

李



Basic

SUTVA

Potential Outcomes. $\gamma_i(1)$. $\gamma_i(0)$

Causal effect = comparison between the potential outcomes under treatment and under control for the same unit or a common set of units

Causal estimands: functions of PO

$$ITE = \tau_i = \gamma_i(1) - \gamma_i(0)$$

(condition) CATE = $\tau(x) = E(\gamma_i(1) - \gamma_i(0) | X=x)$

$$ATE = \tau = E(\gamma_i(1) - \gamma_i(0)) = E_x(\tau(x))$$

(for treated units) ATT = $\tau = E(\gamma_i(1) - \gamma_i(0) | Z_i=1)$

(for control units) ATC = $\tau = E(\gamma_i(1) - \gamma_i(0) | Z_i=0)$

Ratios: $\pi_i = \frac{\gamma_i(1)}{\gamma_i(0)}$

$$\tau^{ATE} = \Pr(Z_i=1)\tau^{ATT} + \Pr(Z_i=0)\tau^{ATC}$$
$$\tau = \frac{E(\gamma_i(1))}{E(\gamma_i(0))}$$

In CRE: $\tau^{ATE} = \tau^{ATT} = \tau^{ATC}$ observe: *

Assumptions \rightarrow assignment mechanism

$$\rightarrow \Pr(Z_i=1 | X_i, Y_i(1), Y_i(0))$$

Unconfounded Assignment: propensity score

$$\Pr(Z_i=1 | X_i, Y_i(1), Y_i(0)) = \Pr(Z_i=1 | X_i)$$

$$Z_i \perp\!\!\!\perp Y_i(1), Y_i(0) | X_i$$

$$\begin{aligned} \textcircled{?} \quad \Pr(Y_i(z) | X_i) &= \Pr(Y_i(z) | X_i, Z_i) \\ &= \Pr(Y_i(z) | X_i, Z_i=2) \\ &= \Pr(Y_i(Z_i) | X_i, Z_i=2) \\ &= \Pr(Y_i^{\text{obs}} | X_i, Z_i=2) \end{aligned}$$

$$Y_i^{\text{obs}} = Y_i = Z_i Y_i(1) + (1-Z_i) Y_i(0) = Y_i(Z_i)$$

$$\text{e.g. } f(\cdot) = (\cdot)^2 + X$$

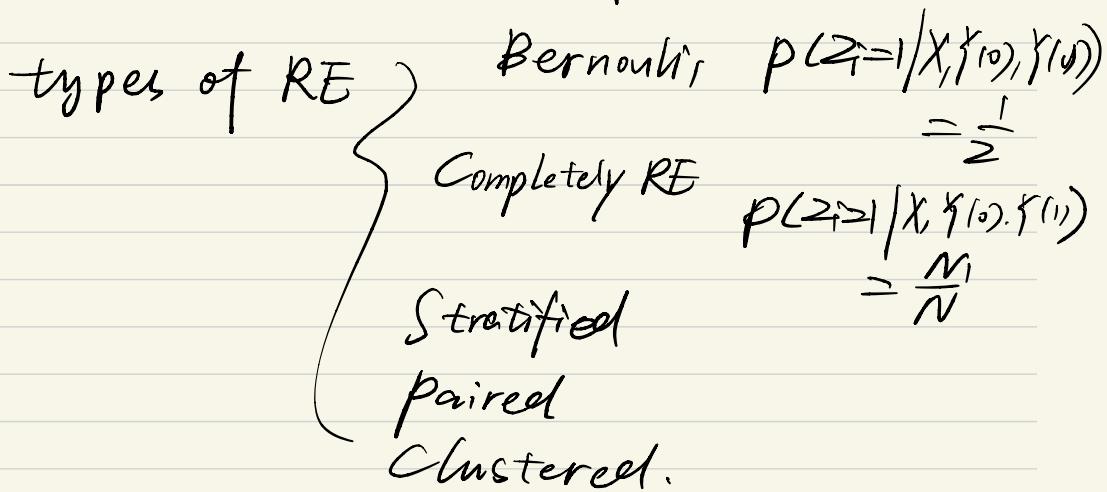
$$Y_i = f(Z_i, Y_i(1), Y_i(0))$$

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

Randomized experiments (RE)

$$Z \perp\!\!\! \perp X_1, Z \perp\!\!\! \perp U, Z \perp\!\!\! \perp Y(1), Y(0)$$

