

Causality I (C.C. 8.9)

Unit level Causal Model:

y : outcome

(x, u) : covariates

g : structural functions.



Suppose $y \in \mathbb{R}$, $x \in \mathcal{X} \subseteq \mathbb{R}$, $u \in \mathcal{U}$

$$y \leftarrow g(x, u).$$

implies. $y = g(x, u)$.

Causal Effects:

1. $\underline{g(x_1, u_0) - g(x_0, u_0)}$.

2. Marginal causal effect:

$$[\nabla_x g](x_0, u_0)$$

where $\nabla_x g = \frac{\partial g}{\partial x}$

Causal Effects for Specific Unit i :

Y_i : outcome variable

(x_i, u_i) : covariates. \rightarrow numbers assigned to i

$$Y_i \leftarrow g(x_i, u_i) \text{ for all } i \in \mathcal{X}.$$

Fix u_i : $\underline{g(x_i, u_i) - g(x_0, u_i)}$.

Potential Outcomes:

$Y_i(\cdot) : \mathcal{X} \rightarrow \mathbb{R}$.

$$Y_i(x) \equiv g(x, u_i) \quad y \leftarrow Y_i(x).$$

Unit level causal effect: $Y_i(x_1) - Y_i(x_0)$.

learning Y_i or some aspect of it is the main goal of causal analysis.

Counterfactuals:

Revised/Factual outcomes: $Y_i = Y_i(X_i) = g(x_i, u_i)$.

Counterfactual outcomes: $Y_i(x)$ for $x \neq X_i$

Status quo treatment effect: $Y_i(x) - Y_i(X_i)$

Population level Causal Effects: Distribution of Effects.

X : random variable

X_i : value drawn from X

x : fixed, $\in \mathcal{X}$.

$$\text{for } i: \Delta_i(x_0 \rightarrow x_i) = g(x_i, u_i) - g(x_0, u_i) = Y_i(x_i) - Y_i(x_0)$$

$$\text{for population: } \Delta(x_0 \rightarrow x) = g(x, u) - g(x_0, u) = Y(x) - Y(x_0)$$

where u is random variable.

Distribution of Treatment Effect:

$$\begin{aligned} DTE(t, x_0 \rightarrow x_1) &= F_{Y(x_1)} - F_{Y(x_0)}(t) \\ &= P(Y(x_1) - Y(x_0) \leq t). \end{aligned}$$

* The proportion who benefit (voting criterion):

$$\text{parameter: } P(Y_1 > Y_0)$$

$$= P(Y_1 - Y_0 > 0)$$

$$= 1 - DTE(0, 0 \rightarrow 1).$$

* Distribution of Marginal Causal Effect

$$\Delta(x_0) = [\nabla_x g](x_0, u)$$

$$\Delta(x_0) = \frac{\partial Y(x_0)}{\partial x}$$

Population Level Causal Effects: Summary Statistics

Average Treatment Effect:

$$ATE \equiv E[\Delta(x_0 \rightarrow x_1)]$$

$$= E[g(x_1, u) - g(x_0, u)]$$

$$= E[g(x_1, u)] - E[g(x_0, u)]$$

Average Structural Function:

$$ASF \equiv E_u[g(x, u)]$$

$$= E[Y(x)]$$

$$\therefore ATE = E[Y(x_1) - Y(x_0)]$$

$$= E[Y(x_1)] - E[Y(x_0)].$$

Quantile Effects:

T^{th} quantile of the distribution of u . l. c. e.

$$Q_{\Delta(x_0 \rightarrow x_1)}(T).$$

T^{th} quantile treatment effect:

$$QTE(T, x_0 \rightarrow x_1) = Q_{g(x_1, u)}(T) - Q_{g(x_0, u)}(T)$$

Usually they are not the same.

Examples of Unit Level Causal Models.

1. Binary:

$(u_i, X_i, Y_i, Y_{i(0)}, Y_{i(1)})$. X_i : binary.

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

2. The Additively Separable Model:

$$(x \in \mathcal{X} \subseteq \mathbb{R})$$

$$Y(x) = m(x) + u.$$

$m: \mathcal{X} \rightarrow \mathbb{R}$ unknown function. u : unknown scalar r.v.

$y = m(x) + u$. $m(\cdot)$ not dependent on i .

$$Y_i(x) = m(x) + U_i$$

$Y_i(x_1) - Y_i(x_0) = m(x_1) - m(x_0)$. homogeneous treatment effects.
(nonlinear in x).

Z (cont). The Linear Constant Coefficient Model:

Assume $m(x) = \theta_0 + \theta_1 x$.

$$Y_i(x) = m(x) + U_i = \theta_0 + \theta_1 x + U_i$$

$Y_i(x) = \theta_0 + \theta_1 x + U_i$ not a "residual"; but an unobserved covariate.

$$Y_i(x_1) - Y_i(x_0) = \theta_1(x_1 - x_0)$$

Structural Interpretation of unobserved random variable U_i :
different X_i across i : observed heterogeneity
different U_i across i : unobserved heterogeneity.

The Difference b/w Fixing and Conditioning:

Fixing: using marginal distribution of U_i , define:

$$\text{potential outcome: } Y(x) = g(x, U)$$

then using \tilde{F} to define causality $Y(x_1) - Y(x_0)$.

it depends only on: (1) $g(x, \cdot)$

(2) marginal dist. F_U

Here we do not say anything about the joint dist of X & U
in the data to define potential outcomes!

Conditioning: implicitly depends on true joint dist of X & U .

$$\begin{aligned} \text{e.g. } \mathbb{E}[Y|X=x] &= \mathbb{E}[g(X, U)|X=x] \\ &= \mathbb{E}[g(x, U)|X=x] \end{aligned}$$

It is not the same thing as.

$$\text{ASF}(x) \equiv \mathbb{E}[g(x, U)] !$$

Causal Effects for Subpopulations:

treated units: $\{i \in \mathcal{X} : X_i = 1\}$

untreated units: $\{i \in \mathcal{X} : X_i = 0\}$.

$$\text{ATT} = \mathbb{E}[Y(1) - Y(0) | X=1]$$

$$\text{ATU} = \mathbb{E}[Y(1) - Y(0) | X=0]$$

by LIE: $\text{ATE} = \text{ATT} \cdot \mathbb{P}(X=1) + \text{ATU} \cdot \mathbb{P}(X=0)$.

Observe another covariate W :

consider units: $\{i \in \mathcal{X} : W_i = w\}$. a subpopulation.

$$\text{CATE} = \mathbb{E}[Y(1) - Y(0) | W=w].$$

Common Identification Proof Techniques.

• Constructive method: parameter is a known function of the data.

$F_{Y,x,u}$: true population distribution.

$F_{Y,x}$: observed population data.

\mathcal{F} : restrictions put on $F_{Y,x,u}$.

$\tilde{F}_{Y,x,u}$: a conjectured value of the population dist.

$\tilde{F}_{Y,x} = \text{MakeData}(\tilde{F}_{Y,x,u})$: corresponding population data that would be observed if $\tilde{F}_{Y,x,u}$ were true.

Theorem: suppose there is a known function g :

such that $\theta(\tilde{F}_{Y,x,u}) = g(\tilde{F}_{Y,x})$,

for any $\tilde{F}_{Y,x,u} \in \mathcal{F}$ and suppose the model is not refuted. Then, θ is point identified.

Pf: WTS: the identified set:

$$\mathcal{D}_I = \{\theta \in \mathbb{H} : \theta = \theta(\tilde{F}_{Y,x,u}) \text{ for some } \tilde{F}_{Y,x,u} \in \mathcal{F} \text{ such that } \text{MakeData}(\tilde{F}_{Y,x,u}) = F_{Y,x}\}.$$

is a singleton.

\because model is not refuted

\therefore non empty

Let $a \in \mathcal{D}_I$. $b \in \mathcal{D}_I$.

by definition of \mathcal{D}_I , $\exists \tilde{F}_{Y,x,u}^a \in \mathcal{F}$ such that $a = \theta(\tilde{F}_{Y,x,u}^a)$
since $\tilde{F}_{Y,x,u}^a \in \mathcal{F}$, $\theta(\tilde{F}_{Y,x,u}^a) = g(\tilde{F}_{Y,x}^a)$ by assumption.
therefore $a = g(\tilde{F}_{Y,x}^a)$.

by definition of \mathcal{D}_I , $\exists \tilde{F}_{Y,x,u}^b \in \mathcal{F}$ such that $b = \theta(\tilde{F}_{Y,x,u}^b)$
since $\tilde{F}_{Y,x,u}^b \in \mathcal{F}$, $\theta(\tilde{F}_{Y,x,u}^b) = g(\tilde{F}_{Y,x}^b)$ by assumption
therefore $b = g(\tilde{F}_{Y,x}^b)$.

Finally, notice that

$$\tilde{F}_{Y,x}^a = \text{MakeData}(\tilde{F}_{Y,x,u}^a) = F_{Y,x}, \tilde{F}_{Y,x}^b = \text{MakeData}(\tilde{F}_{Y,x,u}^b) = F_{Y,x}$$

$$\therefore a = g(\tilde{F}_{Y,x}^a) = b$$

★.

A parameter is point identified if it can be expressed as a known function of the distribution of the observable random variables.

Causality 2. Experiments CC 12.13

Identifying Treatment Effects From the Data Alone.

The fundamental Problem of Causal Inference:

X : the set of logically possible values of treatment.

$X_i \in X$: realized treatment for i .

Y_i : realized outcome. $Y_i = Y_i(X_i)$.

Where $Y_i(x)$ for all $x \in X$ are potential outcomes.

$Y_i(x)$ for all $x \neq X_i$: counterfactuals.

$$\text{Difference: } Y_i(X_1) - Y_i(X_0).$$

Problem: we just observe one value of x in the data.

Theorem: let \mathcal{Y} denote the set of all logically possible values of potential outcomes. For any $i \in \mathcal{X}$, the identified set for $Y_i(x)$ for any $x \in X$ with $x \neq X_i$ is \mathcal{Y} .

Without further assumptions data tells us nothing about counterfactual outcomes.

Learning T : main identification problem.

