

BITS Pilani
Pilani | Dubai | Goa | Hyderabad

MACHINE LEARNING (AIML ZG565)



Session 3&4
(17th August, 2025)

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
- 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 |

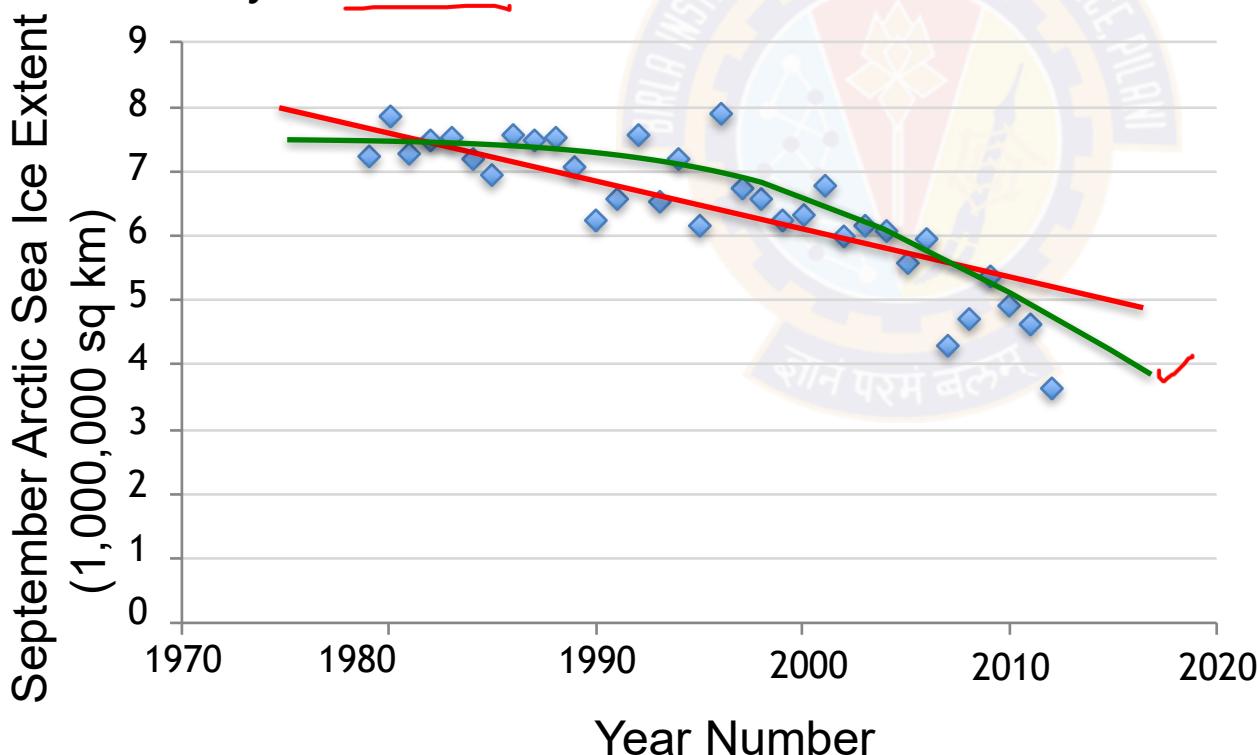
Formal Representation: Interpretation

Supervised Learning: Regression

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

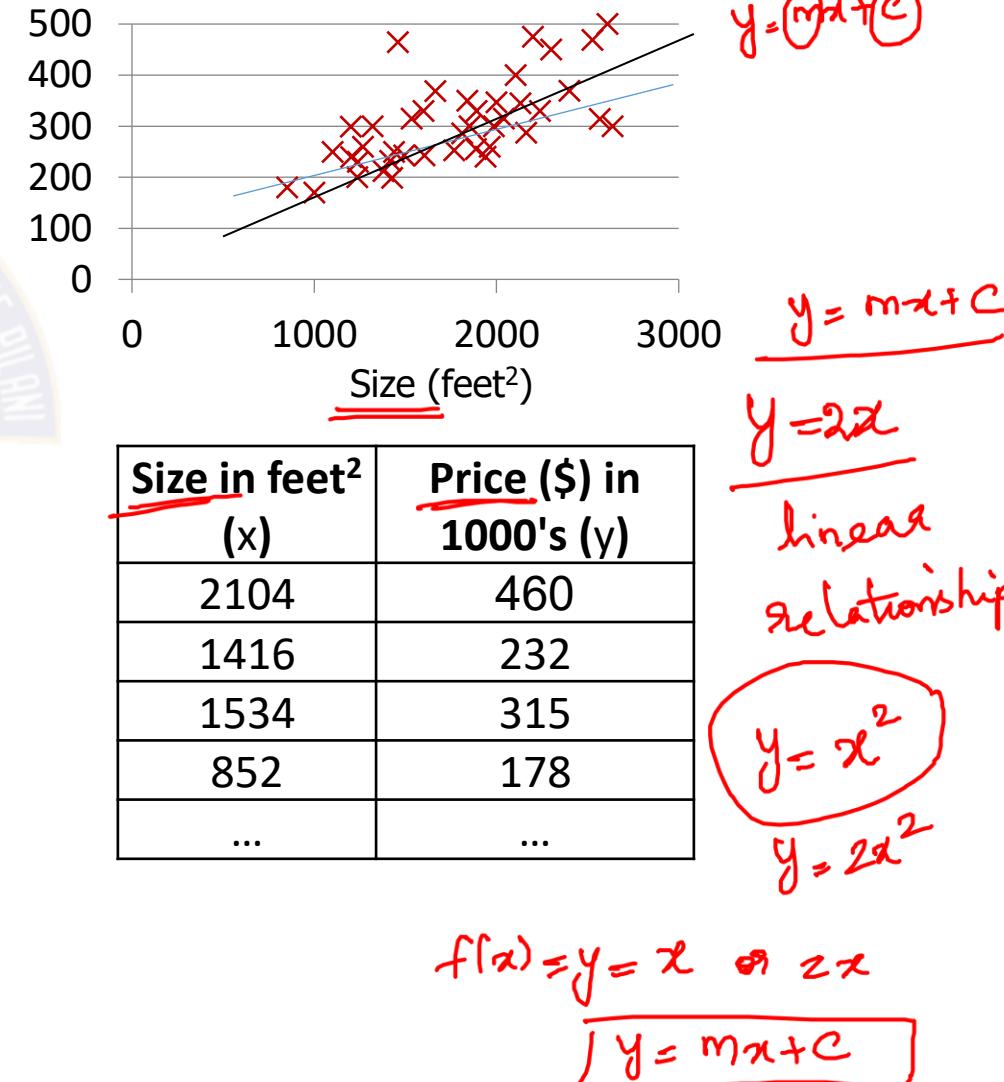
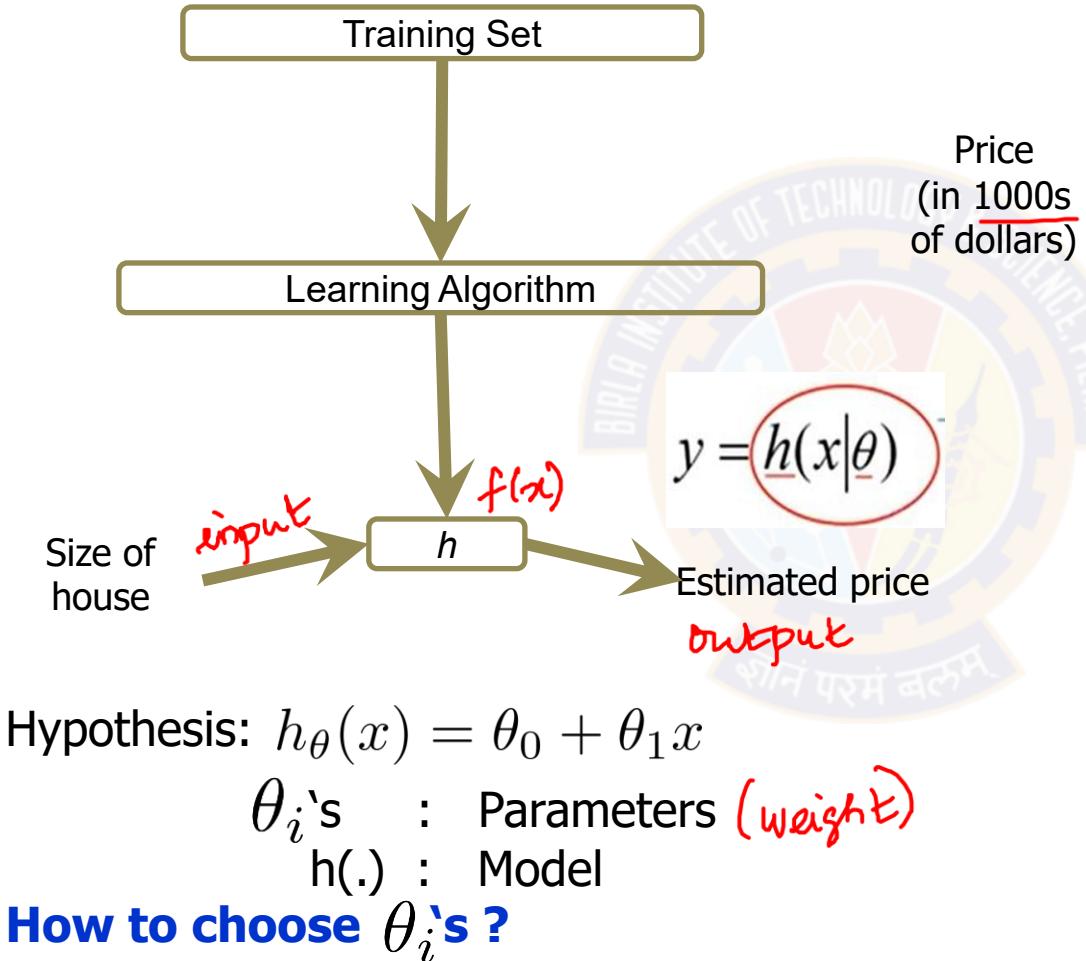
- Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Learn a function $f(x)$ to predict y given x
 - y is real-valued

input value, output value

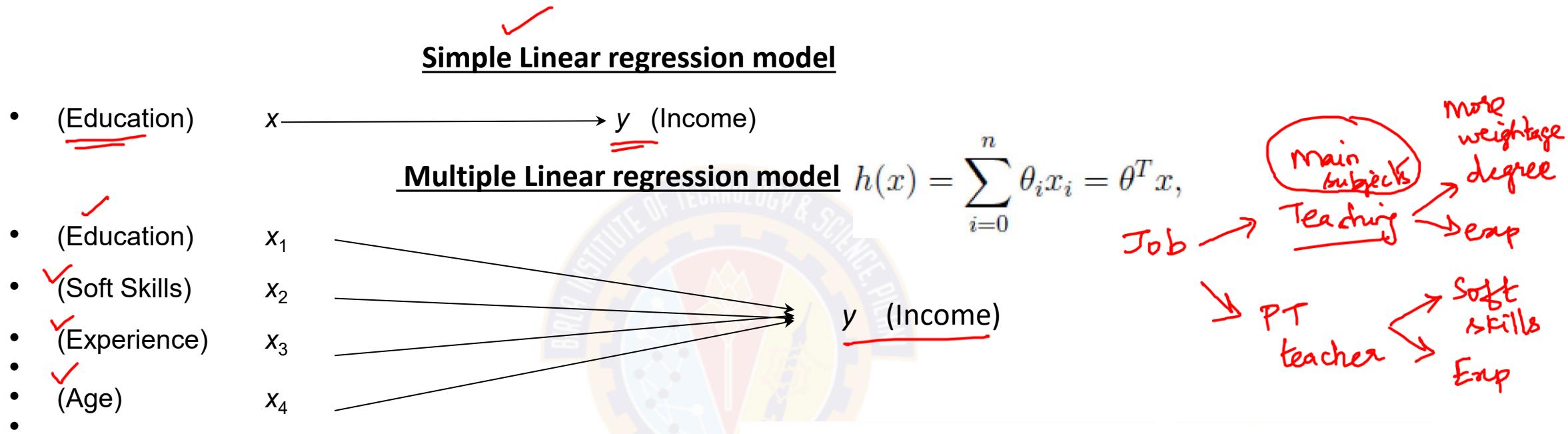


Machine Learning Process : Interpretation

Simple linear regression



Types of Regression Models

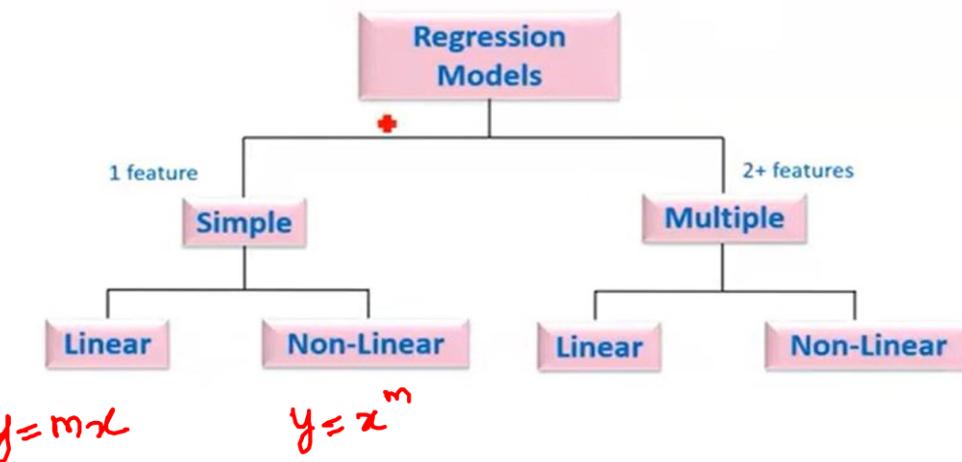


Multiple linear regression model

- 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 $x_0 = 1$

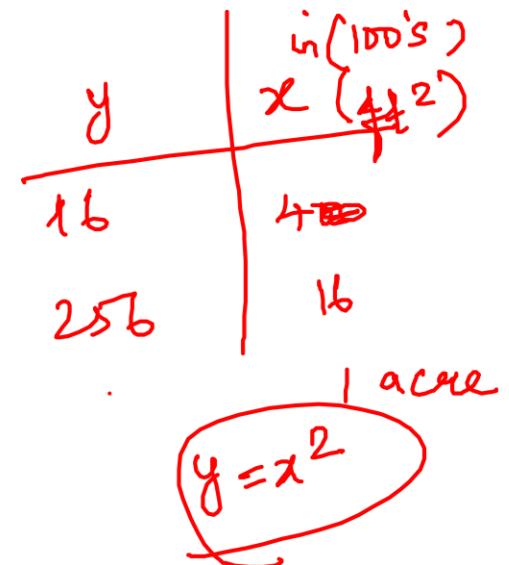


Multiple non-linear functions regression equation

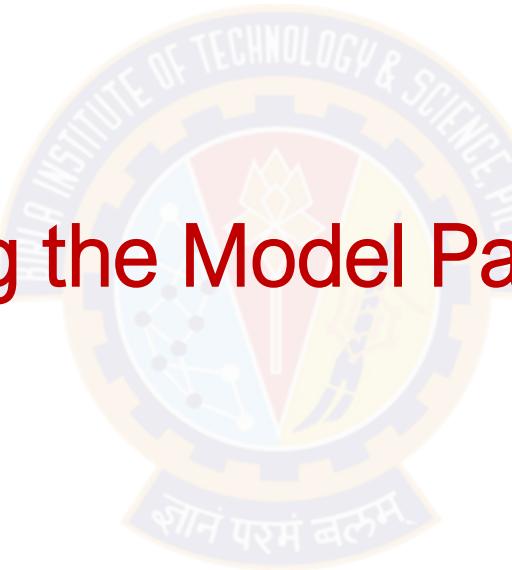
$$y = x^3 + 3x^2 + 4x + C$$

when NOT to use linear regression ?

- 1) When the data is categorical / class
- 2) When the features are highly correlated
- 3) If the relationship is non-linear
- 4) If the data has outliers, strongly affects your results.



Learning the Model Parameters



Analytical

Closed Form Solution Approach

Gradient Descent Approach

How to choose θ_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



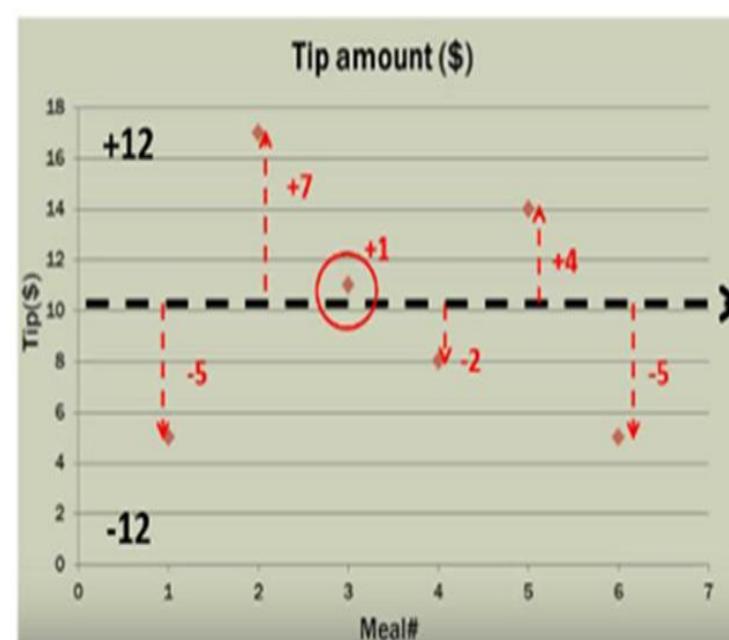
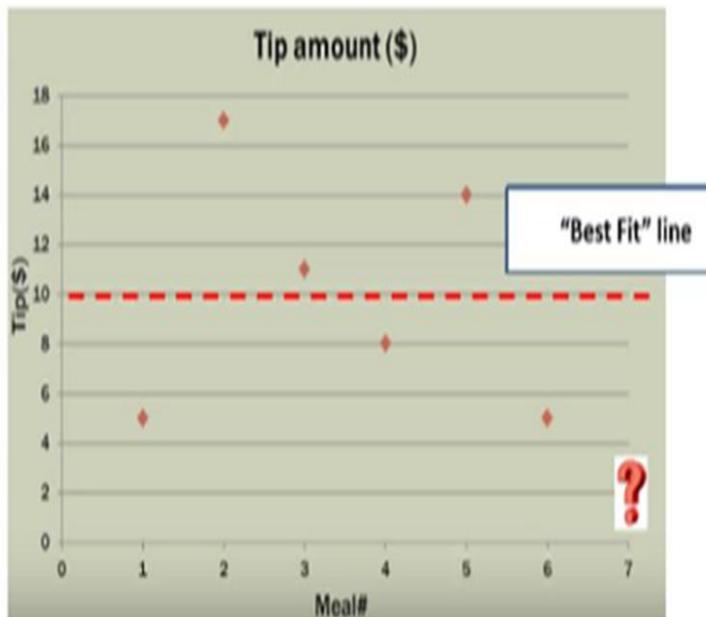
$$y = \theta_0$$

Notion of Cost Function

$$y_{act} - y_{obs}$$

| 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 ² |
|-------|----------|-----------------------|
| 1 | -5 | 25 |
| 2 | +7 | 49 |
| 3 | +1 | 1 |
| 4 | -2 | 4 |
| 5 | +4 | 16 |
| 6 | -5 | 25 |

Sum of Squared Errors (SSE) = 120

↓
minimize this error

Compute θ_1 (slope):

Using least-square method:

$$\theta_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$$

| Meal | Total bill x_i | TIP amt y_i | $x_i - \bar{x}$ | $y_i - \bar{y}$ | $(x_i - \bar{x})(y_i - \bar{y})$ | $(x_i - \bar{x})^2$ |
|------|---------------------|------------------|-----------------|-----------------|----------------------------------|---------------------|
| 1 | 34 | 5 | -40 | -5 | 200 | 1600 |
| 2 | 108 | 17 | 34 | 7 | 238 | 1156 |
| 3 | 64 | 11 | -10 | 1 | -10 | 100 |
| 4 | 88 | 8 | 14 | -2 | -28 | 196 |
| 5 | 99 | 14 | 25 | 4 | 100 | 625 |
| 6 | 51 | 5 | -23 | -5 | 115 | 529 |

$$\bar{x} = 74$$

$$\bar{y} = 10$$

$$\sum(x_i - \bar{x})(y_i - \bar{y}) = 200 + 238 - 10 - 28 + 100 + 115 = 615$$

$$\sum(x_i - \bar{x})^2 = 1600 + 1156 + 100 + 196 + 625 + 529 = 4206$$

$$\theta_1 = \frac{615}{4206} \approx 0.1462 \quad \checkmark$$



Compute θ_0 (intercept):

$$\theta_0 = \bar{y} - \theta_1 \bar{x} = 10 - (0.1462 \times 74)$$

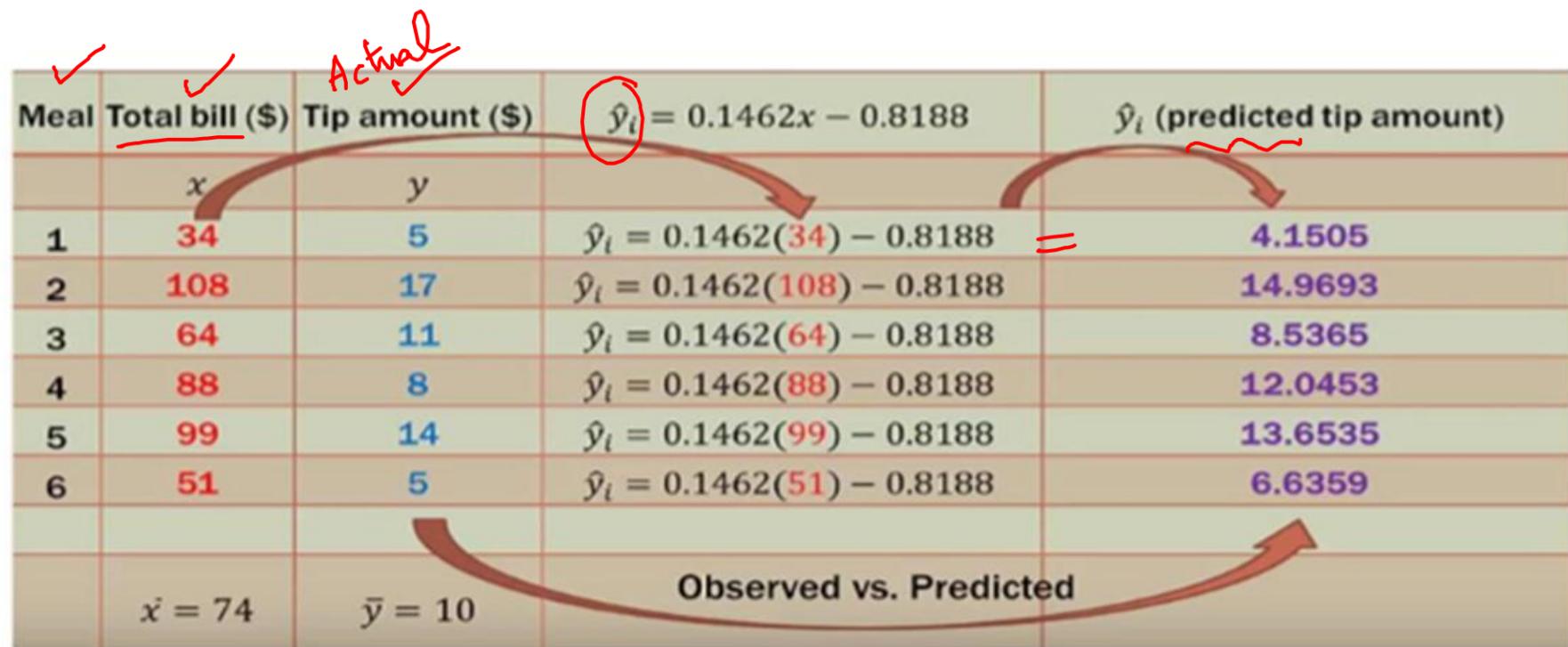
$$\theta_0 = 10 - 10.8188 = -0.8188$$

\therefore By Least-Square method

$$y = -0.8188 + 0.1462 x \quad \checkmark$$

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 | |



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

$$y = mx + c$$

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

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

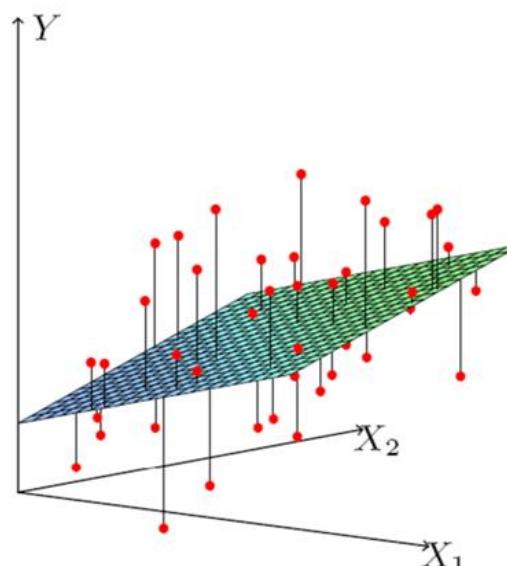
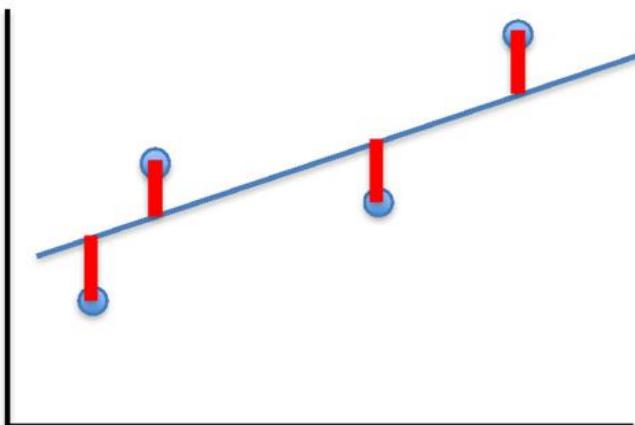
Sum of squared Errors

