

# Houses prices

Size (ft <sup>2</sup> )	# bedrooms	# Age of home	# floors	price
$x_1$	$x_2$	$x_3$	$x_4$	$y$
—	—	—	—	—
—	—	—	—	—
—	—	—	—	—

$n$  = No. of features

$m$  = No. of training examples

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

$$\begin{array}{c}
 \text{cloud} \\
 h_{\theta}(x) \\
 m \times 1
 \end{array}
 \begin{bmatrix}
 h_1(x) \\
 h_2(x) \\
 \vdots \\
 h_m(x)
 \end{bmatrix}
 =
 \begin{array}{c}
 \begin{bmatrix}
 x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\
 x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\
 x_0^{(3)} & x_1^{(3)} & x_2^{(3)} & \dots & x_n^{(3)} \\
 \vdots & \vdots & \vdots & \ddots & \vdots \\
 x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)}
 \end{bmatrix} \\
 m \times n+1
 \end{array}
 \begin{array}{c}
 \begin{bmatrix}
 \theta_0 \\
 \theta_1 \\
 \vdots \\
 \theta_n
 \end{bmatrix} \\
 n+1 \times 1
 \end{array}$$

$$\text{cloud} \quad h_{\theta}(x) = X \theta$$

## \* Multiple Features

$$\begin{aligned}
 h_{\theta}(x) &= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\
 &\stackrel{m}{\leftarrow} \underset{x_0=1}{=} [\theta_0 \quad \theta_1 \quad \dots \quad \theta_n] \underset{(n+1)}{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}} \underset{n \times 1}{=} \theta^T x
 \end{aligned}$$

Multivariate Linear Regression

## \* Gradient Descent for Multiple Features

repeat until Convergence: {

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

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

{

## \* Feature Scaling

→ To speed up gradient descent by modifying the range of our inputs. to work on a smaller scale.

$$x_1 = \frac{\text{Size}}{2000}$$

$$x_1 (0 - 2000)$$

$$x_2 = \frac{\text{No. of rooms}}{5}$$

$$x_2 (1 - 5)$$

→ Divide by the range

So,

$$0 \leq x_i \leq 1$$

Max. Value - Min. Value

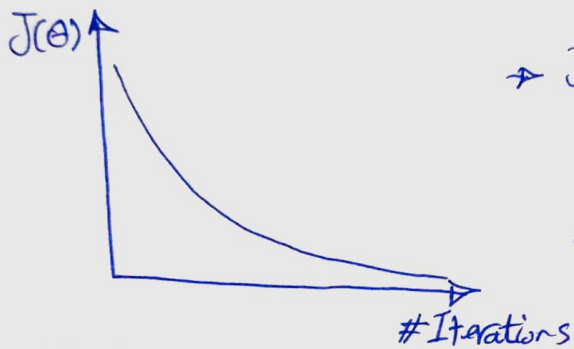
## \* Mean Normalization (Z-Score Normalization)

→ It is another method to reduce the steps of gradient descent by subtracting the feature average value from the input and dividing the result by the ~~std~~ deviation

$$x_i = \frac{x_i - \mu_i}{s_i} \rightarrow \text{avg. value of } x$$

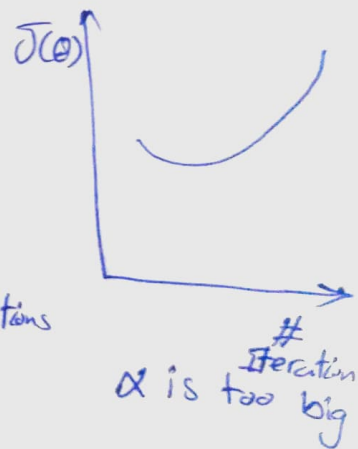
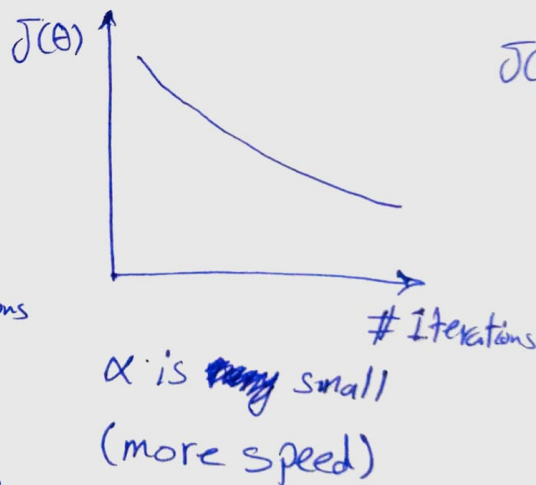
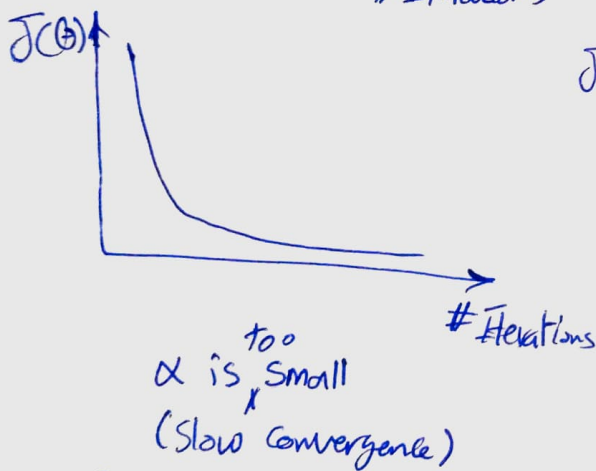
$s_i \rightarrow$  Standard deviation (~~Range = Max - Min~~)

## \* Learning Rate



→  $J(\theta)$  should decrease Every new iteration

→ The speed of  $J(\theta)$  decreasing depends on  $\alpha$



→ But, if  $\alpha$  is too small, gradient descent can be slow to converge.

