

As we know, $A^T A = \sum_i A_i^2$

$$\frac{1}{2} (X\theta - \bar{y})^T (X\theta - \bar{y}) = \frac{1}{2} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$= J(\theta)$$

$$X \in \mathbb{R}^{m \times n}$$

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

$$y \in \mathbb{R}^{m \times 1}$$

To minimize J , we differentiate w.r.t θ

$$\nabla_{\theta} (A^T B A^T C) = B^T A^T C^T + B A^T C$$

[1x1 - scalar]

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

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y}) \rightarrow 1$$

$$= \frac{1}{2} \nabla_{\theta} \text{tr}(\theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y}) \rightarrow 2$$

$$= \frac{1}{2} \nabla_{\theta} (\text{tr}(\theta^T X^T X \theta) - 2 \text{tr}(\bar{y}^T X \theta)) \rightarrow 3$$

$$= \frac{1}{2} (X^T X \theta + X^T X \theta - 2 X^T \bar{y}) \rightarrow 4$$

$$= X^T X \theta - X^T \bar{y}$$

Understandly

$$X \in \mathbb{R}^{m \times n}, \theta \in \mathbb{R}^{n \times 1}, y \in \mathbb{R}^{m \times 1}$$

The term $(X\theta - \bar{y})^T (X\theta - \bar{y})$ is a scalar (1x1 matrix)

we multiplied $\frac{1}{2}$ for convenience

expanding the $(X\theta - \bar{y})^T (X\theta - \bar{y})$

$$= \theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y}$$

So, property $\rightarrow \nabla_{\theta} \text{tr}(\theta^T A \theta) = (A + A^T) \theta$

$= 2 \times \text{tr}(\theta \times \theta)$

Second term, $(\bar{y}^T \times \theta)$

$y^T x \in \mathbb{R}^{1 \times n}$

$\theta \in \mathbb{R}^{n \times 1} \rightarrow y^T x \theta \in \mathbb{R}^{1 \times 1} \rightarrow \text{scalar}$

$\text{tr}(y^T x \theta) = y^T x \theta$

$\nabla_{\theta} \text{tr}(A \theta) = A^T$ (if A is const w.r.t θ)

$A = y^T x$

$\nabla_{\theta} \text{tr} = x^T \bar{y}$

Final gradient

$\nabla_{\theta} J(\theta) = \frac{1}{2} (2 \times \text{tr}(\theta \times \theta) - 2 \times \text{tr}(\bar{y})) = \boxed{x^T x \theta - x^T \bar{y}}$

to minimize J, we set $\nabla = 0$

$x^T x \theta - x^T \bar{y} = 0$

$x^T x \theta = x^T \bar{y}$

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