Linear Regression : Least Squares

- Cost Function

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

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



- Hypothesis

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

- Parameters

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

- Cost Function

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

- 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 $\underline{\boldsymbol{\theta} = [\theta_0, \theta_1]}$



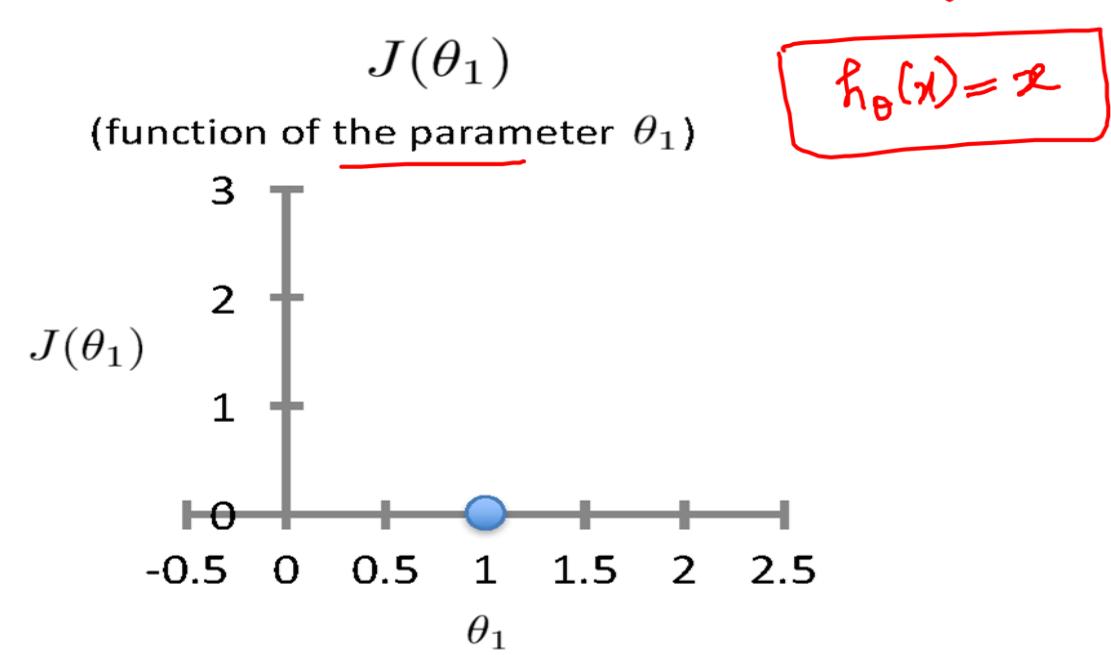
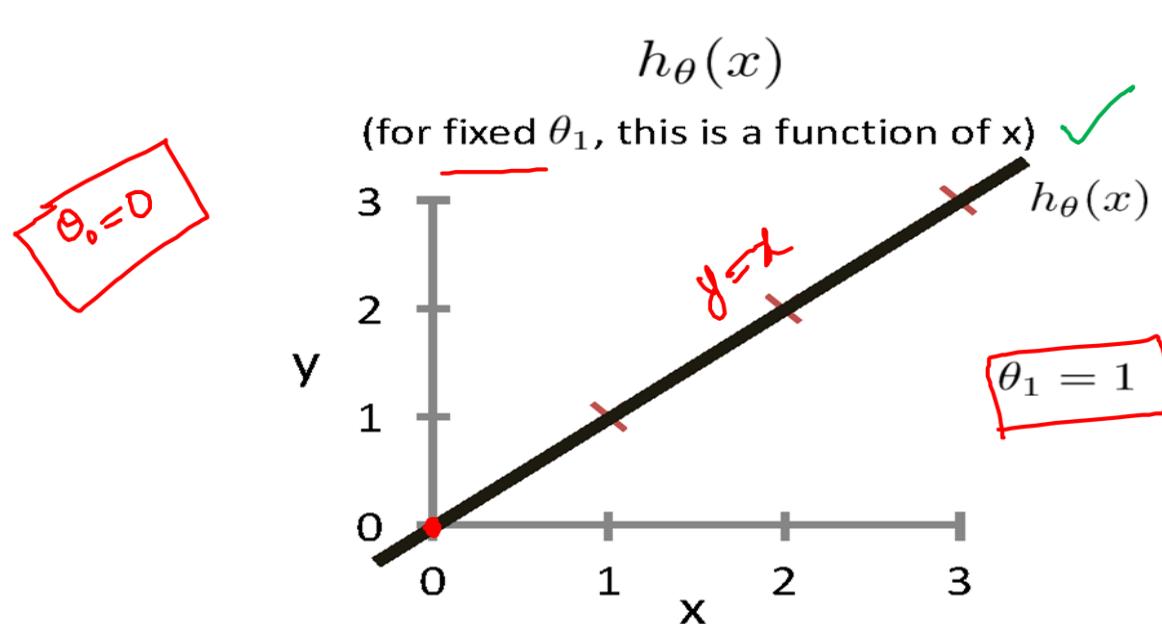
Intuition Behind Cost Function

$\theta_0 \rightarrow$ intercept

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

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

$$\begin{cases} h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x \\ \text{when } \theta_1 = 1 \\ \therefore h_{\boldsymbol{\theta}}(x) = \theta_0 + x \end{cases}$$



Based on example

Source Credit : by Andrew Ng

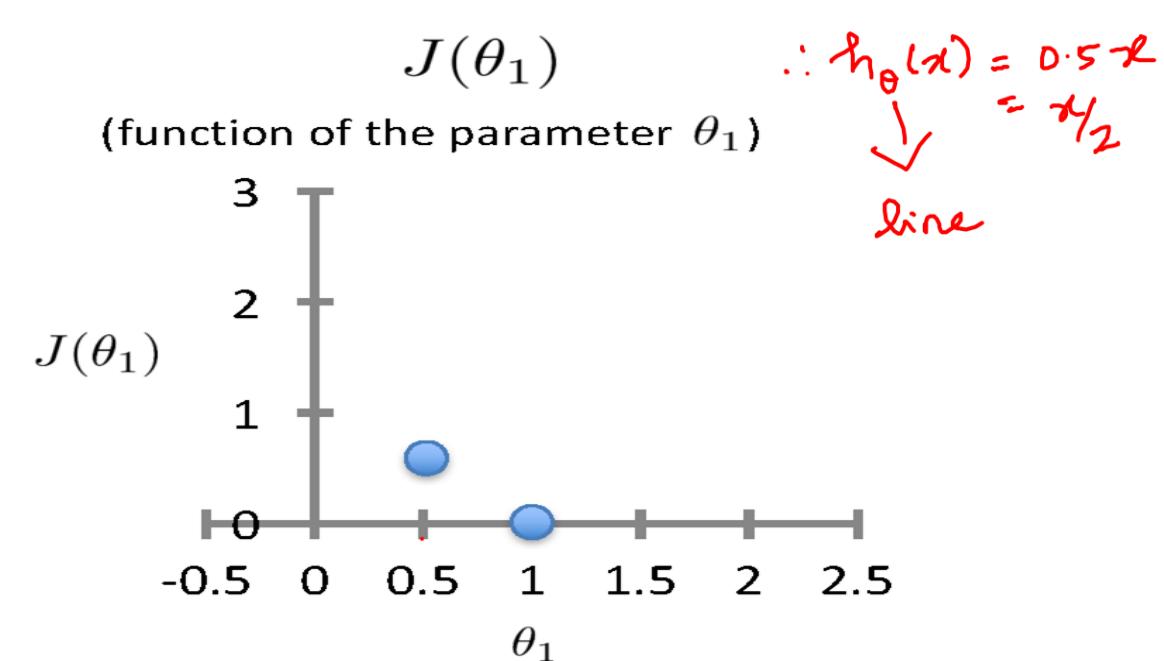
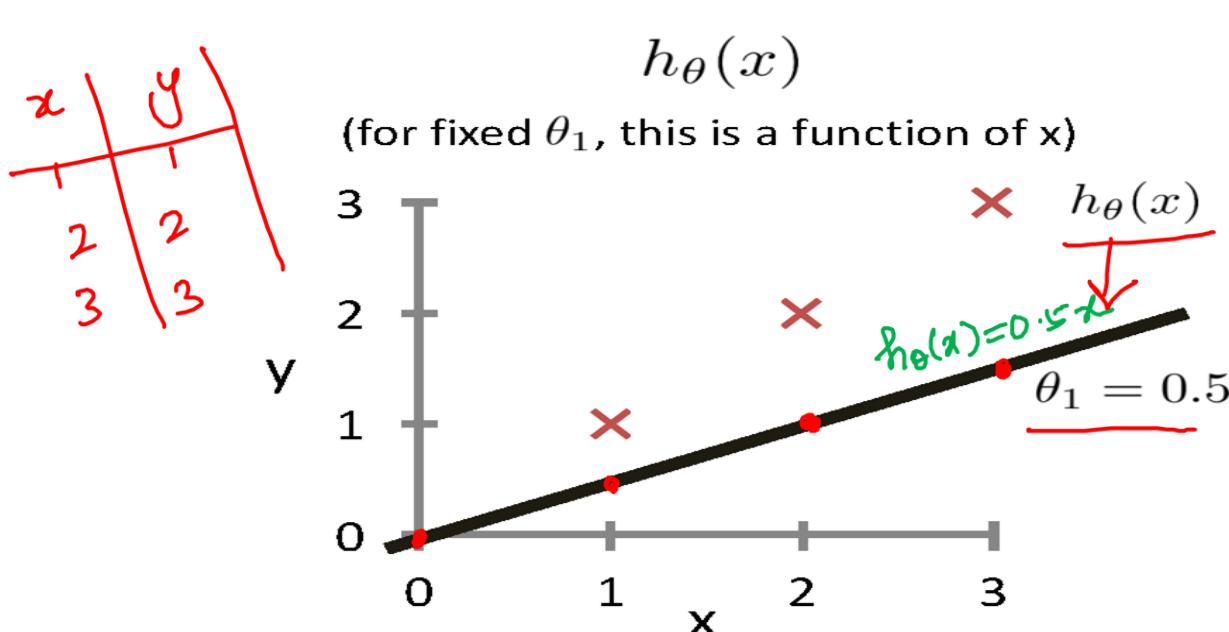
Intuition Behind Cost Function

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

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

so you know

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$ $\theta_0 = 0, \theta_1 = 0.5$



Based on example
by Andrew Ng

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

| Assume $\theta_0 = 0$, $\theta_1 = 0.5$ | | |
|--|-------|--------------------------------|
| x_i | y_i | $h_{\theta}(x) = 0.5x$ |
| 1 | 1 | $h_{\theta}(1) = 0.5$ |
| 2 | 2 | $h_{\theta}(2) = 0.5(2) = 1$ |
| 3 | 3 | $h_{\theta}(3) = 0.5(3) = 1.5$ |

Compute $J(\theta)$ for $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \Rightarrow h_{\theta}(x) = 0.5x$

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n [h_{\theta}(x_i) - y_i]^2$$

sum of squared errors

WKT $n=3$

$$\downarrow \text{no. of data} = \frac{1}{2(3)} \left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right] = 3.5/6 \approx 0.58$$

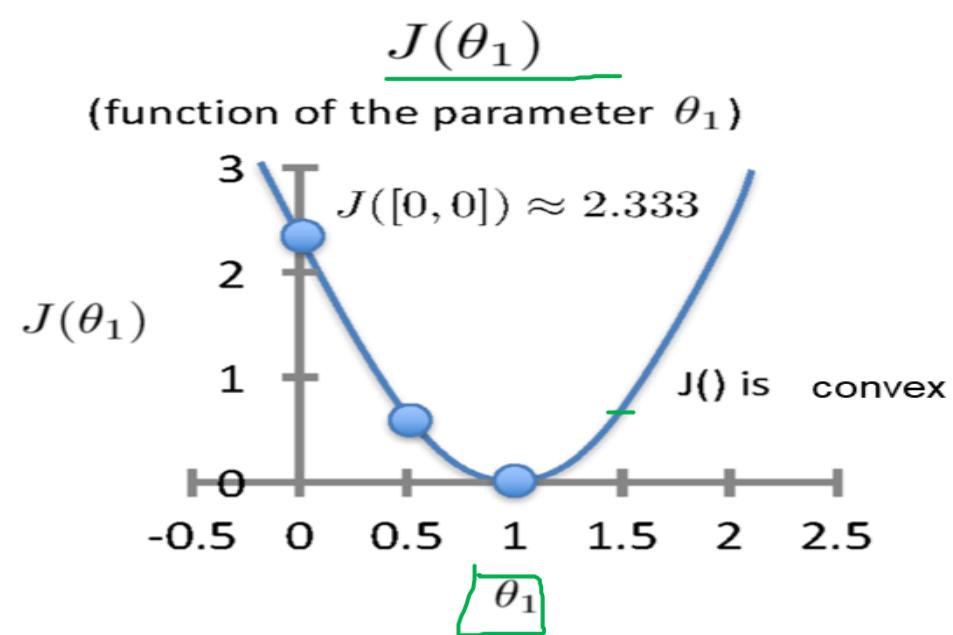
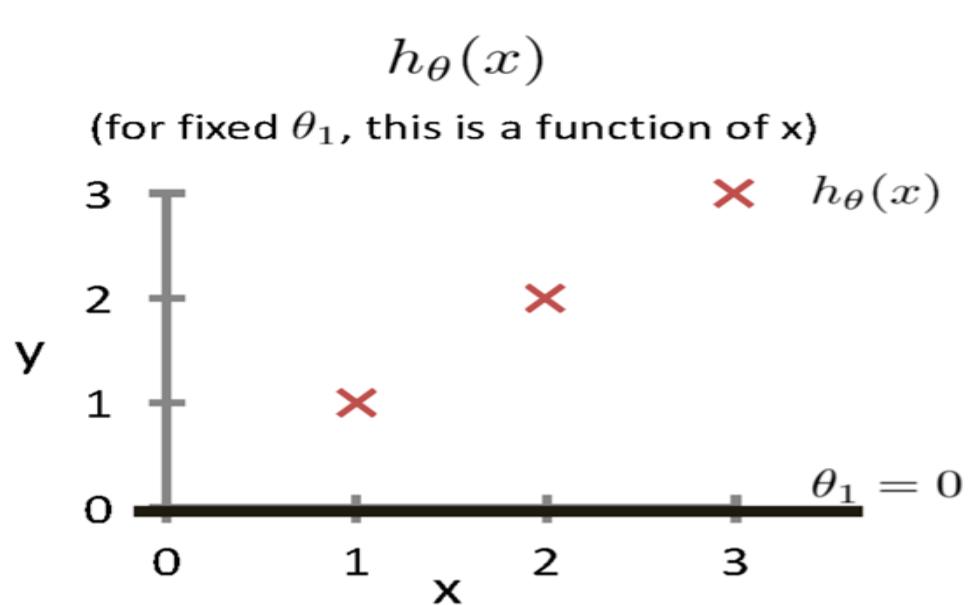
(ii) let $\theta = [0, 1] \Rightarrow h_{\theta}(x) = x$

$$J(\theta) = J(1) = \frac{1}{6} \left[(1-1)^2 + (2-2)^2 + (3-3)^2 \right] = 0$$

Intuition Behind Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$
$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x$$

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



Based on example

Source Credit : by Andrew Ng

$\theta_0 = 0, \theta_1 = 1$
Best fit line = $h_{\boldsymbol{\theta}}(x) = x$

Notion of Maxima and Minima of a function - In General

$$y = \text{slope } m x + C$$

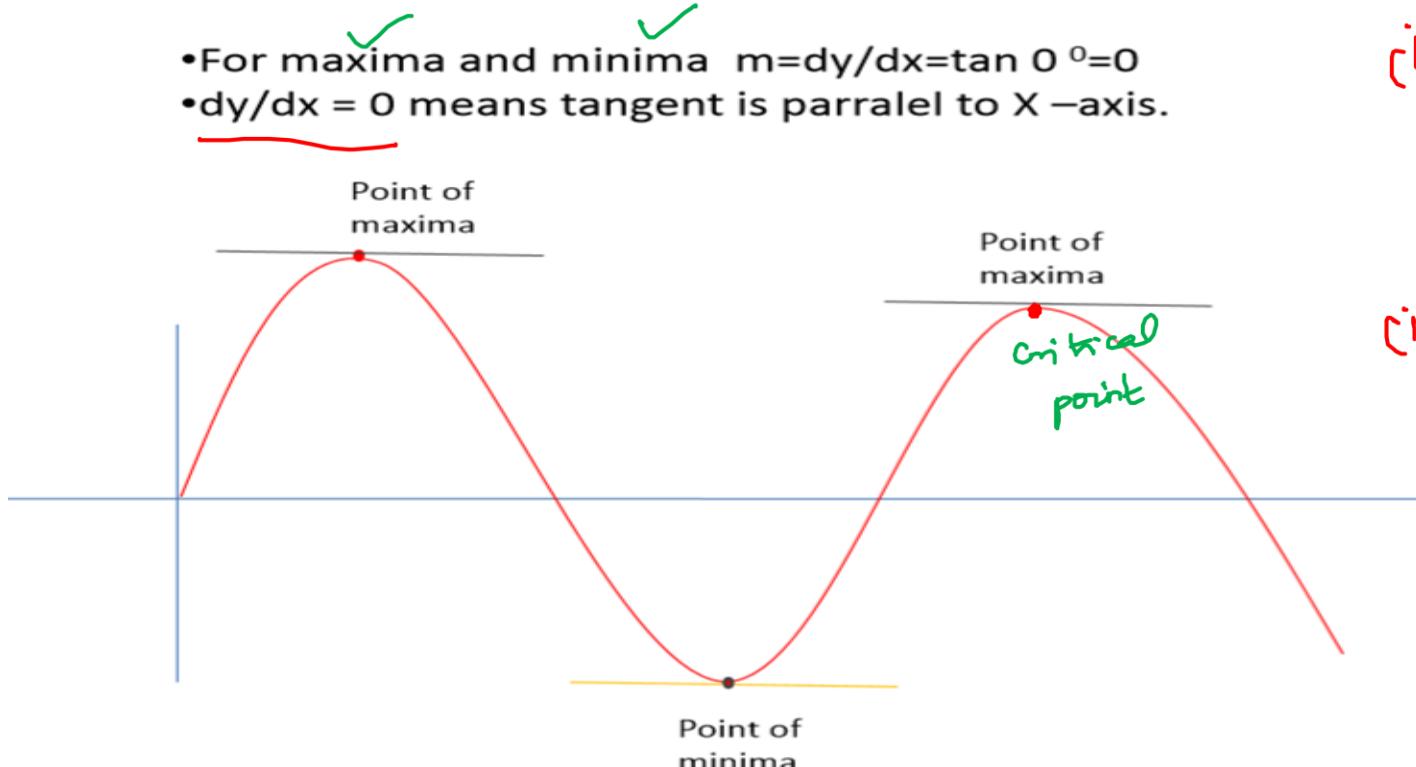
(i) Step 1 : Find derivative

(ii) Step 2 : Find the critical pts

$$\frac{df}{dx} = 0$$

(iii) Step 3 :
 $\frac{d^2f}{dx^2} < 0 \rightarrow \text{Maximum}$

$\frac{d^2f}{dx^2} > 0 \rightarrow \text{Minimum}$



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

Eg: $f(x) = x^2 - 4x + 3$ Find the pt of maxima/minima

Ans: $f'(x) = 2x - 4$

Step 1: $\frac{df}{dx} =$

Step 2: $f'(x) = 0$ find the critical points

$$\therefore 2x - 4 = 0$$

$$\Rightarrow \boxed{x = 2}$$

Step 3: Find the second derivative

$$\frac{d^2 f}{dx^2} = f''(x) = \frac{d}{dx}(2x - 4) = 2$$

$$f''(x) = 2$$

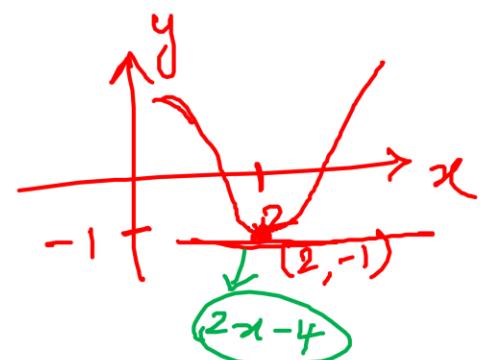
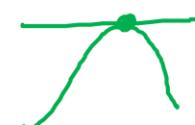
\therefore The function is minimum at $x=2$

What is the value of the fn. at $x=2$?

