

Gauss-Markov Theorem (Reading: F, 2.6)

px1

Q: why is $\hat{\beta} = (X^T X)^{-1} X^T Y$, OLS estimator a good estimator?

criterion for finding estimator

➢ it results from orthogonal projection, makes sense geometrically

LS criterion

➢ (FYI) if $\varepsilon \sim N(0, \sigma^2 I)$, $\hat{\beta}$ is the maximum likelihood estimator (exercise)

multivariate normal

➢ Gauss-Markov thm states $\hat{\beta}$ is BLUE ("Best" Linear Unbiased Estimator)

statistical property

function of parameters

meaning

parameter

- estimable function: a linear combination of the parameters $\psi = c^T \beta$, where c is a known vector, is estimable if and only if there exists a linear combination of y_i 's, i.e., $a^T Y$, such that $E(a^T Y) = c^T \beta$, $\forall \beta \in \mathbb{R}^p$ ($\Rightarrow a^T Y$: an unbiased estimator of $c^T \beta$)

➢ Examples of estimable function

LNP.2-10.
EX2

$$\begin{bmatrix} z_1 \\ z_m \\ w_1 \\ w_n \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + \varepsilon$$

$$\hat{\beta} = \begin{bmatrix} \bar{z} \\ \bar{w} \end{bmatrix}$$

$$\begin{aligned} \text{Say } \psi &= u_1 - u_2 \\ &= [1 \ -1]^T \beta \\ &\uparrow C^T \end{aligned}$$

$$(i) \alpha^T = [1, 0, \dots, 0, -1, 0, \dots, 0]$$

$$E(\alpha^T Y) = E(z_1 - w_1) = u_1 - u_2 \quad \text{C.P.}$$

$$(ii) \alpha^T = [y_1, \dots, y_m, -y_n, \dots, -y_n]$$

$$E(\alpha^T Y) = E(\bar{z} - \bar{w}) = u_1 - u_2$$

$$\text{But, } \text{Var}(z_1 - w_1) > \text{Var}(\bar{z} - \bar{w})$$

(Note $\alpha^T X = C^T$, exercise)

LNP.4-4,
(future
lecture)

prediction:

$$x_0^T = (g_{01}, \dots, g_{0p})$$

$$\text{predict } E(Y|_{x_0}) = x_0^T \beta = \psi$$

$$C^T$$

$$\begin{aligned} \text{predictor: } \hat{x}_0^T \hat{\beta} &= x_0^T (X^T X)^{-1} X^T Y \\ E(x_0^T \hat{\beta}) &= x_0^T \beta \end{aligned}$$

$$\begin{aligned} \because E(\hat{\beta}) &= \beta \\ \therefore E(C^T \hat{\beta}) &= C^T \beta \end{aligned}$$

If not hold,
then $\hat{\beta}$ have
infinite many
solutions.
(future lecture)

➢ If X is of full rank, all $c^T \beta$'s, $\forall c \in \mathbb{R}^p$, are estimable.

- **Theorem.** For a linear model $Y = X\beta + \varepsilon$, suppose ① $E(\varepsilon) = 0$ (i.e., the structural part of the model, $E(Y) = X\beta$, is correct) and ② $\text{Var}(\varepsilon) = \sigma^2 I$ ($\sigma^2 < \infty$). Let ③ $\psi = c^T \beta$ be an estimable function. Then, in the class of all unbiased linear estimators of ψ , $\hat{\psi} = c^T \hat{\beta}$, where $\hat{\beta}$ is the OLS estimator, has the minimum variance and is unique.

proof: Let $a^T Y$ be an unbiased estimator of $C^T \beta$.

