

Modern Methods of Causal Inference

Cause & Effect

The effect of a cause

Cause of an effect X

• Smoking causes lung cancer at the age of 65 for men

$$z_i = \begin{cases} 1 & \text{if the person smokes} \\ 0 & \text{otherwise} \end{cases}$$

$$y_i = \begin{cases} 1 & \text{if the person developed lung cancer} \\ 0 & \text{otherwise} \end{cases}$$

$$(z_i=1, y_i=0)$$

• Potential outcomes model (for mutable causes)

i has two potential outcomes

$y_i(1)$ if they were "treated"

$y_i(0)$ if they were "not treated"

$$(z_i=1, y_i=0) \rightarrow y_i(1)=0 \quad \left. \begin{array}{l} y_i(1)-y_i(0)=0 \end{array} \right\}$$

$$(z_i=0, y_i=0) \rightarrow y_i(0)=0$$

$$(z_j=0, y_j=1) \rightarrow y_j(0)=1 \quad \left. \begin{array}{l} y_j(1)-y_j(0)=-1 \end{array} \right\}$$

$$(z_j=1, y_j=0) \rightarrow y_j(1)=0$$

$$\left. \begin{array}{l} \frac{1}{2}(y_i(1)-y_i(0))+(y_j(1)-y_j(0)) \end{array} \right\}$$

Average of two ICES

Individual causal effect

Assumption: Each individual smokes with probability $\frac{1}{2}$ ($P(z_i)=P(z_j)$) and z_i, z_j are independent of $y_i(0), y_i(1), y_j(0), y_j(1)$,

$$z_i = \begin{cases} 1 & \text{w/p } \frac{1}{2} \\ 0 & \text{w/p } \frac{1}{2} \end{cases}$$

Consider the case where $z_i + z_j = 1$.

$$\begin{aligned} y_i z_i + y_j z_j &= (y_i(1)z_i + y_i(0)(1-z_i)) \cdot z_i + (y_j(1)z_j + y_j(0)(1-z_j)) \cdot z_j \\ &= y_i(1)z_i + y_j(1)z_j \end{aligned}$$

$$E[y_i z_i + y_j z_j | y_i(1), y_i(0), y_j(1), y_j(0), z_i + z_j = 1] = \frac{y_i(1) + y_j(1)}{2}$$

$$\text{Try: } E[y_i(1-z_i) + y_j(1-z_j) | y_i(1), y_i(0), y_j(1), y_j(0), z_i + z_j = 1] = \frac{y_i(0) + y_j(0)}{2}$$

(*) $y_i z_i + y_j z_j - (y_i(1-z_i) + y_j(1-z_j))$ is in expectation equal to $\frac{1}{2}((y_i(1)-y_i(0)) + (y_j(1)-y_j(0)))$

Living in Atlanta:

| | | Lung cancer at 65 | | |
|---------|-----|-------------------|-----|------|
| | | Yes | No | |
| Smoking | Yes | 150 | 5 | 155 |
| | No | 2500 | 122 | 2622 |

$$\frac{5}{150+5} - \frac{122}{2500+122} = -0.014$$

$$\text{proportion of smokers } \frac{155}{155+2622} = 0.056$$

Adding up the two tables above:

| | | Lung cancer at 65 | | |
|---------|-----|-------------------|------|--|
| | | Yes | No | |
| Smoking | Yes | 730 | 355 | |
| | No | 4770 | 2122 | |

$$\frac{355}{355+730} - \frac{2122}{2122+4770} = 0.019 \rightarrow \text{You should smoke (LOL)}$$

Living in Orlando:

| | | Lung cancer at 65 | | |
|---------|-----|-------------------|------|------|
| | | Yes | No | |
| Smoking | Yes | 580 | 350 | 930 |
| | No | 2270 | 2000 | 4270 |

$$\frac{350}{350+580} - \frac{2000}{2000+2270} = -0.092$$

$$\text{proportion of smokers } \frac{930}{930+4270} = 0.178$$

Independence assumption fails

leading to a "ridiculous" result

Covariates: Measurement not subject to the choice of the treatment the individual makes (X_i)

Assumption: Conditioned on X_i , the treatment is independent of the potential outcomes

$$Z_i \perp\!\!\!(1) (Y_{i(1)}, Y_{i(0)}) \mid X_i$$

(" $\perp\!\!\!(1)$ " denotes Philip Dawid's (1979) notation for independence)

Consider K binary covariates. If possible, $2n$ independently chosen people, exactly half of them are in $\{Z=1\}$

$$D_k = \frac{1}{n} \sum_{i=1}^{2n} Z_i X_i^{(k)} - \frac{1}{n} \sum_{i=1}^{2n} (1-Z_i) X_i^{(k)}, k=1, 2, \dots, K$$

$$\Pr(\max_{1 \leq k \leq K} |D_k| \geq t) \leq 2K \exp(-\frac{zt^2}{2t/n})$$

In theory we just need $X_i = (Y_{i(1)}, Y_{i(0)})$

Instead collect information about confounders X_i 's that are informative about your potential outcomes and related to the treatment assignment

i, j exactly one is treated and the other is control

Denote the treated unit by t and the control unit by c $Y_t - Y_c = (\star)$

Pair Matching

$T = \{t_1, \dots, t_T\}$ are my treated units

$C = \{c_1, \dots, c_C\}$ are my control units

Assume $C = |C| > T = |T|$ (For example, more ppl don't smoke)

Then a pair matching is a function $m: T \rightarrow C$ which take t_i to $m(t_i) \in C$.

δ_{tc} is a measure of the distance between X_t and X_c

- $\|X_t - X_c\|_2$
- $(X_t - X_c)^T S^{-1} (X_t - X_c)$ (Mahalanobis distance)
- $|F^{-1}(X_t) - F^{-1}(X_c)|$

$$\min \sum_{t \in T} \delta_{t m(t)}$$

→ Nearest-neighbor approach - one with replacement match

→ Start with the smallest pair (t, c) , then remove it and repeat T times
[Greedy approach]

- one without replacement match

e.g. ε is a small number ($\varepsilon > 0$)

| | c_1 | c_2 | c_3 |
|-------|----------------|----------------|----------------|
| t_1 | ε | 0 | 2ε |
| t_2 | 6ε | ε | 0 |
| t_3 | 6ε | 2ε | ε |

matrix of δ_{tc} values

Greedy approach: (t_1, c_2)

(t_2, c_3)
 (t_3, c_1)

Total $6\varepsilon \rightarrow$ Not ideal

(6ε can be any large number)

600 pair matches $\frac{1}{600} \sum_{t \in T} (Y_t - Y_{m(t)})^2$

$\text{Var}(Y_{m(t)})$ part) without repeat with repeat

$$\frac{\sigma^2}{600} < \frac{\sigma^2}{(600)^2} \sum_{c \in C} r_c^2 \quad r_c \geq 0, \sum_c r_c = T = 600$$

We don't want repeat controls

The optimal pair matching problem

$$a_{tc} = \begin{cases} 0, & \text{if } t \text{ is not matched to } c \\ 1, & \text{if } t \text{ is matched to } c \end{cases}$$

$$\min \sum_{t \in T} \sum_{c \in C} \delta_{tc} \cdot a_{tc} \quad \text{s.t.} \quad \begin{cases} \sum_t a_{tc} \leq 1 \text{ for all } c \\ \sum_c a_{tc} = 1 \text{ for all } t \end{cases}$$

Solving via the minimum network flow problem

① Network of warehouses

② Connection between the warehouses

③ Cost of moving a product between warehouses

④ Capacity constraint of each connection

$$N \subseteq N \times N$$

f_e is the cost per unit in connection $e \in E$

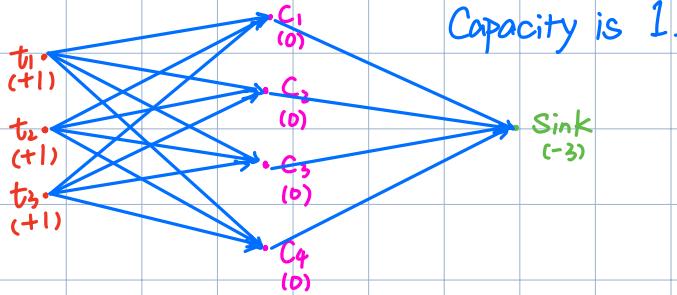
w_e is the capacity constraint

s_n is the supply/demand of warehouse $n \in N$

k_e for $e \in E$ are integers (amount of product to be moved for connection e)

$$\begin{aligned} \min \sum f_e k_e \\ \text{s.t. } \begin{cases} k_e \leq w_e & \text{going out} \\ -s_n = \sum_{e=(n,m)} k_e - \sum_{e=(m,n)} k_e & \text{coming in} \end{cases} \\ \text{for any } m \quad \text{for any } m \end{aligned}$$

$$T = \{t_1, t_2, t_3\} \quad C = \{C_1, C_2, C_3, C_4\}$$

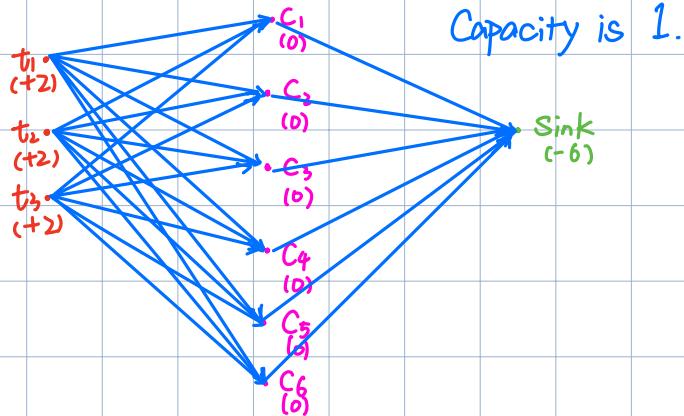


RELAX - IV Optmatch in R
· C call relax

$$\begin{aligned} (|N| \times |E| + (|E| + |N|) \log(|N|) \sum_{n \in N} |S_n|) &\leq \text{const. } C^3 \\ T &\downarrow \\ TC + C &\leq C^2 \\ \leq 4C^3 &\leq \sim C^3 \\ 2T &\downarrow \end{aligned}$$

