



**BITS Pilani**  
Pilani Campus

# Machine Learning

## AIML CZG565

### M3 : Linear Models for Regression

Course Faculty of MTech Cluster  
BITS – CSIS - WILP

## Disclaimer and Acknowledgement



- The content for these slides has been obtained from books and various other source on the Internet
- We here by acknowledge all the contributors for their material and inputs.
- We have provided source information wherever necessary
- Students are requested to refer to the textbook w.r.t detailed content of the presentation deck shared over canvas
- We have reduced the slides from canvas and modified the content flow to suit the requirements of the course and for ease of class presentation

**Source:** “Probabilistic Machine Learning, An Introduction”, Kevin P. Murphy, Slides of Prof. Sugata, Prof. Chetana from BITS Pilani, Prof. Raja vadhana from BITS Pilani , CS109 and CS229 stanford lecture notes and many others who made their course materials freely available online.

# Course Plan

- |     |  |
|-----|--|
| M1  | Introduction                                 |
| M2  | Machine learning Workflow                    |
| M3  | Linear Models for Regression                 |
| M4  | Linear Models for Classification             |
| M5  | Decision Tree                                |
| M6  | Instance Based Learning                      |
| M7  | Support Vector Machine                       |
| M8  | Bayesian Learning                            |
| M9  | Ensemble Learning                            |
| M10 | Unsupervised Learning                        |
| M11 | Machine Learning Model Evaluation/Comparison |

# Agenda

- Linear Model for Regression
  - Direct solution vs Iterative Method
  - Gradient Descent (**Introduction**)
- 
- Linear Basis Function
  - Notion of Bias vs Variance

# Linear Regression

# Inductive Learning Hypothesis : Interpretation

- Target Concept
- Discrete :  $f(x) \in \{\text{Yes, No, Maybe}\}$  Classification
- Continuous :  $f(x) \in [20-100]$  Regression
- Probability Estimation :  $f(x) \in [0-1]$

Sky	AirTemp	Altitude	Wind	Water	Forecast	Humidity
Sunny	Warm	Normal	Strong	Warm	Same	60
Sunny	Warm	High	Strong	Warm	Same	75
Rainy	Cold	High	Strong	Warm	Change	70
Sunny	Warm	High	Strong	Cool	Change	45

## Supervised Learning: Regression

---

Wish to learn a function  $f : X \rightarrow Y$ , where predicted output  $Y$  is real, given the  $n$  real training instances  $\{<x^1, y^1> \dots <x^n, y^n>\}$ .

Examples include

- predict weight from gender, height, age, ...
- Predict house price from locality, area, income, ...
- predict Google stock price today from Google, Yahoo, MSFT prices yesterday
- predict each pixel intensity in robot's current camera image, from previous image and previous action

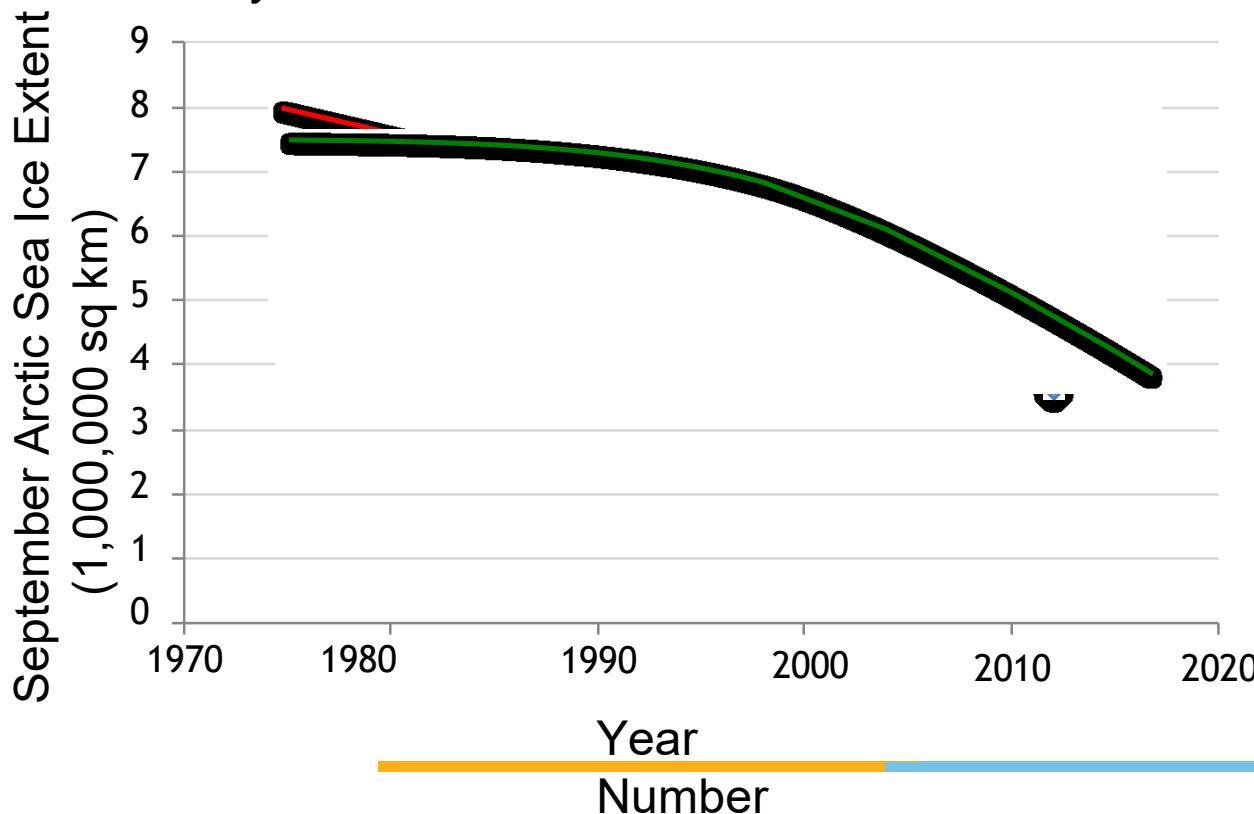
# Formal Representation: Interpretation



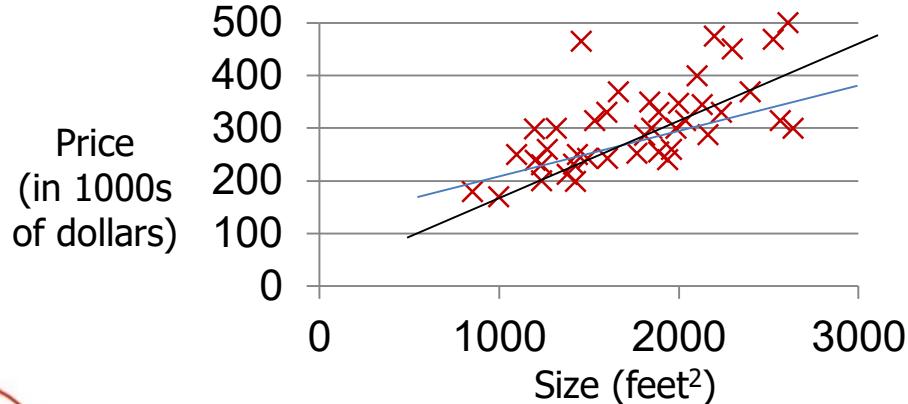
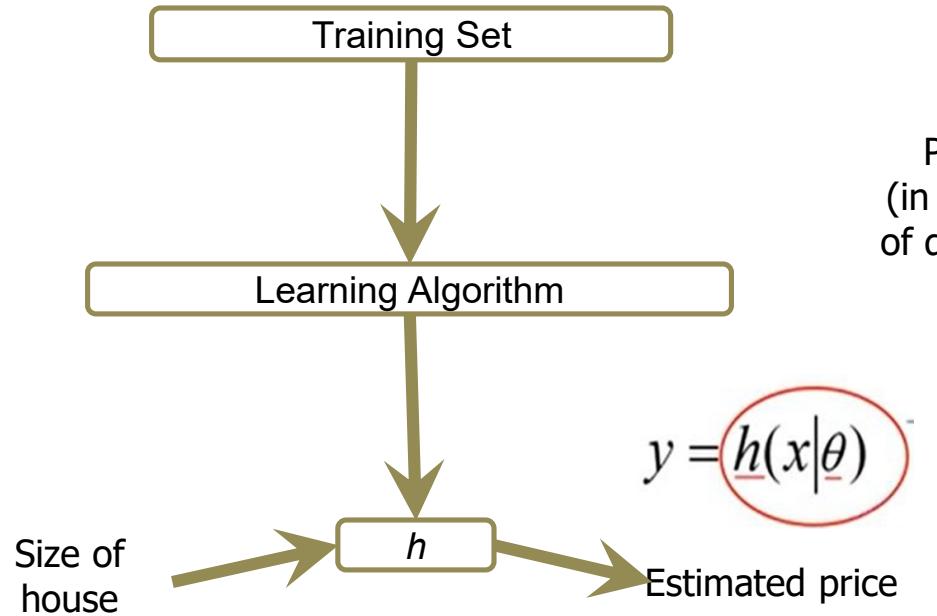
## Supervised Learning: Regression

GOAL : Previously unseen records should be assigned a value as accurately as possible.

- Given  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function  $f(x)$  to predict  $y$  given  $x$ 
  - $y$  is real-valued



# Machine Learning Process : Interpretation



Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

Hypothesis:

$\theta_i$ 's : Parameters  
 $h(\cdot)$  : Model

**How to choose  $\theta_i$ 's ?**

# Types of Regression Models

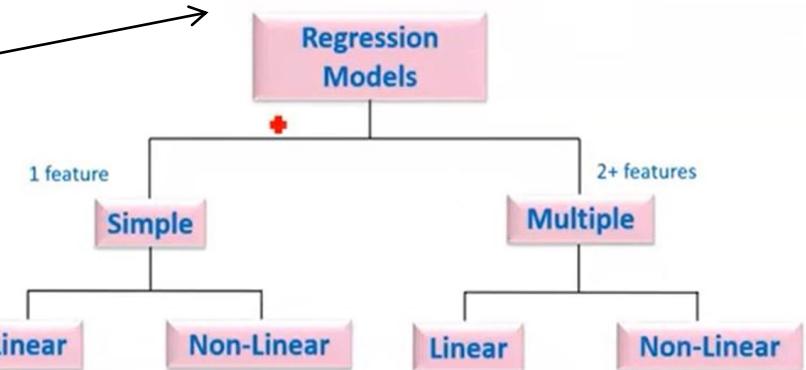
## simple regression model

(Education)  $x \longrightarrow y$  (Income)

## multiple regression model

(Education)  $x_1$   $\longrightarrow$    
 (Soft Skills)  $x_2$   $\longrightarrow$    
 (Experience)  $x_3$   $\longrightarrow$    
 (Age)  $x_4$   $\longrightarrow$   $y$  (Income)

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$



- Hypothesis:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d = \sum_{j=0}^d \theta_j x_j$$

