

Linear Regression Recap

- Training set =  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$   
 features  $\in \mathbb{R}^d$   $\leftarrow$   $\leftarrow$  output  $\in \mathbb{R}$
- Goal: Compute hypothesis  $h_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$  where  

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$
 such that  $h_\theta(x)$  is a good prediction for  
 new feature vectors  $x$ . (i.e.  $x \neq x^{(i)} \forall i$ )
- Under the assumption that the training set is  
 "representative" of future feature vectors,  
 it is reasonable to minimize the cost of  
 prediction of the hypothesis on training set.

Linear Regression under squared loss is

Given training set  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  find  $\theta = (\theta_0, \dots, \theta_n)$   
 such that

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2}_{\text{cost/loss on a single training data}} \quad \text{Total cost on Training} \quad \text{is minimized}$$

$$= \frac{1}{2n} \sum_{i=1}^n \left( \left( \theta_0 + \sum_{j=1}^d \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$$

## Math Recap

(2)

Convention: Vectors in  $\mathbb{R}^d$  thought of as  $(d \times 1)$  matrices  
(column vectors)

$$\|x\|_2^2 = \sum_{i=1}^d x_i^2 = x^T x$$

Partial Derivative  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  i.e.  $f(\theta_1, \dots, \theta_d) \in \mathbb{R}$

$\frac{\partial f}{\partial \theta_i}$  - measures how  $f$  changes as  $\theta_i$  changes  
(while  $\theta_1, \theta_2, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_n$  are fixed).

Gradient:  $\nabla_{\theta} f = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \frac{\partial f}{\partial \theta_d} \end{bmatrix}$

Standard Gradients:

- For  $f(\theta_1, \dots, \theta_d) = \sum_{i=1}^d \theta_i x_i = x^T \theta = \theta^T x$ ,

$$\nabla_{\theta} f = x$$

- For  $f(\theta) = \theta^T A \theta$  where  $A$  is  $d \times d$ ,

$$\nabla_{\theta} f = (A + A^T) \theta.$$

$$X = \begin{bmatrix} 1 & \dots & (x^{(1)})^T \\ \vdots & & \vdots \\ 1 & \dots & (x^{(n)})^T \\ \vdots & & \vdots \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} 1 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}} \right\} n$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}} \right\} d+1$$

$$X\theta = \begin{bmatrix} \theta_0 + \theta_1 x_1^{(1)} + \dots + \theta_d x_d^{(1)} \\ \vdots \\ \theta_0 + \theta_1 x_1^{(n)} + \dots + \theta_d x_d^{(n)} \end{bmatrix} = \begin{bmatrix} h_\theta(x^{(1)}) \\ \vdots \\ h_\theta(x^{(n)}) \end{bmatrix}$$

$$X\theta - y = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ \vdots \\ h_\theta(x^{(n)}) - y^{(n)} \end{bmatrix} \quad \text{where } y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(i)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Linear Regression: Want to find  $\theta \in \mathbb{R}^{d+1}$  s.t.

minimize  $J(\theta) = \frac{1}{2n} \|X\theta - y\|_2^2$   
 $= \frac{1}{2n} (X\theta - y)^T (X\theta - y)$

Find  $\theta$  s.t.  $\nabla_\theta J(\theta) = 0$ .

$$\begin{aligned} \nabla_\theta J(\theta) &= \frac{1}{2n} \left[ \nabla_\theta (X\theta - y)^T (X\theta - y) \right] \\ &= \frac{1}{2n} \left[ \nabla_\theta \left\{ \theta^T X^T X \theta - \underbrace{\theta^T X^T y}_{(X\theta)^T y} - \underbrace{y^T X \theta}_{y^T (X\theta)} + y^T y \right\} \right] \\ &= \frac{1}{2n} \left[ \nabla_\theta \left\{ \theta^T X^T X \theta - 2\theta^T X^T y + y^T y \right\} \right] \quad \left[ \begin{array}{l} (A+B)^T = A^T + B^T \\ (AB)^T = B^T A^T \end{array} \right] \\ &= \frac{1}{2n} \left[ (X^T X + (X^T X)^T) \theta - 2X^T y \right] \quad (x^T \theta = \theta^T x) \\ &= \frac{1}{n} \left[ X^T X \theta - X^T y \right] = \frac{1}{n} X^T (X\theta - y) \end{aligned}$$

$(X^T X)^T = \cancel{(X^T)^T X^T} X^T (X^T)^T = X^T X$

Normal Equations:

$$\nabla_{\theta} J(\theta) = 0.$$

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

Find  $\theta$  s.t.  $X^T X \theta = X^T y$ . (normal equations)

When  $(X^T X)^{-1}$  is ~~the~~ defined,

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

Mooze-Penrose pseudo inverse: For matrix  $A$ ,  $A^+$

Solution  $\theta = (X^T X)^+ X^T y.$

Proposition: Any  $\theta^*$  s.t.  $X^T X \theta^* = X^T y$  has the property

$$J(\theta^*) \leq J(\theta) \quad \forall \theta.$$

Proof: Let  $\theta$  be any values to the parameters.

$$J(\theta) = \|X\theta - y\|_2^2 \geq \|X\theta^* - y\|_2^2 = J(\theta^*)$$

$$= \|(X\theta - X\theta^*) + (X\theta^* - y)\|_2^2$$

$$= [(X\theta - X\theta^*) + (X\theta^* - y)]^T [(X\theta - X\theta^*) + (X\theta^* - y)]$$

$$= \underbrace{\|X\theta - X\theta^*\|_2^2}_{\geq 0} + \underbrace{(X\theta - X\theta^*)^T (X\theta^* - y)}_{J(\theta^*)} + \underbrace{(X\theta^* - y)^T (X\theta - X\theta^*)}_{J(\theta^*)} + \|X\theta^* - y\|_2^2$$

$$= \|X\theta - X\theta^*\|_2^2 + 2 \underbrace{(X\theta - X\theta^*)^T (X\theta^* - y)}_{J(\theta^*)} + \|X\theta^* - y\|_2^2$$

$$(X\theta - X\theta^*)^T (X\theta^* - y) = \theta^T X^T X \theta^* - \theta^{*T} X^T X \theta^* - \theta^T X^T y + \theta^{*T} X^T y.$$

$$= \theta^T [X^T X \theta^* - X^T y] - \theta^{*T} (X^T X \theta^* - X^T y).$$

$$= 0$$

= 0 because normal equations

$$\Rightarrow J(\theta) \geq J(\theta^*)$$

minimize  $J(\theta) = \| X\theta - y \|_2^2$

Each  $\theta$  gives you a linear combination  $a^{(j)}$ .

$X\theta$  - Linear span of  $a^{(j)}$ .

$$X = \begin{bmatrix} 1 & -x^{(1)T} \\ \vdots & \vdots \\ \vdots & -x^{(i)T} \\ \vdots & \vdots \end{bmatrix}$$

$a^{(j)}$  - jth column.