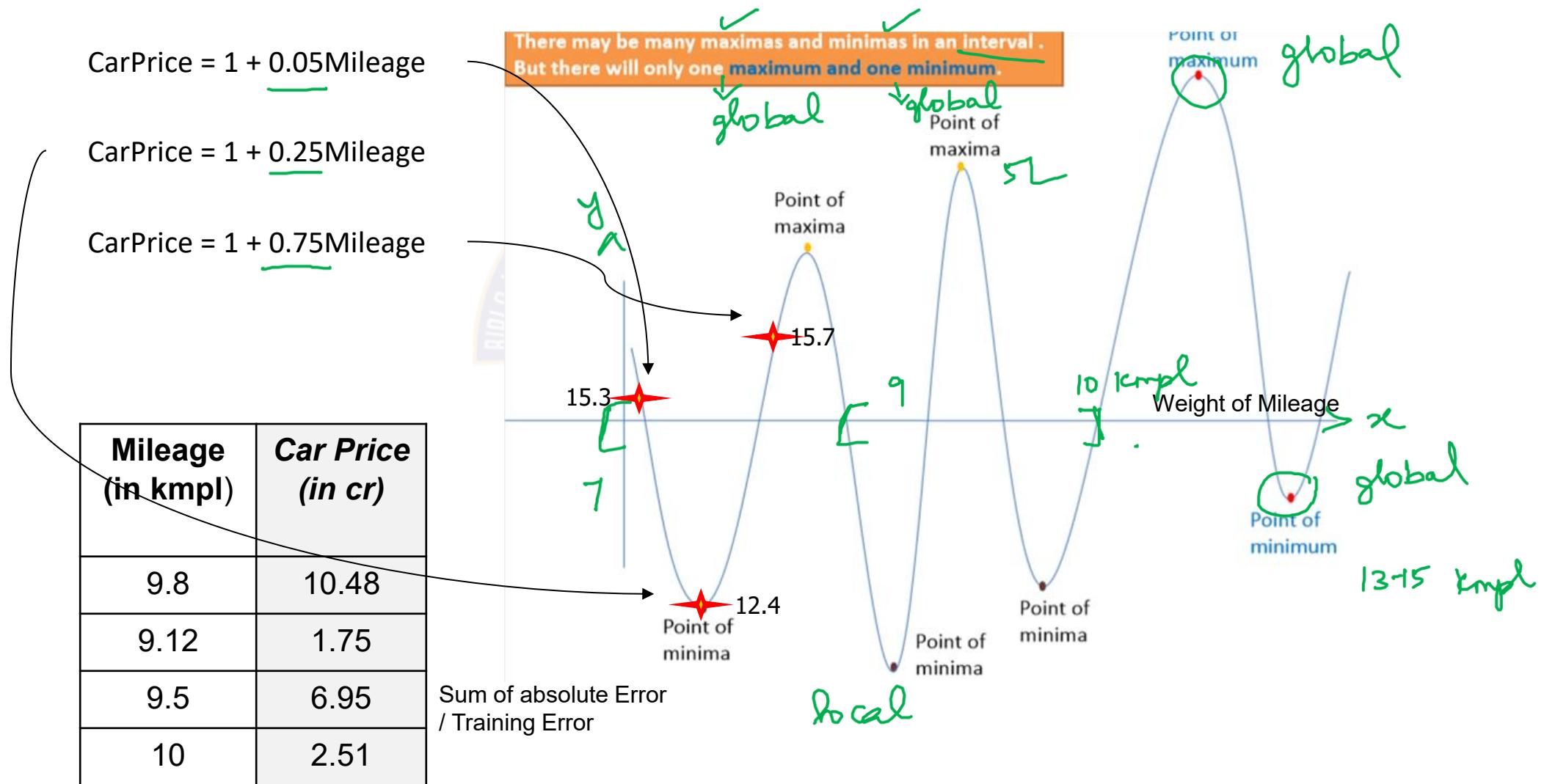
$$\text{Sub. } x=2 \text{ in } f(x) = 2^2 - 4(2) + 3 = -1$$

Suppose $f''(x) = -2$

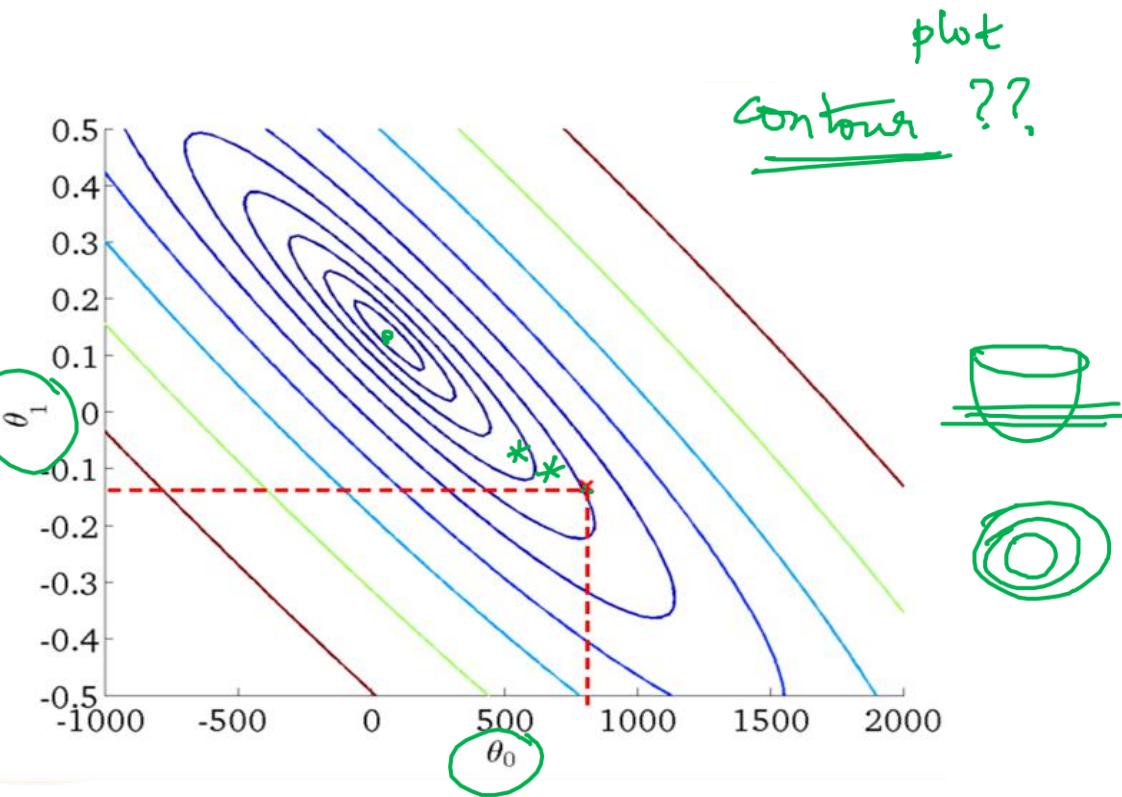
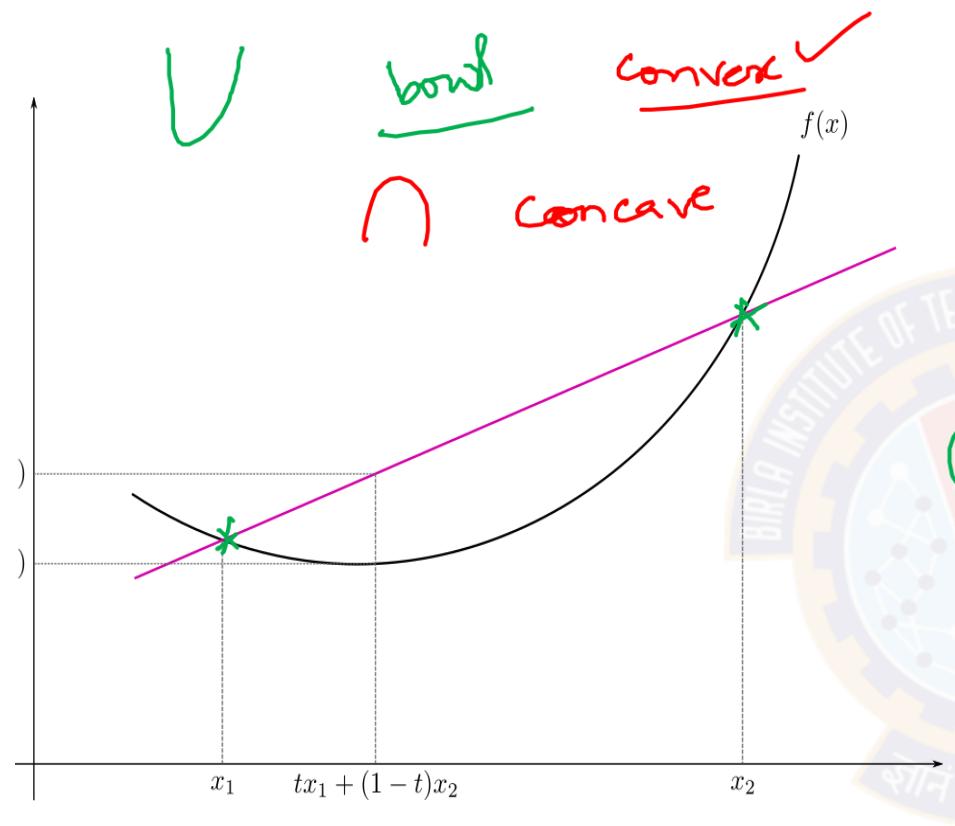
Then $f(x)$ will have maximum



Notion of Maxima and Minima of a function - In Machine Learning



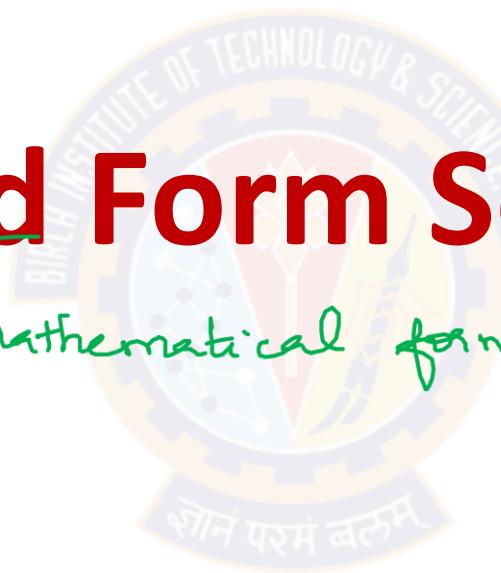
Convex function: Multivariate View



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

Mathematical formula to predict the values



Vectorization

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

- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

- Let

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

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

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

$$y = mx + c$$

| 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

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

column vector

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

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

Let:

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

$n \times 1$

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

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

~~$\boldsymbol{\theta} = n \times (d+1) \cdot (d+1) \times 1 = n \times 1$~~

| 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 |

Vectorization

- For the linear regression cost function:

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

$\sum_{i=1}^n v^2 = v^T v$
 for any vector, v

| 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 |

Let:

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

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

since $\theta^T \mathbf{x} = \mathbf{x}^T \theta$

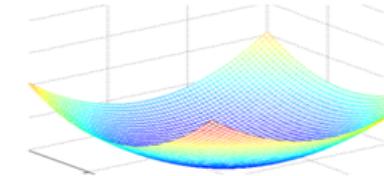
$$\theta^T = \begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \quad X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} \quad y = \begin{bmatrix} 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{bmatrix}$$

Vectorization

- Instead of using GD, solve for optimal θ analytically
 - Notice that the solution is when $\frac{\partial}{\partial \theta} J(\theta) = 0$
- Derivation:

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

1×1 (scalar)



Take derivative and set equal to 0, then solve for θ :

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

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

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

$$\theta = \frac{\mathbf{X}^\top \mathbf{y}}{\mathbf{X}^\top \mathbf{X}}$$

Closed Form Solution:

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

for large 'n'
Computationally expensive

~~invertible~~
invertible

$$(a-b)^T = a^T - b^T$$

$$(ab)^T = b^T a^T$$

$$\begin{aligned}\theta^T x^T y &= (y^T x \theta)^T \\ &= y x^T \theta^T\end{aligned}$$

Note:

$$\frac{\partial}{\partial \theta} (c) = 0$$

$$\frac{\partial}{\partial \theta} (b^T \theta) = b$$

$$\frac{\partial}{\partial \theta} (\theta^T A \theta) = 2A\theta$$

(If A is symmetric)

$$J(\theta) = \frac{1}{2n} [(x\theta - y)^T (x\theta - y)]$$

$$= \frac{1}{2n} [((x\theta)^T - y^T)(x\theta - y)]$$

$$= \frac{1}{2n} [(\theta^T x^T - y^T)(x\theta - y)]$$

$$= \frac{1}{2n} [\theta^T x^T x \theta - y \underbrace{\theta^T x^T - y^T x \theta}_{\text{scalars & equal}} + y y^T]$$

$$J(\theta) = \frac{1}{2n} [\underline{\theta^T x^T x \theta} - 2\theta^T x^T y + yy^T]$$

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{1}{2n} [2x^T x \theta - 2x^T y]$$

$\frac{\partial}{\partial \theta} \rightarrow$ single variable
 $\frac{\partial}{\partial \theta} \rightarrow$ multiple variables

Vectorization

- Can obtain θ by simply plugging \underline{X} and \underline{y} into

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

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

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

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

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

find inverse of this matrix

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

inverse of matrix

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

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

$$\theta_0 = 1, \theta_1 = 1$$

Can obtain θ by simply plugging X and y into

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

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

$$h_{\theta}(x) = 1 + x$$

| x No.of.Years of Experience <u>(in Years)</u> | y Salary Of the Employee <u>(in Lakhs)</u> |
|--|---|
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 5 |
| 5 | 6 |

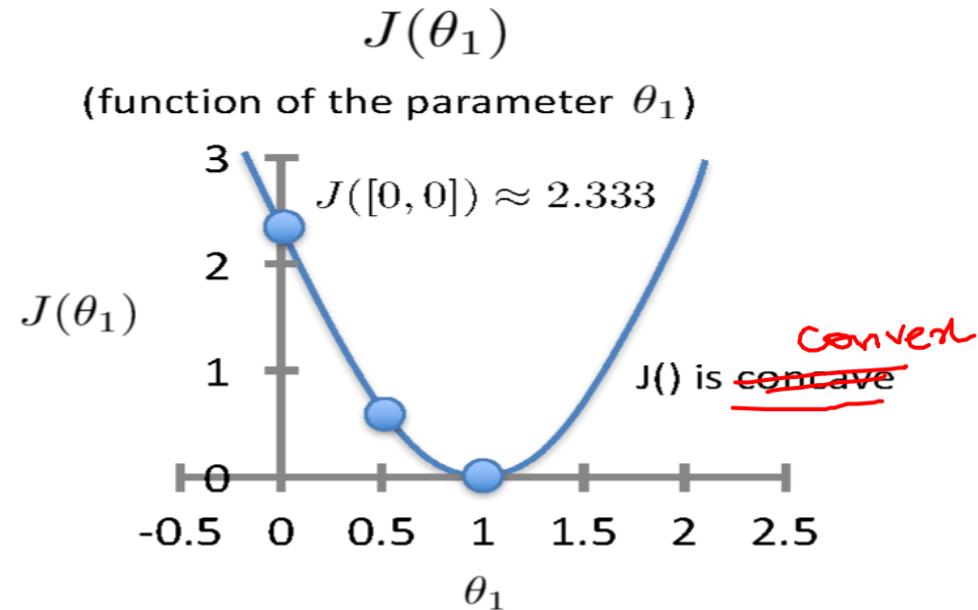
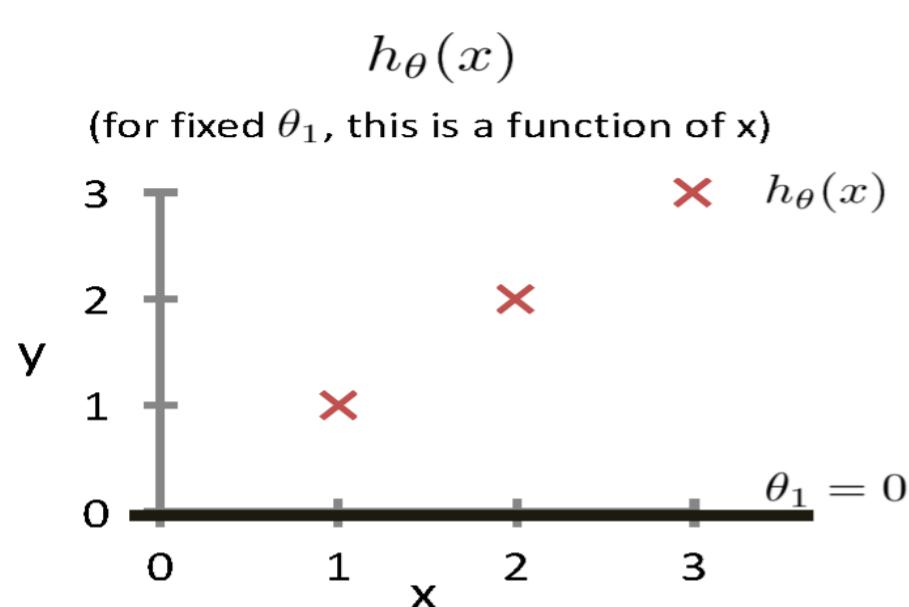
Gradient Descent Approach



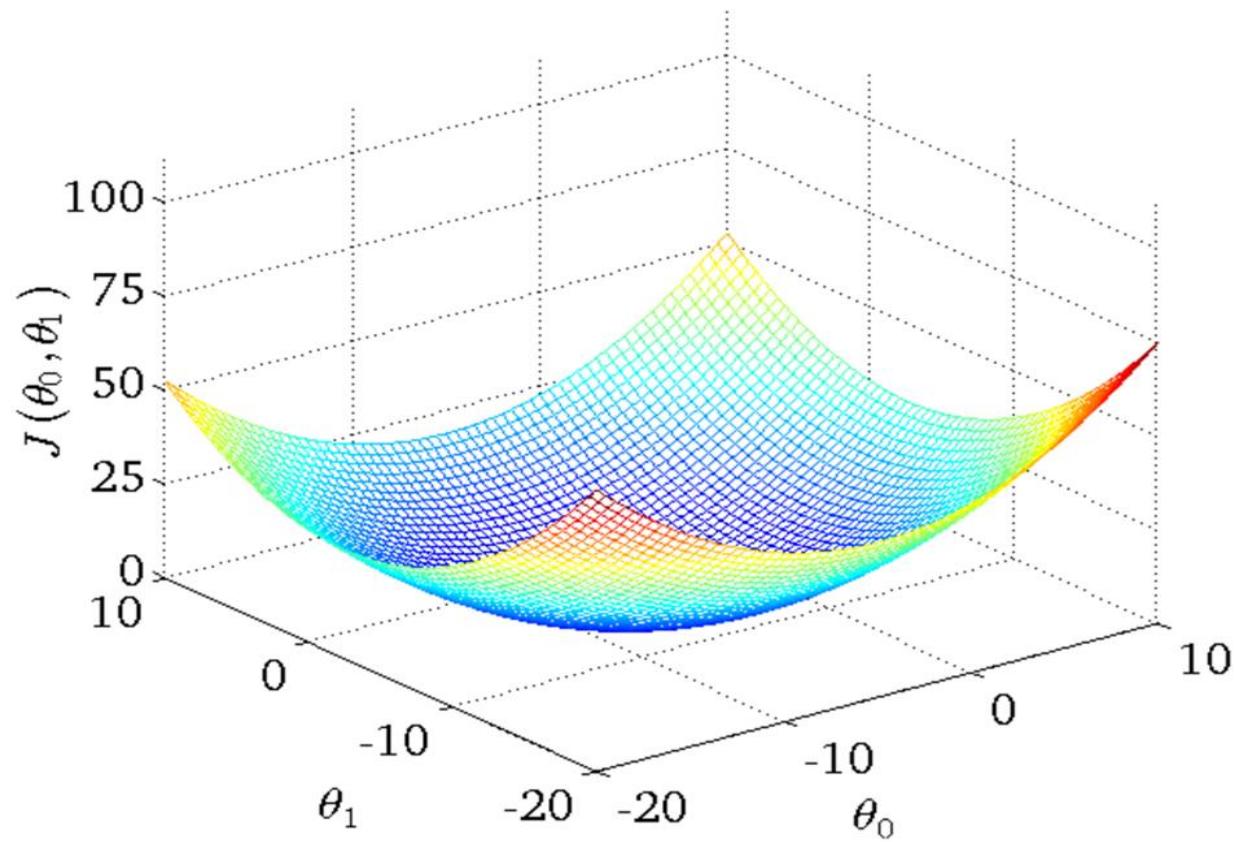
Intuition Behind Cost Function

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

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



Intuition Behind Cost Function

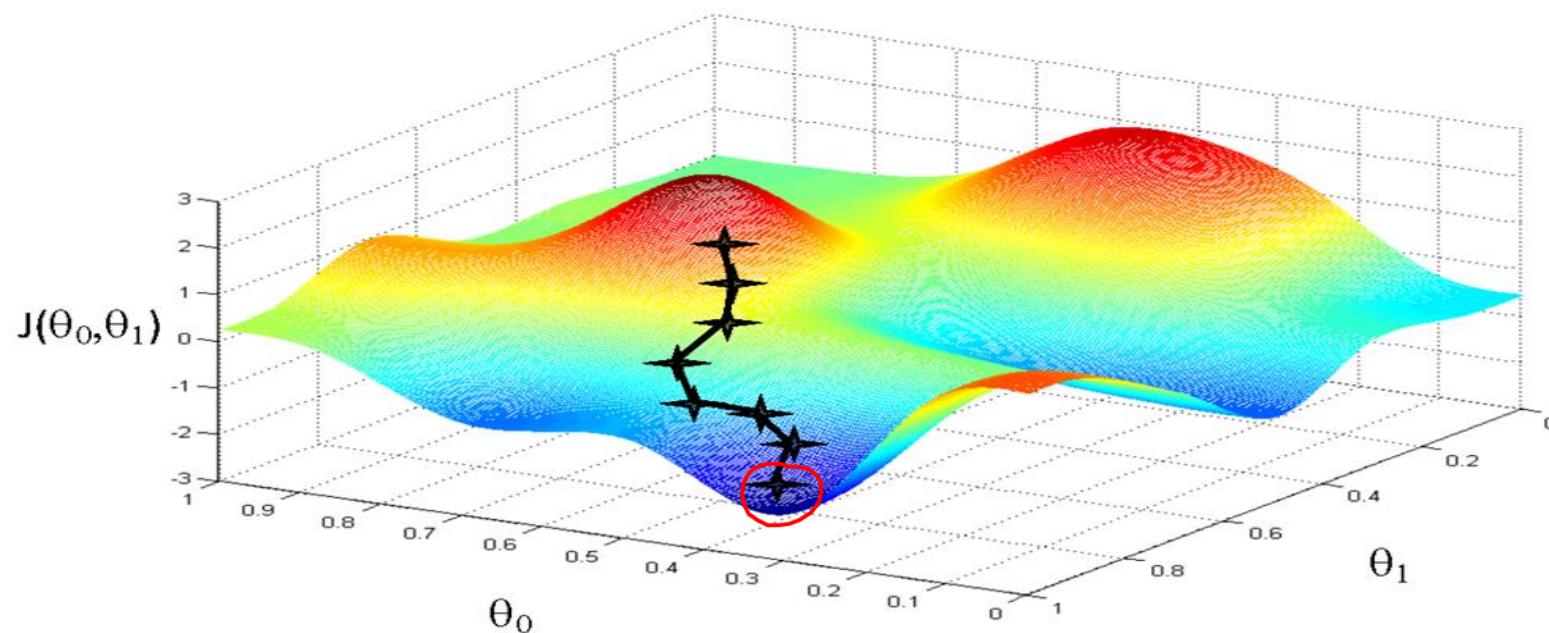


Based on example

Source Credit : by Andrew Ng

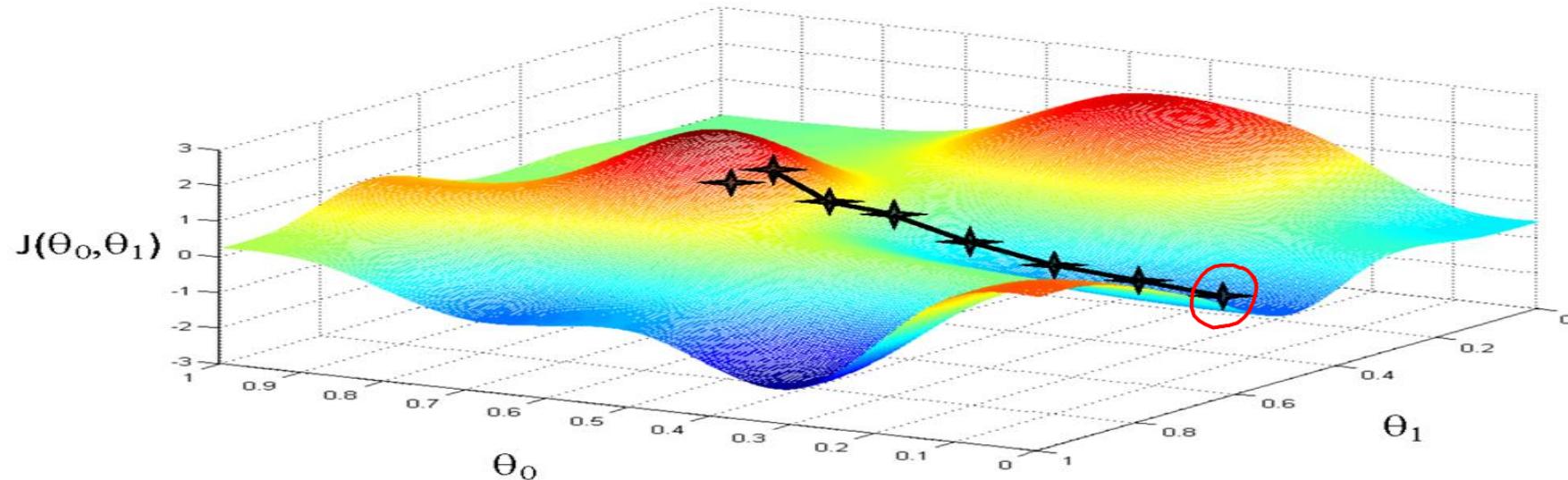
Basic Search Procedure

- Choose initial value for θ $\theta_0=0, \theta_1=1$
- Until we reach a minimum:
 - Choose a new value for θ to reduce $J(\theta)$ $\theta_1=0.5$