1-k matching

$$T = \{t_1, t_2, t_3\} \quad C = \{C_1, C_2, C_3, C_4, C_5, C_6\}$$



Variable ratio matching
(one to various numbers)

$$\hat{P}(z_{t_i=0} | X_t) \times C$$

Full match

$(\mathcal{T}_s, \mathcal{C}_s)$ for $s = 1, \dots, S$

$$|\mathcal{T}_s| = 1 \text{ or } |\mathcal{C}_s| = 1$$

$$\bigcup_{s=1}^S \mathcal{T}_s = \mathcal{T} \quad \bigcup_{s=1}^S \mathcal{C}_s = \mathcal{C}$$

matching is observational design
 matching is non-experimental design

Consider $\mathcal{T} = \{t_1, \dots, t_T\}$, $\mathcal{C} = \{c_1, \dots, c_C\}$

$m: \mathcal{T} \rightarrow \mathcal{C}$ defines a pair matching

- 1-1 match:

$$(t_1, m(t_1)), (t_2, m(t_2)), \dots, (t_T, m(t_T))$$

$$y_{t_1} - y_{m(t_1)}, y_{t_2} - y_{m(t_2)}, \dots, y_{t_T} - y_{m(t_T)}$$

- 1-2 match

$$(t_1, m_1(t_1), m_2(t_1)), \dots, (t_T, m_1(t_T), m_2(t_T))$$

$$y_{t_1} - \frac{1}{2}(y_{m_1(t_1)} + y_{m_2(t_1)}), \dots, y_{t_T} - \frac{1}{2}(y_{m_1(t_T)} + y_{m_2(t_T)})$$

- 1-variable ratio

Potential Outcomes notation

$y_{i(1)}, y_{i(0)} \rightarrow y_{i(1)} - \underbrace{y_{i(0)}}_{\text{by matching, we try to impute this unknown term}}$ individual treatment effect for i

Think of $\hat{y}_{t_i(0)} = y_{m(t_i)}$

$y_{t_i(0)}$ is counterfactual outcome

We constructed a counterfactual estimate using matching

$\frac{1}{T} \sum_{t=1}^T (y_t - y_{m(t)})$ is our final matching estimator. What is it estimating?

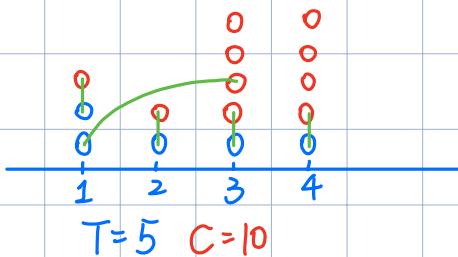
$\frac{1}{T} \sum_{t=1}^T (y_{t(1)} - y_{t(0)})$ corresponds to $E(y_{i(1)} - y_{i(0)} | z_i = 1)$

When the treatment effect is different across the population,

$E(y_{i(1)} - y_{i(0)}), E(y_{i(1)} - y_{i(0)} | z_i = 1), E(y_{i(1)} - y_{i(0)} | z_i = 0)$ are different

$\Delta \frac{1}{T} \sum_{t=1}^T (y_t - y_{mt})$ is NOT a consistent estimator of $E(y_{i(1)} - y_{i(0)} | z_i = 1)$

$T < C, 2T = C$

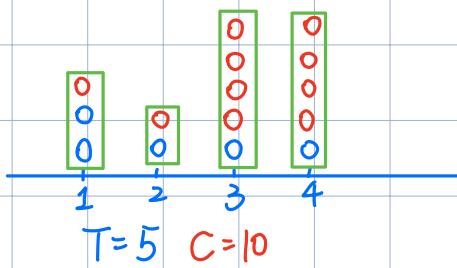


We need $\Pr(z_i = 1 | X) < \frac{1}{2}$ for all X .

Reference: Savje (2022) On inconsistency of matching without replacement Biometrika

Full Matching

In each matched set, there is either exactly one treatment unit or exactly one control unit. It will give consistent estimator for $E(y_{i(1)} - y_{i(0)})$



$E(y_{i(1)} - y_{i(0)}) \rightarrow$ Average treatment effect (ATE)

$E(y_{i(1)} - y_{i(0)} | z_i = 1) \rightarrow$ Average treatment effect on the treated (ATT)

$E(y_{i(1)} - y_{i(0)} | z_i = 0) \rightarrow$ Average treatment effect on the control (ATC)

Balance

$\frac{1}{T} \sum g(X_t) - \frac{1}{T} \sum g(X_{mt})$ is close to zero for many functions g .

$$\cdot g(x) = x \Rightarrow \bar{X}_t - \bar{X}_{mt}$$

Cardinality Matching targeting for balance.

$$f_k : -b_k \sum_{t=1}^T \sum_{c=1}^C a_{tc} \leq \sum_{t=1}^T \sum_{c=1}^C (g(X_t) - g(X_c)) a_{tc} \leq b_k \sum_{t=1}^T \sum_{c=1}^C a_{tc}$$

where $a_{tc} = 1_{\{t \text{ is matched to } c\}}$

Solve for a_{tc}

$$\max \sum_{t=1}^T \sum_{c=1}^C a_{tc} \quad \text{s.t.} \quad \begin{cases} f_k \text{'s} \\ a_{tc} \text{'s are } 0, 1 \\ \sum_{c=1}^C a_{tc} \leq 1, \sum_{t=1}^T a_{tc} \leq 1 \end{cases}$$

Standardized Mean Difference

$$\left| \frac{\bar{X}_t - \bar{X}_{mt}}{\sqrt{\frac{s_e^2 + s_t^2}{2}}} \right|$$

More on matching

Assumption: $(y_{i(1)}, y_{i(0)}) \perp\!\!\!\perp z_i | X_i$ (1)

Suppose I can find a lower dim $b(X_i)$ so that $(y_{i(1)}, y_{i(0)}) \perp\!\!\!\perp z_i | b(X_i)$ (2)

Theorem: There is a one dim covariate $b(X_i)$ for which (1) implies (2).

Lemma 1: Suppose $z_i \perp\!\!\!\perp X_i | b(X_i)$ for some $b(X_i)$, then (1) implies (2).

Lemma 2: Let $\lambda(x) = \Pr(z_i=1 | X_i=x)$. Then $z_i \perp\!\!\!\perp X_i | b(X_i)$ for some function $b(X_i)$
iff $\lambda(x)=f(b(X_i))$ for some f

* In light of the two lemmas, the theorem is proved by showing that

$$(y_{i(1)}, y_{i(0)}) \perp\!\!\!\perp z_i | \lambda(X_i)$$

Pf of lemma 2:

• "if": Suppose $\lambda(X_i)=f(b(X_i))$.

$$\Pr(z_i=1 | X_i, b(X_i)) = \Pr(z_i=1 | X_i) = \lambda(X_i)$$

$$\Pr(z_i=1 | b(X_i)) = E(\Pr(z_i=1 | X_i, b(X_i)) | b(X_i)) = E(\lambda(X_i) | b(X_i)) = f(b(X_i)) = \lambda(X_i)$$

Hence $z_i \perp\!\!\!\perp X_i | b(X_i)$

• "only if": Suppose $z_i \perp\!\!\!\perp X_i | b(X_i)$ but $\lambda(x)$ is NOT $f(b(x))$ for any f .

