

Machine Learning

## Linear Regression with multiple variables

### Multiple features

### Multiple features (variables).

Size (feet²)	Price (\$1000)	
original case, we only consider size of the house one variable ${m x}$	71 🚄	
2104	460	
1416	232	
1534	315	
852	178	
•••	•••	

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

### Multiple features (variables).

<i>_</i>	Size (feet²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)	
_	(>1)	×3)	<b>*</b> 3	<b>₹</b>	<u> </u>	
	2104	5	1	45	460 7	
	1416	3	2	40	232 - M=	47
1	1534	3	2	30	315	
	852	2	1	36	178	
		•••	•••		] ,	
No	tation:	<b>*</b>	*	<b>1</b>	~ (2) =	14167
		mber of fea		n=4		2 €
7	$\Rightarrow x^{(i)} = \inf$	out (feature	es) of $\it i^{th}$ tra	aining example	e. (2)	-407
	$\Rightarrow x_j^{(i)} = va$	lue of featu	re $j$ in $i^{th}j$	training examp	$\frac{1}{2}$ ole. $\frac{1}{2}$	

### Hypothesis:

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

Multivariate linear regression.



Machine Learning

## 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: 
$$\theta_0, \theta_1, \dots, \theta_n$$

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

#### **Gradient descent:**

Repeat 
$$\{$$
  $\Rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$ .  $\Im$  (e)  $\}$  (simultaneously update for every  $j = 0, \dots, n$ )

#### **Gradient Descent**

Previously (n=1):

$$:= \theta_0 - o \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right]$$

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

(simultaneously update  $\theta_0, \theta_1$ )

New algorithm  $(n \ge 1)$ :





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

(simultaneously update  $heta_j$  for  $j=0,\ldots,n$ 

$$\theta_0 := \theta_0 -$$

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

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$



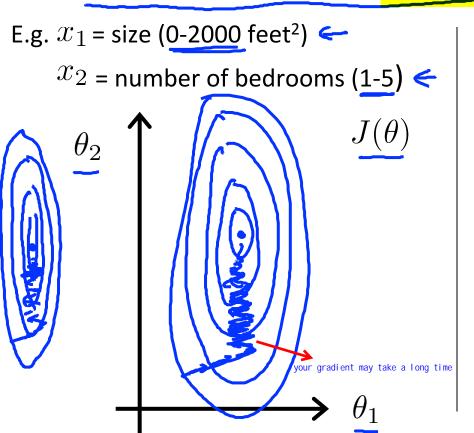
Machine Learning

# Linear Regression with multiple variables

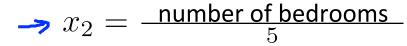
Gradient descent in practice I: Feature Scaling

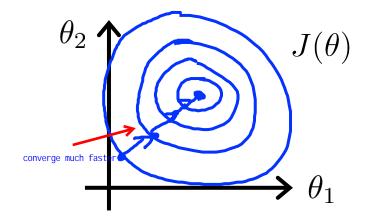
### **Feature Scaling**

Idea: Make sure features are on a similar scale. we scale features, apples to apples comparison



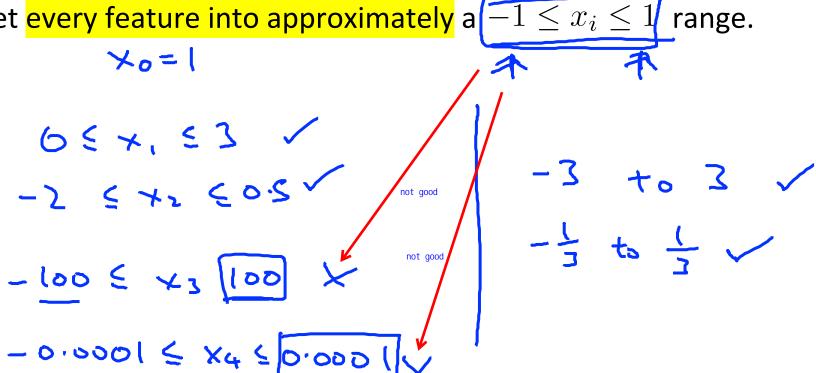
$$x_1 = \frac{\text{size (feet}^2)}{2000}$$





### **Feature Scaling**

Get every feature into approximately a



#### **Mean normalization**

Replace  $\underline{x}_i$  with  $\underline{x}_i - \mu_i$  to make features have approximately zero mean (Do not apply to  $\underline{x}_0 = 1$ ).

