# Machine Learning

## Linear Regression

Dr. Mehran Safayani

safayani@iut.ac.ir

safayani.iut.ac.ir

https://www.aparat.com/mehran.safayani

https://github.com/safayani/machine_learning_course

Department of Electrical and computer engineering,  Isfahan university of technology, Isfahan, Iran

# Supervised Learning

- Regression

- Classification

# example

**Notation:**

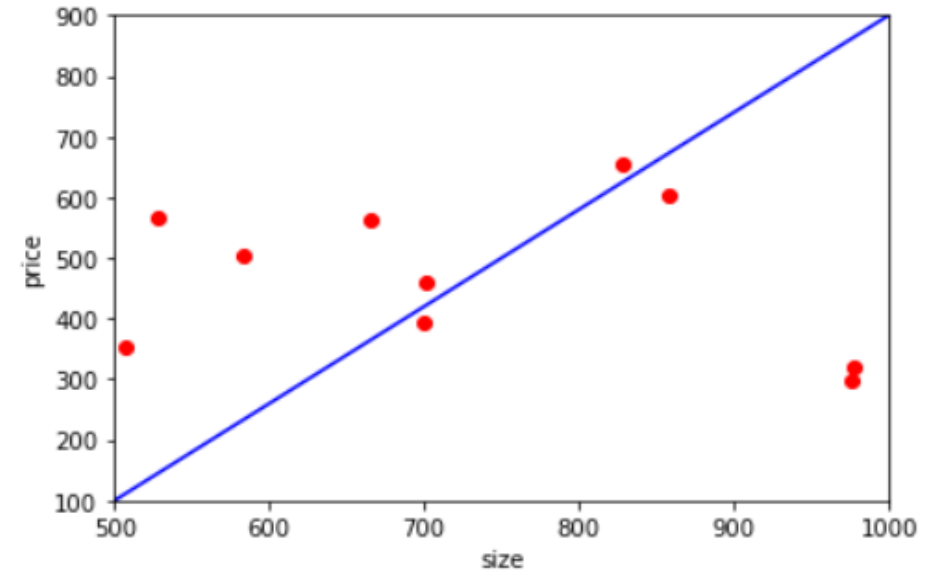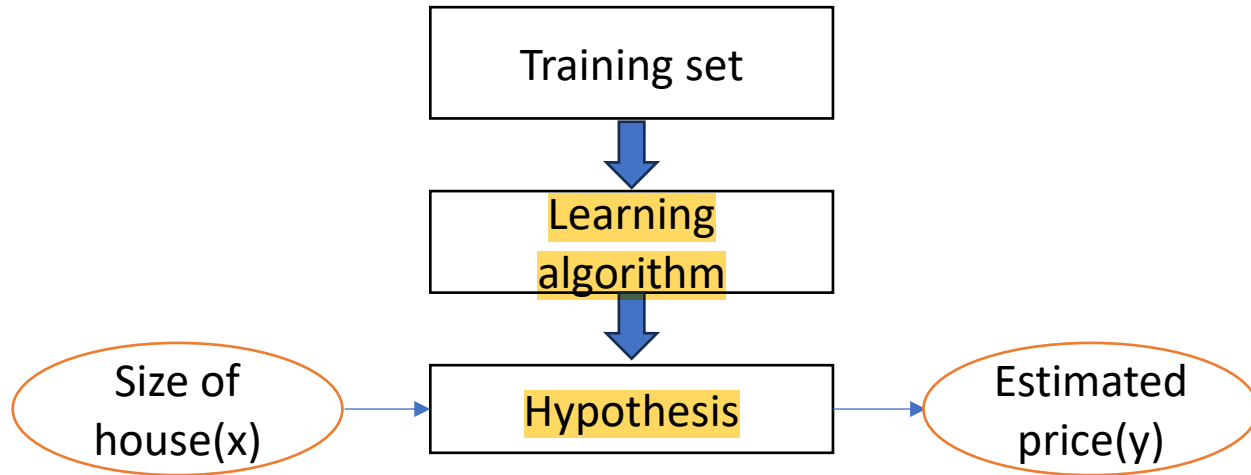m: number of training samples

x: input variable

y: output variable
Or
target variable

$(x_i, y_i)$: i th training sample

| number | Size (x variable) | Price (y variable) | |
|--------|-------------------|--------------------|--------------|
| 1 | 100 | 500 | $(x_1, y_1)$ |
| 2 | 750 | 2000 | $(x_2, y_2)$ |
| 3 | 852 | 178 | $(x_3, y_3)$ |
| | ... | ... | |
| m | 3210 | 870 | $(x_m, y_m)$ |

# example



Training set → Learning algorithm → Hypothesis

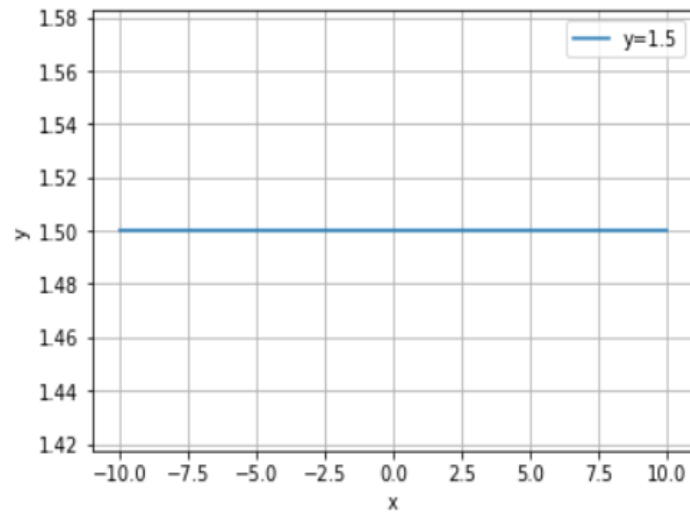Size of house(x) → Hypothesis → Estimated price(y)