estimators with the form $a^T Y$

$$E(a^T Y) = a^T X \beta = c^T \beta \Rightarrow (a^T X - c^T) \beta = 0, \forall \beta \in \mathbb{R}^p \Rightarrow a^T X - c^T = 0 \Rightarrow a^T X = c^T \quad (*)$$

$$\text{Var}(c^T \hat{\beta}) = E^*(c^T \hat{\beta} \hat{\beta}^T c) = c^T E^*(\hat{\beta} \hat{\beta}^T) c = c^T (X^T X)^{-1} \sigma^2 c = a^T X (X^T X)^{-1} X^T a \cdot \sigma^2 = (a^T H a) \sigma^2$$

$$\text{Var}(a^T Y) = E^*(a^T Y Y^T a) = a^T E^*(Y Y^T) a = a^T \sigma^2 I a$$

$$\|a\|^2 = \sigma^2 (a^T a)$$

$$\text{by (*)} \quad \|H\|^2 \quad \text{cf.} \quad \|H\|^2$$

• minimum variance

$$\text{Var}(a^T Y) - \text{Var}(c^T \hat{\beta}) = \sigma^2 a^T a - \sigma^2 a^T H a = \sigma^2 a^T (I - H) a = \sigma^2 a^T (I - H)^T (I - H) a$$

$$= \sigma^2 \| (I - H) a \|^2 \geq 0 \quad (**)$$

in (**), " $=$ " holds iff $a \in \Omega \Leftrightarrow a = X\lambda \Rightarrow$ from (*), $c^T = a^T X = (X\lambda)^T X = \lambda^T X^T X$

$$\Rightarrow C^T \hat{\beta} = \underbrace{\lambda^T (X^T X)}_{C^T} \underbrace{(X^T X)^{-1} X^T Y}_{\hat{\beta}} = \underbrace{\lambda^T X^T Y}_{a^T} = a^T Y$$

Q: Under what assumptions?

Ans. Gauss-Markov conditions

① ② ③

• Implications of the theorem:

➢ the theorem shows that the OLS estimator is a "good" choice Note, OLS est' or $\hat{\beta}$ is MLE if $\varepsilon \stackrel{iid}{\sim} N(0, \sigma^2)$

➢ Note: the theorem does not require normally distributed ε

➢ there may exist non-linear/biased estimators that are "better" ($\text{MSE} = \text{Var} + \text{Bias}^2$)

➢ if some assumptions are not true, there will be better estimators

- Situations where estimators other than the OLS estimator should be considered
 - when ϵ are correlated or have unequal variance, $(\text{Var}(\epsilon)) = \sigma^2 I$ is violated.
 - Q:** what will be seen in data analysis? how to check? diagnostic using residuals
 - use generalized least square estimator **heavy-tailed**, e.g., **Cauchy distribution**.
 - when error distribution is long-tailed (e.g., there exists outliers or $\sigma^2 = \infty$). **Q:** what will be seen in data analysis? how to check?
 - use robust estimators, typically not linear in Y . $\hat{\beta} \cdot \text{a linear function of } Y$
 - when the predictors are highly correlated, i.e., collinear (non-singular $X^T X$) is "partially" violated. **Q:** what will be seen in data analysis? how to check? collinearity detection or PCA
 - use biased estimators such as ridge regression
 - when some important predictors are not included in fitted model ($E(\epsilon) = 0$ is violated, e.g., true model (M_T): $E(Y) = X_1 \beta_1 + X_2 \beta_2$, and fitted model (M_F): $E(Y) = X_1 \beta_1$). **Q:** what will be seen in data analysis? how to check?) $\Rightarrow [X_1 \ X_2] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$
 - Under M_F ,** $\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$ diagnostic using residuals
 - Under M_T ,** $E(\hat{\beta}_1) = (X_1^T X_1)^{-1} X_1^T (X_1 \beta_1 + X_2 \beta_2) = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2 \beta_2$

(Note: From the viewpoint of data analysis, it implies you should examine whether these situations exist in your data when performing regression analysis. It's interesting that the theorem not only shows us how good OLS is, but also indicates what conditions should be concerned and checked in data analysis.)