Then there is a pair X_1, X_2 so that $b(X_1)=b(X_2)$ but $\lambda(X_1) \neq \lambda(X_2)$, i.e.,

$$\Pr(z_i=1 | X_i=X_1) \neq \Pr(z_i=1 | X_i=X_2) \quad \text{Hence the contradiction.}$$

Pf of lemma 1:

Sufficient to prove $\Pr(z_i=1 | y_{i(1)}, y_{i(0)}, b(X_i)) = \Pr(z_i=1 | b(X_i))$.

$$\text{LHS} = E[\Pr(z_i=1 | y_{i(1)}, y_{i(0)}, X_i, b(X_i)) | y_{i(1)}, y_{i(0)}, b(X_i)]$$

$$= E[\Pr(z_i=1 | X_i, b(X_i)) | y_{i(1)}, y_{i(0)}, b(X_i)]$$

$$= E[\lambda(X_i) | y_{i(1)}, y_{i(0)}, b(X_i)] \xrightarrow{\text{Lemma 2}} \lambda(X_i)$$

m data points (X_i, Z_i, Y_i)

Fit a logistic regression of Z_i on X_i

$$\hat{Z}(X_i) = \Pr(Z_i=1|X_i)$$

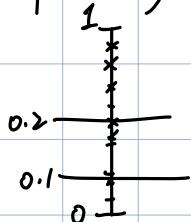
How to incorporate the fitted values in matching?

$$\delta'_{tc} = \delta_{tc} / \mathbb{1}(|\hat{Z}(x_t) - \hat{Z}(x_{m(t)})| < 0.2) = \begin{cases} \delta_{tc}, & \text{if } |\hat{Z}(x_t) - \hat{Z}(x_{m(t)})| < 0.2 \\ \infty, & \text{aw.} \end{cases}$$

$\hat{Z}(X_i) = \Pr(Z_i=1|X_i)$ is "Propensity Score"

δ'_{tc} is called "Propensity Score Caliper"

Propensity Score Stratification



If using 0.2, we have $\frac{1}{0.2} = 5$ groups (strata)

The strata are defined according to the quantiles of propensity score

Weighting

$$ATT = E(Y_i(1) - Y_i(0) | Z_i = 1)$$

$$\frac{1}{T} \sum_{t \in T} Y_t - \frac{1}{T} \sum_{t \in T} Y_{m(t)}$$

$$\hat{E}(Y_i(1) | Z_i = 1) \quad \hat{E}(Y_i(0) | Z_i = 1)$$

$$\frac{1}{T} \sum_{c \in C} w_c Y_c$$

Suppose we want an estimator of $E(Y_i(0) | Z_i = 1)$ of the form $\frac{1}{T} \sum_{c \in C} w_c Y_c$.

$$\begin{aligned} \mu_0 &= E(Z_i Y_i(0)) = E(Z_i Y_i(0) | Z_i = 1) \Pr(Z_i = 1) + E(Z_i Y_i(0) | Z_i = 0) \Pr(Z_i = 0) \\ &= E(Y_i(0) | Z_i = 1) \cdot p + 0 \end{aligned}$$

using $\hat{\mu}_0 = \frac{1}{n} \sum_{c \in C} w_c Y_c = \frac{1}{n} \sum_{i=1}^n (1 - z_i) w_i Y_i$ when $n = T + C$.

Let $m_0(X_i) = E(Y_i(0) | X_i)$.

$$\mu_0 = E(Z_i Y_i(0)) = E[E(Z_i Y_i(0) | X_i)] = E[\Pr(Z_i = 1 | X_i) E(Y_i(0) | X_i)] = E(\hat{Z}(X_i) m_0(X_i))$$

Let $\varepsilon_i = Y_i(0) - m_0(X_i)$.

$$\hat{\mu}_0 - \mu_0 = \frac{1}{n} \sum_{i=1}^n (1-z_i) w_i m_0(x_i) - \frac{1}{n} \sum_{i=1}^n z_i m_0(x_i)$$

$$+ \frac{1}{n} \sum_{i=1}^n (1-z_i) w_i \varepsilon_i \quad \rightarrow \text{has } E(\cdot | x_i) = 0$$

$$+ \frac{1}{n} \sum_{i=1}^n z_i m_0(x_i) - E(z_i y_i(0)) \quad \rightarrow \text{has } E(\cdot | x_i) = 0$$

We can estimate the above bias by minimizing

$$\left| \frac{1}{n} \sum_{i=1}^n (1-z_i) w_i m_0(x_i) - \frac{1}{n} \sum_{i=1}^n z_i m_0(x_i) \right|, \text{ if } m_0(\cdot) \text{ is known}$$

$$\max_{m \in \mathcal{M}} \left| \frac{1}{n} \sum_{i=1}^n (1-z_i) w_i m(x_i) - \frac{1}{n} \sum_{i=1}^n z_i m(x_i) \right|, \text{ o.w.}$$

Suppose $m(x_i) = x_i$ one dimensional.

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (1-z_i) w_i x_i}_{\frac{1}{n} \bar{X}_{w, \text{control}}} - \underbrace{\frac{1}{n} \sum_{i=1}^n z_i x_i}_{\frac{1}{n} \bar{X}_t}$$

Covariate Balancing Methods

We say that $f(x)$ is balanced in a pair match if

$$\frac{1}{T} \sum_{t \in T} f(x_t) \approx \frac{1}{T} \sum_{t \in T} f(x_{mt})$$

In full match we want to balance for the population average.

Balance targeted procedure

Maximum Cardinality Matching

$$B_k : -b_k \sum_{t \in T} \sum_{c \in C} a_{tc} \leq \sum_{t \in T} \sum_{c \in C} a_{tc} g_{tck} \leq b_k \sum_{t \in T} \sum_{c \in C} a_{tc} \quad a_{tc} = \begin{cases} 1 & \text{if } t \text{ is matched to } c \\ 0 & \text{otherwise} \end{cases}$$

$$g_{tck} = f_k(x_t) - f_k(x_c)$$

$$\max \sum_{t \in T} \sum_{c \in C} a_{tc} \quad \text{subject to } B_k \text{ for } k=1, \dots, K \text{ and } a_{tc} \in \{0, 1\}, \sum_c a_{tc} \leq 1, \sum_t a_{tc} \leq 1$$

Recall propensity score $\pi(x)$

We have data $(y_{i(1)}, y_{i(0)}, x_i, z_i)$ i.i.d. from some population. $y_i = \begin{cases} y_{i(1)}, & \text{if } z_i = 1 \\ y_{i(0)}, & \text{if } z_i = 0 \end{cases}$

$$\pi(x_i) = \Pr(z_i = 1 | x_i) \rightarrow \text{propensity score}$$

$$z_i \perp\!\!\!\perp x_i | \pi(x_i)$$

Ideally we can estimate it from the data using some flexible regression model of z_i on x_i .

$$\begin{aligned} E\left(\frac{z_i f(x_i)}{\lambda(x_i)}\right) &= E\left(E\left(\frac{z_i f(x_i)}{\lambda(x_i)} \mid X_i\right)\right) \\ &= E\left(\frac{f(x_i)}{\lambda(x_i)} \Pr(z_i=1 \mid X_i)\right) \\ &= E(f(x_i)) \text{ as long as } \lambda(x_i) > 0 \text{ with prob. 1} \end{aligned}$$

For any function f , $\frac{1}{n} \sum_{i=1}^n \frac{z_i f(x_i)}{\lambda(x_i)} \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$

Consider $f(x_i) = E(y_i \mid X_i, z_i=1)$.

$$\begin{aligned} E\left(\frac{z_i E(y_i \mid X_i, z_i=1)}{\lambda(x_i)}\right) &= E\left(E(y_i \mid X_i, z_i=1)\right) \\ &\stackrel{\text{"}}{=} E\left(\frac{z_i y_i}{\lambda(x_i)}\right) \quad \stackrel{\text{"}}{=} E(y_{i(1)}) \quad (\text{exercise}) \end{aligned}$$

Lemma. $E\left(\frac{z_i y_i}{\lambda(x_i)}\right) = E(y_{i(1)})$ for $\lambda(x_i) > 0$ a.s.

If we want to estimate $E(y_{i(1)})$, we use $\frac{1}{n} \sum_{i=1}^n \frac{z_i y_i}{\lambda(x_i)}$

Similarly we can estimate $E(y_{i(0)})$ by $\frac{1}{n} \sum_{i=1}^n \frac{(1-z_i)y_i}{1-\lambda(x_i)}$ for $\lambda(x_i) < 1$ a.s.

$\frac{1}{n} \sum_{i=1}^n \left[\frac{z_i y_i}{\lambda(x_i)} - \frac{(1-z_i)y_i}{1-\lambda(x_i)} \right]$ is a reasonable estimator for ATE = $E(y_{i(1)} - y_{i(0)})$ when $\lambda(x_i) \in (0, 1)$ a.s.

Not stable when $\lambda(x_i)$ is close to 0 or 1.

