

Test:  $Q$ : how about test  $H_0: \beta_1=0$  (or  $c$ ) in models 1 and 2 when orthogonality exists between  $\{x_1, I\}$  and  $x_2$ ? will the test results be identical?

model 1:  $\omega_1: y = \beta_0 + \varepsilon$  vs.  $\Omega_1: y = \beta_0 + \beta_1 x_1 + \varepsilon$   $\leftarrow$   $x_2: \beta_1$   $x_2 \perp \Omega_1$   
 $\varepsilon \in \Omega_1$

model 2:  $\omega_2: y = \beta_0 + \beta_2 x_2 + \varepsilon$  vs.  $\Omega_2: y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$   $\leftarrow$   $x_2: \beta_2$

$F = [(RSS_{\omega} - RSS_{\Omega}) / (df_{\omega} - df_{\Omega})] / [RSS_{\Omega} / df_{\Omega}] \sim F_{1, df_{\Omega}}$

$RSS_{\omega_1} - RSS_{\Omega_1} = RSS_{\omega_2} - RSS_{\Omega_2}$  ( $Q$ : why?)  $\rightarrow$  same projection space

but,  $RSS_{\Omega_1} \neq RSS_{\Omega_2}$ , and  $df_{\Omega_1} \neq df_{\Omega_2}$ , i.e.,  $\hat{\sigma}_{\Omega_1}^2 \neq \hat{\sigma}_{\Omega_2}^2$ .

$x_2$ : an important effect

$Q$ : when will the test results be consistent? ( $\hat{\sigma}_{\Omega_1}^2 \approx \hat{\sigma}_{\Omega_2}^2$ ) when will be very different?

Note: although the tests do depend on the presence of  $x_2$ , the dependence is usually not as strong as in non-orthogonal cases.  $\leftarrow$  Check graph in LN p. 4-13  $Q$ : If orthogonal, can we simultaneously remove effects with insignificant  $t$ 's?

orthogonality is very unlikely to achieve in observational data (it's a feature of experimental data from a good design. In experimental case, orthogonal design is an important criterion). At best, predictors are almost uncorrelated and "near" orthogonality holds.

cf. orthogonality

fitted model:  $Y = X\beta + \varepsilon$

centering  $x_1, x_2$  (LN p. 3-7)

Randomization: In an exp't, suppose that true model is  $Y = X\beta + Z\gamma + \varepsilon$ , but  $Z$  cannot be measured or may not even be suspected  $\Rightarrow E(\hat{\beta}) = \beta + (X^T X)^{-1} X^T Z \gamma \neq 0$

$Q$ : what's the best way of controlling  $X$  to make  $X$  and  $Z$  as orthogonal as possible?

- Reading: Faraway (2005, 1st ed.), 3.6
- Further reading: D&S, Appendix 6A

$\{ -1, +1 \}$   $\leftarrow$  cf.  $\{ a, b \}$   $\rightarrow$   $X$  &  $Z$  independent

### Identifiability

model:  $Y = X\beta + \varepsilon$ , where  $X$  is an  $n \times p$  matrix  $\Rightarrow$  OLS estimator  $\hat{\beta} = (X^T X)^{-1} X^T Y$

$Q$ : what if the inverse of  $X^T X$  does not exist?  $\leftarrow$  i.e.  $\text{rank}(X^T X) = \text{rank}(X) < p$   $\leftarrow$   $p \times p$  matrix  $\text{rank}(X^T X) = p$

$\beta$  (or  $X$ ) is called *unidentifiable* when  $X^T X$  is singular ( $\Leftrightarrow \text{rank}(X) < p \Leftrightarrow \dim(\Omega) < p \Leftrightarrow$  at least one column of  $X$  is a linear combination of other columns)

the normal equation  $X^T X \beta = X^T Y$  has infinite many solutions. Any  $\hat{\beta} = (X^T X)^{-} X^T Y$ , is a solution, but should not be regarded as an estimate of  $\beta$ .

generalized inverse of  $A: A A^- A = A$   $\leftarrow$  not unique

$\hat{Y}$  and  $\hat{\varepsilon}$  are still unique

$Q$ : Why does unidentifiability happen?

