

## Introduction

- Using economic model of agent behaviour to extrapolate outside of the support of the data ex-ante is known as structural analysis. Otherwise ex-post analysis of policies usually falls under reduced form analysis.
- This paper focuses on discussing program evaluation for ex post reduced form analysis.

## Causality and Potential Outcome Model

- POM writes the realized outcome in terms of the potential outcome and treatment

$$\hookrightarrow Y_i = D_i Y_{1i} + (1 - D_i) Y_{0i}, \begin{cases} D_i \in \{0,1\} \text{ is binary treatment} \\ (Y_{0i}, Y_{1i}) \text{ are potential outcomes} \end{cases}$$

- SUTVA rules out interference between units. This is violated in educational setting with peer effects as realized outcome depends on treatment status of other units.
- The unit level causal effect is not identified from the data.

$$\hookrightarrow Y_{1i} - Y_{0i} \text{ where } (Y_{0i}, Y_{1i}) \text{ are potential outcomes}$$

- The causal estimands are generally joint distributions of potential outcomes, treatment, and covariates. The estimand is causal when it just depends on the distribution of the potential outcomes.
- The ATE and ATT only depend on the distribution of the potential outcomes, and hence they are causal estimands.

$$\hookrightarrow \begin{cases} \text{ATE} = E(Y_{1i} - Y_{0i}) \\ \text{ATT} = E(Y_{1i} - Y_{0i} | D_i = 1) \end{cases}, \text{ where } D_i \in \{0,1\}$$

- If potential outcomes are not independent of the treatment, than neither the ATE or ATT is identified from the mean outcome differences across control and treatment.
- Confounding arises when treatment selection is dependant on potential outcome.
- Let us consider the difference in mean outcome between control and treatment:

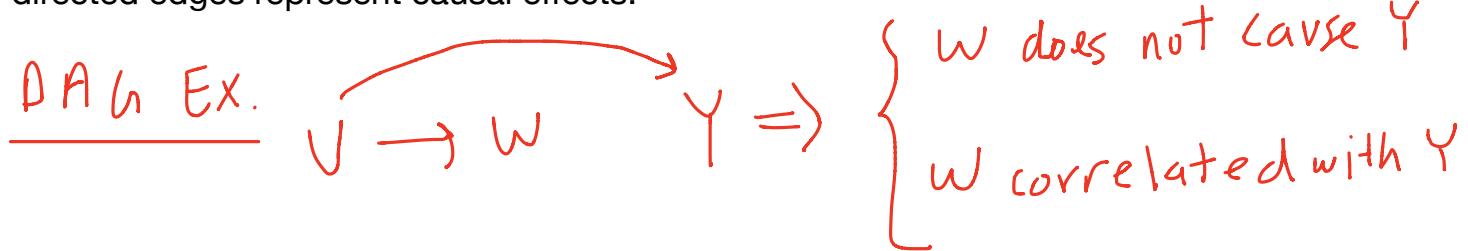
$$\tau = E(Y_i | D_i = 1) - E(Y_i | D_i = 0) \Rightarrow$$

$$\begin{aligned} \textcircled{1} \quad \tau &= E(Y_{1i} - Y_{0i} | D_i = 1) + [E(Y_{0i} | D_i = 1) - E(Y_{0i} | D_i = 0)] \\ &= \text{ATT} + \text{Bias}_{\text{ATT}} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \tau &= E(Y_{1i} - Y_{0i}) + E(Y_{0i}) - E(Y_{1i}) + E(Y_{1i} | D_i = 1) - E(Y_{0i} | D_i = 0) \\ &= \text{ATE} + \Pr(D_i = 1) \text{ATT} + (1 - \Pr(D_i = 1)) \text{ATT} \\ &= \text{ATE} + \text{Bias}_{\text{ATE}} \end{aligned}$$

- In general, neither the ATE and ATT are identified in observational data as shown above.

- A DAG is a visual summary of all the causal effects of interest, nodes are random variables, directed edges represent causal effects.



- Above DAG is a good example of the cliche "correlation does not imply causation". W does not cause Y, but W is correlated to Y because the confounder U effects both W and Y.