Inverse Propensity Weighted estimator (IPW estimator)

Suppose we estimate $\hat{\lambda}(x_i)$ to make sure that $\frac{1}{n} \sum_{i=1}^n \frac{z_i f(x_i)}{\hat{\lambda}(x_i)} \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$

Choose a model $\lambda_p(x)$ for $\Pr(z_i=1 \mid x)$

The corresponding likelihood is $\prod_{i=1}^n [\lambda_p(x_i)]^{z_i} [1 - \lambda_p(x_i)]^{1-z_i}$

$$\text{s.t. } \left| \frac{1}{n} \sum_{i=1}^n \frac{z_i f(x_i)}{\lambda_p(x_i)} - \frac{1}{n} \sum_{i=1}^n f(x_i) \right| < b_k, \quad k = 1, \dots, K$$

→ Covariate Balancing Propensity Score

• A solution to very small or very large PS estimator is trimming the estimated propensity score.

$$\sum_{i=1}^n z_i w_i f_k(x_i) \approx \frac{1}{n} \sum_{i=1}^n f_k(x_i) \text{ for } w_i = \frac{1}{n} z_i(x_i)$$

Stable weight: $\underline{w} = (\frac{1}{T}, \dots, \frac{1}{T})$ is a target stable weight

If we can find w that is NOT too far from the stable weights and it satisfies

$$Bkw: |\sum_{i=1}^n z_i w_i f_k(x_i) - \frac{1}{n} \sum_{i=1}^n f_k(x_i)| \leq b_k \text{ for } k=1, \dots, K.$$

Minimize $\sum_{i=1}^n (w_i - w'_i)^2$ such that Bkw holds for $k=1, \dots, K$

Then estimate $E(y_{i(1)})$ by $\sum_{i=1}^n z_i w_i y_i$

Similarly, we can estimate $E(y_{i(0)})$

An alternative objective is to maximize $\sum_{i=1}^n -w_i \log w_i$

* 09/14/23 • Identification in causal inference

• Instrumental variables

n individuals drawn i.i.d. from a population. We want to estimate a population quantity using our data

Full data: $(y_{i(1)}, y_{i(0)}, X_i, z_i) \stackrel{iid}{\sim} P$

CATE
↑

$$\tau = E(y_{i(1)} - y_{i(0)} | z_i = 1) \text{ or } E(y_{i(1)} - y_{i(0)} | z_i = 0) \text{ or } E(y_{i(1)} - y_{i(0)} | X_i = n) \equiv E(y_{i(1)} - y_{i(0)})$$

Observed data: $(y_i, X_i, z_i) \stackrel{iid}{\sim} P_{obs}$ where $y_i = y_{i(1)}z_i + y_{i(0)}(1-z_i)$

Identification: A parameter is identifiable if it can be expressed as a function of the observed data distribution P_{obs} .

$\tau = g(P_{obs})$ In that case, g is the identifier function. Then we can give a point estimate of τ as $\hat{\tau} = g(\hat{P}_{obs})$

We can get the sampling distribution of $\hat{\tau}$ by bootstrapping.

In general, none of those are identifiable without additional assumptions.

Assumption: $(y_{i(1)}, y_{i(0)}) \perp\!\!\!\perp z_i | X_i$ (Ignorability, Exchangeability, No-?? confounders)

$$\tau = E(y_{i(1)} - y_{i(0)}) = E(y_{i(1)}) - E(y_{i(0)})$$

$\left\{ \begin{array}{l} y_{i(1)} \perp\!\!\!\perp z_i | X_i \\ \Pr(z_i = 1 | X_i) > 0 \text{ a.s.} \end{array} \right.$

$$E(y_{i(1)}) = E[E(y_{i(1)} | X_i)] \xrightarrow{\text{assumption}} E[E(y_{i(1)} | X_i, z_i = 1)] = E[E(y_i | X_i, z_i = 1)]$$

Lemma: $\tau = E(y_{i(1)} - y_{i(0)})$ is identifiable under the assumptions $E(y_{i(1)} | X_i, z_i = 1) = E(y_{i(1)} | X_i)$

$$| E(y_{i(0)} | X_i, z_i = 0) = E(y_{i(0)} | X_i)$$

$$\text{as } \tau = E[E(y_i | X_i, z_i=1) - E(y_i | X_i, z_i=0)]$$

* $0 < \Pr(z_i=1 | X_i) < 1$ a.s. overlap Since it is equivalent to say that $f(x_i | z_i=1)$ & $f(x_i | z_i=0)$ have the same support.

g -formula, back-door adjusted, outcome model best identification

$$\begin{aligned} E\left[\frac{z_i y_i}{\pi(x_i)}\right] &= E\left[E\left(\frac{z_i y_i}{\pi(x_i)} | X_i\right)\right] = E\left[\frac{1}{\pi(x_i)} E(z_i y_i | X_i)\right] \\ &= E\left[\frac{1}{\pi(x_i)} E(z_i | X_i) E(y_i | z_i=1, X_i)\right] \\ &= E\left[E(y_i | z_i=1, X_i)\right] \end{aligned}$$

Proposition: $\tau = E(y_i | z_i=1) - E(y_i | z_i=0)$ is identifiable under the assumption $(y_i | z_i, X_i) \perp\!\!\!\perp z_i | X_i$ and

$$0 < \Pr(z_i=1 | X_i) < 1 \text{ as } E\left[E(y_i | z_i=1, X_i) - E(y_i | z_i=0, X_i)\right]$$

$$\text{or } E\left(\frac{z_i y_i}{\pi(x_i)} - \frac{(1-z_i)y_i}{1-\pi(x_i)}\right)$$

$$* E(y_i | z_i=1) = E\left[E(y_i | z_i=1, X_i) + \frac{z_i(y_i - E(y_i | z_i=1, X_i))}{\pi(x_i)}\right] \rightarrow \text{the third formula}$$

Theorem: Consider any function $h(x)$ and $e(x)$,

$$E(y_i | z_i=1) = E\left[h(X_i) + \frac{z_i(y_i - h(X_i))}{e(X_i)}\right] \rightarrow \text{Doubly Robust Identification}$$

as long as one of the following is true,

$$(i) \quad h(X_i) = E(y_i | z_i=1, X_i); \quad (ii) \quad e(X_i) = \pi(X_i)$$

$$\begin{aligned} \text{Pf. Assume (ii), } E\left[h(X_i) + \frac{z_i(y_i - h(X_i))}{\pi(x_i)}\right] &= E\left[\frac{h(X_i)(\pi(x_i) - z_i)}{\pi(x_i)} + \frac{z_i y_i}{\pi(x_i)}\right] \\ &= E\left[\frac{h(X_i)(\pi(x_i) - z_i)}{\pi(x_i)}\right] + E\left[\frac{z_i y_i}{\pi(x_i)}\right] \\ &= 0 + E(y_i | z_i=1) \end{aligned}$$

Proposition: $\tau = E(y_i | z_i=1) - E(y_i | z_i=0) | \tilde{X}_i = \tilde{x}$ where \tilde{X}_i is a subvector of X_i . Under Ignorability and overlap, we can identify τ by

$$\int_{x:\tilde{x}} [E(y_i | z_i=1, X_i=x) - E(y_i | z_i=0, X_i=x)] dP_{\substack{\text{marginal} \\ (\text{integrating out } \tilde{x})}}(x)$$

$(Y_{i(1)}, Y_{i(0)}) \perp\!\!\!\perp Z_i \mid X_i$ might fail to hold:

Example: Effect of NICU on mortality after birth

X_i = mother choose to give birth in a hospital with nicu

U_i = the intent of the mother to go hospital with a nicu

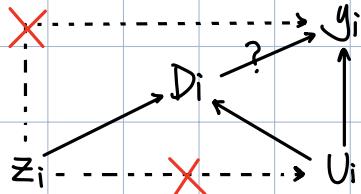
$(Y_{i(1)}, Y_{i(0)}) \perp\!\!\!\perp Z_i \mid X_i, U_i$ $\xrightarrow{\text{unmeasurable}}$

$(Y_{i(1)}, Y_{i(0)}) \not\perp\!\!\!\perp Z_i \mid X_i$

Example: Effect of strong/harsh sentencing in a criminal case on recitivization

Instrumental Variables Method

(D_i, X_i, Y_i) $Z_i \leftarrow$ an instrument



- $\left. \begin{array}{l} \text{(a) the instrument is associated with the treatment} \\ \text{(b) the instrument is NOT associated with the unmeasured covariates} \\ \text{(c) the instrument has no direct effect on the outcome} \end{array} \right\}$

Y_i = whether the baby was healthy

X_i = observed covariates

U_i = mother's intent to go to a hospital with nicu

D_i = whether the baby was treated in a nicu

Z_i = whether the mother lived more than 30 min away from the nearest hospital with nicu

Notation: $Z_i = 0$ or 1

$$D_i(Z_i) \begin{cases} D_i(1) & \text{if } Z_i = 1 \\ D_i(0) & \text{if } Z_i = 0 \end{cases}$$

$$Y_i(Z_i, d_i) \begin{cases} Y_{i(0,0)} \\ Y_{i(0,1)} \\ Y_{i(1,0)} \\ Y_{i(1,1)} \end{cases}$$

$$(Z_i, D_i(1), D_i(0), Y_{i(1,1)}, Y_{i(1,0)}, Y_{i(0,1)}, Y_{i(0,0)}, X_i) \stackrel{\text{iid}}{\sim} P$$

$$(Z_i, D_i, Y_i, X_i) \stackrel{\text{iid}}{\sim} P_{\text{obs}}$$

Assumptions: ① If $z_i = z'_i$, then $D_i(z_i) = D_i(z'_i)$

If $z_i = z'_i$, $d_i = d'_i$, then $y_i(z_i, d_i) = y_i(z'_i, d'_i)$

② $E(D_i(1)|X_i) \neq E(D_i(0)|X_i)$ a.s.