$$h(x) = \theta_0 + \theta_1 x$$

$$\text{parameters} = \left\{ \theta_0, \theta_1 \right\}$$

۴

$$h(x) = \theta_0 + \theta_1 x$$

parameters=


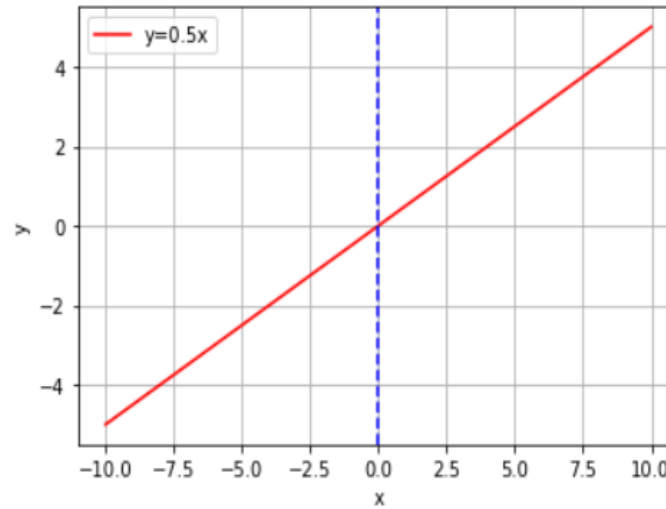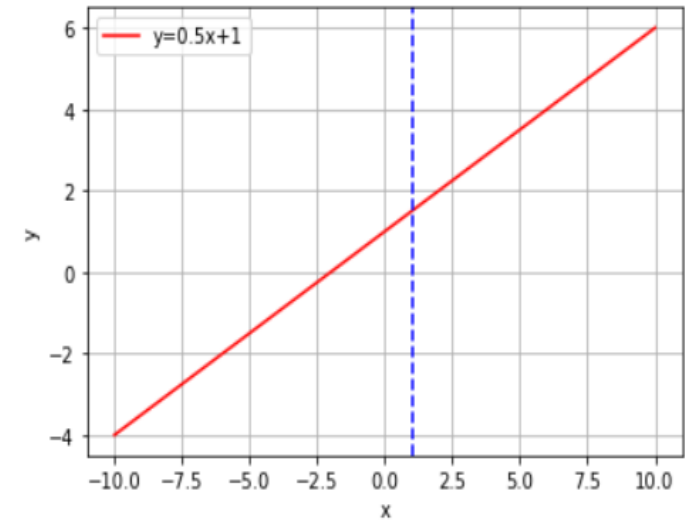
h(x) = 1.5

$\theta_0 = 1.5$
$\theta_1 = 0$

h(x) = 0.5x

$\theta_0 = 0$
$\theta_1 = 0.5$

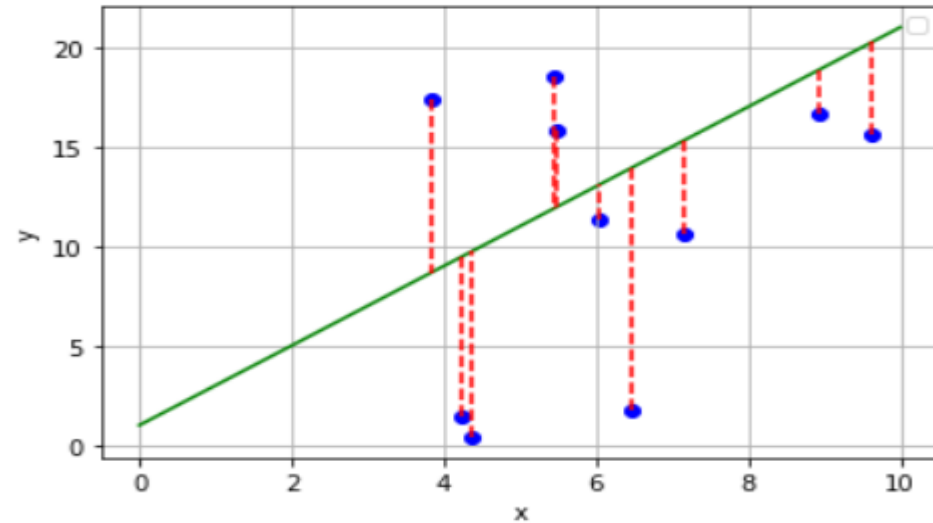h(x)=0.5x+1

$\theta_0 = 1$
$\theta_1 = 0.5$

# Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h(x_i) - y_i)^2$$

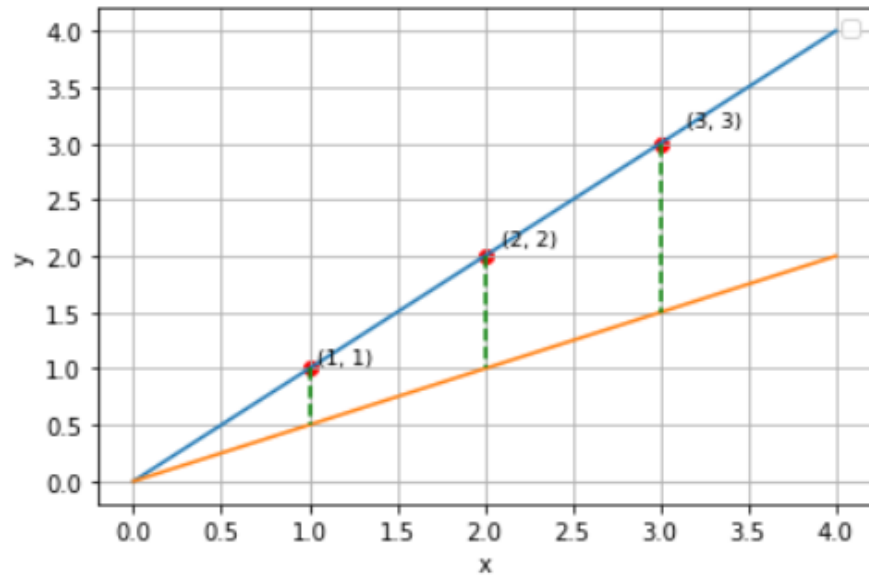Mean square error(MSE)

