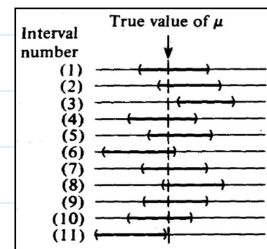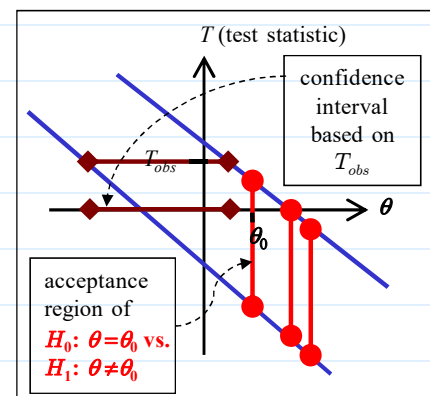# Confidence intervals and regions

- **Q**: Why need <u>interval/region</u> estimation? <u>What more information</u> can it provide compared to <u>point estimation</u>?

  e.g., estimate of $\beta = 3.5$, but <u>accept $H_0$: $\beta = 0$</u>. <u>How</u> to give such result an <u>explanation</u>? <u>Why point estimation</u> cause such <u>confusing</u>?

- An <u>interval/region</u> estimation provides

  ➢ <u>plausible values</u> for <u>parameter</u>
  ➢ <u>uncertainty</u> in parameter estimator

  ➢ <u>information</u> about <u>its length</u> and the <u>values it covers</u> may be <u>helpful</u>
  ➢ information related to <u>testing</u>



- <u>meaning</u> of $100(1-\alpha)\%$ confidence interval or region, e.g., <u>95%</u> confidence interval

- *duality* of <u>interval/region</u> estimation and <u>hypothesis test</u>: For a <u>$100(1-\alpha)\%$</u> confidence region, any <u>point</u> $\theta$ that *lies within* the region represents a <u>null hypothesis</u> that *would not be rejected* at the <u>$100\alpha\%$</u> significance level while every <u>point</u> $\theta$ *outside* represents a <u>null hypothesis</u> that *would be rejected*.

- Model: $\underline{Y = X\beta + \varepsilon}$, $\underline{\varepsilon \sim N(0, \sigma^2 I)}$; $\hat{\beta}$ : <u>OLS estimator</u> $\Rightarrow \hat{\beta} \sim N(\beta, (X^TX)^{-1}\sigma^2)$

  ➢ <u>Confidence region</u> for $\underline{A\beta}$, where $\underline{A}$ is a <u>full rank</u> $\underline{d \times p}$ matrix and $\underline{d \le p}$

  $$A\hat{\beta} \sim N(A\beta, \underline{A(X^TX)^{-1}A^T\sigma^2}) \Rightarrow [(A\hat{\beta}-A\beta)^T[A(X^TX)^{-1}A^T]^{-1}(A\hat{\beta}-A\beta)]/\sigma^2 \sim \chi^2_{\underline{d}},$$

  $$(n-p)\,\hat{\underline{\sigma}}^2 / \sigma^2 \sim \chi^2_{\underline{n-p}},$$

  and they are **independent**.

  $$[(A\hat{\beta}-A\beta)^T[A(X^TX)^{-1}A^T]^{-1}(A\hat{\beta}-A\beta)] / (d\,\hat{\sigma}^2) \sim F_{\underline{d,n-p}}$$

  ➢ <u>$100(1-\alpha)\%$ confidence region of $\underline{A\beta}$</u>: <u>collection</u> of $\underline{A\beta}$'s (or $\beta$) that satisfy

  **general form** $[(A\hat{\beta}-\underline{A\beta})^T[A(X^TX)^{-1}A^T]^{-1}(A\hat{\beta}-\underline{A\beta})] / (d\,\hat{\sigma}^2) \le F_{d,n-p}^{(\alpha)}$

  The regions are often <u>ellipsoidally</u> shaped (**Q**: <u>why?</u>).

- <u>Examples</u>:

  ➢ confidence region <u>for $\beta$</u>, i.e, $\underline{A=I_{p \times p}}$

  $$(\hat{\beta}-\beta)^T X^T X(\hat{\beta}-\beta) \le (\underline{p}\,\hat{\sigma}^2)\,F_{\underline{p,n-p}}^{(\alpha)}$$

  (**Q**: What's the <u>confidence region</u> for <u>all effects</u>?)