Relationship btw realized & potential outcomes.

$$Y_i = Y_i(X_i)$$

$$\Leftrightarrow Y_i = \sum_{x \in X} \mathbb{I}[X_i=x] Y_i(x)$$

$$\text{In binary Case: } Y_i = Y_{i(1)} X_i + Y_{i(0)} (1-X_i)$$

Marginal Distribution of Potential Outcomes.

Consider some treatment $x \in X$.

$$\begin{aligned} P(Y(x) \leq y) &= P(Y(x) \leq y | X=x) P(X=x) + P(Y(x) \leq y | X \neq x) P(X \neq x) \\ &= \underbrace{P(Y \leq y | X=x)}_{\text{known}} \underbrace{P(X=x)}_{\text{unknown}} + \underbrace{P(Y \leq y | X \neq x)}_{\text{A dist. of counterfactual outcomes.}} \underbrace{P(X \neq x)}_{\text{unknown}} \end{aligned}$$

Theorem. $\text{supp}(Y(x)) \subseteq S \quad \forall x \in X$.

For fixed $y \in \mathbb{R}$, the no assumptions identified set for $P(Y(x) \leq y)$:

$$[P(Y \leq y | X=x) P(X=x), P(Y \leq y | X=x) P(X=x) + P(X \neq x)]$$

Conditional on $X=x$, $P(Y(x) \leq y | X=x) - P(Y(x) \leq y)$

usually $P(Y \leq y)$ and $P(Y \neq y) = P(Y \leq y | X=x) - P(Y(x) \leq y)$

are not the same. Selection bias

The Joint Identified Set:

We want to know Treatment Effect $Y(X_1) - Y(X_0)$

must know $P(Y(x) \leq y_1, Y(x) \leq y_0)$

Theorem (*) Let $x_1, x_0 \in X$. The identified set for $F_{Y(x_1), Y(x_0)}(\cdot, \cdot)$ is the set of all joint cdfs whose marginal cdfs are identified.

→ We cannot identify $P(Y(x) \leq y)$ for all $x \in X$ because we can't observe the counterfactual.

That would require $P(X=x) = 1 \quad \forall x \in X$ which is impossible since we have $\sum_{x \in X} P(X=x) = 1$.

The Average Treatment Effect:

Suppose $\text{supp}(Y(x)) \subseteq \mathcal{Y} \quad \forall x \in X$.

$$\text{Normalization: } Y^{\text{norm}}(x) = \frac{Y(x) - Y_{\min}}{Y_{\max} - Y_{\min}} \quad Y_{\max} = \sup \mathcal{Y} \quad Y_{\min} = \inf \mathcal{Y}.$$

The identified set of mean potential outcomes:

$$\mathbb{E}[Y(x)] \in [LB(x), UB(x)]$$

$$= [\mathbb{E}[Y | X=x] P(X=x), \mathbb{E}[Y | X=x] P(X=x) + P(X \neq x)]$$

The identified set for ATE:

$$\mathbb{E}[Y(x_1)] \in [LB(x_1), UB(x_1)]$$

$$\mathbb{E}[Y(x_0)] \in [LB(x_0), UB(x_0)].$$

By Theorem (*):

$$\text{Joint identified set} = [UB(x_1), UB(x_0)] \times [LB(x_0), UB(x_0)].$$

$$\therefore \underline{\mathbb{E}[Y(x_1)] - \mathbb{E}[Y(x_0)]} \in [LB(x_1) - UB(x_0), UB(x_1) - LB(x_0)].$$

ATE.

$$WIDTH = 2 - [\mathbb{P}(X=x_1) + \mathbb{P}(X=x_0)]$$

$$\text{proof: ATE} = \mathbb{E}[Y(x_1)] - \mathbb{E}[Y(x_0)]$$

$$= \mathbb{E}[Y(x_1) | X=x_1] P(X=x_1) + \mathbb{E}[Y(x_0) | X \neq x_1] P(X \neq x_1)$$

$$\text{By LIE. } -[\mathbb{E}[Y(x_0) | X=x_0] P(X=x_0) + \mathbb{E}[Y(x_1) | X \neq x_0] P(X \neq x_0)]$$

$$= \mathbb{E}[Y(x_0) | X=x_1] P(X=x_1) - \mathbb{E}[Y(x_1) | X=x_0] P(X=x_0).$$

$$+ \underline{\mathbb{E}[Y(x_1) | X \neq x_1] P(X \neq x_1)} - \underline{\mathbb{E}[Y(x_0) | X \neq x_0] P(X \neq x_0)}.$$

Unknown. Unknown.

$$\therefore \widehat{\text{ATE}} = \mathbb{E}[Y(x_0) | X=x_1] P(X=x_1) + \mathbb{E}[Y(x_0) | X \neq x_1] P(X=x_1) + 1 \cdot \mathbb{P}(X \neq x_1) - 0 \cdot \mathbb{P}(X \neq x_0)$$

$$\widehat{\text{ATE}} = \mathbb{E}[Y(x_0) | X=x_1] P(X=x_1) + \mathbb{E}[Y(x_0) | X \neq x_1] P(X=x_1) + 0 \cdot \mathbb{P}(X \neq x_1) - 1 \cdot \mathbb{P}(X \neq x_0)$$

$$\therefore \text{width} = \mathbb{P}(X \neq x_1) + \mathbb{P}(X \neq x_0)$$

$$= 1 - \mathbb{P}(X=x_1) + 1 - \mathbb{P}(X=x_0)$$

$$= 2 - [\mathbb{P}(X=x_1) + \mathbb{P}(X=x_0)].$$

Randomized Experiments.

$x \in X$: all logically feasible Treatment Values.

$Y(x)$: potential outcomes.

$Y = Y(X)$: realized outcomes.

Identification of Marginal Distributions of Potential Outcomes:

Theorem: Suppose $X \perp\!\!\!\perp \{Y(x) : x \in X\}$

randomization \Leftrightarrow
statistical independence
assumption.

marginal Suppose the joint dist. of (Y, X) is known. Then,
distribution $F_{Y(x)}$ is point identified for all $x \in X$.
of any potential proof: $\forall y \in \mathbb{R}, x \in \text{supp}(X)$:

$$\text{known in data.} \quad P(Y \leq y | X=x) = P(Y(x) \leq y | X=x)$$

$$\text{by randomization:} \quad = P(Y(x) \leq y)$$

$$= F_{Y(x)}(y).$$

$$\text{Means: } P(Y \leq \cdot | X=x) = P(Y(x) \leq \cdot | X=x) = P(Y(x) \leq \cdot | X=x')$$

$\forall x, x' \in \mathbb{R}$, know one, know the others.

Look at the dist. of observed outcomes among people with treatment values $X=x$ to learn about the counterfactual outcomes of the people with $x \neq x$.

Likewise, any parameter that depends only on the marginal distributions of $Y(x)$ is point identified.

$$\text{ATE}(x_0 \rightarrow x_1) = E[Y(x_1)] - E[Y(x_0)]$$

$$\text{QTE}(z, x_0 \rightarrow x_1) = Q_{Y(x_1)}(z) - Q_{Y(x_0)}(z).$$

Identification of the joint distribution of potential outcomes: partially.

Since we only observe one value of $(Y_{i(1)}, Y_{i(0)})$ for each i .

Interpreting Randomization:

Randomization \Leftrightarrow statistical independence assumption

$$X \perp\!\!\!\perp \{Y(x) : x \in X\}$$

Measuring:

Fix some $x_{i(t)} \in X$, look at dist. of $Y(x_{i(t)})$ in diff. groups: the set of units with $X=x_i$ for $x_i \in X, x_i \neq x_j$ in each group, dist. of $Y(x_{i(t)})$ is the same.

In binary setting:

$Y(x) | X=1$ has the same dist. of $Y(x) | X=x_0$.

$$\forall x \in \{0, 1\}.$$

Covariates in experiments:

W : vector of covariates.

Distinguish: ① W^{pre} : Not causally affected by Treatment

$$X \perp\!\!\!\perp \{Y(x) : x \in X\}, W^{\text{pre}}$$

② W^{post} : Might be affected.

$$X \not\perp\!\!\!\perp \{Y(x) : x \in X\}, W^{\text{post}}.$$

5 reasons to include:

To test randomization

Suppose treatment randomly assigned.

$$\Rightarrow X \perp\!\!\!\perp \{Y(x) : x \in X\}, W^{\text{pre}} \text{ (falsifiable)}$$

Compare $P(W^{\text{pre}} | X=1)$ with $P(W^{\text{pre}} | X=0)$

$$\text{If different, then: } X \not\perp\!\!\!\perp \{Y(x) : x \in X\}, W^{\text{pre}}$$

To Improve Precision:

Suppose Treatment randomly assigned.

$$X \perp\!\!\!\perp \{Y(x) : x \in X\}, W^{\text{pre}}.$$

Suppose X is binary treatment, then:

$$\text{ATE} = E[Y | X=1] - E[Y | X=0]$$

Consider Two OLS estimands:

① $\text{On } (1, X)$: asymptotic variance of OLS estimator is

$$\frac{\text{Var}(Y^{\perp X})}{\text{Var}(X)}.$$

② $\text{On } (1, X, W^{\text{pre}})$, $X \perp\!\!\!\perp W^{\text{pre}} \Rightarrow$ coefficient on X same. Asymptotic covariance of OLS on estimator is:

$$\frac{\text{Var}(Y^{\perp X, W^{\text{pre}}})}{\text{Var}(X^{\perp W^{\text{pre}}})}.$$

Since $X \perp\!\!\!\perp W^{\text{pre}}$, $X^{\perp W^{\text{pre}}} = X$

$$\therefore \frac{\text{Var}(Y^{\perp X, W^{\text{pre}}})}{\text{Var}(X)}.$$

Since $\text{Var}(A^{\perp B}) \leq \text{Var}(A)$, $(A^{\perp B})^{\perp C} = A^{\perp B, C}$,

$$\frac{\text{Var}(Y^{\perp X, W^{\text{pre}}})}{\text{Var}(X)} \leq \frac{\text{Var}(Y^{\perp X})}{\text{Var}(X)}.$$

Stratified Randomized Experiments:

researchers assign treatments randomly within groups defined by pre-treatment covariates W .