Estimation of a subset of parameters

p. 3-12

 $\beta_1 \leftarrow$ of main interest

- Consider the model $\mathbf{Y} = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 + \epsilon$ with constant variance. Let $\mathbf{X} = [\mathbf{X}_1 \ \mathbf{X}_2]$ and $\underline{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$. Then, $\mathbf{Y} = \mathbf{X} \underline{\beta} + \epsilon$ uncorrelated (*)
- Suppose $\hat{\underline{\beta}}$ is the OLS estimator of $\underline{\beta}$ under (*). Then, $\mathbf{X}^T = \begin{bmatrix} \mathbf{X}_1^T \\ \mathbf{X}_2^T \end{bmatrix} \times [\mathbf{X}_1 \ \mathbf{X}_2] = \mathbf{X}$

$$\hat{\underline{\beta}} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{bmatrix} \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{X}_2 \\ \mathbf{X}_2^T \mathbf{X}_1 & \mathbf{X}_2^T \mathbf{X}_2 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}_1^T \mathbf{Y} \\ \mathbf{X}_2^T \mathbf{Y} \end{bmatrix}$$
Check orthogonality (future lecture) cf.
 - Clearly, $\hat{\beta}_1 \neq (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{Y}$ in general (equal if $\mathbf{X}_1^T \mathbf{X}_2 = 0$) (□)
 - Note. $(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{Y}$ is the OLS estimator of β_1 under $\mathbf{Y} = \mathbf{X}_1 \beta_1 + \epsilon$.
- Q:** What is the close form of $\hat{\beta}_1$? (useful when β_2 are nuisance parameters, which are not of the main interest but required to be in the model)
- Normal equation $\mathbf{X}^T \mathbf{X} \hat{\underline{\beta}} = \mathbf{X}^T \mathbf{Y}$: LN p. 3-5 check (*) in LN p. 3-11
 - $(\mathbf{X}_1^T \mathbf{X}_1) \hat{\beta}_1 + (\mathbf{X}_1^T \mathbf{X}_2) \hat{\beta}_2 = \mathbf{X}_1^T \mathbf{Y}$ substituting
 - $(\mathbf{X}_2^T \mathbf{X}_1) \hat{\beta}_1 + (\mathbf{X}_2^T \mathbf{X}_2) \hat{\beta}_2 = \mathbf{X}_2^T \mathbf{Y} \Rightarrow \hat{\beta}_2 = (\mathbf{X}_2^T \mathbf{X}_2)^{-1} [\mathbf{X}_2^T \mathbf{Y} - (\mathbf{X}_2^T \mathbf{X}_1) \hat{\beta}_1]$
- Do the substitution, and get

$$(\mathbf{X}_1^T \mathbf{X}_1) \hat{\beta}_1 + (\mathbf{X}_1^T \mathbf{X}_2) (\mathbf{X}_2^T \mathbf{X}_2)^{-1} [\mathbf{X}_2^T \mathbf{Y} - (\mathbf{X}_2^T \mathbf{X}_1) \hat{\beta}_1] = \mathbf{X}_1^T \mathbf{Y} \quad (\Delta)$$

$$\Rightarrow [\mathbf{X}_1^T \mathbf{X}_1 - \mathbf{X}_1^T \mathbf{X}_2 (\mathbf{X}_2^T \mathbf{X}_2)^{-1} \mathbf{X}_2^T \mathbf{X}_1] \hat{\beta}_1 = [\mathbf{X}_1^T - \mathbf{X}_1^T \mathbf{X}_2 (\mathbf{X}_2^T \mathbf{X}_2)^{-1} \mathbf{X}_2^T] \mathbf{Y}$$

• Let $\underline{H}_2 = \underline{X}_2(\underline{X}_2^T \underline{X}_2)^{-1} \underline{X}_2^T$ be the hat matrix of \underline{X}_2 . $\underline{\mathcal{Y}} = \underline{P}_{\Omega_2} \underline{\mathcal{Y}} + \underline{P}_{\Omega_2^\perp} \underline{\mathcal{Y}}$ p. 3-13
 $= \underline{H}_2 \underline{\mathcal{Y}} + (\underline{I} - \underline{H}_2) \underline{\mathcal{Y}}$

- The matrix \underline{H}_2 is the orthogonal projection matrix onto $\underline{\Omega}_2 \equiv \text{span}\{\underline{X}_2\}$.
- The matrix $\underline{I} - \underline{H}_2$ is the orthogonal projection matrix onto $\underline{\Omega}_2^\perp$.

