

MACHINE LEARNING (ML-6)

Dr. NEERAJ GUPTA, Department of CEA, GLA University, Mathura

AGENDA

- Linear regression with multiple variable
- Gradient descent for multiple variables
- Gradient descent in practice I: Feature Scaling

SRC : * Andrew NG

Multiple features (variables).

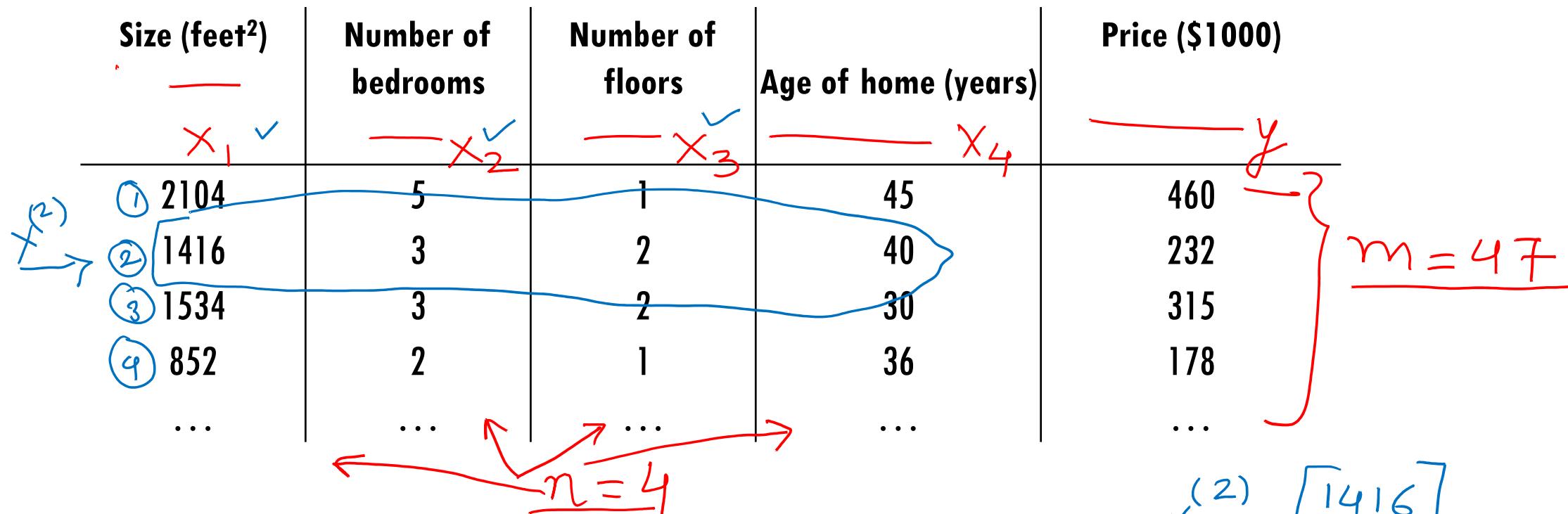
Size (feet ²)	Price (\$1000)
x ↙	y ↙
2104	460
1416	232
1534	315
852	178
...	...

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

Multiple features (variables).

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

Multiple features (variables).



Notation:

$\rightarrow n$ = number of features

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

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

$$\underline{x}^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$x_3^{(2)} = ? = 2$$

Hypothesis:

Previously:

$$h_{\theta}(x) = \underline{\theta_0} + \underline{\theta_1 x} \quad \swarrow \text{single variable}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \rightarrow \underline{\text{multiple}} \underline{\text{features}}$$

Eg. $h_{\theta}(x) = 80 + 0.2 \underline{x_1} + 0.08 \underline{x_2} + 5 \underline{x_3} - 2 \underline{x_4} \quad \swarrow$

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

n-features

For convenience of notation, define $\underline{x_0} = 1.$

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

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

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

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

$$\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

θ^T
 $(n+1) \times 1$ matrix
 $\theta^T \cdot x$

Multivariate linear regression.