Assume  $x_0 = 1$

---

# Learning the Model Parameters

Closed Form Solution Approach

Gradient Descent Approach

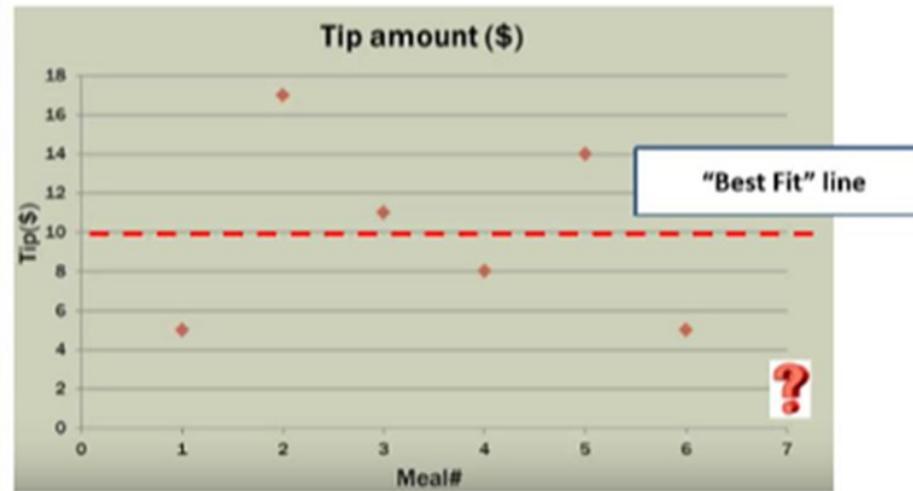
**How to choose  $\theta_i$ 's ?**

---

# Notion of Cost Function

Meal #	Tip amount (\$)
1	5.00
2	17.00
3	11.00
4	8.00
5	14.00
6	5.00

mean = \$10

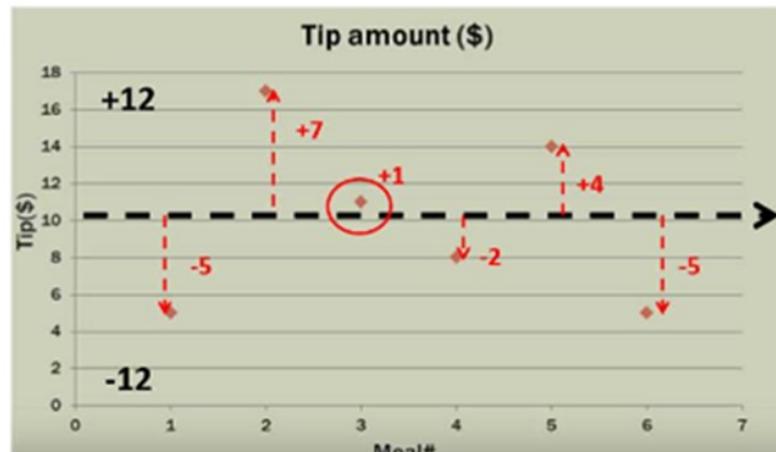


$$y = \theta_0$$

# Notion of Cost Function

Meal #	Tip amount (\$)
1	5.00
2	17.00
3	11.00
4	8.00
5	14.00
6	5.00

mean = \$10



RESIDUALS (error): Distance from best fit line and the actual values

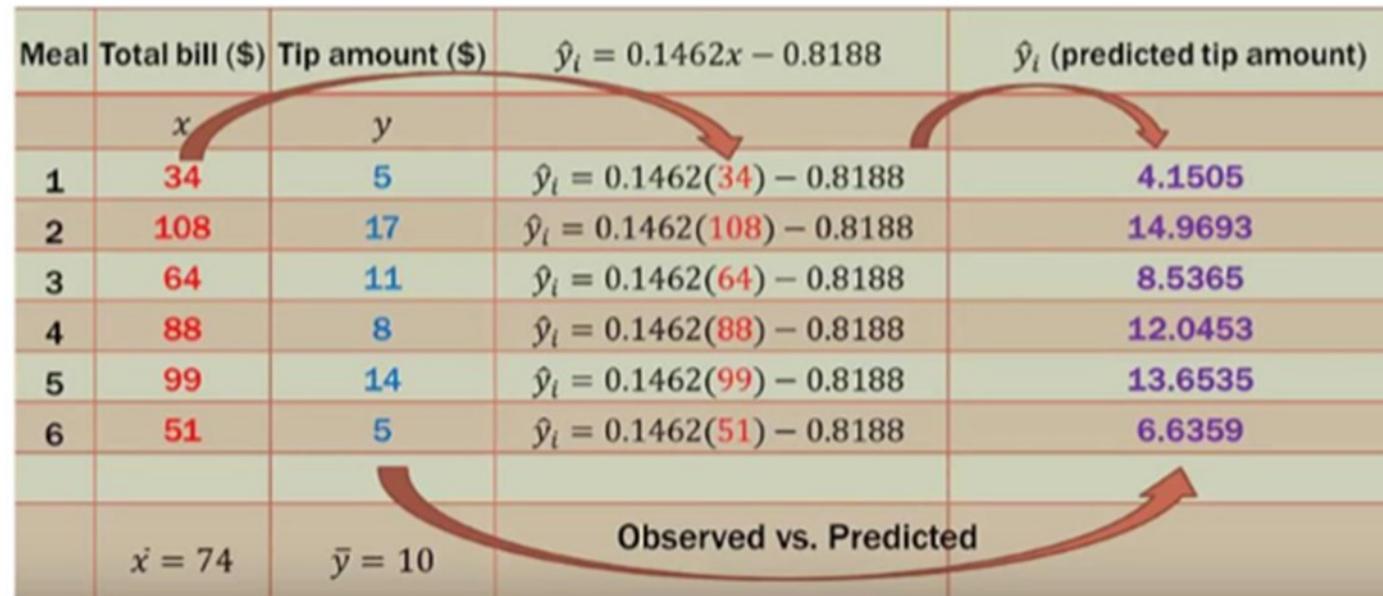
Meal#	Residual	Residual <sup>2</sup>
1	-5	25
2	+7	49
3	+1	1
4	-2	4
5	+4	16
6	-5	25

Sum of Squared Errors (SSE) = 120

# Notion of Cost Function

Meal #	Tip amount (\$)
1	5.00
2	17.00
3	11.00
4	8.00
5	14.00
6	5.00

mean = \$10



$$y = \theta_0 + \theta_1 x$$

$$J(\Theta) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

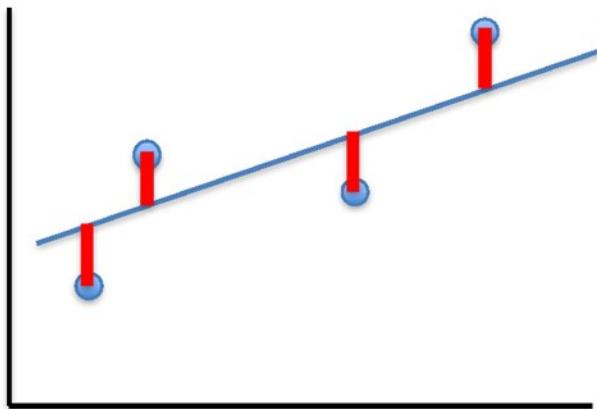
$$J(\Theta) = \frac{1}{2n} \sum_{i=1}^n (h_\Theta(x^{(i)}) - y^{(i)})^2$$

# Linear Regression : Least Squares

- Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

- Fit by solving  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$



- Hypothesis

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

- Parameters

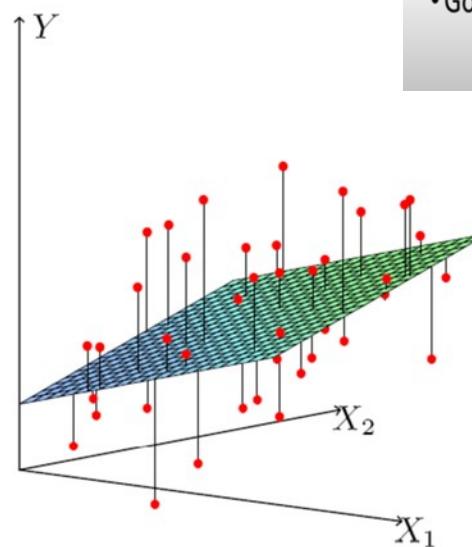
$$\theta_0 \text{ and } \theta_1$$

- Cost Function

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

- Goal

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



# Intuition Behind Cost Function

---

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [\theta_0, \theta_1]$

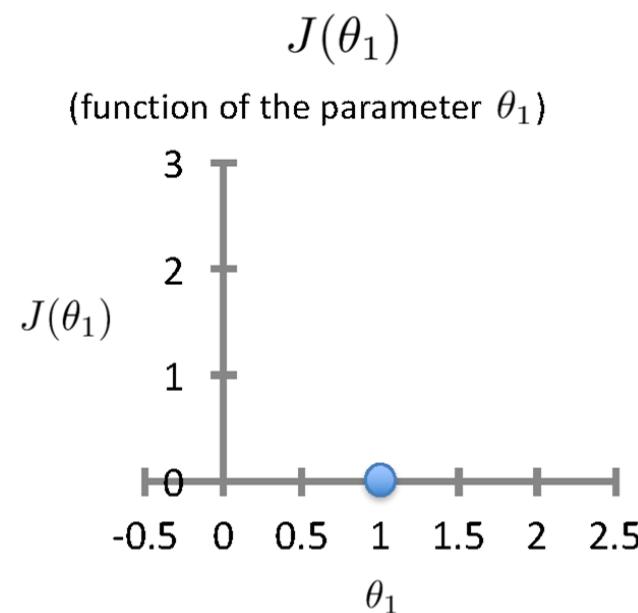
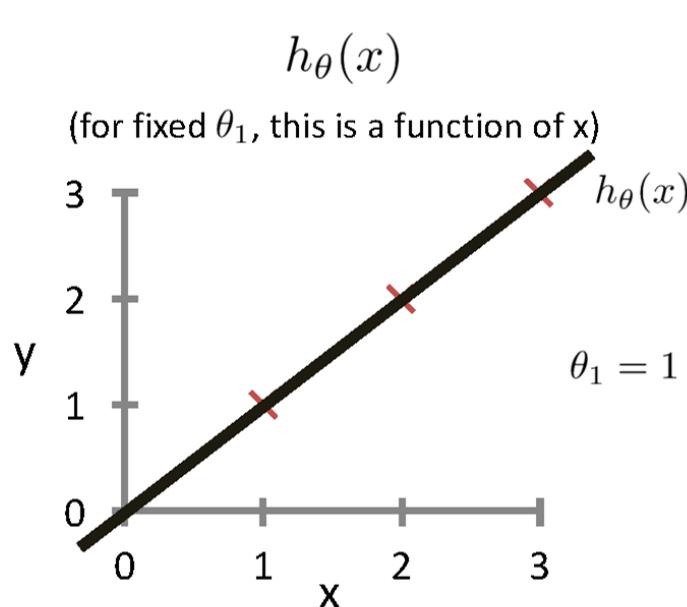
Cost function for linear regression:

Source Credit : Based on example  
by Andrew Ng

# Intuition Behind Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [\theta_0, \theta_1]$

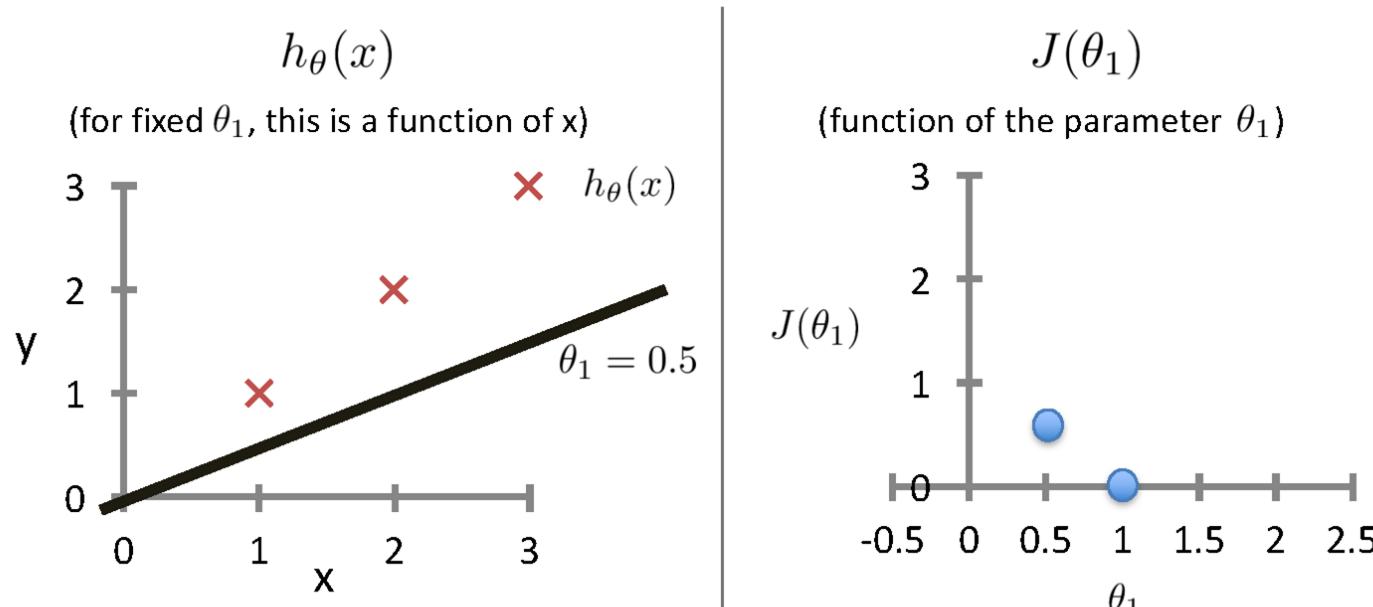


Source Credit : Based on example by Andrew Ng

# Intuition Behind Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [\theta_0, \theta_1]$



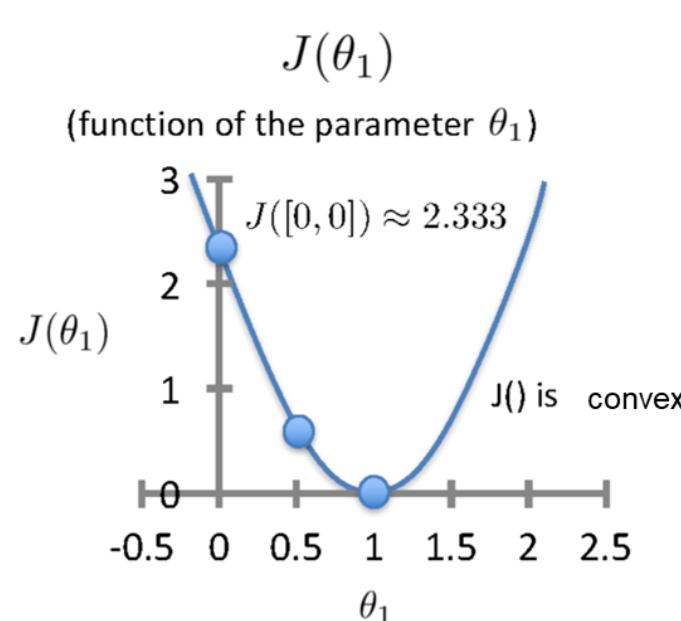
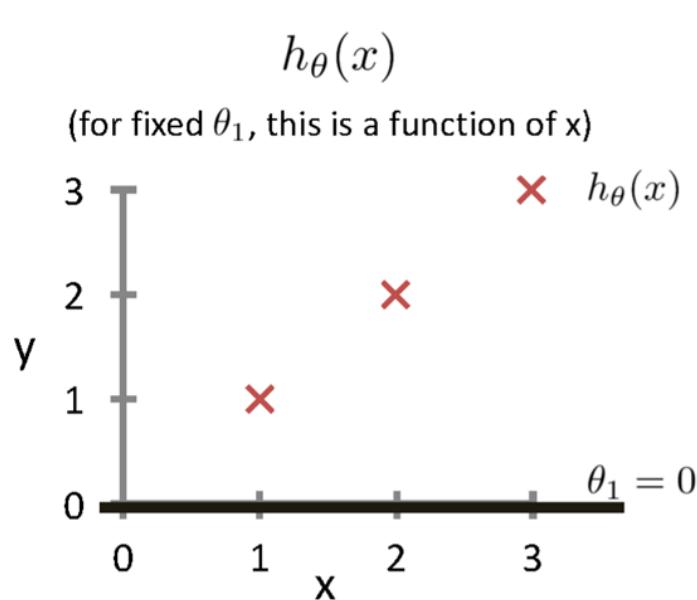
Based on example  
by Andrew Ng

$$J([0, 0.5]) = \frac{1}{2 \times 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 0.58$$

# Intuition Behind Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [\theta_0, \theta_1]$

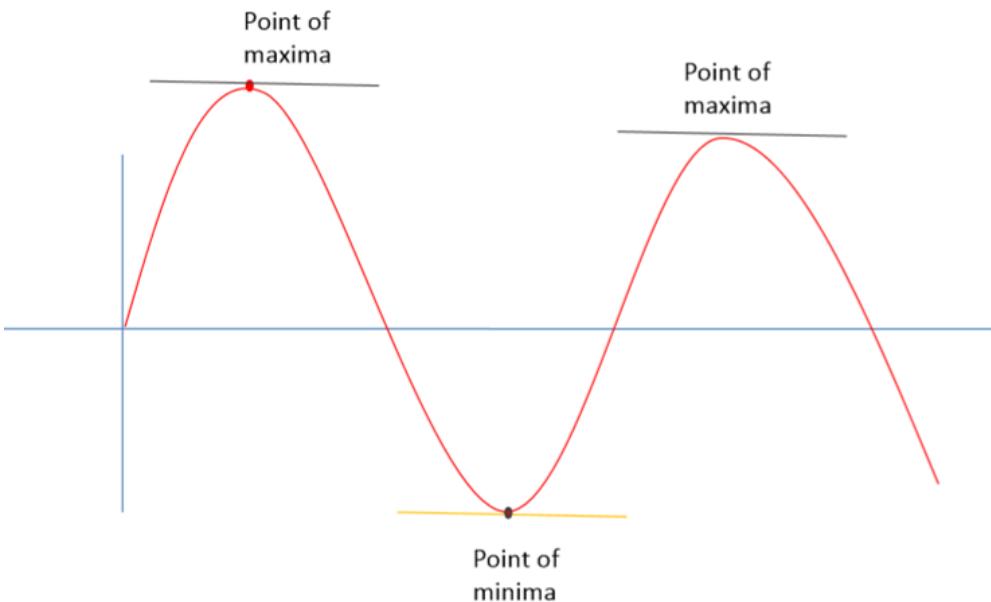


Source Credit : Based on example by Andrew Ng

# Notion of Maxima and Minima of a function

## In General

- For maxima and minima  $m = dy/dx = \tan \theta = 0$
- $dy/dx = 0$  means tangent is parallel to X – axis.



# Notion of Maxima and Minima of a function

## In Machine Learning

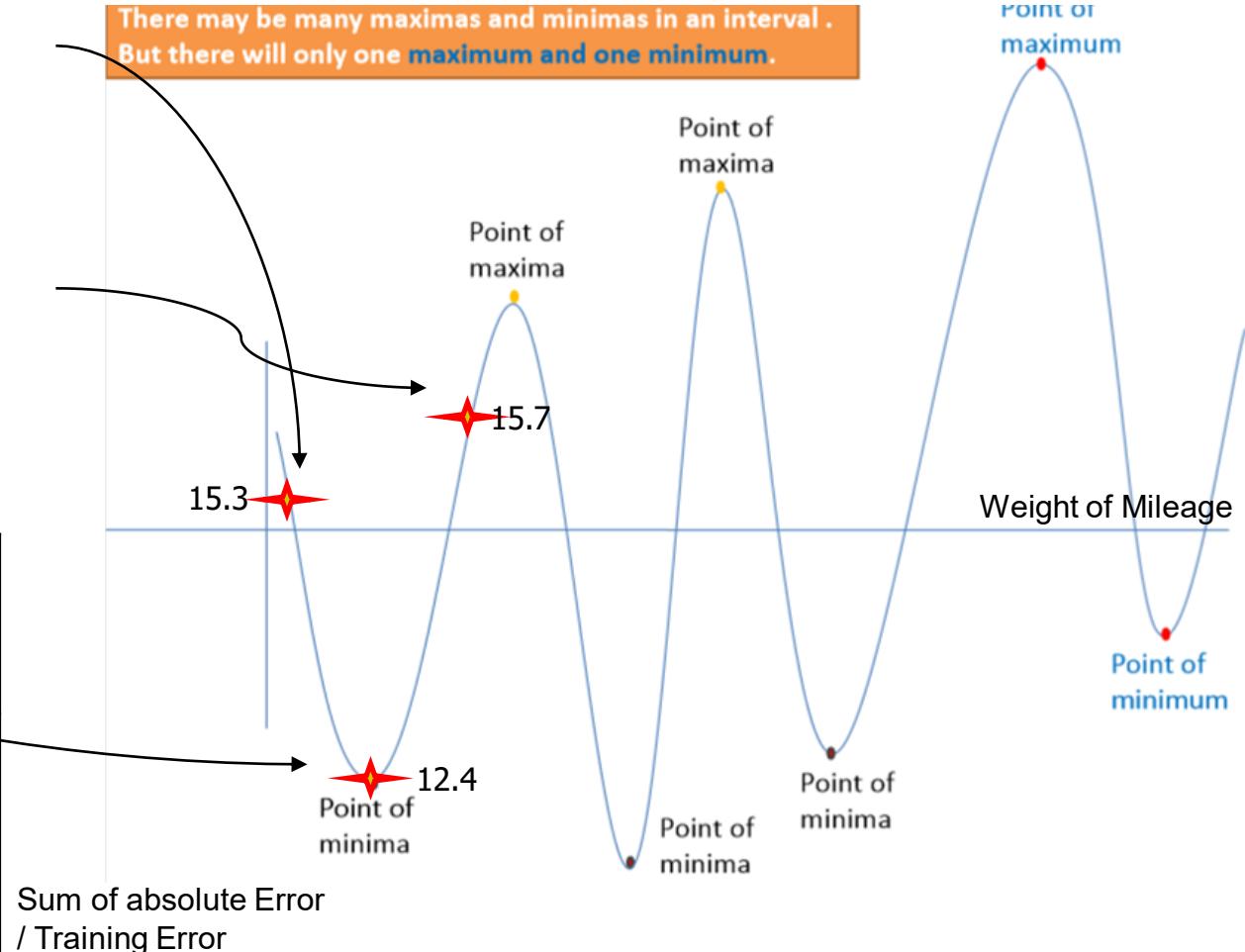
$$\text{CarPrice} = 1 + 0.05\text{Mileage}$$

$$\text{CarPrice} = 1 + 0.25\text{Mileage}$$

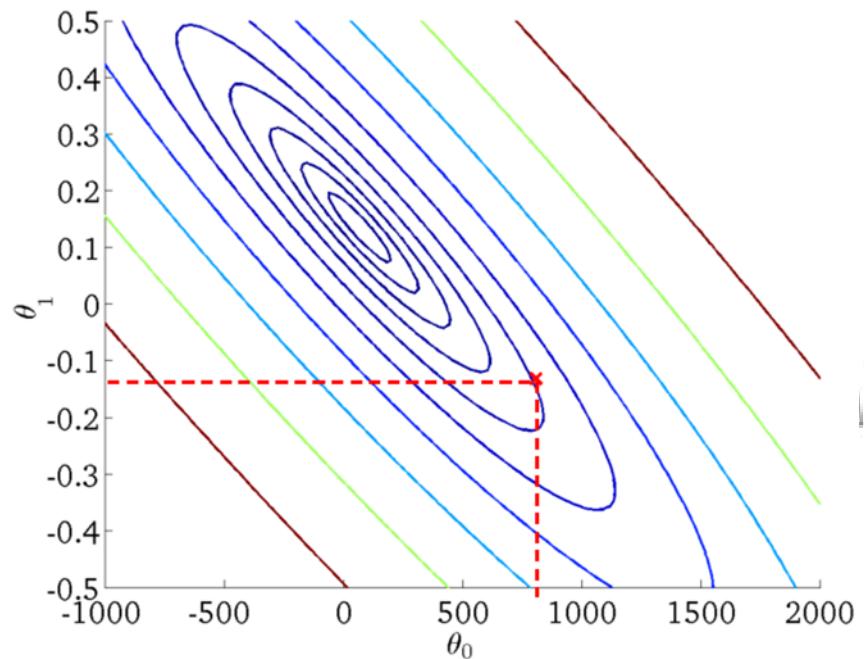
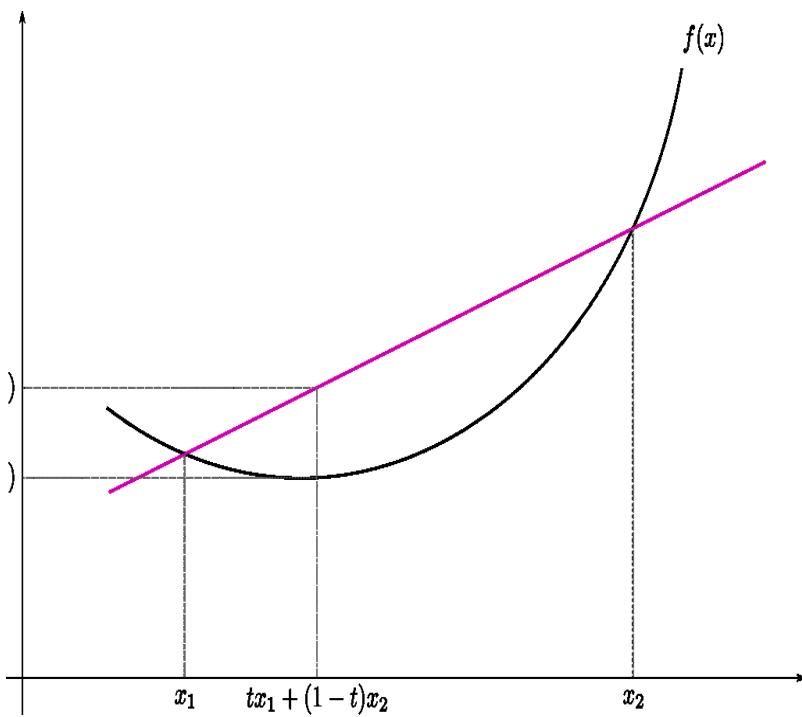
$$\text{CarPrice} = 1 + 0.75\text{Mileage}$$

There may be many maximas and minimas in an interval.  
But there will only one maximum and one minimum.