• Then, we have $\underline{X}_1^T \xrightarrow{(I-H_2)^T} \underline{R}^n = \underline{\Omega}_2 \oplus \underline{\Omega}_2^\perp$
 $(\Delta) \Rightarrow [\underline{X}_1^T (\underline{I} - \underline{H}_2) \underline{X}_1] \hat{\beta}_1 = \underline{X}_1^T (\underline{I} - \underline{H}_2) \underline{Y} = \text{span}\{\underline{x}_1, \underline{x}_2\} \oplus \text{span}\{\underline{x}_1, \underline{x}_2\}^\perp$
 $\Rightarrow [\underline{X}_1^T (\underline{I} - \underline{H}_2)^T (\underline{I} - \underline{H}_2) \underline{X}_1] \hat{\beta}_1 = \underline{X}_1^T (\underline{I} - \underline{H}_2)^T (\underline{I} - \underline{H}_2) \underline{Y}$
 $\xrightarrow{P \times P} (\underline{X}_1^T \underline{X}_1) \hat{\beta}_1 = \underline{X}_1^T \underline{\tilde{Y}} \quad (\Leftarrow \text{normal equation for } \hat{\beta}_1) \xrightarrow{\text{cf.}} \text{normal equation for } \hat{\beta}_1 \text{ (LNp.3-12)}$

Check Ex.1 (LNp.3-7) where $\underline{\tilde{X}}_1 = \underline{(I - H_2)} \underline{X}_1 = \underline{X}_1 - \underline{H}_2 \underline{X}_1$
 $\rightarrow \text{cov}(\hat{\beta}_1) = \sigma^2 (\underline{X}_1^T \underline{X}_1)^{-1}$

• From the normal equation for $\hat{\beta}_1$, we get $\hat{\beta}_1 = (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T \underline{\tilde{Y}} = (\underline{X}_1^T \underline{X}_1)^{-1} \underline{X}_1^T \underline{Y}$, which is the OLS estimator of the linear model

$\underline{\tilde{Y}} = \underline{X}_1^T \hat{\beta}_1 + \underline{\tilde{\epsilon}}$, where $\underline{\tilde{\epsilon}} = (\underline{I} - \underline{H}_2) \underline{\epsilon}$.

model (□) in LNp.3-12

Estimating σ^2

$$Y = X\beta + \epsilon = X\hat{\beta} + \hat{\epsilon}$$

$$\sigma^2 = \text{Var}(\epsilon), \quad \hat{\epsilon} : \text{surrogate of } \epsilon$$

p. 3-14

- Q: Which part in data contains information about σ^2 ?
- Q: What is a suitable function (statistics) of $\hat{\epsilon}$ for estimating σ^2 ?

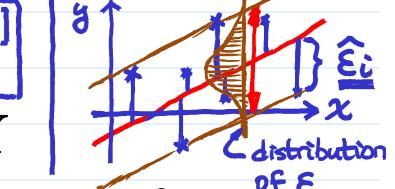
Ans: residuals ($\hat{\epsilon}$)

$\boxed{\text{RSS}} = \sum_{i=1}^n (\hat{\epsilon}_i - \bar{\hat{\epsilon}})^2 = \hat{\epsilon}^T \hat{\epsilon}$ length² of $\hat{\epsilon}$

$\rightarrow = 0, \text{ if the model space } \Omega \text{ contains } \mathbb{I}$

$E(\hat{\epsilon}^T \hat{\epsilon}) = E[(I - H)\mathbf{Y}] = (I - H)X\beta = \Omega$

$= [Y^T(I - H)^T][(I - H)Y] = Y^T(I - H)Y$ a quadratic function of y_i 's



$E(\hat{\epsilon}^T \hat{\epsilon}) = (n-p)\sigma^2$, where $n = \# \text{ of observations}$, $p = \# \text{ of parameters in } \beta$

(Note: $\hat{\epsilon}$ is located in the $(n-p)$ -dim residual space Ω^\perp)

Hint. $E(Y^T A Y) = \text{trace}[A \cdot \text{cov}(Y)] + [E(Y)]^T A [E(Y)]$

Let $A = I - H$

$$E(Y^T(I - H)Y) \quad \text{trace}[(I - H)\sigma^2 I]$$

$$E(\hat{\epsilon}^T \hat{\epsilon}) \quad (n-p)\sigma^2$$

$$(X\beta)^T(I - H)(X\beta)$$

$$\beta^T X^T(I - H)X\beta = 0$$

① trace of matrix = sum of eigenvalues
 ② $I - H$ has eigenvalues: $(n-p) 1's \& p 0's$