Linear Regression with multiple variables

GRADIENT DESCENT FOR MULTIPLE VARIABLES

SRC : * Andrew NG

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n \rightarrow \Theta$ | $n+1 - \text{dim-vector}$

Cost function:

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

$J(\theta)$

Gradient descent:

Repeat {

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

$J(\theta)$ ↘
 α . of feature ↑
 $j = 0, \dots, n$

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

Gradient Descent

Previously ($n=1$):

Repeat {

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

$\downarrow = 0, \dots, n$

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

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

$\frac{\partial}{\partial \theta_1} J(\theta)$

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$) :

Repeat {

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

\uparrow \uparrow

$j=0 \quad j=0$

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

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

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

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

⋮

Linear Regression with multiple variables

GRADIENT DESCENT IN PRACTICE I: FEATURE SCALING

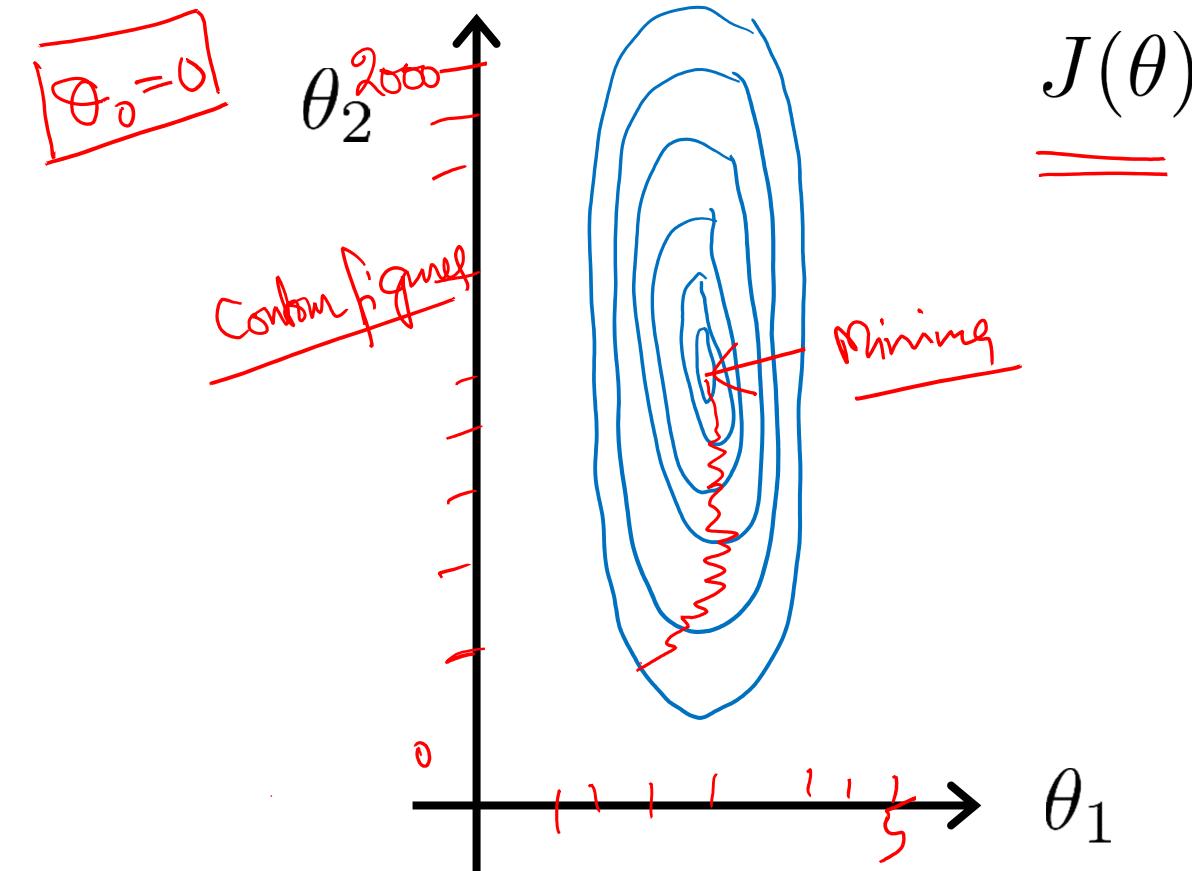
SRC : * Andrew NG

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. x_1 = size (0-2000 feet²) \rightarrow

x_2 = number of bedrooms (1-5)