## Randomized Experiments

- Randomization of treatment implies treatment is independent of potential outcomes.

$$\hookrightarrow D_i \perp (Y_{0i}, Y_{1i})$$

- With randomization, the bias terms for ATT and ATE are both 0 and hence the ATE is identified from comparing mean outcomes across control and treatment.
- Note that the ATT = ATE under randomization of the treatment.

$$\hookrightarrow ATE = ATT = E(Y_i | D_i = 1) - E(Y_i | D_i = 0) = E(Y_{1i} - Y_{0i})$$

$$D_i \perp (Y_{0i}, Y_{1i})$$

- The ATE is a function of the marginal distribution of the potential outcomes, they are both identified under randomization.

$$\hookrightarrow \begin{cases} F_{Y_i}(y) = \Pr(Y_i \leq y) = \Pr(Y \leq y | D=1) \\ F_{Y_0}(y) = \Pr(Y_0 \leq y) = \Pr(Y \leq y | D=0) \end{cases}$$

Identified since  $D$  is randomized

- Since the potential outcomes marginal distributions are identified from the data, the quantile treatment effect is also identified.

Identified quantile dbn

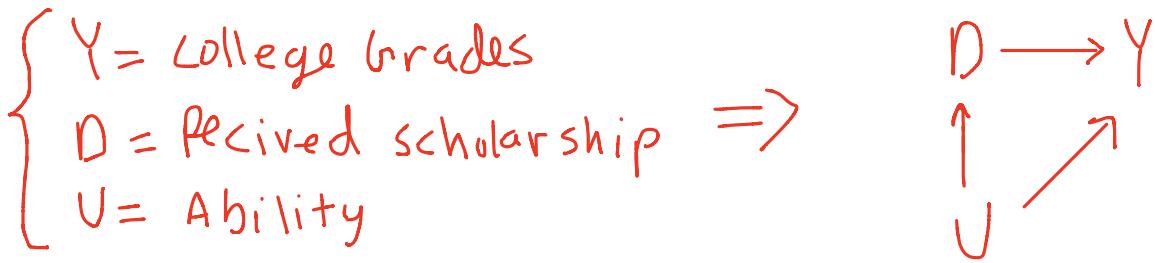
$$\begin{cases} Q_{Y_i}(\theta) = Q_{Y|D=1}(\theta), \theta \in (0,1) \text{ quantile of } Y \\ Q_{Y_0}(\theta) = Q_{Y|D=0}(\theta) \end{cases}$$

$$\Rightarrow \text{Quantile Treat. Effect} = Q_{Y_i}(\theta) - Q_{Y_0}(\theta)$$

## Selection on Observables

- Suppose we are interested in learning the impact of receiving a scholarship on college grades. This is difficult to learn as those who receive scholarships will be of higher ability than those who don't and may receive higher grades even if they had not gotten the scholarship.

- The following DAG represents the above problem:



- The direct causal effect we are interested in is from D to Y, but it is confounded by U.

$\hookrightarrow \text{Paths} \begin{cases} \text{Direct: } D \rightarrow Y \\ \text{Backdoor: } D \leftarrow U \rightarrow Y \end{cases}$

- We can block the backdoor path by conditioning on U. We can only do this if U is observed.
- Above example motivates the conditional independence assumption (CIA):

$\hookrightarrow \text{CIA: } (Y_{1i}, Y_{0i}) \perp D_i \mid X_i$

- Under the CIA the conditional average treatment effect is identified:

$$E(Y_{1i} - Y_{0i} \mid X_i) = E(Y_i \mid X_i, D_i=1) - E(Y_i \mid X_i, D_i=0)$$

↓  
CIA

- CIA along the common support assumption can identify the ATT and ATE:

Assume  $\begin{cases} \text{CIA: } (Y_{0i}, Y_{1i}) \perp D_i \mid X_i \\ \text{Common support: } P(D_i=1 \mid X_i) \in (0, 1) \end{cases} \Rightarrow$

$$\textcircled{1} \text{ ATT} = E_{X_i \mid D_i=1} [E(Y_i \mid X_i, D_i=1) - E(Y_i \mid X_i, D_i=0)]$$

$$\textcircled{2} \text{ ATE} = E_x [E(Y_i \mid X_i, D_i=1) - E(Y_i \mid X_i, D_i=0)]$$