Under randomization. confounding ✓



(GRE)

Fisher's Randomized Test

assumptions: $\gamma^{(0)}, \gamma^{(1)}$ fixed

sharp null hypotheses

$$H_0: \gamma_i(1) = \gamma_i(0) \text{ for all units } i=1, \dots, n$$

test statistic

$$T = T(Z, \gamma) \quad \left. \begin{array}{l} \text{T Statistic} \\ \text{Wilcoxon Rank} \end{array} \right\}$$

exact p-value $\Pr(T \geq T^{\text{obs}})$

CRE

Neyman's Repeated Sampling Approach

estimand : ATE

1. find unbiased estimator of ATE
2. obtain variance and interval estimates of the unbiased estimator (derive the estimator distribution under repeated sampling based on the (randomized) distribution of Z)

estimator:

$$\text{DIM } \hat{\tau}^{\text{unadj}} = \sum_{i=1}^N \left(\frac{Z_i Y_i}{N_1} \right) - \sum_{i=1}^N \left(\frac{(1-Z_i) Y_i}{N_0} \right) = \bar{Y}_1 - \bar{Y}_0$$

unbiasedness $E(\hat{\tau}^{\text{unadj}} / \text{Sampling}) = T$

$$V(\hat{\tau}^{\text{unadj}}) = \frac{S_0^2}{N_0} + \frac{S_1^2}{N_1} - \frac{S_{01}^2}{N}$$

\checkmark Neyman

$$= \frac{S_0^2}{N_0} + \frac{S_1^2}{N_1}$$

$$\check{V}_{\text{Neyman}} = \frac{S_0^2}{N_0} + \frac{S_1^2}{N_1}$$

CRE

unbiasedness 被 RE 保证. covariance balance 是★
covariance balance 保证了 variance

Covariate Balance and Adjustment

a secondary analysis

Option 1. by design

- { stratified RE
- paired RE
- rerandomization

Option 2. by analysis

- { regression (ANCOVA)
- propensity score weighting
- model-based imputation

Observational Studies

In observational studies, some structural (often untestable) assumptions must be made.

Assumption 1. Unconfoundedness (untestable)

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

$$\Pr(Z_i \mid Y_{i(0)}, Y_{i(1)}, X_i) = \Pr(Z_i \mid X_i)$$

Assumption 2. Overlap (directly checked from the data)

$$0 < \Pr(Z_i=1 \mid X_i) < 1 \text{ for all } i$$

(Assumptions for estimating ATT : 1. $Z_i \perp\!\!\!\perp Y_{i(0)} \mid X_i$)

$$2. \Pr(Z_i=1 \mid X_i) < 1$$

Under Assumptions, we have

$$\Pr(Y_{i(z)} \mid X) = \Pr(Y \mid X, Z=z)$$

the distribution of $Y_{i(z)}$ = observed distribution of Y in
p.o. treatment arm $Z=z$

(DS)

so we have

$$\begin{aligned} T^{ATE} &= E_x \{ \tau(x) \} \\ &= E \{ M_1(x) - M_0(x) \} \\ &= E \{ E(Y | Z=1, X) - E(Y | Z=0, X) \} \end{aligned}$$

$$\begin{aligned} \hat{M}_z(x) &= E(Y(z) | X) \\ &= E(Y | Z=z, X) \end{aligned}$$

$$\begin{aligned} \text{ATE} &= \int \{ E(Y | Z=1, X=x) - E(Y | Z=0, X=x) \} f(x) dx \\ \text{离散} &= \sum_x \{ E(Y | Z=1, X=x) - E(Y | Z=0, X=x) \} p(x) \end{aligned}$$

$$\Pr(Y(z) | X) = \Pr(Y | Z=z, X) = f_z(x)$$

$$\begin{cases} \Pr(Y(z)) = \sum_x \Pr(Y(z) | X=x) \cdot \Pr(X=x) \\ \Pr(Y | Z=z) = \sum_x \Pr(Y | Z=z, X=x) \Pr(X=x | Z=z) \end{cases}$$

Two estimation strategies for τ^{ATE}

① Outcome modeling / regression

$$\hat{\tau}^{ATE} = \hat{\tau}^{reg} = \frac{1}{N} \sum_{i=1}^N \left\{ \hat{\mu}_1(x_i) - \hat{\mu}_0(x_i) \right\}$$

