

# Week 1

## Linear Regression with one variable (Model Representation)

### Supervised Learning

Given the "Right Answer" for each example in the data.

### Regression Problem

Predict real-value output

Training set of housing prices

Six infant<sup>2</sup> ( $x$ )

price (\$) in 1000's ( $y$ )

2104

460

1416

232

1536

315

852

178

...

...

Notation:

$m$  = Number of training Examples

$x$ 's = "Input" variable/features

$y$ 's = "output" variable / "target" variable

$(x, y)$  — one training example

$(x^{(i)}, y^{(i)})$  —  $i$ th training example

Training Set

How do we represent  $h$ ?

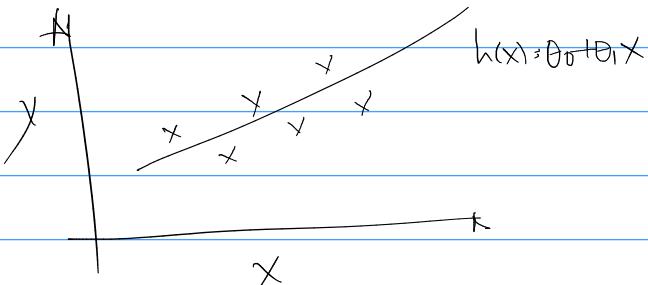
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Shorthand:  $h(x)$

Size of house

$\underline{x}$   
hypothesis function

Estimated price  
(estimated value of  $y$ )



Linear regression with one variable ( $x$ )

Univariate linear regression

## — Linear regression with one variable (Cost function)

Hypothesis:  $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$ 's : Parameters

How do we choose  $\theta_i$ 's?

Idea: Choose  $\theta_0, \theta_1$  so that  $h_\theta(x)$  is close to  $y$  for our training examples  $(x_i, y_i)$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ & h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)} \\ & J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ & \text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1) \\ & \text{Cost function} \\ & \text{Square error function} \end{aligned}$$

# Linear Regression with one variable (Cost function intuition I)

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

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

Goal:

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

Simplified

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

$$\theta_0 = 0$$

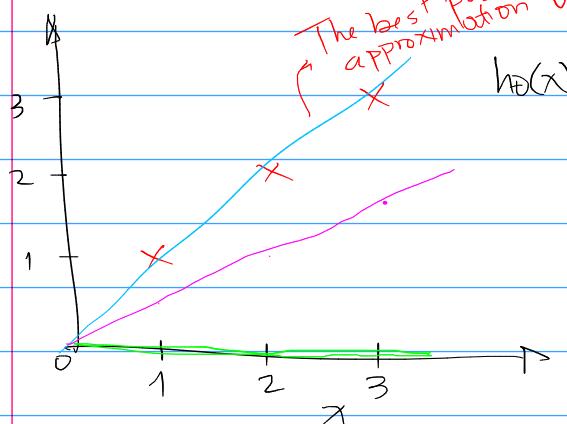
$$\theta_1$$

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

$$\underset{\theta_1}{\text{minimize}} \quad J(\theta_1)$$

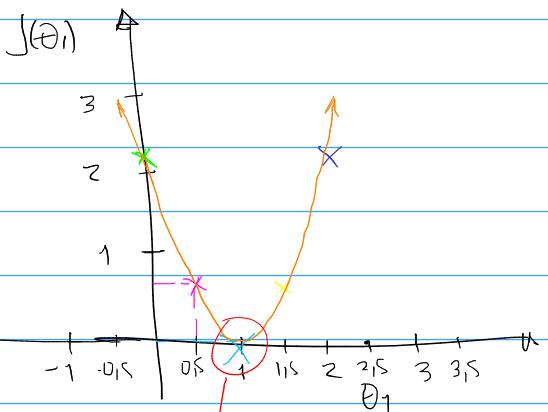
$$h_{\theta}(x)$$

(for fixed  $\theta_1$ , this is a function of  $x$ )



$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )



If  $\theta_1 = 1$  when  $m=3$

$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0 \end{aligned}$$

If  $\theta_1 = 0.5$  when  $m=3$

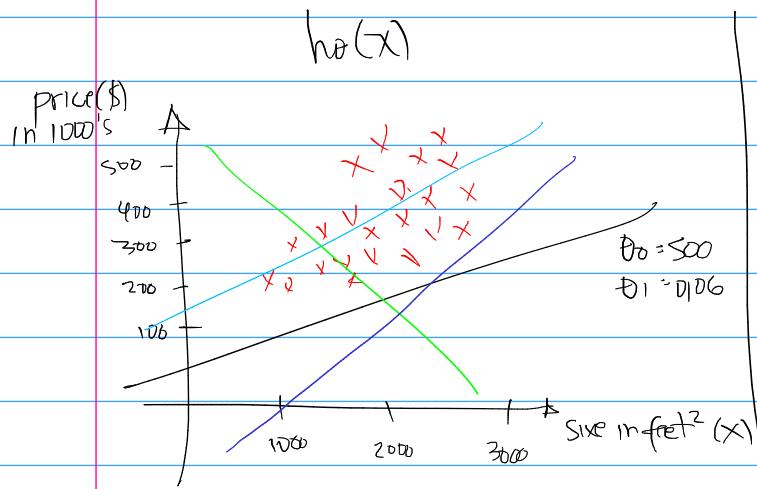
$$\begin{aligned} J(0.5) &= \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \\ &= \frac{1}{2 \cdot 3} \cdot (3.5) \approx 0.58 \end{aligned}$$