  ➢ confidence region of $\underline{\beta_i, \beta_j}$, i.e, $\underline{A} = \begin{pmatrix} 0,\cdots,0,1,0,\cdots,0,0,0,0,\cdots,0 \\ 0,\cdots,0,0,0,0,\cdots,0,1,0,\cdots,0 \end{pmatrix}$

  $$[(A\hat{\beta}-\underline{A\beta})^T[\underline{A(X^TX)^{-1}A^T}]^{-1}(A\hat{\beta}-\underline{A\beta})] \le (\underline{2}\,\hat{\sigma}^2)F_{\underline{2,n-p}}^{(\alpha)}$$

$$\mathbf{Z} \equiv \begin{pmatrix} \hat{\beta}_i \\ \hat{\beta}_j \end{pmatrix} \sim \mathrm{N}\left( \begin{pmatrix} \beta_i \\ \beta_j \end{pmatrix}, \sigma^2 \begin{pmatrix} (\mathbf{X}^T\mathbf{X})_{ii}^{-1} & (\mathbf{X}^T\mathbf{X})_{ij}^{-1} \\ (\mathbf{X}^T\mathbf{X})_{ji}^{-1} & (\mathbf{X}^T\mathbf{X})_{jj}^{-1} \end{pmatrix} \right) \equiv \mathrm{N}\left(\mu, \sigma^2\Sigma\right) \Leftarrow \Sigma = A(\mathbf{X}^T\mathbf{X})^{-1}A^T$$

confidence region of $\underline{\beta_i}$ and $\underline{\beta_j}$: $\left\{ \mu \mid \underline{(\mathbf{Z}-\mu)^T\Sigma^{-1}(\mathbf{Z}-\mu) \leq c} \right\}$ for some $c$

- example: <u>confidence region</u> and <u>intervals</u> of $\underline{\beta_{p15}}$ and $\underline{\beta_{p75}}$

  □ **Q1**: why the <u>straight lines</u> <u>not tangential</u> to the ellipse?

  $$\begin{aligned} 1-\alpha &= P\left(\{(\beta_{p15}, \beta_{p75}) \in \text{C. Region}\}\right) \\ &= P\left(\{(\beta_{p15}, \beta_{p75}) \in \text{C.I.}_{p15} \times \mathbb{R}\}\right) \\ &= P\left(\{(\beta_{p15}, \beta_{p75}) \in \mathbb{R} \times \text{C.I.}_{p75}\}\right) \end{aligned}$$



  □ **Q2**: <u>what</u> can you say, based on the <u>plot</u>, about the <u>results</u> of testing $\underline{H_0^1}{:}\beta_{p15}{=}0$, $\underline{H_0^2}{:}\beta_{p75}{=}0$, and $\underline{H_0^3}{:}\beta_{p15}{=}\beta_{p75}{=}0$?

  □ **Q3**: <u>where</u> will be the point $\underline{(0,0)}$ <u>located</u> if the data <u>accept</u> $H_0^1$, $H_0^2$, <u>reject</u> $H_0^3$? <u>how</u> to <u>explain</u> the result if $(0,0)$ falls <u>in other regions</u>? (<u>exercise</u>)

  □ **Q4**: what is the <u>correlation</u> between $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$? how will the <u>shape</u> of <u>ellipse</u> <u>change</u> when the <u>correlation</u> becomes <u>larger</u> or <u>smaller</u>?

  □ **Q5**: can you see <u>why</u> the <u>situation in</u> **Q3** will happen <u>more</u> <u>frequently</u> when the <u>correlation</u> between $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$ gets <u>larger</u>?

  □ **Q6**: if $\hat{\beta}_{p15}$ and $\hat{\beta}_{p75}$ are <u>uncorrelated</u>, what would be the <u>shape</u> of the <u>confidence region</u>? <u>why</u> situation in **Q3** <u>less possible</u> to occur?