要估计 $\hat{\mu}_2(x)$

② Inverse probability weighting IPW

$$\hat{\tau}^{ATE} = \hat{\tau}^{IPW} = \frac{\sum_{i=1}^N z_i y_i / \hat{e}(x_i)}{\sum_{i=1}^N z_i / \hat{e}(x_i)} - \frac{\sum_{i=1}^N (1-z_i) y_i / \{1-\hat{e}(x_i)\}}{\sum_{i=1}^N (1-z_i) / \{1-\hat{e}(x_i)\}}$$

(高斯的估计见 BP Book. 10.4.1)

③ Matching

④ Stratification (高斯的 x)

Outcome Modeling | Regression

$$\hat{T}_{ATE} = \frac{1}{N} \sum_{i=1}^N \left\{ \hat{\mu}_1(x_i) - \hat{\mu}_0(x_i) \right\}$$

estimands

都是 $\mu_1(x)$

的函数

$$\hat{T}_{ATT} = \frac{1}{N_1} \sum_{i=1}^N z_i \left\{ y_i - \hat{\mu}_0(x_i) \right\}$$

$$\hat{T}_{ATC} = \frac{1}{N_0} \sum_{i=1}^N (1-z_i) \left\{ \hat{\mu}_1(x_i) - y_i \right\}$$

$$\hat{T}_{CATE}(x) = \hat{\mu}_1(x) - \hat{\mu}_0(x)$$

$$CATE = E(y_1 - y_0 | x=x)$$

估计 $\mu_2(x)$

To estimate $\mu_2(x)$, we have many choices of model.

① Linear regression, we assume

(注: 这种方法
不一定可行, 因

$$E(Y | Z=1, X) = \mu_1(x) = \beta_0 + \beta_2 + \beta_X x$$

为 Z 和 X 有
可能共线性)

$$E(Y | Z=0, X) = \mu_0(x) = \beta_0 + \beta_X x$$

$$E(Y | Z, X) = \beta_0 + \beta_2 \cdot Z + \beta_X^T X$$

$$Y \sim \beta_2 + X$$

$$\therefore T(X) = \mu_1(x) - \mu_0(x)$$

$$= \beta_0 + \beta_2 + \beta_X^T X - (\beta_0 + \beta_X^T X) = \beta_2$$

$$\therefore T_{ATE} = E\{T(X)\} = \beta_2$$

的系数 β_2

② We can also assume

$$E(Y|Z=1, X) = \mu_1(x) = \beta_0 + \beta_Z + \beta_X^T X + \beta_{ZX}^T X$$
$$= (\beta_0 + \beta_Z) + (\beta_X^T + \beta_{ZX}^T) X$$

$$\nwarrow E(Y|Z=0, X) = \mu_0(x) = \beta_0 + \beta_X^T X$$

$$E(Y|Z, X) = \beta_0 + \beta_Z Z + \beta_X^T X + \beta_{ZX}^T X Z$$

$$\therefore T(x) = \mu_1(x) - \mu_0(x)$$

$$= \beta_Z + \beta_{ZX}^T X$$

$$\therefore \bar{T}^{ATE} = E\{T(x)\} = \beta_Z + \beta_{ZX}^T E(X)$$

$$\therefore \bar{T}^{ATE} = \beta_Z + \beta_{ZX}^T \bar{X}$$

$$(或者无中为 X 为 \tilde{X}, Y \sim 1 + Z + \tilde{X} + Z\tilde{X})$$

$$\bar{T}^{ATE}$$

Why balance is important?

"If the imbalance of the covariates between the two groups is large, the model-based results heavily relies on extrapolation in the region with little overlap, which is sensitive to the model specification assumption."

