*Handwritten annotations at top:*

Causal Quantity

Study design

$\longrightarrow$ Est

Inter

(Least) identifiability

Valid,
efficiency,
robustness

# Potential Outcomes and Causal Effects

Xinzhou Guo

HKUST

Spring 2024

(Credited to Zhichao Jiang)

# Causation:

Treatment A                    Treatment B

Is treatment A better than treatment B

(in helping patients control blood pressure)?

e.g. job training program, gender.

Association: correlation, regression coef.

treatment ✓ blood pressure reduction

# Simpson's paradox

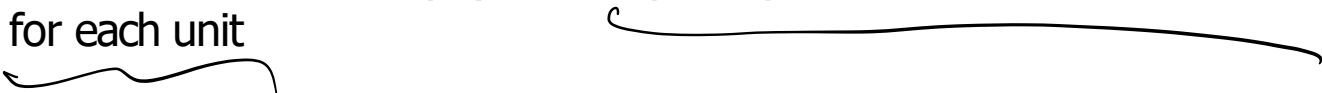| | small stones | | large stones | |
|---|---|---|---|---|
| | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment A: open surgical procedures
- Treatment B: a minimally-invasive procedure

- Success rate for small stones: **93%** (81/87) > 87% (234/270)
  Success rate for large stones: **73%** (192/263) > 69% (55/80)
- Overall success rate: 78% (273/350) < **83%** (289/350)
  -- Why and Is treatment A better than treatment B?

# Potential outcomes framework (Neyman 1923; Rubin 1974)

- Success rate (A > B) $\longrightarrow$ positive association between stone removal and treatment A

- Association ≠ Causation; <span style="color:red">The comparison between treatment A and treatment B is about association or causation?</span>

- Causation: comparison between <span style="color:red">potential</span> outcomes under treatment and control for the same unit(s) $\longrightarrow$ What if xxx?

- Defining causal quantities by potential outcomes requires a <span style="color:red">thought experiment</span>; neither data nor actual experimentation needed
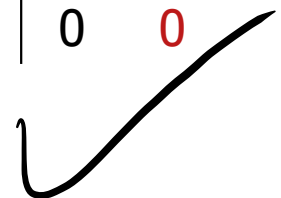
# Potential outcome and observed outcome

- Observed data: treatment $Z_i$, outcome $Y_i$
- Potential outcomes: $Y_i(1)$ and $Y_i(0)$
    - categorical: $Y_i(0)$, $Y_i(1),\ldots,$ $Y_i(K-1)$
    - continuous: $Y_i(z)$ for any $z \in \mathbb{R}$
    - observed outcome: $Y_i(Z_i) \longrightarrow$ only one potential outcome is observed for each unit

# Potential outcome and observed outcome

- Observed data: treatment $Z_i$, outcome $Y_i$
- Potential outcomes: $Y_i(1)$ and $Y_i(0)$
    - categorical: $Y_i(0)$, $Y_i(1)$,..., $Y_i(K-1)$
    - continuous: $Y_i(z)$ for any $z \in \mathbb{R}$
    - observed outcome: $Y_i(Z_i)$ ⟶ only one potential outcome is observed for each unit

| Unit $i$ | $Z_i$ | $Y_i(1)$ | $Y_i(0)$ |
|----------|-------|----------|----------|
| 1 | 1 | 0 | 1 |
| 2 | 0 | 0 | 1 |
| 3 | 1 | 0 | 0 |
| 4 | 1 | 1 | 1 |
| 5 | 0 | 1 | 0 |

$Y_i(1) - \frac{1}{n}(0)$

| Unit $i$ | $Z_i$ | $Y_i$ |
|----------|-------|-------|
| 1 | 1 | 0 |
| 2 | 0 | 1 |
| 3 | 1 | 0 |
| 4 | 1 | 1 |
| 5 | 0 | 0 |