Theorem: Suppose $X \perp\!\!\!\perp \{Y(x) : x \in X\} | W$

Then $F_{Y(x)(\cdot | W)}(\cdot | w)$ is point identified for all $x \in \text{supp}(X | W=w)$, $w \in \text{supp}(W)$.

p.f.: $\forall y \in \mathbb{R}$ and $(w, x) \in \text{supp}(W, X)$,

$$P(Y \leq y | W=w, X=x) = P(Y(x) \leq y | W=w, X=x)$$

$$\text{known in the data.} \quad = P(Y(x) \leq y | W=w, X=x)$$

$$= P(Y(x) \leq y | W=w)$$

$$= F_{Y(x)(\cdot | W)}(y | w)$$

Proposition Suppose $X \perp\!\!\!\perp \{Y(x) : x \in \mathcal{X}\} | W$.

Then $F_{Y(x)}(\cdot)$ not necessarily $= F_{Y|X}(\cdot|x)$.

• Use Experiments to Learn about Subgroup Treatment Effect

Suppose treatment unconditionally randomly assigned.

$$X \perp\!\!\!\perp \{Y(x) : x \in \mathcal{X}\}, W_{pre}.$$

Lemma: $A \perp\!\!\!\perp (B, C) \Rightarrow A \perp\!\!\!\perp B | C$

Proposition: Suppose $X \perp\!\!\!\perp \{Y(x) : x \in \mathcal{X}\}, W_{pre}$, Then $F_{Y(x)|W_{pre}}(w)$ is point identified for any $x \in \mathcal{X}$, $w \in \text{supp}(W_{pre})$.

$$\text{CATE} = \mathbb{E}[Y_{(1)} - Y_{(0)} | W_{pre} = w]$$

Selection on Observables.

Randomly Assigned Treatments with Observational Data.
not experimental data.

In observational data setting, $X \perp\!\!\!\perp \{Y(x) : x \in \mathcal{X}\} | W$ is the unconfoundedness assumption. (= conditional independence).
→ Binary Treatment.

Assumption ①: Unconfoundedness: non falsifiable.
 $X \perp\!\!\!\perp (Y_{(1)}, Y_{(0)}) | W$

Assumption ②: Overlap: (sufficient variation).

$P(X=1 | W=w) \in (0, 1)$ for all $w \in \text{supp}(W)$
falsifiable.

Main Identification Results.

$Y = Y(x)$ potential outcomes.

Theorem: Suppose the unconfoundedness and overlap hold.

Suppose the distribution of (Y, X, W) is observed.

Then: ① $F_{Y(x)}|W(\cdot|w)$ is point identified for all $x \in \{0, 1\}$ and $w \in \text{supp}(W)$.

② $F_{Y(x)}$ is point identified for all $x \in \{0, 1\}$.

Proof: Part 1:

$$\forall y \in \mathbb{R}, w, x \in \text{supp}(W, X):$$

$$\begin{aligned} P(Y \leq y | W=w, X=x) &= P(Y(x) \leq y | W=w, X=x) \\ &= P(Y_{(x)} \leq y | W=w, X=x) \\ &= P(Y(x) \leq y | W=W) \\ &= F_{Y(x)}|W(y | W) \end{aligned}$$

Overlap: $\forall x \in \{0, 1\}$, LHS is known for all $w \in \text{supp}(W)$.

Thus the RHS is known for all $x \in \{0, 1\}$ and $w \in \text{supp}(W)$.

Part 2:

$$\begin{aligned} F_{Y(x)>y} &= P(Y(x) > y) \\ &= \mathbb{E}[P(Y(x) > y | W)] \\ &= \mathbb{E}[F_{Y(x)|W}(y | W)] \end{aligned}$$

RHS: ① Marginal distribution of W (\mathbb{E}). Known from data.

② $F_{Y(x)|W}(\cdot | W)$ for all $w \in \text{supp}(W)$. Known from Part 1.

The theory implies that any parameters that only depends on $F_{Y(x)|W}(\cdot | W)$ and $F_{Y(x)}$ is point identified.

Including:

$$\text{① CATE: } \text{CATE}(w) = \mathbb{E}[Y_{(1)} - Y_{(0)} | W=w]$$

$$\begin{aligned} \text{② ATE: } \text{ATE}(w) &= \mathbb{E}[Y_{(1)} - Y_{(0)}] \\ &= \mathbb{E}[\text{CATE}(W)] \end{aligned}$$

$$\text{③ CATT: } \text{CATT}(w) = \mathbb{E}[Y_{(1)} - Y_{(0)} | W=w, X=1]$$

$$\text{④ ATT: } \text{ATT} = \mathbb{E}[Y_{(1)} - Y_{(0)} | X=1]$$

$$\text{⑤ T-th QTE: } \text{QTE}(T) = Q_{Y_{(1)}}(T) - Q_{Y_{(0)}}(T)$$

$$\text{⑥ T-th QTT: } \text{QTT}(T) = Q_{Y_{(1)}}|_{X(T=1)} - Q_{Y_{(0)}}|_{X(T=0)}$$

Accessing Unconfoundedness: What to control for?

Set up: Y : outcome variable of interest.

X : treatment

$W = (W_1, \dots, W_K)$ a vector of covariates.

U : a vector of unobserved variables.

$$(t) \begin{cases} Y \leftarrow g(X, W, U) \\ X \leftarrow h(Y, W, U) \\ W_k \leftarrow m_k(Y, X, W_{-k}, U) \text{ for } k=1 \dots K. \end{cases}$$

A graph G : a set of nodes and a set of edges between them.

A directed G : have \rightarrow . ordered pair of nodes.

A causal graph: nodes are variables.

An walk from v_i to v_j is a sequence of nodes x_0, \dots, x_k such that ① $x_0 = v_i$ ② $x_k = v_j$ ③ there is an edge btw x_i and x_{i+1} , $i \in \{0, 1, \dots, k-1\}$.

A directed walk: all arrows are in same direction.

A path: all the nodes in a walk are different.
+ directed walk = directed path.

A directed cycle: a directed walk that starts and ends at the same node, all others are distinct.

A acyclic graph: does not contain any directed cycles.

DAG: directed acyclic graph.

Baseline Assumptions:

① Initial conditions: U in (*) is external.

No $\rightarrow U$, only $U \rightarrow$.

② No simultaneity: causal graph of (*) is acyclic.
not possible for two variables to affect each other

③ \rightarrow for any joint dist. of U , (*) yields a unique distribution of (Y, X, W) .

Simple example:

W : a scalar.

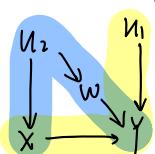
$$Y \leftarrow g(X, W, U_1)$$

$$X \leftarrow h(U_2)$$

$$W \leftarrow m(U_2)$$

U_1, U_2 independent

Causal Graph:



path: Y collider

path: Y not collider

$X \& Y$ d-separated:
every path has a collider

$X \& Y$ d-connected:
there is a path with no
collider.

Conditional Independence in Causal Graphs.

structure of the causal graph \rightarrow unconfoundedness assumption

d-separation

Definition: Consider a directed graph between V_1 and V_2 . A

collider on this path is any node on the path with two incoming arrows. (common effect) d: directional.

Definition: A path is d-connected if it contains no collider.

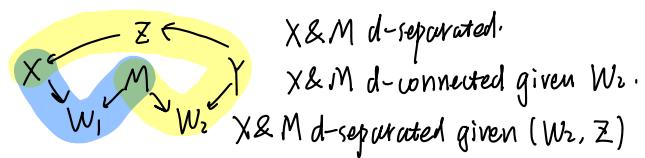
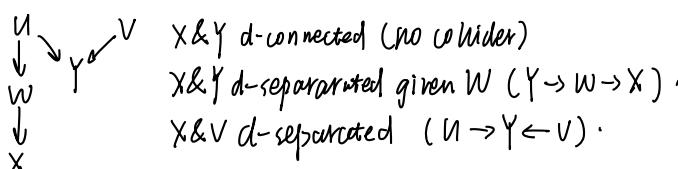
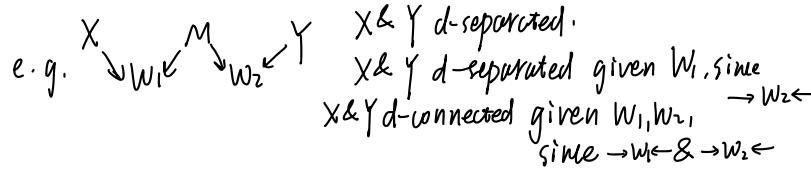
Definition: The nodes V_1 and V_2 are d-connected if there is at least one path between them that is d-connected.

If no path between them is d-connected, then V_1 and V_2 are d-separated.

Including Conditioning Variables.

Definition: Let W be a set of variables. V_1 and V_2 are d-connected given W if there is at least one path from V_1 to V_2 that satisfies all of the following criteria:

1. there are no chains with $W \in W$ as a middle node.
that is: no $A \rightarrow W \rightarrow B$.
2. there are no forks with $W \in W$ as a middle node.
that is: no $A \leftarrow W \rightarrow B$.
3. for all colliders C on the path, ($A \rightarrow C \leftarrow B$):
Either: a. C is in W
b. at least one of the descendants of C is in W .

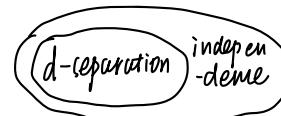


d-separation implies independence

Theorem: Let G be a causal graph corresponding to (*).

Suppose assumption ①&② holds. Let X, Y, W be subsets of variables (Y, X, W, U) . Then for any joint dist of U :

$$\begin{aligned} X \text{ and } Y \text{ are d-separated given } W \\ \Downarrow \\ X \perp\!\!\!\perp Y \mid W. \end{aligned}$$



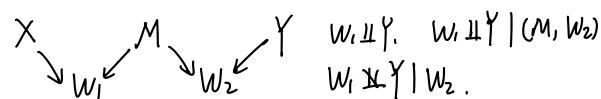
However,
 $X \perp\!\!\!\perp Y \mid W \not\Rightarrow$ d-separated.

