

# Normal equation

we will know that the normal equation give us the best fit model parameter (or feature weights)  $\theta$  through the equation

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

- Our target is to prove and understand this formula
- This require you to know a little bit about Linear Algebra. Don't worry, I will go through some important things before I come to prove.

$$1) (AB)^T = B^T \cdot A^T$$

2) If  $A$  is symmetric :

$$\rightarrow A^T = A$$

3) give you a scalar  $\alpha = y^T A x$

we know that because  $\alpha$  is scalar, we have :  $\alpha = \alpha^T$

$\rightarrow$  From that :

$$\bullet y^T A x = x^T A^T y$$

• partial derivative :

$$\frac{\partial \alpha}{\partial x} = y^T \cdot A$$

$$\frac{\partial \alpha}{\partial y} = x^T A^T$$

4) given scalar :  $\alpha = x^T A x$

•  $A$  is not symmetric

$$\frac{\partial \alpha}{\partial x} = x^T (A + A^T)$$

•  $A$  is symmetric

$$\frac{\partial \alpha}{\partial x} = 2 x^T A$$

Revise:

$$y = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T x$$

with  $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1}$  ;  $x = \begin{bmatrix} x_0 = 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1}$

→ with many data, we write it in the form

$$\hat{y} = X\theta$$

$$X = \begin{bmatrix} x_0^1 & x_1^1 & \dots & x_n^1 \\ x_0^2 & x_1^2 & \dots & x_n^2 \\ \vdots & \vdots & & \vdots \\ x_0^m & x_1^m & \dots & x_n^m \end{bmatrix}_{m \times n} \quad \cdot \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1} \quad \cdot \quad \hat{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}$$

•  $n$  is number of features

•  $m$  is number of samples

•  $x_j^i$  is the  $j^{\text{th}}$  feature in  $i^{\text{th}}$  sample

⇒ We want to minimize the cost function.

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

**Prove** • we can drop  $\frac{1}{2}$  in  $J(\theta)$  for easier to minimize bec se it's constant and do not have impact when you want to minimize  $J(\theta)$

• Rewritten  $\sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$  in matrix form:

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

$$= \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$= \begin{bmatrix} (\hat{y}^{(1)} - y^{(1)}) & (\hat{y}^{(2)} - y^{(2)}) & \dots & (\hat{y}^{(m)} - y^{(m)}) \end{bmatrix}_{1 \times m} \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \\ \vdots \\ \hat{y}^{(m)} - y^{(m)} \end{bmatrix}_{m \times 1}$$

$$= (X\theta - y)^T \cdot (X\theta - y)$$

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

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

→ From now we have  $J(\theta) = \theta^T X^T X \theta - \theta^T X^T y - y^T X \theta + y^T y$

• We want to find  $\theta$  minimize  $J(\theta)$  or  $\theta = \underset{\theta}{\operatorname{argmin}} J(\theta)$

take derivative of  $J(\theta)$  with respect to  $\theta$

$$\begin{aligned} \frac{\partial J}{\partial \theta} &= \frac{\partial}{\partial \theta} \left( \underbrace{\theta^T X^T X \theta}_{(u)} - \underbrace{\theta^T X^T y}_{(3)} - \underbrace{y^T X \theta}_{(3)} + y^T y \right) \\ &= 2\theta^T X^T X - y^T X - y^T X + 0 \\ &= 2\theta^T X^T X - 2y^T X \end{aligned}$$

this symmetric because  $X^T X = (X^T X)^T$

→ check out the Linear Algebra above to understand this step

Find the critical point to have the  $\theta$  that minimize  $J(\theta)$

$$\frac{\partial J}{\partial \theta} = 0$$

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

$$\Rightarrow (\theta^T X^T X)^T = (y^T X)^T$$

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

$$\Rightarrow \theta = (X^T X)^{-1} X^T y \quad (\text{prove})$$

Conclusion:  $\theta = (X^T X)^{-1} X^T y$  that is best fit for model