If  $\theta_1 = 0$

$$\begin{aligned} J(0) &= \frac{1}{2m} [(0-1)^2 + (0-2)^2 + (0-3)^2] \\ &= \frac{1}{2 \cdot 3} [1+4+9] = \frac{1}{2 \cdot 3} \cdot [14] = \frac{14}{6} = 2.33 \end{aligned}$$

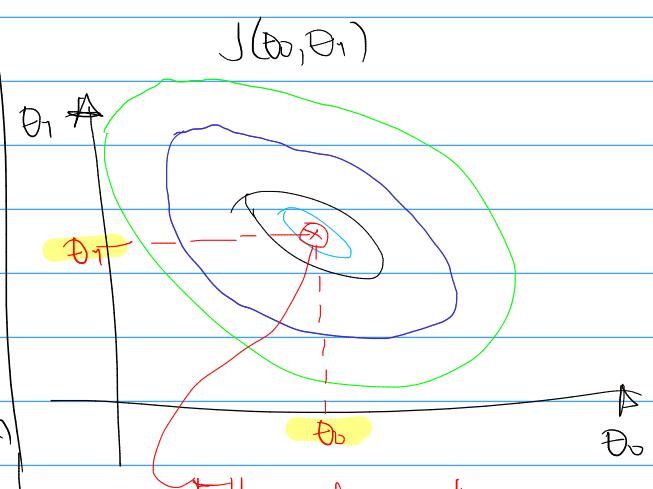
$\star \theta_1 = 1$  we found the minimum and therefore we find the best possible approximation

## Linear regression with one variable (Cost function Intuition II)



$$h_0(x) = 500 + 0.06x$$

Each coordinate  $(\theta_0, \theta_1)$  correspond to an hypothesis function  $h_0(x) = \theta_0 + \theta_1 x$



This value is the minimum and therefore, its  $\theta_0$  and  $\theta_1$  coordinates are the one that makes a better approximation to the data.

## Linear Regression with one Variable [Gradient descent]

Gradient descent - The most common algorithm

Have some function  $J(\theta_0, \theta_1)$

Want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

### Outline

- Start with some  $\theta_0, \theta_1$  (Say  $\theta_0=0, \theta_1=0$ )
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at a minimum

### Gradient descent Algorithm

Repeat until convergence h

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

} learning rate

derivative term

Simultaneously update  $\theta_0, \theta_1$ .

Correct: Simultaneous Update

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

Incorrect

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

$$\theta_0 := \text{temp}_0$$

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

$$\theta_1 := \text{temp}_1$$

### Exercise:

Suppose  $\theta_0=1, \theta_1=2$  We simultaneously update  $\theta_0, \theta_1$  using the rule  $\theta_j := \theta_j + \sqrt{\theta_0 \theta_1}$   
(for  $j=0$  and  $j=1$ )

$$\text{temp}_0 := \theta_0 + \sqrt{\theta_0 \theta_1} = 1 + \sqrt{2}$$

$$\text{temp}_1 := \theta_1 + \sqrt{\theta_0 \theta_1} = 2 + \sqrt{2}$$

$$\theta_0 = 1 + \sqrt{2}$$

$$\theta_1 = 2 + \sqrt{2}$$

## — Linear Regression with one variable (Gradient descent Intuition)

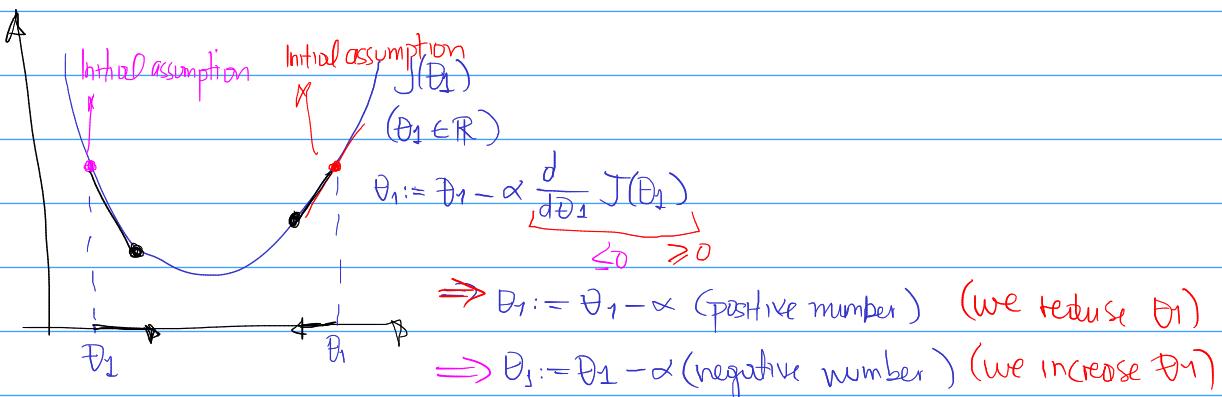
Repeat until convergence h

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{Simultaneously update } \theta_0 \text{ and } \theta_1)$$