Intuition behind d-separation:

① Mediators: $Y \leftarrow \beta_Y M + U_Y$ $X \rightarrow M \rightarrow Y$
Structural System: $M \leftarrow \beta_M x + U_M$ $X \& Y$ d-connected $X \perp\!\!\!\perp Y$
 $X \leftarrow U_X$ $X \& Y$ d-separated given M
all U mutually indep. $X \perp\!\!\!\perp Y \mid M$

② Common Causes: $Y \leftarrow \beta_Y M + U_Y$ $X \leftarrow M \rightarrow Y$
Structural System: $X \leftarrow \beta_{X,M} M + U_X$ $X \& Y$ d-connected $X \perp\!\!\!\perp Y$
 $M \leftarrow U_M$ $X \& Y$ d-separated given M
all U mutually indep. $X \perp\!\!\!\perp Y \mid M$

③ Colliders: $M \leftarrow \beta_{M,X} X + \beta_{M,Y} Y + U_M$ $X \leftarrow M \rightarrow Y$
Structural System $X \leftarrow U_X$ $X \& Y$ d-separated $X \perp\!\!\!\perp Y$
 $Y \leftarrow U_Y$ $X \& Y$ d-connected given M
all U mutually indep. $X \perp\!\!\!\perp Y \mid M$



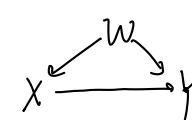
Deriving Sufficient Conditions for Unconfoundedness

Start from:

$$Y \leftarrow g(X, W, U_Y)$$

$$X \leftarrow h(W, U_X)$$

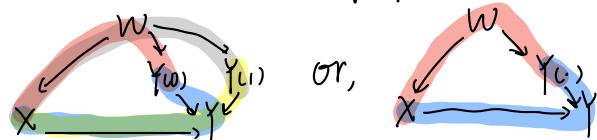
$$W \leftarrow U_W$$



U 's are mutually indep.

Define potential outcome: Then: $Y \leftarrow g(X, W, U_Y)$
 $Y(x) = g(x, W, U_Y)$ $Y(x) \leftarrow g(x, W, U_Y)$
 $X \leftarrow h(W, U_X)$
 $W \leftarrow U_W$

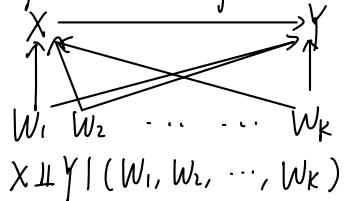
Then we have the following graph



Using the graph theory, $X \& Y_{(0)}$ are desegregated condition on W . Thus, $X \perp\!\!\!\perp Y(x) \mid W$. (unconfoundedness).

In general, condition on post-treatment outcomes is dangerous.

The Unconfoundedness Temporale Causal Graph:



Generally, should control for all observed pre-treatment variables.

Should avoid imprecision from over control.

It has two drawbacks:

- ①. increase standard errors in the estimators.



But, control for W reduces the remaining variation in X .

- ②. make overlap worse

$$p(w) = P(X=1 \mid W=w) \rightarrow 0 \text{ or } \rightarrow 1 \text{ when } W \text{ includes too many.}$$

Proposition: There are always at least two different graphs in \mathcal{G} , where \mathcal{G} is the set of all causal graphs where $X \perp\!\!\!\perp Y(x) \mid W$.

- pf: ① $X \rightarrow Y$
② $X \leftarrow W \rightarrow Y$

Thus, assumptions on causal graphs provides sufficient, not necessary, conditions for unconfoundedness.

Benefits of Causal Graphs: ① arguably more interpretable than unconfoundedness alone ② make the model falsifiable.

The Classical Linear Model for Potential Outcomes.

Functional form of potential outcomes are linear in a set of known functions of the treatment with constant coefficients.

Benefits: ① treatment effect parameter equivalent to certain OLS estimands
② more subtle identification results ③ solve extrapolation identifications.

$X \in \mathcal{X}$: a vector of treatment variables.

$x \in \mathbb{R}^k$: all logical values of these variables.

$Y(x)$: potential outcomes $\forall x \in \mathcal{X}$

Assumption: $\forall x \in \mathcal{X}, Y(x) = p(x)' \beta + u$. $p(x)$: J vector.

restrictions: ① linear functional form

② homogeneous treatment effect β .

$$Y_i(x_{i1}) - Y_i(x_0) = (p(x_i) - p(x_0))' \beta.$$

Exogeneity.

① full independence: $X \perp\!\!\!\perp Y(x)$

under linear functional form: $X \perp\!\!\!\perp U$

② mean independence: $\mathbb{E}[U \mid X=x] = \mathbb{E}[U]$

③ uncorrelatedness: $\text{Corr}(U, p_j(X)) = 0$.

Theorem: Identification of the classical linear model:

Suppose $F_{Y,X}$ satisfies:

① Finite moments on observables: $\mathbb{E}[p(X)Y] < \infty, \mathbb{E}[p(X)p(X)'] < \infty$

② Sufficient Variation: $\mathbb{E}[p(X)p(X)']$ invertible.

③ Linear potential outcomes

④ finite moments on unobservables: $\mathbb{E}[p(X)U] < \infty, \mathbb{E}[U] < \infty$.

⑤ Exogeneity: $\mathbb{E}[p(X)U] = \mathbb{E}[p(X)] \mathbb{E}[U]$.

$$\text{Corr}(U, p_j(X)) = 0$$

⑥ normalization: $\mathbb{E}[U] = 0$.

Then β is point identified.

$$\beta = \mathbb{E}[p(X)p(X)']^{-1} \mathbb{E}[p(X)Y].$$

⑦ ⑧ ⑨ not falsifiable, then the identified set for β is a singleton. $B_2 = \{\mathbb{E}[p(X)p(X)']^{-1} \mathbb{E}[p(X)Y]\}$.

Proof: $Y = Y(x)$

$$= p(X)' \beta + U$$

$$p(X)Y = p(X)p(X)' \beta + p(X)U.$$

$$\mathbb{E}[p(X)Y] = \mathbb{E}[p(X)p(X)'] \beta + \mathbb{E}[p(X)U].$$

$$= \mathbb{E}[p(X)] \mathbb{E}[U] = 0.$$

$$\beta = \mathbb{E}[p(X)p(X)']^{-1} \mathbb{E}[p(X)Y]$$

Suppose $p(U)$ does not contain a constant.

Let $\tilde{p}(x) = [p(x)']$ then: $\mathbb{E}[\tilde{p}(x)\tilde{p}(x)']$ invertible.

define $\tilde{\beta} = [\frac{\mathbb{E}[U]}{\beta}]$, $\tilde{U} = U - \mathbb{E}[U]$

then $Y = \tilde{p}(x)' \tilde{\beta} + \tilde{U}$ holds.

To prove B_2 is a singleton. Have to construct $(\tilde{p}, \tilde{F}_{X,Y})$ consistent with $F_{Y,X}$ such that $\tilde{\beta} = \mathbb{E}[p(X)p(X)']^{-1} \mathbb{E}[p(X)Y]$.

$$\text{① define } \tilde{p} = [\frac{\mathbb{E}[U]}{\beta}].$$

$$\text{② define } \tilde{U} = Y - p(X)' \tilde{\beta}$$

Thus, $\tilde{F}_{X,U} = F_{X,\tilde{U}}$. $(\tilde{p}, \tilde{F}_{X,U})$ consistent with linear ...

$$\mathbb{E}[p(X)\tilde{U}] = 0, \mathbb{E}[\tilde{U}] = 0 \text{ if } p(X) \text{ consists a constant}$$

∴ this construction is consistent with all model assumptions.

$$\begin{aligned} P(Y \in y | X=x) &= P(p(x)' \tilde{\beta} + \tilde{U} \in y | X=x) \\ &= P(\tilde{U} \in y - p(x)' \tilde{\beta} | X=x) \\ &= F_{\tilde{U}|X}(y - p(x)' \tilde{\beta} | X=x). \end{aligned}$$

Alternative proof:

$$\begin{aligned} \text{Suppose } p(x) = \begin{pmatrix} 1 \\ x \end{pmatrix}, Y(x) = \beta_0 + \beta_1 x + U. \\ \text{cov}(Y, X) &= \text{cov}(p(x)' \beta + U, X) \\ &= 0 + \beta_1 \text{var}(x) + \underbrace{\text{cov}(U, X)}_{= E[UX] - E[U]E[X]} \\ &= \beta_1 \text{var}(x) \\ \therefore \beta_1 &= \frac{\text{cov}(Y, X)}{\text{var}(x)} \\ &= E[UX] - E[U]E[X] \\ &= 0 \end{aligned}$$

Connecting Regression and Causality

The linear potential outcomes model with uncorrelatedness gives a causal interpretation of OLS estimand.

Corollary. Suppose all the assumptions hold.

$$\text{Let } E = Y - p(x)' \beta^{\text{OLS}}$$

where $\beta^{\text{OLS}} = E[p(x)p(x)']^{-1} E[p(x)Y]$

$$\text{Then, } U \equiv E.$$

Interpretation: Under uncorrelatedness assumption about $X \perp\!\!\!\perp U$ in linear potential outcome outcomes, U (unobserved heterogeneity) equals E (regression residual), which is by construction uncorrelated to the regressors $p_j(x)$.

Identification without the sufficient variation assumption

Lemma. Suppose Sufficient Variation does not hold.

① β is not point identified.

② The identified set: $\beta_2 = \{b \in \mathbb{R}^J : E[p(x)p(x)']b = E[p(x)Y]\}$

Identification under mean independence.

Lemma. Suppose $F_{Y|X}$ satisfies: ① sufficient variation ② linear potential outcomes ③ finite moments ④ exogeneity.

(If $p(x)$ contains a constant, then Normalization holds)

Then, β and $F_{Y|X}$ are point identified.

Proposition: Suppose all the assumptions hold. If $p(x)$ does not contain a constant,

$$E[Y|X=x] = p(x)' \beta + E[U]$$

If $p(x)$ does contain a constant,

$$E[Y|X=x] = p_1(x)' \beta_1 + (\beta_0 + E[U]).$$

This implies that the model is falsifiable since the conditional mean function must satisfy this specific form.

$$\begin{aligned} \text{Proof: } E[Y|X=x] &= E[p(x)' \beta + U | X=x] \\ &= E[p(x)' \beta | X=x] + E[U | X=x] \\ &= p(x)' \beta + E[U] \end{aligned}$$

Connecting Regression and Causality

Step ①. linear potential outcome model implies

$$A_{\text{OLS}}(x) = p(x)' \beta + E[U]$$

Step ②. add exogeneity assumption $\rightarrow \beta$ point identified.
under mean independence, $E[Y|X=x] = p(x)' \beta + E[U]$.