*unobservable*

$$Y_i = Z_i \, Y_i(1) + (1 - Z_i) \, Y_i(0)$$

$$= \begin{cases} Y_i(1) & Z_i = 1 \\ \\ Y_i(0) & Z_i = 0 \end{cases}$$

# Hidden assumptions on potential outcomes

- The notation of $Y_i(z)$ implies three assumptions
  - no interference between units:

$$Y_i(Z_1, \ldots, Z_n) = Y_i(Z_i)$$

  - same version of treatment

  - treatment occurs before outcomes

- Stable Unit Treatment Value Assumption (SUTVA)
  - no interference
  - only one version of treatment

# Violation of SUTVA

**No interference** can be violated in infectious diseases or network experiments. For instance, if some of my friends receive flu shots, my chance of getting the flu decreases even if I do not receive the flu shot; if my friends see an advertisement on Facebook, my chance of buying that product increases even if I do not see the advertisement directly. It is an active research area to study situations with interfering units in modern causal inference literature (e.g., Hudgens and Halloran, 2008).

**Same treatment version** can be violated for treatments with complex components. For instance, when studying the effect of cigarette smoking on lung cancer, the type of cigarettes may matter; when studying the effect of college education on income, the type and major of college education may matter.

# Causal quantity

- Any causal quantity is a function of potential outcomes

$$\log Y_i(1) - \log Y_i(0), \quad \frac{Y_i(1)}{Y_i(0)}, \quad \mathbf{1}\{Y_i(1) > Y_i(0)\}, \quad \text{etc.}$$

① outcome

② practical meaning

③ estimable

# Causal quantity

- Any causal quantity is a function of potential outcomes

$$\log Y_i(1) - \log Y_i(0), \quad \frac{Y_i(1)}{Y_i(0)}, \quad \mathbf{1}\{Y_i(1) > Y_i(0)\}, \quad \text{etc.}$$

- A causal effect is defined to be the comparison of the potential outcomes on the same units

# Causal quantity

- Any causal quantity is a function of potential outcomes

$$\log Y_i(1) - \log Y_i(0), \quad \frac{Y_i(1)}{Y_i(0)}, \quad \mathbf{1}\{Y_i(1) > Y_i(0)\}, \quad \text{etc.}$$

- A causal effect is defined to be the comparison of the potential outcomes on the same units

- Fundamental problem of causal inference
  → only one potential outcome is observed
  - we never see both $Y_i(1)$ and $Y_i(0)$
  - most features of $Y_i(1) - Y_i(0)$ are not point identified, e.g., $\text{pr}\{Y_i(1) - Y_i(0) \leq 0\}$

  //
  estimable

$$p(\, \xi_i(1) - \xi_i(0) \leq 0 \,)$$

$$\uparrow$$

$$(\, \xi_i(1), \; \xi_i(0) \,)$$

marginal dist.

# Average causal effect

- Individual causal effect: $Y_i(1) - Y_i(0) \rightarrow$ difficult to estimate

# Average causal effect

- Individual causal effect: $\underline{Y_i(1) - Y_i(0)}$ -> difficult to estimate
- Average causal effect (ACE): $E\{Y_i(1) - Y_i(0)\}$

$$= E\, Y_i(1) - E\, Y_i(0)$$

$$= \int x\, dF_{Y_i(1)}(x) - \int x\, dF_{Y_i(0)}(x)$$

$$F(Y_i(1), Y_i(0))$$

# Average causal effect

- Individual causal effect: $Y_i(1) - Y_i(0)$ -> difficult to estimate
- Average causal effect (ACE): $E\{Y_i(1) - Y_i(0)\}$
- $E\{Y_i(1)\} \neq E\{Y_i \mid Z_i = 1\}$ (when will they be the same?)
  - $E\{Y_i(1)\}$: average of $Y_i(1)$ for units 1 to 5
  - $E(Y_i \mid Z_i = 1) = E\{Y_i(1) \mid Z_i = 1\}$: average of $Y_i(1)$ for units 1,3,4