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

$\theta_0, \theta_1$

# example
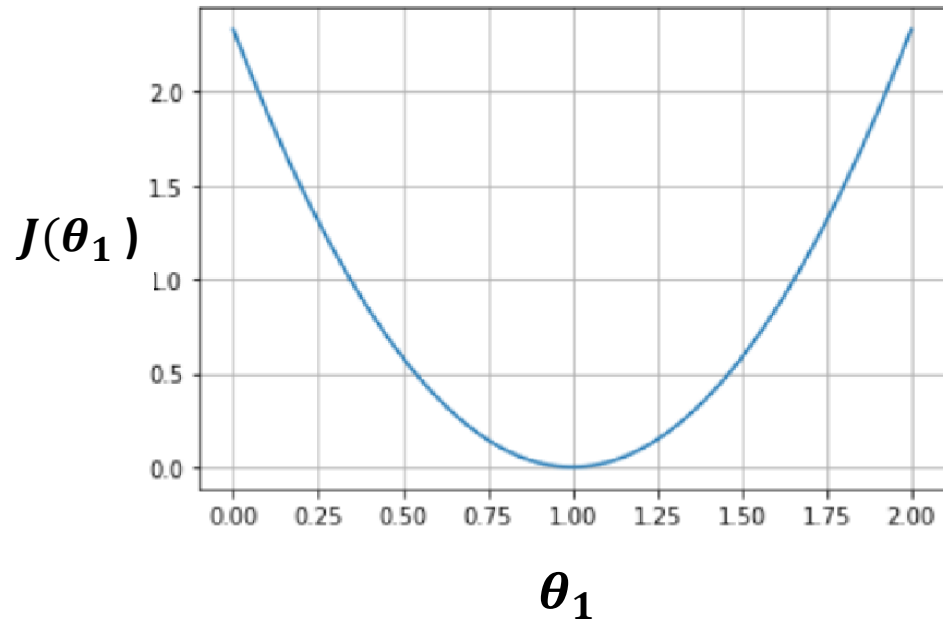


$$J(\theta_0 = 0, \theta_1 = 0.5) = \frac{1}{2m} \sum_{i=1}^{m} (0.5x_i - y_i)^2$$

$$= \frac{1}{2*3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{6} (3.5) = 0.58$$

--------------------------------------------------------

$$J(\theta_0 = 0, \theta_1 = 1) = \frac{1}{2m} \sum_{i=1}^{m} (x_i - y_i)^2$$

$$= \frac{1}{2*3} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$= \frac{1}{6} (0) = 0$$

٧

# example



$\theta_1$

| $\boldsymbol{\theta_1}$ | $\boldsymbol{J(\theta_1)}$ |
|---|---|
| 0 | 14/6 |
| 0.5 | 0.58 |
| 1 | 0 |
| 1.5 | 0.58 |
| 2 | 14/6 |

- Plotting the cost for each value of $\theta_1$

- The minimum point: $\theta_1 = 1$

- Using Grid Search to find best values of parameters
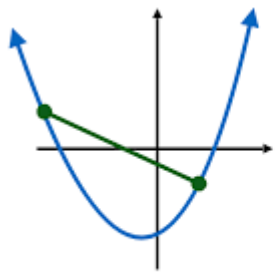
# Cost Function

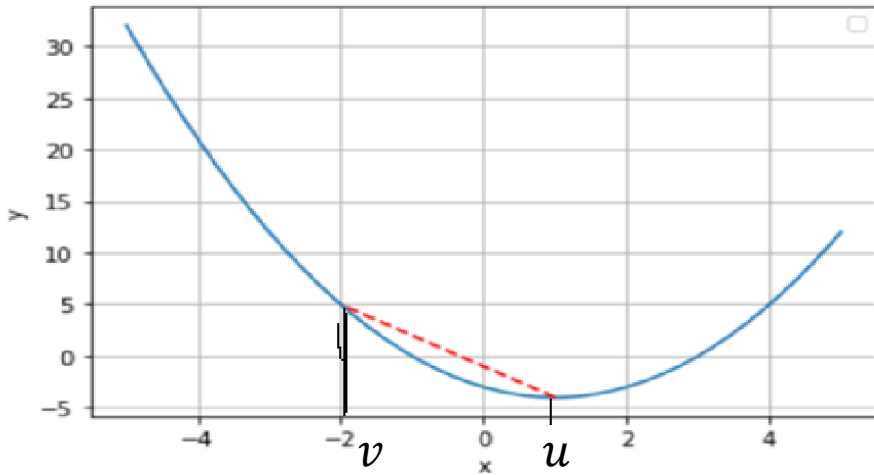- $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^{m} |h(x_i) - y_i|$    <span style="color:red">Mean absolute error(MAE)</span>

- Better for outliers compared with MSE
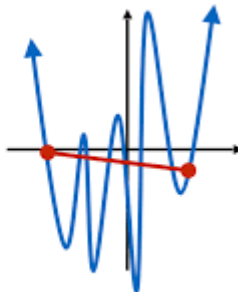
# Convexity



Function h(u) with u∈ X is convex if for any u, v ∈ X and for any $0 \leq \lambda \leq 1$ we have:

**h($\lambda$u +(1- $\lambda$)v) $\leq$ $\lambda$ h(u) + (1- $\lambda$) h(v)**

برای توابع محدب هر بهینه محلی یک بهینه سراسری است.