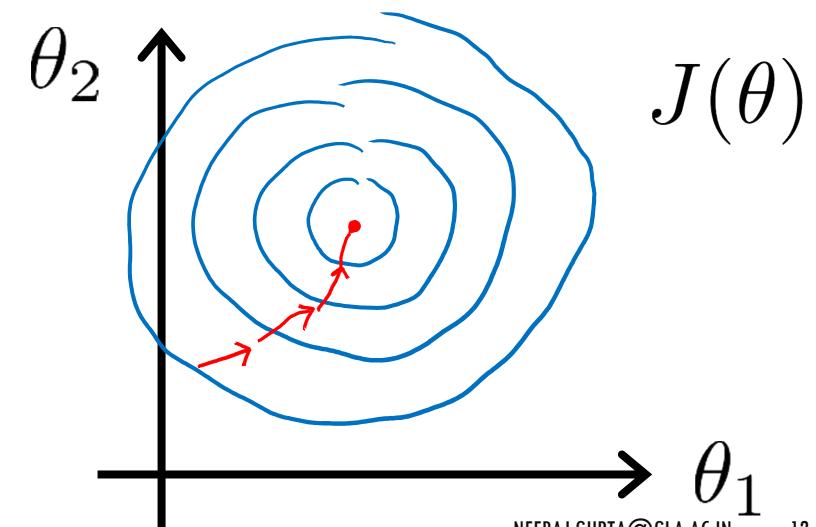


$$x_1 = \frac{\text{size (feet}^2)}{2000}$$

0 - 2000
1 - 5

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

$$0 \leq x_1 \leq 1$$
$$0 \leq x_2 \leq 1$$



Feature Scaling

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

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq 100 \quad \times$$

$$-0.0001 \leq x_4 \leq 0.0001 \quad \times$$

Mean normalization

Replace x_i with $\frac{x_i - \mu_i}{s_i}$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

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

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

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

avg value
of x_i
in training set

Range (Max-Min)