]

We assume

$$\min_{\theta_1} J(\theta_1) \quad (\text{just one parameter})$$



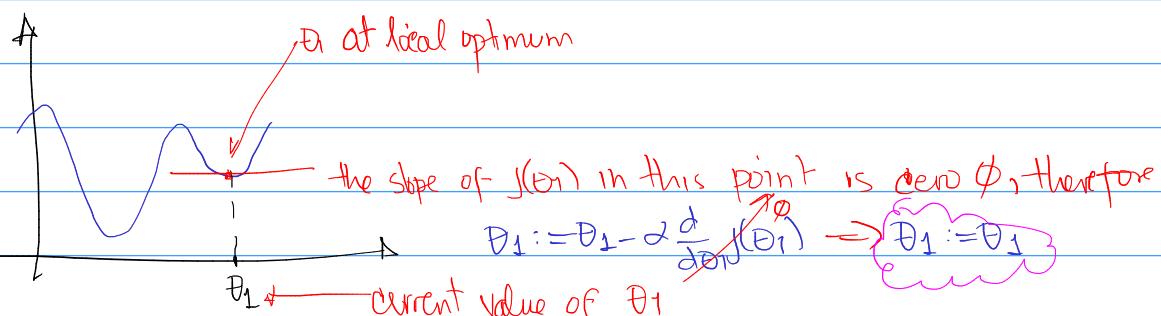
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

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

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

### Exercise

Suppose  $\theta_1$  is at a local optimum of  $J(\theta_1)$ , such as shown. What will one step of gradient descent  $\theta_1 := \theta_1 - \alpha \frac{d}{d \theta_1} J(\theta_1)$  do?



Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed. As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.

## Linear Regression with one variable (Gradient descent for linear regression)

### Gradient Descent Algorithm

Repeat until convergence {

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

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

}

### Linear Regression Model

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

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

We are going to use gradient descent to  
minimize  $J(\theta_0, \theta_1)$

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

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

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

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

### Gradient Descent Algorithm for Linear Regression

Repeat until convergence {

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

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

The function  $J(\theta_0, \theta_1)$  for linear regression is always convex, that's why  
the minimization is always a global optimum.

Gradient Descent also called "Batch" Gradient Descent

"Batch" : Each step of gradient descent uses all the training examples.

## — Questionario: Calificando (Linear Regression with One Variable)

①

X: Number of "A" grades in first year	
5	
3	
0	
4	

Y: The number of "A" grades in second year	
4	
4	
1	
3	

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

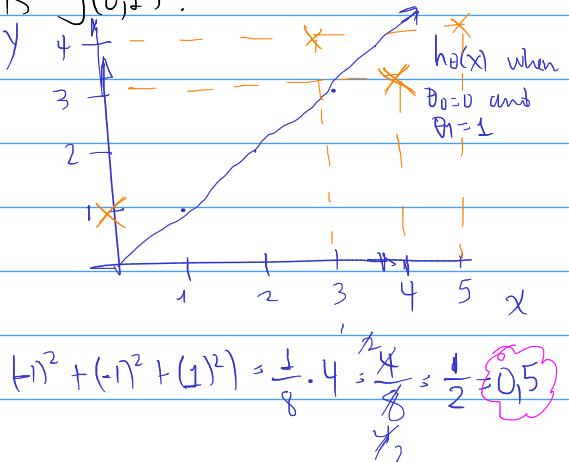
$\boxed{\theta_0=4}$  # of training examples

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

$$\text{If } \theta_0 = 0 \wedge \theta_1 = 1 \Rightarrow h_{\theta}(x) = 0 + 1x = x$$

$$\begin{aligned} J(\theta_0, \theta_1) &= \frac{1}{2(4)} \sum_{i=1}^4 (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{8} \cdot ((1x^{(1)} - 4)^2 + (1x^{(2)} - 4)^2 + (1x^{(3)} - 1)^2 + (1x^{(4)} - 3)^2) \\ &= \frac{1}{8} \cdot ((-3)^2 + (-3)^2 + (0)^2 + (1)^2) = \frac{1}{8} \cdot (9 + 9 + 0 + 1) = \frac{1}{8} \cdot 19 = \frac{19}{8} = 2.375 \end{aligned}$$

What is  $J(0, 1)$ .