③ $z_i \perp\!\!\!\perp (y_{i(1,1)}, y_{i(1,0)}, y_{i(0,1)}, y_{i(0,0)}) | X_i$

④ $y_i(z_i, d_i) = y_i(z'_i, d_i)$

Examples:

• Randomized Encouragement

The effect of drinking while pregnant on infant birthweight

Instrument $\{1, \dots, n\}$ $z_i = s_i^1$

• Natural Experiment

The effect of going to war on future wages.

Vietnam war lottery

$z_i = \begin{cases} 1 & \text{if their name came up in lottery} \\ 0 & \text{o.w.} \end{cases}$

• Distance to speciality care

The effect of being treated in a NICU on infant health.

$z_i = \begin{cases} 1 & \text{if the mother lived less than 30 mins from a nicu hospital} \\ 0 & \text{o.w.} \end{cases}$

• Preference for the treatment

The effect of harsh sentencing on recidivism.

$z_i = \begin{cases} 1 & \text{if the judge in the case had a history of harsh sentencing} \\ 0 & \text{o.w.} \end{cases}$

• The effect of food stamp on CVD

$z_i = \begin{cases} 1 & \text{if the person is from Florida} \\ 0 & \text{if the person is from California} \end{cases}$

Mendallian Randomization

z_i = gene expression

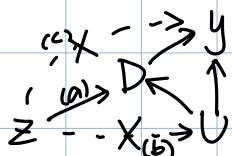
The effect of drinking on infant birthweight

$$z_i = \begin{cases} 1, & \text{if ALDH2 has null expression} \\ 0, & \text{o.w.} \end{cases}$$

Compliance Classification

| z_i | $D_i(z_i)$ | | | |
|-------|---------------------|-----------------|---------------|--------------------|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| | Always takers at | Compliers co | Defiers de | Never takers nt |

9/21/23 · Instrumental variables



- Identification theorem

- Violation of assumptions

- Additional discussion

Potential Outcomes

$$y(z, d) = y(z', d) \text{ for all } z, z', d.$$

We simplify our notation to $y(d)$

$$\text{ATE} = E(y(1) - y(0))$$

$$= \Pr(\text{nt}) E(y(1) - y(0) | \text{nt})$$

$$+ \Pr(\text{co}) E(y(1) - y(0) | \text{co})$$

$$+ \Pr(\text{de}) E(y(1) - y(0) | \text{de})$$

$$+ \Pr(\text{at}) E(y(1) - y(0) | \text{at})$$

"—" unable to measure b/c changing
 z doesn't change D for nt and at

Observed data:

y, z, D all are binary. There is no X .

We have data based on a multinomial distribution on $\{0, 1\}^3$. We have 8 parameters we can estimate.

Parameters of the full data

| | | | | | |
|------------------------------------|------------------------------------|------------|----------------------------------|--|---|
| $P(z=1)$ | | | | | } |
| $P_r(c_0)$ | $P_r(d_e)$ | $P_r(a_t)$ | $P_r(n_t)$ | | |
| $E(y^1 a_t)$ | $E(y^0 a_t)$ | | | | |
| $E(y^1 c_0)$ | $E(y^0 c_0)$ | | | | |
| $E(y^1 n_t)$ | $E(y^0 n_t)$ | | | | |
| $E(y^1 d_e)$ | $E(y^0 d_e)$ | | | | |

→ are ruled out by the monotonicity assumption

• Monotonicity, $D(z) \geq D(z')$ for any $z \geq z'$

For binary z , $D(1) \geq D(0)$. This implies no defiers.

Our estimand $E(y_{(1)} - y_{(0)}|c_0) \rightarrow$ Complier ATE (CATE) \rightarrow Local ATE (LATE)

Identification Theorem, Assume there is no X .

$$E(y_{(1)} - y_{(0)}|c_0) = \frac{E(y|z=1) - E(y|z=0)}{E(D|z=1) - E(D|z=0)}$$

$$\text{Pf: } E(y|z=1) - E(y|z=0)$$

$$\begin{aligned} &= E(y_{(1)}D_{(1)} + y_{(0)}(1-D_{(1)})|z=1) - E(y_{(1)}D_{(0)} + y_{(0)}(1-D_{(0)})|z=0) \\ &\stackrel{z \perp\!\!\!\perp y_{(1)}, y_{(0)}}{=} E(y_{(1)}D_{(1)} + y_{(0)}(1-D_{(1)})) - E(y_{(1)}D_{(0)} + y_{(0)}(1-D_{(0)})) \\ &= E[(y_{(1)} - y_{(0)})(D_{(1)} - D_{(0)})] \end{aligned}$$

$$\text{Monotonicity } E(y_{(1)} - y_{(0)} | D_{(1)} - D_{(0)} = 1) P_r(D_{(1)} - D_{(0)} = 1)$$

$$\text{Using association of } z \text{ and } D \quad E(y_{(1)} - y_{(0)} | c_0) \cdot (E(D|z=1) - E(D|z=0)) \quad \square$$

For any given person i , we only see $z_i, D_i(z_i)$.

$D_{(1)}$ and $D_{(0)}$ are not simultaneously observable. Hence we do not know their compliance group.

Understanding LATE

Option 1: If the treatment effect is homogeneous across compliance groups

Option 2: If we can explain who compliers are.

Lemma: For any function $h(y, X)$,

$$E(h(y_{(1)}, X) | C_0) = \frac{E(h(y, X)D | z=1) - E(h(y, X)D | z=0)}{E(D | z=1) - E(D | z=0)}$$

$$E(h(y_{(0)}, X) | C_0) = \frac{E(h(y, X)(1-D) | z=1) - E(h(y, X)(1-D) | z=0)}{E(1-D | z=1) - E(1-D | z=0)}$$

Pf, $E(h(y, X)D | z=1) - E(h(y, X)D | z=0)$

$$= E(h(y_{(1)}, X)D(1) | z=1) - E(h(y_{(0)}, X)D(0) | z=0)$$

$$= E[E(h(y_{(1)}, X)D(1) | z=1, X) - E(h(y_{(0)}, X)D(0) | z=0, X)]$$

$$= E(h(y_{(1)}, X)D(1)) - E(h(y_{(0)}, X)D(0))$$

$$= E(h(y_{(1)}, X)(D(1) - D(0)))$$

$$= E(h(y_{(1)}, X) | C_0) \Pr(D=1 | z=1) - \Pr(D=1 | z=0)) \quad \square$$

Using the lemma, we can estimate

$$E(y_{(1)} | C_0)$$

$$E(y_{(0)} | C_0)$$

$$E(y_{(1)} | \text{at}) = E(y | z=1, D=1)$$

$$E(y_{(0)} | \text{nt}) = E(y | z=0, D=0)$$

Consider a linear regression model,

$$y_d = \alpha_0 + \alpha_1 d + \alpha_2 X + \alpha_3 U + \varepsilon_d$$

Assume $\alpha_2 = 0$.

$$y_d = \alpha_0 + \alpha_1 d + \alpha_3 U + \varepsilon_d$$

What we see: $y = \alpha_0 + \alpha_1 D + \alpha_3 U + \varepsilon_0 + (\varepsilon_1 - \varepsilon_0)D$

Assume $E(\varepsilon_d | U) = 0$. $E(\varepsilon_d | D, U) = 0$

$$E(y | D=1) = \alpha_0 + \alpha_1 + \alpha_3 E(U | D=1) + E(\varepsilon_1 | D=1)$$

$$E(y | D=0) = \alpha_0 + \alpha_3 E(U | D=0) + E(\varepsilon_0 | D=0)$$

$$E(Y|D=1) - E(Y|D=0) = \alpha_1 + \underbrace{\alpha_3 [E(U|D=1) - E(U|D=0)]}_{\text{bias}} + \underbrace{E(\varepsilon_1|D=1) - E(\varepsilon_0|D=0)}_{=0}$$

Z is an instrument $E(U|Z=1) - E(U|Z=0) = 0$ (*)

$$E(\varepsilon_d|Z=1) = E(\varepsilon_d|Z=0) = 0$$

Then, $E(Y|Z=1) - E(Y|Z=0)$

$$= \alpha_1 [E(D|Z=1) - E(D|Z=0)] + \alpha_3 [E(U|Z=1) - E(U|Z=0)]$$

In this linear model, if assumption (b) does not hold (so that (*) does not necessarily hold) the bias of LATE estimator is $\alpha_3 \cdot \frac{E(U|Z=1) - E(U|Z=0)}{E(D|Z=1) - E(D|Z=0)}$

The bias of the previous estimator is $\alpha_3 [E(U|D=1) - E(U|D=0)]$

$$Y_{(d,z)} = \alpha_0 + \alpha_1 d + \alpha_2 z + \alpha_3 U + \varepsilon_d$$

(a) $E(D|Z=1) \neq E(D|Z=0)$

$$\overline{\sum_i 1_{(Z_i=1)}} \sum_i 1_{(D_i=1, Z_i=1)} - \overline{\sum_i 1_{(Z_i=0)}} \sum_i 1_{(D_i=1, Z_i=0)}$$

Theorem. (Y_i, D_i, Z_i) are iid from a population.