➢ confidence interval <u>for $\beta_i$</u>, i.e, $\underline{A}=(0,\dots,0,1,0,\dots,0)$

$$(\hat{\beta}_i-\beta_i)^2/(\mathbf{X}^T\mathbf{X})^{-1}_{ii} \leq \hat{\sigma}^2 F_{\underline{1,n-p}}^{(\alpha)} \Rightarrow \left| (\hat{\beta}_i-\beta_i)/\left(\hat{\sigma}\sqrt{(\mathbf{X}^T\mathbf{X})^{-1}_{ii}}\right) \right| \leq t_{\underline{n-p}}^{(\alpha/2)}$$

<u>alternative method:</u>

① $\hat{\beta}_i \sim N(\beta_i, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}_{ii})$, ② $(n-p)\hat{\sigma}^2/\sigma^2 \sim \chi^2_{n-p}$, and ③ they are <u>independent</u>

$$\Rightarrow \underline{(\hat{\beta}_i-\beta_i)/\left(\hat{\sigma}\sqrt{(\mathbf{X}^T\mathbf{X})^{-1}_{ii}}\right) \sim t_{n-p}} \qquad \Rightarrow \text{C.I.:} \ \underline{\hat{\beta}_i} \ \pm \ \underline{t_{n-p}^{(\alpha/2)}} \times \underline{\left(\hat{\sigma}\sqrt{(\mathbf{X}^T\mathbf{X})^{-1}_{ii}}\right)}.$$

➢ confidence interval for <u>prediction of mean response</u> at $\underline{x_0}$

$$x_0^T\hat{\beta} - x_0^T\beta \sim N(0, (x_0^T(\mathbf{X}^T\mathbf{X})^{-1}x_0)\sigma^2) \Rightarrow (x_0^T\hat{\beta} - x_0^T\beta)/\left(\hat{\sigma}\sqrt{x_0^T(\mathbf{X}^T\mathbf{X})^{-1}x_0}\right) \sim t_{n-p}$$

$$\Rightarrow \text{C.I.:} \ \underline{x_0^T\hat{\beta}} \ \pm \ t_{n-p}^{(\alpha/2)} \times \underline{\left(\hat{\sigma}\sqrt{x_0^T(\mathbf{X}^T\mathbf{X})^{-1}x_0}\right)}$$

- **Q**: for a given dataset and $\underline{\alpha}$, the <u>length</u> of the <u>C.I.</u> is <u>related to $x_0$ only</u>. What $\underline{x_0}$ will cause a <u>wider C.I.</u>? **Ans**: $\underline{x_0}$ that is <u>away from</u> the "<u>center</u>" of $\underline{X}$

- <u>interpolation and extrapolation</u>
  - □ <u>interpolation</u>: $\underline{x_0}$ lie "<u>within the range</u>" of $\underline{X}$
  - □ <u>extrapolation</u>: $\underline{x_0}$ lie "<u>outside the range</u>" of $\underline{X}$
    (**Q**: <u>fitted model</u> still <u>hold outside the range</u>?)
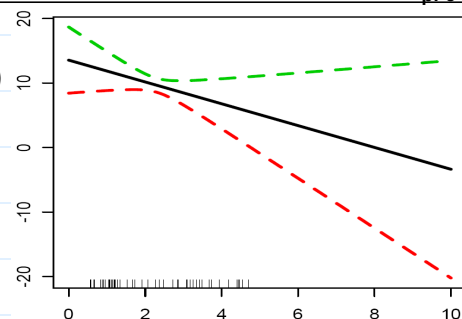    ✓ <u>quantitative</u> $x_0$ and <u>qualitative</u> $x_0$

- Example: 95% *pointwise* confidence band for prediction of mean responses (model: $y = \beta_0 + \beta_1 x + \varepsilon$)

  

  - **Q1**: why the confidence intervals get wider when we move away from the range of data?

    _____

  - **Q2**: what's the danger of extrapolation?
  - **Q3**: does the widening reflect the possibility that the mean structure of the model may change outside the range?
  - **Q4**: does the plot represent a *simultaneous* confidence band for **all** prediction of mean response?

➤ C.I. for prediction of future observation at $x_0$

$$x_0^T\hat{\beta} - (x_0^T\beta + \varepsilon) \sim N(\mathbf{0}, (x_0^T(X^TX)^{-1}x_0 + \underline{1})\sigma^2)$$

$$\Rightarrow \text{C.I.:} \quad \underline{x_0^T\hat{\beta}} \pm t_{n-p}^{(\alpha/2)} \times \left( \hat{\sigma}\sqrt{1 + x_0^T(X^TX)^{-1}x_0} \right)$$

➤ a general form for confidence interval:

$$\textbf{estimate} \pm (\textbf{critical value}) \times (\textbf{standard error of estimate})$$