③ Suppose we set  $\theta_0 = -2$ ;  $\theta_1 = 0.5$ . What is  $h_{\theta}(6)$ ?

$$\text{If } \theta_0 = \cancel{-2} \wedge \theta_1 = 0.5 \Rightarrow h_{\theta}(x) = -2 + 0.5x \Rightarrow h_{\theta}(6) = -2 + 0.5(6) = -2 + 3 = 1$$

$\cancel{-2}$  You idiot!

④ a. and b are ~~truth~~

$$⑤ J(\theta_0, \theta_1) = 0$$

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

For  $J(\theta_0, \theta_1) = 0 \Leftrightarrow$  Our training set can fit perfectly by a straight line.

$$⑥ \text{for a) } y(1) = -569.6 + 530.9(1) = -38.7$$

$$\text{for b) } y(1) = -569.6 - 530.9(1) = -1100.5$$

$$⑦ \theta_0 = -1; \theta_1 = 2$$

$$h_{\theta}(x) = -1 + 2x \Rightarrow h_{\theta}(6) = -1 + 2(6) = -1 + 12 = 11$$

⑧ c and d are correct

# Week 2

## Linear Regression with multiple variables (multiple features)

Size (feet <sup>2</sup> ) $x_1$	Number of Bedrooms $x_2$	Number of floors $x_3$	Age of home (years) $x_4$	Price (\$1000) $y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

Notation: for example  $\vec{x}^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^4$  Now  $\vec{x}^{(i)}$  is a vector of the features for each  $y^{(i)}$

$\Rightarrow x_2^{(2)} = 2$

$n = \text{number of features}$

$\vec{x}^{(i)} = \text{input (features) of } i^{\text{th}} \text{ training example}$

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

Hypothesis:

Previously:  $h_{\theta}(x) = \theta_0 + \theta_1 x$  → We can't use this form because now we have multiple features.

$h_{\theta}(x) \rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots + \theta_n x_n$ . → New form of hypothesis

$$\text{E.g. } h_{\theta}(x) = 80 + 0.1x_1 + 0.1x_2 + 3x_3 - 2x_4$$

For convenience of notation, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\underbrace{[\theta_0 \cdot \theta_1 \dots \theta_n]}_{\theta^T} \quad \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

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

Multivariate Linear Regression

## — Linear Regression with multiple variables (Gradient descent for multiple variables)

Hypothesis:  $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters:  $\theta_0, \theta_1, \dots, \theta_n = \theta$   $n+1$ -dimensioned vector

Cost function:  $J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$  or  $J(\theta) =$   
 $\downarrow \theta^T x^{(i)}$

Gradient Descent:

Repeat {

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

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

New algorithm ( $n \geq 1$ )

Repeat {

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

}

+ just the feature that  
multiplies  $\theta_j$

(Simultaneously update  
 $\theta_j$  for  $j=0, \dots, n$ )

Linear Regression with multiple variables (Gradient descent in practice 1:  
Feature Scaling)

## Feature Scaling

Idea: Make sure features are on a similar scale

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

→ We can scale the features

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

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

That way gradient  
descent converge quickly.

$$0 \leq x_1 \leq 1$$

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

$$0 \leq x_2 \leq 1$$

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

## Mean Normalization

Replace  $x_i$  with  $\frac{x_i - \mu_i}{\sigma_i}$  to make features have approximately zero mean

(Do not apply to  $x_0 = 1$ )

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

avg value

$$-0,5 \leq x_1 \leq 0,5$$

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

$$-0,5 \leq x_2 \leq 0,5$$

range (max-min)

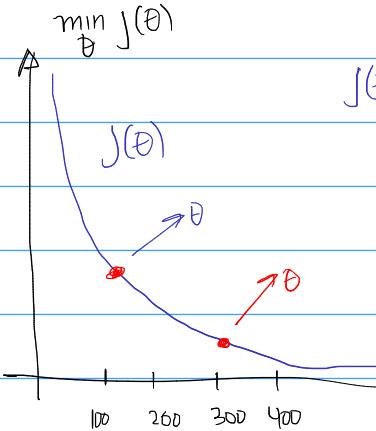
(standard deviation)

- Linear Regression with multiple variables (Gradient descent in practice II: learning Rate)

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

- "Debugging": How to make sure gradient descent is working correctly
- How to choose learning rate  $\alpha$ .

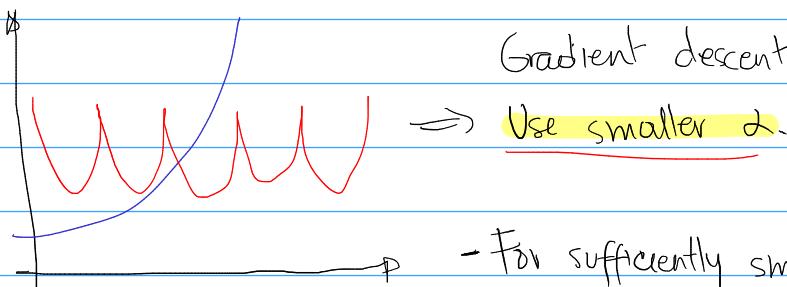
- Making Sure gradient descent is working correctly -



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

Example automatic convergence test:

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



Gradient descent not working.