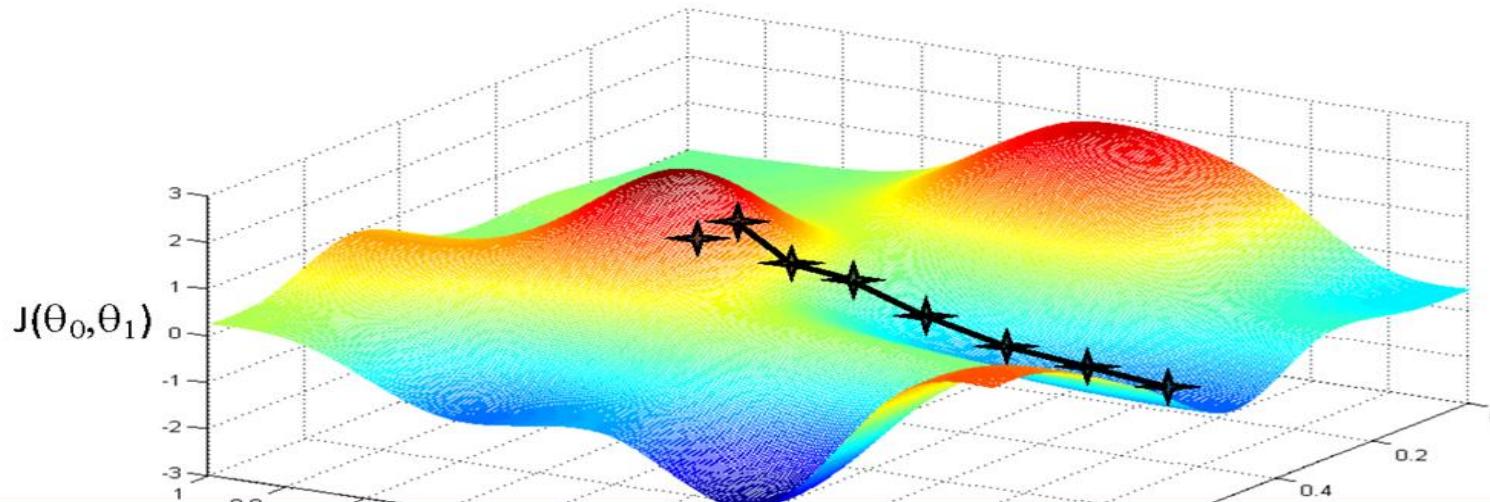
Basic Search Procedure

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



Basic Search Procedure

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



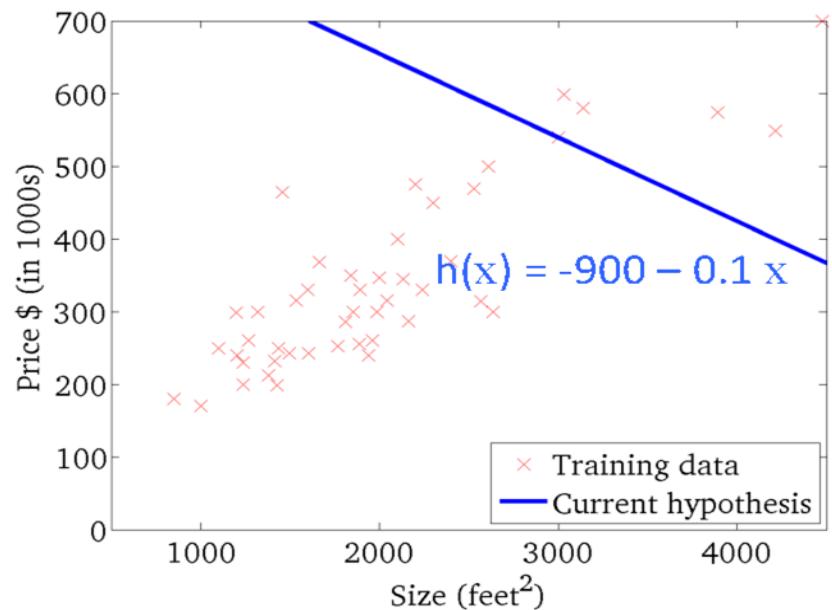
Since the least squares objective function is convex (concave),
we don't need to worry about local minima

Gradient Descent

| Size in feet ² (x) | Price (\$) in 1000's (y) |
|-------------------------------|--------------------------|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

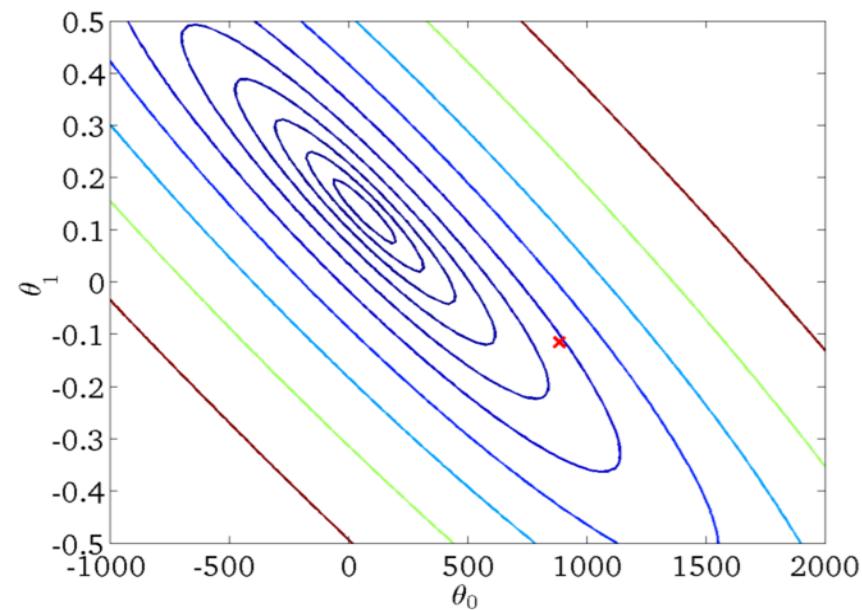
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)

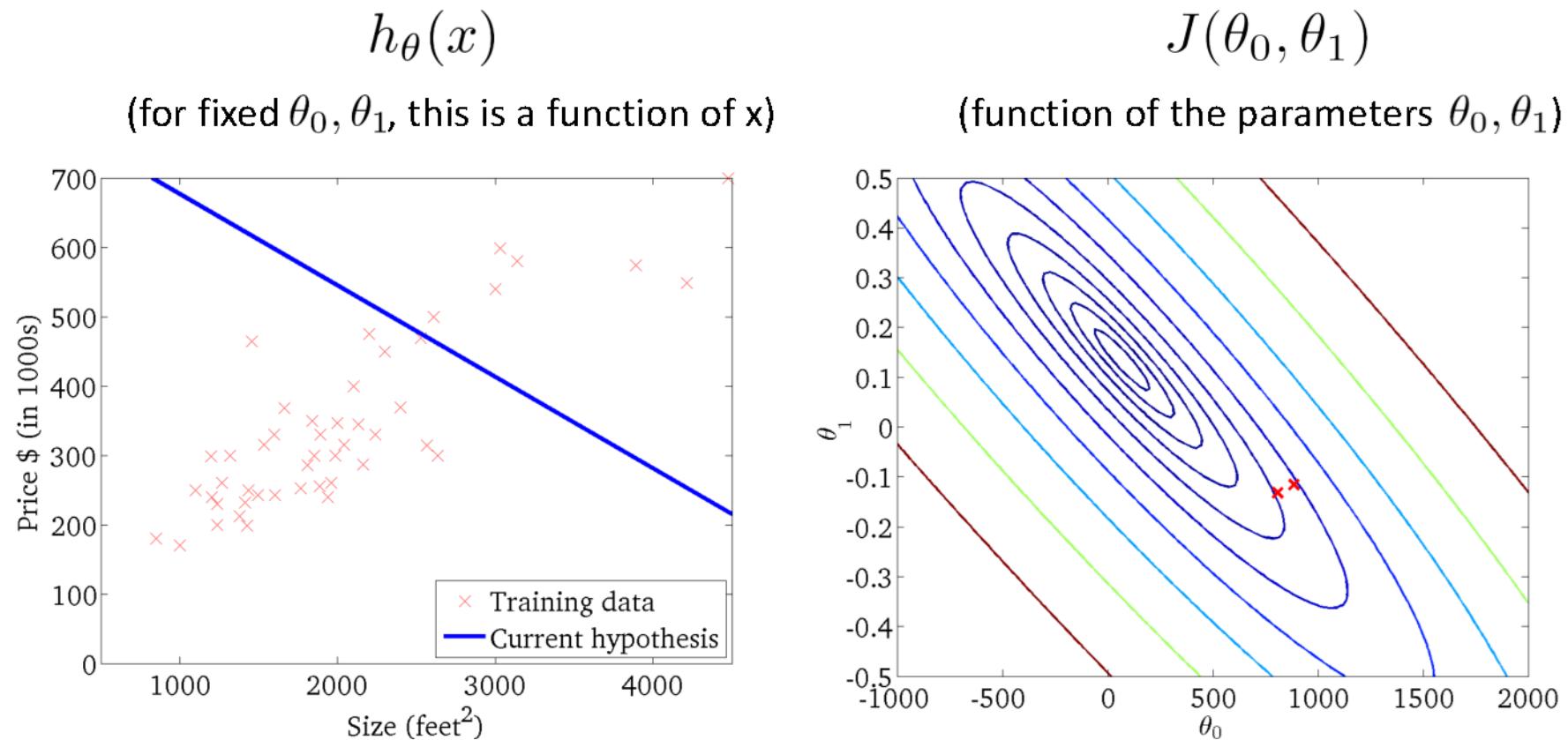


$$J(\theta_0, \theta_1)$$

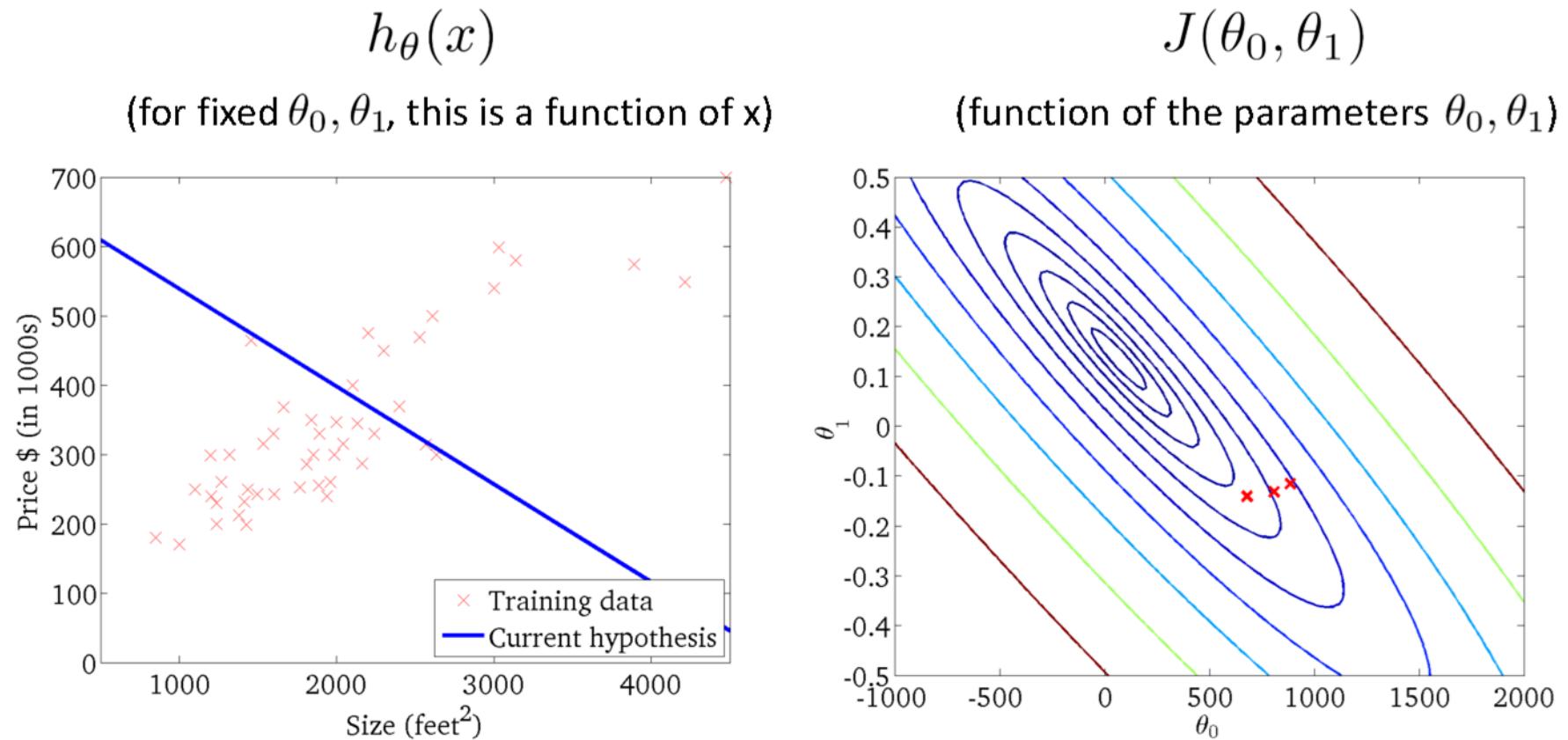
(function of the parameters θ_0, θ_1)