- Note that both the ATT and ATE are just weighted averages of the within covariate cell treatment effects. We can see this more clearly by letting X take on a finite number of values  $X_1, \dots, X_m$ :

$$\textcircled{1} \text{ ATT} = \sum_{R=1}^m [E(Y_i \mid X_i=x_R, D_i=1) - E(Y_i \mid X_i=x_R, D_i=0)] P(X_i=x_R \mid D_i=1)$$

$$\textcircled{2} \text{ ATE} = \sum_{R=1}^m [E(Y_i | X_i = x_R, D_i = 1) - E(Y_i | X_i = x_R, D_i = 0)] P(X_i = x_R)$$

- It is noteworthy to mention that regressing Y on D and X and assuming CIA will allow us to recover either the ATT or ATE.

$$Y_i = \alpha + \gamma D_i + \beta X_i + \varepsilon_i \Rightarrow \gamma \notin \{\text{ATE}, \text{ATT}\}$$

- The coefficient on D will represent the average of the within covariate cell treatment effects weighted by the variance of the treatment.

$$\gamma = \sum_{R=1}^m [E(Y_i | X_i = x_R, D_i = 1) - E(Y_i | X_i = x_R, D_i = 0)] w_R,$$

$$\text{where } w_R = \frac{V(D_i | X_i = x_R) P(X_i = x_R)}{\sum_{r=1}^m V(D_i | X_i = x_r) P(X_i = x_r)},$$

$$\text{with } V(D_i | X_i = x_R) = P(D_i = 1 | X_i = x_R) (1 - P(D_i = 1 | X_i = x_R))$$

- The OLS will associate covariate cells that have the most common support across control and treatment group with the largest weights. These weights lower the variance of estimator.
- We can show the ATT within covariate cell effects are weighted by the probability of treatment:

$$\text{ATT} = \sum_{R=1}^m [E(Y_i | X_i = x_R, D_i = 1) - E(Y_i | X_i = x_R, D_i = 0)] w_R^{\text{ATT}},$$

$$w_R^{\text{ATT}} = P(X_i = x_R | D_i = 1) \stackrel{\text{Bayes}}{=} \frac{P(D_i = 1 | X_i = x_R) P(X_i = x_R)}{\sum_{r=1}^m P(D_i = 1 | X_i = x_r) P(X_i = x_r)}$$

### Matching Estimators

- Although OLS cannot be used to estimate the ATT or ATE under CIA, matching estimators can estimate these causal parameters of interest.
- Matching involves pairing each observation i with another observation j(i) that is of the opposite treatment group but with very similar covariates.

For each i, find j(i) such that  $\begin{cases} D_{j(i)} = 1 - D_i \\ X_{j(i)} \approx X_i \end{cases}$

- The matching pairs can be used to identify both the ATT and ATE:

$$\textcircled{1} \text{ ATT} = \frac{1}{n} \sum_{i=1}^n D_i (Y_i - Y_{j(i)})$$

$$\textcircled{2} \text{ ATE} = \frac{1}{n} \sum_{i=1}^n D_i (Y_i - Y_{j(i)}) - \frac{1}{n} \sum_{i=1}^n (1 - D_i) (Y_i - Y_{j(i)})$$

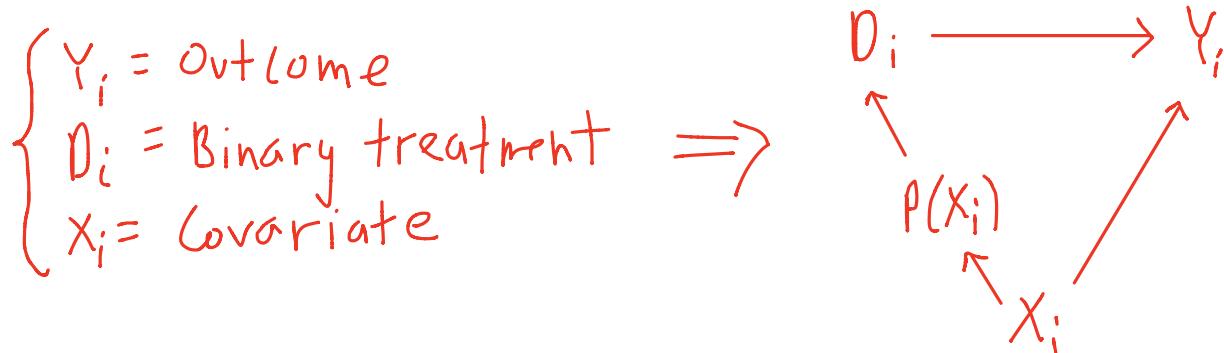
- Another type of matching is known as "prosperity score matching". The propensity score is:

$$\hookrightarrow P(X_i) = P(D_i = 1 | X_i)$$

- The prosperity score essentially summarizes the information contained in the covariates. That is if CIA is attained from conditioning on  $X$  then CIA is also obtained from conditioning on  $P(X)$ :

$$\text{If } (Y_{ii}, Y_{oi}) \perp\!\!\!\perp D_i | X_i \Rightarrow (Y_{ii}, Y_{oi}) \perp\!\!\!\perp D_i | P(X_i)$$

- The DAG below illustrates the above implication:



- The ATE can be identified using propensity score matching:

$$\text{ATE} = \frac{1}{n} \sum_{i=1}^n \frac{D_i Y_i}{P(X_i)} - \frac{1}{n} \sum_{i=1}^n \frac{(1 - D_i) Y_i}{1 - P(X_i)}$$

- Observations with large  $P(X)$  are overrepresented in the treatment group, so they are weighted down when in treatment group, and weighted up when in control group. The vice-versa is true for small  $P(X)$ .

### Differences-in-Differences

- DD model is identified after making certain assumptions on how unobserved confounders effects the outcome over time.
- Consider the standard panel data model:

$$Y_{it} = D_{it} \gamma_{it} + \mu_i + \delta_t + \varepsilon_{it}, \quad D_{it} \in \{0, 1\}$$

- The potential outcomes in this panel data setting are:

$$\hookrightarrow \begin{cases} Y_{0it} = \mu_i + \delta_t + \varepsilon_{it} \\ Y_{1it} = \gamma_{it} + \mu_i + \delta_t + \varepsilon_{it} \end{cases} \Rightarrow \gamma_{it} = Y_{1it} - Y_{0it}$$

$$\Rightarrow \text{POM: } Y_{it} = D_{it} Y_{1it} + (1 - D_{it}) Y_{0it}$$

- For simplicity assume there are two periods: ( $t = 0$ , pre-treatment) and ( $t = 1$ , post-treatment).

$\hookrightarrow t \in \{0, 1\}$  and  $D_{it0} = 1$  for all  $i$

- The unit FE is the confounders and is related to the treatment, and assume changes in unobservable are mean independent of treatment status.

Assume  $\begin{cases} \text{Cov}(\mu_i, D_{it}) \neq 0 & \varepsilon_{it1} - \varepsilon_{it0} \\ E(\varepsilon_{it1} | D_{it}) = E(\varepsilon_{it0}) \text{ or } E(\Delta \varepsilon_{it} | D_{it}) = E(\Delta \varepsilon_{it0}) \end{cases}$

