

$$MSE(w) = \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2 \quad \text{ενενων}$$

$$h_w(x) = w \cdot x \Rightarrow MSE(w) = \frac{1}{m} \sum_{i=1}^m (w x^{(i)} - y^{(i)})^2 = C(w)$$

And ~~Property~~:

$$\vec{a}^T \cdot \vec{a} = \sum_i a_i^2 \quad (a)$$

Property:

$$\frac{1}{m} (\vec{w}x - \vec{y})^T (\vec{w}x - \vec{y}) = \frac{1}{m} \sum_{i=1}^m (w x^{(i)} - y^{(i)})^2 = C(w)$$

Derivatives with respect to w :

$$\nabla_w C(w) \stackrel{(a)}{=} \nabla_w \left\{ \frac{1}{m} (\vec{w}x - \vec{y})^T (\vec{w}x - \vec{y}) \right\}$$

$$\stackrel{(b)}{=} \frac{1}{m} \nabla_w \left\{ \vec{w}^T \vec{x}^T \vec{w}x - \vec{w}^T \vec{x}^T \vec{y} - \vec{w}x \vec{y}^T + \vec{y}^T \vec{y} \right\}$$

$$\stackrel{(c)}{=} \frac{1}{m} \nabla_w \left\{ \text{tr} \vec{w}^T \vec{x}^T \vec{w}x - \text{tr} \vec{w}^T \vec{x}^T \vec{y} - \text{tr} \vec{w}x \vec{y}^T + \cancel{\text{tr} \vec{y}^T \vec{y}} \right\}$$

$$\stackrel{(d)}{=} \frac{1}{m} \nabla_w \left\{ \text{tr} \vec{w}^T \vec{x}^T \vec{w}x - 2 \text{tr} \vec{w}x \vec{y}^T \right\} \quad \left(\begin{array}{l} \vec{w}x \vec{y}^T = \\ = \vec{w}(\vec{x}^T \vec{y})^T \end{array} \right)$$

$$\stackrel{(e)}{=} \frac{1}{m} \nabla_w \left\{ \text{tr} \vec{w}^T \vec{x}^T \vec{w}x - 2 \vec{x}^T \vec{w} - 2 \vec{x}^T \vec{y} \right\}$$

$$2 \vec{x}^T (\vec{x} \vec{w})$$

'Aps

→ ...

And derivative:

$$\vec{a}^T \cdot \vec{a} = \sum_i a_i^2 \quad (a)$$

Property:

$$\frac{1}{n} (\mathbf{W}\mathbf{X} - \vec{y})^T (\mathbf{W}\mathbf{X} - \vec{y}) = \frac{1}{n} \sum_{i=1}^m (w x^{(i)} - y^{(i)})^2 = C(w)$$

Derivatives with respect to w :

$$\nabla_w C(w) \stackrel{(a)}{=} \nabla_w \left\{ \frac{1}{n} (\mathbf{W}\mathbf{X} - \vec{y})^T (\mathbf{W}\mathbf{X} - \vec{y}) \right\}$$

$$\stackrel{(b)}{=} \frac{1}{n} \nabla_w \left\{ \mathbf{W}^T \mathbf{X}^T \mathbf{W} \mathbf{X} - \mathbf{W}^T \mathbf{X}^T \vec{y} - \mathbf{W} \mathbf{X} \vec{y}^T + \vec{y}^T \vec{y} \right\}$$

$$\stackrel{(c)}{=} \frac{1}{n} \nabla_w \left\{ \text{tr} \mathbf{W}^T \mathbf{X}^T \mathbf{W} \mathbf{X} - \text{tr} \mathbf{W}^T \mathbf{X}^T \vec{y} - \text{tr} \mathbf{W} \mathbf{X} \vec{y}^T + \text{tr} \vec{y}^T \vec{y} \right\}$$

$$\stackrel{(d)}{=} \frac{1}{n} \nabla_w \left\{ \text{tr} \mathbf{W}^T \mathbf{X}^T \mathbf{W} \mathbf{X} - 2 \text{tr} \mathbf{W} \mathbf{X} \vec{y}^T \right\} \quad \left(\begin{array}{l} \mathbf{W} \mathbf{X} \vec{y}^T = \\ = \mathbf{W} (\mathbf{X}^T \vec{y})^T \end{array} \right)$$

$$\stackrel{(e)}{=} \frac{1}{n} \left\{ \text{tr} \mathbf{W}^T \mathbf{X}^T \mathbf{W} \mathbf{X} (2 \mathbf{X} \mathbf{X}^T \mathbf{W} - 2 \mathbf{X}^T \vec{y}) \right\}$$

$$= \frac{2}{n} \left(\mathbf{X}^T (\mathbf{X} \mathbf{W} - \vec{y}) \right)$$

\Rightarrow $\begin{array}{l} \text{Απns} \\ \text{Ταλ φωνs} \\ 1115201700170 \end{array}$

(b) we simplified the cost function and used the linearity of the gradient operator.

(c) we used the fact that a trace of a real number is the number itself and we eliminated the $\vec{y}^T \vec{y}$ term because it has no dependency on w .

(d) we used the fact that $\text{tr} a = \text{tr} a^T$.

(e) We used the following rules of matrix calculus.

$$\nabla_A \text{tr. } A^T B A = (B + B^T) A \text{ where } A = w, B = x x^T$$

$$\nabla_A \text{tr} B^T A = B \text{ where } A = w, B = x^T y$$