*differed cohorts*

| Unit $i$ | $Z_i$ | $Y_i(1)$ | $Y_i(0)$ |
|----------|-------|----------|----------|
| 1 | 1 | 0 | 1 |
| 2 | 0 | 0 | 1 |
| 3 | 1 | 0 | 0 |
| 4 | 1 | 1 | 1 |
| 5 | 0 | 1 | 0 |

| Unit $i$ | $Z_i$ | $Y_i$ |
|----------|-------|-------|
| 1 | 1 | 0 |
| 2 | 0 | 1 |
| 3 | 1 | 0 |
| 4 | 1 | 1 |
| 5 | 0 | 0 |

$\bar{E} Y_{i(1)}$ : average potential outcome

receiving $1$ for $\boxed{all \quad patients}$

$E [ Y_i | z_i = 1 ]$ : average response

for $\boxed{all \quad patients \quad receiving \quad 1}$

$Y_i = z_i Y_{i(1)} + (1 - z_i) Y_{i(0)}$

$\longrightarrow E [ \quad Y_{i(1)} \qquad\qquad | z_{i=1})$

$= E [ Y_{i(1)} | z_i = 1 ]$

$\neq E [ Y_{i(1)} ]$

$z_i \perp Y_{i(1)}$

# Other causal quantities of interest

$$ATT \stackrel{?}{=} ATU \stackrel{?}{=} ATE$$

- Average treatment effect on the treated (ATT) and on the untreated (ATU): $E\{Y_i(1) - Y_i(0) \mid Z_i = 1\}$, $E\{Y_i(1) - Y_i(0) \mid Z_i = 0\}$

- Heterogeneous effects:
  - conditional average causal effect: $ACE(\mathbf{x}) = E\{Y_i(1) - Y_i(0) \mid \mathbf{X}_i = \mathbf{x}\}$
  - applications to precision medicine

- Non-additive effects:
  - quantile treatment effects, e.g.,

    $\text{median}\{Y_i(1) - Y_i(0)\}$ or $\text{median}\{Y_i(1)\} - \text{median}\{Y_i(0)\}$

  - odds ratio

    $$\frac{\text{pr}\{Y_i(1) = 1\}/\text{pr}\{Y_i(1) = 0\}}{\text{pr}\{Y_i(0) = 1\}/\text{pr}\{Y_i(0) = 0\}}$$

# Causal effect is comparison of potential outcomes

- Let Z = **1**(Take Aspirin at 3 pm). Which of the following qualifies/qualify as a causal effect?

  *not potential outcome*

  (A)  E(temperature | Z = 1) − E(temperature | Z = 0)  ✗

  (B)  E(potential pain scale at 4 pm with Aspirin | Z = 1) − E(potential pain scale at 4 pm without Aspirin | Z = 0)  ✗

  *On the same*

  (C)  E(potential pain scale at 2 pm with Aspirin) − E(potential pain scale at 2 pm without Aspirin)  ✗  *not potential outcome*

  (D)  my body temperature after taking Aspirin - my body temperature before taking Aspirin  ✗  *not potential outcome*

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.
- Can immutable characteristics have meaningful causal effects?

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.
- Can immutable characteristics have meaningful causal effects?

- Strategies:

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.
- Can immutable characteristics have meaningful causal effects?

- Strategies:
  1. causal effects of perceived characteristics:
     - Causal effect of a job applicant's gender/race on call-back rates (Bertrand and Mullainathan, 2004)

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.
- Can immutable characteristics have meaningful causal effects?

- Strategies:
  1. causal effects of perceived characteristics:
     - Causal effect of a job applicant's gender/race on call-back rates (Bertrand and Mullainathan, 2004)
  2. reinterpretation
     - Causal effect of having a female politician on policy outcomes (Chattopadhyay and Duflo, 2004)

# Causal effects of immutable characteristics

- "No causation without manipulation" (Holland, 1986)
- Immutable characteristics or attributes: gender, race, age, etc.
- Can immutable characteristics have meaningful causal effects?