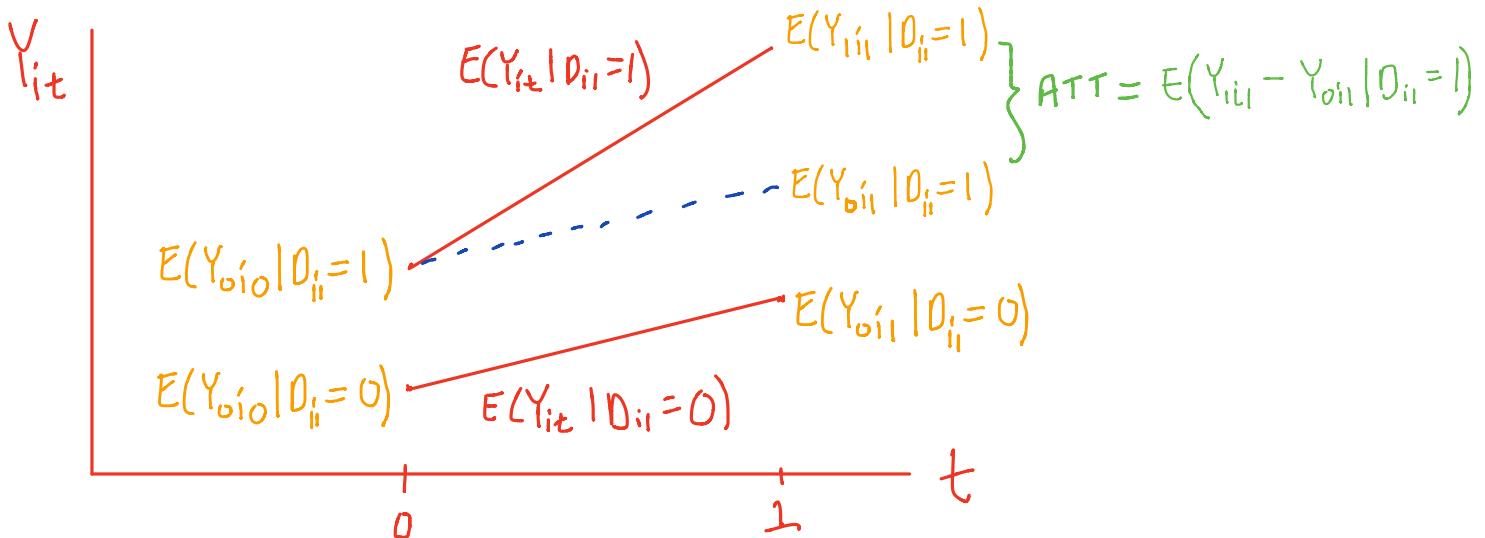
- Under the above assumptions, the ATT is identified using a diff-in-diff estimator:

$$\begin{array}{ll} \text{Post} & \left\{ \begin{array}{l} E(Y_{it1} | D_{it1}=1) = E(\gamma_{it1} | D_{it1}=1) + E(\mu_i | D_{it1}=1) + \delta_1 + E(\varepsilon_{it1} | D_{it1}=1) \\ E(Y_{it1} | D_{it1}=0) = E(\mu_i | D_{it1}=0) + \delta_1 + E(\varepsilon_{it1} | D_{it1}=0) \end{array} \right. \\ \text{Treat} & \left. \begin{array}{l} E(Y_{it0} | D_{it1}=1) = E(\mu_i | D_{it1}=1) + \delta_0 + E(\varepsilon_{it0} | D_{it1}=1) \\ E(Y_{it0} | D_{it1}=0) = E(\mu_i | D_{it1}=0) + \delta_0 + E(\varepsilon_{it0} | D_{it1}=0) \end{array} \right. \end{array}$$

$$\begin{array}{ll} \text{Pre} & \left\{ \begin{array}{l} E(Y_{it0} | D_{it1}=1) = E(\mu_i | D_{it1}=1) + \delta_0 + E(\varepsilon_{it0} | D_{it1}=1) \\ E(Y_{it0} | D_{it1}=0) = E(\mu_i | D_{it1}=0) + \delta_0 + E(\varepsilon_{it0} | D_{it1}=0) \end{array} \right. \\ \text{Treat} & \left. \begin{array}{l} E(Y_{it1} | D_{it1}=1) = E(\gamma_{it1} | D_{it1}=1) + E(\Delta \varepsilon_{it1} | D_{it1}=1) - E(\Delta \varepsilon_{it1} | D_{it1}=0) \\ E(Y_{it1} | D_{it1}=0) = E(\gamma_{it1} | D_{it1}=0) + E(\Delta \varepsilon_{it1} | D_{it1}=0) \end{array} \right. \end{array}$$

$$\begin{aligned} \Rightarrow \text{DD} &= E(\Delta Y_{it1} | D_{it1}=1) - E(\Delta Y_{it1} | D_{it1}=0) \\ &= E(\gamma_{it1} | D_{it1}=1) + E(\Delta \varepsilon_{it1} | D_{it1}=1) - E(\Delta \varepsilon_{it1} | D_{it1}=0) \\ &= E(Y_{it1} - Y_{it0} | D_{it1}=1) = \text{ATT} \end{aligned}$$

- The differences-in-differences intuition is illustrated in the figure below:



- We can use the above figure to define the common trends assumption:

Common trends:  $E(Y_{0i1} - Y_{0i0} | D_{i1} = 0) = E(Y_{0i1} - Y_{0i0} | D_{i1} = 1)$

↳ Implied by  $E(\Delta \varepsilon_{i1} | D_{i1}) = E(\Delta \varepsilon_{i0})$

- Under absence of the treatment, the average change in the outcome across the control and treatment group is the same over time.
- The common trends assumption is not invariant in nonlinear transformation of the outcome.
- That is, the identification of differences-in-differences is dependant on the functional form of the outcome. The DD identified in  $Y$  may not be identified in  $\log(Y)$ .
- Plausibility of common trends assumption can be evaluated using pre-treatment periods by checking for identical trends across control and treatment group.

### Synthetic Control

- Let  $j$  denote the unit, and  $t$  be the time period. Assume first unit is exposed to intervention at  $T_0$ , rest of the units are control group. Also suppose all units have  $k$  pre-treat characteristics.

$X_1$  = pre-treat charac. for treated unit ( $R \times 1$ )

$X_0$  = pre-treat charac. for control unit ( $R \times J$ )

- Synthetic control is essentially a weighted average of untreated units such that they match up with the treated unit prior to the treatment.

$$\min_w \|X_1 - X_0 w\| = \min_w \sqrt{(X_1 - X_0 w)' V (X_1 - X_0 w)}$$

↳  $V$  weights match  $X_1 \approx X_0 w$

- The ATT for  $T > T_0$  is estimated as follows (compare observed outcome to synthetic control):

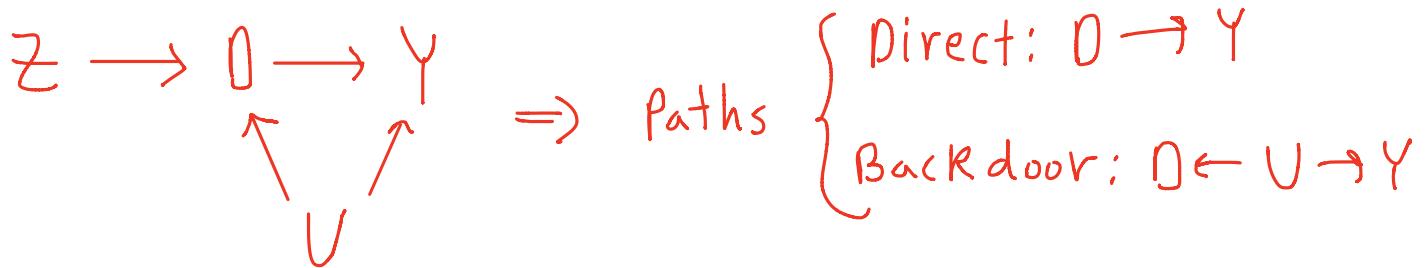
$$\hat{\alpha}_{1t} = Y_{1t} - \sum_{j=2}^{J+1} w_j^* Y_{jt} \quad \text{for } t \geq T_0$$

### Instrumental Variables

- IVs are best motivated for experiments without perfect compliance. Although the treatment assignment is randomized, treatment participation is still endogenous.
- Assignment effects by themselves are known as Intent-to-treat (ITT) effects:

$$\text{ITT} = E(Y_i | Z_i = 1) - E(Y_i | Z_i = 0)$$

- The ITT are relevant for interventions that are focused on evaluating the effect of providing units access to the treatment.
- The following DAG describes the standard IV setting:



- In the above DAG U is confounder that makes it difficult to measure the direct causal effect of the treatment on the outcome. Z is an instrument that effects the outcome only through the treatment.
- The IV uses exogenous variation in D that is induced by Z to estimate the direct effect of D on Y.
- The binary treatment can be expressed using a potential treatment model:

$$\hookrightarrow D_i = Z_i D_{oi} + (1 - Z_i) D_{ni}, \text{ where } (D_{oi}, D_{ni}) \text{ are potential treat.}$$

- One sided compliance is the control group cannot get the treatment, but not everyone in the treatment group participates in the assigned treatment.

$$\hookrightarrow \begin{cases} \text{Perf. Compliance: } \Pr(D_{ii} = 1) = \Pr(D_{oi} = 0) = 1 \\ \text{One sided: } \Pr(D_{oi} = 0) = 1 \text{ but } \Pr(D_{ii} = 1) < 1 \end{cases}$$