Linear Regression with multiple variables

GRADIENT DESCENT IN PRACTICE II: LEARNING RATE

SRC : * Andrew NG

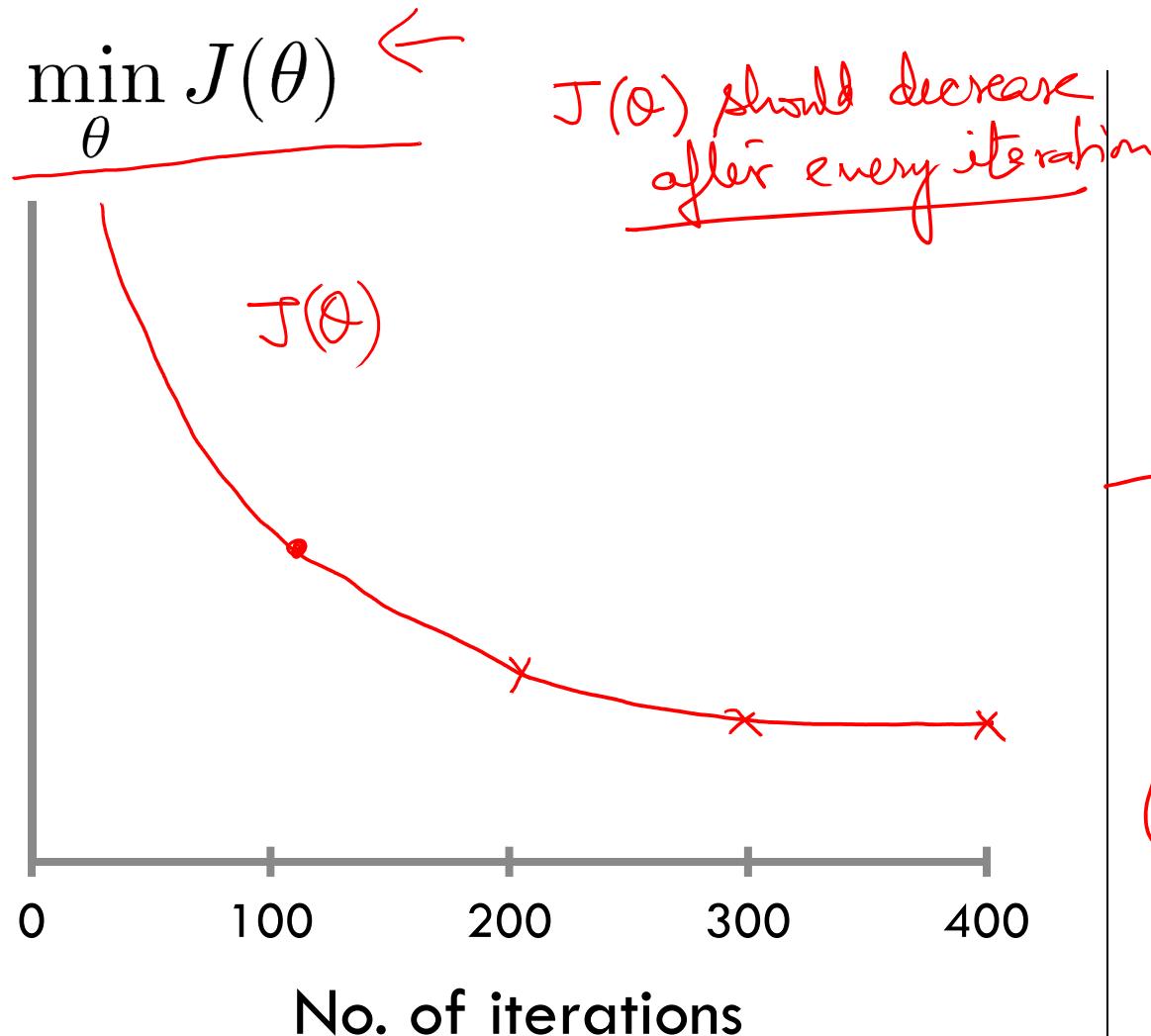
Gradient descent

↓ Learning Rate

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

- “Debugging”: How to make sure gradient descent is working correctly. 
- How to choose learning rate α

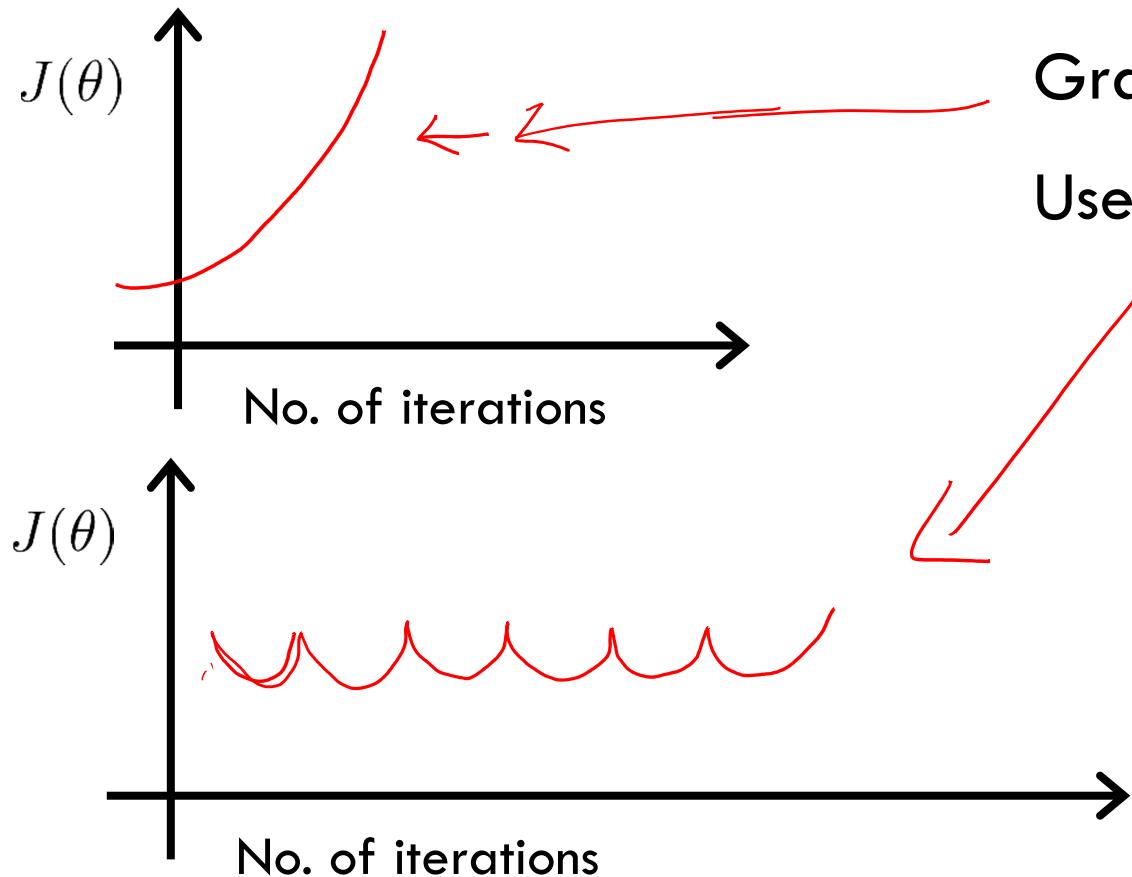
Making sure gradient descent is working correctly.



