

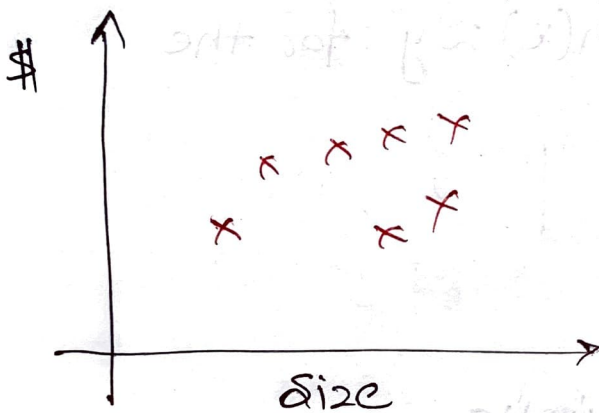
# CS229: Linear Regression

①

↳ Supervised learning

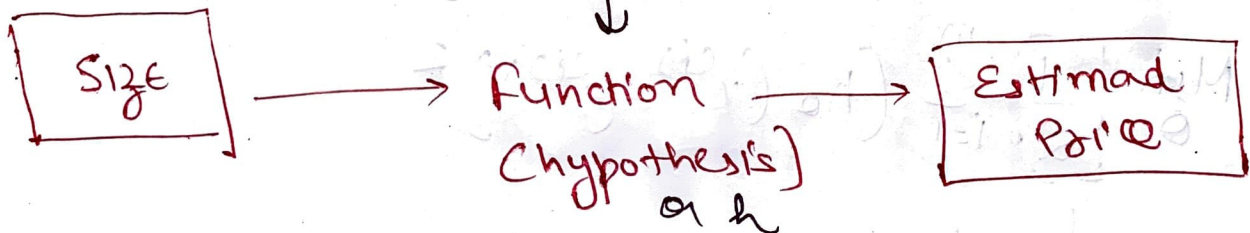
$X \rightarrow y$

Size	# Bedrooms	Price (\$)
2104	3	400
1416	2	232
1534	3	315
852	2	178



TRAINING SET

↓  
LEARNING ALGORITHM



Q. How to represent  $h$ ?

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$x_1 = \text{Size}$

$x_2 = \text{\# bedrooms}$

$$h(x) = \sum_{j=0}^2 \theta_j x_j$$

where  $x_0 = 1$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

$\theta$  = parameters

$m$  = # training sample

$x$  = inputs / features

$y$  = output (target variable)

$$h(x) = \sum_{j=0}^2 \theta_j x_j \text{ where } x_0 = 1$$

Choose  $\theta$  such that  $h(x) \approx y$  for the training samples

$$h_0(x) = h(x)$$

Ordinary least squares

$$(h_0(x) - y)^2 \rightarrow \text{minimize}$$

$$\min_{\theta} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$\min_{\theta} \left\{ \frac{1}{2} \right\} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$\hookrightarrow$  for easier derivative

\* Cost function

$$J(\theta) = \min_{\theta} \frac{1}{2} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

## \* GRADIENT DESCENT :

(3)

1. Start with some  $\theta$  (say  $\theta = \vec{0}$ ) [Zero vector]
2. Keep changing the  $\theta$  to reduce  $J(\theta)$ .

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

$\rightarrow \alpha = \text{learning rate}$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$= 2 \times \frac{1}{2} (h_{\theta}(x) - y) \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$\Rightarrow (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n - y)$$

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

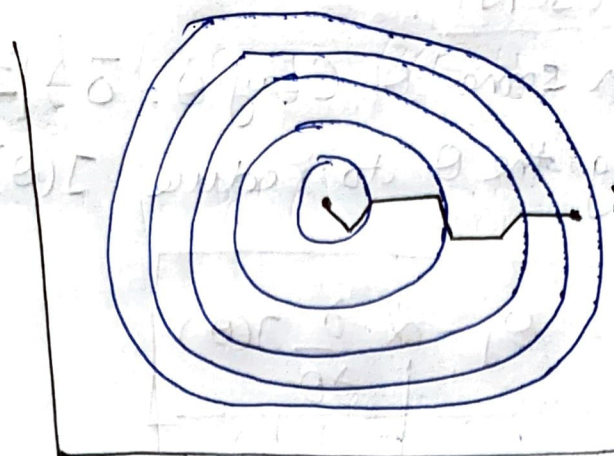
$$\theta_j := \theta_j - \alpha (h_{\theta}(x) - y) \cdot x_j$$

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

Repeat until convergence  
for  $(j = 0, 1, \dots, n)$



$\theta_1$



$\theta_0$

### BATCH GRADIENT DESCENT

↳ You look at entire dataset as a batch of data

Repeat  $\ell$

for  $j=1$  to  $m$

$$\theta_j := \theta_j - \alpha (h(x^{(j)}) - y^{(j)}) \cdot x_j^{(j)}$$

### STOCHASTIC GRADIENT DESCENT

↳ using single samples

$\theta_1$



$\theta_0$

} follows a noisy path

# # NORMAL EQUATION

(3)

↳ Only for linear Regression

↳ Reaches Global Optimum in onestep.

$\nabla_{\theta} J(\theta)$  — derivative

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

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

$$\begin{bmatrix} \text{Design} \\ \text{matrix} \end{bmatrix} X\theta = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$= \begin{bmatrix} x^{(1)T} \theta \\ x^{(2)T} \theta \\ \vdots \\ x^{(m)T} \theta \end{bmatrix} = \begin{bmatrix} h_0(x^{(1)}) \\ h_0(x^{(2)}) \\ \vdots \\ h_0(x^{(n+1)}) \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

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

$$x\theta - y = \begin{bmatrix} h(x^{(1)}) - y^{(1)} \\ \vdots \\ h(x^{(m)}) - y^{(m)} \end{bmatrix}$$

$$\because Z^T Z = Z^2$$

$$(x\theta - y)^T (x\theta - y)$$

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

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

$$= \frac{1}{2} \nabla_{\theta} [\theta^T x^T x\theta - \theta^T x^T y - y^T x\theta + y^T y]$$

$$\left[ (ax - b)(ax - b) = a^2 x^2 - axb - bax + b^2 \right]$$

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

$$= x^T x\theta - x^T y \stackrel{\text{set}}{=} 0$$

$$x^T x\theta = x^T y$$

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

→ Normal equation