## 2.3 Projection and geometric view

**[Proj 1]** Definition of projection & least squares fit

**[Proj 2]** (1) Properties of the projection map $\Rightarrow$ induces a matrix

$\left\{ \begin{array}{l} \text{(1) Uniqueness} \quad \text{(2) Linearity} \quad \text{(3) Idempotent} \\[1em] \text{(4) Map} \Rightarrow \text{Matrix: specific form} \quad P_{Col(x)} = x(x^Tx)^{-1}x^T \\[1em] \text{(5) Relationship with OLS} \quad \text{(6) } I-P \text{ is also projection} \end{array} \right.$

(2) Properties of the projection matrix in OLS

(Use the properties to prove OLS conclusions)

# 3. Statistical properties for OLS

**[Stat 1]** From only moment structure $\quad E(Y \mid x) = x^T\beta$

$\Rightarrow$ Additive model $\quad Y = x^T\beta + \epsilon \quad$ (Assumptions on $\epsilon$)

**[Stat 2]** Properties: Mean and variance of $\hat{\beta}$ and residuals

12) Residual

Mean $\quad E(\hat{R} \mid x) = E(Y - \hat{Y} \mid x) \qquad \underset{n \times 1}{\hat{Y}} = P\underset{n \times 1}{Y}$

$$= (I-P) \, E(Y \mid x)$$

$$= (I-P) \, X\beta \quad = 0 \qquad (I-P) \, X = 0$$

Covariance $\quad cov(\underset{n \times 1}{\hat{R}} \mid x) = cov(Y - \hat{Y} \mid x)$

$$= E\left[ (Y-\hat{Y})(Y-\hat{Y})^T \mid x \right]$$

$$= E\left\{ (I-P) \, Y \, \{(I-P)Y\}^T \mid x \right\}$$

$$= (I-P) \, E(YY^T \mid x) \, (I-P)^T \quad ①$$

By $Y = X\beta + \epsilon,$ ✫ $E(YY^T \mid x) = E\left\{ (X\beta+\epsilon)(X\beta+\epsilon)^T \mid x \right\}$

$$= E\left\{ X\beta\beta^T X^T + \epsilon(X\beta)^T + X\beta\epsilon + \epsilon\epsilon^T \mid x \right\}$$

$$= X\beta\beta^T X^T + \sigma^2 I$$

① $= (I-P) \left\{ X\beta\beta^T X^T + E(\epsilon \mid x)(X\beta)^T + X\beta \, E(\epsilon \mid x) + E(\epsilon\epsilon^T \mid x) \right\} (I-P)^T$

$$= (I-P) \, \sigma^2 I \, (I-P)^T$$

$\qquad\qquad\qquad\qquad\qquad$ By $\qquad (I-P) X = 0$

$$= \sigma^2 (I-P)$$

$\qquad\qquad\qquad\qquad\qquad\qquad E(\epsilon \mid x) = 0$

$\qquad\qquad\qquad\qquad\qquad\qquad E(\epsilon\epsilon^T \mid x) = \sigma^2 I$

$cov(\hat{R} \mid x) = \sigma^2 (I-P)$

## Residuals sum of squares (RSS)

$$RSS = \| Y - \hat{Y} \|^2 \qquad \hat{R} = Y - \hat{Y}$$

$$= \| (I-P) Y \|^2 \qquad Y - \hat{Y} = (I-P) Y$$

$$= \{ (I-P) Y \}^T (I-P) Y$$

$$= Y^T (I-P)^T (I-P) Y \qquad (I-P)^T (I-P) = I-P$$

$$= Y^T (I-P) Y$$

This is a quadratic form in $Y$

mean RSS $\qquad E( RSS \mid x )$

$$= E\{ Y^T (I-P) Y \mid x \} \qquad (E1)$$

$$= E[ tr\{ (I-P) Y Y^T \} \mid x] \qquad (E2)$$

$$= \sigma^2 \, tr(I-P) \qquad (E3)$$

$$= \sigma^2 (n-p) \qquad ( \text{By property on } P )$$

(Sep. 21st Notes)