# example

$$if\ \theta_1 = -1:$$

$$MAE = \frac{1}{3}\ [\ |1 - (-1)| + |2 - (-2)| + |3 - (-3)|\ ] = 4$$

-----------------------------------------------------------------

$$if\ \theta_1 = 0:$$

$$MAE = \frac{1}{3}\ [\ |1 - 0| + |2 - 0| + |3 - 0|\ ] = 2$$

-----------------------------------------------------------------

$$if\ \theta_1 = 1:$$

$$MAE = \frac{1}{3}\ [\ |1 - 1| + |2 - 2| + |3 - 3|\ ] = 0$$

-----------------------------------------------------------------

$$if\ \theta_1 = 2:$$

$$MAE = \frac{1}{3}\ [\ |1 - 2| + |2 - 4| + |3 - 6|\ ] = 2$$

# example



MAE is convex

| $\theta_1$ | $J(\theta_1)$ |
|:---:|:---:|
| -1 | 4 |
| -0.5 | 3 |
| 0 | 2 |
| 0.5 | 1 |
| 1 | 0 |
| 1.5 | 1 |
| 2 | 2 |
| 2.5 | 3 |
| 3 | 4 |

# Cost Function

$h_\theta(x_i) = \theta_0 + \theta_1 \, x_i$

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

**Minimize** $J(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})$

$\boldsymbol{\theta_0}, \boldsymbol{\theta_1}$

If $J(\theta_1) = (\theta_1 - 2)^2$

$\dfrac{dJ(\theta_1)}{d\theta_1} = 0$ $\longrightarrow$ $\dfrac{dJ(\theta_1)}{d\theta_1} = 2(\theta_1 - 2) = 0$ $\longrightarrow$ $\theta_1 = 2$

١٣

# Gradient Descent

**Minimize** $J(\theta_0, \theta_1)$
$\quad \theta_0, \theta_1$

**Minimize** $J(\theta_0, \theta_1, \dots, \theta_n)$
$\quad \theta_0, \theta_1, \dots, \theta_n$

Repeat until convergence:

$\quad$ For j=0,…,n

$$\theta_j = \theta_j - \alpha \frac{dJ(\theta_0, \theta_1, \dots, \theta_n)}{d\theta_j}$$

$\alpha$ **is learning rate**

**Updating all $\theta_j$ $Simultaneous$ly**

Convergence condition:
$$\|\theta^{t+1} - \theta^t\|_2 \leq \varepsilon$$

# Gradient Descent

Correct form

$\text{temp0} = \boldsymbol{\theta_0} - \alpha \dfrac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_0}}$

$\text{temp1} = \boldsymbol{\theta_1} - \alpha \dfrac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_1}}$

$\boldsymbol{\theta_0} = \text{temp0}$
$\boldsymbol{\theta_1} = \text{temp1}$

✔

Incorrect form

$\boldsymbol{\theta_0} = \boldsymbol{\theta_0} - \alpha \dfrac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_0}}$

$\boldsymbol{\theta_1} = \boldsymbol{\theta_1} - \alpha \dfrac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_1}}$

✖

# Gradient Descent

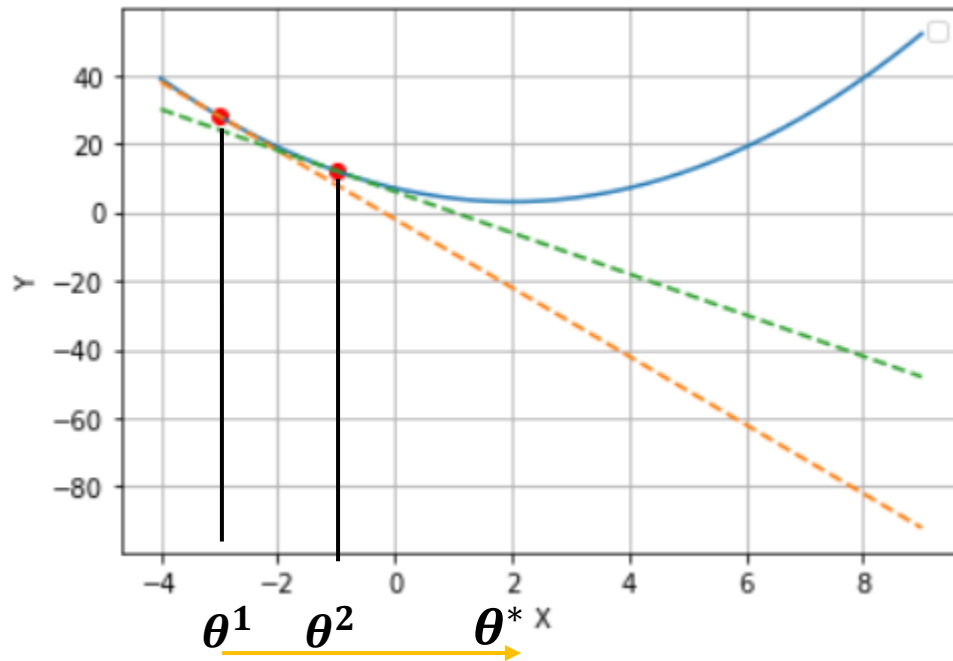خطوط مماس نشان داده شده دارای شیب یا مشتق مثبت هستند. درنتیجه:

$$\frac{dJ(\theta^1)}{d\theta^1} > 0 \; , \; \alpha > 0 \implies \alpha \frac{dJ(\theta^1)}{d\theta^1} > 0$$

$$\implies \theta^2 = \theta^1 - \alpha d\theta^1$$

$\theta$ کوچکتر میشود و به سمت چپ حرکت میکنیم.