- In general, the population can be partitioned into the following groups:

$$\hookrightarrow \text{pop}^n \begin{cases} \text{Compliers: } D_{ii} > D_{oi} \quad (D_{ii} = 1 \Rightarrow D_{oi} = 0) \\ \text{Always-takers: } D_{ii} = D_{oi} = 1 \\ \text{Never-takers: } D_{ii} = D_{oi} = 0 \\ \text{Defiers: } D_{ii} < D_{oi} \quad (D_{ii} = 0 \Rightarrow D_{oi} = 1) \end{cases}$$

- The exclusion restriction says the potential outcome don't directly depend on the treatment assignment. The IV only effects the outcome through its impact on the treatment.
- The independence assumption is implied by the randomization of the IV. The exclusion and independence assumption implies:

$$\hookrightarrow \begin{cases} \text{Exclusion} \\ \text{Indep.} \end{cases} \Rightarrow (Y_{oi}, Y_{ii}, D_{ii}, D_{oi}) \perp \!\!\! \perp Z_i$$

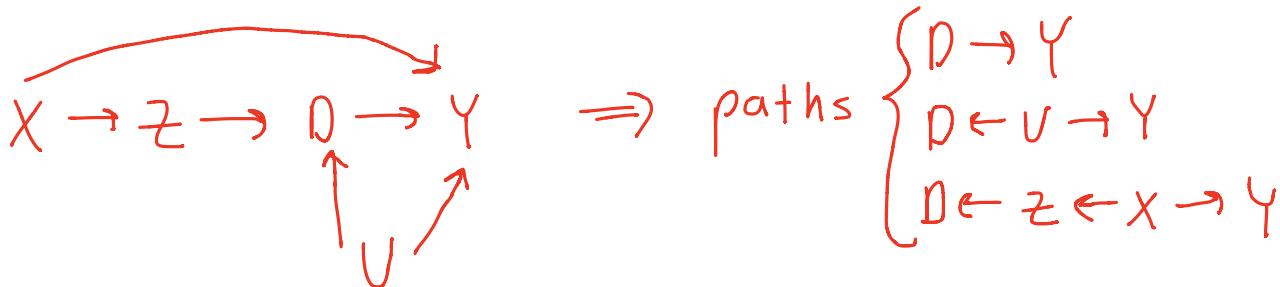
- The IV parameter is:

$$\hat{\tau}_{IV} = \frac{\text{Cov}(Y, Z)}{\text{Cov}(D, Z)} = \frac{E[(Y_{ii} - Y_{oi})(D_{ii} - D_{oi})]}{E(D_{ii} - D_{oi})}$$

- Under heterogeneous treatment effects and additionally imposing monotonicity (assignment to treatment makes unit more likely to be treated) identifies the average treatment effects for compliers.

Monotonicity:  $\Pr(D_{oi} \geq D_{0i}) \Rightarrow \text{LATE} = E(Y_{1i} - Y_{0i} | D_{1i} > D_{0i})$

- Under one sided non-compliance, the LATE is equivalent to the ATT.
- There may be observed variables that effect the instrument and outcome. This model is represented in the DAG below.



- The IV Z above is not valid unless we condition on X. The LATE after condition on X is:

$$\text{LATE} = \frac{E_x [E(Y_i | X_i, Z_i = 1) - E(Y_i | X_i, Z_i = 0)]}{E_x [\Pr(D_i = 1 | X_i, Z_i = 1) - \Pr(D_i = 1 | X_i, Z_i = 0)]}$$

### Regression Discontinuity Design

- The RD design relies on their being discontinuous change in the treatment after the value of the running variable crosses some threshold.

$$\begin{cases} X = \text{running variable} \\ C = \text{threshold} \end{cases} \Rightarrow D_i = I(X_i \geq C) = \begin{cases} 1, & \text{if } X_i \geq C \\ 0, & \text{if } X_i < C \end{cases}$$