Balance

$$\frac{\Pr(X|Z=1)}{\Pr(X|Z=0)}$$

metrics of balance:

ASD

$$= \frac{\frac{\Pr(Z=1|X) \Pr(X)}{\Pr(Z=1)}}{\frac{\Pr(Z=0|X) \Pr(X)}{\Pr(Z=0)}}$$

visualize balance:

Love plot

$$= \frac{\Pr(Z=1|X)}{\Pr(Z=0|X)} \cdot \frac{\Pr(Z=0)}{\Pr(Z=1)}$$

Overlap

Matching

often be applied in settings where

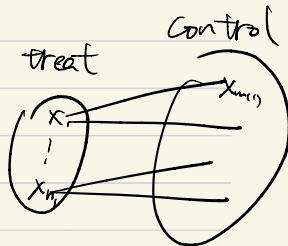
the interest is ATT

} a large reservoir of potential controls

Nearest-Neighbor Matching

$$\hat{y}_i(0) = \begin{cases} \sum_{j \in M_i} y_j / M, & z_j = 1 \\ y_i, & z_i = 0 \end{cases}$$

$$\hat{y}_i(1) = \begin{cases} y_i, & z_i = 1 \\ \sum_{j \in M_i} y_j / M, & z_j = 0 \end{cases}$$



$$\hat{T}^{ATE} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i(1) - \hat{y}_i(0))$$

$$\hat{T}^{ATT} = \frac{1}{N_1} \sum_{i=1}^N z_i (\hat{y}_i(1) - \hat{y}_i(0))$$

can be used as a pre-processing step

Stratification

Suppose we have a k-level covariate X

$$E(Y_{(1)}) = \sum_k E(Y|X=k, Z=1) \Pr(X=k)$$

$$E(\hat{Y}_{(1)}) = \sum_k \bar{Y}_{k,1} \frac{n_k}{n}$$

$$E(\hat{Y}_{(0)}) = \sum_k \bar{Y}_{k,0} \frac{n_k}{n}$$

$$\hat{\tau} = \sum_k (\bar{Y}_{k,1} - \bar{Y}_{k,0}) \frac{n_k}{n}$$

Propensity Score

a probability

mechanism

a summary statistic of assignment

$$e(x) = \Pr(Z=1 | X) = E(Z|X)$$

a summary score of the covariates

Property 1. $Z \perp X | e(x)$



Property 2. $Z_i \perp\!\!\!\perp \{Y_{i(0)}, Y_{i(1)}\} \mid X_i \Rightarrow Z_i \perp\!\!\!\perp \{Y_{i(0)}, Y_{i(1)}\} \mid e(X_i)$

Propensity score analysis

Stage 1. estimate the propensity score $\hat{e}(x)$

Stage 2. given the $\hat{e}(x)$, estimate the causal effects through stratification-weighting.

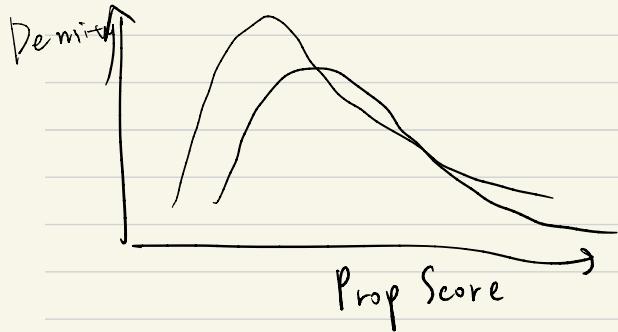
matching, regression and etc.

"ensure overlap and balance"

estimate propensity score

initial fit \rightarrow discard outliers \rightarrow check balance \rightarrow re-fit

$$\text{logit } \Pr(z_i=1 | X_i) = \beta X_i$$



Propensity Score Weighting

IPW Inverse Probability Weighting

$$E\left(\frac{z_i Y}{e(x)} - \frac{(1-z_i) Y}{1-e(x)}\right) = \tau_{ATE}$$

$$\text{proof: } E\left(\frac{z_i Y}{e(x)}\right) = E\left(E\left(\frac{z_i Y}{e(x)} | x\right)\right)$$

IPW creates a weighted population

$$= E\left(\frac{1}{e(x)} E(z_i Y | x)\right) \quad Y = z_i Y_{(1)} + (1-z_i) Y_{(0)}$$

$e(x) = E(z_i | x)$

$$= E\left(\frac{1}{e(x)} \cdot E\left(z_i^2 Y_{(1)} + z_i(1-z_i) Y_{(0)} | x\right)\right)$$

the covariates distributions of 2 groups are balanced

$$= E\left(\frac{1}{e(x)} E(z_i^2 | x) E(Y_{(1)} | x)\right) = E\left(E(Y_{(1)} | x)\right) = Y_{(1)}$$