Example automatic convergence test:

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

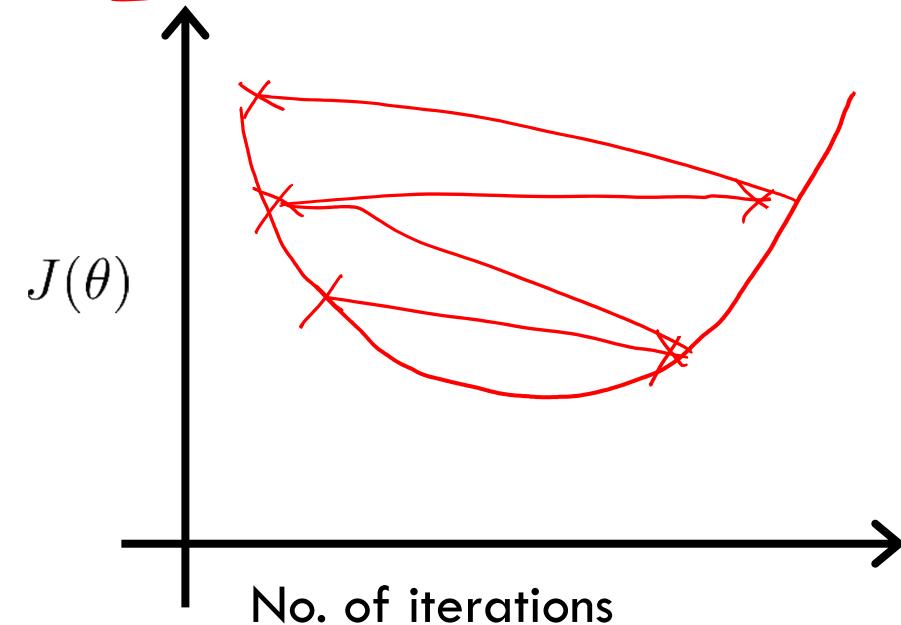
Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller α

learning rate



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

Summary:

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

To choose α , try

..., 0.001, , 0.01, , 0.1, , 1, ...