- Strategies:
  1. causal effects of perceived characteristics:
     - Causal effect of a job applicant's gender/race on call-back rates (Bertrand and Mullainathan, 2004)
  2. reinterpretation
     - Causal effect of having a female politician on policy outcomes (Chattopadhyay and Duflo, 2004)
  3. redefinition:
     - Race as a "bundle of sticks": skin color, neighborhood, socio-economic status, etc. (Sen and Wasow, 2016)

# Resolving Simpson's paradox

|  | small stones | | large stones | |
| --- | --- | --- | --- | --- |
|  | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for ~~large~~ stones)

  *Small*

# Resolving Simpson's paradox

|  | small stones | | large stones | |
| --- | --- | --- | --- | --- |
|  | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for ~~large~~ stones)
  *small*

- Simpson's paradox:
  - $\hat{E}(Y_i \mid Z_i = 1, X_i = x) > \hat{E}(Y_i \mid Z_i = 0, X = x)$ for $x = 0, 1$
  - $\hat{E}(Y_i \mid Z_i = 1) < \hat{E}(Y_i \mid Z_i = 0)$

# Resolving Simpson's paradox

|  | small stones | | large stones | |
| --- | --- | --- | --- | --- |
|  | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for large stones)

- Simpson's paradox:
  - $\hat{E}(Y_i \mid Z_i = 1, X_i = x) > \hat{E}(Y_i \mid Z_i = 0, X = x)$ for $x = 0, 1$
  - $\hat{E}(Y_i \mid Z_i = 1) < \hat{E}(Y_i \mid Z_i = 0)$
  - the sign of association may be reversed when adding covariates

# Why Association Fail?

Small –Stone
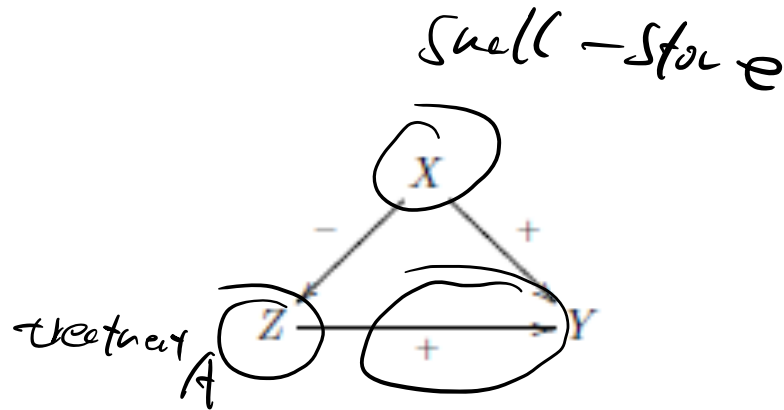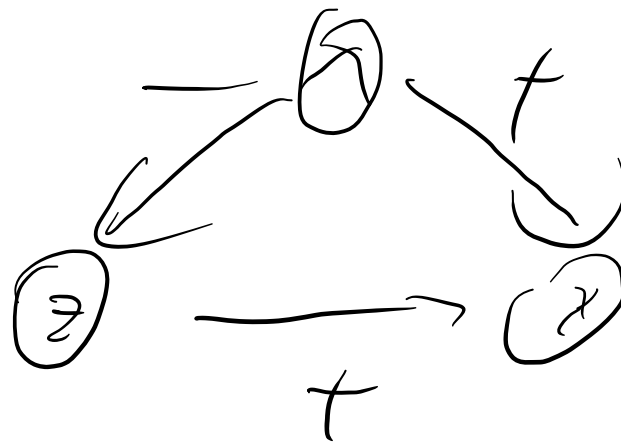
X

$-$      $+$
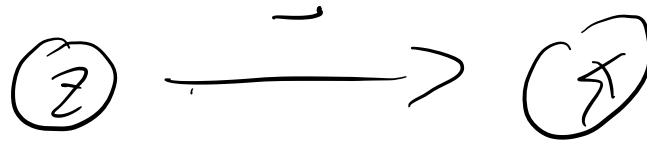
treatment A (Z) ——→ Y

$+$

FIGURE 1.1: A diagram for the kidney stone example. The signs indicate the associations of two variables, conditioning on other variables pointing to the downstream variable.