$E\left(\frac{XZ}{e(x)}\right) = E\left(\frac{X(1-z_i)}{1-e(x)}\right)$

Define inverse probability weights

$$\{ w_1(x_i) = \frac{1}{e(x_i)}, \text{ for } z_i = 1 \}$$

$$\{ w_0(x_i) = \frac{1}{1-e(x_i)}, \text{ for } z_i = 0 \}$$

moment estimator

$$\hat{\tau}_{ipw,1} = \frac{1}{N} \left\{ \sum_{i=1}^N \frac{Y_i z_i}{e(x_i)} - \sum_{i=1}^N \frac{Y_i (1-z_i)}{1-e(x_i)} \right\}$$

$$= \frac{1}{N} \sum_{i=1}^N \left\{ Y_i z_i w_1(x_i) - Y_i (1-z_i) w_0(x_i) \right\}$$

Normalize weights (sum of weights within group should be 1)

$$\hat{\tau}_{ipw,2} = \frac{\sum_{i=1}^N Y_i z_i w_1(x_i)}{\sum_{i=1}^N z_i w_1(x_i)} - \frac{\sum_{i=1}^N Y_i (1-z_i) w_0(x_i)}{\sum_{i=1}^N (1-z_i) w_0(x_i)}$$

$$V(\hat{\tau}_{ipw,2}) < V(\hat{\tau}_{ipw,1})$$

Propensity Score Weighting

Extension Beyond IPW: Move the Goal post

use specified target populations.

"often units who are "on the fence" are of interest"

Balancing Weights

$f(x)$: density of the observed covariates

$g(x)$: density of a target population's covariates

$$h(x) = \frac{g(x)}{f(x)} \quad \text{a tilting function}$$

A new class of estimands: ATE over the target population

$$\begin{aligned} T_h &= E_g [\tau_{i(1)} - \tau_{i(0)}] = \frac{\int \tau(x) g(x) \mu(dx)}{\int g(x) \mu(dx)} \\ &= \frac{\int \tau(x) f(x) h(x) \mu(dx)}{\int f(x) h(x) \mu(dx)} = \frac{E\{h(x) \tau(x)\}}{E\{h(x)\}} \end{aligned}$$

choice of $h(x)$ determines the target population, estimand, weights

Let $f_z(x) = \Pr(X=x | Z=z)$ covariate distribution within each group z

$$f_1(x) = \Pr(X|Z=1) = \frac{\Pr(Z=1|x) \Pr(x)}{\Pr(Z=1)}$$

$$\propto e(x) f(x)$$

$$f_0(x) = \Pr(X|Z=0) = \frac{\Pr(Z=0|x) \Pr(x)}{\Pr(Z=0)}$$

$$\propto (1-e(x)) f(x)$$

weight $f_z(x)$ to the target population g

$$w_1(x) f_1(x) \propto g(x) = f(x) h(x)$$

$$w_0(x) f_0(x) \propto g(x) = f(x) h(x)$$

$$\Rightarrow w_1(x) \propto \frac{f(x) h(x)}{f_1(x)} \propto \frac{f(x) h(x)}{e(x) f(x)} = \frac{h(x)}{e(x)}$$

balancing weights

$$w_0(x) \propto \frac{f(x) h(x)}{f_0(x)} \propto \frac{f(x) h(x)}{(1-e(x)) f(x)} = \frac{h(x)}{1-e(x)}$$

Sample estimator of WATE

$$\hat{T}_n = \frac{\sum_i w_1(x_i) Z_i Y_i}{\sum_i w_1(x_i) Z_i} - \frac{\sum_i w_0(x_i) (1-Z_i) Y_i}{\sum_i w_0(x_i) (1-Z_i)}$$

"each unit is weighted by its probability of being assigned to the opposite group."