❖ **Reading**: Faraway (2005, 1st ed.), 3.4, 3.5
❖ **Futher reading**: D&S, 5.3, 5.4, 5.5

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

# **Sampling model**

population



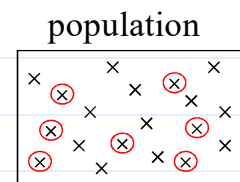- experimental data vs. observational data

  It depends on whether we have *control* over predictors

  examples: yield of crop.

  experimental: fertilizer, …,
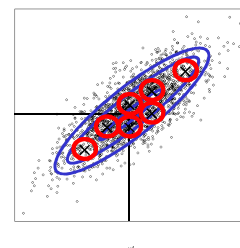  observational: exposure, weather, …

  

  ➤ **Q**: What difference between inferences based on experimental data and observational data?

  **Ans**: experimental data: causation,

  observational data: often only association (**Note**. lurking variable)

  ➤ **Q**: Is this model description, $Y = X\beta + \varepsilon$, $\varepsilon \sim N(\mathbf{0}, \sigma^2 I)$, appropriate for observational data? **Note** that

  (1) observational $X$ are random variables
  (2) in LM, $X$ are treated as fixed values, i.e., no distribution assigned for $X$

  ➤ **Q**: difference between "$X$ is random" and "$X$ measured with (random) error"

  example:

- for some data sets, we can regard the data as a *sample* drawn from a *population*. In the case, we want to say something about the unknown population values using estimated values that are obtained from the sampled data. (example?)

- the data should be generated using a "(simple) *random* sample" of the population so that they can be representative
- conditional distribution of multivariate normal: If

$$\mathbb{Z} = \left[\frac{Z_1}{Z_2}\right] \sim N\left(\left[\frac{\mu_1}{\mu_2}\right], \left[\begin{array}{cc}\Sigma_{11} & \Sigma_{12}\\ \Sigma_{21} & \Sigma_{22}\end{array}\right]\right),$$

then

$$Z_1 \Big| Z_2 = z_2 \sim N\left(\mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(z_2 - \mu_2),\ \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\right)$$

- an alternative view of regression: data $(y_i, x_i)$, $i=1,\ldots,n$, are randomly sampled from a multivariate Normal population,
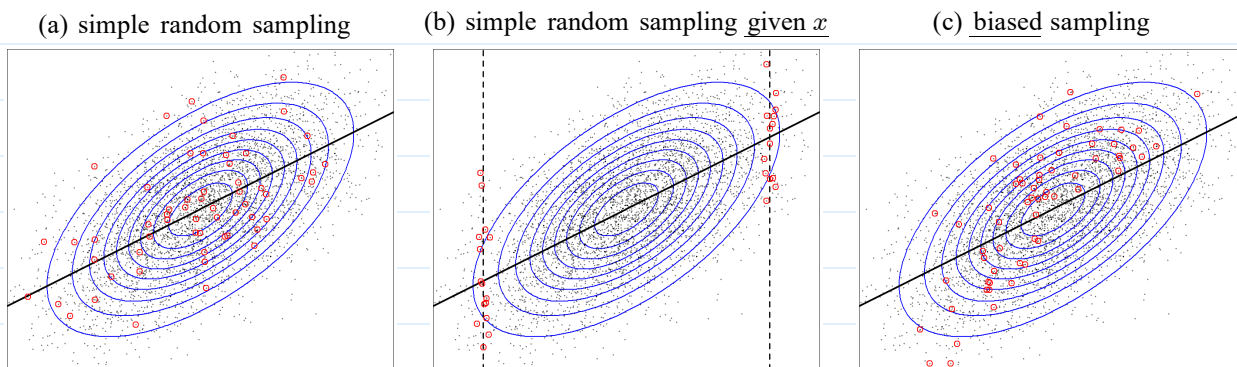
$$\left[\begin{array}{c}Y_i\\ X_i\end{array}\right] \sim N\left(\left[\begin{array}{c}\mu_Y\\ \mu_X\end{array}\right], \left[\begin{array}{cc}\sigma_Y^2 & \Sigma_{XY}^T\\ \Sigma_{XY} & \Sigma_{XX}\end{array}\right]\right) \Rightarrow y_i | X_i = x_i \sim N((\mu_Y - \beta^T\mu_X) + \beta^T x_i,\ \sigma^2).$$

It is a linear model with $\beta = \Sigma_{XX}^{-1}\Sigma_{XY}$, $\sigma^2 = \sigma_Y^2 - \Sigma_{XY}^T\Sigma_{XX}^{-1}\Sigma_{XY} \equiv \sigma_Y^2(1 - r^2)$.