(a) $\frac{E(Y_i|Z_i=1) - E(Y_i|Z_i=0)}{E(D_i|Z_i=1) - E(D_i|Z_i=0)} = \frac{\text{Cov}(Y_i, Z_i)}{\text{Cov}(D_i, Z_i)}$

(b) $\sqrt{n} \left[\frac{\sum_i Z_i (Y_i - \bar{Y})}{\sum_i Z_i (D_i - \bar{D})} - \frac{\text{Cov}(Y_i, Z_i)}{\text{Cov}(D_i, Z_i)} \right] \xrightarrow{d} N \left(0, \frac{E[\eta_i^2 (Z_i - E Z_i)^2]}{[\text{Cov}(D_i, Z_i)]^2} \right)$

where $\eta_i = Y_i - E Y_i - \frac{\text{Cov}(Y_i, Z_i)}{\text{Cov}(D_i, Z_i)} (D_i - E D_i)$.

09/26/23 • Classical IV methods

$$y_i = D_i \beta + X_i' \alpha + \varepsilon_i, \quad E(\varepsilon_i) = 0$$

β ← estimand

$$\begin{cases} E[y_{(1)}|X_i] = \beta + X_i' \alpha \\ E[y_{(0)}|X_i] = X_i' \alpha \end{cases} \quad \text{assuming } E(\varepsilon_i|X_i) = 0$$

Run a least square of y on D_i, X_i

Get the coefficient of D_i from the regression

Is the coefficient consistent for β ?

Only consistent if $E(\varepsilon_i|D_i, X_i) = 0$ (Exogeneity)

If not, aka $E(\varepsilon_i|D_i, X_i) \neq 0$, it is called Endogeneity.

$$\begin{cases} (1) y_i = D_i \beta + X_i' \alpha + \varepsilon_i \\ (2) D_i = Z_i \lambda + X_i \tau + \eta_i \end{cases} \quad \text{assuming } \begin{cases} E(\varepsilon_i|Z_i, X_i) = 0 \\ E(\eta_i|Z_i, X_i) = 0 \end{cases}$$

First stage model

Goal is to estimate or make inference about β .

Procedure 1:

Get $\hat{\lambda}$ and $\hat{\tau}$ using LS of the first equation.

Write $y_i = Z_i \delta + X_i \theta + V_i$ where $\delta = \lambda \beta$, $\theta = \alpha + \tau \beta$, $V_i = \varepsilon_i + \eta_i \beta$

Now we have $E(V_i|Z_i, X_i) = 0$, so we can estimate $\hat{\delta}$ and $\hat{\theta}$ consistently by LS of the reduced form equation.

$$\hat{\beta} = \left(\hat{\lambda}' \left(\frac{1}{n} \sum Z_i Z_i' \right) \hat{\lambda} \right)^{-1} \hat{\lambda}' \left(\frac{1}{n} \sum Z_i Z_i' \right) \hat{\delta} \quad (\text{Weighted LS of } \hat{\delta} \text{ on } \hat{\lambda})$$

Procedure 2:

From $\hat{\lambda}$ and $\hat{\tau}$ calculate $\hat{D}_i = Z_i \hat{\lambda} + X_i \hat{\tau}$

$$y_i = \hat{D}_i \beta + X_i' \alpha$$

and estimating β using LS.

The estimators from the two procedures are equivalent.

This is called 2SLS (2 stage least square estimator)

$$\sqrt{n} \begin{pmatrix} \hat{\delta} - \delta \\ \hat{\lambda} - \lambda \end{pmatrix} \xrightarrow{d} N \left[0, \begin{pmatrix} Q_z^{*-1} & 0 \\ 0 & Q_z^{*-1} \end{pmatrix} \Lambda^* \begin{pmatrix} Q_z^{*-1} & 0 \\ 0 & Q_z^{*-1} \end{pmatrix} \right]$$

where $Q_z^* = E(z_i z_i')$, $\Lambda^* = \lim_{n \rightarrow \infty} \text{Var} \left(\frac{1}{\sqrt{n}} \sum_i \varepsilon_i z_i, \frac{1}{\sqrt{n}} \sum_i \eta_i z_i' \right)$

$$\delta = \lambda \beta$$

① We can test for $H_0: \lambda = 0$

Reject H_0 , then use 2SLS

Check the F statistic value is very large. Then proceed to 2SLS. (>40)

② Consider testing $H_0: \beta = \beta_0$

Equivalently, $H_0: \delta - \lambda \beta_0 = 0$

$$g(\beta_0) := \delta - \lambda \beta_0 \quad \hat{g}(\beta_0) = \hat{\delta} - \hat{\lambda} \beta_0$$

$$\sqrt{n} (\hat{g}(\beta_0) - g(\beta_0)) \xrightarrow{d} N(0, \Sigma g(\beta_0))$$

$$\text{Under } H_0, n \hat{g}(\beta_0) \Sigma g(\beta_0)^{-1} \hat{g}(\beta_0) \xrightarrow{d} \chi_k^2$$

Invert this test to get confidence set for β .

- * 09/28/23 • Causal inference in Panel data
 • Difference-in-differences

$$t = 1, \dots, T \quad i = 1, \dots, N$$

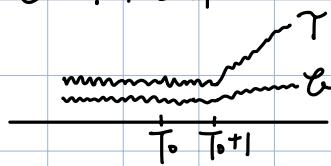
$$y_{it} \quad D_{it} = \begin{cases} 0 & \text{if unit } i \text{ is not treated at time } t \\ 1 & \text{if unit } i \text{ is treated at time } t \end{cases}$$

Treatment is staggered - $D_{it} = 1$ if $D_{i(t-1)} = 1$

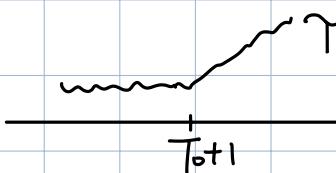
It may happen that the units that are treated are all treated at the same time

$T = \{i : D_{iT} = 1\}$ $T_0 =$ the last time period when units in T are not treated

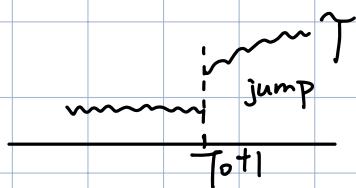
$$\mathcal{C} = \{i : D_{iT} = 0\}$$



Difference-in-difference



Before-after analysis



Structure break or
interrupted time series model

| Dit | | 1 2 T | | | | | | | | | | | | | | | | | |
|-----|---|---------------------|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
| 1 | i | 0 0 ... 0 1 1 ... 1 | | | | | | | | | | | | | | | | | |
| 2 | | 0 0 ... 0 1 1 ... 1 | | | | | | | | | | | | | | | | | |
| | : | | | | | | | | | | | | | | | | | | |
| N | | 0 0 0 | | | | | | | | | | | | | | | | | |

units might be treated at different time

Two units who are observed for two time periods.

One of the two units is treated in the second time period.

The other unit remains untreated.

$y_{it}(1)$ for $t = 1, 2$ - potential outcomes for unit i at time 1 and 2 if the unit was treated

$y_{it}(0)$ for $t = 1, 2$ - potential outcomes if unit i remained untreated

$$y_{it} = \begin{cases} y_{it}(1) & \text{if } D_{iz} = 1 \\ y_{it}(0) & \text{if } D_{iz} = 0 \end{cases}$$

Treatment anticipation: y_{ii} depends on whether the unit gets treatment at time 2

Assumption: No anticipation of the treatment, i.e., $y_{ii}(0) = y_{ii}(1)$.