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

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

$$-0.5 \le x_1 \le 0.5, -0.5 \le x_2 \le 0.5$$

$$x_1 = \frac{x_1 - y_1}{2000}$$

$$y_2 = \frac{y_2 - y_2}{5}$$

$$y_3 = \frac{y_3 - y_4}{5}$$

$$y_4 = \frac{y_4 - y_5}{5}$$

$$y_5 = \frac{y_5 - y_5}{5}$$



Machine Learning

# Linear Regression with multiple variables

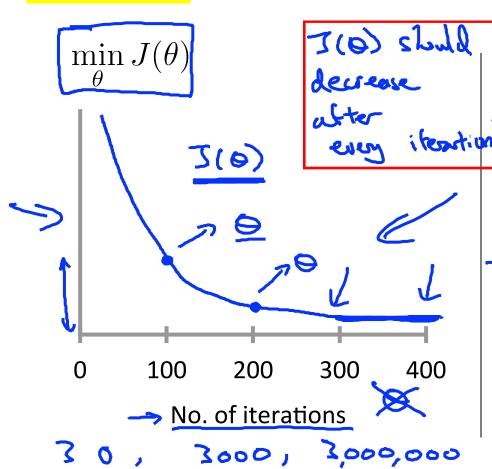
Gradient descent in practice II: Learning rate

### **Gradient descent**

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

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

### Making sure gradient descent is working correctly.

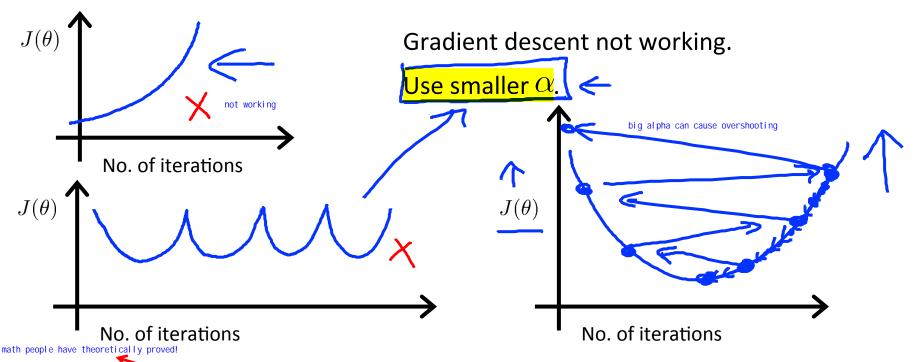


Example automatic convergence test:

should decrease

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

### Making sure gradient descent is working correctly.

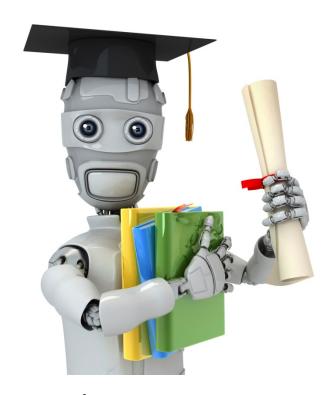


- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.
- But if lpha is <mark>too small</mark>, gradient descent can be slow to converge.