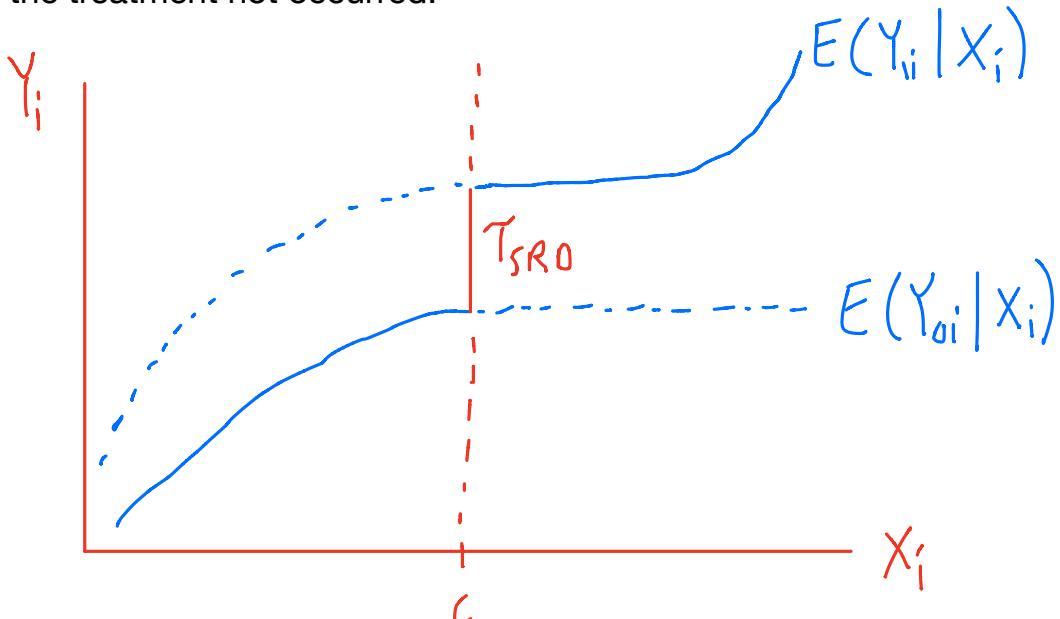
- There are three types of RD designs 1) Sharp RD, 2) Fuzzy RD, and 3) Kink RD. The first two are the most popular.
- The idea behind the RD design is that the treatment status is essential random around the threshold at which the treatment status changes.
- The sharp RD parameter of interest is:

$$\gamma_{SRD} = E(Y_{1i} - Y_{0i} | X_i = C) = \lim_{x \rightarrow C^-} E(Y_i | X_i = x) - \lim_{x \rightarrow C^+} E(Y_i | X_i = x)$$

- The above sharp RD parameter is nonparametrically identified if the conditional expectation of the potential outcomes are continuous at cutoff.

Assume  $\begin{cases} E(Y_{1i} | X_i = x) \\ E(Y_{0i} | X_i = x) \end{cases}$  are cont. at  $x = C$

- The requirement of the continuity assumption becomes clear in the figure below. It allows us to compare observed mean outcome post-treatment (just to the right of  $c$ ) to the counterfactual had the treatment not occurred:



- The key RD identifying assumptions can be supported by empirical tests. A common falsification test is to show there is no discontinuity in covariates at the cutoff.
- RD assumes there is no manipulation of individuals at the threshold (individuals can't control which side they end up near the cutoff). This assumption can be supported by showing there is no discontinuity in the distribution of the running variable at the cutoff.
- The most common estimation strategy is a weighted local-linear regression. This involves 1) using data points only near the cutoff and 2) regression that weights points closer to the cutoff.

$$\min_{\beta_0, \beta_1, \beta_2, \gamma} \sum_{i=1}^n \left( Y_i - \beta_0 - \beta_1 T_i - (X_i - c) \beta_2 - (X_i - c) D_i \beta_3 \right)^2 K\left(\frac{X_i - \bar{X}}{h}\right)$$

Subject to  $X_i \in [c-h, c+h]$ ,  $h = \text{bandwidth}$

## Conclusion

- Paper discusses the most common empirical tools for program evaluation: 1) Experiments, 2) Conditional Independence, 2) Instrumental Variables, 3) Differences-in-differences, and 4) Regression discontinuity.
- However paper is not exhaustive of all policy tools, see Athey and Imbens (2017) as a nice complement to this paper.