observational data, some examples:

$x_1$  &  $x_2 = ax_1 + b$

- same predictor measured in different scales, and both are in the model
- $x_1 + x_2 = x_3$ , or  $x_1 + x_2 + x_3 = c$ , and all three are in the model with intercept
- $X$  is *supersaturated*:  $p > n$ , i.e., more effects than observations

$Q$ : What is the hat matrix  $H$  under a saturated model?  $H = I_n$

(Note. *saturated*  $X$ : when  $p = n$  and  $X^T X$  is nonsingular  $\Rightarrow \hat{\beta}$  is identifiable, but no degrees of freedom left for estimation of  $\sigma$

$X^T X$ :  $p \times p$  matrix  $\text{rank}(X^T X) \leq \min(n, p) = n < p$

because  $Y = \hat{Y}$  and  $\hat{\varepsilon} = 0 \Rightarrow$  cannot do testing or C.I.)

$\therefore \dim(Y) = \dim(\Omega) = n$

- such problems can be avoided by paying attention.

➤ experimental data, e.g., two-sample case:

treatment data:  $y_1, \dots, y_n$  <sup>with mean  $\mu_1$</sup>  control data:  $y_{n+1}, \dots, y_{m+n}$  <sup>with mean  $\mu_2$</sup>  Suppose we model the response by an overall mean  $\mu$  and group effects  $\alpha_1$  and  $\alpha_2$ :

$$y_i = \mu + \alpha_1 + \epsilon_i, \quad i=1, \dots, n; \quad y_i = \mu + \alpha_2 + \epsilon_i, \quad i=n+1, \dots, n+m,$$

$(\mu, \alpha_1, \alpha_2)$	$(\mu_1, \mu_2)$
10 5 -1	15 9
5 10 4	15 9
0 15 9	15 9
12 3 -3	15 9

$$\begin{matrix} \text{treatment} \\ \text{control} \end{matrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \\ y_{n+1} \\ \vdots \\ y_{m+n} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \\ \epsilon_{n+1} \\ \vdots \\ \epsilon_{m+n} \end{pmatrix} \Rightarrow X \text{ (or } \beta) \text{ is unidentifiable}$$

⇒ over-parameterized: some constraint must be imposed on  $(\mu, \alpha_1, \alpha_2)$ , say

$\mu=0$  or  $\alpha_1+\alpha_2=0$

$\mu=0$

$\alpha_1 = \mu_1$ : treatment mean  
 $\alpha_2 = \mu_2$ : control mean

$\alpha_1 + \alpha_2 = 0$

$\mu = \frac{\mu_1 + \mu_2}{2}$ : average of 2 means  
 $\alpha_1 = -\alpha_2 = \frac{\mu_1 - \mu_2}{2}$ : (mean difference)/2

• “unidentifiable” means

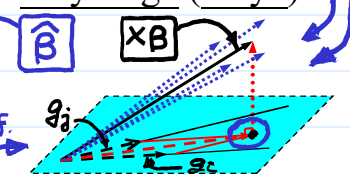
1. insufficient data to estimate the parameters of interest, or
2. more parameters than are necessary to model the data

• an eigen-decomposition of  $X^T X$  will reveal the linear combinations that gave rise to the unidentifiability (check lab)

$\exists \mathbf{q}$  s.t.  $a_1 x_1 + \dots + a_p x_p \approx 0$   
Let  $\mathbf{q} = (a_1, \dots, a_p)$ , then  $(X^T X) \mathbf{q} = 0 = 0 \cdot \mathbf{q}$ .  
⇒  $\mathbf{q}$  is an eigenvector of  $X^T X$ , whose corresponding eigenvalue is 0.  
 $\text{cor}(\beta_i, \beta_j) + \text{cor}(g_i, g_j) = 0$

• what causes problem is data close to “unidentifiable,” (i.e., strong collinearity) ⇒ model is still identifiable, but standard error of estimates can be very large (why?)