- Patient with larger stones tends to take treatment A
- Patients with smaller stones have higher success probability.

# Resolving Simpson's paradox

| | small stones | | large stones | |
|---|---|---|---|---|
| | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for large stones)

- Can Simpson's paradox happen using ACE instead of success rate?

$$A C \bar{E} = \bar{E} ( Y_{(1)} - Y_{(0)} )$$

$$A C E_{X=0} = \bar{E} ( Y_{(1)} - Y_{(0)} \mid X = \textcircled{0} )$$

$$ACE = E(Y(1) - Y(0))$$

$$= P(X=1) \, E(Y(1) - Y(0) \mid X=1)$$

$$+ P(X=0) \, E(Y(1) - Y(0) \mid X=0)$$

$$= \sum_{i=0,1} ACE_{X=i} \; \boxed{P(X=i)}$$

$$ACE_{X=i} > 0$$

$$\Rightarrow ACE > 0$$

# Resolving Simpson's paradox

|  | small stones | | large stones | |
|---|---|---|---|---|
|  | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for large stones)

- Can Simpson's paradox happen using ACE instead of success rate?
  - $E\{Y_i(1) \mid X_i = x\} > E\{Y_i(0) \mid X = x\}$ for $x = 0, 1$
  - $E\{Y_i(1)\} < E\{Y_i(0)\}$?

$$E(Y_i(1) - Y_i(0))$$

$$= \sum_{x=0,1} P(X_i = x) \, E(Y_i(1) - Y_i(0) \mid X = x)$$

# Resolving Simpson's paradox

|  | small stones | | large stones | |
|---|---|---|---|---|
|  | success | fail | success | fail |
| Treatment A | 81 | 6 | 192 | 71 |
| Treatment B | 234 | 36 | 55 | 25 |

- Treatment $Z_i$ (1 for A); outcome $Y_i$ (1 for success); covariate $X_i$ (1 for large stones)

- Can Simpson's paradox happen using ACE instead of success rate?
  - $E\{Y_i(1) \mid X_i = x\} > E\{Y_i(0) \mid X = x\}$ for $x = 0, 1$
  - $E\{Y_i(1)\} < E\{Y_i(0)\}$?
- Simpson's paradox cannot happen for ACE; Is treatment A better than treatment B?

$$E(\, Y_{i(0)} \,) \qquad vs \qquad E(\, Y_i \mid Z_i = 1)$$

$$\|$$

$$E\big(\, Z_i\, Y_{i(1)} + (1 - Z_i)\, Y_{i(0)} \mid Z_i \big)$$
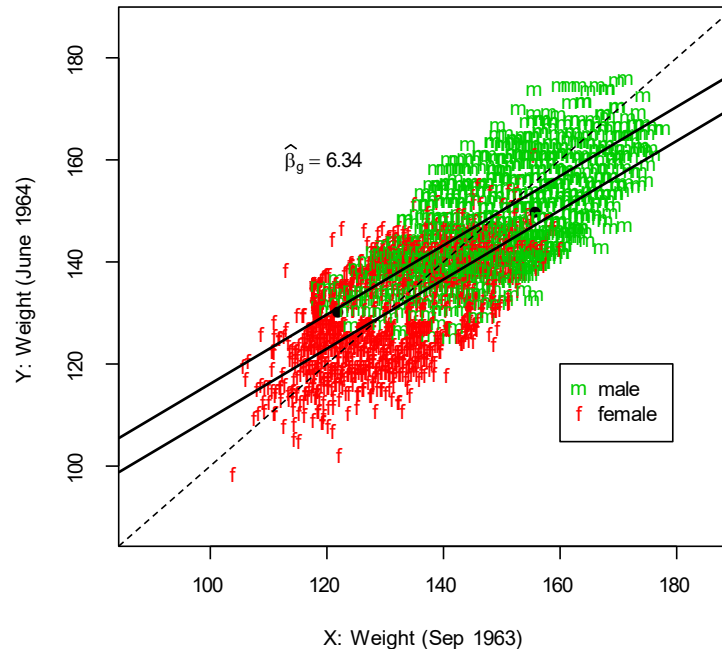
$$\|$$