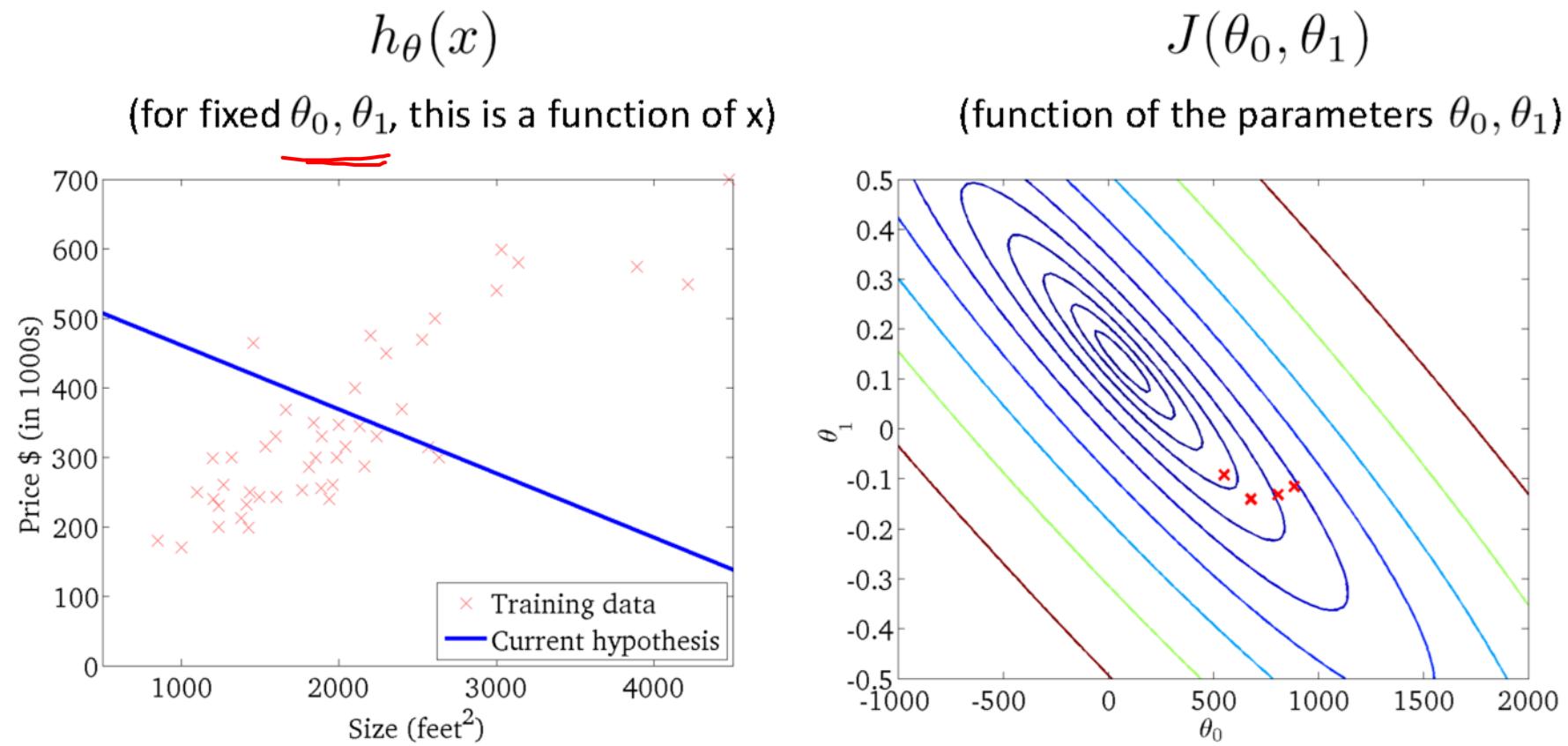
Gradient Descent



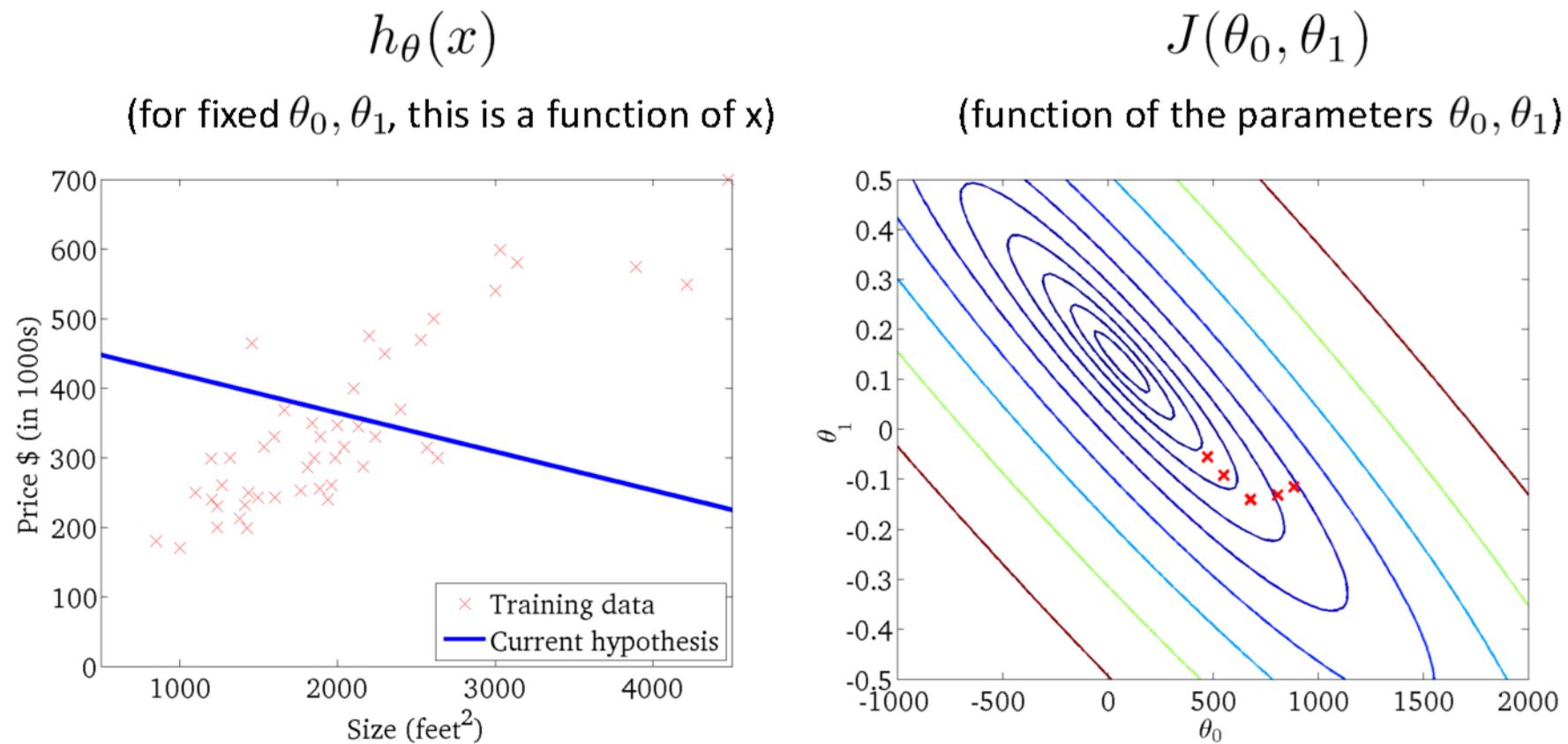
Gradient Descent



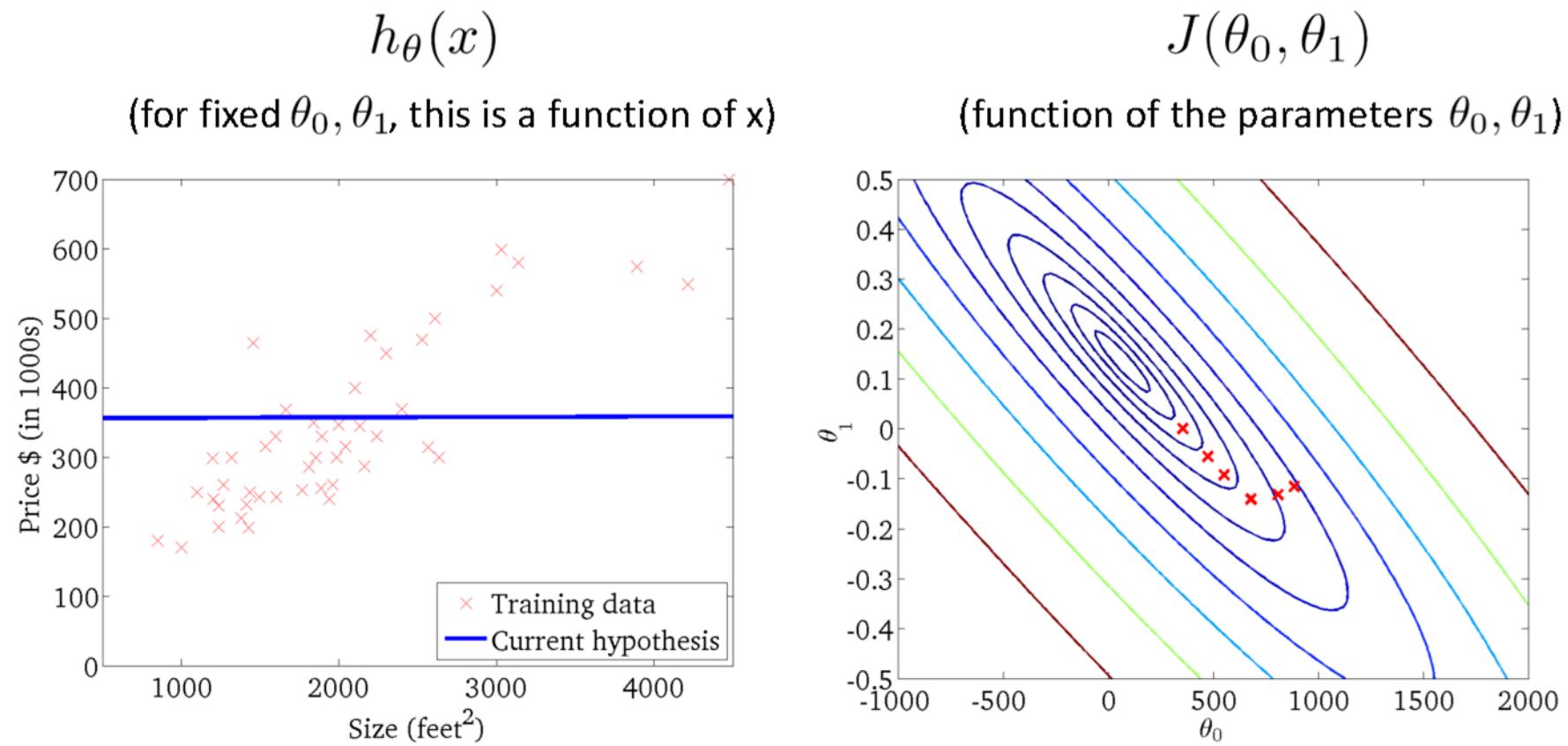
Gradient Descent



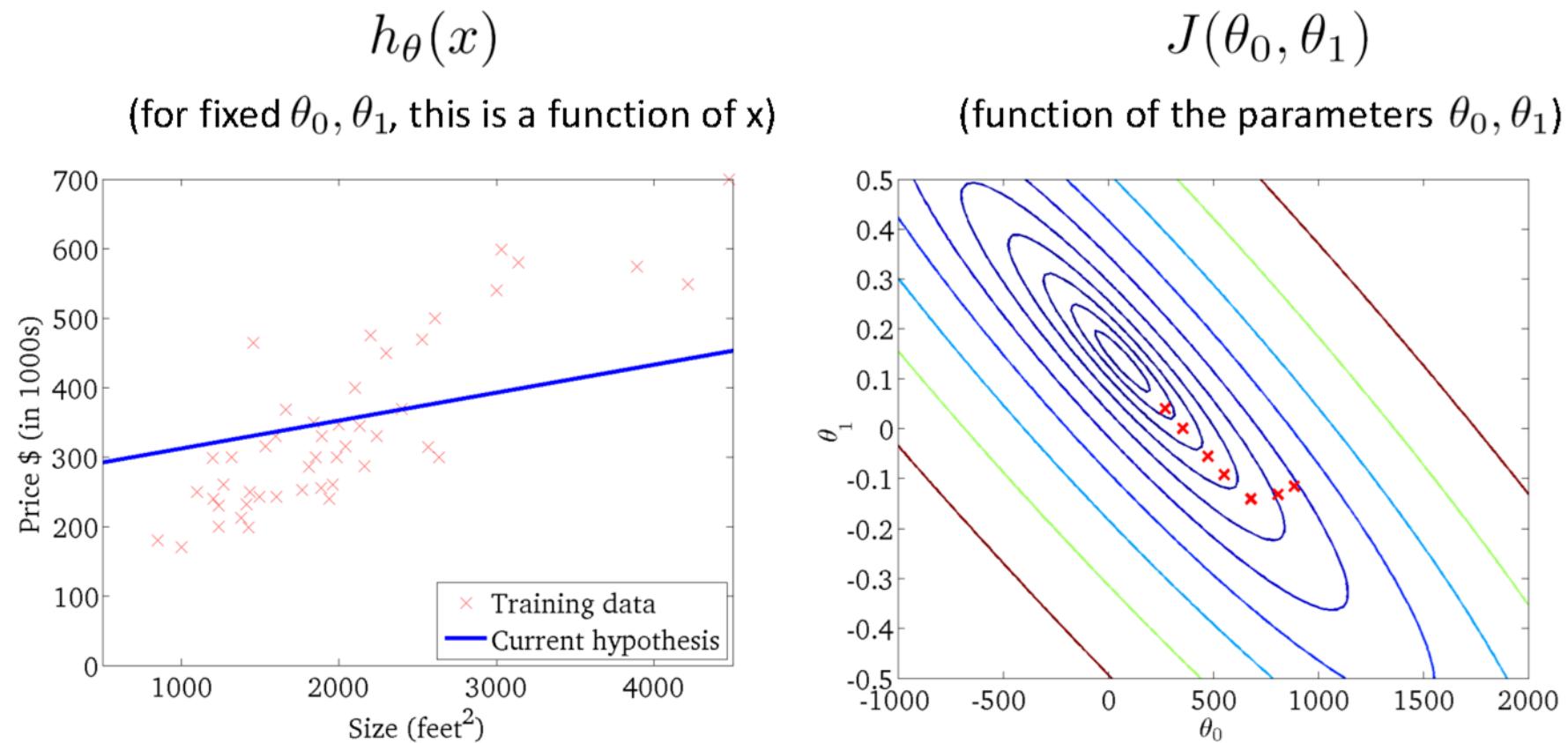
Gradient Descent



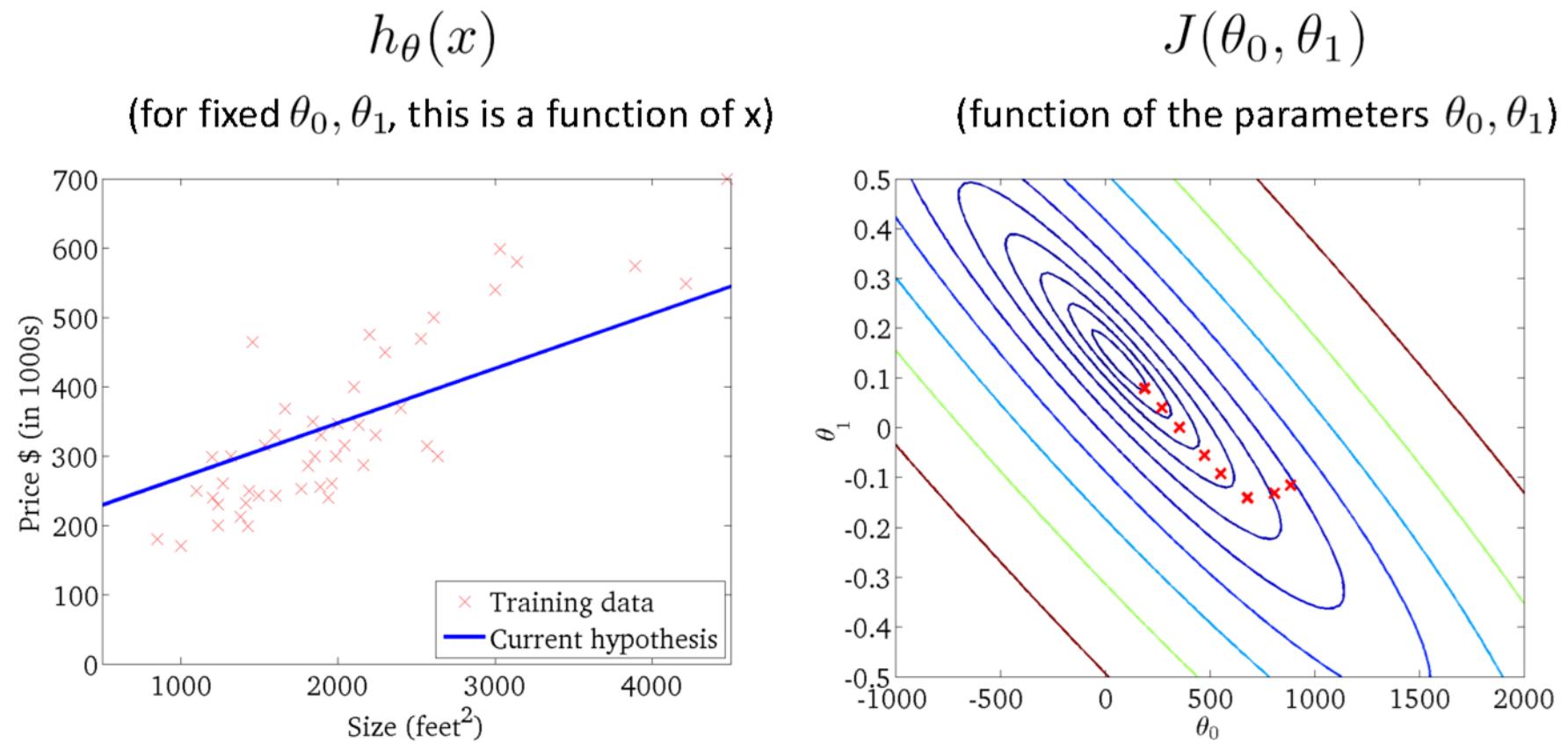
Gradient Descent



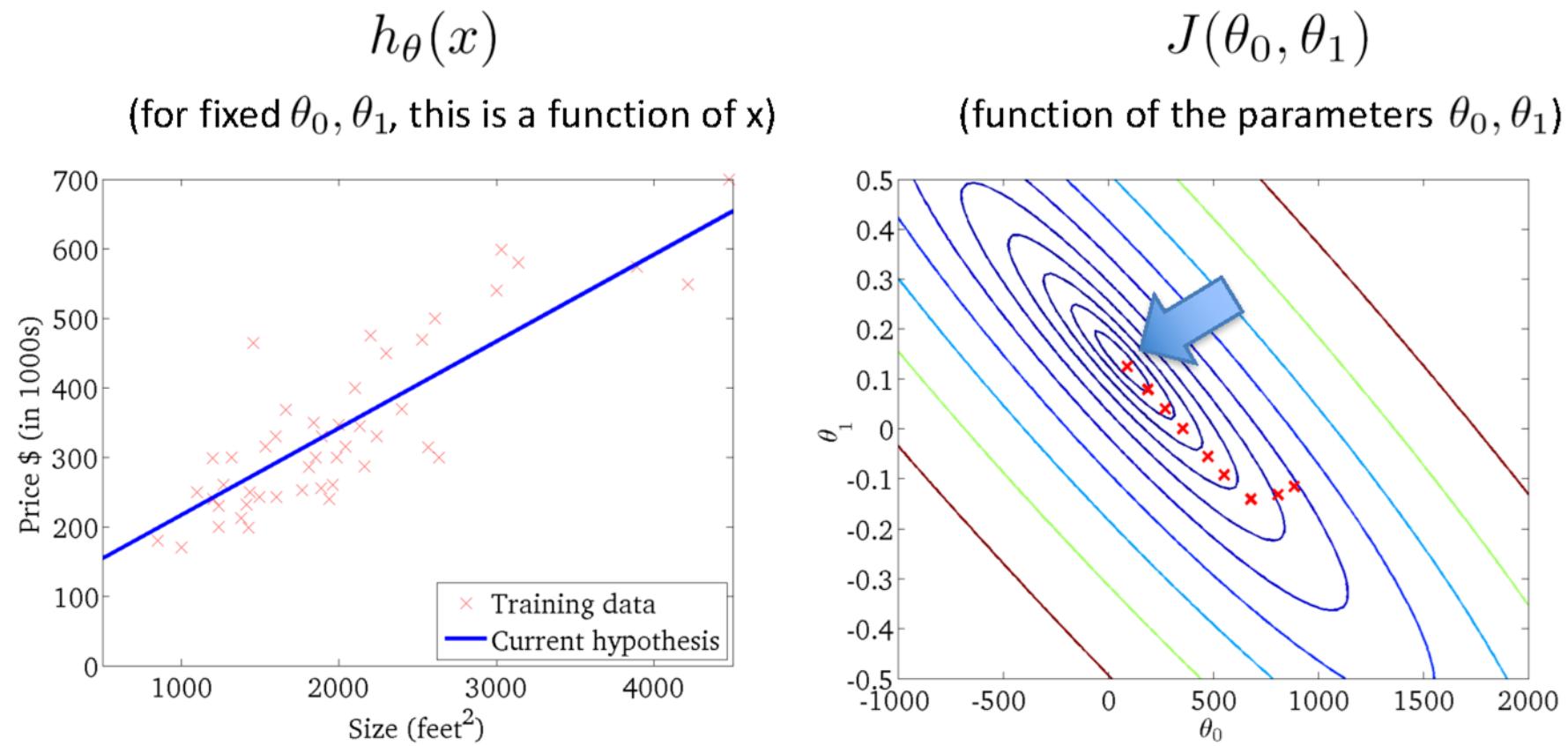
Gradient Descent



Gradient Descent



Gradient Descent



Gradient Descent

minimize the
cost function

$$(h_{\theta}(x) - y_i)^2$$

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

function should be convex

how to find ??

Hessian Matrix

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}$$

should be
positive
definite

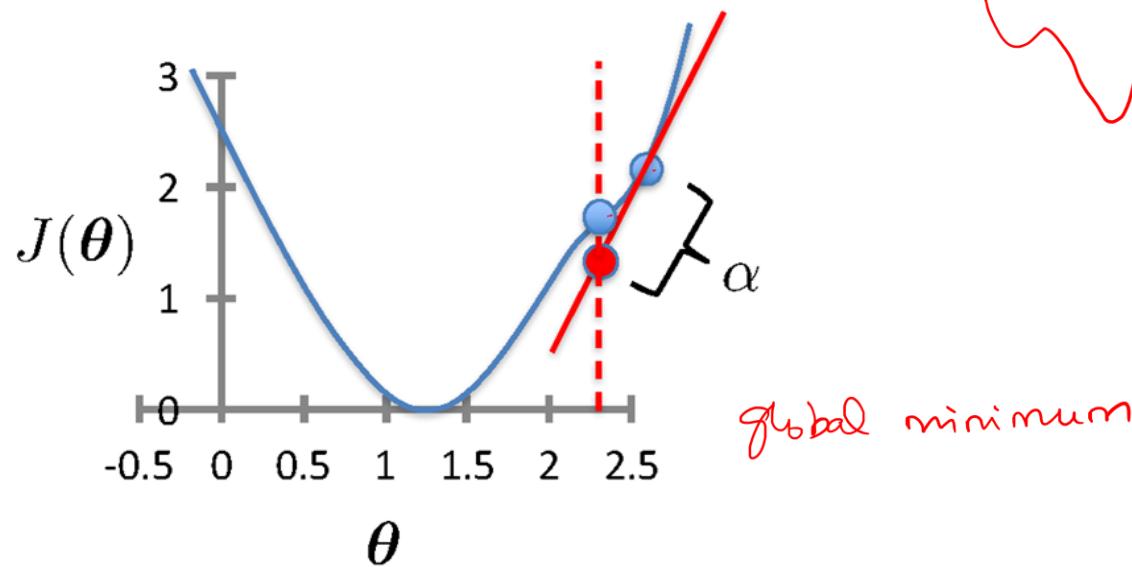
- Initialize $\theta \rightarrow$ weight
- Repeat until convergence

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

simultaneous update
for $j = 0 \dots d$

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

step size



global minimum

Gradient Descent

- Initialize θ
- 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: $\frac{\partial}{\partial \theta_j} J(\theta) = \cancel{\frac{\partial}{\partial \theta_j}} \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta} (x^{(i)}) - y^{(i)} \right)^2$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2$$

$$\cancel{\frac{\partial}{\partial \theta_j}} (y^{(i)}) = 0$$

$$= \cancel{\frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)} \times \frac{\partial}{\partial \theta_j} \left(\sum_{k=0}^d \theta_k \cancel{x_k^{(i)}} - y^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$

$$\sum_{k=0}^d \theta_k x_k = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

$$\frac{d}{dx} (x^2) = 2x (1)$$

$$(a-b)^2 = 2(a-b) \frac{d}{da} (a-b)$$

$$\frac{d}{da}$$

$$\sum_{k=0}^d \theta_k x_k = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_j x_j + \dots$$

Substitute values for $k=0, 1, 2, \dots, d$

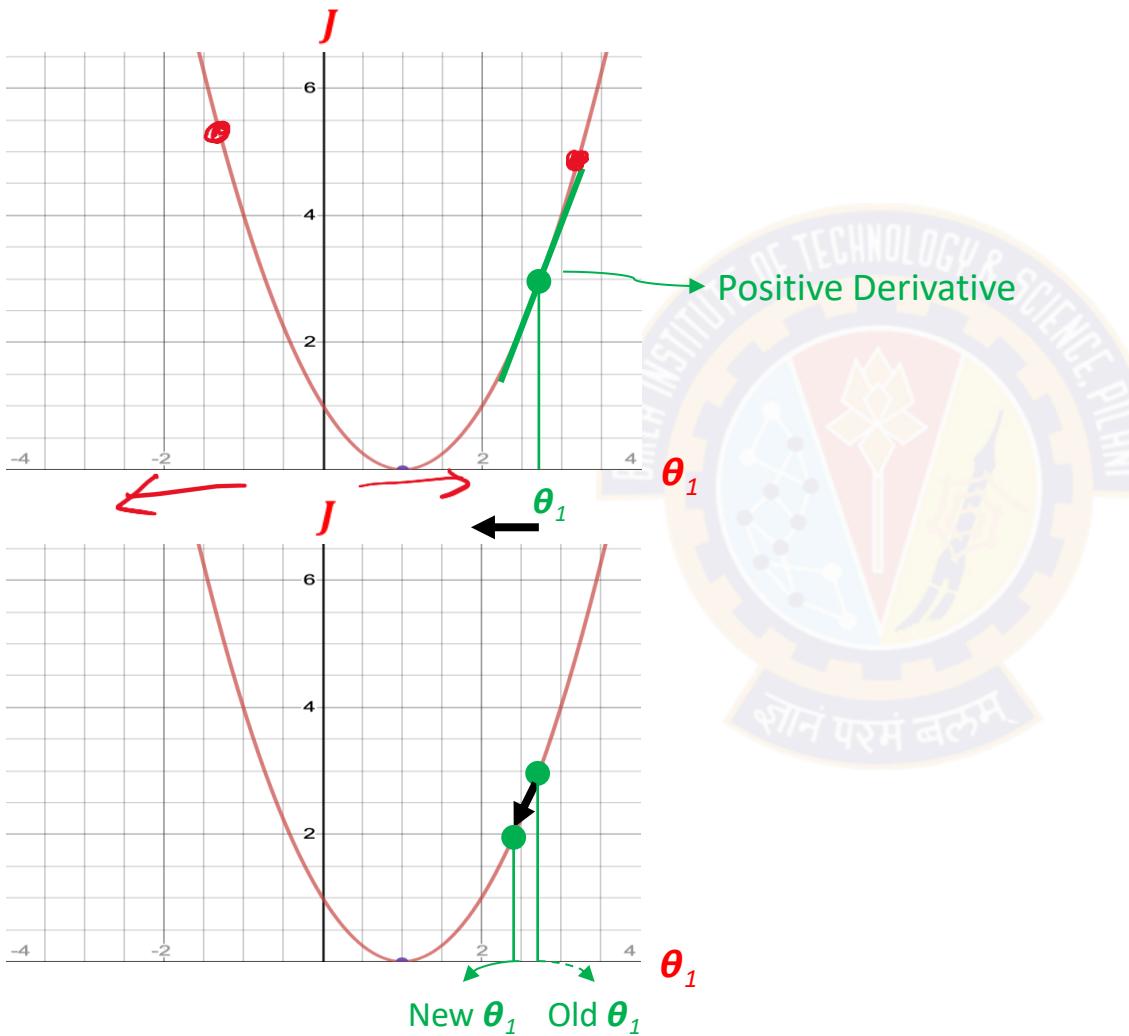
$$\frac{\partial}{\partial x} (\theta^0 w) = 2$$

$$\frac{\partial}{\partial \theta_0} (\theta_j x_j) = x_j$$

$$\frac{\partial}{\partial x_j} (\theta_j x_j) = \theta_j$$

$$\frac{\partial}{\partial \theta_j} (y^i) = 0 \text{ since } y^i \text{ is independent of } \theta$$