Linear Regression with multiple variables

FEATURES AND POLYNOMIAL REGRESSION

Housing prices prediction

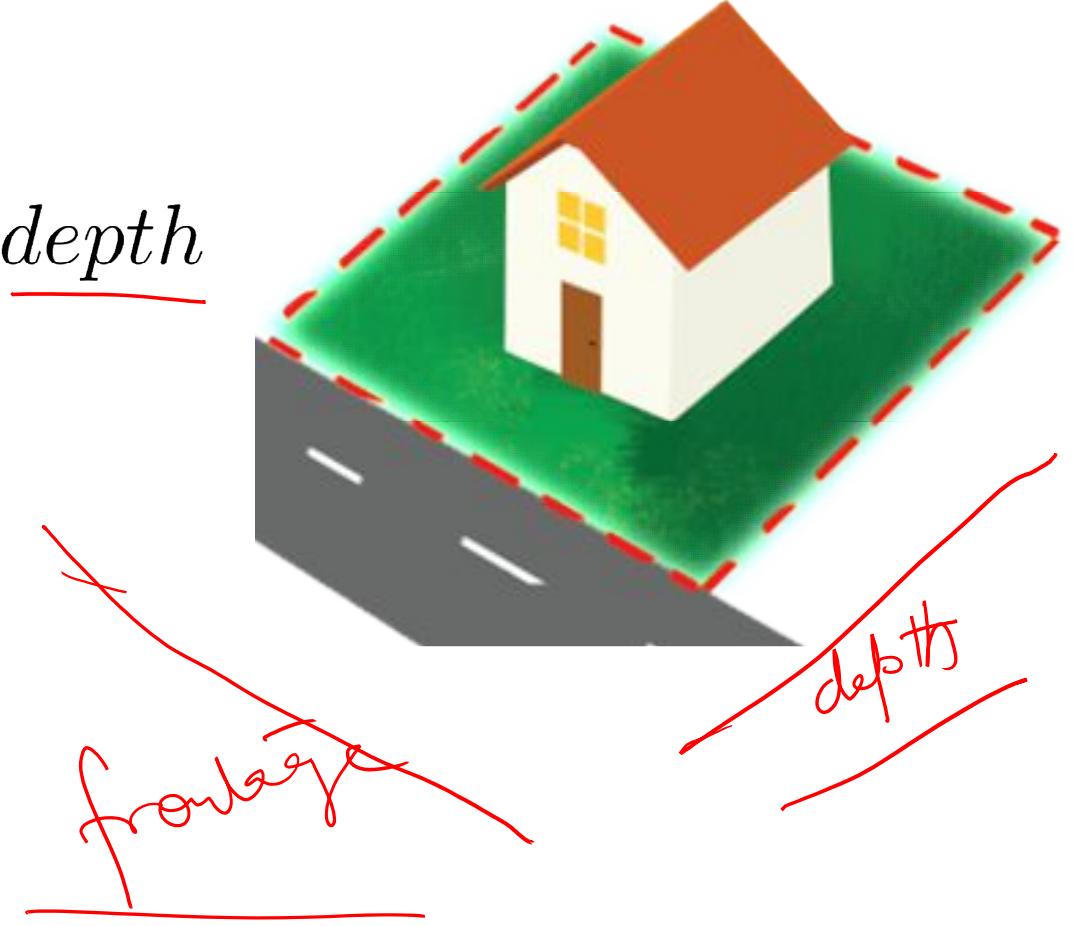
$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underline{\text{frontage}} + \theta_2 \times \underline{\text{depth}}$$

✓ Area = ?

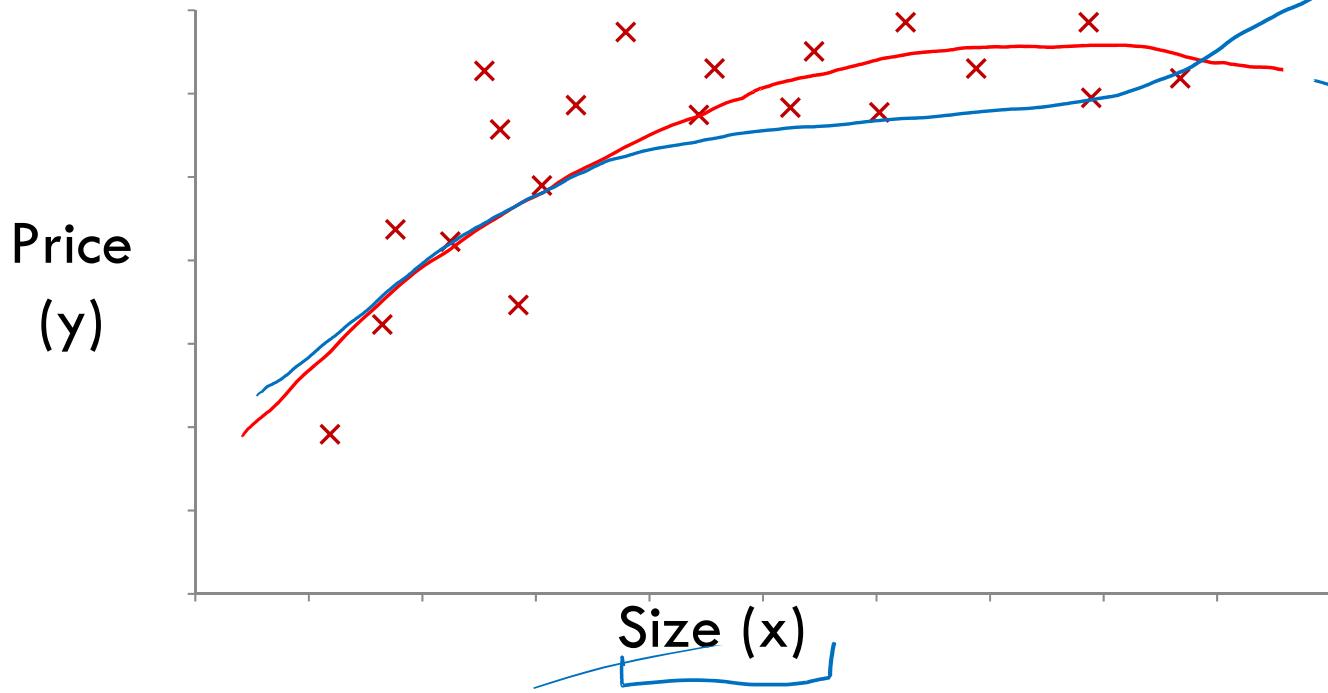
$$X = \text{frontage} \times \text{depth}$$