• statistical softwares handle unidentifiability differently. R will automatically fit a reduced model when  $X$  is unidentifiable.



- ❖ Reading: Faraway (2005, 1st ed.), 2.9
- ❖ Further reading: D&S, 4.2, 20.4, Appendix 20A.

graph in LNp.4-13

LNp.1-1 ← Interpreting parameter estimates → prediction (Recall. In LNp.5-9 p. 5-13)

• Q:  $Y = X\beta + \epsilon$ , what does  $\hat{\beta}$  mean?

What? Some matters needing attention about  $\hat{\beta}$ :  $\hat{Y} = X_i \hat{\beta}_i$

	$X_1, X_2$ orthogonal	$X_1, X_2$ collinear
$\hat{\beta}_1$	$E(\hat{\beta}_1) = \beta_1$	$E(\hat{\beta}_1) = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2 \beta_2$
$E(\hat{Y})$	$E(\hat{Y}) = X_1 \beta_1$	$E(\hat{Y}) = X_1 \beta_1 + [X_1 (X_1^T X_1)^{-1} X_1^T] (X_2 \beta_2)$

Why? ➤  $\hat{\beta}$  have units [e.g., fuel consumption data, fitted model: true:  $Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$   
fitted:  $Y = X_1 \beta_1 + \epsilon$ ]  
fuel = 154.19 + (-4.23)Tax + (0.47)Dlic + (-6.14)Income + (18.54)log<sub>2</sub>(Miles)]

its unit = unit of y      its unit = (unit of y)/(unit of  $x_i$ )

- sign of  $\hat{\beta}$ : direction of the relationship between the term and the response
- interpretation of estimated value (see next two slides)
- better to also consider values contained in its confidence interval
- causality or association
- the parameters  $\beta$

1. what if the C.I. contains 0?
2. what if (upper bound - lower bound) of the C.I. is very large?
3. what if (upper bound/lower bound) of the C.I. is very large or almost 1 or 0?

- some  $\beta_i$ 's have physical interpretation, especially those from a conceptual model [e.g., attach weights  $x$  to a spring and measure the extension  $y$ ]  
⇒ unfortunately, such cases are rare

check LNp.1-2

- usually,  $\beta_i$ 's do not have such physical interpretation  
⇒ in the case, the model  $Y = X\beta + \epsilon$  is only an empirical model, i.e., a convenience for representing a complex reality within the range of  $X$  ⇒

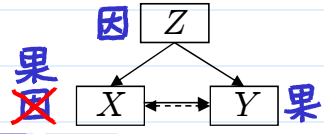
local approximation the real meaning of a particular  $\beta_i$  is not obvious, interpretation is difficult

Some interpretations of parameter estimates

a naive interpretation:

“A unit increase in  $X_i$  will cause an average change of  $\hat{\beta}_i$  in  $Y$ ” ← causality statement

- Q: what if there exist lurking variables? e.g.  $Z$ .  
[e.g.,  $X$ : shoe size,  $Y$ : reading abilities,  $Z$ : age of child]



⇒ causal conclusion is doubtful ← But, can be OK for prediction purpose

- Q: what if the roles of predictor and response are mistakenly switched?  
[e.g.,  $Y$ : fire damage, and  $X$ : numbers of firefighters called out]

All three (X,Y) scatter plots show  $X \uparrow Y \uparrow$

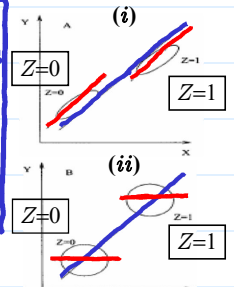
Q: what if some important effects are not included in model?

Even association is questionable

$X$  fixed.  $E(\hat{\beta}_1) = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2 \beta_2$  ← LN p. 5-9

$X$  random. true model:  $E(Y | X_1, X_2) = X_1 \beta_1 + X_2 \beta_2$

Beware of information found in scatter plot (one y, one x)



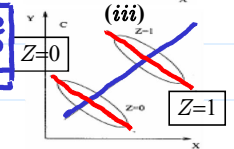
fitted model:  $E(Y | X_1) = X_1 \beta_1$

$E(Y | X_1) \stackrel{!}{=} X_1 \beta_1 + E(X_2 | X_1) \beta_2$

$Var(Y | X_1) \stackrel{!}{=} \sigma^2 + \beta_2^T Var(X_2 | X_1) \beta_2$

- even though we have all important variables in the model and no lurking variables, there still are problems, e.g.:

how to define effects?