### **Summary:**

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

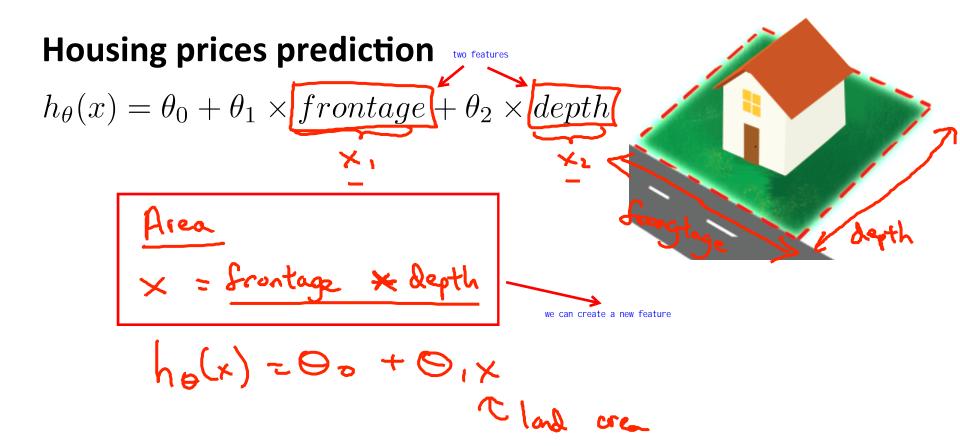
To choose  $\alpha$ , try



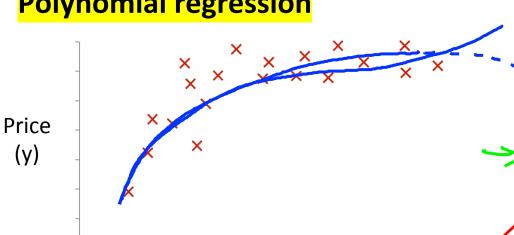
Machine Learning

# Linear Regression with multiple variables

Features and polynomial regression



### **Polynomial regression**



Size (x)

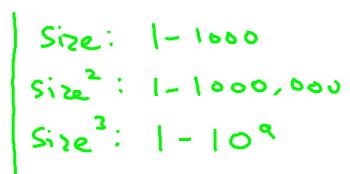
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$$

$$\Rightarrow x_1 = (size)$$

$$\rightarrow x_2 = (size)^2$$

$$x_3 = (size)^3$$



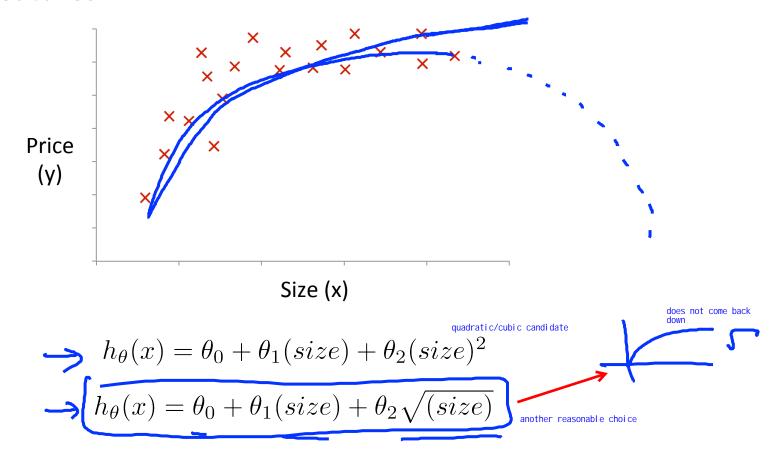
 $\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$ 

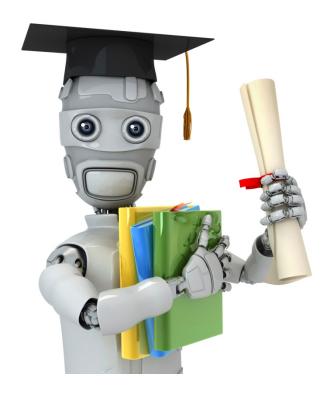
 $\theta_0 + \theta_1 x + \theta_2 x^2$ 

fit the cubic model

choose your features

#### **Choice of features**



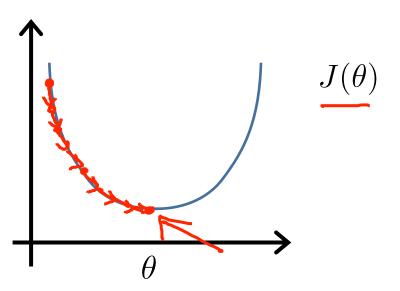


Machine Learning

### Linear Regression with multiple variables

### Normal equation Analytically solve the problem/pseudo inverse

### **Gradient Descent**

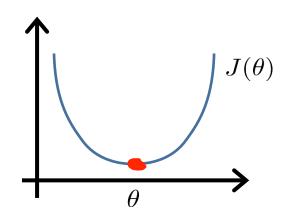


Normal equation: Method to solve for  $\theta$  analytically.