$$\Rightarrow E\left( \frac{RSS}{n-p} \right) = \sigma^2$$

$$\Rightarrow \frac{RSS}{n-p} \text{ is an unbiased estimator of } \sigma^2$$

$(E1) \Rightarrow (E2)$

random    fixed        random

$$Y^T \; (I-P) \; Y \qquad \text{scalar}$$

$1 \times n \quad n \times n \quad n \times 1$

$$= tr\{ \; Y^T (I-P) Y \}$$

$$= tr\{ \; (I-P) \, Y Y^T \}$$

fixed    random

$$(E2) = \text{tr} \{ (I-P) \; \underline{E(YY^T | X)} \}$$

$$= \text{tr} \{ (I-P) ( \underwave{X\beta\beta^T X^T} + \sigma^2 I ) \} \qquad \text{By } (I-P)X = 0$$

$$= \text{tr} \{ (I-P) \sigma^2 I \}$$

**Remark :** ① $X \in \mathbb{R}^{n \times p}$ $\qquad \sigma^2 = E\left( \dfrac{RSS}{n-p} \right)$

② $p$ covariates and 1 intercept $\qquad X \in \mathbb{R}^{n \times (p+1)} \qquad \sigma^2 = E\left( \dfrac{RSS}{n-(p+1)} \right)$

**Discussions on the SLR** $\qquad p = 2 \qquad$ 1 intercept + 1 covariate

$$Y = \alpha + \beta X + \epsilon \qquad (\alpha, \beta \in \mathbb{R})$$

**(1)** OLS estimates: unbiased $\qquad E(\hat{\alpha} | X) = \alpha$

$$E(\hat{\beta} | X) = \beta$$

Follows from MLR conclusion.

**Exercise :** $\checkmark \begin{cases} \hat{\beta} = \dfrac{\widehat{\text{cov}}(X, Y)}{\widehat{\text{var}}(X)} \\ \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x} \end{cases} \Rightarrow \begin{cases} E(\hat{\alpha} | X) = \alpha \\ E(\hat{\beta} | X) = \beta \end{cases}$

OLS estimate covariance $\text{cov}\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \sigma^2 (X^TX)^{-1}$

$$\underset{n\times 2}{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

Exercise 1. Plug-in $X$ into $\sigma^2(X^TX)^{-1}$

2. $\hat{\beta}, \hat{\alpha}$ formula    ( after lecture )

Hint: $\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad-bc}\begin{pmatrix} d & -b \\ -b & a \end{pmatrix}$

$$\text{var}(\hat{\beta}|x) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$\text{var}(\hat{\alpha}|x) = \sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n\sum_{i=1}^{n}(x_i-\bar{x})^2} = \sigma^2 \times \frac{\bar{x}^2 + \frac{\sum_{i=1}^{n}(x_i-\bar{x})^2}{n}}{\sum_{i=1}^{n}(x_i-\bar{x})^2}$$

$$= \sigma^2 \left\{ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i-\bar{x})^2} \right\}$$

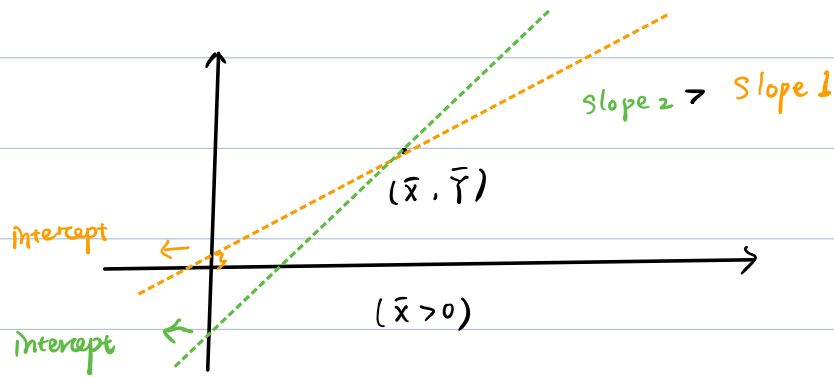$\geq 0$    equality holds if $\bar{x}=0$

$\text{Var}(\hat{\alpha}|x)$ is minimized if $\bar{x}=0$.

$$\text{cov}(\hat{\alpha},\hat{\beta}|x) = -\sigma^2 \frac{\bar{x}}{\sum_{i=1}^{n}(x_i-\bar{x})^2}$$

If $\bar{x}=0$, $\hat{\alpha}$ and $\hat{\beta}$ are uncorrelated.