Guarantee of Convergence



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

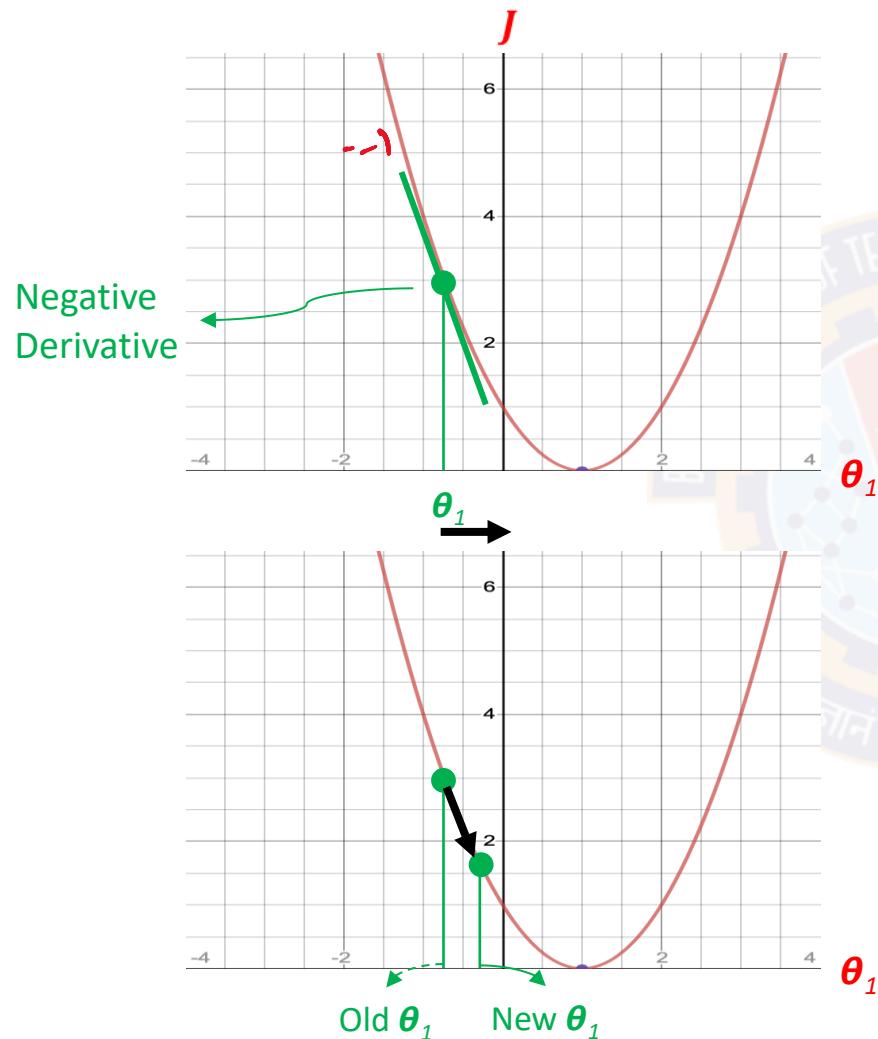
Decrease θ_1 by a certain value

$$\alpha = D \cdot D\backslash$$



Source Credit: Prof. Mohammad Hammoud

Guarantee of Convergence



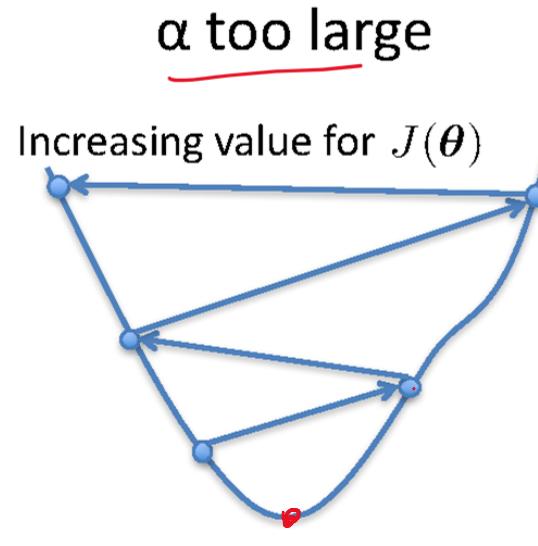
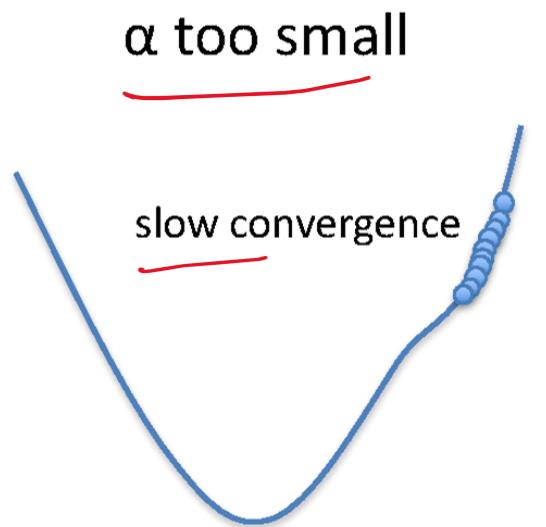
$$\begin{aligned}\theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\ &= \theta_1 - \alpha (\text{Negative Number})\end{aligned}$$

Increase θ_1 by a certain value

$$\theta_1 - \alpha (-0.5) = \theta_1 + 0.5\alpha$$

Source Credit: Prof. Mohammad Hammoud

Choosing Learning Rate



- 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 α

Gradient Descent algorithm : Effect of Feature Scaling

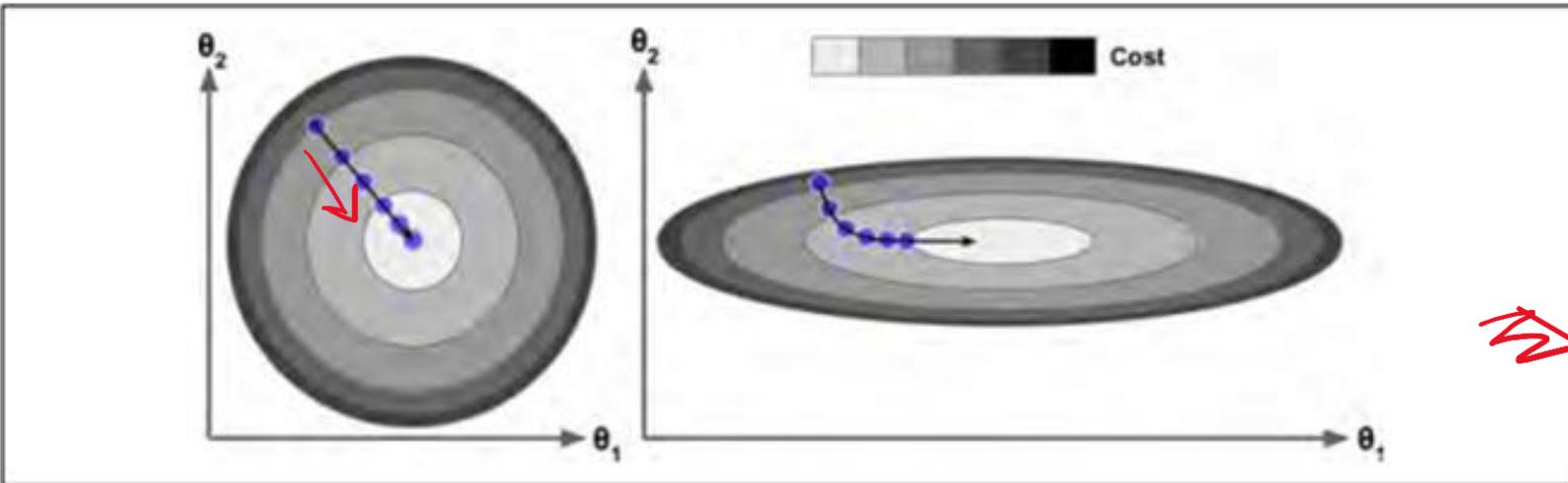


Figure 4-7. Gradient Descent with and without feature scaling

| Aspect | With Scaling | Without Scaling ✓ |
|-----------------------|---------------------|---------------------------|
| Contours shape | Circular ✓ | Elongated |
| Convergence speed | Fast | Slow |
| Gradient descent path | Straight to minimum | Zigzagging, inefficient ✓ |

Why this happens?

Gradient descent uses the slope (partial derivatives) to update each parameter θ_j . When features have very different scales:

- One direction (e.g., θ_1) might change much faster than another (e.g., θ_2).
- This causes inefficient updates, as the algorithm overcorrects in one direction and undercorrects in another.
- The result is a zigzagging path instead of a direct descent.

Solution: Feature Scaling

Feature Scaling ensures:

- All features are on a similar scale, typically by:
 - Standardization (zero mean, unit variance)
 - Min-max normalization (rescale to [0, 1])
- This transforms the cost function contours into circular shapes.
- Gradient descent moves more directly toward the optimum.

Gradient Descent algorithm : Effect of Learning Rate

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

$\eta \rightarrow$ learning rate

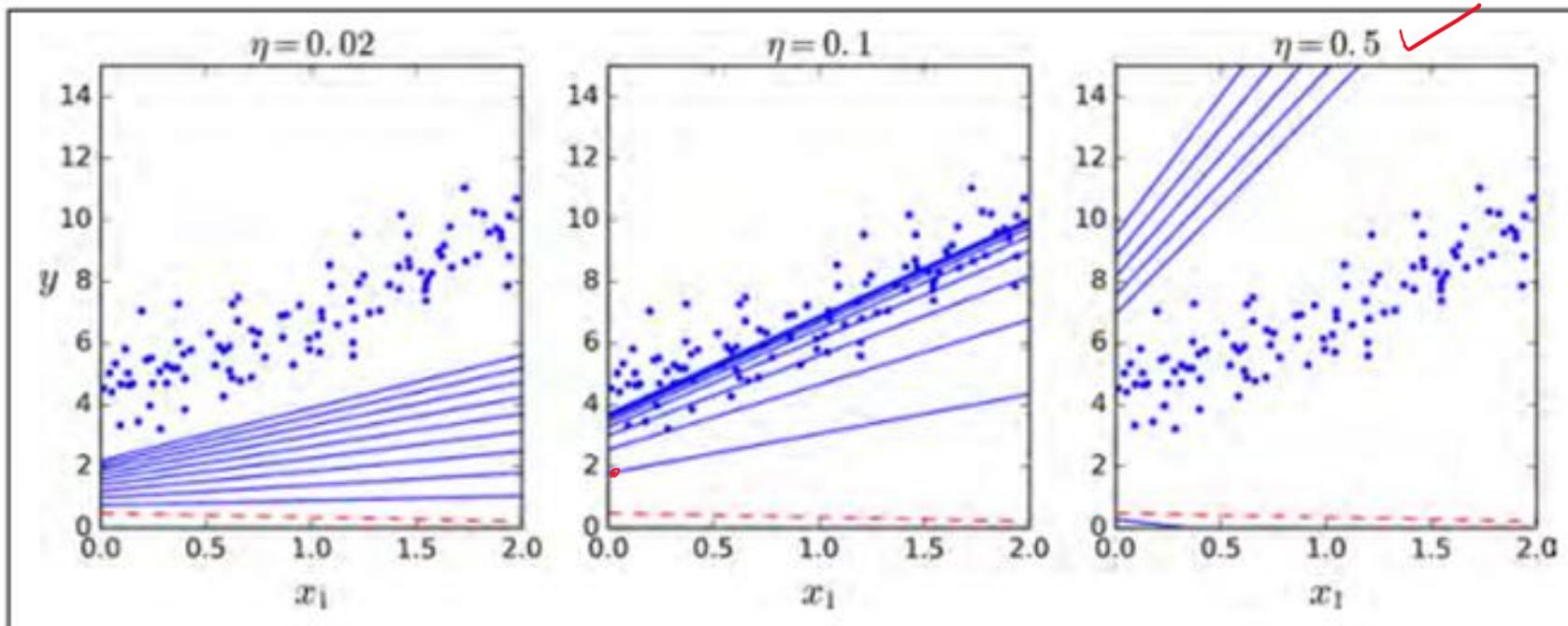


Figure 4-8. Gradient Descent with various learning rates

The learning rate (η) controls how big the steps are in gradient descent.

| Learning Rate | Behavior | Result |
|---------------|--------------------------|-----------------------------------|
| Too Small ✓ | Slow, steady updates | <u>Very slow convergence</u> |
| Just Right | Efficient updates ✓ | Fast and stable convergence |
| Too Large ✓ | Overshooting, divergence | <u>Unstable or no convergence</u> |

Gradient Descent for Linear Regression

- Initialize θ *weights*
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(\underline{h_\theta(x^{(i)})} - y^{(i)} \right) \underline{x_j^{(i)}} \quad \begin{matrix} \text{simultaneous} \\ \text{update} \\ \text{for } j = 0 \dots d \end{matrix}$$

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

$$L_2 \text{ norm: } \|\mathbf{v}\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$$

$$\mathbf{v}_i = (v_1, v_2, v_3, \dots, v_n)$$

$\|\mathbf{v}\|_2$

$\|\mathbf{v}\|_2$

L_2

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 $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

2. Cost function & derivative

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

$$\underbrace{\theta_0 + \theta_1 x_1 + \theta_2 x_2}_{y}$$

$$\theta_j$$

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

$$\partial \theta_0 = 1$$

$$\theta'_0 = \theta_0 - \frac{1}{n} * \text{learning rate} * (\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2 - y))$$

$$\theta'_1 = \theta_1 - \frac{1}{n} * \text{learning rate} * (\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2 - y) * x_1)$$

$$\theta'_2 = \theta_2 - \frac{1}{n} * \text{learning rate} * (\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2 - y) * x_2)$$

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).

$$y = mx + c$$

$$n = 3$$

$$\theta_1 = -0.03$$

$$\theta_2 = -0.03$$

$$h_{\theta}(x) = 5 - 0.03x_1 - 0.03x_2$$

$$\theta_0$$

| 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 |

actual ~~estimated~~ y_i

- Identification of the equations: RRCHD = 5 - 0.03 * BMI - 0.03 * Diastolic Pressure
- Cost function & derivative
 - $\theta'_0 = \theta_0 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}))$
 $= 5 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}))$
 - $\theta'_1 = \theta_1 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}) * \text{BMI})$
 $= -0.03 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}) * \text{BMI})$
 - $\theta'_2 = \theta_2 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}) * \text{Diastolic Pressure})$
 $= -0.03 - 1/3 * 0.02 * (\text{sum } (5 - 0.03 * \text{BMI} - 0.03 * \text{Diastolic Pressure} - \text{RRCHD}) * \text{Diastolic Pressure})$
- Apply the equations : Answer at the end of first iteration:
 $\theta_0 = 5.0016$, $\theta_1 = 0.0476$, $\theta_2 = 0.179$

given y values

$$\theta_0' = \theta_0 - \frac{\alpha}{n} \left[\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2) - \hat{y}^i \right]$$

$$\theta_1' = \theta_1 - \frac{\alpha}{n} \left[\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2) - \hat{y}^i \right] \cdot x_1$$

$$\theta_2' = \theta_2 - \frac{\alpha}{n} \left[\sum (\theta_0 + \theta_1 x_1 + \theta_2 x_2) - \hat{y}^i \right] \cdot x_2$$

First iteration: For patient 1:

$$\hat{y} = 5 - 0.03(35) - 0.03(80) = 1.55$$

$$\text{Actual } y = 1.81, \hat{y} = 1.55$$

$$\text{Error}_1 = (\hat{y} - y) = 1.55 - 1.81 = -0.26$$

For patient 2 :

$$\hat{y} = 5 - 0.03(25) - 0.03(80) = 1.85$$

$$\text{Error}_2 = 1.85 - 1.22 = 0.63$$

For patient 3 : $\hat{y} = 5 - 0.03(30) - 0.03(100) = 1.1$

$$\text{Error}_3 = 1.1 - 1.71 = -0.61$$

Update the parameters :

$$\theta_0 = 5 - \frac{0.02}{3} \sum_{i=1}^3 \text{Errors} = 5 - \frac{0.02}{3} [-0.26 + 0.63 - 0.61] \\ < 5.0016$$

Update θ_1 [BMI] :

$$\theta_1 = -0.03 - \frac{0.02}{3} \sum \text{error} \cdot \text{BMI}$$

$$\begin{aligned}\underline{\text{new}} \theta_1 &= -0.03 - \frac{0.02}{3} [(-0.26)(35) + (0.63)(25) - 0.61(30)] \\ &= -0.03 + 0.07766 = 0.04766 \quad \checkmark\end{aligned}$$

Update θ_2 (Diastolic Pressure)

$$\begin{aligned}\theta_2 &= -0.03 - \frac{0.02}{3} [\sum \text{Error. diastolic pressure}] \\ &= -0.03 + 0.20934 = 0.17934\end{aligned}$$

$$\begin{array}{l}\varepsilon = 0.01 \\ \varepsilon = 0.05\end{array}$$

Convergence criteria:

$$\|\theta_{\text{new}} - \theta_{\text{old}}\|_2 < \varepsilon$$

Closed Form Solution Vs. Gradient Descent

Gradient Descent

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

Closed Form Solution

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



Types of Gradient Descent Algorithms



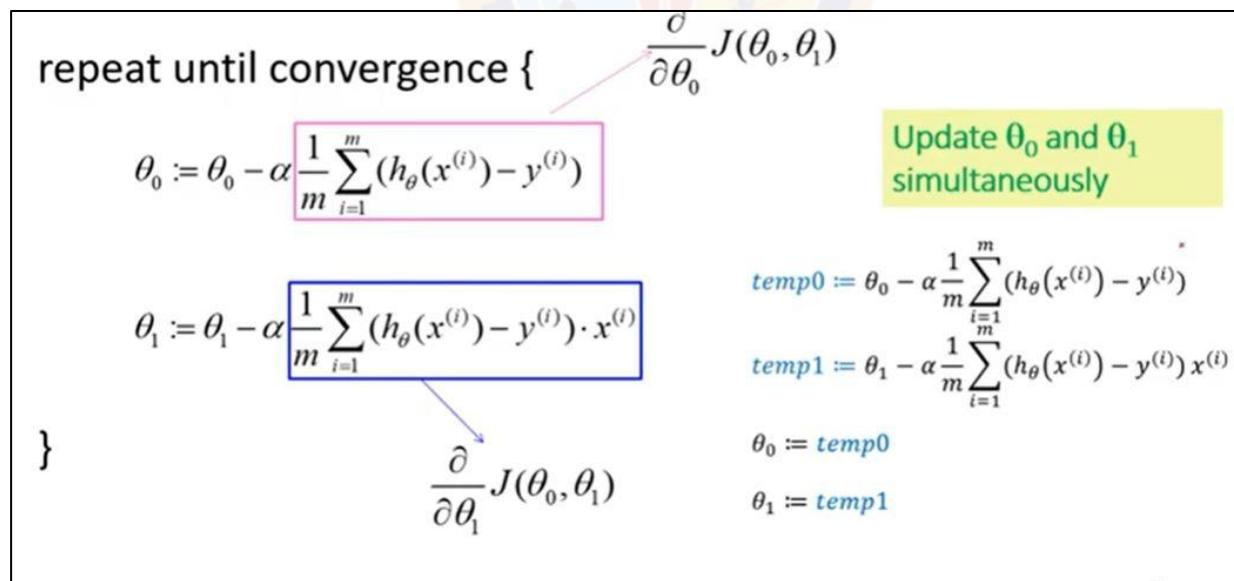
Gradient Descent: Variants

150

n = 10 groups

1 group = 15 data

- Batch gradient descent refers to calculating the derivative from all training data before calculating an update.
- Minibatch refers to calculating derivative of mini groups of training data before calculating an update.
- Stochastic gradient descent refers to calculating the derivative from each training data instance and calculating the update immediately



Gradient Descent: Variants

- **Batch** gradient descent refers to calculating the derivative from all training data before calculating an update.

Initialize the Parameters $(\theta_0^1, \theta_1^1, \dots)$

K=1

Repeat until Convergence {

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

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

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

.....

K=k+1

}

return $(\theta_0^1, \theta_1^1, \dots)$

Hn

Gradient Descent: Variants

- **Minibatch** refers to calculating derivative of mini groups of training data before calculating an update.

Divide the training instances into “N” batches each of size “m”

Initialize the Parameters $(\theta_0^1, \theta_1^1, \dots)$

K=1

Repeat until Convergence {

Repeat for every batch in 1 : N , each with ‘m’ instances {

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

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

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

.....

K=k+1

}

}

return $(\theta_0^1, \theta_1^1, \dots)$

Gradient Descent: Variants

- **Stochastic gradient** descent refers to calculating the derivative from each training data instance and calculating the update immediately

Randomly shuffle training instances

Initialize the Parameters $(\theta_0^1, \theta_1^1, \dots)$

K=1

Repeat until Convergence {

Sample with replacement, only one random training instance “ i ” at a time

$$\theta_0^{k+1} = \theta_0^k - \alpha(h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1^{k+1} = \theta_1^k - \alpha(h_\theta(x^{(i)}) - y^{(i)}) * x_1^{(i)}$$

$$\theta_2^{k+1} = \theta_2^k - \alpha(h_\theta(x^{(i)}) - y^{(i)}) * x_2^{(i)}$$

.....