• estimate σ^2 by $\hat{\sigma}^2 = \hat{\epsilon}^T \hat{\epsilon} / (n-p) = \text{RSS} / (n-p) \Rightarrow$ an unbiased estimator

• actually, $\hat{\sigma}^2$ has the minimum variance among all quadratic unbiased estimators of σ^2 $\xleftarrow{\text{c.f.}} \text{Gauss-Markov Thm.}$

• $\hat{\sigma} = \sqrt{\text{RSS} / (n-p)} \rightarrow \text{se.}(\hat{\beta}_i) = \sqrt{(X^T X)^{-1}_{ii} \cdot \frac{\text{RSS}}{n-p}}$ (LNp.3-8)

not unbiased

• (FYI) if $\epsilon \sim N(0, \sigma^2 I)$, the maximum likelihood estimator of σ^2 is $\hat{\epsilon}^T \hat{\epsilon} / n = \text{RSS} / n$ (exercise)

Note. OLS estimator of β is also MLE

c.f.
 $n-p$

❖ Reading: Faraway (2005, 1st ed.), 2.4

$\hat{Y} \neq Y$ Goodness-of-Fit: how well does the model fit the data?

- Q: What's "goodness"-of-fit? Why we need it? What can it imply and what cannot?

model 1 (Ω): $y = \beta_0 + \beta_1 g_1(x) + \cdots + \beta_{p-1} g_{p-1}(x) + \epsilon$, Recall LN p.3-7 $\rightarrow Y - \hat{Y}$

model 2 (ω): $y = \beta_0 + \epsilon$. model in one-sample problem (LN p.2-10)

- R^2 , coefficient of determination or percentage of variance explained

$\delta^2 = RSS_{n-p}$ RSS_Ω is the RSS calculated from model 1 (Ω , with all effects g_i 's), $\sum_i (y_i - \bar{y})^2 = TSS_\omega$ is the RSS calculated from model 2 (ω , without any effects g_i 's).

$$\begin{aligned} R^2 &= 1 - \frac{RSS_\Omega}{TSS_\omega} = \frac{TSS_\omega - RSS_\Omega}{TSS_\omega} \\ &= 1 - \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\ &= \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\ &= \left(\frac{\sum (y_i - \bar{y})(\hat{y}_i - \bar{y})}{\sqrt{\sum (y_i - \bar{y})^2 \sum (\hat{y}_i - \bar{y})^2}} \right)^2 = (\text{cor}(Y, \hat{Y}))^2 \end{aligned}$$

$\text{average of } \hat{y}_i \text{'s is } \bar{y} (\because y_i = \hat{y}_i + \epsilon_i)$

$\text{Var}(Y) = E[\text{Var}(Y|X)] + \text{Var}[E(Y|X)]$

Interpretation of R^2 : "proportion of total variation in Y that can be explained by the effects g_i 's." (\rightarrow concept: source of variation in Y) without observing x_i 's

規律 $\rightarrow R = \text{correlation between } \hat{Y} \text{ and } Y$; for simple regression (i.e., Ω only has one effect g_1), $R = \text{correlation between } g_1 \text{ and } Y$ (from the geometry viewpoint, ...)

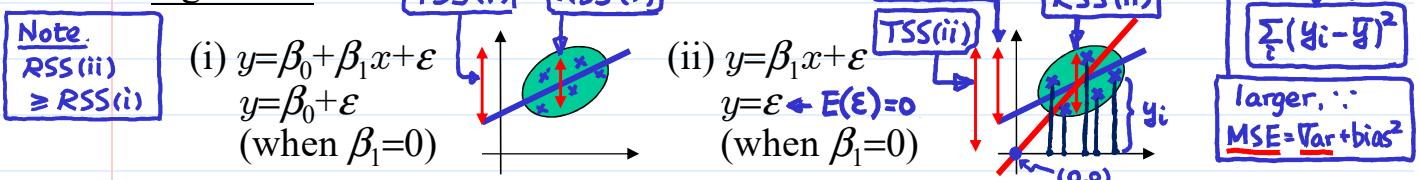
$$D = \hat{\beta}_0 + \hat{\beta}_1 g_1$$

$\triangleright 0 \leq R^2 \leq 1$, a value closer to 1 indicates a better fit. (what if $n \approx p$?)