Overlap Weights

$$h(x) = e(x) (1 - e(x))$$

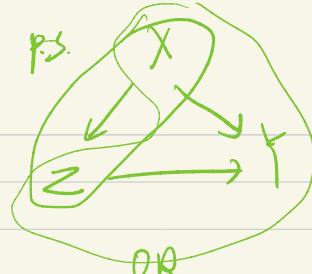
$$\begin{cases} w_1(x) \propto 1 - e(x), & \text{for } z=1 \\ w_0(x) \propto e(x), & \text{for } z=0 \end{cases}$$

$$h(x) = e(x) (1 - e(x))$$

$$= \left(\frac{1}{e_1} + \frac{1}{1-e_1} \right)^{-1} \quad \left(\frac{1}{e_1} + \frac{1}{e_2} + \frac{1}{e_3} \right)^{-1}$$

Doubly Robust Estimation

estimand: ATE



QR estimator:

$$E(Y|z) | X) = m_z(x)$$

$$\hat{T}_{\text{Doub}} = \frac{1}{N} \sum_{i=1}^N \{ z_i (Y_i - \hat{m}_0(x_i)) + (1-z_i) (\hat{m}_1(x_i) - Y_i) \}$$

① can be sensitive when overlap is weak

PW Weighting estimator:

$$\hat{T}_{\text{PW}} = \frac{\sum Y_i z_i / e(x_i)}{\sum z_i / e(x_i)} - \frac{\sum Y_i (1-z_i) / (1-e(x_i))}{\sum (1-z_i) / (1-e(x_i))}$$

DR estimator:

$$\begin{aligned} T &= E(Y_{(1)}) - E(Y_{(0)}) = E \left\{ \frac{2Y}{e(x)} - \frac{2 - e(x)}{e(x)} m_0(x) \right\} - \\ &\quad E \left\{ \frac{(1-z)Y}{1 - e(x)} + \frac{2 - e(x)}{1 - e(x)} m_1(x) \right\} \\ &= E \left\{ m_1(x) + \frac{z_i (Y_i - m_1(x_i))}{e(x_i)} \right\} \end{aligned}$$

用 \hat{m}_0 , \hat{m}_1 代替插式

$$- E \left\{ m_0(x_i) + \frac{(1-z_i) (Y_i - m_0(x_i))}{1 - e(x_i)} \right\}$$

$$\begin{aligned} \hat{T}_{\text{DR}} &= \frac{1}{N} \sum_{i=1}^N \left\{ \frac{z_i Y_i}{\hat{e}(x_i)} - \frac{z_i - \hat{e}(x_i)}{\hat{e}(x_i)} \hat{m}_1(x_i) \right\} \\ &\quad - \frac{1}{N} \sum_{i=1}^N \left\{ \frac{(1-z_i) Y_i}{1 - \hat{e}(x_i)} + \frac{z_i - \hat{e}(x_i)}{1 - \hat{e}(x_i)} \hat{m}_0(x_i) \right\} \\ &= \dots \end{aligned}$$

$$\text{let's see } \hat{M}_{1,\text{dr}} = \frac{1}{N} \sum_{i=1}^N \left\{ \frac{\hat{Y}_i - \hat{Y}_i}{\hat{e}_i} - \frac{(Z_i - \hat{e}_i) \hat{m}_i}{\hat{e}_i} \right\}$$

$$\hat{M}_{1,\text{dr}} - E(Y_{(1)}) = E \left\{ \frac{\geq(Y - m_1(x))}{e(x)} + m_1(x) \right\} - E(Y_{(1)})$$

$m_1(x) = E(Y_{(0)}|x)$
 $e(x) = E(z=1|x)$

$$= E \left\{ \frac{\geq(Z Y_{(1)} + (1-Z) Y_{(0)} - m_1(x))}{e(x)} + m_1(x) - Y_{(1)} \right\}$$

$$= E \left\{ \frac{\geq(Y_{(1)} - m_1(x))}{e(x)} - (Y_{(1)} - m_1(x)) \right\}$$

$$= E \left\{ \frac{Z - e(x)}{e(x)} (Y_{(1)} - m_1(x)) \right\} \quad \checkmark$$

$$= E \left(E \left\{ \frac{Z - e(x)}{e(x)} (Y_{(1)} - m_1(x)) \mid X \right\} \right) \quad \xrightarrow{\text{Z} \perp \text{Y}_{(1)} \text{Y}_{(0)} | X}$$

$$= E \left(E \left\{ \frac{Z - e(x)}{e(x)} \mid X \right\} \times E \left\{ Y_{(1)} - m_1(x) \mid X \right\} \right)$$

$$= E \left(\underbrace{\frac{E(Z|x)}{e(x)}}_{\text{true}(x)} \times \underbrace{\left(E(Y_{(1)}|x) - m_1(x) \right)}_{m_1,\text{true}(x)} \right)$$

so $\hat{M}_{1,\text{dr}} - E(Y_{(1)}) = 0$ if either $E(Z|x) = e(x)$ or

(asymptotically unbiased) $E(Y_{(1)}|x) = m_1(x)$

\hat{T}_{dr} is a consistent estimator of ATE if either
model is correctly specified. (large sample property)