$$E(\, Y_{i(1)} \mid Z_i = 1)$$

Association: difference after
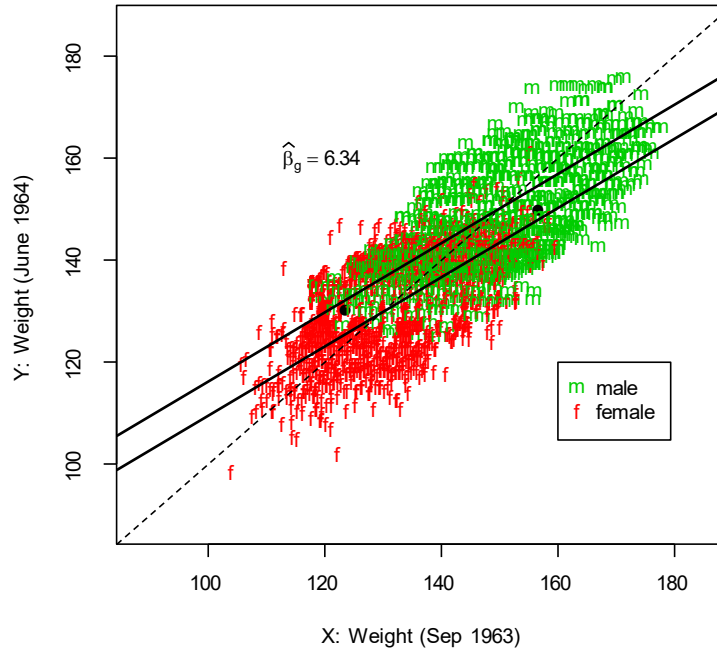
.. before treatment

# Lord's paradox (Lord, 1967)

- Question: are the effects of the diet provided in the dining hall different for males and females?
- Data: gender $G_i$; weight in 1963 $X_i$; weight in 1964 $Y_i$

# Lord's paradox (Lord, 1967)

- Question: are the effects of the diet provided in the dining hall different for males and females?
- Data: gender $G_i$; weight in 1963 $X_i$; weight in 1964 $Y_i$
- $E(Y_i \mid G_i = 1) = E(X_i \mid G_i = 1) = 150$
- $E(Y_i \mid G_i = 0) = E(X_i \mid G_i = 0) = 130$

# Lord's paradox (Lord, 1967)



- Statistician A: average weights unchanged for both males and females
- Statistician B: $Y_i = \beta_0 + \beta_g G_i + \beta_X X_i + E_i \quad \beta_g = 6.34$
- What is the interpretation of $\beta_g$
- Who is correct?

$$E(Y \mid X, G=1)$$
$$- E(Y \mid X, G=0)$$
$$= 6.34$$

# Resolving Lord's paradox

- Formulation
    - treatment $Z_i$ (1 for dining)
    - pre-treatment: gender $G_i$ (1 for male); weight in 1963 $X_i$
    - post-treatment: weight in 1964 $Y_i$
    - potential outcomes: $Y_i(1)$ and $Y_i(0)$

# Resolving Lord's paradox

- Formulation
  - treatment $Z_i$ (1 for dining)
  - pre-treatment: gender $G_i$ (1 for male); weight in 1963 $X_i$
  - post-treatment: weight in 1964 $Y_i$
  - potential outcomes: $Y_i(1)$ and $Y_i(0)$

- Causal quantity: $\Delta_g = E\{Y_i(1) - Y_i(0) \mid G_i = g\}$ for $g = 0, 1$
  - difference between males and females: $\Delta_1 - \Delta_0$

# Resolving Lord's paradox

- $E\{Y_i(1) \mid G_i = g\} = E(Y_i \mid G_i = g)$, $E\{Y_i(0) \mid G_i = g\} =$???