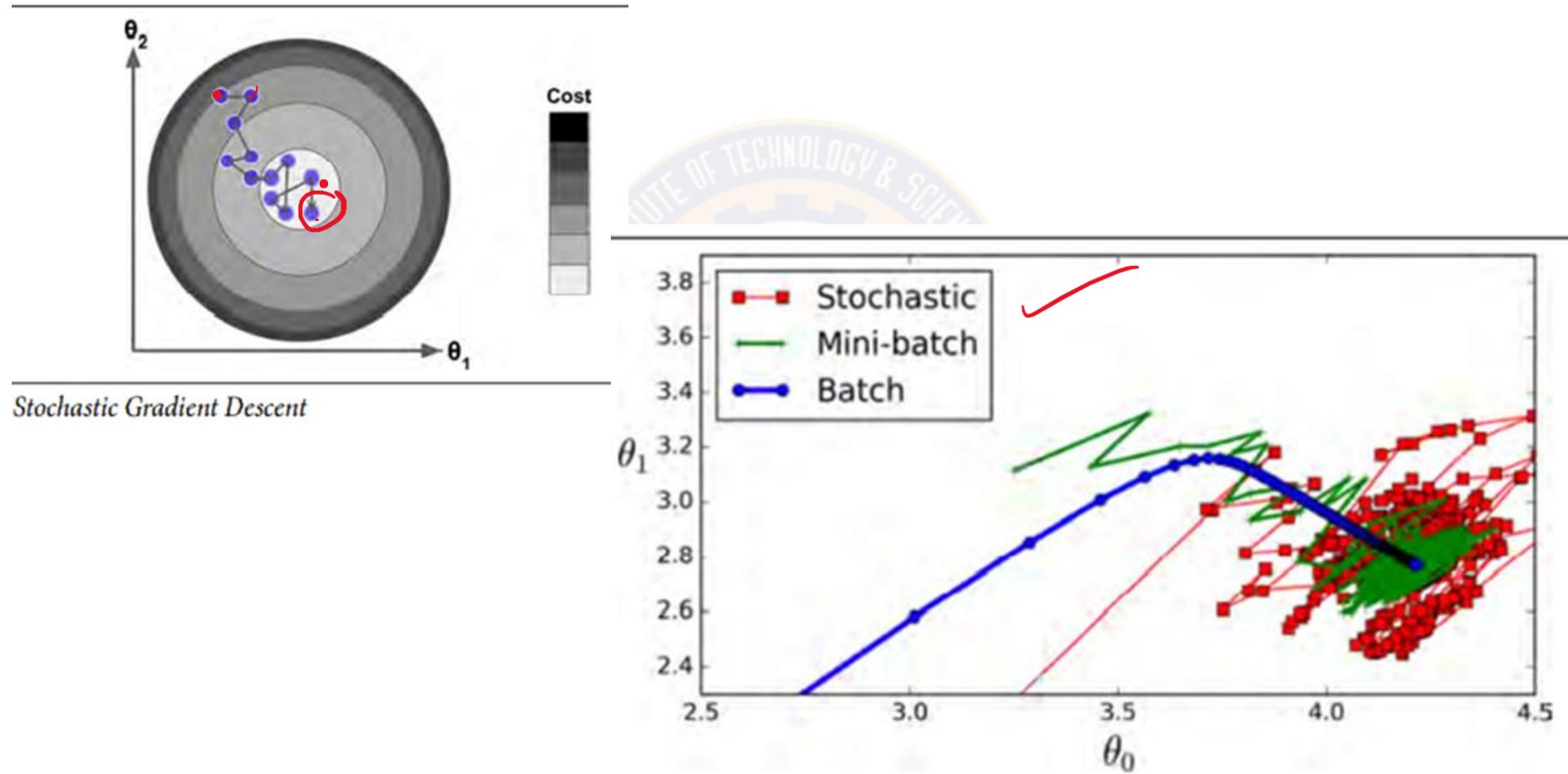
K=k+1

}

}

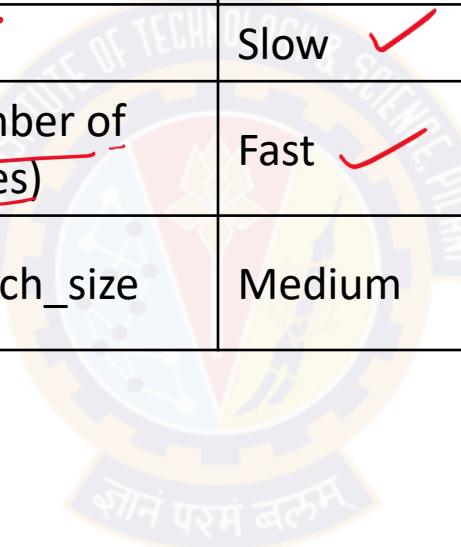
return $(\theta_0^1, \theta_1^1, \dots)$

Gradient Descent: Variants - Effect on Convergence



1. Gradient Descent paths in parameter space

| Type | Data Used per Update | Updates per Epoch | Speed | Stability | Common Use Case |
|------------|----------------------|-----------------------|--------|-----------|-------------------------|
| Batch | All data | 1 ✓ | Slow ✓ | High | Small datasets |
| Stochastic | 1 sample | n (number of samples) | Fast ✓ | Low | Online learning ✓ |
| Mini-Batch | Small batches | n / batch_size | Medium | Medium | Deep learning (default) |


 ✓
 chatgpt → gives one answer
 Yes | No
 Second time Yes ✗ Not repeated

Evaluation Metrics



Evaluation of Linear Regression Model

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

$x \uparrow$ $y \uparrow$
 $x \uparrow$ $y \downarrow$
 $x \uparrow$ $y \uparrow$

\rightarrow Strength of relationship (x, y)

$R^2 = 0.9$

Model 1

$$\text{CarPrice} = 8.5 + 0.5 \text{Mileage} - 1.5 \text{Mileage}^2$$

base rate

Mean Square Error

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

Mean Absolute Error

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}|$$

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2}$$

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{Total}}}$$

| Unseen Data | |
|----------------------|----------------------|
| Mileage (in kmpl) | Car Price (in cr) |
| 7.5 | 9.25 |
| 10 | 6.5 |
| | |

Model 2

$$\text{CarPrice} = -5.5 + 1.5 \text{Mileage}$$

Evaluation of Linear Regression Model

R-squared

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

variation in 'y' that is explained by a regression model

$$\underline{\text{explained variation}} = \hat{y} - \bar{y}$$

variation in 'y' that is not captured/explained by a regression model

$$\underline{\text{unexplained variation}} = y - \hat{y}$$

$$\text{total variation} = (y - \hat{y}) + (\hat{y} - \bar{y}) = (y - \bar{y})$$

$\hat{y} \rightarrow \text{predicted } y\text{-value}$

| Car Price (in cr) |
|----------------------|
| 10.48 |
| 1.75 |
| 6.95 |
| 2.51 |
| Mean Y |

$$R^2 = 1 - \frac{SS_{residual}}{SS_{Total}}$$

Model 1

$$\text{CarPrice} = 8.5 + 0.5 \underline{\text{Mileage}} - 1.5 \underline{\text{Mileage}^2}$$

- **variation that is explained by a regression model**
- **measures the goodness of fit of a regression model**

$$SS_{explained} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SS_{residual} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SS_{Total} = \sum_{i=1}^n (y_i - \bar{y})^2$$

where :

$SS_{explained}$ = explained variation sum of squares

$SS_{residual}$ = unexplained variation sum of squares

SS_{Total} = total variation sum of squares

$$(a_0x_0 + a_1x_1 + a_2x_2^2 + a_3x_3^3 + \dots + a_nx_n^n)$$

Linear Basis Models

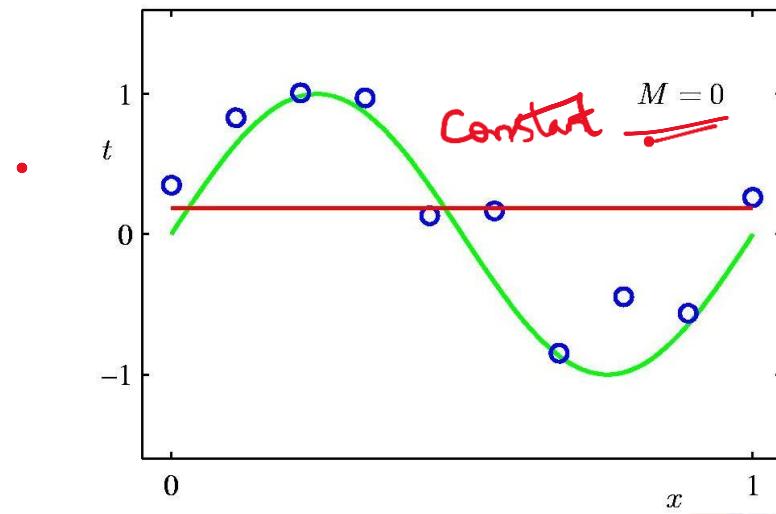


$$x, \quad y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

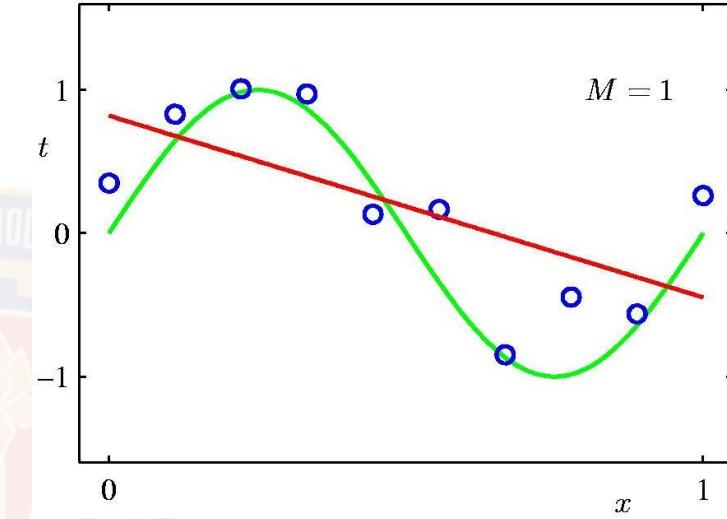
What if output is a non-linear function of input vector?

Polynomial Regression

degree of the polynomial



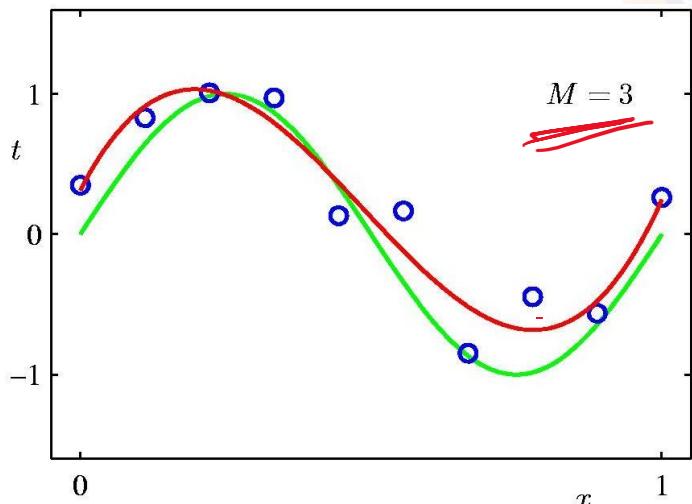
Constant
 $M = 0$



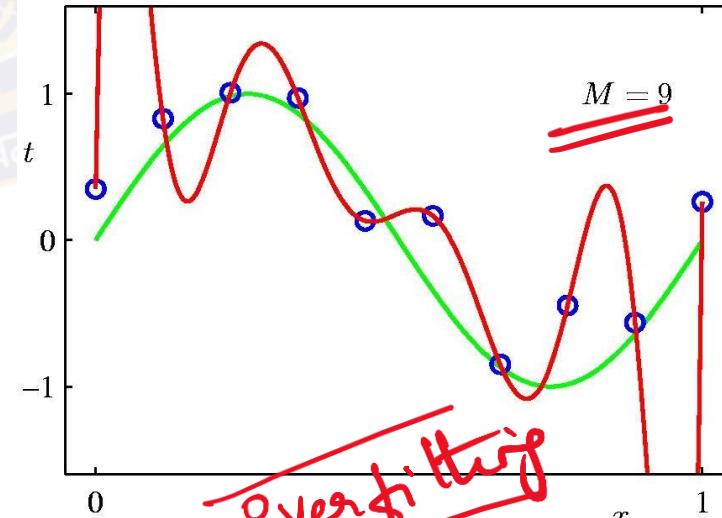
$M = 1$

$$y = ax^M + b$$

when $M = 1$



$M = 3$



$M = 9$

overfitting

new instance
??

Linear Basis Function Models

- The inputs \mathbf{X} for linear regression can be:
 - Original quantitative inputs
 - Transformation of quantitative inputs
 - e.g. log, exp, square root, square, etc.
 - Polynomial transformation
 - example: $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot x^3$
 - Basis expansions
 - Dummy coding of categorical inputs
 - Interactions between variables
 - example: $x_3 = x_1 \cdot x_2$

This allows use of linear regression techniques
to fit non-linear datasets.

| Input X | Output Y |
|----------|----------|
| exp(2) | |
| exp(4) | |
| exp(6.3) | |
| exp(9.2) | |

$y = ax + bx^2$

| X No.of.Years of Experience (in Years) | X^2 | Y Salary Of the Employee (in Lakhs) |
|---|-------|--|
| 1 | 1 | 2 |
| 2 | 4 | 3 |
| 3 | 9 | 4 |
| 4 | 16 | 5 |
| 5 | 25 | 6 |

| X1 = Graduate | X2 = PostGraduate | X3 = Others | Y Salary Of the Employee |
|------------------|----------------------|----------------|-----------------------------------|
| 0 | 0 | 1 | 2 |
| 1 | 0 | 0 | 3 |
| 0 | 0 | 1 | 4 |
| 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 6 |

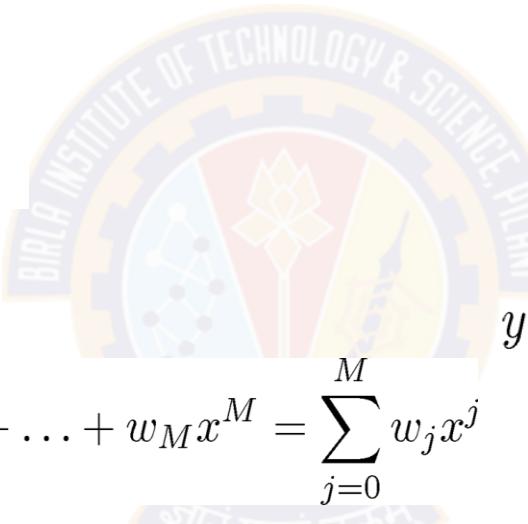
Linear Basis Function Models

Example: an M-th order polynomial function of one dimensional feature x :

$$y(x, \mathbf{w}) = w_0 + \sum_{j=1}^M w_j x^j$$

where x^j = j-th power of x

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$



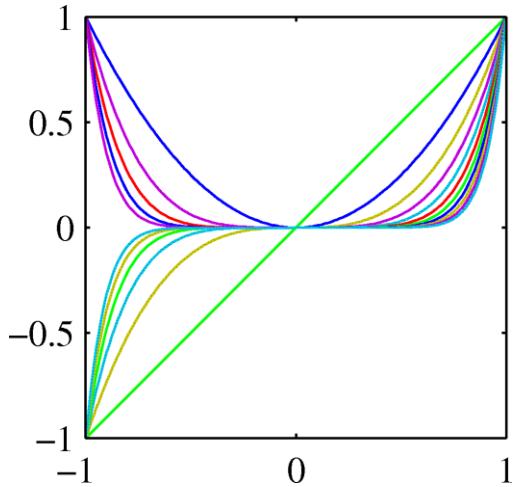
| X No.of Years of Experience (in Years) | X^2 | Y Salary Of the Employee (in Lakhs) |
|---|-------|--|
| 1 | 1 | 2 |
| 2 | 4 | 3 |
| 3 | 9 | 4 |
| 4 | 16 | 5 |
| 5 | 25 | 6 |

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

- $\phi_j(x)$ are known as *basis functions*. Typically, $\phi_0(x) = 1$, so that w_0 acts as a *bias*.
- In the simplest case, we use *linear basis functions* : $\phi_d(x) = x^d$.
- They are called *linear models* because this function is linear in \mathbf{w} .

$$\phi_d(x) = x, x^2, x^3, \dots, x^d$$

Linear Basis Function Models - Examples

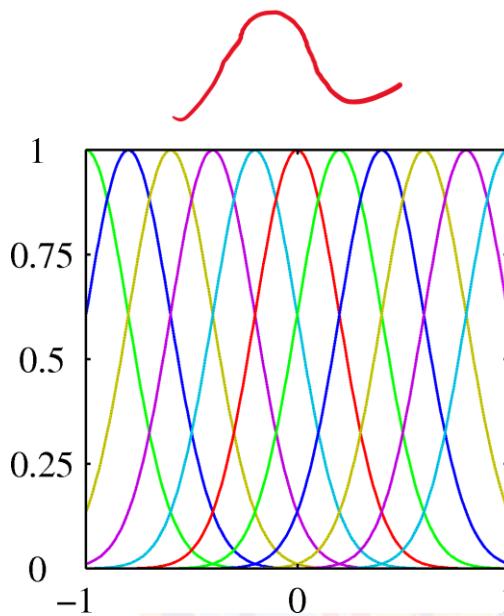


$x=1$ x^1, x^2
 $x=1.01$ x^1, x^2, x^3

Polynomial basis functions:

$$\phi_j(x) = x^j.$$

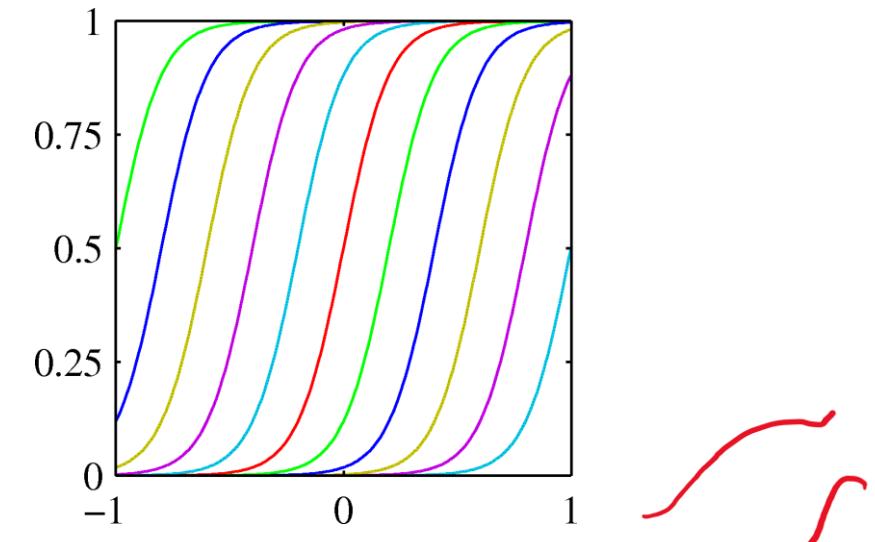
These are global; a small change in x affect all basis functions.



Gaussian basis functions:

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$

These are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (width).



Sigmoidal basis functions:

$$\text{where } \phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

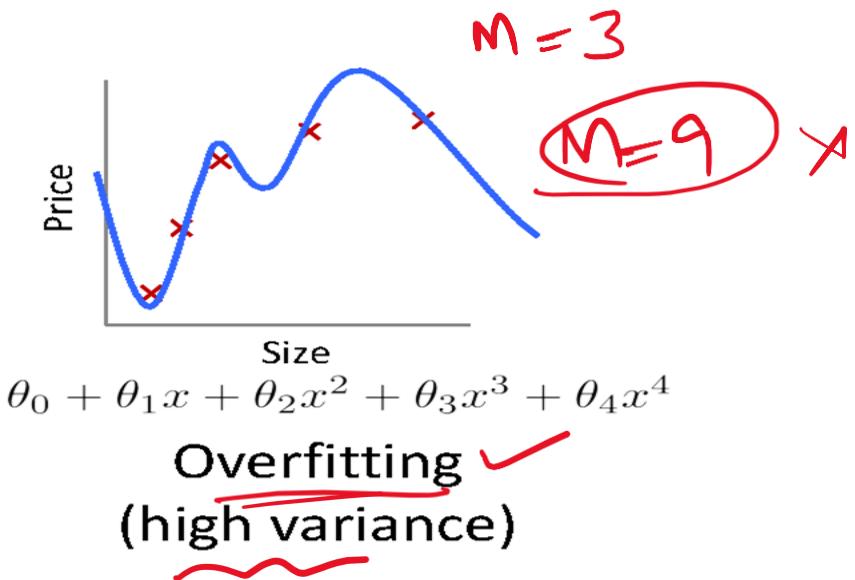
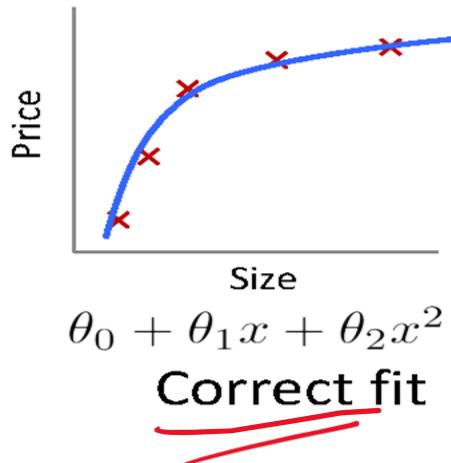
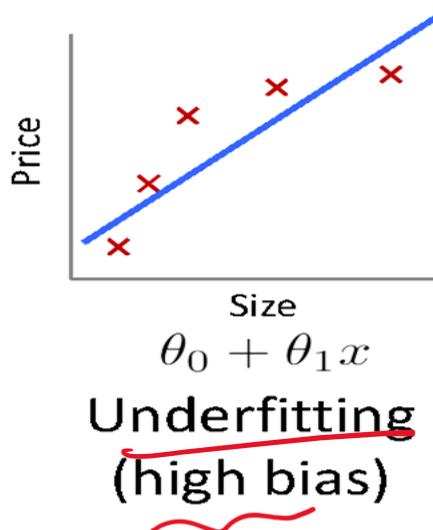
$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Also these are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (slope).

Notion of Bias - Variance



Quality of Fit

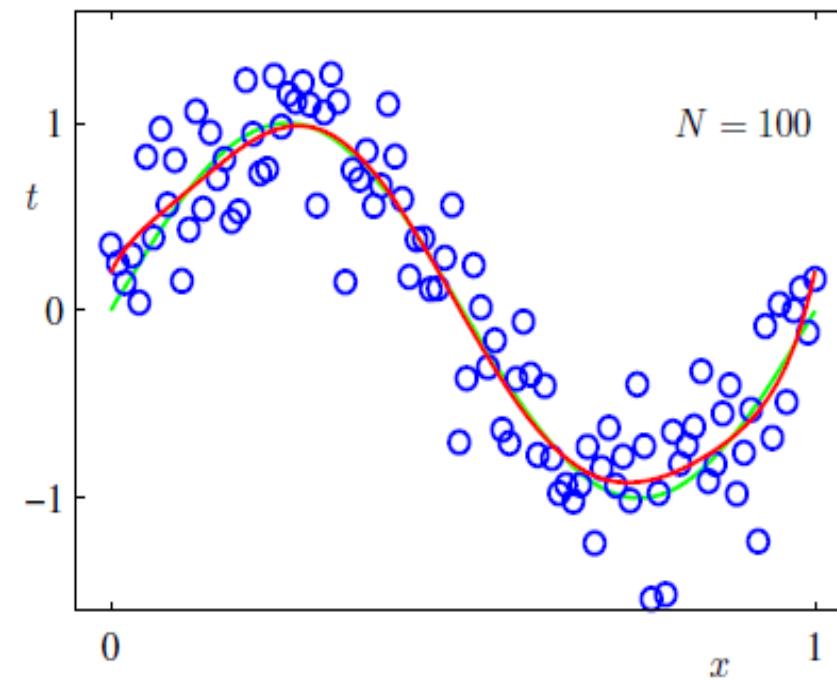
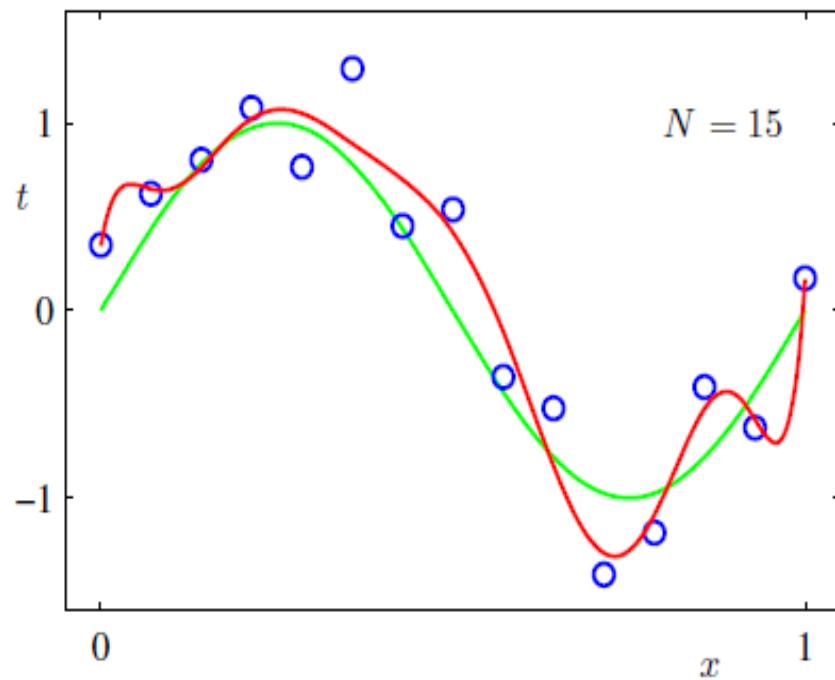


Overfitting:

- The learned hypothesis may fit the training set very well ($J(\theta) \approx 0$)
- ...but fails to generalize to new examples

Handling Overfitting – Way 1

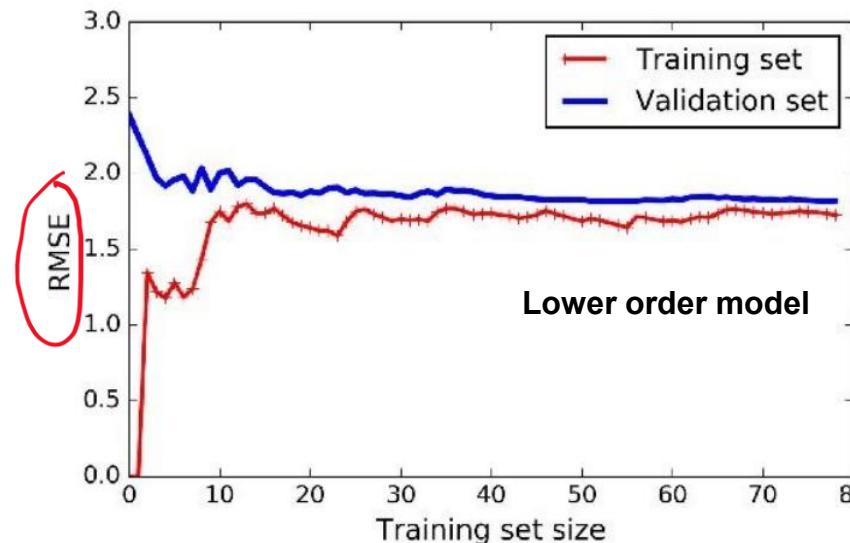
- Increase in Size of the data set reduces the over-fitting problem



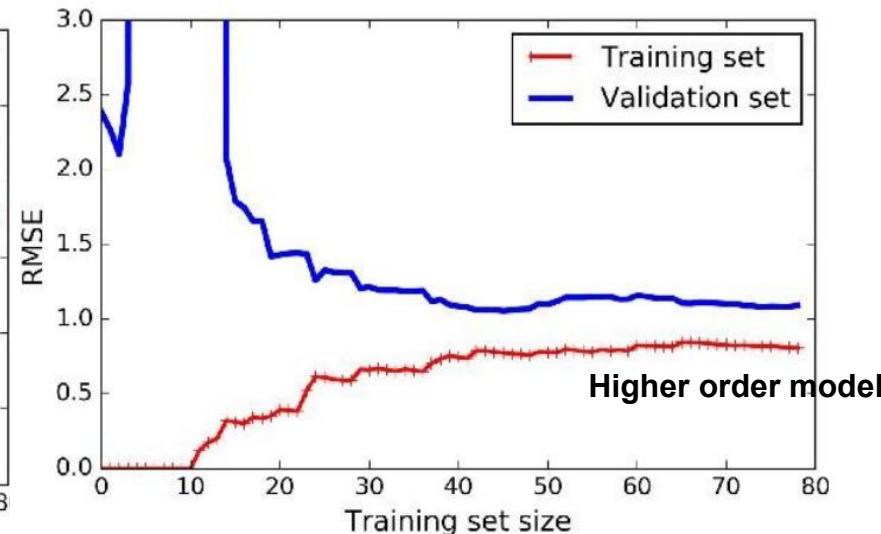
Effect of Training Size on Over fitting

Problem Type 4 : Interpretation of the Model Fit

- Size of training dataset needs to be large to prevent when higher order model is used.