# Gradient Descent


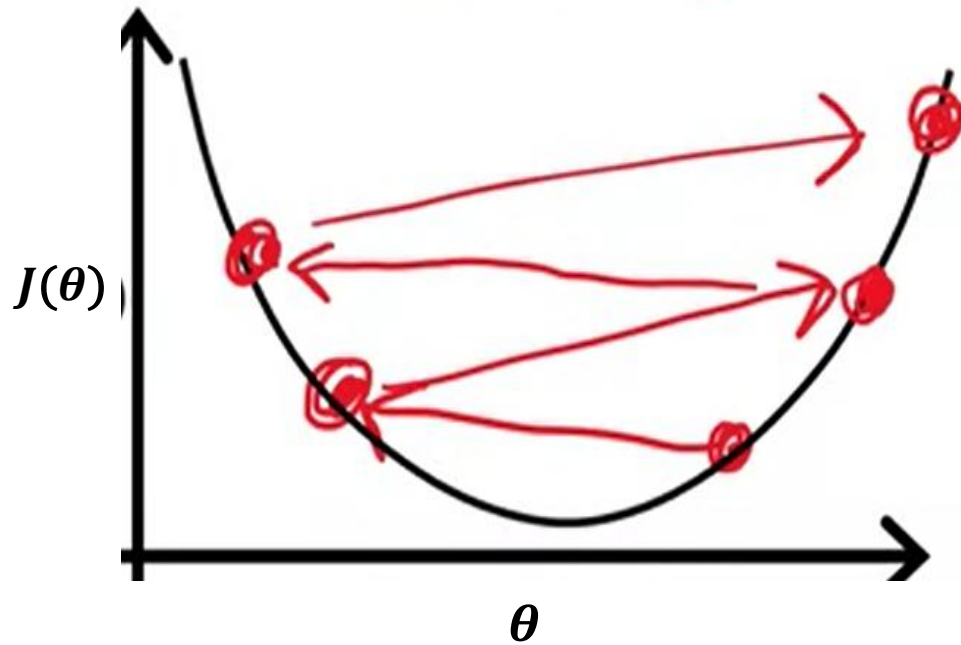
خطوط مماس نشان داده شده دارای شیب یا مشتق مثبت هستند.
درنتیجه:

$$\frac{dJ(\theta^1)}{d\theta^1} < 0 \ , \ \alpha > 0 \implies \alpha\frac{dJ(\theta^1)}{d\theta^1} < 0$$

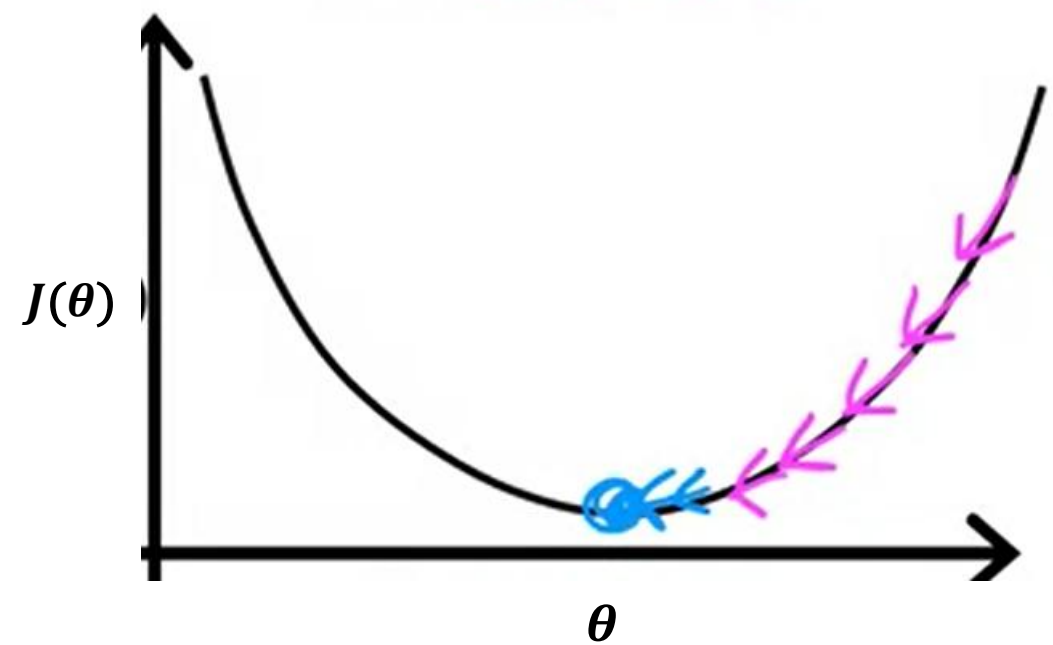$$\implies \theta^2 = \theta^1 - \alpha d\theta^1$$

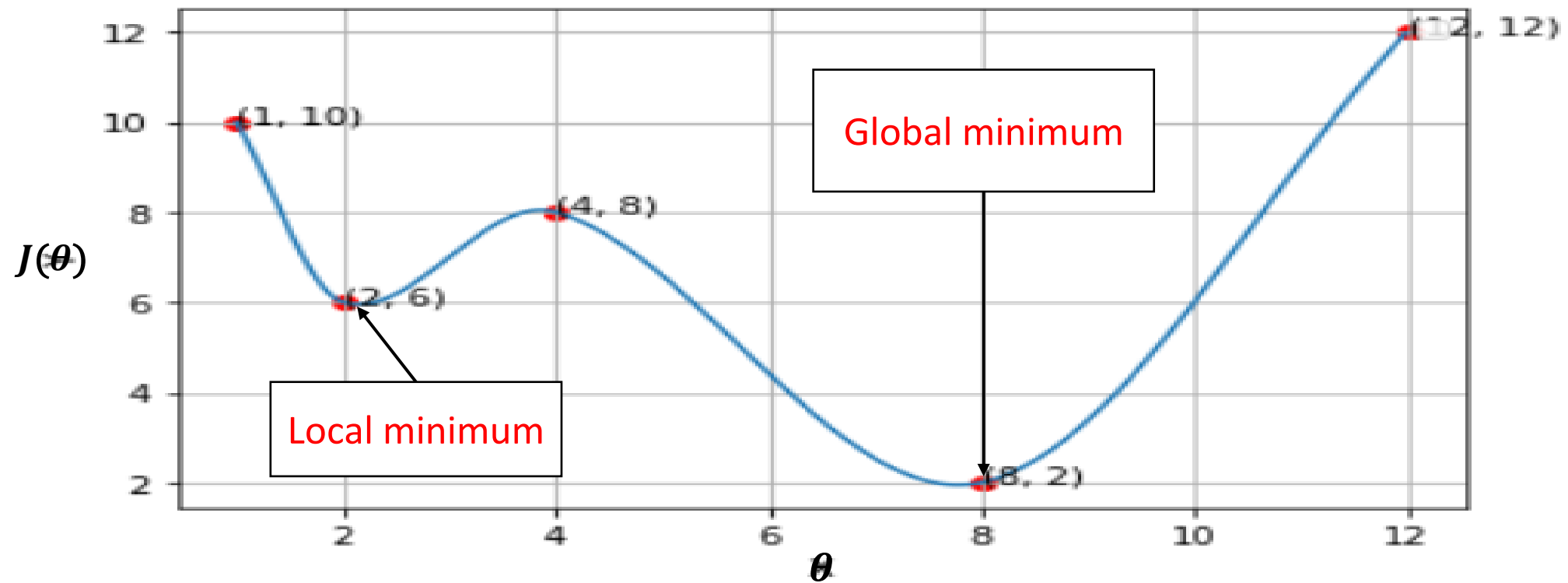$\theta$ بزرگتر میشود و به سمت راست حرکت میکنیم.