Mileage (in kmpl)	Car Price (in cr)
9.8	10.48
9.12	1.75
9.5	6.95
10	2.51



# Convex Function : Multivariate View



Real-valued function defined on an  $n$ -dimensional interval is called **convex** if the line segment between any two points on the graph of the function lies above or on the graph

# Closed Form Solution

# Vectorization

- Benefits of vectorization

- More compact equations
- Faster code (using optimized matrix libraries)

- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = [ 1 \quad x_1 \quad \dots \quad x_d ]$$

- Can write the model in vectorized form as  $h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad \mathbf{x}^\top = [ 1 \quad \text{No.of.Years Experience} ]$$

X No.of.Years of Experience (in Years)	Y Salary Of the Employee (in Lakhs)
1	2
2	3
3	4
4	5
5	6

# Vectorization

- Consider our model for  $n$  instances:

$$h(\mathbf{x}^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$\mathbb{R}^{(d+1) \times 1}$

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}$$

$\mathbb{R}^{n \times (d+1)}$

Let:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- Can write the model in vectorized form as  $h_{\boldsymbol{\theta}}$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

# Vectorization

- For the linear regression cost function:

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

$$= \frac{1}{2n} \sum_{i=1}^n \left( \theta^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y})$$

Dimensions:   
 $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$  (blue arrow)   
 $\mathbf{y} \in \mathbb{R}^{(d+1) \times 1}$  (blue arrow)   
 $\mathbf{X} \in \mathbb{R}^{1 \times n}$  (blue arrow)   
 $\mathbf{y} \in \mathbb{R}^{n \times 1}$  (blue arrow)

Let:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

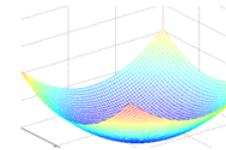
$$\theta^T = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

X	Y
No.of.Years of Experience (in Years)	Salary Of the Employee (in Lakhs)
1	2
2	3
3	4
4	5
5	6

# Vectorization

- Instead of using GD, solve for optimal  $\theta$  analytically

– Notice that the solution is when  $\frac{\partial}{\partial \theta} J(\theta) = 0$



- Derivation:

$$\begin{aligned} J(\theta) &= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - \boxed{\mathbf{y}^\top \mathbf{X} \theta} - \boxed{\theta^\top \mathbf{X}^\top \mathbf{y}} + \mathbf{y}^\top \mathbf{y} \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \end{aligned}$$

Take derivative and set equal to 0, then solve for  $\theta$ :

$$\frac{\partial}{\partial \theta} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \cancel{\mathbf{y}^\top \mathbf{y}}) = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y} = 0$$

$$(\mathbf{X}^\top \mathbf{X})\theta = \mathbf{X}^\top \mathbf{y}$$

Closed Form Solution:

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

# Vectorization

- Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- If  $X^T X$  is not invertible (i.e., singular), may need to:
  - Use pseudo-inverse instead of the inverse
    - In python, `numpy.linalg.pinv(a)`
  - Remove redundant (not linearly independent) features
  - Remove extra features to ensure that  $d \leq n$

# Fit the Linear Regression Model :

## Using Closed Form – Problem Type 1

X No.of.Years of Experience (in Years)	Y Salary Of the Employee (in Lakhs)
1	2
2	3
3	4
4	5
5	6

$$\left( X^T \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix} * X \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \end{bmatrix} \right) - 1 * X^T \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix} * y \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

$$\left( \begin{bmatrix} 5 & 15 \\ 15 & 55 \end{bmatrix} \right) - 1 * X^T \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix} * y \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

$$\left( \begin{bmatrix} 1.1 & -0.3 \\ -0.3 & 0.1 \end{bmatrix} \right) * X^T \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix} * y \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

$$\begin{bmatrix} 0.8 & 0.5 & 0.2 & -0.1 & -0.4 \\ -0.2 & -0.1 & 0 & 0.1 & -0.2 \end{bmatrix} * y \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

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

Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

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

$$\theta^T = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} y = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

# Gradient Descent Approach

# Notion of Gradient Descent

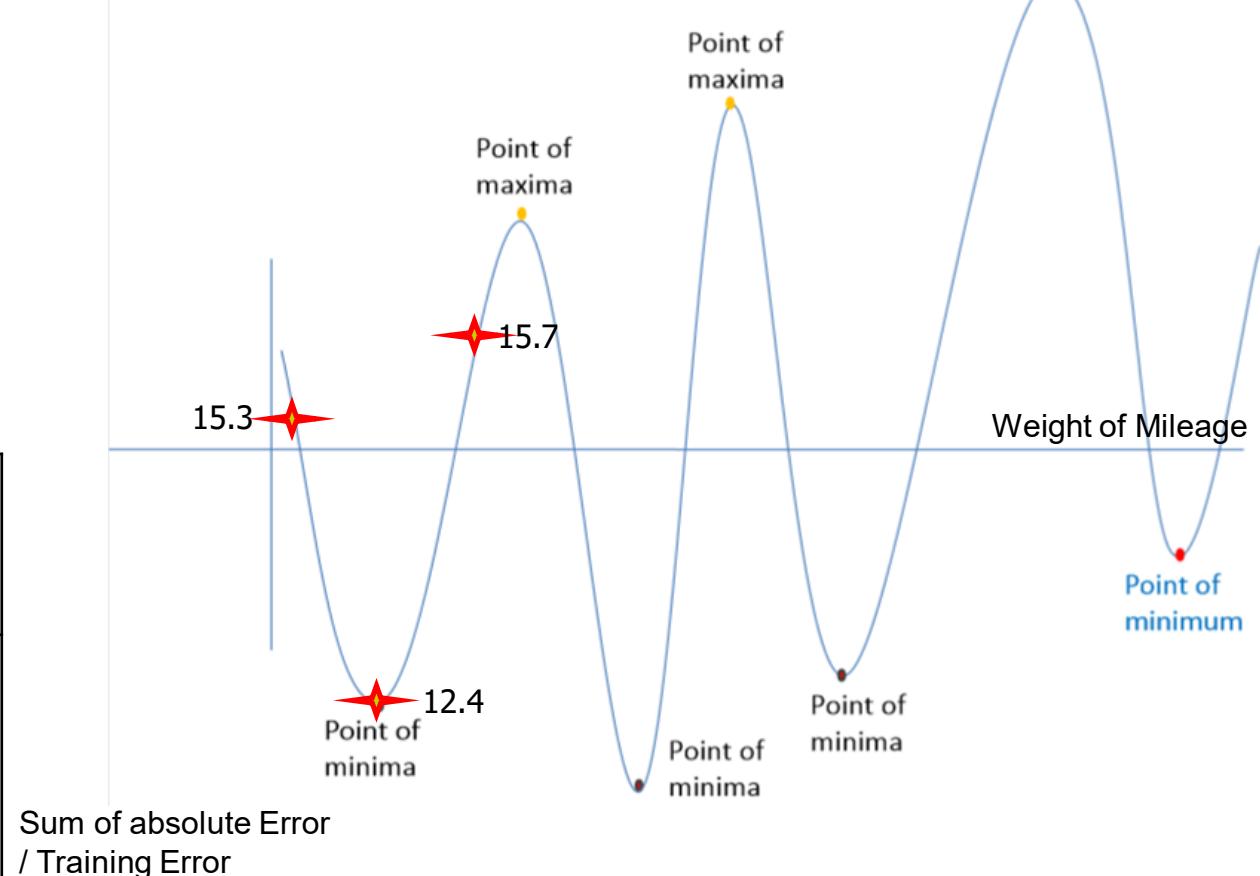
CarPrice = 1 + 0.05Mileage

CarPrice = 1 + 0.25Mileage

CarPrice = 1 + 0.75Mileage

Mileage (in kmpl)	Car Price (in cr)
9.8	10.48
9.12	1.75
9.5	6.95
10	2.51

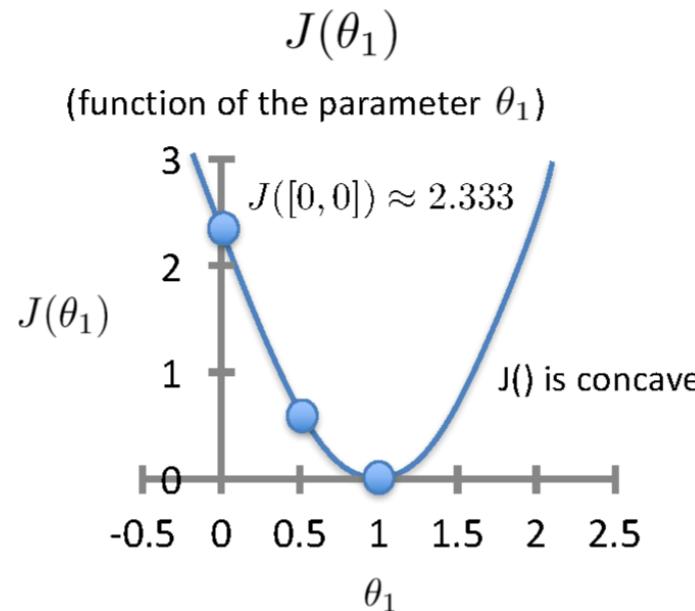
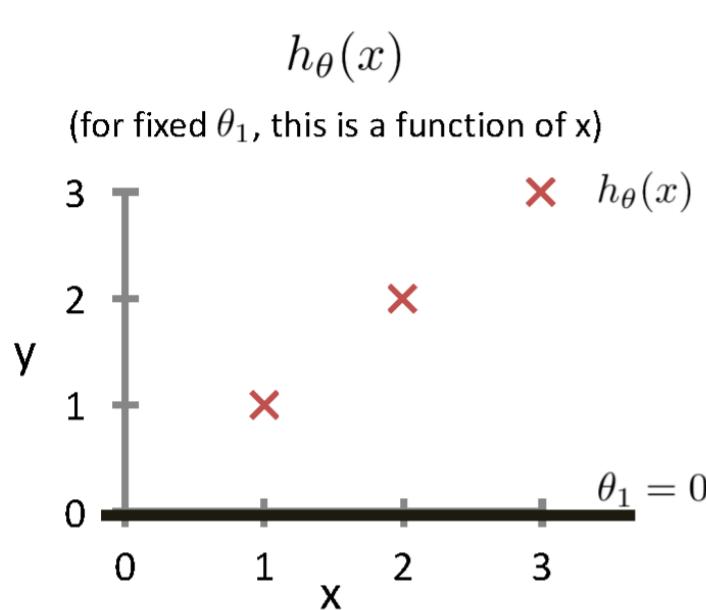
There may be many maximas and minimas in an interval.  
But there will only one maximum and one minimum.



