

$$SSE = \underbrace{\tilde{y}^T}_{1 \times n} \underbrace{\tilde{y}}_{n \times 1} - 2 \underbrace{\tilde{w}^T}_{1 \times (p+1)} \underbrace{X^T}_{(p+1) \times n} \underbrace{\tilde{y}}_{n \times 1} + \underbrace{\tilde{w}^T}_{1 \times (p+1)} \underbrace{X^T}_{(p+1) \times n} \underbrace{X}_{n \times (p+1)} \underbrace{\tilde{w}}_{(p+1) \times 1}$$

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

let $\tilde{x} \in \mathbb{R}^n$. let $a \in \mathbb{R}$ be a constant wrt \tilde{x} . $\Rightarrow \frac{\partial}{\partial \tilde{x}}[a] = \vec{0}_n$ (0)

let $\vec{a} \in \mathbb{R}^n$ constant wrt \tilde{x}

$$\frac{\partial}{\partial \tilde{x}} [\vec{a}^T \tilde{x}] = \begin{bmatrix} \frac{\partial}{\partial x_1} [a_1 x_1 + a_2 x_2 + \dots + a_n x_n] \\ \vdots \\ \frac{\partial}{\partial x_n} [a_1 x_1 + a_2 x_2 + \dots + a_n x_n] \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \vec{a} \neq \vec{a}^T$$

let $a, b \in \mathbb{R}$ constants wrt \tilde{x}

$$\frac{\partial}{\partial \tilde{x}} [a f(\tilde{x}) + b g(\tilde{x})] = \begin{bmatrix} \frac{\partial}{\partial x_1} [a f(\tilde{x}) + b g(\tilde{x})] \\ \vdots \\ \frac{\partial}{\partial x_n} [a f(\tilde{x}) + b g(\tilde{x})] \end{bmatrix} = \begin{bmatrix} a \frac{\partial}{\partial x_1} [f(\tilde{x})] + b \frac{\partial}{\partial x_1} [g(\tilde{x})] \\ \vdots \\ a \frac{\partial}{\partial x_n} [f(\tilde{x})] + b \frac{\partial}{\partial x_n} [g(\tilde{x})] \end{bmatrix}$$

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

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

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

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

$$\tilde{x}^T (A \tilde{x}) = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} \vec{a}_1^T \tilde{x} \\ \vec{a}_2^T \tilde{x} \\ \vdots \\ \vec{a}_n^T \tilde{x} \end{bmatrix} = x_1 \vec{a}_1^T \tilde{x} + x_2 \vec{a}_2^T \tilde{x} + \dots + x_n \vec{a}_n^T \tilde{x}$$

$$= x_1 (a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n) + x_2 (a_{21} x_1 + a_{22} x_2 + \dots + a_{2n} x_n) + \dots + x_n (a_{n1} x_1 + a_{n2} x_2 + \dots + a_{nn} x_n)$$

$$\frac{\partial}{\partial x_1} \left\{ \begin{array}{l} \text{---} \end{array} \right\} = 2a_{11} x_1 + 2a_{12} x_2 + \dots + 2a_{1n} x_n = 2 \vec{a}_1^T \tilde{x}$$

$$\frac{\partial}{\partial x_2} \left\{ \begin{array}{l} \text{---} \end{array} \right\} = 2a_{21} x_1 + 2a_{22} x_2 + \dots + 2a_{2n} x_n = 2 \vec{a}_2^T \tilde{x}$$

$$\vdots$$

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

$$\begin{aligned} \frac{\partial}{\partial \tilde{w}} [\tilde{y}^T \tilde{y} - 2 \tilde{w}^T X^T \tilde{y} + \tilde{w}^T X^T X \tilde{w}] &\stackrel{\text{rule \#2}}{=} \frac{\partial}{\partial \tilde{w}} [\cancel{\tilde{y}^T \tilde{y}}] - 2 \frac{\partial}{\partial \tilde{w}} [\tilde{w}^T (X^T \tilde{y})] + \frac{\partial}{\partial \tilde{w}} [\tilde{w}^T X^T X \tilde{w}] \\ &\stackrel{\text{rule \#1}}{=} -2 X^T \tilde{y} + \frac{\partial}{\partial \tilde{w}} [\tilde{w}^T (X^T X) \tilde{w}] \stackrel{\text{rule \#3}}{=} -\cancel{2} X^T \tilde{y} + \cancel{2} X^T X \tilde{w} \stackrel{\text{set } \vec{0}_{p+1}}{=} \vec{0}_{p+1} \quad \text{and solve for } \vec{b} \\ &\Rightarrow (X^T X)^{-1} X^T X \tilde{w} = (X^T X)^{-1} X^T \tilde{y} \Rightarrow \boxed{\vec{b} = (X^T X)^{-1} X^T \tilde{y}} \Rightarrow \hat{y}_i = g(\hat{x}_i) = \tilde{x}_i^T \vec{b} \end{aligned}$$

predictions

In order to compute the OLS coefficients (vector b), you need $X^T X$, a $(p+1) \times (p+1)$ square matrix, to be invertible. Equivalently, $\text{rank}[X^T X] = p + 1$ i.e. "full rank" i.e. all columns of $X^T X$ are linearly independent. Since there's a thm: $\text{rank}[X^T X] = \text{rank}[X]$, this means $\text{rank}[X] = p + 1$, i.e. the columns of X are linearly indep.

$$X = \begin{bmatrix} \vdots & \uparrow & \uparrow & \uparrow \\ & \tilde{x}_{\cdot 1} & \tilde{x}_{\cdot 2} & \dots & \tilde{x}_{\cdot p} \\ & \downarrow & \downarrow & & \downarrow \\ \vdots & & & & \end{bmatrix}$$

feature measurements
on all n subjects

If X is full rank that means there is no exact data duplication e.g. 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.

$$\tilde{y} = \tilde{\hat{y}} + \tilde{e} \Rightarrow \tilde{e} = \tilde{y} - \tilde{\hat{y}}, \quad SSE = \sum_{i=1}^n e_i^2 = \tilde{e}^T \tilde{e}$$

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

you sometimes say the model has $p+1$ "degrees of freedom" (i.e. the number of parameters, w_0, w_1, \dots, w_p , is $p + 1$) and $p + 1 = \text{dim}[\text{colsp}[X]]$.