⇒ Use smaller  $\lambda$ .

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

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

To choose  $\alpha$ , try:

..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...

— Linear Regression with multiple variables (Features and polynomial regression)

Housing prices prediction —

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

$$x = \text{frontage} * \text{depth}$$

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

↓  
land area.

Polynomial regression —

If we think we should model our hypothesis these ways:

Then  $\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$  ✓  $\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

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

where  $x_1 = (\text{size})$   $x_3 = (\text{size})^3$

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$
$$x_2 = (\text{size})^2$$

$$\text{size: } 1 - 1000$$

$$\text{size}^2: 1 - 1000000$$

$$\text{size}^3: 1 - 10^9$$

## Linear Regression with multiple variables (Normal Equation)

Method to solve for  $\theta$  analytically.

Intuition: If 1D ( $\theta \in \mathbb{R}$ )  $\Rightarrow \frac{d}{d\theta} J(\theta) = \dots \stackrel{\text{set } 0}{=} 0$  } This way we found  
 $J(\theta) = a\theta^2 + b\theta + c$  solve for  $\theta$  an optimum.

Now:

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

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

Solve for  $\theta_0, \theta_1, \dots, \theta_n$ .

For example  $m=4$

$$X = \begin{bmatrix} x_0^{(i)} & 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \quad m \times (n+1)$$

each column is a feature  $x_j$

Each row is a feature  $x^{(i)}$

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

$m$ -dimensional vector

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

Optimum values for  $\theta$

Feature Scaling is not necessary with this method.

$m$  training examples,  $n$  features.

### Gradient Descent

- Need to choose  $\alpha$
- Needs many iterations
- Works well even when  $n$  is large

### Normal Equation

- No need to choose  $\alpha$
- Don't need to iterate
- Need to compute  $(X^T X)^{-1}$
- Slow if  $n$  is very large.  $n \geq 10,000$

so much computational cost

Submit()  $\rightsquigarrow$  function to submit your solution.

## Quiz

①  ~~$x_2^{(1)}$~~  = 4761  $\Rightarrow x_2 = \begin{bmatrix} 7921 \\ 5184 \\ 8836 \\ 4761 \end{bmatrix}$

$$\mu_2 = \frac{7921 + 5184 + 8836 + 4761}{4} = 6700.5$$

$$S_2 = \max - \min = 8836 - 4761 = 4075 \rightarrow$$

$$\text{Normalizing } x_2^{(1)} = \frac{4761 - 6700.5}{4075} = -0.47595 \approx -0.48$$

②  $m=28 \quad n=4$

$$X = \begin{bmatrix} x_0^{(1)} & \dots & x_4^{(1)} \\ x_0^{(2)} & & \vdots \\ \vdots & & \ddots \\ x_0^{(28)} & & x_4^{(28)} \end{bmatrix} \quad 28 \times 5 \text{ matrix} \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(28)} \end{bmatrix} \quad 28 \text{-dimensional vector} \quad 28 \times 1 \text{ matrix}$$

$$X^T \cdot X = 5 \times 28 \text{ matrix} \cdot 28 \times 5 \text{ matrix} \rightarrow 5 \times 5 \text{ matrix}$$

$$(X^T \cdot X)^{-1} = 5 \times 5 \text{ matrix}$$

$$(X^T \cdot X)^{-1} X^T = 5 \times 5 \text{ matrix} \cdot 5 \times 28 \text{ matrix} = 5 \times 28 \text{ matrix}$$

$$(X^T \cdot X)^{-1} X^T \cdot y = 5 \times 28 \text{ matrix} \cdot 28 \times 1 \text{ matrix} = 5 \times 1 \text{ matrix}$$

④ The answer is ①

⑤ The answer is ②

①  $x_1^{(3)} = \begin{bmatrix} 94 \\ 8836 \end{bmatrix}$   $x_1^{(3)} = 94$   $x_1 = \begin{bmatrix} 89 \\ 72 \\ 94 \\ 69 \end{bmatrix}$

$$\mu_1 = \frac{89 + 72 + 94 + 69}{4} = 81$$

$$\mu_1 = \frac{94 - 69}{2} = 12.5$$

$$S_1 = 94 - 69 = 25 \Rightarrow \text{Normalizing } x_1^{(3)}, \frac{94 - 12.5}{25} = 3.26 \times \text{value}$$

$$\text{Normalizing } x_1^{(3)} = \frac{94 - 81}{25} = 0.52$$

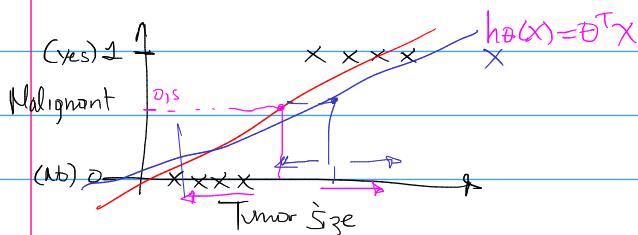
## Week 3

### - Logistic Regression (Classification)

Email: Spam / Not Spam?