Theorem Suppose all the assumptions satisfy. then

$$A_{\text{OLS}}(x) = p(x)' \beta + E[U] = E[Y|X=x]$$

And it is point identified.

Thus $E[Y|X=x]$ has a causal interpretation.

Identification under full independence

Assumption: Exogeneity: $X \perp\!\!\!\perp U$ (implies $Y(x) \perp\!\!\!\perp X$).

Full independence has further falsifying power.

Proposition. Suppose $F_{Y|X}$ satisfies sufficient variation

$$\text{Let } E[Y|X=x] = p(x)' \beta + E[U].$$

$$\text{Define } \tilde{U} = Y - (p(x)' \beta + E[U])$$

Then, linear potential outcomes, finite moments and exogeneity are falsified iff $\tilde{U} \perp\!\!\!\perp X$.

$$\begin{aligned} \text{Var}[Y|X=x] &= \text{Var}(p(x)' \beta + U | X=x) \\ &= \text{Var}(p(x)' \beta + \tilde{U} | X=x) \\ &= \text{Var}(\tilde{U} | X=x) \\ &= \text{Var}(\tilde{U}) \quad (\text{since } X \perp\!\!\!\perp U). \end{aligned}$$

\rightarrow homoskedasticity

Covariates in the Linear Model

Reasons to include covariates.

- ① to test assumptions like unconfoundedness.
- ② to improve the statistical precision of estimators.
- ③ necessary for causality identification analysis. Here!
- ④ need them to define other parameter of interest.
- ⑤ use them to solve identification problems like missing data.

Nonparametric unconfoundedness:

$$\{Y(x) : x \in X\} \perp\!\!\!\perp X | W$$

Also assume $Y(x) = p(x)' \beta + U$, then:

$$\{Y(x) : x \in X\} \perp\!\!\!\perp X | W \Leftrightarrow U \perp\!\!\!\perp X | W$$

$$E[U | X=x, W=w] = E[U | W=w]$$

$$E[Y | X=x, W=w] = E[Y(x) | W]$$

$$\begin{array}{c} \uparrow \\ \text{known} \end{array} = p(x)' \beta + E[U | W=w].$$

\uparrow to identify this,

need some version of sufficient variation assumption.

Linear Unconfoundedness \Leftrightarrow Nonparametric Unconfoundedness

$$Y(x) = p(x)' \beta + u.$$

Here, we use "unrelatedness" after partitioning out W .

Assumption: $\text{corr}(p_j(X)^{\perp W}, U^{\perp W}) = 0$ for all $j \in \{1, \dots, J\}$

Remember: $\text{cov}(A^L, B^L) = \text{cov}(A^L, B)$, $\text{cov}(A^L, D^L)$

then: $\text{cov}(p_j(X)^{\perp W}, U) = 0$, $\text{cov}(p_j(X), U^{\perp W}) = 0$.

Identification:

Suppose: ① linear potential outcomes ② finite moments, sufficient variation, $\text{var}(Y, p(x), W)$ exists, positive definite

③ linear unconfoundedness.

Then: 1. β and X, u point identified.

$$2. \beta = \mathbb{E}[p(x)^{\perp W} (p(x)^{\perp W})' \mathbb{E}[p(x)^{\perp W} Y^{\perp W}]]$$

by FWL, β is the coefficient on $p(x)$ in OLS of Y on $(p(x), W)$.

$$3. \beta_2 = \mathbb{E}[p(x)^{\perp W} (p(x)^{\perp W})' \mathbb{E}[p(x)^{\perp W} Y^{\perp W}]].$$

Suppose that $p(x) = \begin{bmatrix} 1 \\ x \end{bmatrix}$. Let $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$

then: $Y(x) = \beta_0 + \beta_1 x + u$. since $(A+B)^L = A^L + B^L$

$$\text{cov}(Y^{\perp W}, X^{\perp W}) = \text{cov}(\beta_0 + \beta_1 x^{\perp W} + u^{\perp W}, X^{\perp W})$$

$$= 0 + \beta_1 \text{cov}(X^{\perp W}, X^{\perp W}) + \text{cov}(u^{\perp W}, X^{\perp W})$$

$$= \beta_1 \text{var}(X^{\perp W}).$$

$$\text{sufficient variation} \Rightarrow \beta_1 = \frac{\text{cov}(Y^{\perp W}, X^{\perp W})}{\text{var}(X^{\perp W})}$$

* With linear unconfoundedness assumptions, coefficients on $p(x)$ can be interpreted causally. Coefficient on W cannot!

In practice most covariates are thought to be endogenous.

The bias from including outcomes as control variables

① endogenous themselves

② destroys the exogeneity of X . *

Proof 1st.

Set up.

Suppose we have two outcomes, Y_1, Y_2

$$\begin{aligned} Y_1(x) &= \alpha_1 + \beta_1 x + u_1 & u_1 \longrightarrow Y_1 \\ Y_2(x) &= \alpha_2 + \beta_2 x + u_2 & u_2 \perp u_1 \\ \text{Suppose } X &\perp\!\!\!\perp (u_1, u_2), \quad X \longrightarrow Y_1 \\ \text{then, } \beta_1 &\text{ is point identified.} \end{aligned}$$

If include Y_2 as control, do OLS of Y_1 on $(1, X, Y_2)$, then

$$\text{coeff. on } X = \frac{\text{cov}(Y_1, X^{\perp Y_2})}{\text{var}(X^{\perp Y_2})}$$

let $(\hat{\beta}_0, \hat{\beta}_1)$ be the coefficient of X on $(1, Y_2)$, residual = $X^{\perp Y_2}$

$$\text{cov}(Y_1, X^{\perp Y_2}) = \text{cov}(\beta_1 x + u_1, X^{\perp Y_2})$$

$$= \beta_1 \text{cov}(x, X^{\perp Y_2}) + \text{cov}(u_1, X^{\perp Y_2})$$

$$= \beta_1 \text{var}(X^{\perp Y_2}) + \text{cov}(u_1, X - [\hat{\beta}_0 + \hat{\beta}_1 Y_2])$$

$$= \beta_1 \text{var}(X^{\perp Y_2}) - \hat{\beta}_1 \text{cov}(u_1, Y_2).$$

$$\text{Hence } \frac{\text{cov}(Y_1, X^{\perp Y_2})}{\text{var}(X^{\perp Y_2})} = \beta_1 - \hat{\beta}_1 \frac{\text{cov}(u_1, Y_2)}{\text{var}(X^{\perp Y_2})}$$

$$\begin{aligned} \text{Next, } \text{cov}(Y_2, u_1) &= \text{cov}(\alpha_2 + \beta_2 x + u_2, u_1) \\ &= \text{cov}(u_2, u_1) \end{aligned}$$

$$\text{Moreover, } \hat{\beta}_1 = \frac{\text{cov}(X, Y_2)}{\text{var}(Y_2)}$$

$$= \frac{\text{cov}(X, \alpha_2 + \beta_2 x + u_2)}{\text{var}(Y_2)}$$

$$= \beta_2 \frac{\text{var}(X)}{\text{var}(Y_2)}$$

$$= \beta_2 \frac{\text{var}(X)}{\text{var}(\alpha_2 + \beta_2 x + u_2)}$$

$$= \beta_2 \frac{\text{var}(X)}{\beta_2^2 \text{var}(X) + \text{var}(u_2)}$$

$$\text{Therefore, } \frac{\text{cov}(Y_1, X^{\perp Y_2})}{\text{var}(X^{\perp Y_2})} = \beta_1 - \beta_2 \frac{\text{var}(X)}{\beta_2^2 \text{var}(X) + \text{var}(u_2)} \frac{\text{cov}(u_2, u_1)}{\text{var}(X^{\perp Y_2})}$$

If either ① $\beta_2 = 0$ (Y_2 actually not a control) or

② $\text{cov}(u_2, u_1) = 0 \rightarrow$ generally will not hold.

then, adding Y_2 to be a covariate will not cause bias.

prop 2nd.

Basic idea in using regression to get causality:

Variation in the conditional mean function yields causal effects.

$$\mathbb{E}[Y_1 | X=x, Y_2=y_2] = \mathbb{E}[x + \beta_1 x + u_1 | X=x, Y_2=y_2]$$

$$= \alpha_1 + \beta_1 x + \mathbb{E}[u_1 | X=x, Y_2=y_2]$$

$$= \alpha_1 + \beta_1 x + \mathbb{E}[u_1 | X=x, u_2=y_2 - \alpha_2 - \beta_2 x]$$

$$= \alpha_1 + \beta_1 x + \mathbb{E}[u_1 | u_2=y_2 - \alpha_2 - \beta_2 x]$$

$$\frac{\partial \mathbb{E}[Y_1 | X=x, Y_2=y_2]}{\partial x} = \beta_1 + \underbrace{\frac{\partial \mathbb{E}[u_1 | u_2=y_2 - \alpha_2 - \beta_2 x]}{\partial x}}_{*}$$

If ① $\beta_2 = 0$, then * = 0 or

② u_1 mean independent of u_2

then, adding Y_2 to be a covariate will not cause bias.

Mean Independence with Endogenous Controls.

$$Y(x, w) = \beta_0 + \beta_1 x + \beta_2 w + u.$$

Mean independence assumption: $\mathbb{E}[u | X=x, W=w] = \mathbb{E}[u]$

* Rarely possible! Instead, we plot:

Assumption: Conditional mean independence (control function).

$$\mathbb{E}[u | X=x, W=w] = \mathbb{E}[u | W=w].$$

Assumption: Linearity.

$$\mathbb{E}[u | W=w] = \gamma_0 + \gamma_1 w . \text{ allow for } \gamma_1 \neq 0 . *$$

do OLS of Y on $(1, X, W)$:

$$\begin{aligned}
\mathbb{E}[Y|X=x, W=w] &= \mathbb{E}[\beta_0 + \beta_1 X + \beta_2 W + U | X=x, W=w] \\
&= \beta_0 + \beta_1 x + \beta_2 w + \mathbb{E}[U | X=x, W=w] \\
&= \beta_0 + \beta_1 x + \beta_2 w + \mathbb{E}[U | W=w] \\
&= \beta_0 + \beta_1 x + \beta_2 w + \gamma_0 + \gamma_1 w \\
&= (\beta_0 + \gamma_0) + \beta_1 x + (\beta_2 + \gamma_1)w.
\end{aligned}$$

$\frac{\partial \mathbb{E}[Y|X=x, W=w]}{\partial x} = \beta_1$, even if W is endogenous!
i.e. β_1 is point identified.

