

$$SSE = \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} |x_{ij} - (p+1)x_{ij}|^p + \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} |x_{ij} - (p+1)x_{ij}|^{p+1}$$

2/24/21

$$\frac{\partial SSE}{\partial \vec{w}} := \begin{bmatrix} \frac{\partial SSE}{\partial w_0} \\ \frac{\partial SSE}{\partial w_1} \\ \vdots \\ \frac{\partial SSE}{\partial w_p} \end{bmatrix} = \vec{0}_{p+1} \text{ and solve for } b_0, b_1, \dots, b_p$$

Let  $\vec{x} \in \mathbb{R}^n$ , let  $a \in \mathbb{R}$  be a constant with  $\vec{x} = \frac{\partial}{\partial \vec{x}} [a] = \vec{0}_n$

(0)

Let  $\vec{a} \in \mathbb{R}^n$  constant with  $\vec{x}$

$$(1) \quad \frac{\partial}{\partial \vec{x}} [\vec{a}^\top \vec{x}] = \begin{bmatrix} \frac{\partial}{\partial x_1} [a_1 x_1 + a_2 x_2 + \dots + a_n x_n] \\ \vdots \\ \frac{\partial}{\partial x_n} [a_1 x_1 + a_2 x_2 + \dots + a_n x_n] \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \vec{a} \neq \vec{a}^\top$$

let  $a, b \in \mathbb{R}$  constant with  $\vec{x}$

$$(2) \quad \frac{\partial}{\partial \vec{x}} [af(\vec{x}) + bg(\vec{x})] = \begin{bmatrix} \frac{\partial}{\partial x_1} [af(\vec{x}) + bg(\vec{x})] \\ \vdots \\ \frac{\partial}{\partial x_n} [af(\vec{x}) + bg(\vec{x})] \end{bmatrix} = \begin{bmatrix} a \frac{\partial}{\partial x_1} [f(\vec{x})] + b \frac{\partial}{\partial x_1} [g(\vec{x})] \\ \vdots \\ a \frac{\partial}{\partial x_n} [f(\vec{x})] + b \frac{\partial}{\partial x_n} [g(\vec{x})] \end{bmatrix}$$

$$= a \frac{\partial}{\partial \vec{x}} [f(\vec{x})] + b \frac{\partial}{\partial \vec{x}} [g(\vec{x})]$$

let  $A \in \mathbb{R}^{n \times n}$ , symmetric, constant with  $\vec{x}$

$$\frac{\partial}{\partial \vec{x}} [\vec{x}^T A \vec{x}], A \vec{x} = \begin{bmatrix} \leftarrow \vec{a}_1 \rightarrow \\ \leftarrow \vec{a}_2 \rightarrow \\ \vdots \\ \leftarrow \vec{a}_n \rightarrow \end{bmatrix} \begin{bmatrix} \uparrow \\ \vec{x} \\ \downarrow \end{bmatrix} = \begin{bmatrix} \vec{a}_1 \cdot \vec{x} \\ \vec{a}_2 \cdot \vec{x} \\ \vdots \\ \vec{a}_n \cdot \vec{x} \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n \end{bmatrix}$$

This scalar expression,  $\vec{x}^T A \vec{x}$  is called a "quadratic form" and it's a common expression and very well-studied

$$\vec{x}^T (A \vec{x}) = [x_1, x_2, \dots, x_n] \begin{bmatrix} \vec{a}_1 \cdot \vec{x} \\ \vec{a}_2 \cdot \vec{x} \\ \vdots \\ \vec{a}_n \cdot \vec{x} \end{bmatrix} = x_1 a_{11} \vec{x} + x_2 a_{21} \vec{x} + \dots + x_n a_{n1} \vec{x}$$

$$= x_1 (a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n) + x_2 (a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n) + \dots +$$

$$x_n (a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n)$$

$$\frac{\partial}{\partial x_1} \left[ \right] = 2a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n + a_{21}x_2 + \dots + a_{n1}x_n = 2a_{11}x_1 + 2a_{12}x_2 + \dots + 2a_{1n}x_n$$