Online transactions: Fraudulent (Yes/No)?

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



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

If  $h_{\theta}(x) \geq 0.5$ , predict ' $y=1$ '

If  $h_{\theta}(x) \leq 0.5$ , predict ' $y=0$ '

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

Classification:  $y = 0$  or  $1$

$h_{\theta}(x)$  can be  $> 1$  or  $\leq 0$

## — Logistic Regression (Hypothesis Representation)

Logistic Regression Model

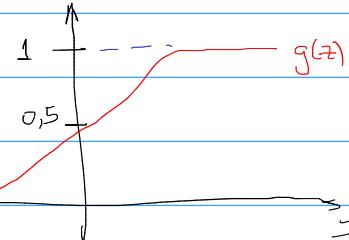
Want  $0 \leq h_{\theta}(x) \leq 1$

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

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

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

{ Sigmoid function  
} Logistic function



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

Example: If  $X = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix}$

We need to fit the parameter  $\theta$

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

Tell the patient that 70% chance of tumor being malignant

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

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

$y=0$  or  $1$

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

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

## - Logistic Regression (Decision boundary)

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

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

$\theta^T x \geq 0$

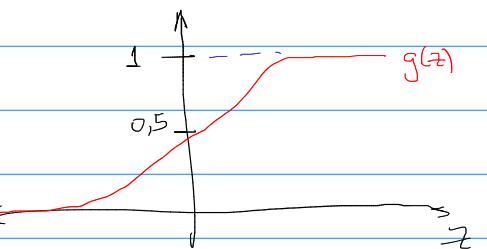
Suppose predict ' $y=1$ ' if  $h_{\theta}(x) \geq 0.5$

predict ' $y=0$ ' is  $h_{\theta}(x) < 0.5$

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

$$\theta^T x < 0$$

$$g(z) < 0.5$$



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

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

## Decision Boundary

$$x_2$$

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

$$x_1$$

$$x_1, x_2$$

Predict ' $y=1$ ' if  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$$\theta^T x$$

$$x_1 + x_2 \geq 3$$

$$x_1 + x_2 < 3$$

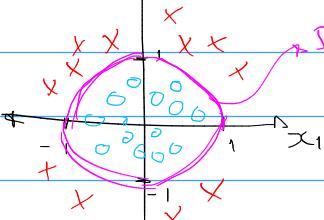
$$\begin{cases} h_{\theta}(x) = 0.5 \\ x_1 + x_2 = 3 \end{cases}$$

## Non-linear decision boundaries

$$x_2$$

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

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



Decision boundary

Predict ' $y=1$ ' if  $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 \geq 0$

$$x_1^2 + x_2^2 = 1$$

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

→ Logistic regression (Cost function)

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

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

$$\text{cost}(h_\theta(x^{(i)}) - y^{(i)}) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2 \rightsquigarrow \text{Non convex function for } \frac{1}{1+e^{-\theta}x}$$

Logistic Regression cost function

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

$h_\theta(x) \rightarrow 0 \rightsquigarrow \text{cost} \rightarrow \infty$

$h_\theta(x)=1$

$h_\theta(x) \rightarrow 1 \rightsquigarrow \text{cost} \rightarrow \infty$

— Logistic Regression (Simplified cost function and gradient descent)

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) \quad \text{Better form to express cost function}$$

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

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

To fit parameters  $\theta$

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

To make a prediction given new  $x$ :

$$\text{output } h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \quad P(y=1|x, \theta)$$

### 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  $\min_{\theta} J(\theta)$  Replacing (2) in (1)

Repeat  $\tilde{\theta}_j$

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

$$(1) \quad \tilde{\theta}_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{simultaneously update all } \theta_j)$$

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

Trying to obtain the  $\frac{\partial}{\partial \theta_j} J(\theta) \Rightarrow \frac{\partial}{\partial \theta_j} y^{(i)} \log(h_\theta(x^{(i)})) = \frac{y^{(i)}}{h_\theta(x^{(i)})} \cdot \frac{\partial}{\partial \theta_j} \left( \frac{1}{1+e^{-\theta^T x}} \right)$

$$\Rightarrow \frac{y^{(i)}}{h_\theta(x^{(i)})} \cdot \left( -\frac{1}{(1+e^{-\theta^T x})^2} \right) \cdot \frac{\partial}{\partial \theta_j} (1+e^{-\theta^T x}) = \frac{y^{(i)}}{h_\theta(x^{(i)})} \cdot \frac{1}{(1+e^{-\theta^T x})^2} \cdot e^{-\theta^T x} \cdot (-x_j) \rightarrow h_\theta(x^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) = \frac{(1-y^{(i)})}{1-h_\theta(x^{(i)})} \cdot \frac{\partial}{\partial \theta_j} (-h_\theta(x^{(i)})) = \frac{(1-y^{(i)})}{1-h_\theta(x^{(i)})} \cdot \frac{\partial}{\partial \theta_j} (-(1+e^{-\theta^T x})^{-1})$$

$$= \frac{(1-y^{(i)})}{1-h_\theta(x^{(i)})} \cdot (1) \cdot (1+e^{-\theta^T x})^2 \cdot \frac{\partial}{\partial \theta_j} (1+e^{-\theta^T x}) = \frac{(1-y^{(i)})}{1-h_\theta(x^{(i)})} \cdot \frac{1}{(1+e^{-\theta^T x})^2} \cdot e^{-\theta^T x} \cdot (-x_j) \rightarrow h_\theta(x^{(i)})^2$$