# Intuition Behind Cost Function

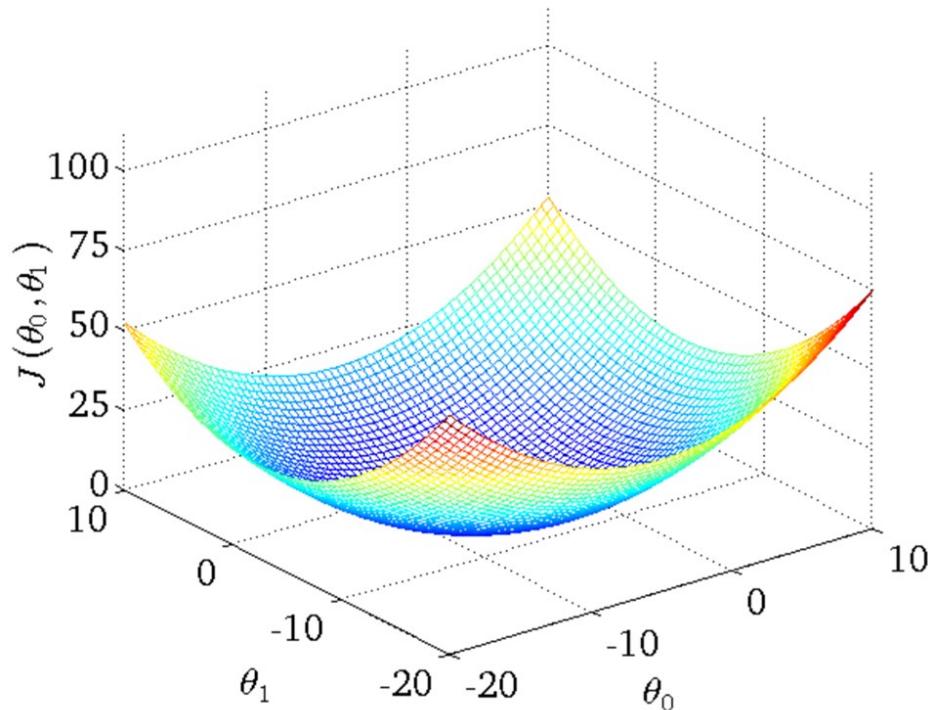
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

For insight on  $J()$ , let's assume  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [\theta_0, \theta_1]$



Source Credit : Based on example by Andrew Ng

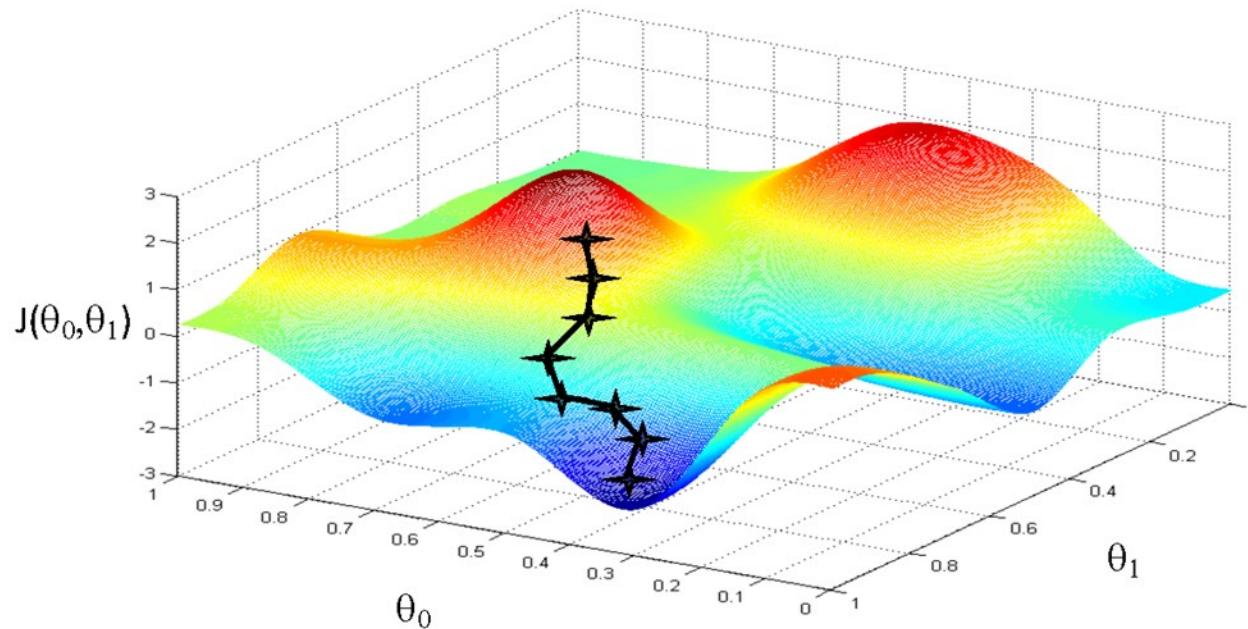
# Intuition Behind Cost Function



Source Credit : Based on example  
by Andrew Ng

# Basic Search Procedure

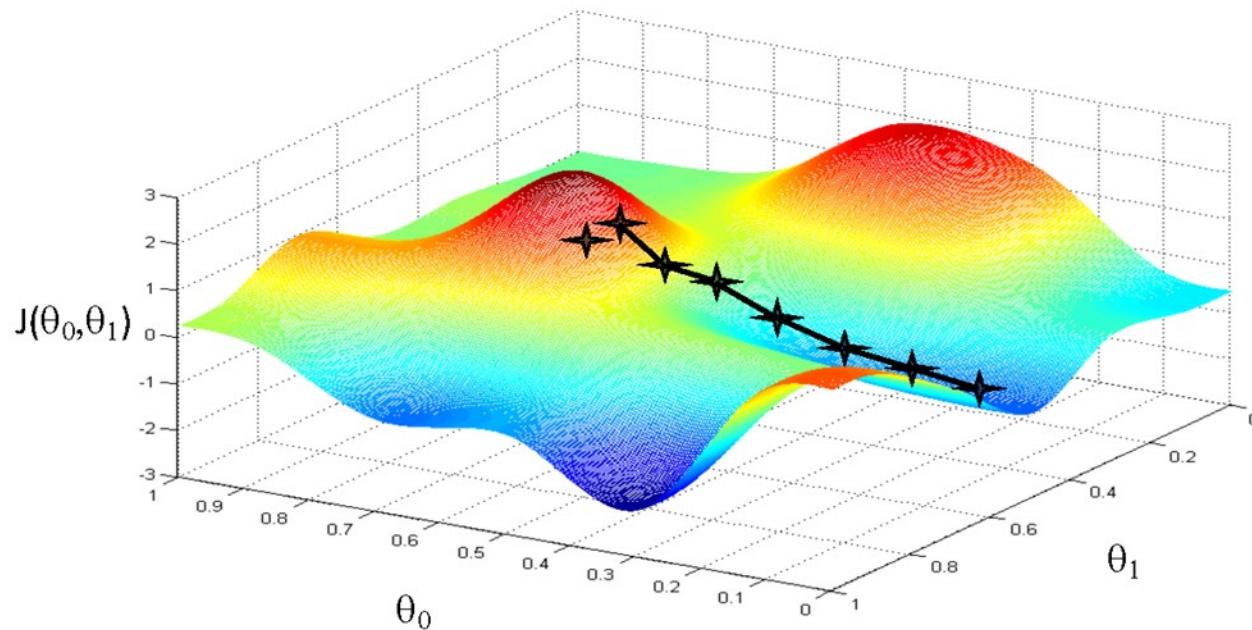
- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$



Source Credit : Figure by Andrew Ng

# Basic Search Procedure

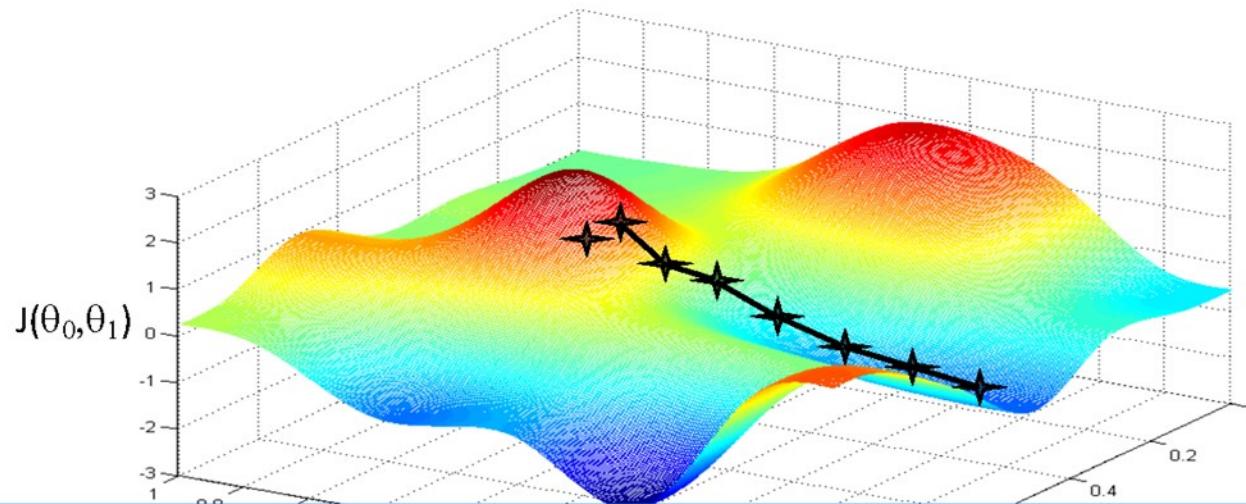
- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$



Source Credit : Figure by Andrew Ng

# Basic Search Procedure

- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$



Since the least squares objective function is convex (concave),  
we don't need to worry about local minima

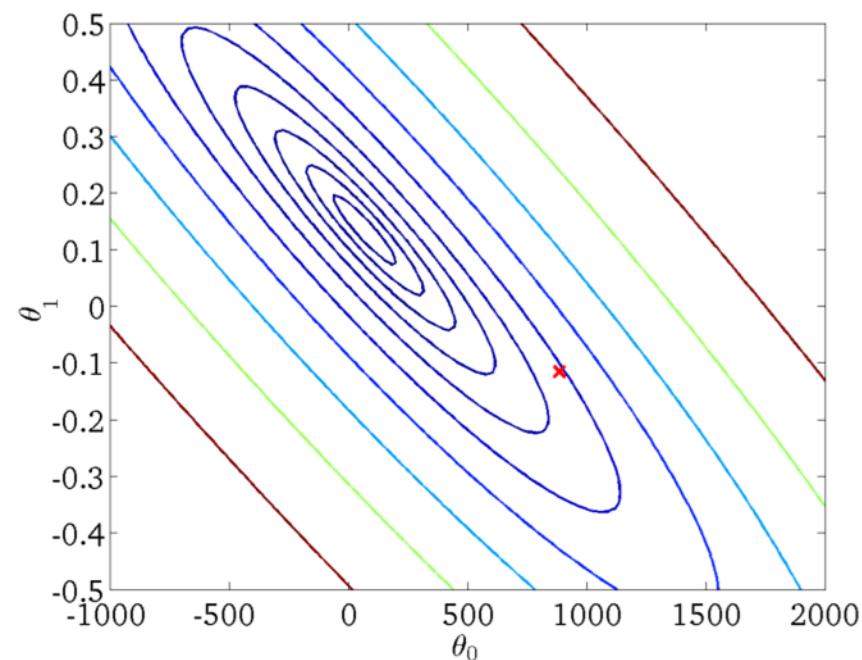
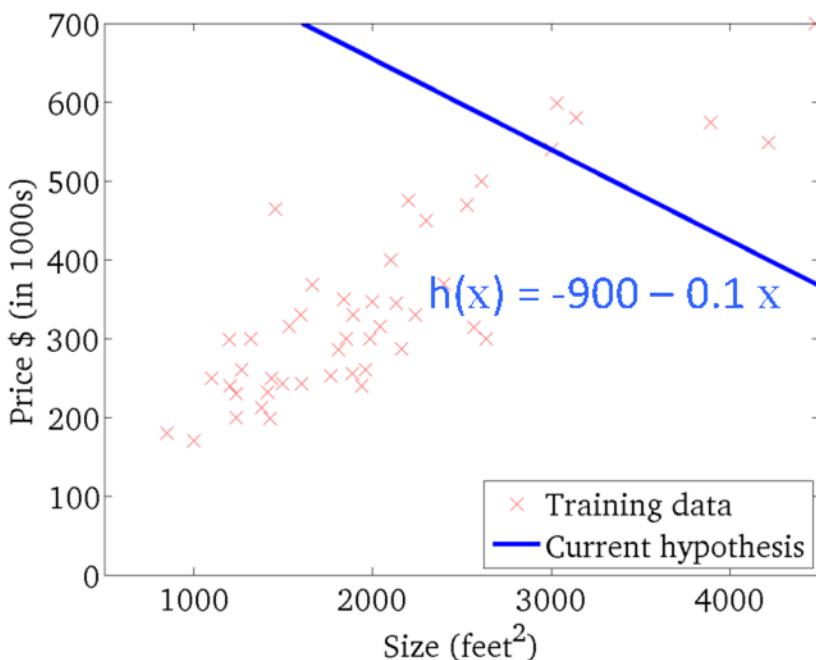
Source Credit : Figure by Andrew Ng

# Gradient Descent

$h_{\theta}(x)$   
 (for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )  
 ...

Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

$J(\theta_0, \theta_1)$   
 (function of the parameters  $\theta_0, \theta_1$ )

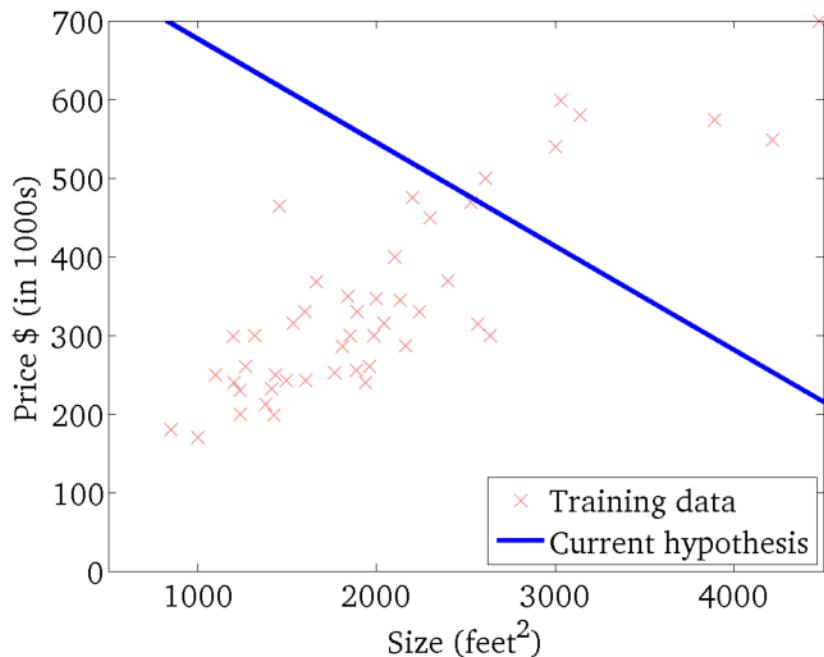


Source Credit : Slide by Andrew Ng

# Gradient Descent

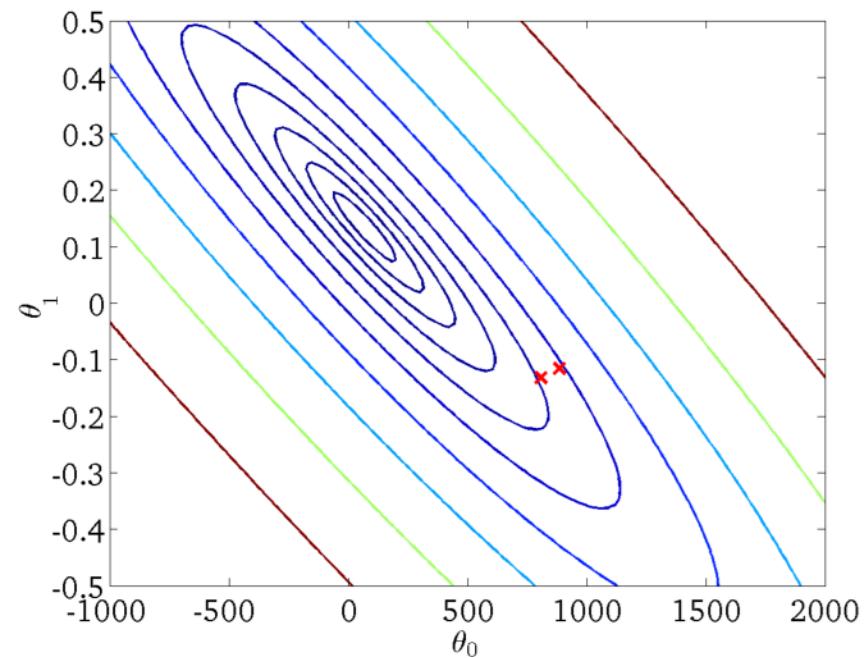
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )

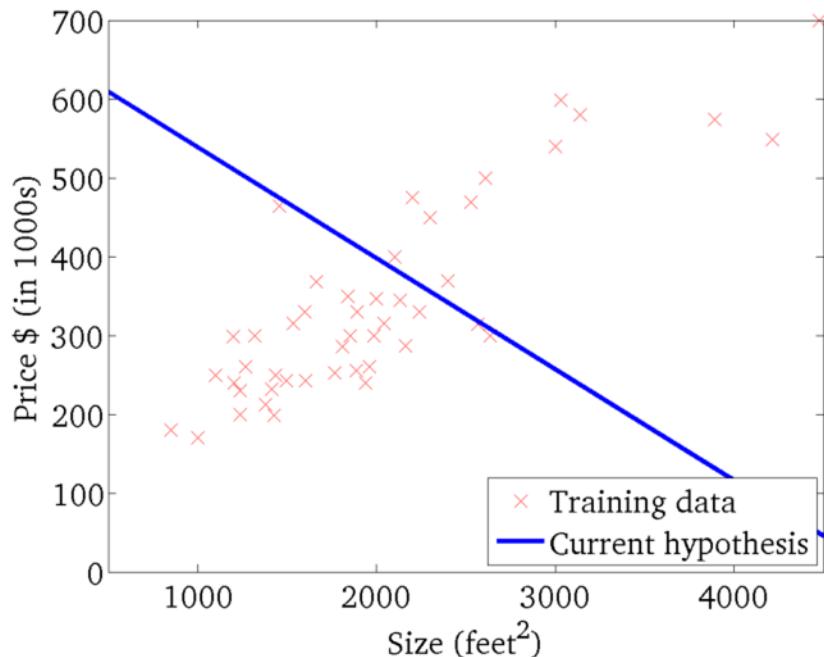


Source Credit : Slide by Andrew Ng

# Gradient Descent

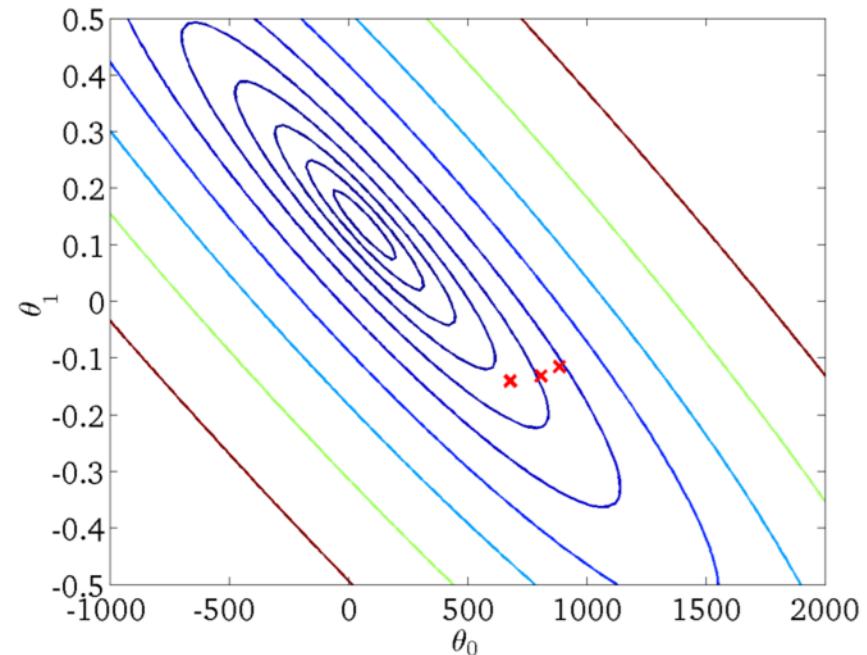
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



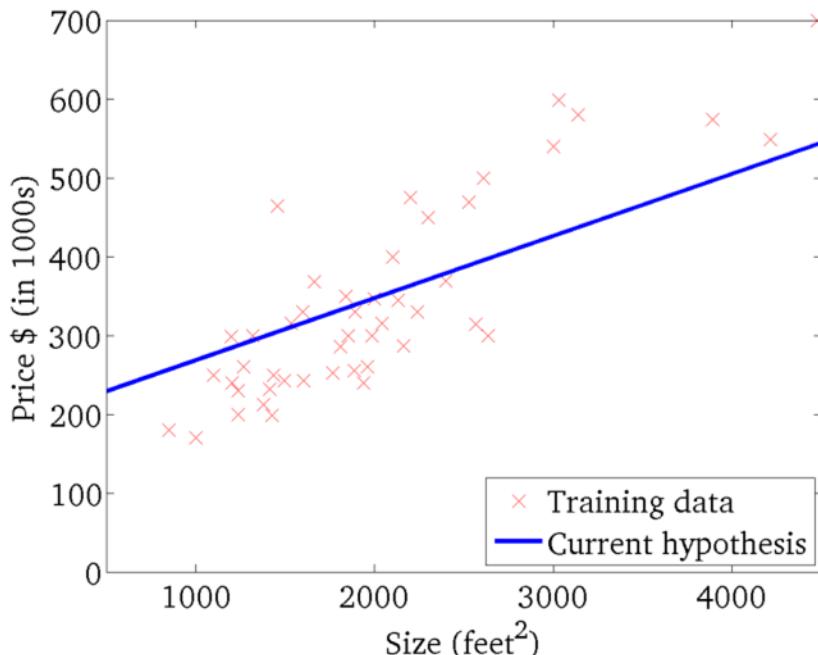
Source Credit : Slide by Andrew Ng

# Gradient Descent

(....After Few Iteration or SEARCH!!)