If $\bar{x}>0$, covariance is negative.

Slope 2 > Slope 1

$(\bar{x}, \bar{Y})$

Intercept

Intercept

$(\bar{x} > 0)$

OLS fit always
pass $(\bar{x}, \bar{Y})$ data center.

## (2) Residuals

$$\hat{R}_i = Y_i - \hat{Y}_i$$

$$= Y_i - \hat{\alpha} - \hat{\beta} X_i$$

$$= Y_i - \bar{Y} - \hat{\beta}(X_i - \bar{x})$$

$$= (\alpha + \beta X_i + \epsilon_i) - (\alpha + \beta \bar{x} + \bar{\epsilon}) - \hat{\beta}(X_i - \bar{x})$$

$$= (\beta - \hat{\beta})(X_i - \bar{x}) + \epsilon_i - \bar{\epsilon} \quad ✳$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x}$$

$$Y_i = \alpha + \beta X_i + \epsilon_i$$
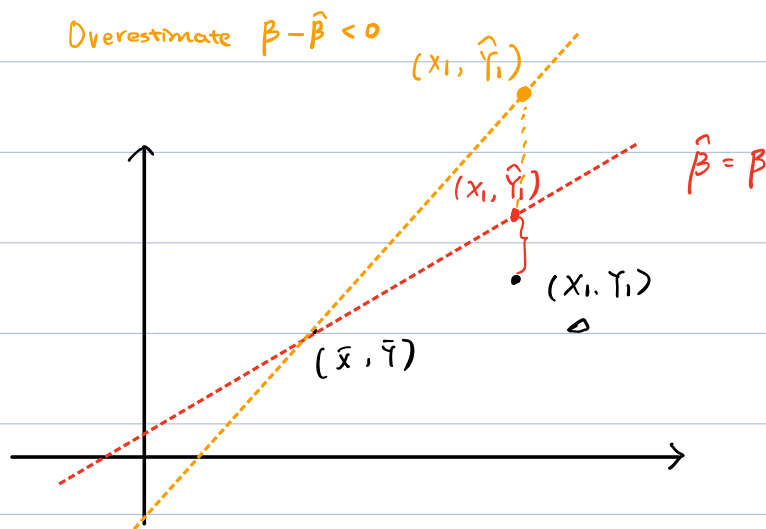$$\bar{Y} = \alpha + \beta \bar{x} + \bar{\epsilon}$$
$$\bar{\epsilon} = \frac{1}{n} \sum_{i=1}^{n} \epsilon_i$$

## Interpretation

The residuals $\hat{R}_i$ are not only influenced by errors $\epsilon$
but also $(\beta - \hat{\beta})(X_i - \bar{x})$, i.e. how well we
recover the true slope $\beta$.

( If $\hat{\beta} = \beta$, view $\epsilon_i - \bar{\epsilon}$ as baseline errors.)

Overestimate $\beta - \hat{\beta} < 0$

$(x_1, \hat{Y_1})$

$\hat{\beta} = \beta$

$(x_1, \hat{Y_1})$

$(X_1, Y_1)$

$(\bar{x}, \bar{Y})$

$\hat{R}_{1, \hat{\beta}=\beta} = Y_1 - \hat{Y}_{1, \hat{\beta}=\beta}$
$\left.\begin{array}{l}\\ \\ \end{array}\right\} \Rightarrow \hat{R}_{2, \hat{\beta}>\beta}$ is more negative than $\hat{R}_{1, \hat{\beta}=\beta}$
$\hat{R}_{2, \hat{\beta}>\beta} = Y_1 - \hat{Y}_{1, \hat{\beta}>\beta}$

1. Overestimate $\beta$: $\quad \beta - \hat{\beta} < 0 \qquad$ as $\hat{\beta} \to +\infty$

For $i$ with $X_i - \bar{x} > 0$ (right of mean) $\hat{R_i} \searrow -\infty$

$X_i - \bar{x} < 0$ (left of mean) $\hat{R_i} \nearrow +\infty$

2. Underestimate $\beta \qquad \beta - \hat{\beta} > 0$

For $i$ with $X_i - \bar{x} > 0$, $\quad \hat{R_i} \nearrow +\infty$

$X_i - \bar{x} < 0$, $\quad \hat{R_i} \searrow -\infty$

Covariance of residuals: $\quad \text{cov}(\hat{R} | x) = \sigma^2 (I - P)$
$\qquad\qquad\qquad\qquad\qquad\quad n\times 1$