# Choosing Learning Rate



$\alpha$ is too large

$\alpha$ is small

# Gradient Descent Weakness

# Linear regression model

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

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

$\frac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_0}} = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i)$

$\frac{dJ(\boldsymbol{\theta_0}, \boldsymbol{\theta_1})}{d\boldsymbol{\theta_1}} = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i) x_i$

# Linear regression model

Repeat until convergence:

$$\boldsymbol{\theta_0} = \boldsymbol{\theta_0} - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i)$$

$$\boldsymbol{\theta_1} = \boldsymbol{\theta_1} - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i) x_i$$

بروز رسانی همزمان

$$\theta^t = \begin{bmatrix} \boldsymbol{\theta_0} \\ \boldsymbol{\theta_1} \end{bmatrix} \quad , \qquad \theta^{t+1} = \begin{bmatrix} \boldsymbol{\theta_0} \\ \boldsymbol{\theta_1} \end{bmatrix} \quad , \qquad d\theta = \begin{bmatrix} d\boldsymbol{\theta_0} \\ d\boldsymbol{\theta_1} \end{bmatrix}$$

Convergence condition:

- $\|\theta^{t+1} - \theta^t\|_2 = \sqrt[2]{(\theta_0^{t+1} - \theta_0^t)^2 + (\theta_1^{t+1} - \theta_1^t)^2} < \varepsilon$

- $\|d\theta\|_2 < \varepsilon$

# <mark>Batch</mark> Gradient Descent

$$\frac{dJ(\boldsymbol{\theta})}{d\boldsymbol{\theta}} = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x_i) - y_i \right)$$

# Batch Gradient Descent

$\theta_0 \longleftarrow random,\quad \theta_1 \longleftarrow random$

Repeat until convergence: $\Big\{$

J $\longleftarrow$ 0 , $d\theta_1 \longleftarrow$ 0 , $d\theta_0 \longleftarrow$ 0

For i = 1 to m:

$$h_\theta(x_i) = \theta_0 + \theta_1\, x_i$$
$$j\ += (h_\theta(x_i) - y_i)^2$$
$$d\theta_1\ += 2\,(h_\theta(x_i) - y_i)\, x_i$$
$$d\theta_0\ += 2\,(h_\theta(x_i) - y_i)$$

J /= 2m

$d\theta_1$ /= 2m

$d\theta_0$ /= 2m

$\theta_1 = \theta_1 - \alpha d\theta_1$

$\theta_0 = \theta_0 - \alpha d\theta_0$

$\Big\}$

# Gradient Descent

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

$$\boldsymbol{\theta_1} = \boldsymbol{\theta_1} - \alpha \frac{dJ(\boldsymbol{\theta_1})}{d\boldsymbol{\theta_1}}$$

$J(\theta)$

$\theta$

| number | size | #bedrooms | # floors | Price(y) |
|--------|------|-----------|----------|----------|
| 1 | 100 | 2 | 1 | 10000 |
| 2 | 150 | 3 | 2 | 175000 |
| ... | ... | ... | ... | ... |
| m | ... | ... | ... | ... |

n: #features = 3

m: #training data

$x_i$: i th data in training set

$x_j^i$: j th feature of i th data in training set

$$h_\theta(x^i) = \theta_0 + \theta_1 x_1^i + \theta_2 x_2^i + \ldots + \theta_n x_n^i$$

$$y = [y^1, y^2, \ldots, y^m]^T \in R^{m*1}$$

$$X = [x^1, x^2, \ldots, x^m]^T \in R^{m*(n+1)}$$

$$\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ . \\ . \\ . \\ x_n \end{bmatrix} \in R^{n+1} \quad x_0 = 1 \quad , \quad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ . \\ . \\ . \\ \theta_n \end{bmatrix} \in R^{n+1} \quad \theta_0 \text{ is bias}$$

$$\Longrightarrow \quad h_\theta(x) = x^T\theta = \theta^T x$$

# Cost function

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

$$e^i = (x^i)^T \theta - y^i \quad \Longrightarrow \quad \underset{m \times 1}{e} = \underset{m \times 1}{X\theta} - \underset{m \times 1}{y} \quad \Longrightarrow \quad J(\theta) = \frac{1}{2m} e^T e$$

$$e \, , \, X\theta \, , \, y \in R^m$$

# Gradient Descent

For j=0,…,n

$$\boldsymbol{\theta_j} = \boldsymbol{\theta_j} - \alpha \frac{d\boldsymbol{J(\theta_0, \theta_1, \ldots, \theta_n)}}{d\boldsymbol{\theta_j}}$$

②

$$\frac{d\boldsymbol{J(\theta)}}{d\boldsymbol{\theta}} = \frac{1}{m} X^T \text{e}$$

(n+1)*1

m*1

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

$$\frac{d\boldsymbol{J(\theta_0, \theta_1, \ldots, \theta_n)}}{d\boldsymbol{\theta_j}} = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i) \, x_j^i$$

(j=0,…,n , $x_0^i = 1$)