When we are interested in the "transformed" parameters, regression can be applied.

- **Q**: what information in these samples is proper?

(a) simple random sampling     (b) simple random sampling given $x$     (c) biased sampling



| | (a) | (b) | (c) |
|---|---|---|---|
| Joint distribution $(X, Y)$ | ✓ | ✗ | ✗ |
| Conditional distribution $Y|X$ | ✓ | ✓✓ | ✗ |
| Marginal distribution $X$ | ✓ | ✗ | ✓ or ✗ |

- **Q**: what will happen if the sample is not random?
  (i) biased sample
  (ii) sample of convenience
  (iii) sample = population
  these nonrandom samples can cause problems in the inference (e.g., $R^2$, LNp.3-18)

❖ **Reading**: Faraway (2005, 1st ed.), 3.8, nonrandom samples
       Weisberg (2005), *Applied Linear Regression*, 3rd Ed., 4.2, 4.3

# Orthogonality

- **Q**: consider the <u>two models</u>:
  <u>model 1</u>: $y = \beta_0 + \beta_1 x_1 + \varepsilon$,
  <u>model 2</u>: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

  In general, $\hat{\beta}_1$, in the <u>two models</u> are <u>not</u>
  identical (**Q**: <u>why?</u>)
  [also, test <u>$H_0: \beta_1=0$</u> (or <u>c</u>) <u>not identical</u>]
  an *exception*: when $x_1$ and $x_2$ are <u>orthogonal</u>

<div style="border:1px solid">

<u>fitted model</u>=<u>model 1</u>: $Y=X_1\beta_1+\varepsilon$
<u>true model</u>=<u>model 2</u>: $Y=X_1\beta_1+X_2\beta_2+\varepsilon$

- $E(\hat{\beta}_1) = \beta_1 + \underline{(X_1^T X_1)^{-1} X_1^T X_2 \beta_2}$
- $E(X_1\hat{\beta}_1) = X_1\beta_1$
  $\quad + \underline{X_1(X_1^T X_1)^{-1} X_1^T (X_2\beta_2)}$

Note. If <u>fitted model</u>=<u>model 2</u>
- $E(\hat{\beta}_1) = \beta_1$
</div>

- $Y=\underline{X\beta}+\varepsilon=\underline{X_1\beta_1+X_2\beta_2}+\varepsilon$, where $\underline{\beta=[\beta_1\ \beta_2]^T}$ and $\underline{X=[X_1\ X_2]}$ with the property
  $\underline{X_1^T X_2=0} \Rightarrow \underline{X_1}$ and $\underline{X_2}$ are <u>orthogonal</u> (<u>generalization?</u>)

$$X^T X = \begin{pmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{pmatrix} = \begin{pmatrix} X_1^T X_1 & \boxed{0} \\ 0 & X_2^T X_2 \end{pmatrix}$$

$$\Rightarrow \quad (X^T X)^{-1} = \begin{pmatrix} (X_1^T X_1)^{-1} & \boxed{0} \\ 0 & (X_2^T X_2)^{-1} \end{pmatrix}$$

- <u>Estimation</u>: $\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T Y$, $\hat{\beta}_2 = (X_2^T X_2)^{-1} X_2^T Y$, and $\hat{\beta}_1, \hat{\beta}_2$ <u>independent</u>
  $\Rightarrow$ note that $\hat{\beta}_1$ will be the <u>same regardless of</u>
  whether $X_2$ is in the model or not (and vise versa).