\triangleright Q: What is a good value of R^2 ? Ans: It depends.

\triangleright Warning. R^2 as defined here does not make any sense if an intercept is not in model. Consider simple regression:

$$\begin{aligned} R^2 &= 1 - \frac{RSS}{TSS} \\ &= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \end{aligned}$$



- Beware of high R^2 reported from models without an intercept.
- Formulate exists for no-intercept R^2 (e.g, compare the RSS s of the models $\Omega: y = \beta_1 x + \epsilon$ and $\omega: y = \beta_0 + \epsilon$), but same graphical intuition is not available.
- Resist throwing out intercept term in model (Note: when intercept is not in model (or $1 \notin \Omega$), sum of residuals might not be zero) Note: $X^T \hat{\epsilon} = 0$

\triangleright Note. R^2 does not indicate whether the model

$E(Y) = X\beta$ is correct, i.e.,

high $R^2 \nrightarrow$ fitted model is correct

low $R^2 \nrightarrow$ fitted model is incorrect

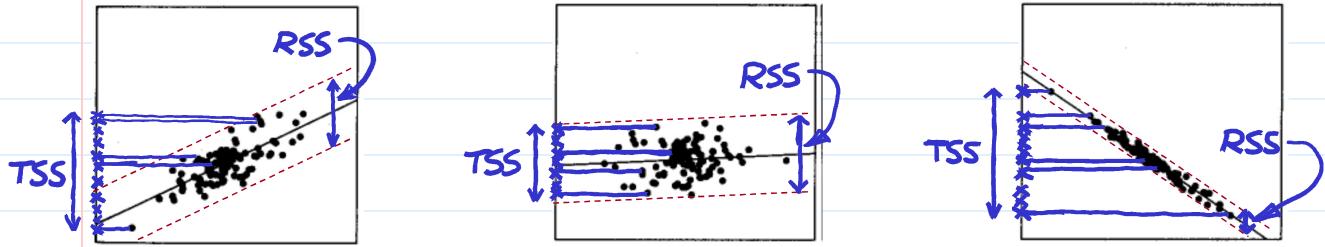
(e.g., a model with $n \approx p$ will have an $R^2 \approx 1$; however, such high R^2 might only indicate over-fitting)

$$Y = X\beta + \epsilon = X^* \hat{\beta}^* + \hat{\epsilon}$$

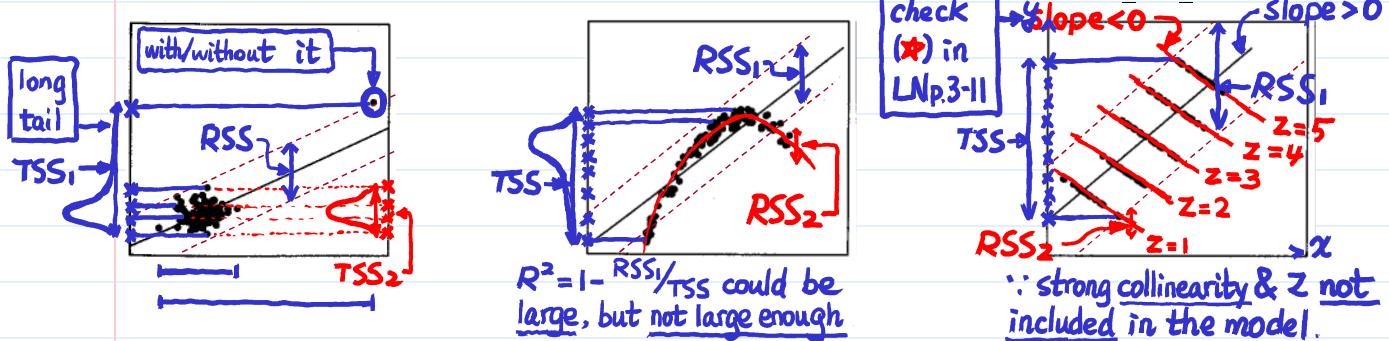
規律 隨機

Q: What happen if
 $\text{cor}(Y, X\beta)^2 \gg \text{cor}(Y, \hat{Y})^2$

Cases for which R^2 seems appropriate under the model: $y = \beta_0 + \beta_1 x + \varepsilon$
 (Q: Which R^2 is higher?)