①

$$\frac{d\boldsymbol{J(\theta)}}{d\boldsymbol{\theta}_j} = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i) \, x_j^i$$

# حجم محاسبات ضرب ماتریس

$$A \in R^{a*b} \quad , \quad B \in R^{b*c} \quad \Longrightarrow \quad AB \in R^{a*c} \qquad (2b-1)\,ac\ flops$$

$$O(abc)$$

$Calculating:\ \textcolor{red}{e = \mathrm{X}\theta - y}$

a=m
b=n+1
c=1

$m(2n+1)$ (ضرب و جمع)
$m$ (تفریق)

2m(n+1) (مجموع) $\Longrightarrow$ $\textcolor{red}{O(mn)}$

$$\frac{dJ(\theta)}{d\theta} : (2\mathrm{m} - 1)(n+1) + (n+1) = 2m(n+1) \qquad \Longrightarrow$$

تقسیم بر m

$\textcolor{red}{O(mn)}$ در مجموع

$$\frac{1}{m}X^T \mathrm{e}$$

a=n+1
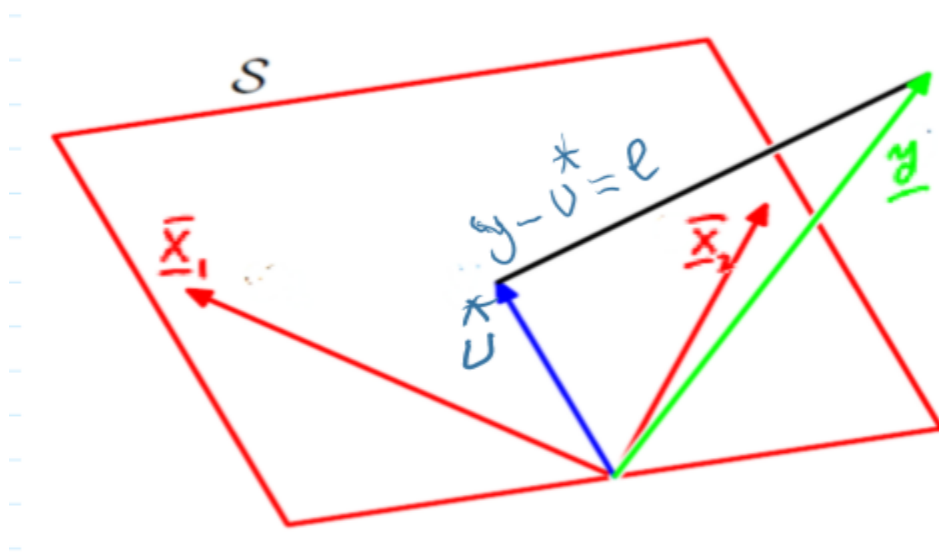b=m
c=1

# مفهوم هندسی

$$\text{Min}_{W} \| y - XW \|_2 = \text{Min} \| e \|_2$$

<span style="color:red">Span of X:</span>

فضایی که توسط ستون های X پوشش داده می شود. هر بردار در این فضا به صورت U = XW نشان داده می شود. و به آن span(X) می گویند. U بهینه که به صورت *U نشان داده می شود برداری است که e = y – U* بر span(X) عمود باشد. یا به عبات دیگر *Uای باید انتخاب شود که برابر با نگاشت y در span(X) باشد.

تعداد ویژگی * تعداد داده

۳۰

# Feature Scaling

$$x_1^i \ , x_2^i \ , \ldots , x_n^i$$

$$-1 \leq x_j \ \leq 1$$
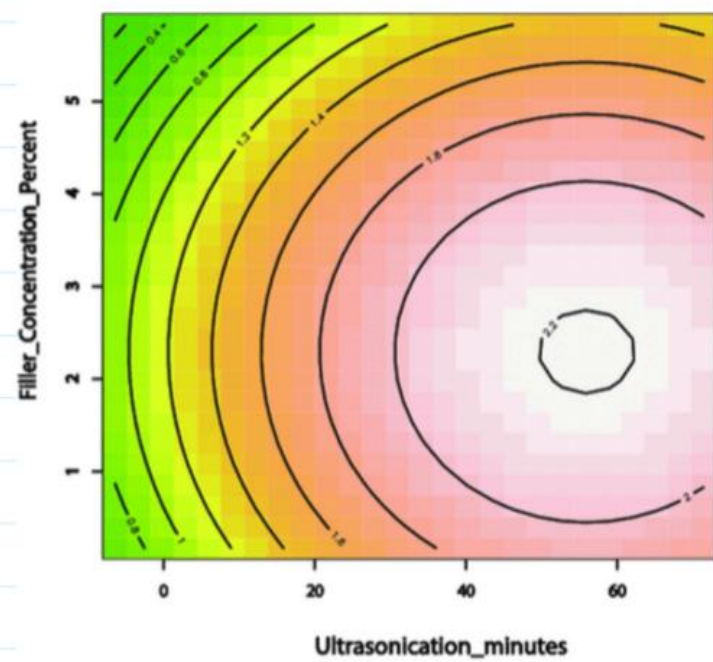
$0 < \ x_1 \ < 1000$