- **Q**: what if <u>only two predictors</u>, say <u>some $x_i$ in $X_1$</u>
  and <u>some $x_j$ in $X_2$</u>, are <u>orthogonal?</u>

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- Test: **Q**: how about test <u>$H_0: \beta_1=0$</u> (or <u>c</u>) in <u>models 1 and 2</u> when <u>orthogonality</u>
  <u>exists</u> between $\{\underline{x_1, 1}\}$ and $\underline{x_2}$? will the <u>test results</u> be <u>identical?</u>
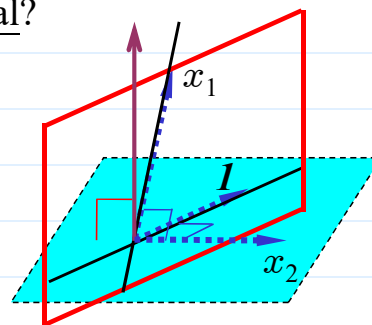
  <u>model 1</u>: $\omega_1: y=\beta_0+\varepsilon$ vs. $\Omega_1: y=\beta_0+\beta_1 x_1+\varepsilon$

  <u>model 2</u>: $\omega_2: y=\beta_0+\beta_2 x_2+\varepsilon$ vs. $\Omega_2: y=\beta_0+\beta_1 x_1+\beta_2 x_2+\varepsilon$
  $F=[(\underline{RSS_\omega-RSS_\Omega})/(df_\omega-df_\Omega)]/[\underline{RSS_\Omega/df_\Omega}]\sim F_{1,\underline{df_\Omega}}$
  $\underline{RSS_{\omega_1}-RSS_{\Omega_1}}=\underline{RSS_{\omega_2}-RSS_{\Omega_2}}$ (**Q**: <u>why?</u>)
  but, $\underline{RSS_{\Omega_1}\neq RSS_{\Omega_2}}$, and $\underline{df_{\Omega_1}\neq df_{\Omega_2}}$, i.e., $\hat{\sigma}^2_{\Omega_1} \neq \hat{\sigma}^2_{\Omega_2}$.

  **Q**: <u>when</u> will the <u>test results</u> be <u>consistent?</u> ($\hat{\sigma}^2_{\Omega_1} \approx \hat{\sigma}^2_{\Omega_2}$) <u>when</u> will be <u>very different?</u>
  **Note**: although the <u>tests</u> do <u>depend on</u> the <u>presence</u> of $x_2$, the <u>dependence</u> is usually
  <u>not as strong as</u> in <u>non-orthogonal</u> cases.



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

- <u>Randomization</u>: In an <u>exp't</u>, suppose that true model is $Y=\underline{X\beta}+\underline{Z\gamma}+\varepsilon$, but $\underline{Z}$ <u>cannot</u>
  <u>be measured</u> or <u>may not even be suspected</u> $\Rightarrow E(\hat{\beta})=\beta+\underline{(X^T X)^{-1} X^T Z\gamma}$
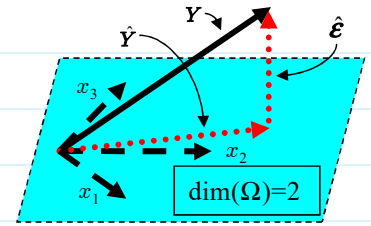  **Q**: what's the <u>best way</u> of <u>controlling $X$</u> to make $X$ and $Z$ as <u>orthogonal</u> as possible?

❖ **Reading**: Faraway (2005, 1st ed.), 3.6
❖ **Futher reading**: D&S, Appendix 6A

# Identifiability

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

  **Q**: what if the <u>inverse</u> of $\underline{X}^T\underline{X}$ does <u>not</u> exist?

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

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

  ➤ $\hat{Y}$ and $\hat{\boldsymbol{\varepsilon}}$ are still <u>unique</u>

- **Q**: <u>Why</u> does <u>unidentifiability</u> <u>happen</u>?
  ➤ <u>observational data</u>, some <u>examples</u>:
    - <u>same predictor</u> measured in <u>different scales</u>, and <u>both</u> are <u>in the model</u>
    - $\underline{X_1}{+}\underline{X_2}{=}\underline{X_3}$, or $\underline{X_1}{+}\underline{X_2}{+}\underline{X_3}{=}c$, and <u>all three</u> are <u>in the model</u> with intercept
    - $\underline{X}$ is *supersaturated*: $p{>}n$, i.e., <u>more effects</u> than <u>observations</u>

      (**Note**. *saturated* $X$: when $p{=}n$ and $\underline{X}^T\underline{X}$ is <u>nonsingular</u> $\Rightarrow \hat{\boldsymbol{\beta}}$ is <u>identifiable</u>, but <u>no degrees of freedom</u> left for <u>estimation of</u> $\sigma$ because $\underline{Y}{=}\hat{Y}$ and $\hat{\boldsymbol{\varepsilon}}{=}\boldsymbol{0} \Rightarrow$ <u>cannot</u> do <u>testing</u> or <u>C.I.</u>)
    - <u>such problems</u> can be <u>avoided</u> by <u>paying attention</u>.