Under the assumption $y_{ii} = y_{ii}(1) = y_{ii}(0)$

$$y_{iz} = D_{iz} y_{iz}(1) + (1 - D_{iz}) y_{iz}(0)$$

| | $y_{ii}(\cdot)$ | $y_{iz}(\cdot)$ |
|---|-----------------|-----------------|
| 1 | ✓ | ✓ |
| 0 | ✓ | ✗ |

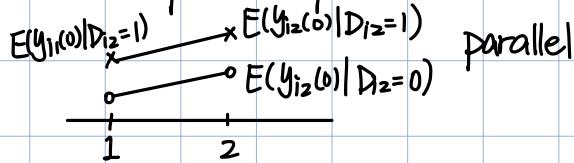
if $D_{iz} = 1$

$$\text{ATT} = E(y_{iz}(1) - y_{iz}(0) | D_{iz} = 1)$$

Assumption: $D_{iz} \perp \underbrace{(y_{iz}(0) - y_{ii}(0))}_{\text{the trend of unit } i \text{ under no treatment}}$

the trend of unit i under no treatment

This assumption implies $E(y_{iz}(0) - y_{ii}(0) | D_{iz} = 1) = E(y_{iz}(0) - y_{ii}(0) | D_{iz} = 0)$



parallel trends in mean between the groups

$$\text{ATT} = E(y_{iz}(1) - y_{iz}(0) | D_{iz} = 1)$$

$$= E(y_{iz}(1) | D_{iz} = 1) - E(y_{iz}(0) | D_{iz} = 1)$$

$$= E(y_{iz} | D_{iz} = 1) - [E(y_{ii}(0) | D_{iz} = 1) + E(y_{iz}(0) - y_{ii}(0) | D_{iz} = 0)]$$

$$= E(y_{iz} - y_{ii} | D_{iz} = 1) - E(y_{iz} - y_{ii} | D_{iz} = 0)$$

→ Difference-in-difference identification

For iid data (y_{ii}, y_{iz}, D_{iz}) , the estimator is

$$\frac{1}{\sum(D_{iz}=1)} \sum_i 1(D_{iz}=1) (y_{iz} - y_{ii}) - \frac{1}{\sum(D_{iz}=0)} \sum_i 1(D_{iz}=0) (y_{iz} - y_{ii})$$

Challenges with DiD:

- ① Only restricted to two time periods
- ② It assumes no anticipation
- ③ Parallel trend only holds given covariates
- ④ The scale dependence of the parallel trend assumption

$\exists_i \perp (y_{i2}(1), y_{i2}(0)) \mid X_i \rightarrow$ This is invariant to the scale

$D_{i2} \perp y_{i2}(0) - y_{i1}(0) \mid X_i \rightarrow$ But this is NOT (for example, $\log y_{i2}(0) - \log y_{i1}(0)$)

③ says it might only hold that $E(y_{i2}(0) - y_{i1}(0) \mid D_{i2}=1, X_i) = E(y_{i2}(0) - y_{i1}(0) \mid D_{i2}=0, X_i)$ a.s.
 → Parallel trend given covariates

Theorem: Under no-anticipation assumption and parallel trends given covariates,

$$\text{ATT} = E \left[\left(\frac{D_{i2}}{E(D_{i2})} - \frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} \right) (y_{i2} - y_{i1}) \right]$$

assuming $\lambda(X_i) = \Pr(D_{i2}=1 \mid X_i) \in (0, 1)$ almost surely.

$$\begin{aligned} \text{Pf. } E\left(\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)}\right) &= E\left[E\left(\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} \mid X_i\right)\right] \\ &= E\left[\frac{\lambda(X_i)}{1-\lambda(X_i)} E(1-D_{i2} \mid X_i)\right] \\ &= E(\lambda(X_i)) = E(D_{i2}) \end{aligned}$$

Hence, RHS becomes

$$\underbrace{\frac{1}{E(D_{i2})} E\left[D_{i2} - \frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} (y_{i2} - y_{i1})\right]}$$

$$\cdot E[D_{i2}(y_{i2} - y_{i1})] = E[y_{i2} - y_{i1} \mid D_{i2}=1] \cdot E(D_{i2})$$

$$\cdot E\left[\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} (y_{i2} - y_{i1})\right] = E\left[\frac{\lambda(X_i)}{1-\lambda(X_i)} E((1-D_{i2})(y_{i2} - y_{i1}) \mid X_i)\right]$$

$$= E\left[\frac{\lambda(X_i)}{1-\lambda(X_i)} \cdot E(y_{i2} - y_{i1} \mid X_i, D_{i2}=0) \cdot \Pr(D_{i2}=0 \mid X_i)\right]$$

$$\text{ATT} = E(y_{i2}(1) - y_{i2}(0) \mid D_{i2}=1)$$

$$= E_x\left[E(y_{i2}(1) - y_{i2}(0) \mid X_i, D_{i2}=1) \mid D_{i2}=1\right]$$

$$= E_x\left[E(y_{i2}(1) - y_{i2}(0) + y_{i1}(0) - y_{i1}(0) \mid X_i, D_{i2}=1) \mid D_{i2}=1\right]$$

$$= E_x\left[E(y_{i2}(1) - y_{i1}(0) \mid X_i, D_{i2}=1) - E(y_{i2}(0) - y_{i1}(0) \mid X_i, D_{i2}=1) \mid D_{i2}=1\right]$$

$$= E_x\left[E(y_{i2}(1) - y_{i1}(0) \mid X_i, D_{i2}=1) - E(y_{i2}(0) - y_{i1}(0) \mid X_i, D_{i2}=0) \mid D_{i2}=1\right]$$

↓ Parallel trend given covariates

$$= E_x \left[E(y_{i2} - y_{i1} | X_i, D_{i2}=1) - E(y_{i2} - y_{i1} | X_i, D_{i2}=0) \mid D_{i2}=1 \right]$$

exercise: complete the last two steps of the calculation.

In light of "~~" in the proof, the estimator \widehat{ATT} is

$$\text{1st item} = \frac{1}{\frac{\sum(D_{i2}=1)}{n}} \cdot \frac{1}{n} \sum_i D_{i2} (y_{i2} - y_{i1}) = \frac{1}{\sum(D_{i2}=1)} \sum_{i: D_{i2}=1} (y_{i1} - y_{i2})$$

(minus)

$$\text{2nd item} = \frac{1}{\frac{\sum(D_{i2}=1)}{n}} \cdot \frac{1}{n} \sum_i \frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} (y_{i2} - y_{i1}) = \frac{1}{\sum(D_{i2}=1)} \sum_i \frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)} (y_{i2} - y_{i1})$$

10/03/23

Issues of Difference-in-difference

- ① multiple time periods
- ② possible anticipation
- ③ parallel trends conditioned on the X_i
- ④ Scale restriction on parallel trends

Assumptions:

$$(1) \quad y_{i1}(1) = y_{i1}(0)$$

$$(2) \quad y_{i2}(0) - y_{i1}(0) + D_{i2} \mid X_i \quad \checkmark \quad (3)$$

$$E(y_{i2}(0) - y_{i1}(0) \mid D_{i2}=1, X_i) = E(y_{i2}(0) - y_{i1}(0) \mid D_{i2}=0, X_i) \text{ a.s.}$$

$$(3) \quad \Pr(D_{i2}=1 \mid X_i) \in (0, 1)$$

$\tau = E(y_{i2}(1) - y_{i2}(0) \mid D_{i2}=1)$ is identified as

$$E \left[\frac{D_{i2}}{E(D_{i2})} (y_{i2} - y_{i1}) - \frac{\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)}}{E(\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)})} (y_{i2} - y_{i1}) \right]$$

We can also identify τ by

$$E \left[\frac{D_{i2}}{E(D_{i2})} (y_{i2} - y_{i1} - E(y_{i2} - y_{i1} \mid X_i, D_{i2}=0)) \right]$$

exercise

$$\text{or } E \left[\frac{D_{i2}}{E(D_{i2})} (y_{i2} - y_{i1} - E(y_{i2} - y_{i1} \mid X_i, D_{i2}=0)) - \frac{\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)}}{E(\frac{\lambda(X_i)(1-D_{i2})}{1-\lambda(X_i)})} (y_{i2} - y_{i1} - E(y_{i2} - y_{i1} \mid X_i, D_{i2}=0)) \right]$$

For Issue ①:

$$Y_{i1}(1), Y_{i2}(1), \dots, Y_{iT}(1)$$

$$Y_{i1}(0), Y_{i2}(0), \dots, Y_{iT}(0)$$

$$D_{i1}, D_{i2}, \dots, D_{iT}$$

Assume T_0 is the last time the unit remains untreated

$$T_0 < T \quad \text{or} \quad T_0 = T$$

$$\text{ATT}_g = E(Y_{ig}(1) - Y_{ig}(0) \mid D_{iT}=1), g=1, \dots, T$$

$$\text{Assumption: } E[Y_{ig}(0) - Y_{i1}(0) \mid D_{iT}=1] = E[Y_{ig}(0) - Y_{i1}(0) \mid D_{iT}=0]$$

$$\text{Assumption: } Y_{i1}(1) = Y_{i1}(0), Y_{iT_0}(1) = Y_{iT_0}(0)$$

$$\text{Proposition: } \text{ATT}_g = E(Y_{ig} - Y_{i1} \mid D_{iT}=1) - E(Y_{ig} - Y_{i1} \mid D_{iT}=0)$$

$$E(Y_{it}(0) - Y_{it(t-1)}(0) \mid D_{iT}=1) = E(Y_{it}(0) - Y_{it(t-1)}(0) \mid D_{iT}=0) \text{ for all } t \geq 1$$

$$\Rightarrow E(Y_{ig}(0) - Y_{ig(g-1)}(0) \mid D_{iT}=1) = E(Y_{ig}(0) - Y_{ig(g-1)}(0) \mid D_{iT}=0)$$

$$\text{In fact, } \text{ATT}_g = E(Y_{ig} - Y_{ig'} \mid D_{iT}=1) - E(Y_{ig} - Y_{ig'} \mid D_{iT}=0) \text{ for } g' = g-1, \dots, T_0$$