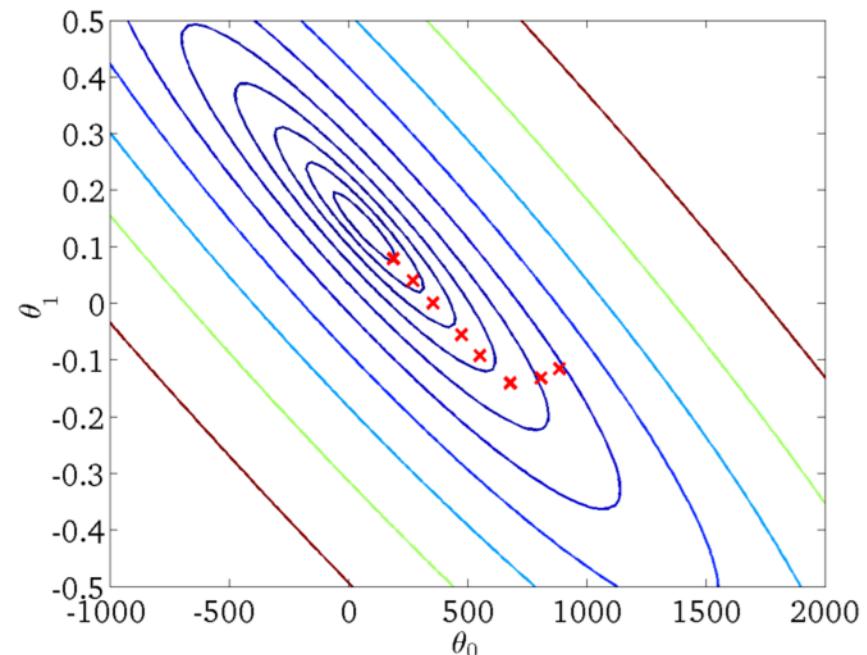
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )

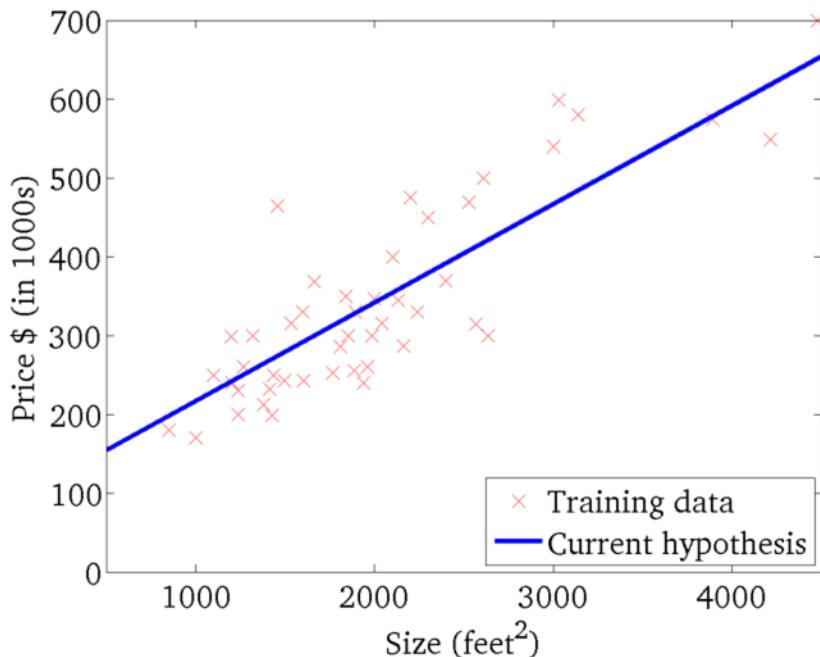


Source Credit : Slide by Andrew Ng

# Gradient Descent

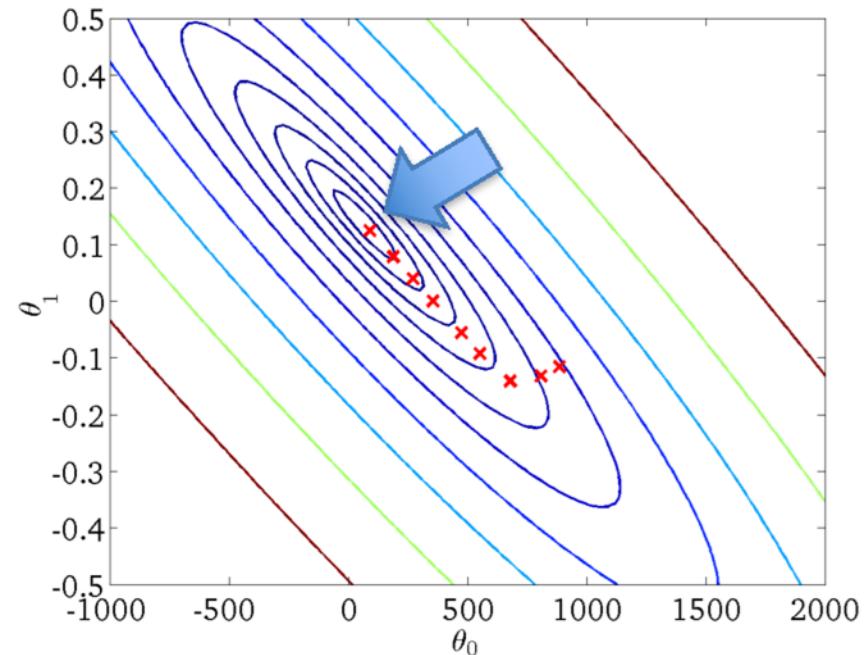
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



Source Credit : Slide by Andrew Ng

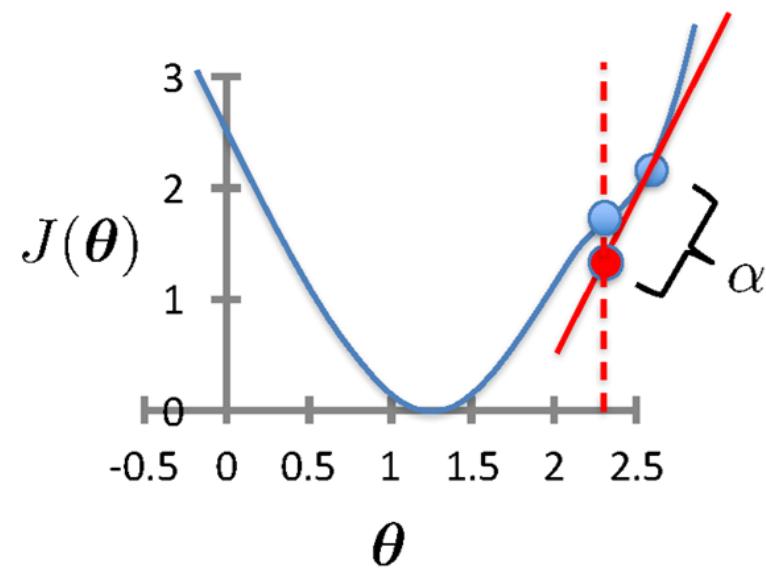
# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

learning rate (small)  
e.g.,  $\alpha = 0.05$



# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

For Linear Regression:

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\
 &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \\
 &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}
 \end{aligned}$$

# Gradient Descent for Linear Regression

- Initialize  $\theta$
- Repeat until convergence

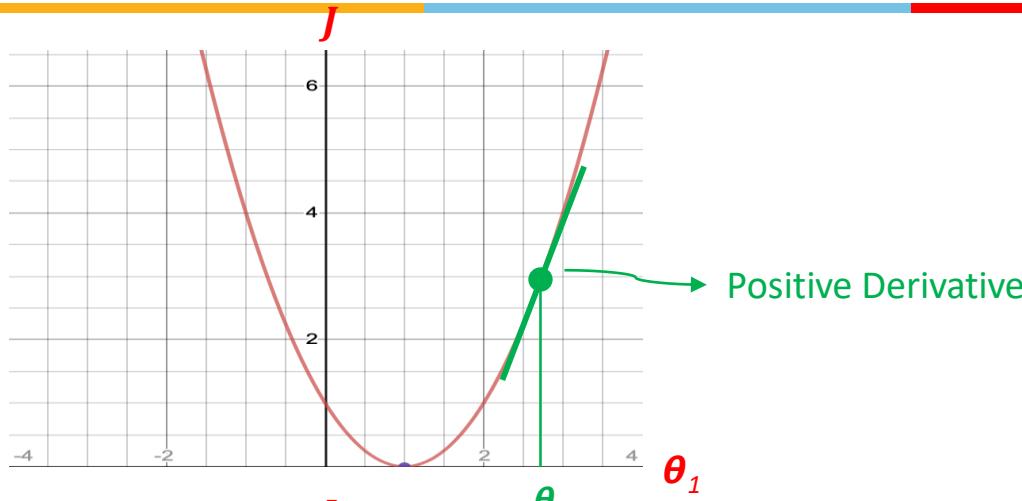
$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\theta} \left( \mathbf{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)}$$

simultaneous update for  $j = 0 \dots d$

- To achieve simultaneous update
  - At the start of each GD iteration, compute  $h_{\theta} \left( \mathbf{x}^{(i)} \right)$
  - Use this stored value in the update step loop
- Assume convergence when  $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

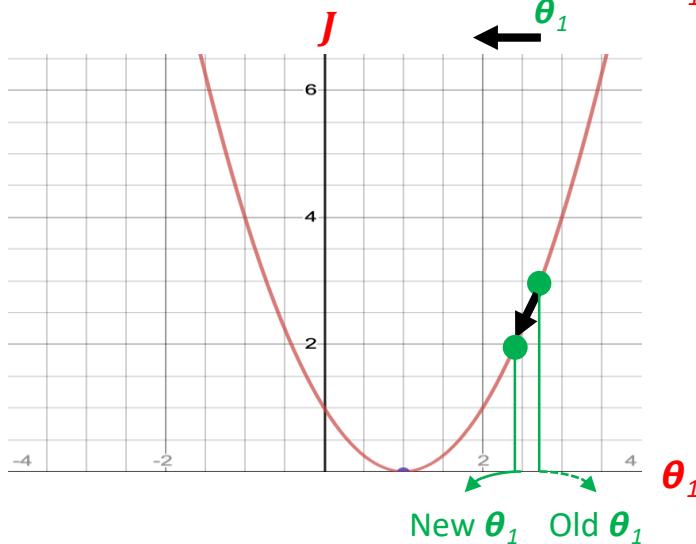
L<sub>2</sub> norm:  $\|\mathbf{v}\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$

# Guarantee of Convergence



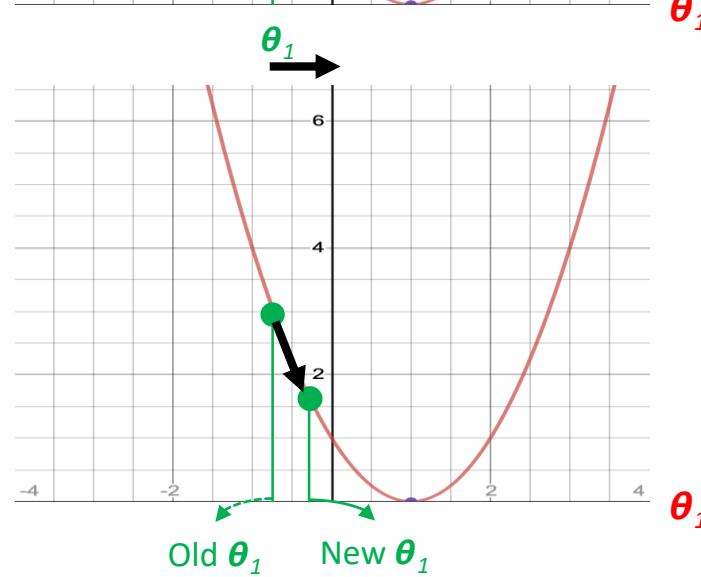
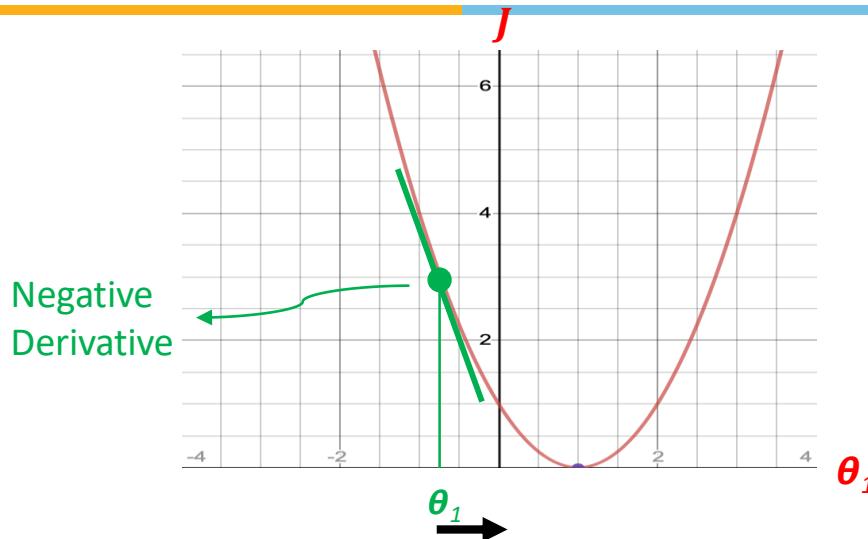
$$\begin{aligned}\theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\ &= \theta_1 - \alpha (\text{Positive Number})\end{aligned}$$