$$\hat{Y} \quad Y \quad \hat{\boldsymbol{\varepsilon}}$$
$$x_3 \quad x_2 \quad x_1 \quad \boxed{\dim(\Omega){=}2}$$

  ➤ <u>experimental data</u>, e.g., <u>two-sample</u> case:
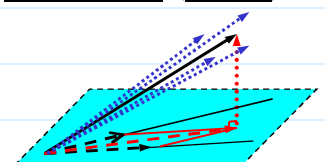
  <u>treatment</u> data: $y_1, ..., y_n$,  <u>control</u> data: $y_{n+1}, ..., y_{m+n}$.  Suppose we <u>model</u> the <u>response</u> by an <u>overall mean</u> $\mu$ and <u>group effects</u> $\alpha_1$ and $\alpha_2$:

  $$y_i{=}\underline{\mu}{+}\underline{\alpha_1}{+}\varepsilon_i,\ i{=}1,...,n;\quad y_i{=}\underline{\mu}{+}\underline{\alpha_2}{+}\varepsilon_i,\ i{=}n{+}1,...,n{+}m,$$

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

  $\Rightarrow$ *over-parameterized*: some <u>constraint</u> must be <u>imposed</u> on $(\underline{\mu}, \alpha_1, \alpha_2)$, say $\underline{\mu{=}0}$ or $\underline{\alpha_1{+}\alpha_2{=}0}$

- "unidentifiabile" means
  1. <u>insufficient data</u> to estimate the <u>parameters of interest</u>, or
  2. <u>more parameters</u> than are <u>necessary to model</u> the data

- an <u>eigen-decomposition</u> of $\underline{X}^T\underline{X}$ will <u>reveal</u> the <u>linear combinations</u> that <u>gave rise</u> to the <u>unidentifiability</u> (check <u>lab</u>)

- what causes problem is <u>data</u> *close to* "unidentifiable," (i.e., <u>strong collinearity</u>) $\Rightarrow$ model is still <u>identifiable</u>, but <u>standard error</u> of <u>estimates</u> can be <u>very large</u> (why?)

- <u>statistical softwares</u> handle unidentifiability <u>differently</u>. R will <u>automatically fit</u> a <u>reduced model</u> when $\underline{X}$ is unidentifiable.

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

# Interpreting parameter estimates

- **Q**: $Y = X\beta + \varepsilon$, what does $\hat{\beta}$ mean?

  Some <u>matters</u> needing attention <u>about</u> $\hat{\beta}$ :

  - $\hat{\beta}$ have <u>units</u> [e.g., <u>fuel consumption</u> data, fitted model:
    <u>fuel</u> = <u>154.19</u> + (<u>− 4.23</u>)Tax + (0.47)Dlic + (−6.14)Income + (18.54)$\log_2$(Miles)]

  - <u>sign</u> of $\hat{\beta}$ : <u>direction</u> of the <u>relationship</u> between the <u>term</u> and the <u>response</u>
  - <u>interpretation</u> of <u>estimated value</u> (see <u>next two slides</u>)
  - better to also <u>consider values</u>
    contained <u>in its confidence interval</u>
  - <u>causality</u> or <u>association</u>
  - the parameters $\beta$
    - some $\beta_i$'s have <u>physical interpretation</u>, especially those from a <u>conceptual</u>
      <u>model</u> [e.g., attach <u>weights</u> $x$ to a <u>spring</u> and measure the <u>extension</u> $y$]
      $\Rightarrow$ unfortunately, <u>such cases</u> are <u>rare</u>
    - usually, $\beta_i$'s do <u>not have</u> such <u>physical interpretation</u>
      $\Rightarrow$ in the case, the model $Y = X\beta + \varepsilon$ is only an *empirical model*, i.e., a
      <u>convenience</u> for <u>representing a complex reality</u> <u>within</u> the <u>range of</u> $X$ $\Rightarrow$
      the <u>real meaning</u> of a particular $\beta_i$ is <u>not obvious</u>, <u>interpretation is difficult</u>