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

$$m = h_\theta(x^{(i)})^2 \cdot e^{-\theta^T x} \cdot (-x_j)$$

$$\frac{-y^{(i)}[1-h_\theta(x^{(i)})] + h_\theta(x^{(i)})[1-y^{(i)}]}{1-h_\theta(x^{(i)})} \cdot h_\theta(x^{(i)}) \cdot e^{-\theta^T x} \cdot (-x_j)$$

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

$$\Rightarrow \left[ -y^{(i)} + y^{(i)} h_\theta(x^{(i)}) + h_\theta(x^{(i)}) - y^{(i)} h_\theta(x^{(i)}) \right] (-x_j) = - (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)} \quad (2) \quad \text{replace (2) in (1)}$$

## - Logistic Regression (Advanced Optimization)

### - Optimization Algorithms

- Gradient Descent

- Conjugate gradient }

- BFGS }

- L-BFGS }

### Advantage

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

### Disadvantage

- More complex

### Examples

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

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

$$\theta_1 = 5$$

$$\theta_2 = 5$$

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

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

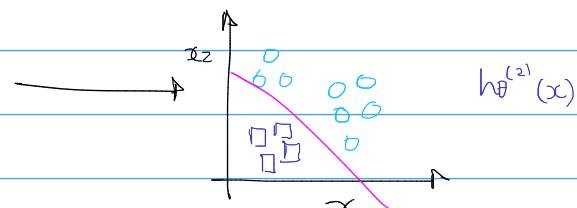
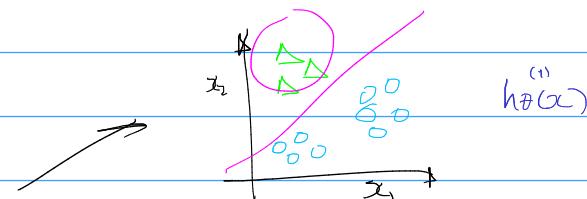
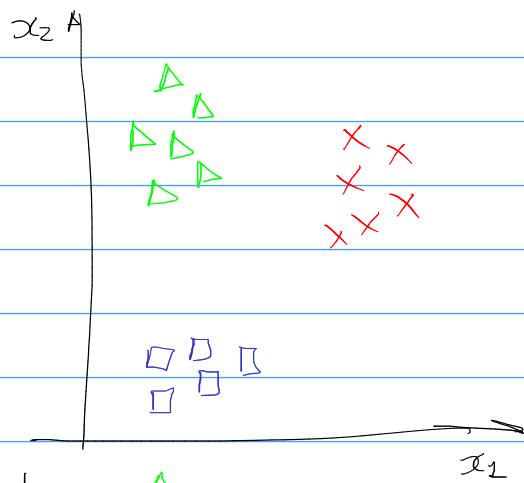
$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

### Optimal Example

## - Logistic Regression [Multi-class classification: one-vs-all]

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

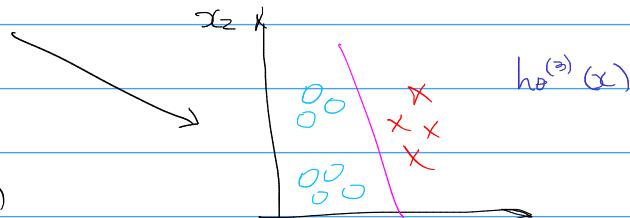
$$y=1 \quad y=2 \quad y=3 \quad y=4$$



Class 1: ▲

Class 2: □

Class 3: ×



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

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.

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

## - Qn 3 -

(5)  $\theta_0 + \theta_1 x_1 + \theta_2 x_2$  when  $\theta_0 = -6$ ;  $\theta_1 = 1$ ;  $\theta_2 = 0$

$$-6 + x_1 = 0$$

$$\boxed{x_1 = 6}$$

$$6 - x_1 = 0$$

$$\boxed{x_1 = 6}$$

$$6 - x_1 > 0$$

$$\boxed{6 > x_1}$$

(1)  $P(y=0|x; \theta) = 0.4$   $\text{and } h_{\theta}(x) = 0.4$   $P(y=1|x; \theta) = 0.4$   
 ~~$P(y=1|x; \theta) = 0.6$~~   $P(y=0|x; \theta) = 0.6$

(2)  $\circ$  adding polynomial features could increase how well we can fit the training data  $\circ$  At optimal value of  $\theta$ , we will have  $J(\theta) \geq 0$