For Issue ②:

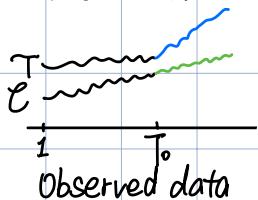
$$\text{Assumption: } Y_{i1}(0) = Y_{i1}(1), \dots, Y_{ik}(0) = Y_{ik}(1) \text{ for some } k < T_0 \quad (\text{Limited anticipation})$$

Under limited anticipation,

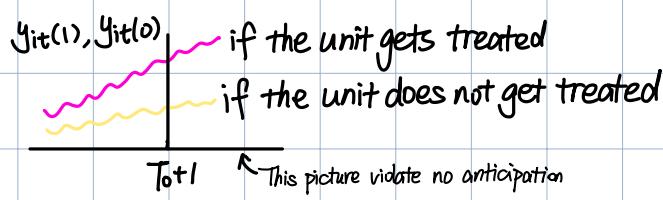
$$\text{ATT}_g = E(Y_{ig} - Y_{ig'} \mid D_{iT}=1) - E(Y_{ig} - Y_{ig'} \mid D_{iT}=0) \text{ for } g' = g-1, \dots, k$$

* 10/05/23

$T \leftarrow$ treated unit



$C \leftarrow$ control unit

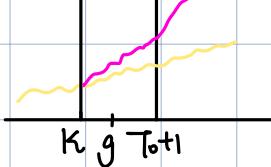


This picture violate no anticipation

No anticipation:



Partial anticipation:



We have two time periods = 1, 2, and y_{i1}, y_{i2}, D_{i2}

Potential outcomes: $y_{i1}(0), y_{i2}(0), y_{i2}(1) [y_{i1}(1) = y_{i1}(0)]$

Standard parallel trend:

$$E(y_{i2}(0) - y_{i1}(0) | D_{i2}=1) = E(y_{i2}(0) - y_{i1}(0) | D_{i2}=0)$$

We say that a parallel trend that is invariant under scale transformation would satisfy $E(g(y_{i2}(0)) - g(y_{i1}(0)) | D_{i2}=1) = E(g(y_{i2}(0)) - g(y_{i1}(0)) | D_{i2}=0)$ for all strictly monotone function g . (1)

Proposition, (1) is true if and only if

$$(2) F_{y_{i2}(0)|D_{i2}=1}(y) - F_{y_{i1}(0)|D_{i2}=1}(y) = F_{y_{i2}(0)|D_{i2}=0}(y) - F_{y_{i1}(0)|D_{i2}=0}(y) \text{ for all } y.$$

Pf. (2) \Rightarrow (1): Take the derivative by d operation then multiply by $g'(y)$ on both side.

Then integrate.

$$(1) \Rightarrow (2): \text{ Fix } \bar{y}. \quad g_1(y) = y, \quad g_2(y) = y - 1_{(y \leq \bar{y})}$$

Apply g_1 and g_2 to (1), and then subtract to get (2).

Proposition: (2) is true if and only if

$$(3) F_{y_{it}(0)|D_{i2}=d}(y) = \theta C_t(y) + (1-\theta) H_d(y)$$

where $t=1, 2$, $d=0, 1$, θ is a number between 0 and 1, C_t and H_d are CDF.

Pf. (3) \Rightarrow (2) by simple subtraction of the distribution function.

Lemma: Given two distribution functions F_1 and F_2 (with densities f_1 and f_2).

$$F_1 = \theta G + (1-\theta)\tilde{F}_1 \quad \text{and} \quad F_2 = \theta G + (1-\theta)\tilde{F}_2$$

where θ and \tilde{F}_1, \tilde{F}_2 only depends on $f_1 - f_2$. $G, \tilde{F}_1, \tilde{F}_2$ are CDFs, $\theta \in [0, 1]$.

$$(2) \Rightarrow (3): \begin{cases} F_{Y_{i2}(0)|D_{i2}=1} = \theta_1 G_1 + (1-\theta_1) \tilde{F}_{Y_{i2}(0)|D_{i2}=1} \\ F_{Y_{i1}(0)|D_{i2}=1} = \theta_1 G_1 + (1-\theta_1) \tilde{F}_{Y_{i1}(0)|D_{i2}=1} \\ F_{Y_{i2}(0)|D_{i2}=0} = \theta_0 G_0 + (1-\theta_0) \tilde{F}_{Y_{i2}(0)|D_{i2}=0} \\ F_{Y_{i1}(0)|D_{i2}=0} = \theta_0 G_0 + (1-\theta_0) \tilde{F}_{Y_{i1}(0)|D_{i2}=0} \end{cases}$$

By (2) and the lemma, $\theta_0 = \theta_1$, $\tilde{F}_{Y_{i2}(0)|D_{i2}=1} = \tilde{F}_{Y_{i2}(0)|D_{i2}=0}$ & $\tilde{F}_{Y_{i1}(0)|D_{i2}=1} = \tilde{F}_{Y_{i1}(0)|D_{i2}=0}$

So, $F_{Y_{it}|D_{i2}=d} = \theta G_d + (1-\theta) \tilde{F}_t$ where $\tilde{F}_t = \tilde{F}_{Y_{it}(0)|D_{i2}=d}$

- $i \rightarrow Y_{ii}, D_{i2}$ Panel Data

- $i \rightarrow (Y_i, D_i, T_i)$ Repeated Cross-sectional Data

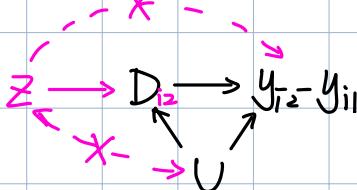
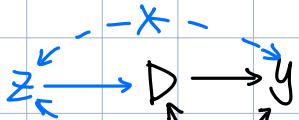
$$Y_i = h(D_i, T_i, X_i)$$

$$Y_i = \alpha + \beta D_i + \gamma T_i + \varepsilon_i$$

$$Y_{i2}(0) - Y_{i1}(0) \perp D_{i2} \leftarrow (Y_{i2}(0), Y_{i1}(0), Y_{i2}(1), Y_{i1}(1)) \perp D_{i2}$$



Instrumented parallel trends $\leftarrow \bar{z}_i$



$$\text{LATE} = E(Y_{i2}(1) - Y_{i1}(0) | \text{Complier})$$

$$= \frac{E(Y_{i2} - Y_{i1} | z_i=1) - E(Y_{i2} - Y_{i1} | z_i=0)}{E(D_{i2} | z_i=1) - E(D_{i2} | z_i=0)}$$

David Campbell:

$$Y_i = \alpha + \tau C_i + \beta I_i + \gamma T_i + \text{error}$$

$$C_i = \begin{cases} 1, & \text{if it is in the control group} \\ 0, & \text{if it is in the treatment group} \end{cases}$$

Pre treatment Post treatment

$$I_i = \begin{cases} 1, & \text{if } C_i=1 \text{ and time = 2} \\ 0, & \text{if } C_i=0 \text{ and time is 1 or 2.} \end{cases}$$

$$I_i = C_i \times T_i$$

$$T_i = \begin{cases} 0, & \text{at time 1} \\ 1, & \text{at time 2} \end{cases}$$

$$i \rightarrow (Y_i, C_i, T_i)$$

Note that this model implies

$$E(Y_i | T_i=1, C_i=1) - E(Y_i | T_i=0, C_i=1) = (\alpha + \tau + \beta + \gamma) - (\alpha + \tau) = \beta + \gamma$$

$$E(Y_i | T_i=1, C_i=0) - E(Y_i | T_i=0, C_i=0) = (\alpha + \gamma) - \alpha = \gamma$$

$(Y_{i(0)}, Y_{i(1)}, C_i, T_i) \leftarrow \text{Full data}$

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

Interested in estimating $\text{ATT} = E(Y_{i(1)} - Y_{i(0)} | C_i=1, T_i=1)$

U_i : unmeasured confounders

$$Y_{i(0)} = h(U_i, T_i), U_i \text{ is also related to } C_i$$

For particular h and U , this model reduces to our previous linear model

[e.g. $U_i = \alpha + \tau C_i$, h a linear function]