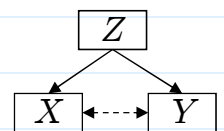- <u>Some interpretations</u> of parameter estimates
  - a <u>naive</u> interpretation:
    **"A <u>unit increase</u> in $X_i$ will *cause* an <u>average change</u> of $\hat{\beta}_i$ in $Y$"** $\Leftarrow$ | causality statement |
    [e.g., $Y$: <u>annual income</u>, and $X$: <u>years of education</u>]
    - **Q**: what if there exist <u>lurking variables</u>?
      [e.g., $X$: <u>shoe size</u>, $Y$: <u>reading abilities</u>, $Z$: <u>age of child</u>]
      $\Rightarrow$ <u>causal conclusion</u> is <u>doubtful</u>

      

    - **Q**: what if the <u>roles of predictor and response</u> are <u>mistakenly switched</u>?
      [e.g., $Y$: <u>fire damage</u>, and $X$: <u>numbers of firefighters</u> called out]
    - **Q**: what if some <u>important effects</u> are <u>not included</u> in <u>model</u>?
      - $\underline{X}$ <u>fixed</u>. $E(\hat{\beta}_1) = \beta_1 + (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \beta_2$
      - $\underline{X}$ <u>random</u>. true model: $E(Y \mid \mathbf{X}_1, \mathbf{X}_2) = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2$ ,
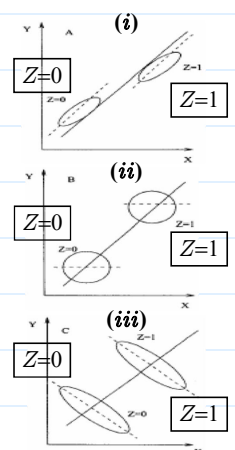        <u>fitted</u> model: $E(Y \mid \mathbf{X}_1) = \mathbf{X}_1 \beta_1$
        $E(Y \mid \mathbf{X}_1) = \mathbf{X}_1 \beta_1 + E(\mathbf{X}_2 \mid \mathbf{X}_1) \beta_2$
        $Var(Y \mid \mathbf{X}_1) = \sigma^2 + \beta_2^T \, Var(\mathbf{X}_2 \mid \mathbf{X}_1) \, \beta_2$

      

    - even though we have <u>all important variables</u> in the model
      and <u>no lurking variables</u>, there still are <u>problems</u>, e.g.:
      $y = \beta_0 + \underline{\beta_1} \, \underline{X_1} + \beta_2 \, \underline{X_2} + \varepsilon = \beta_0 + (\underline{\beta_1} - \underline{\beta_2}) \, \underline{X_1} + \beta_2 (\underline{X_1 + X_2}) + \varepsilon$
    - in a <u>properly designed</u> <u>experiment</u>, the <u>naive interpretation</u> is
      <u>more reasonable</u> (because of its <u>use of orthogonal designs</u> and
      <u>randomization</u>); but for <u>observational data</u>, it's <u>often questionable</u>.

➢ an alternative interpretation

" **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** $Y$ "

- **Q**: can other terms be held constant? e.g.
  - ▫ $X_1$ 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.

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

➢ an interpretation from prediction viewpoint

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

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

- 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 $X$

- **Q**: how to make a stronger case for causality (be associated with → cause)?

  ➢ include all relevant variables/effects ⇒ however, even though you try hard to do so, the possibility of an unsuspected lurking variables will always exists

  ➢ 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?

  ➢ use non-statistical knowledge of the physical nature of the relationship
    ⇒ 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.

    [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.]

  ➢ 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 ⇒ it takes decades of studies to go from association to causality

❖ **Reading**: Faraway (2005, 1st ed.), 3.6, 3.7

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

$Y=X\beta+\varepsilon$, $\varepsilon\sim N(0, \sigma^2 I)$   p. 5-17

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

  ➢ 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

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

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

  ➢ $\varepsilon$ may have unequal variance

  ➢ $\varepsilon$ may be correlated

  ➢ $\varepsilon$ may not be normally distributed

    ▪ this is less serious when sample size is large. Notice that even if $\varepsilon$ 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

    ▪ for small sample sizes, bootstrap method offers a solution

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

p. 5-18

- structural component [ $E(Y)=X\beta$ ]

  ➢ errors in $X$

  ➢ serious collinearity in $X$

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

  **Q**: where does the full model come from?

    1. physical theory may suggest a model --- wonderful but relatively uncommon
    2. experience from past data --- may help suggesting a reliable model
    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

  ➢ 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. 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*

❖ **Reading**: Faraway (2005, 1st ed.), 3.8