(3) for  $\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

↳

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{1+e^{-\theta_j x_i}} - y^{(i)} \right) x_j^{(i)} \quad \theta := \theta - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

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

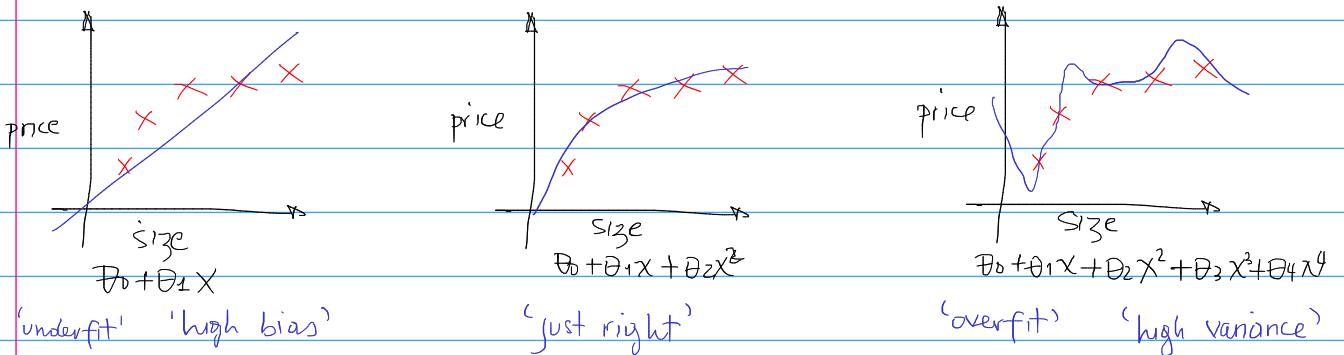
$$\theta := \theta - \alpha \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{1+e^{-\theta x_i}} - y^{(i)} \right) x^{(i)}$$

(4)  $\circ$  sigmoid is never greater than one

$\circ$  The cost  $J(\theta)$  for logistic is always  $\uparrow$  or equal to zero

$\circ$  The one vs all allows to use logistic regression for problems in which each  $y^{(i)}$  comes from a fixed, discrete set of values.

## → Regularization (The problem of overfitting)



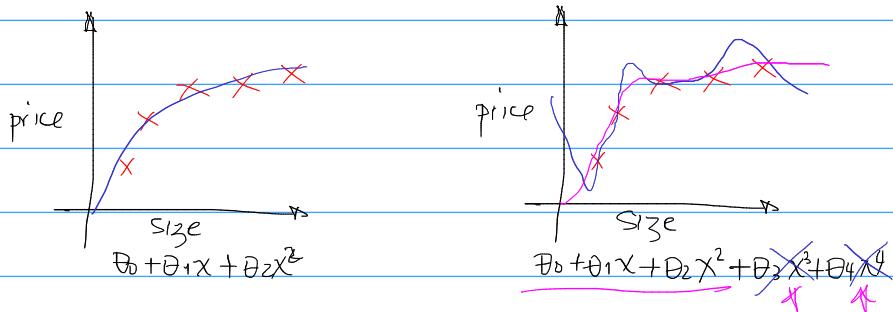
Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples)

### Addressing 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  $y$ .

## → Regularization (Cost function)



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

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2 \Rightarrow \underline{\theta_3 \approx 0} \wedge \underline{\theta_4 \approx 0}$$

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

- 'Simpler' hypothesis
- Less prone to overfitting

Example.

Housing

- Features:  $x_1, x_2, \dots, x_{100}$

- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

Regularized cost function

↳ regularization parameters

If  $\lambda$  is extremely large value  $\Rightarrow$  'under-fit' prediction

## — Regularization (Regularized linear regression)

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

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

Gradient descent:

Repeat h

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

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

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

Normal Equation

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

$$\Rightarrow \theta = (X^\top X + \lambda \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix})^{-1} X^\top y$$

— Regularization (Regularized Logistic Regression)

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

Gradient descent

repeat {

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

$\nearrow \frac{\partial}{\partial \theta_0} J(\theta)$  and  $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

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

## Quiz Regularization

- ① - adding a new feature to the model always results in equal or better performance on the training set ✓  
- introducing regularization to the model always results in equal or better performance on examples not in the training set. \*① X
- ②  $\theta = \begin{bmatrix} 1.37 \\ 0.51 \end{bmatrix}$  80%.
- ③ - Consider a classification problem. Adding regularization may cause your classifier to incorrectly classify some training examples (which it had correctly classified when not using regularization, i.e. when  $\lambda=0$ )
- ④, ⑤ overfit and underfit graphics.
- ① - Adding many new features to the model make it more likely to overfit the training set.  
- \*① X
- ②  $\theta = \begin{bmatrix} 2.75 \\ 1.32 \end{bmatrix}$
- ③ using too large values of  $\lambda$  can cause your hypothesis to underfit the data.
- ④⑤ graphics