But, $\frac{\partial \mathbb{E}[Y|X=x, W=w]}{\partial W} = \beta_2 + \gamma_1 \neq \beta_2$, biased!

Proposition: Under the above assumptions, the coefficient on X from OLS of Y on $(1, X, W)$ is β_1 .

Proof: by FWL: coeff = $\frac{\text{cov}(Y, X^{\perp W})}{\text{var}(X^{\perp W})}$.

$$\begin{aligned}
\text{cov}(Y, X^{\perp W}) &= \text{cov}(\beta_0 + \beta_1 X + \beta_2 W + U, X^{\perp W}) \\
&= \beta_1 \text{cov}(X, X^{\perp W}) + \beta_2 \text{cov}(W, X^{\perp W}) + \text{cov}(U, X^{\perp W}) \\
&= \beta_1 \text{var}(X^{\perp W}) + \text{cov}(U, X^{\perp W}).
\end{aligned}$$

$$\begin{aligned}
\text{cov}(U, X^{\perp W}) &= \mathbb{E}[X^{\perp W}(U - \mathbb{E}[U])] \\
&= \mathbb{E}[X^{\perp W}] \quad \text{assume } \mathbb{E}[U] = 0 \\
&= \mathbb{E}_{x,w}[\mathbb{E}[X^{\perp W}U | x, w]] \\
&= \mathbb{E}_{x,w}[\mathbb{E}(U | x, w)X^{\perp W}] \\
&= \mathbb{E}_{x,w}[\mathbb{E}[U | w]X^{\perp W}] \\
&= \mathbb{E}_{x,w}[(\gamma_0 + \gamma_1 w)X^{\perp W}] \\
&= \text{cov}(\gamma_0 + \gamma_1 w, X^{\perp W}) \quad \text{since } \mathbb{E}[X^{\perp W}] = 0 \\
&= 0
\end{aligned}$$

Non-parametric endogenous controls

Without $\mathbb{E}[U | W] = \gamma_0 + \gamma_1 W$,

$$\mathbb{E}[Y | X=x, W=w] = \beta_0 + \beta_1 x + \beta_2 w + \mathbb{E}[U | W=w]$$

$\beta_1 = \frac{\partial \mathbb{E}[Y | X=x, W=w]}{\partial x}$ is identified!

Important: $Y(x) \perp\!\!\!\perp X | W \Leftrightarrow Y(x) \perp\!\!\!\perp W | X$

Using conditional mean independence \rightarrow allow endogenous control.

Assessing Linearity

The non-parametric additively separable model

Mean independence:

$$m(x) = \mathbb{E}[Y | X=x] = p(x)' \beta + \mathbb{E}[U]$$

Suppose that linear model is false.

Theorem: Suppose the joint distribution of (Y, X) is observed. And: Additively Separable Potential Outcomes, Finite Moments, Mean Independence. Then: $g(x) + \mathbb{E}[U]$ is point identified for all $x \in \text{supp}(X)$. model is not falsifiable.

Proof: $\forall x \in \text{supp}(X)$:

$$\begin{aligned}
\mathbb{E}[Y | X=x] &= \mathbb{E}[g(x) + U | X=x] \\
&= \mathbb{E}[g(x) + \mathbb{E}[U] | X=x] \\
&= g(x) + \mathbb{E}[U] \\
&= g(x) + \mathbb{E}[U]. \text{ identified.}
\end{aligned}$$

Corollary 7: Suppose all the relevant assumption hold.

1. Let $\tilde{U} = U - \mathbb{E}[U]$. Then $F_{\tilde{U}}$ is point identified.
2. $F_{Y|x}$ is point identified for all $x \in \text{supp}(X)$
3. The joint distribution of any finite set of potential outcomes $(Y_{(x_1)}, \dots, Y_{(x_k)})$ is point identified, so long as $x_1, \dots, x_k \in \text{supp}(X)$.

Proof: 1. we have $U = Y - g(x)$

$$\text{so } \tilde{U} = Y - (g(x) + \mathbb{E}[U]).$$

\uparrow then known $\underbrace{\mathbb{E}[U]}$ known

$$\begin{aligned}
\text{That is: } F_{\tilde{U}}(u) &= P(\tilde{U} \leq u) \\
&= P(Y - g(x) - \mathbb{E}[U] \leq u) \\
&= \int_{\text{supp}(Y|x)} (y - g(x) - \mathbb{E}[U]) dF_{Y|x}(y|x)
\end{aligned}$$

It is known.

$$2. Y(x) = g(x) + \mathbb{E}[U] + \tilde{U}$$

$$F_{Y|x}(y) = P(Y(x) \leq y)$$

$$\begin{aligned}
&\uparrow \\
&= P(g(x) + \mathbb{E}[U] + \tilde{U} \leq y) \\
&= P(\tilde{U} \leq y - g(x) - \mathbb{E}[U])
\end{aligned}$$

then known $= F_{\tilde{U}}(y - g(x) - \mathbb{E}[U]) \leftarrow$ known.

3. for example, just consider two values $x_1, x_2 \in \text{supp}(X)$

$$F_{Y|x_1, Y|x_2}(y_1, y_2) = P(Y(x_1) \leq y_1, Y(x_2) \leq y_2)$$

$$= P(\tilde{U} \leq y_1 - g(x_1) - \mathbb{E}[U], \tilde{U} \leq y_2 - g(x_2) - \mathbb{E}[U])$$

$$= P(\tilde{U} \leq \min\{y_1 - g(x_1) - \mathbb{E}[U], y_2 - g(x_2) - \mathbb{E}[U]\})$$

$$= F_{\tilde{U}}(\min\{y_1 - g(x_1) - \mathbb{E}[U], y_2 - g(x_2) - \mathbb{E}[U]\})$$

then point identified. \uparrow known

Cost for allowing a fully non-parametric function form:
 cannot extrapolate off the support.

Full independence

Proposition: $m(x) = \mathbb{E}[Y | X=x]$. Define $\tilde{U} = Y - m(x)$

Then. Additively separable + finite moments + full independence are falsified iff $\tilde{U} \not\perp\!\!\!\perp X$.

The non-parametric potential outcomes model

replace $X \perp\!\!\!\perp U$ with: $Y(x) \perp\!\!\!\perp X$ for all $x \in \text{supp}(X)$

Thus $F_{Y|x}(y) = P(Y \leq y | X=x)$ is point identified.

Best linear approximation ASF = $m(x)$.

Let $\beta^{ou} = \mathbb{E}[p(x)p(x)']^{-1} \mathbb{E}[p(x)Y]$. Then within a class of functions $f_p(x) = p(x)' b$, $p(x)' \beta^{ou}$ is the best approximation to ASF using weighted L_2 distance.

Assessing Exogeneity

Reasons for exogeneity failure ($\mathbb{E}[p(x)u] \neq \mathbb{E}[p(x)]\mathbb{E}[u]$)

- ① Simultaneity
- ② Omitted variables.

Interpreting OLS without exogeneity

Dominic Variable Bias X, W scalar

$$Y(x, w) = \beta_0 + \beta_1 x + \beta_2 w + V$$

$$\text{Let } Y(x) = Y(x, W) = \beta_0 + \beta_1 x + (\beta_2 w + V) = \beta_0 + \beta_1 x + U; U \equiv \beta_2 w + V.$$

$$\text{Suppose } \text{cov}(X, V) = 0, \text{ cov}(W, V) = 0.$$

$$\text{do OLS on } (I, X, W) \rightarrow \beta = (\beta_0, \beta_1, \beta_2)$$

$$\text{do OLS on } (I, X): \text{cov}(Y, X) = \text{cov}(\beta_0 + \beta_1 X + \beta_2 W + V, X)$$

$$\stackrel{\text{true}}{=} \beta_1 \text{Var}(X) + \beta_2 \text{Cov}(W, X)$$

$$\text{then } \frac{\text{cov}(Y, X)}{\text{Var}(X)} = \beta_1 + \beta_2 \frac{\text{cov}(W, X)}{\text{Var}(X)}$$

$$\text{Bias} = \frac{\text{cov}(Y, X)}{\text{Var}(X)} - \beta_1 = \beta_2 \frac{\text{cov}(W, X)}{\text{Var}(X)}$$

$$\text{Depend on: } ① \beta_2 \quad ② \frac{\text{cov}(W, X)}{\text{Var}(X)}$$

$$\text{Different way of seeing: } \text{cov}(X, U) = \text{cov}(X, \beta_2 w + V) \\ = \beta_2 \text{cov}(X, W).$$

For mean independence: $\mathbb{E}[V|X=x] = \mathbb{E}[V]$.

$$\text{Then } \mathbb{E}[U|X=x] = \mathbb{E}[\beta_2 w + V|X=x] \\ = \beta_2 \mathbb{E}[w|X=x] + \mathbb{E}[V|X=x] \\ = \beta_2 \mathbb{E}[w|X=x]$$

if $\beta_2 \neq 0$ and mean of W depend on X , then

$$\mathbb{E}[U|X=x] \neq \mathbb{E}[U]$$

$$\text{Consider } \mathbb{E}[Y|X=x] = \mathbb{E}[\beta_0 + \beta_1 x + U|X=x]$$

$$= \beta_0 + \beta_1 x + \mathbb{E}[U|X=x]$$

$$\frac{\partial \mathbb{E}[Y|X=x]}{\partial x} = \beta_1 + \frac{\partial \mathbb{E}[U|X=x]}{\partial x}$$

$$\text{Suppose } \mathbb{E}[W|X=x] = S, \text{ then } \frac{\partial \mathbb{E}[U|X=x]}{\partial x} = \beta_2 S$$

$$\mathbb{E}[Y|X=x] = \beta_0 + (\beta_1 + \beta_2 S)x, \frac{\partial \mathbb{E}[Y|X=x]}{\partial x} = \beta_1 + \beta_2 S$$

Heterogeneous Treatment Effects: Random Coefficients

Motivation: the unit level causal effects change from person to person.