$\Rightarrow$

$x_1 : \dfrac{size}{1000}$

$0 < \ x_2 \ < 5$

$x_2 : \dfrac{\#bedrooms}{5}$
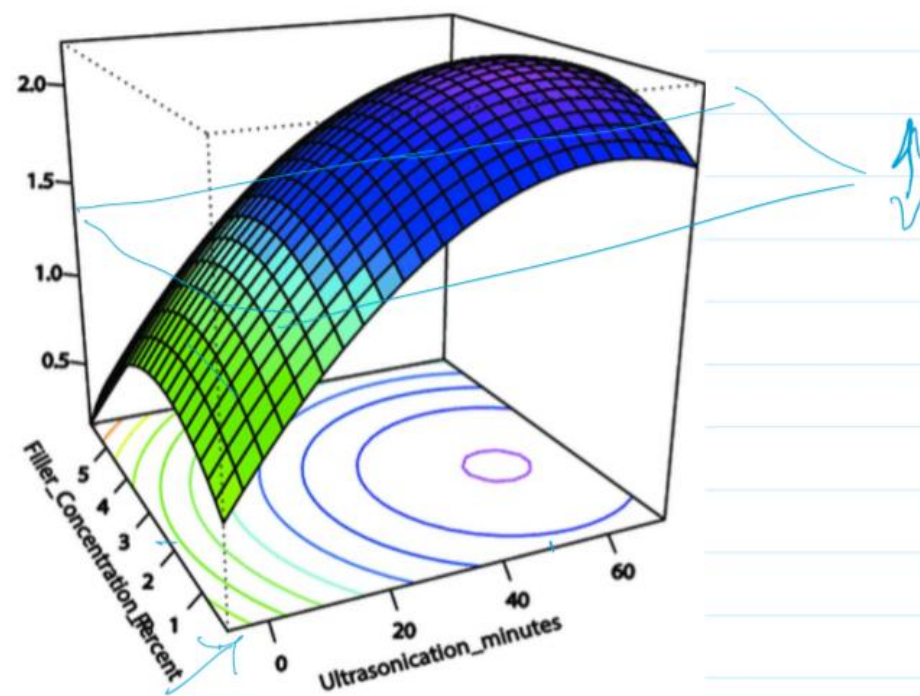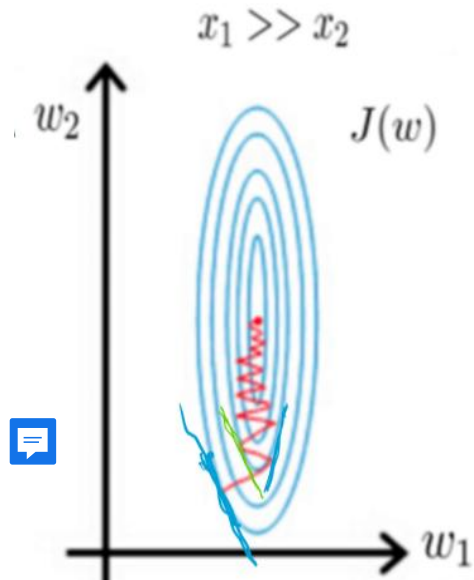
# Contour Plot

$$\frac{w_1^2}{b^2} + \frac{w_2^2}{a^2} = 1$$

2a: قطر بزرگ
2b: قطر کوچک

Gradient descent
without scaling

$x_1 \gg x_2$

$w_2$      $J(w)$

$w_1$

Gradient descent
after scaling variables

$0 \le x_1 \le 1$
$0 \le x_2 \le 1$

$w_2$      $J(w)$

$w_1$

$$\frac{w_1^2}{a^2} + \frac{w_2^2}{a^2} = 1$$

# Feature Scaling

**Scaled features:**

- $0 \leq x_1 \leq 3$ ✔
- $-3 \leq x_1 \leq 3$ ✔
- $-2 \leq x_2 \leq 0.5$ ✔
- $-\dfrac{1}{3} \leq x_2 \leq \dfrac{1}{3}$ ✔

**Need scaling:**

$-100 \leq x_3 \leq 100$ ✘

$-0.001 \leq x_4 \leq 0.001$ ✘

# Feature Scaling

$$x_1^* = \frac{x_1 - \mu_1}{standard\_deviation}$$

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$$\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_j^i$$
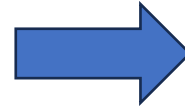
$$bedroom^* = \frac{bedroom - 2.5}{5}$$

$$size^* = \frac{size - 300}{2000}$$

# Creating New Features

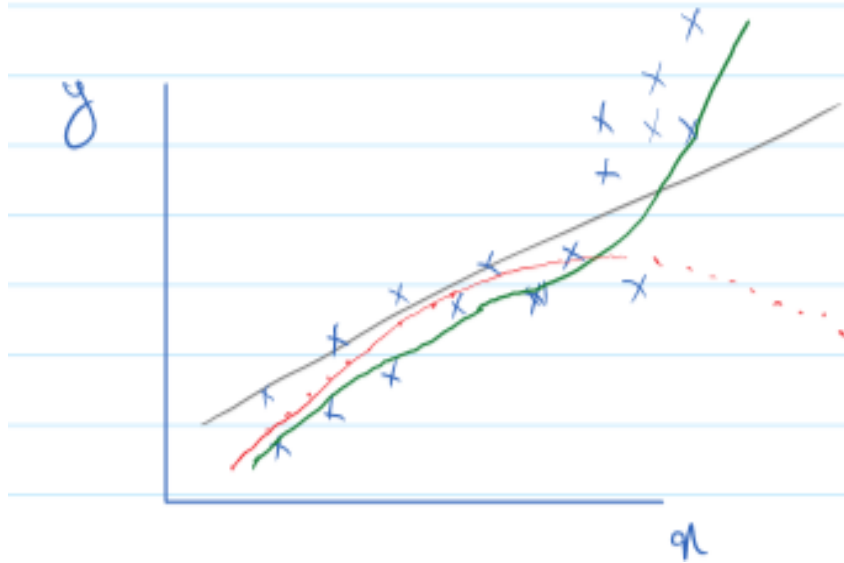$$h_\theta(x) = \theta_0 + \theta_1 \, x_1 + \theta_2 \, x_2$$

| قیمت خانه | طول خانه | عرض خانه |
|---|---|---|

( مساحت خانه ) $\quad x^* = x_1 * x_2$

$$h_\theta(x) = \theta_0 + \theta_1 \, x^*$$

# Creating New Features



درجه 2:

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

درجه 3:

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

Need scaling:

x: 0,...,1000
$x^2$: 0,..., $10^6$
$x^3$: 0,..., $10^9$

We can use:
$$x , x^2 , x^3 , \sqrt{x}$$
$$\theta_0 + \theta_1 \, x + \theta_2 \sqrt{x}$$