Some  $R^2$  same

$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon = \beta_0 + (\beta_1 - \beta_2) X_1 + \beta_2 (X_1 + X_2) + \varepsilon$

- in a properly designed experiment, the naive interpretation is more reasonable (because of its use of orthogonal designs and randomization); but for observational data, it's often questionable.

model:  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$   
(i)  $\hat{\beta}_1 > 0$  (ii)  $\hat{\beta}_1 = 0$   
(iii)  $\hat{\beta}_1 < 0$

an alternative interpretation

after adjusted for the other terms

“A unit increase in  $X_i$  with all the other (specified) terms held constant will be associated with an average change of  $\hat{\beta}_i$  in  $Y$ ”

if specified terms are changed,  $\hat{\beta}_i$  & its interpretation could be different

- Q: can other terms be held constant? e.g.

obs'nal data: hard to fix  $X_2$  & change  $X_1$   
exp'tal data: sliding level

$X_1$  and  $X_2$  are highly correlated

consider the model  $E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 = \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2$

- it requires the specification of the other terms/effects

e.g.  $\hat{\beta}_1 \leftarrow X_1, X_2$   
 $\hat{\beta}_1 \leftarrow X_1, X_2 + X_3$   
 $\hat{\beta}_1 \leftarrow X_1, X_2, X_3$

Q: what will happen in the analysis when strong collinearity exists between effects?

⇒ estimates and tests of  $\beta_i$ 's may significantly change according to what other effects are included. It makes the interpretation almost impossible (check lab).

principal components (future lecture)

In some cases, the problem can be removed by redefining the terms into new linear combinations that may be easier to interpret.

i.e., regard a LM as nothing but a local approximation

an interpretation from prediction viewpoint regarding the parameters and their estimates as fictional quantities, and concentrating on prediction enable a rather cautious interpretation of  $\hat{\beta}$ :

an important objective in the analysis of LM

Recall. unidentifiable many  $\beta$  same

given  $(g_{1,0}, \dots, g_{i,0}, \dots, g_{p-1,0}) \rightarrow \hat{y}_0$ , observe  $(g_{1,0}, \dots, g_{i,0} + 1, \dots, g_{p-1,0}) \rightarrow \hat{y}_0 + \hat{\beta}_i$

avoid "unit increase in  $X_i$  & held constant in others"

- prediction is more stable than parameter estimation (check lab)

- directly interpretable and success may be measured in future

- dangers of extrapolation, be cautious when  $x_0$  is outside the range of  $X$

avoid causality conclusion



• **Q:** how to make a stronger case for causality (be associated with  $\rightarrow$  cause)?

- include all relevant variables/effects  $\Rightarrow$  however, even though you try hard to do so, the possibility of an unsuspected lurking variable will always exist
- fit a variety of models and see if a similar effect is observed, i.e., whether the estimates of  $\beta_i$  similar no matter what the fitted models are?

Recall. Level of evidence (LNp.1-12)

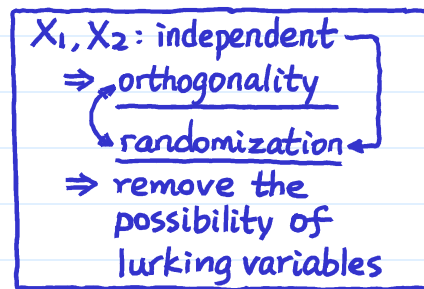
Note. when orthogonality exists  $\Rightarrow$  identical  $\hat{\beta}_i$

e.g. smoking  $\rightarrow$  lung cancer

- use non-statistical knowledge of the physical nature of the relationship  $\Rightarrow$  conceptual model is more persuasive than empirical model
- multiple studies under different conditions can help confirm a relationship.
- in a few cases, one can infer causality from an observational study.

biased sample?