Non-parametric potential outcomes model

$$Y(x) = g(x, u).$$

$$Y_i(x_i) - Y_j(x_j) = g(x_i, u_i) - g(x_j, u_j)$$

Random Coefficient Models β random

$$Y(x) = \beta x + u$$

$$Y_i(x) = \beta_i x + u_i$$

The OLS estimand:

$$Y(x) = \beta x + u$$

$$\text{cov}(Y, X) = \mathbb{E}[(Y - \mathbb{E}[Y])(X - \mathbb{E}[X])] \\ = \mathbb{E}[Y(X - \mathbb{E}[X])] - \mathbb{E}[Y]\mathbb{E}[X - \mathbb{E}[X]]$$

$$= \mathbb{E}[Y(X - \mathbb{E}[X])].$$

$$= \mathbb{E}[(\beta X + u)(X - \mathbb{E}[X])] \quad \text{assume } \mathbb{E}[u|X=x] = \mathbb{E}[u]$$

$$= \mathbb{E}[\beta X(X - \mathbb{E}[X])]$$

$$\text{Also assume } \mathbb{E}[\beta|X=x] = \mathbb{E}[\beta]$$

$$\text{then } \text{cov}(Y, X) = \mathbb{E}[\beta X(X - \mathbb{E}[X])]$$

$$= \mathbb{E}[\mathbb{E}[\beta|X]X(X - \mathbb{E}[X])]$$

$$= \mathbb{E}[\beta]\mathbb{E}[X(X - \mathbb{E}[X])]$$

$$= \mathbb{E}[\beta]\text{Var}(X) \quad \text{then,}$$

$$\text{Theorem } Y(x) = \beta x + u, \mathbb{E}[u|x] = \mathbb{E}[u], \mathbb{E}[\beta|x] = \mathbb{E}[\beta]$$

$$\text{then: } \mathbb{E}[\beta] = \frac{\text{cov}(Y, X)}{\text{Var}(X)} = \text{average random effect}$$

$$\text{ATE}(x_1 \rightarrow x_2) = \mathbb{E}[Y(x_1) - Y(x_2)] = \mathbb{E}[(\beta x_1 + u) - (\beta x_2 + u)]$$

$$= \mathbb{E}[\beta(x_2 - x_1)] = \mathbb{E}[\beta](x_2 - x_1)$$

Heterogeneous Treatment Effects: Quantile Regression (most common approach).

Theorem (Equivalence): $h: \mathbb{R} \rightarrow \mathbb{R}$ weakly left-continuous.

$$Y: \mathbb{R} \text{ v. } \text{then } Q_{h(Y)}(z) = h(Q_Y(z)).$$

Proof: suppose F_Y strictly ↑ everywhere, h strictly ↑. h continuous.

$$\mathbb{P}(Y \leq Q_Y(z)) = z.$$

$$\mathbb{P}(Y \leq y) = \mathbb{P}(h(Y) \leq h(y)) \Rightarrow z = \mathbb{P}(Y \leq Q_Y(z)) \\ = \mathbb{P}(h(Y) \leq h(Q_Y(z)))$$

$\therefore h(Q_Y(z))$ is the z^{th} quantile of $h(Y)$, $= Q_{h(Y)}(z)$.

Model 1: The additively separable model with quantile independence

Theorem potential outcome: $Y = g(x) + u$.

observed outcome $Y = Y(x)$. (Y, X) observed.

Assume: ① $Q_{h(Y)}(z|x) = Q_{h(Y)}(z)$, $Q_{h(Y)}(z)$ is known

Then: $g(x), F_x, u$ are point identified

$$\text{proof: } Q_{Y|X}(z|x) = Q_{g(x)+u|X}(z|x)$$

$$= Q_{g(x)+u}(z|x)$$

$$= g(x) + Q_{u|X}(z|x) \quad (\text{Equivalence})$$

$$= g(x) + Q_u(z)$$

$$g(x) = \underbrace{Q_{Y|X}(z|x) - Q_u(z)}_{\text{known from the data}} \text{ identified}$$

Corollary: point identification of the classical linear model.

Suppose $Y(x) = p(x)\beta + u$.

Assume: $\text{Med}(u|Y=x) = \text{Med}(u)$, $\text{Med}(u) = 0$,

then $\theta = (\beta, F_x, u)$ is point identified.

Model 2: Non-separable Unobservables

If $Y(x) = g(x) + u$, then $\frac{\partial Y(x)}{\partial x} = g'(x)$. not heterogeneous.

To model heterogeneous treatment effect: $Y(x) = g(x, u)$, $\frac{\partial Y(x)}{\partial x} = g'(x, u)$

Theorem: Suppose $Y(x) = g(x, u)$

$$Y = Y(x), \text{ observe } (Y, x).$$

Assume: ① U scalar ② $p(x)' \beta(u)$ weakly ↑, left cont. ③ $X \perp\!\!\!\perp U$ ④ $U \sim \text{unif}[0, 1]$

Then: $g(x, u)$ is point identified.

Pf:
$$\begin{aligned} Q_{Y|X}(z|x) &= Q_{g(X, u)|x}(z|x) \\ &\stackrel{\substack{\text{known from} \\ \text{data.}}}{=} Q_{g(x, u)|x}(z|x) \\ &= Q_{g(x, u)}(z) \\ &= g(x, Q_u(z)) \text{ since } U \sim \text{unif}[0, 1]. \\ &= g(x, z) \rightarrow \text{identified}. \end{aligned}$$

Model 3. The linear quantile regression model.

Theorem: Identification of linear quantile regression model.

$$Y(x) = p(x)' \beta(u).$$

$Y = Y(x)$. (Y, X) observed. Assumes:

① U scalar ② $p(x)' \beta(u)$ weakly increasing, left continuous
③ $X \perp\!\!\!\perp U$. ④ $U \sim \text{unif}[0, 1]$. ⑤ no multicollinearity.

Then: $\beta: [0, 1] \rightarrow \mathbb{R}$ is point identified.

Pf: let $g(x, u) = p(x)' \beta(u)$, then by previous Theorem, ✓.

$$Q_{Y|X}(z|x) = p(x)' \beta(z).$$

$$ASF(x) = p(x)' \mathbb{E}[\beta(u)]$$

$$= p(x)' \int_0^1 \beta(u) du. \text{ since } U \sim \text{unif}[0, 1].$$

$$\begin{aligned} ATE(x_1 \rightarrow x_2) &= ASF(x_2) - ASF(x_1) \\ &= p(x_2)' \mathbb{E}[\beta(u)] - p(x_1)' \mathbb{E}[\beta(u)]. \\ &= (p(x_2) - p(x_1))' \mathbb{E}[\beta(u)]. \end{aligned}$$

$$\begin{aligned} Q_{Y|X}(z) &= Q_{p(x)' \beta(u)}(z) \\ &= p(x)' \beta(Q_u(z)) \\ &= p(x)' \beta(z). \end{aligned}$$

$$\begin{aligned} QTE(x_1 \rightarrow x_2) &= Q_{Y|X}(z) - Q_{Y|X}(z) \\ &= p(x_2)' \beta(z) - p(x_1)' \beta(z) \\ &= (p(x_2) - p(x_1))' \beta(z) \end{aligned}$$

Instrumental Variables.

Simultaneity.

Motivation: to identify supply and demand curve.

Observe in data: (P, Q, Z_1, Z_2) .

Unknown parameters: $(\alpha_1, \alpha_2, \gamma_1, \gamma_2, \beta_1, \beta_2)$.

When are these parameters point identified?

Model:
$$\begin{cases} P = \alpha_1 + \gamma_1 Q + \beta_1 Z_1 + U_1 \\ Q = \alpha_2 + \gamma_2 P + \beta_2 Z_2 + U_2 \end{cases}$$

Let Z_1 be the supply shifter, Z_2 be the demand shifter.

Then:

Theorem: $\begin{cases} P = \alpha_1 + \gamma_1 Q + \beta_1 Z_1 + \beta_1 Z_2 + U_1 \text{ (Supply)} \\ Q = \alpha_2 + \gamma_2 P + \beta_2 Z_2 + \beta_2 Z_1 + U_2 \end{cases}$

Assume: ① unique equilibrium: $\gamma_1 \gamma_2 \neq 1$
an IV $\begin{cases} \text{② exclusion: } \beta_1 = 0, \beta_2 = 0 \\ \text{satisfies } \text{③ exogeneity: } \text{cov}(Z_1, U_1) = 0, \text{cov}(Z_1, U_2) = 0 \\ \text{④ relevant: } \beta_1 \neq 0, \beta_2 \neq 0. \\ \text{⑤ normalization: } \mathbb{E}[U_1] = 0, \mathbb{E}[U_2] = 0. \end{cases}$

Then, (γ_1, γ_2) are point identified.

Proof #1: Assume $\alpha_1 = 0, \alpha_2 = 0$:

$$\begin{cases} P = \gamma_1 Q + \beta_1 Z_1 + \beta_1 Z_2 + U_1 \\ Q = \gamma_2 P + \beta_2 Z_1 + \beta_2 Z_2 + U_2 \end{cases} \Leftrightarrow \begin{bmatrix} 1 & -\gamma_1 \\ -\gamma_2 & 1 \end{bmatrix} \begin{bmatrix} P \\ Q \end{bmatrix} = \begin{bmatrix} \beta_1 & 0 \\ 0 & \beta_2 \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} + \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}$$

$$\text{Solve it, } P = \frac{\beta_1}{1 - \gamma_1 \gamma_2} Z_1 + \frac{\gamma_2 \beta_1}{1 - \gamma_1 \gamma_2} Z_2 + \frac{U_1 + \gamma_1 U_2}{1 - \gamma_1 \gamma_2}$$

$$\text{We get: } Q = \frac{\gamma_2 \beta_1}{1 - \gamma_1 \gamma_2} Z_1 + \frac{\beta_2}{1 - \gamma_1 \gamma_2} Z_2 + \frac{U_2 + \gamma_2 U_1}{1 - \gamma_1 \gamma_2}$$

Reduced form: endo... = f(exo...)

$$\begin{aligned} \text{let } (\pi_{11}, \pi_{12}) &= \left(\frac{\beta_1}{1 - \gamma_1 \gamma_2}, \frac{\gamma_2 \beta_1}{1 - \gamma_1 \gamma_2} \right), \quad V_1 = \frac{U_1 + \gamma_1 U_2}{1 - \gamma_1 \gamma_2} \\ (\pi_{21}, \pi_{22}) &= \left(\frac{\gamma_2 \beta_1}{1 - \gamma_1 \gamma_2}, \frac{\beta_2}{1 - \gamma_1 \gamma_2} \right), \quad V_2 = \frac{U_2 + \gamma_2 U_1}{1 - \gamma_1 \gamma_2} \end{aligned}$$

then $P = \pi_{11} Z_1 + \pi_{12} Z_2 + V_1$

$$\begin{aligned} Q &= \pi_{21} Z_1 + \pi_{22} Z_2 + V_2 \\ \mathbb{E}[Z_1 V_1] &= \mathbb{E}\left[Z_1 \frac{U_1 + \gamma_1 U_2}{1 - \gamma_1 \gamma_2}\right] \\ &= \frac{\mathbb{E}[Z_1 U_1] + \gamma_1 \mathbb{E}[Z_1 U_2]}{1 - \gamma_1 \gamma_2} \\ &= 0 \text{ by ③ and ⑤.} \end{aligned}$$

$\mathbb{E}[Z_2 V_1] = 0$ also.