# Resolving Lord's paradox

- $E\{Y_i(1) \mid G_i = g\} = E(Y_i \mid G_i = g)$, $E\{Y_i(0) \mid G_i = g\} =$ ???

- $Y_i(0)$ is missing for all units -> no conclusion without assumptions about $Y_i(0)$ (identifiability issue)

# Difference in Difference

1963     diet $\longrightarrow$     1964

1963     no diet $\longrightarrow$     1964

- Statistician A: $Y_i(0) = X_i \rightarrow \Delta_1 - \Delta_0 = 0$

$$E( Y_i - X_i \mid G_i = g )$$

$$= \bar{E}( Y_i(0) \mid G_i = g )$$

$Y_i(0)$
$= X_i$

$$- E( X_i \mid G_i = g )$$

$$\overset{?}{=} E( Y_i(1) - Y_i(0) \mid G_i = g )$$

# Resolving Lord's paradox

- Statistician A: $Y_i(0) = X_i$ -> $\Delta_1 - \Delta_0 = 0$

- Statistician B: $Y_i = \beta_0 + \beta_g G_i + \beta_X X + E_i$

- Statistician A: $Y_i(0) = X_i$ -> $\Delta_1 - \Delta_0 = 0$

- Statistician B: $Y_i = \beta_0 + \beta_g G_i + \beta_X X + E_i$
    - $E\{Y_i(1) \mid X_i, G_i = g\} = a_g + bX_i$ -> $a_1 - a_0 = \beta_g$
    - $Y_i(0) = a + bX_i$

$$\Delta_1 - \Delta_0 \,?$$

$$E\left\{ Y_{i}^{(c)} \mid X, G_i = g \right\} = a_g + bX_i$$

$$\Delta_g \left[ E(Y_{i(1)} \mid G_i = g) - E(Y_{i(0)} \mid G_i = g) \right]$$

$$\parallel$$

$$E(a_g + b X_i \mid G_i = g)$$

$$\Delta_1 - \Delta_0$$

$$\Rightarrow = a_1 - a_0 + \left[ E(b X_i - Y_{i(0)} \mid G_i = g) - E(b X_i - Y_{i(0)} \mid G_i = 0) \right]_g$$

regression coef

$$\xrightarrow{?} \quad \boxed{\beta_g = a_1 - a_0} \doteq \underbrace{\left( \Delta_1 - \Delta_0 \right)}_{\text{difference}}$$

$$ACE$$

# Resolving Lord's paradox

- Statistician A: $Y_i(0) = X_i \to \Delta_1 - \Delta_0 = 0$

- Statistician B: $Y_i = \beta_0 + \beta_g G_i + \beta_X X + E_i$
  - $E\{Y_i(1) \mid X_i, G_i = g\} = a_g + b X_i \to a_1 - a_0 = \beta_g$
  - $Y_i(0) = a + b X_i$

- Both statisticians' conclusions depend on untestable assumptions

# Identification links thought experiment and data

- The target parameters, as defined by potential outcomes, is a function of the unobservables

- Question of identification: what can we learn about this function from the observed data?

- Identification maps assumptions and data to information about target parameters; Which causal quantity is identifiable?

- A parameter is identified if, under the stated assumptions, alternative values of the parameter implies different distributions of observable data

- Identification is a binary property

- In order to achieve identification, assumptions are unavoidable, but we need to figure out what assumptions are plausible in practice

# Statistical inference links population and sample

- In practice, we only see a finite sample of the observables
- We do not know the population distribution of data
- Statistical inference: using the sample to infer about the population

- It is useful to separate identification from statistical inference

Sample $\quad\longrightarrow\quad$ Population $\quad\longrightarrow\quad$ Target parameters

statistical inference $\qquad\qquad$ identification

- Identification: how much can you learn about the quantities of interest if you had an infinite amount of data?

- We will keep returning to these two steps in the whole semester

# Summary

- Causation: comparison of potential outcomes for the same unit(s)
- Causal quantity is a function of potential outcomes

- Fundamental problem of causal inference: only one potential outcome is observed
- Identification links thought experiment and data
- Statistical inference links population and sample