### Intuition: If 1D $(\theta \in \mathbb{R})$

$$J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \underbrace{\text{Set}}_{\theta} O$$
Solve by



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

$$\underline{\frac{\partial}{\partial \theta_j} J(\theta)} = \cdots \stackrel{\text{set}}{=} 0 \quad \text{(for every } j\text{)}$$

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

analytical solution from calculus

### Examples: $\underline{m} = 4$ .

J	Size (feet²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)	)
$\rightarrow x_0$	$x_1$	$x_2$	$x_3$	$x_4$	y	
1	2104	5	1	45	460	7
1	1416	3	2	40	232	
1	1534	3	2	30	315	
1	852	2	_1	<b>3</b> 6	178	7
	$X = \begin{bmatrix} 1 & 1 \end{bmatrix}$	2104 5 1 1416 3 2 1534 3 2 852 2 1 $(x)^{-1}X^{T}y$	2 40 2 30 1 36 <b>y</b> -	y =	460 232 315 178	√est or

pseudo inverse of X

### $m = xamples (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}); n = xamples (x^{(m)}, y^{(m)}); n =$

Andrew Ng

$$\underbrace{\theta = \underbrace{(X^T X)^{-1} X^T y}}_{(Y^T Y)^{-1} \dots}$$

 $(X)^{-1}$  is inverse of matrix  $X^TX$ .

$$\frac{A: \times^{T} \times}{\left( \times^{T} \times \right)^{-1}} = A^{-}$$

### m training examples, n features.

very useful!!! Gradient Descent

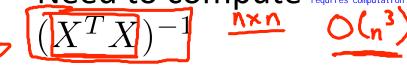


- $\rightarrow$  Need to choose  $\alpha$ .
- → Needs many iterations.
  - Works well even when n is large.



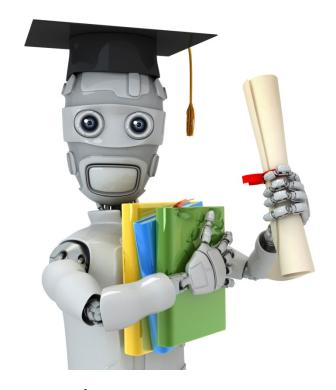
### **Normal Equation**

- $\rightarrow$  No need to choose  $\alpha$ .
- Don't need to iterate.
  - Need to compute requires computation



Slow if n is very large.

$$N = 1000$$
 $N = 10000$ 
 $M = 100000$ 
 $M = 1000000$ 
 $M = 100000$ 
 $M = 100000$ 



Machine Learning

# Linear Regression with multiple variables

Normal equation and non-invertibility (optional)

### Normal equation

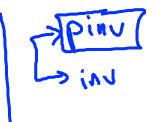
$$\theta = (X^T X)^{-1} X^T y$$



- What if  $X^TX$  is non-invertible? (singular/degenerate)

- Octave: pinv (X'\*X) \*X'\*y

pseudo inverse



### What if $X^TX$ is non-invertible?

Redundant features (linearly dependent).

E.g. 
$$x_1 = \text{size in feet}^2$$

$$x_2 = \text{size in m}^2$$

$$x_1 = (3.28)^2 \times 1$$

$$x_2 = (3.28)^2 \times 2$$

$$x_3 = (3.28)^2 \times 2$$

$$x_4 = (3.28)^2 \times 2$$

- Too many features (e.g.  $m \leq n$ ).
  - Delete some features, or use regularization.