$\qquad\qquad\qquad\qquad\qquad P = X(X^TX)^{-1}X^T$

① Plug in $\quad \underset{n\times 2}{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \checkmark$

② Direct calculate through $\quad \hat{R_i} = Y_i - \hat{\alpha} - \hat{\beta}X_i$

**Exercise:** $\quad var(\hat{R}_i \mid X) = \sigma^2 \times (I - P)_{(i,j)} = \sigma^2 \times (1 - P_{ii})$

$$P = X(X^TX)^{-1}X^T = \frac{1}{\sum\limits_{i=1}^{n}(x_i-\bar{x})^2} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \frac{\sum\limits_{i=1}^{n}x_i^2}{n} - x_1\bar{x} & \frac{\sum\limits_{i=1}^{n}x_i}{n} - x_2\bar{x} & \cdots & \frac{\sum\limits_{i=1}^{n}x_i^2}{n} - x_n\bar{x} \\ -\bar{x} + x_1 & \cdots\cdots & & -\bar{x} + x \end{pmatrix}$$

$$\underbrace{\phantom{XXXXX}}_{n \times 2} \qquad\qquad \underbrace{\phantom{XXXXXX}}_{2 \times n}$$

**$i$-th diagonal** $\quad P_{ii} = \dfrac{1}{\sum\limits_{i=1}^{n}(x_i-\bar{x})^2}\left( \dfrac{1}{n}\sum\limits_{i=1}^{n}x_i^2 - 2x_i\bar{x} + x_i^2 \right)$

$$var(\hat{R}_i \mid X) = \sigma^2 \times (1 - P_{ii})$$

$$= \sigma^2 \times \left( 1 - \frac{1}{\sum\limits_{i=1}^{n}(x_i-\bar{x})^2}\left( \frac{1}{n}\left(\sum\limits_{i=1}^{n}x_i^2 - n\bar{x}^2\right) + (x_j-\bar{x})^2 \right\} \right)$$

$$= \sigma^2 \times \left( 1 - \frac{1}{n} - \underbrace{\frac{(x_j-\bar{x})^2}{\sum\limits_{i=1}^{n}(x_i-\bar{x})^2}}_{\geqslant 0} \right)$$

$$\leqslant \sigma^2 \times \left(1 - \frac{1}{n}\right) \underset{\sim}{\leqslant} \sigma^2 = var(\epsilon_i \mid X) \qquad n \geqslant 2$$

$\hat{R}_i$ has less variability than $\epsilon_i$

## Optimality of OLS ⇒ Gauss–Markov Theorem

Why OLS estimates $\hat{\beta}$ not other estimates?

Nice properties: Unbiasedness $E(\hat{\beta} \mid x) = \beta$

$$E\left(\frac{RSS}{n-p} \mid x\right) = 6^2$$

△ Goal: For a reasonable class of estimates of $\beta$, OLS $\hat{\beta_j}$

is an unbiased estimate of $\beta_j$ with the smallest variance.

△
$$\beta_j = \underset{1 \times p}{e_j^T} \underset{p \times 1}{\beta} \qquad e_j = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \rightarrow j\text{th position indicator}$$

$e_j^T \beta$    unbiased    small variance

We can generalize this to $\underset{p \times 1}{c^T \beta}$ for any given $c \in \mathbb{R}^P$.

## ▲ Gauss–Markov Theorem

When columns of $x$ are linearly independent $(\hat{\beta} = (x^T x)^{-1} x^T y)$

among the class of linear unbiased estimates of $c^T \beta$,

$c^T \hat{\beta}$ is the unique estimate with the minimum variance.

△ We say $c^T \hat{\beta}$ is the <u>best</u> <u>linear</u> <u>unbiased</u> estimate of $c^T \beta$.

<span style="color:red">BLUE</span>

Linear unbiased estimate: any estimate in the form

$$m^T y \qquad (Y \in \mathbb{R}^n, \; m \in \mathbb{R}^{n \times P}) \quad E(m^T Y \mid x) = \beta$$