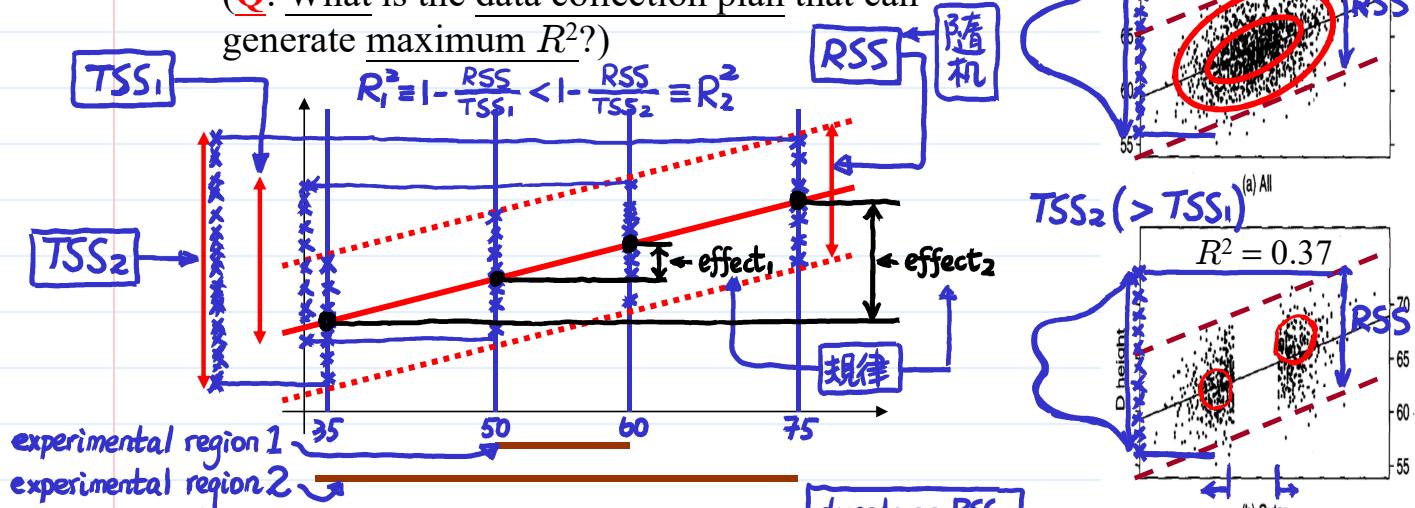
Cases for which R^2 is inappropriate under the model: $y = \beta_0 + \beta_1 x + \varepsilon$



- e.g.,
 Y: weight
 X: height
- When (Y, X) follows a multivariate Normal distribution, R^2 has a connection with population parameters (more details will be given in sampling model, future lecture).
 - mean, variance, correlation of X & Y
 - each member of the population has equal chance of being selected
 - However, if (simple) random sampling is not used in data collection stage, the R^2 no longer estimates a population value (examples?)

- R^2 value can be manipulated merely by changing the data collection plan

(Q: What is the data collection plan that can generate maximum R^2 ?)



Alternative measure for goodness of fit: $\hat{\sigma}$ depends on RSS_2 irrelevant to TSS_w

It is related to standard error of estimates of β and prediction $s.e.(\hat{\beta}_i) = \sqrt{(\mathbf{x}^T \mathbf{x})_{ii}^{-1}} \cdot \hat{\sigma}$

It is measured in the unit of the response (cf: R^2 is free of unit)

check testing for lack of fit (future lecture)

❖ Reading: Faraway (2005, 1st ed.), 2.7

❖ Further reading: D&S, 5.2, 11.2, 21.5

Note: almost identical fitted lines