## \* Features Reformation [Feature Engineering]

→ We Can Reform our features to improve Regression

$$h(x) = \theta_0 + \theta_1 \underset{\text{Width}}{x_1} + \theta_2 \underset{\text{length}}{x_2}$$

$$\Rightarrow h(x) = \theta_0 + \theta_1 \underset{\text{Area} = x_1 x_2}{x}$$

$$\text{let } J(\theta) = a\theta^2 + b\theta + c$$

To minimize  $J(\theta)$  we take the derivative  $\frac{dJ(\theta)}{d\theta}$  and set it to zero to solve for  $\theta$

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

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

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

$$\rightarrow X\theta = y \quad \text{where } \begin{array}{l} X \in \mathbb{R}^{m \times n+1} \\ \theta \in \mathbb{R}^{n+1} \\ y \in \mathbb{R}^m \end{array}$$

$$X^T X \theta = X^T y$$

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

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

We want to minimize the least-square cost function

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

$$h_\theta(x) = \theta^T X$$

$$\text{so, } J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y)$$

Matrix Notation

We can simply throw  $\frac{1}{2m}$  part away as we will make  $\frac{d}{d\theta} J(\theta) = 0$

$$\text{so, } J(\theta) = (X\theta)^T X\theta - (X\theta)^T y - y^T X\theta + y^T y$$

$$= \theta^T X^T X \theta - 2(X\theta)^T y + y^T y$$

$$\frac{dJ(\theta)}{d\theta} = 0 \Rightarrow$$

$$\cancel{2X^T X \theta} = \cancel{2X^T y}$$

$$2X^T X \theta - 2X^T y = 0 \Rightarrow X^T X \theta = X^T y$$

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

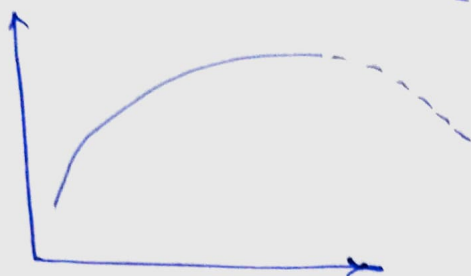


## Polynomial Regression

→ Linear Regression is not the best to fit the data well.

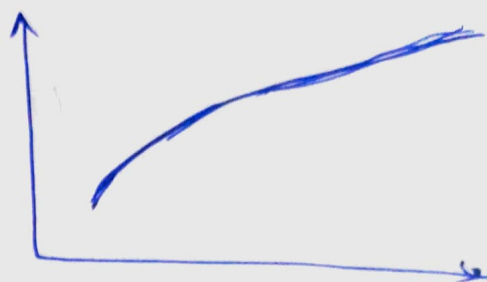
→ We need to change the hypothesis to be a quadratic, Cubic, or Square root ---- or any other form

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

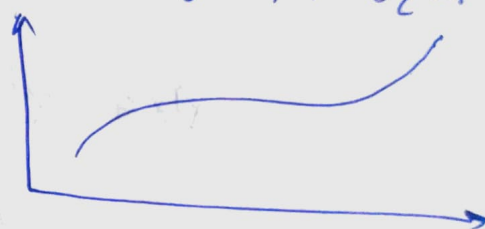


→ It goes down after awhile, so it is not good for our example (house pricing)

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$$



$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$$



→ Using the Square root function may simplify the work more than the Cubic function as the latter needs more variables

→ The range will be changed according to the new function

$$\text{if } 0 \leq x \leq 10$$

$$\rightarrow 0 \leq x^2 \leq 100$$

$$\rightarrow 0 \leq \sqrt{x} \leq \sqrt{10}$$

$$0 \leq x^3 \leq 1000$$

## Normal Equation

→ To solve for  $\theta$  analytically

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

→ No need for feature scaling

Design Matrix

$$X = \begin{bmatrix} 1 & x_1 & x_2 & \dots \\ 1 & x_1 & x_2 & \dots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}, y = \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix}$$

m-vector

$m \times (n+1)$   
↓  
No. of features  
→ No. of training examples

## Gradient Descent

→ Need to choose  $\alpha$

→ Need several iterations

→  $O(n^2)$

→ Works well when  $n$  is big  
10,000 ↑

## Normal Equation

→ No need to choose  $\alpha$

→ No need to iterate

→  $O(n^3)$ , need to calculate  $(X^T X)^{-1}$

→ Slow if  $n$  is very large

\* What if  $(X^T X)$  isn't invertible? [Singular = Noninvertible]

→ This may occur due to,

① Redundant features (linearly dependent)

$x_1 = \text{area in ft}^2$

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

→ delete one

② Too many features ( $m \ll n$ )

→ delete some features

→ Regularization

⇒ For OLS: P-inv → pseudo inverse will compute  $(X^T X)$  even if it is Noninvertible