$$= 2 \vec{a}_1 \cdot \vec{x}$$

$$\frac{\partial}{\partial x_2} \left[ \right] = a_{12}x_1 + 2a_{21}x_1 + 2a_{22}x_2 + \dots + a_{2n}x_n + \dots + a_{n2}x_n = 2a_{12}x_1 + 2a_{22}x_2 + \dots + 2a_{n2}x_n$$

$$= 2 \vec{a}_2 \cdot \vec{x}$$

$$\frac{\partial}{\partial \vec{x}} [\vec{x}^T A \vec{x}] = \begin{bmatrix} 2 \vec{a}_1 \cdot \vec{x} \\ 2 \vec{a}_2 \cdot \vec{x} \\ \vdots \\ 2 \vec{a}_n \cdot \vec{x} \end{bmatrix} = 2 A \vec{x}$$

*b is the optimal / "best"  
"w"*

$$\begin{aligned}
 & \text{sse} \quad \text{rule #2} \quad \text{by rule #0} \\
 & \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y} - 2 \vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}] \stackrel{\text{rule #2}}{=} \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y}] - 2 \frac{\partial}{\partial \vec{w}} [\vec{w}^T (X^T \vec{y})] + \frac{\partial}{\partial \vec{w}} [\vec{w}^T X^T X \vec{w}] \\
 & \downarrow \text{rule #1} \quad \text{rule #3} \quad \text{rule #3} \quad \text{symm.} \\
 & = -2 X^T \vec{y} + \frac{\partial}{\partial \vec{w}} [\vec{w}^T (X^T X) \vec{w}] = -2 X^T \vec{y} + 2 X^T X \vec{w} \xrightarrow{\text{set } \vec{w} = \vec{0}_{p+1} \text{ solve for } \vec{b}} \\
 & = 2(X^T X)^{-1} X^T \vec{y} = (X^T X)^{-1} X^T \vec{y} = \boxed{\vec{b} = (X^T X)^{-1} X^T \vec{y}} \quad \Rightarrow \hat{y}_* = g(\vec{x}_*) = \underbrace{\vec{x}_*^T \vec{b}}_{\text{predictions}}
 \end{aligned}$$

In order to compute the OLS coefficients (vector  $\vec{b}$ ), you need  $X^T X$ , a  $(p+1) \times (p+1)$  square matrix, to be invertible. Equivalently  $\text{rank}[X^T X] = p+1$  ie "full rank" ie all columns are linearly independent. Since there are  $n$  rows,  $\text{rank}[X^T X] = \text{rank}[X]$ , this means  $\text{rank}[X] = p+1$ , ie the columns of  $X$  are linearly indep.

$$X = \left[ \begin{array}{c|ccccc|c}
 1 & & & & & & & 1 \\
 \vdots & \uparrow & \uparrow & & & & & \vdots \\
 1 & x_{1,1} & x_{1,2} & \dots & x_{1,p} & & & 1 \\
 \vdots & \downarrow & \downarrow & & \downarrow & & & \downarrow \\
 \vdots & & & & & & & \vdots
 \end{array} \right]$$

feature measurements  
on all  $n$  subjects

If  $X$  is full rank that means.... there is no exact data duplication eg  $x_1$ : height measured in inches and  $x_2$ : height measured in centimeters. What if you do have a feature that is linearly dependent with the other features in  $X$ ? You just drop it. Then  $X$  will be full rank and you're good to estimate the OLS coefficients.

$$\vec{y} = \vec{\hat{y}} + \vec{e} \Rightarrow \vec{e} = \vec{y} - \vec{\hat{y}}, \text{ sse} = \sum_{i=1}^n e_i^2 = \vec{e}^T \vec{e}$$

$$\text{MSE} = \frac{1}{n-(p+1)} \text{SSE}, \text{RMSE} = \sqrt{\text{MSE}}, R^2 = \frac{\text{SST}-\text{SSE}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$$

$$= \frac{s_y^2 - s_e^2}{s_y^2} \text{ (same)}$$

You sometimes say the model has  $p+1$  "degrees of freedom" (ie. the number of parameters,  $w_0, w_1, \dots, w_p$ , is  $p+1$ ) and  $p+1 = \dim [{}^t \mathbf{0}] \text{sp} [\mathbf{X}]$ .