[e.g., Dahl and Moretti (2003): parents of a single girl are 5% more likely to divorce than parents of a single boy. This observational study functions like an experimental design because the sex of a child is a purely random matter.]



caused by some lurking variables or unobserved important predictors?

$\rightarrow X_1$ : sex,  $X_2$ : other possible variables  
 even if these steps are accomplished, one can never be 100% sure of the causality relationship purely based on a statistical analysis. For example, consider the history of the study of the link between smoking and lung cancer  $\Rightarrow$  it takes decades of studies to go from association to causality

❖ Reading: Faraway (2005, 1st ed.), 3.6, 3.7 *FYI. a new field in statistics: causal inference*

## What can go wrong? many many things ...

$$Y = X\beta + \epsilon, \epsilon \sim N(0, \sigma^2 I)$$

source and quality of the data (**Q:** how was the data collected?)

Note. the information contained in the data might be limited or wrong from the beginning.  
 ➢ data may not be a random sample of the population. Situations such as biased sample, a sample of convenience, or sample=population

➢ important predictors may not have been observed (**Q:** how may you find out?)

➢ observational data often make causal conclusions problematic, reason: lack of orthogonality, collinearity, lurking variables, ...

➢ the range of X and qualitative nature of some predictors may limit effective predictions, it's unsafe to extrapolate too much

e.g.  
 • too low  $R^2$   
 • too large  $\sigma^2$   
 • check with expertise

➢ Key: data collected should be representative of the population of interest

• error component [we hope  $\epsilon \sim N(0, \sigma^2 I)$ ]

隨機 the objective of data collection - DOE survey

➢  $\epsilon$  may have unequal variance

WLS

➢  $\epsilon$  may be correlated

GLS

e.g.,  $\epsilon \sim$  heavy-tailed distribution like Cauchy  $\rightarrow$  outlier  $\rightarrow$  robust regression

➢  $\epsilon$  may not be normally distributed

Note  $\hat{\beta} = (X^T X)^{-1} X^T Y$

need it in testing & C.R., not estimation

this is less serious when sample size is large. Notice that even if  $\epsilon$  is not normal,  $\hat{\beta}$  might tend to normality due to CLT. With large sample size, normality of data is not much of a problem

Inference based only on normality of  $\hat{\beta}$  might still be valid.

for small sample sizes, bootstrap method offers a solution

treat empirical cdf of  $\hat{\epsilon}$  as the true cdf of  $\epsilon$

• structural component [  $E(Y)=X\beta$  ] → 規見律

➤ errors in  $X$  → measurement error model

➤ serious collinearity in  $X$

- principal component regression
- ridge regression
- shrinkage method (like LASSO, ...)

➤ some inferences strongly rely on the choice of full model,  $X\beta$  (example?)

Q: where does the full model come from? → conceptual model

Lab5-4

confidence decrease ↓

1. physical theory may suggest a model --- wonderful but relatively uncommon

2. experience from past data --- may help suggesting a reliable model [cf.]

3. no prior experience --- explore current data to find an empirical model

▪ confidence in inference will depend on confidence in the model

▪ an empirical model can be regarded as a local approximation of the underlying true system on some "safe" range of  $X$



• many statistical theory rests on the assumption that the model (error and structural components) is correct. In practice, the best one can hope for is often "empirical model ≈ underlying system". [Box: "all models are wrong but some are useful"]

• publication and experimenter bias / confirmation bias

$H_0: \beta_i = 0$  is true

➤ significant level, say 5% ⇒ keep studying, sooner or later, one will come up with a significant result (about 5% chance) even if one really does not exist.

5% ↑

Problem: significant results get published but not insignificant results

➤ experimenter bias ⇒ many ways of analyzing data, experimenters may be tempted to pick the one that gives them the results they want/expect

不想要則視而不見

看不見的大猩猩 ←

❖ Reading: Faraway (2005, 1<sup>st</sup> ed.), 3.8