Decrease  $\theta_1$  by a certain value



Source Credit: Prof. Mohammad Hammoud

# Guarantee of Convergence



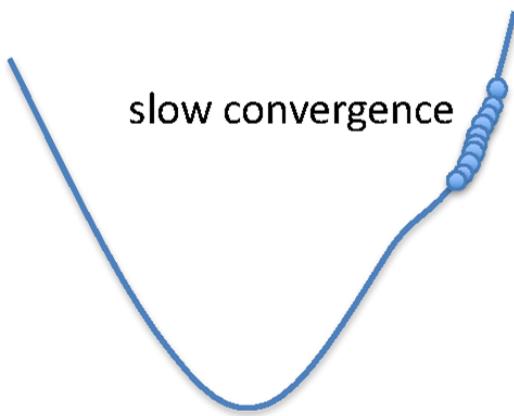
$$\begin{aligned}\theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\ &= \theta_1 - \alpha (\text{Negative Number})\end{aligned}$$

Increase  $\theta_1$  by a certain value

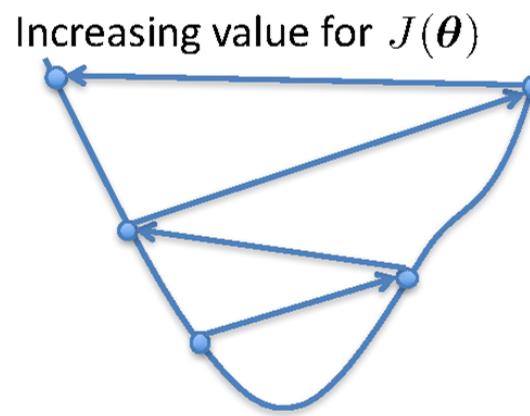
Source Credit: Prof. Mohammad Hammoud

# Choosing Learning Rate

$\alpha$  too small



$\alpha$  too large



- May overshoot the minimum
- May fail to converge
- May even diverge

To see if gradient descent is working, print out  $J(\theta)$  each iteration

- The value should decrease at each iteration
- If it doesn't, adjust  $\alpha$

# Gradient Descent algorithm : Effect of Feature Scaling

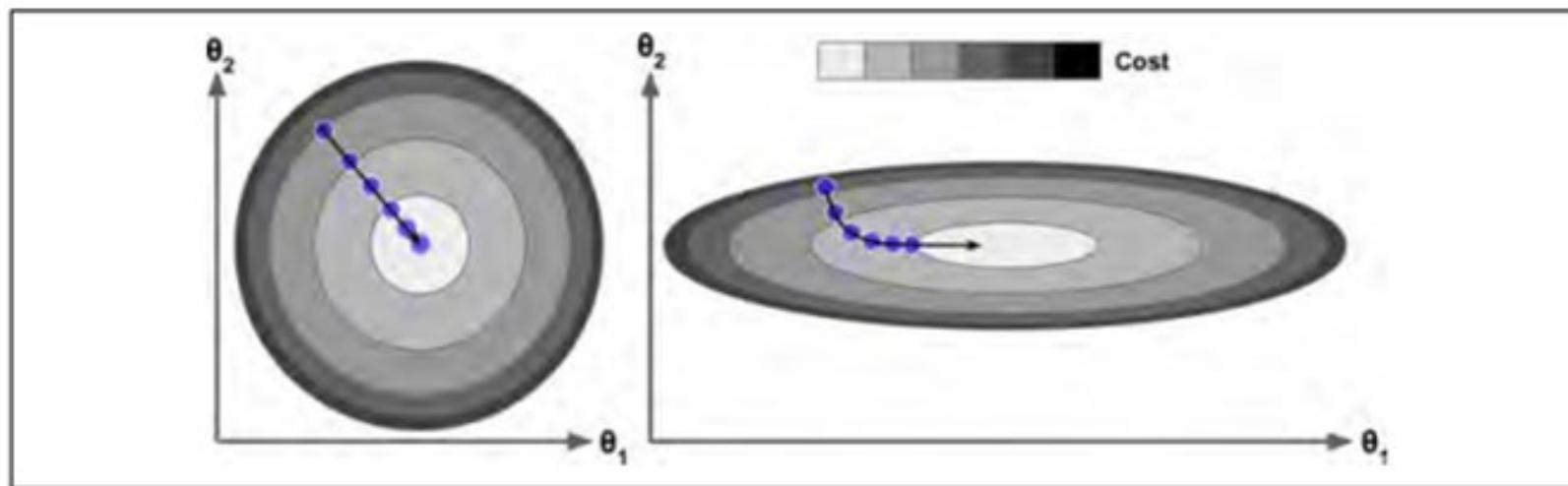


Figure 4-7. Gradient Descent with and without feature scaling

# Gradient Descent algorithm : Effect of Learning Rate :



## Problem Type 2 – Interpretation of Convergence or Effect of Hyper parameters

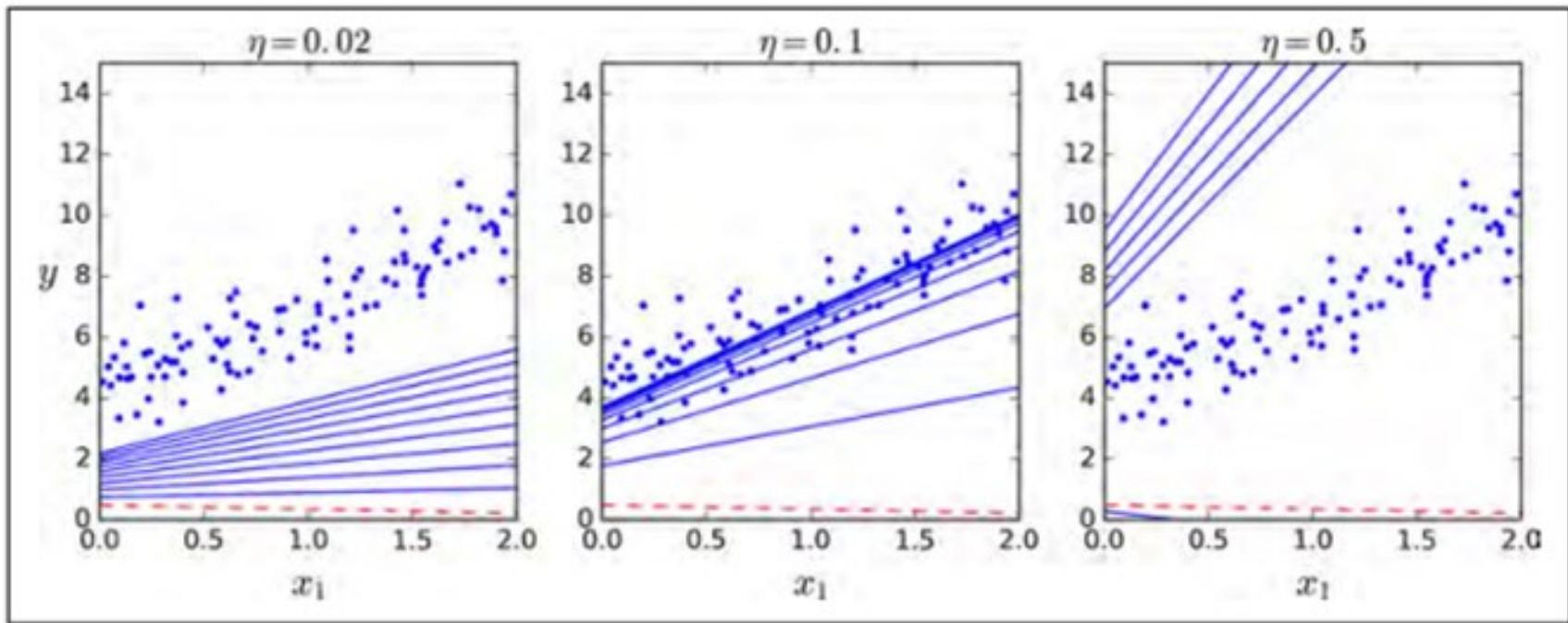


Figure 4-8. Gradient Descent with various learning rates

# Fit a linear Regression Line :

## Gradient Descent : Problem Type 3

---

### Steps :

(Assuming : 'n' no.of.instances and two predictors { $x_1, x_2$ } and linear regression)

1. Identification of the equations  $y = w_0 + w_1X_1 + W_2X_2$
2. Cost function & derivative
  1.  $W_0` = w_0 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y))$
  2.  $W_1` = w_1 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y) * x_1)$
  3.  $W_2` = w_2 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y) * x_2)$
3. Apply the equations

# Fit a linear Regression Line :

## Gradient Descent

Fit a linear regression. Show only the first iteration of Gradient descent algorithm using learning rate of **0.02** for the following data , if the Relative Risk of Coronary Heart Disease is believed to be only linearly dependent on BMI as well as Diastolic Pressure. Assume the intercept of the regression model as **5** and the slope of independent variables as **-0.03 (negative)**.

Patient	Systolic Pressure mm Hg	Diastolic Pressure mm Hg	BMI	Waist Circumference Threshold cm	RR-CHD (Relative Risk of Coronary Heart Disease)
1	140	80	35	100	1.81
2	120	80	25	80	1.22
3	130	100	30	60	1.71

### Steps :

1. Identification of the equations  $y = w_0 + w_1X_1 + w_2X_2$
2. Cost function & derivative
  1.  $w_0' = w_0 - \frac{1}{3} * \text{learning rate} * (\sum (w_0 + w_1X_1 + w_2X_2 - y))$
  2.  $w_1' = w_1 - \frac{1}{3} * \text{learning rate} * (\sum (w_0 + w_1X_1 + w_2X_2 - y) * x_1)$
  3.  $w_2' = w_2 - \frac{1}{3} * \text{learning rate} * (\sum (w_0 + w_1X_1 + w_2X_2 - y) * x_2)$
3. Apply the equations

# Fit a linear Regression Line :

## Gradient Descent

Fit a linear regression. Show only the first iteration of Gradient descent algorithm using learning rate of **0.02** for the following data , if the Relative Risk of Coronary Heart Disease is believed to be only linearly dependent on BMI as well as Diastolic Pressure. Assume the intercept of the regression model as **5** and the slope of independent variables as **-0.03 (negative)**.

Patient	Systolic Pressure mm Hg	Diastolic Pressure mm Hg	BMI	Waist Circumference Threshold cm	RR-CHD (Relative Risk of Coronary Heart Disease)
1	140	80	35	100	1.81
2	120	80	25	80	1.22
3	130	100	30	60	1.71

### Steps :

1. Identification of the equations:  $RR-CHD = 5 - 0.03 * BMI - 0.03 * DiastolicPressure$
2. Cost function & derivative
  1.  $W0' = w0 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD))$   
 $= 5 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD))$
  2.  $W1' = w1 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD) * BMI)$   
 $= -0.03 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD) * BMI)$
  3.  $W2' = w2 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD) * DiastolicPressure)$   
 $= -0.03 - 1/3 * 0.02 * (\text{sum } (5-0.03BMI-0.03DiastolicPressure - RRCHD) * DiastolicPressure)$
3. Apply the equations : Answer at the end of first iteration:  
 $W0 = 5.0016$  ,  $W1 = 0.0476$ ,  $w2= 0.179$

# Closed Form Solution Vs. Gradient Descent

---

## Gradient Descent

- Requires multiple iterations
- Need to choose  $\alpha$
- Works well when  $n$  is large
- Can support incremental learning

## Closed Form Solution

- Non-iterative
- No need for  $\alpha$
- Slow if  $n$  is large
  - Computing  $(X^T X)^{-1}$  is roughly  $O(n^3)$

# References

---

- T1 - Chapter 1 – Machine Learning, Tom Mitchell
- Chapter 1, 2 – Introduction to Machine Learning, 2<sup>nd</sup> edition, Ethem Alpaydin
- R1 – Chapter #1, # 3,#4 (Christopher M. Bishop, Pattern Recognition & Machine Learning) & Refresh your MFDS course basics

---

# Thank you !

## Next Session Plan:

- Variants of Gradient Descent
- Frequent Evaluation Metrics for Linear Regression Models
- Notion of Linear Basis Functions
- Notion of Bias vs Variance
- Ways to Handle Overfitting
- Types of Regularization