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

↖ Land Area



Polynomial regression



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

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

$$h_{\theta}(x) = \theta_0 + \theta_1(x) + \theta_2x_2 + \theta_3x_3 + \theta_2x_2 + \theta_3x_3$$

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

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

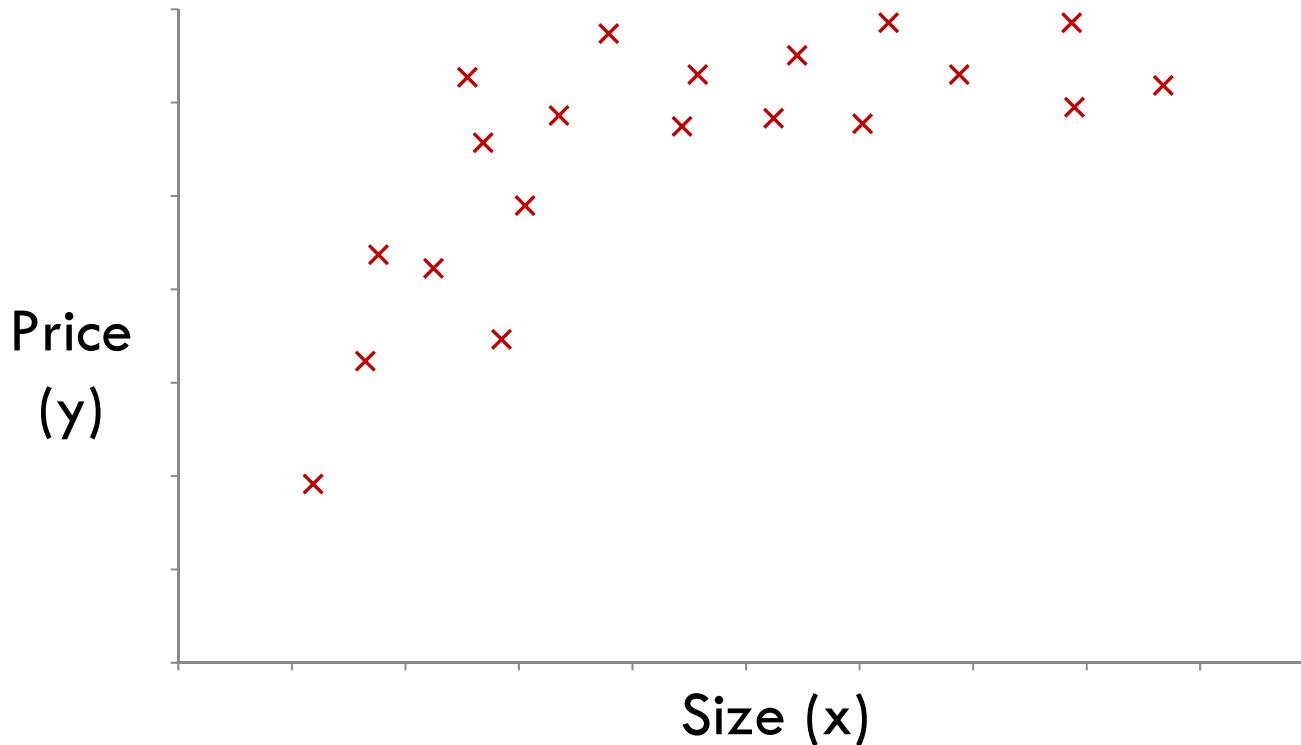
$$x_1 = (\text{size})$$

$$x_2 = (\text{size})^2$$

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

$1 - 1000, 1 - 1000, 000, 1 - 10^9$

Choice of features



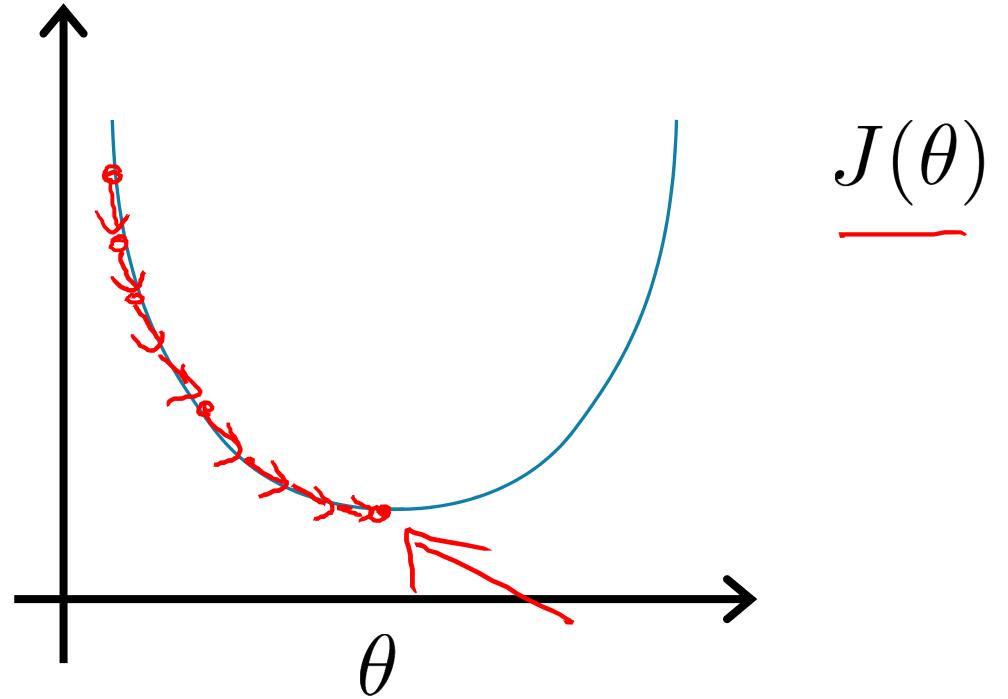
$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2 \sqrt{(\text{size})}$$

Linear Regression with multiple variables

NORMAL EQUATION

Gradient Descent



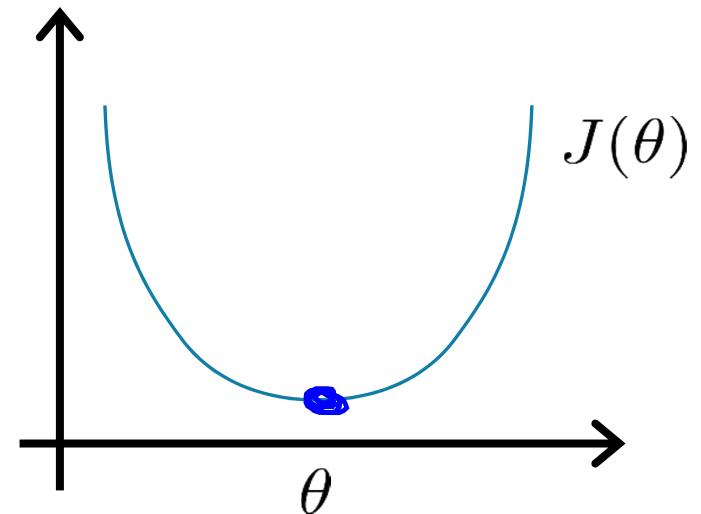
Normal equation: Method to solve for θ analytically.

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

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

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for θ



$$\theta \in \mathbb{R}^{n+1}$$

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

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

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

Examples: $m = 4$.

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

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

$m \times (n+1)$

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

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

m -dimensional vector

Examples: $m = 5$.

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	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
1	3000	4	1	38	540

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \\ 1 & 3000 & 4 & 1 & 38 \end{bmatrix} \quad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \\ 540 \end{bmatrix}$$

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

m training examples, n features.

Gradient Descent

- Need to choose α .
- Needs many iterations.
- Works well even when n is large.

$$\underline{n = 10^6}$$

$\leftarrow -$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$
- Slow if n is very large.

$$\underline{n = 100}$$

$$\underline{n = 1000}$$

$$\dots \underline{n = 10000}$$

$O(n^3)$

THANKS

Keep Learning
Keep Growing



Dr. Neeraj Gupta
Assistant Professor, Dept. of CEA
neeraj.gupta@gla.ac.in