This model is suffering from high bias.
It is underfitting the data



• This model has **high variance**
• It is overfitting



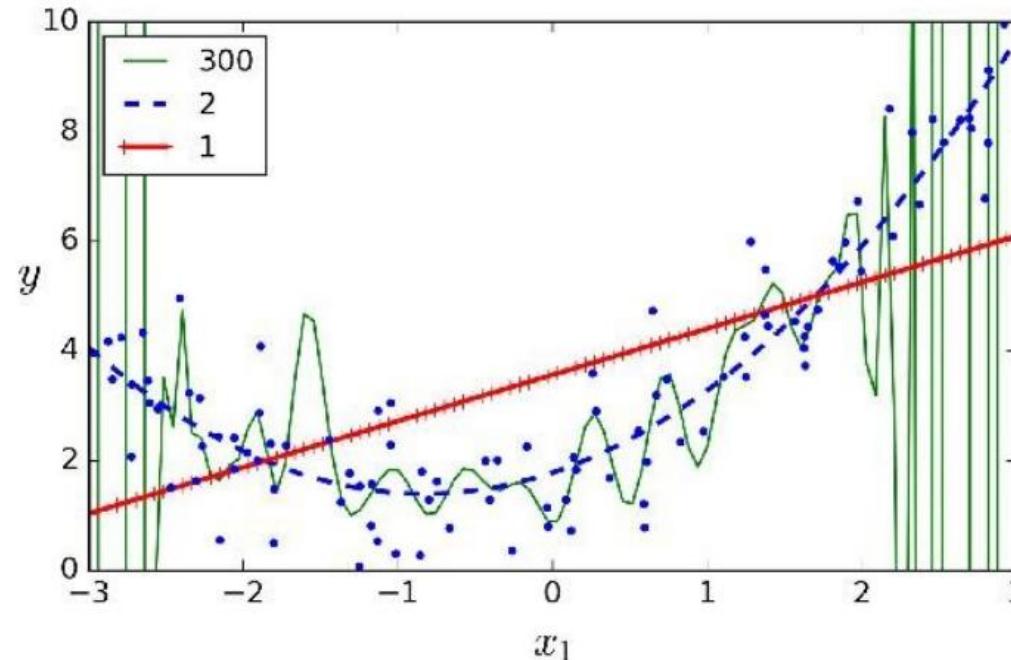
Polynomial Fitting can lead to Overfitting

Handling Overfitting – Way 2 – Reduce the complexity of the model

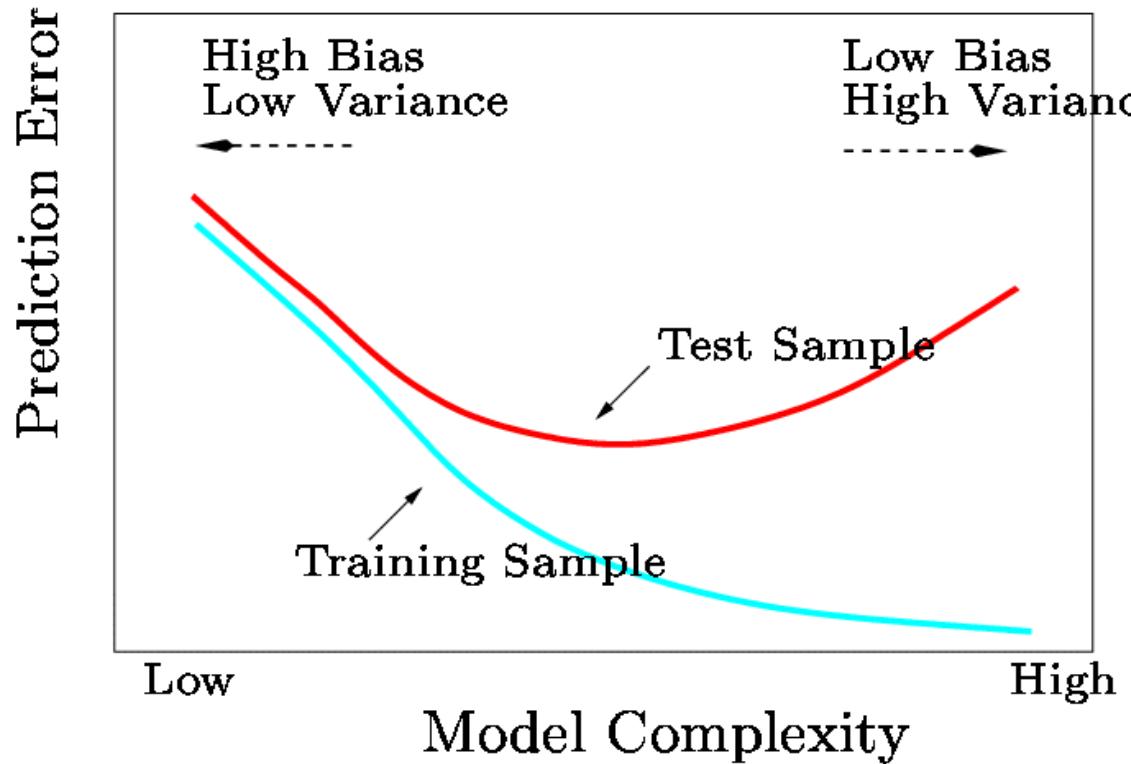
Underlying target function is quadratic

features scaling

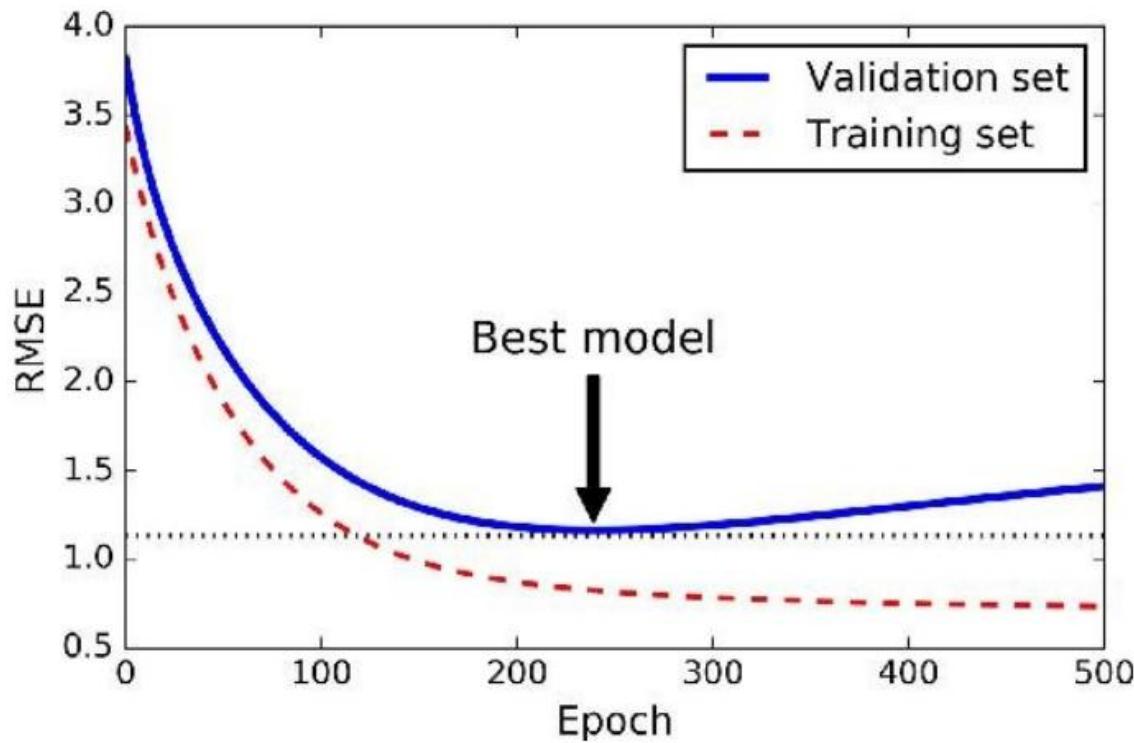
- Linear model results in under fitting with large bias
- Polynomial of order 300 results in a large variance



How to experiment on Model Complexity ?



Handling Overfitting – Way 3



Regularization

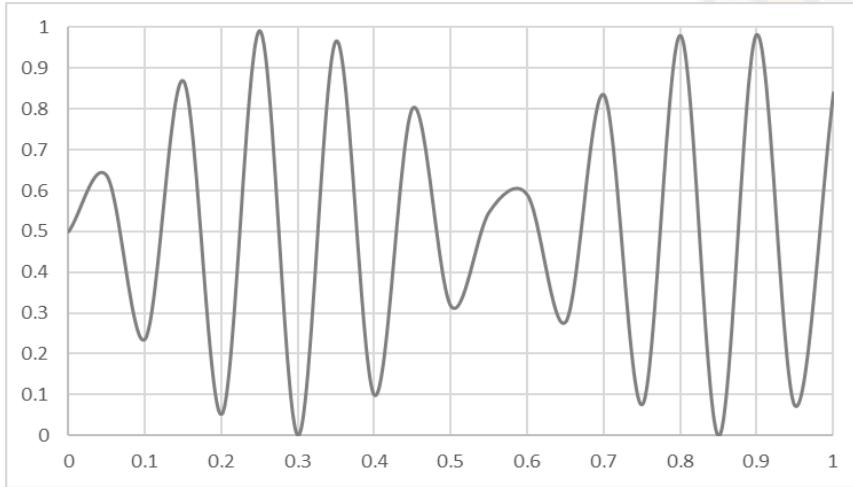
(*This notion is common for both Linear Regression – Module 3 and Logistic Regression – Module 4)



Overfitting vs Underfitting

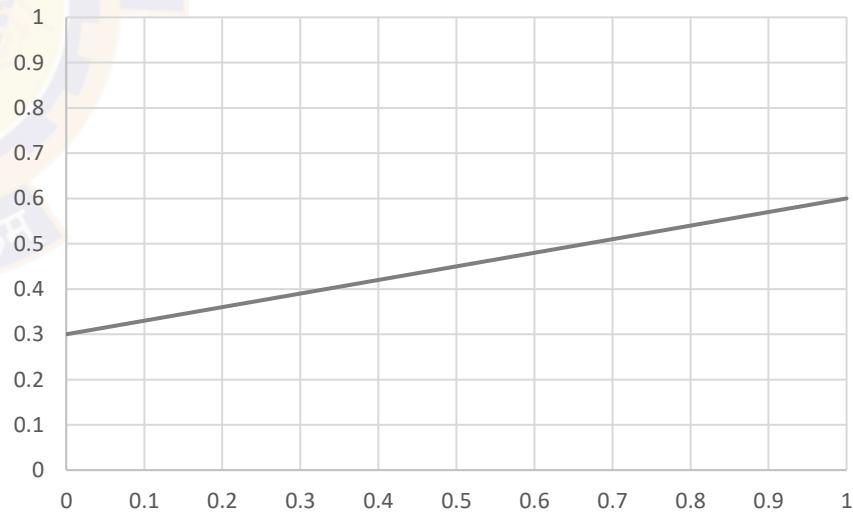
Overfitting ✓

- Fitting the data too well
 - Features are noisy / uncorrelated to concept



Underfitting ✓

- Learning too little of the true concept
 - Features don't capture concept
 - Too much bias in model



Regularization ✓

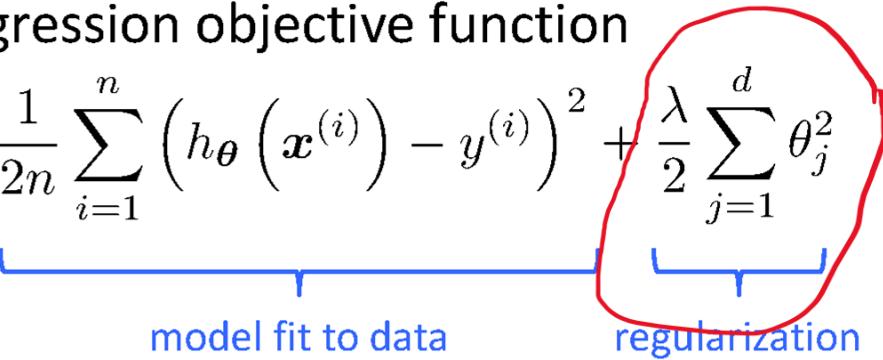
- A method for automatically controlling the complexity of the learned hypothesis
- Idea: penalize for large values of θ_j
 - Can incorporate into the cost function
 - Works well when we have a lot of features, each that contributes a bit to predicting the label
- Can also address overfitting by eliminating features (either manually or via model selection)



Regularization

- Linear regression objective function

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



- λ is the regularization parameter ($\lambda \geq 0$)
- No regularization on θ_0 !



Understanding Regularization

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

- Note that $\sum_{j=1}^d \theta_j^2 = \|\boldsymbol{\theta}_{1:d}\|_2^2 = \theta_1^2 + \theta_2^2 + \theta_3^2 + \dots + \theta_d^2$
 - This is the magnitude of the feature coefficient vector!

- We can also think of this as:

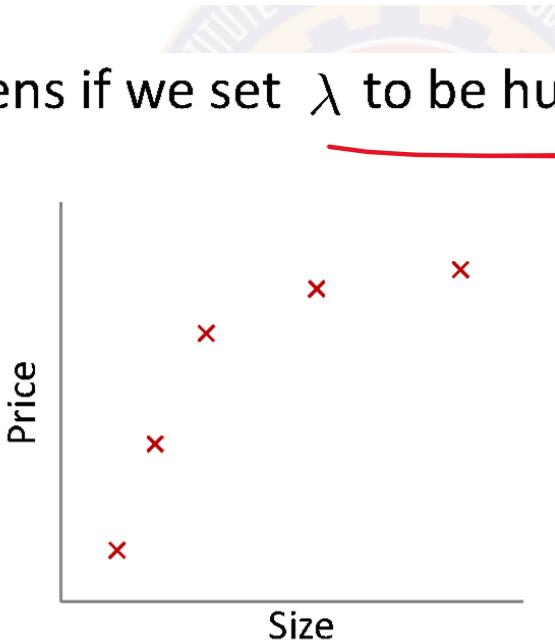
$$\sum_{j=1}^d (\theta_j - 0)^2 = \|\boldsymbol{\theta}_{1:d} - \vec{0}\|_2^2$$

- L₂ regularization pulls coefficients toward 0

Understanding Regularization

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

- What happens if we set λ to be huge (e.g., 10^{10})?



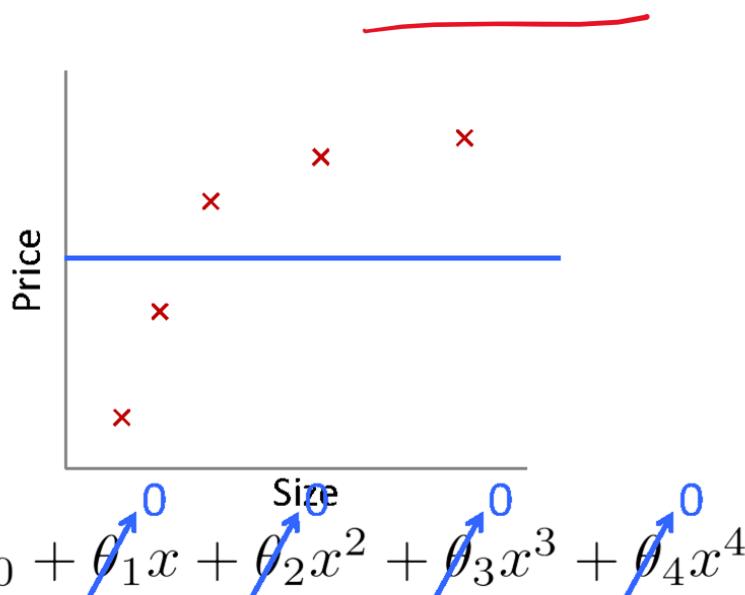
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Based on example by Andrew Ng

Understanding Regularization

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

- What happens if we set λ to be huge (e.g., 10^{10})?



Based on example by Andrew Ng

Regularization **Ridge Regression / Tikhonov regularization**

- Cost Function

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

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

- Gradient update:

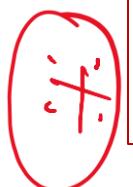
$$\frac{\partial}{\partial \theta_0} J(\boldsymbol{\theta})$$

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

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

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

regularization ✓



- We can rewrite the gradient step as:

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

Regularization

Lasso Regression (Least Absolute Shrinkage and Selection Operator Regression)

L_1 Regularisation

- Cost Function

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

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

- Gradient update:

$$\frac{\partial}{\partial \theta_0} J(\theta) \quad \theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

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

regularization

where $\text{sign}(\theta_i) = \begin{cases} -1 & \text{if } \theta_i < 0 \\ 0 & \text{if } \theta_i = 0 \\ +1 & \text{if } \theta_i > 0 \end{cases}$

Regularization Elastic Net



- Cost Function

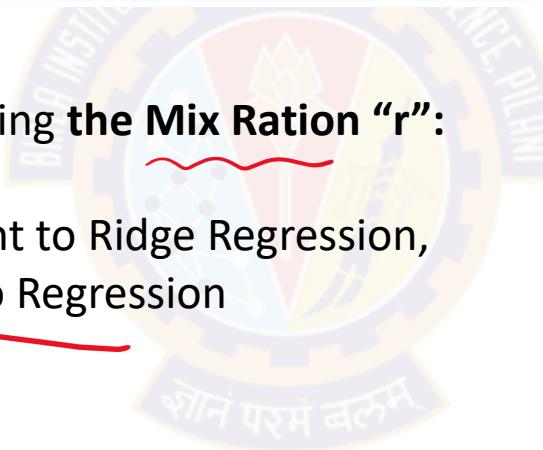
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + r \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2 + \frac{(1-r)}{2} \lambda \sum_{j=1}^d |\theta_j| .$$

Control the regularization using **the Mix Ration “r”**:

When,

$r = 0$, Elastic Net is equivalent to Ridge Regression,

$r = 1$, it is equivalent to Lasso Regression



```
from sklearn.linear_model import ElasticNet  
ElasticNet(alpha=0.1, l1_ratio=0.5)
```

How to choose the right Regularization? Common Usage & Observation

- L1 regularization has the ability to set some coefficients to 0 exactly leading to a sparse model
- L1 regularization helps in feature selection by eliminating the features that are not important
- **L1 cannot be used efficiently in gradient-based approaches since it is not-differentiable unlike L2**
- L2 will in general lead to small magnitudes of weights but not exactly 0.
- Elastic net – Prefer in Highly correlated features

X
Study hours

Y
Exam score

$$\|x\|^2$$
$$x = 0.01 \quad \|x\|^2$$
$$=$$

Numerical Exercise – For Student Practice

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 |

Apply a regularization on the same problem with regularization constant 5 and apply GD for 2 iterations & interpret the results. Try both ridge regression as well as lasso regression. Below equation is changed only for ridge regression. Students must apply appropriate GD update equation changes for this problem.

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. $w_1' = w_1 * (1 - \text{learning rate} * \text{regularization constant}) - \frac{1}{3} * \text{learning rate} * (\sum (w_0 + w_1X_1 + w_2X_2 - y) * x_1)$
 3. $w_2' = w_2 * (1 - \text{learning rate} * \text{regularization constant}) - \frac{1}{3} * \text{learning rate} * (\sum (w_0 + w_1X_1 + w_2X_2 - y) * x_2)$
3. Apply the equations

Python Lab Exercise – For Student Practice

Go to your virtual lab file under : Machine Learning LabCapsule 2 Linear _
Polynomial Regression 3 Polynomial Regression

Download ML_Lab 5 PolynomialRegression.ipynb

Try to change the degree to {2, 5, 7, 10} in the below function line :
PolynomialFeatures(degree=3)

Observe and interpret on the performance of the training data

Note: Training vs Test Data split is not coded in this implementation. Below are the suggestions to experiment further:

- Add few hundreds of data (Use synthetic data generation python libraries available : Refer here)
- Split the data into 80% train vs 20% test set
- Built the polynomial model using above four different degree on training set
- For each of the model apply in the test set
- Find the MSE for both training data and separately for test data
- Plot the four experiment results in a plot and interpret the notion of overfit vs underfit
- Suggest which degree is a best fit .



Thank you