Then, π_{11} and π_{12} are point identified.

$\mathbb{E}[Z_1 V_2] = \mathbb{E}[Z_2 V_2] = 0$ by similar calculation

Then π_{21} and π_{22} are point identified.

Then, $\frac{\pi_{12}}{\pi_{22}} = \gamma_1, \frac{\pi_{21}}{\pi_{11}} = \gamma_2$ point identified.

Then, $\beta_1 = \pi_{11}(1 - \gamma_1 \gamma_2), \beta_2 = \pi_{22}(1 - \gamma_1 \gamma_2)$ point identified.

Proof 2:

$$\begin{aligned} \text{cov}(P, Z_2 | Z_1) &= \text{cov}(\alpha_1 + \gamma_1 Q + \beta_1 Z_1 + U_1, Z_2 | Z_1) \\ &= \gamma_1 \text{cov}(Q, Z_2 | Z_1) + \beta_1 \text{cov}(Z_1, Z_2 | Z_1) + \text{cov}(U_1, Z_2 | Z_1) \\ &= \gamma_1 (\text{cov}(Q, Z_2 | Z_1) + \text{cov}(U_1, Z_2 | Z_1)) \end{aligned}$$

If strengthen assumption:

$\mathbb{E}[U_1 | Z_1, Z_2] = 0, \mathbb{E}[U_2 | Z_1, Z_2] = 0$, then:

$$\begin{aligned} \text{cov}(U_1, Z_2 | Z_1) &= \mathbb{E}[U_1 Z_2 | Z_1] - \mathbb{E}[U_1 | Z_1] \mathbb{E}[Z_2 | Z_1] \\ &= \mathbb{E}[Z_2 \mathbb{E}[U_1 | Z_2, Z_1] | Z_1] - \mathbb{E}[\mathbb{E}[U_1 | Z_2, Z_1] | Z_1] \mathbb{E}[Z_2 | Z_1] \\ &= 0 \text{ by ③ and ⑤.} \end{aligned}$$

Then $\gamma_1 = \frac{\text{cov}(P, Z_2 | Z_1)}{\text{cov}(Q, Z_2 | Z_1)}$ point identified.

Next: we show that $\text{cov}(Q, Z_2 | Z_1) \neq 0$.

$$\begin{aligned}\text{cov}(Q, Z_2 | Z_1) &= \text{cov}(\alpha_2 + \gamma_2 P + \beta_2 Z_2 + U_2, Z_2 | Z_1) \\ &= \gamma_2 \text{cov}(P, Z_2 | Z_1) + \beta_2 \text{var}(Z_2 | Z_1) + \text{cov}(U_2, Z_2 | Z_1) \\ &= \gamma_2 \text{cov}(P, Z_2 | Z_1) + \beta_2 \text{var}(Z_2 | Z_1) \text{ by } ③\end{aligned}$$

Then, $\text{cov}(P, Z_2 | Z_1) = \gamma_1 \text{cov}(Q, Z_2 | Z_1)$
 $= \gamma_1 \gamma_2 \text{cov}(P, Z_1 | Z_1) + \gamma_1 \beta_2 \text{var}(Z_2 | Z_1)$

then, $\text{cov}(P, Z_2 | Z_1) = \frac{\gamma_1 \beta_2 \text{var}(Z_2 | Z_1)}{1 - \gamma_1 \gamma_2}$

then, $\text{cov}(Q, Z_2 | Z_1) = \gamma_2 \frac{\gamma_1 \beta_2 \text{var}(Z_2 | Z_1)}{1 - \gamma_1 \gamma_2} + \beta_2 \text{var}(Z_2 | Z_1)$
 $= \frac{\beta_2}{1 - \gamma_1 \gamma_2} \text{var}(Z_2 | Z_1) \neq 0 \text{ since } \beta_2 \neq 0, \text{var}(Z_2 | Z_1) > 0$

Similarly: $\gamma_2 = \frac{\text{cov}(Q, Z_1 | Z_2)}{\text{cov}(P, Z_1 | Z_2)}$, point identified.

$$\beta_2 = (1 - \gamma_1 \gamma_2) \frac{\text{cov}(Q, Z_2 | Z_1)}{\text{var}(Z_2 | Z_1)} \text{ point identified.}$$

Similarly, $\beta_1 = (1 - \gamma_1 \gamma_2) \frac{\text{cov}(P, Z_1 | Z_2)}{\text{var}(Z_1 | Z_2)}$ point identified.

$V_1 = E[P] - \gamma_1 E[Q] - \beta_1 E[Z_1]$, α_2 similarly.

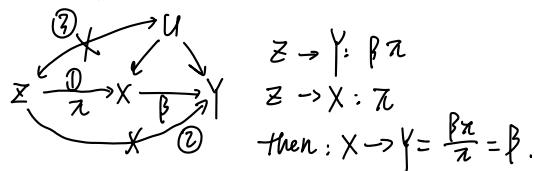
then, $(\alpha_1, \alpha_2, \beta_1, \beta_2, \gamma_1, \gamma_2)$ all identified

$$Y(X, Z) = \beta X + \delta Z + U \quad X, Z \text{ scalar.}$$

$$\text{cov}(Y, Z) = \text{cov}(\beta X + \delta Z + U, Z)$$

$$= \beta \text{cov}(X, Z) + \delta \text{var}(Z) + \text{cov}(U, Z)$$

$$\therefore \beta = \frac{\text{cov}(Y, Z)}{\text{cov}(X, Z)} - \delta \frac{\text{var}(Z)}{\text{cov}(X, Z)} - \frac{\text{cov}(U, Z)}{\text{cov}(X, Z)}$$



Assumption: ① relevant: $\text{cov}(X, Z) \neq 0$.

② exclusion: $S = 0$, Z no direct causal effect on Y .

③ instrument exogeneity: $\text{cov}(U, Z) = 0$.
 (use Z randomly assigned)

Then: $\beta = \frac{\text{cov}(Y, Z)}{\text{cov}(X, Z)}$ (2SLS estimand)

$$\beta = \left[\frac{\text{cov}(Y, Z)}{\text{cov}(X, Z)} \right] = \left[\frac{\text{cov}(Y, Z)}{\text{var}(Z)} \right] / \left[\frac{\text{cov}(X, Z)}{\text{var}(Z)} \right]$$

Second stage
need exclusion
Causal effect of $Z \rightarrow Y$. need exogeneity
ITT

reduced form
 \downarrow
falsifiable $\text{cov}(X, Z) = 0$

$$Y = \beta X + U \quad \text{first stage, get } X^{\text{pred}} = \pi Z$$

$$= \beta (\pi Z + X^{\perp \beta}) + U$$

$$= \underbrace{\beta \pi Z}_{\text{reduced form}} + \underbrace{\beta X^{\perp \beta}}_{\frac{\text{cov}(Y, Z)}{\text{var}(Z)}} + U$$

$$Y = \beta X^{\text{pred}} \Rightarrow \beta = \frac{\text{cov}(Y, X^{\text{pred}})}{\text{var}(X^{\text{pred}})} = \frac{\beta \pi}{\pi}$$

Generally, $Y(x) = p(x)' \beta + u$

now: $\text{corr}(p_k(x), u) \neq 0$.

Theorem: Consider $Y(x, z) = p(x)' \beta + q(z)' \delta + u$

Assume: ① exclusion: $S = 0$ K -vector L -vector

② exogeneity: $E[q(z)u] = E[q(z)]E[u]$

③ relevance: $\text{rank}(E[q(z)p(x)']) = K$

④ Normalization: $p(x)$ consist a constant and $E[u] = 0$

Then: β is point-identified.

When $L = k$: $\beta = E[q(z)p(x)']^{-1} E[q(z)Y]$

proof: we have $0 = E[q(z)u]$

$$= E[q(z)(Y - p(x)'\beta)]$$

$$= E[q(z)Y] - E[q(z)p(x)']\beta.$$

$$E[q(z)p(x)']\beta = E[q(z)Y]$$

since $\text{rank}(E[q(z)p(x)']) = K$, full rank.

And if $L = k$, then:

$$\beta = E[q(z)p(x)']^{-1} E[q(z)Y].$$

if $L > K$, then: choose a $K \times L$ Π :

$$\beta = E[\Pi q(z)p(x)']^{-1} E[\Pi q(z)Y]$$

(assume invertible, new instrument $\Pi q(z)$).

Pick optimal Π : same to WLS to minimize variance.

Under homoskedasticity, optimal Π = residual of $\pi Z + X^2$.

If linear model does not hold, different instrument.

Triangular Models

$$\begin{cases} Y(x) = g(x, u) & \text{observed outcome, } Y = g(x, u) \\ X(z) = h(z, u) & X = h(z, u) \end{cases}$$

Assume additively separable & linear.

$$Y(x) = \gamma_0 + \gamma_1 x + u$$

$$X(z) = \tau_0 + \tau_1 z + v$$

Assume Z exogenous. $E[ZV] = E[Z]E[V]$, $\tau_1 \neq 0$, $E[zv] = E[z]E[v]$

then: $(\gamma_0, \gamma_1, \tau_0, \tau_1)$ identified.

Heterogeneous Treatment Effects with IV: Random Coefficients