Overlap and Balance

$$Y_{(0)}, Y_{(1)} \perp Z \mid X=x$$

$$f^{(0)}(x) = E(Y \mid z=0, x=x)$$

$$= E(Z Y_{(0)} + (1-Z) Y_{(1)} \mid z=0, x=x)$$

$$= E(Y_{(0)} \mid z=0, x=x) = E(Y_{(0)} \mid x=x)$$

Treatment Effect Heterogeneity

Machine Learning Approaches

$$Y = f(z, x) + \epsilon$$

$$f(z, x) = E(Y | z=z, x=x)$$

\leftarrow model this

$$\text{CATE } \tau(x) = f(1, x) - f(0, x)$$

S-Learner (single)

use all data to fit model $f(z, x)$

$$\text{e.g. BART } f(z, x) = \sum_{t=1}^T g_t(x, z; J_t, M_t)$$

T-Learner (two)

fit separate models to the treated/control groups

$$E(Y | z=1, x=x) = f_1(x)$$

$$E(Y | z=0, x=x) = f_0(x)$$

$$\text{CATE } \tau(x) = f_1(x) - f_0(x).$$

R-Learner

$$Y_i = \mu(x_i) + \tau(x_i) Z_i + \epsilon_i \quad f(z, x) = \mu(x) + \tau(x) \cdot z$$

$$\mu(x_i) = E(Y_i | X_i) = \mu(x_i) + \tau(x_i) e(x_i)$$

$$Y_i - \mu(x_i) = (Z_i - e(x_i)) \tau(x_i) + \epsilon_i$$

$$\tau(\cdot) = \arg \min_{\tau} \mathbb{E} \left[(Y_i - \mu(x_i)) - (Z_i - e(x_i)) \tau(x_i) \right]$$

Causal Forests

*γ -learner
for CATE, estimate $f_{\gamma}(x)$*

→ split to maximize heterogeneity of the treatment effect

$$\frac{|\bar{z}_c - \bar{z}_r|}{\sqrt{\text{Var}(\bar{z}_c) + \text{Var}(\bar{z}_r)}}$$

hope: withinleaf K.
as if randomized
similar covariates

leaves = $L_1(x), \dots, L_K(x)$

In leaf K : estimated treatment effect

$$\hat{t}_K = \frac{1}{\sum_{\{i: x_i \in L_K(x)\}} 1} \sum_{\{i: x_i \in L_K(x)\}} y_i - \frac{1}{\sum_{\{i: x_i \notin L_K(x)\}} 1} \sum_{\{i: x_i \notin L_K(x)\}} y_i$$

→ tuning cross-validation

$$\text{MSE } L(\hat{t}) = E[(y_{i(1)} - y_{i(0)} - \hat{t}(x_i))^2]$$

unknown

Honesty criterion: divide samples into
three subsamples (train, estimate causal
effect, test)

b-learner for CATE. estimate $f(z, x)$

BART (Bayesian Additive Regression Trees)

partition the space of x and z

$$f(z, x) = \sum_{t=1}^T g(x, z; J_t, M_t)$$

J_t represents the tree structure

M_t parameters for predictions in each leaf

$$M_t = (M_{t1}, \dots, M_{tL_t}) \quad L_t: \text{number of leaves}$$

Bayesian approach

prior probability of splitting decreases with depth

shrinkage of mean parameters $M_t \sim N(0, \frac{\sigma^2}{T})$

Do we expect CATE to be as complex as $f(z, x)$?

→ re-parameterize

$$f(z, x) = \mu(x) + \tau(x)z$$

further we can include p.s. score

$$f(z, x) = \mu(x, \hat{e}(x)) + \tau(x)z$$

use ML for both outcome model and propensity model

Double Machine Learning (DML) for high-dimensional data

→ An early example : use ML in DR estimator

→ A canonical example : Partial Linear Model.

intuition: relation between Y and X is usually more complex than the relation between Y and Z

idea: use ML model for $Y \sim X$ and linear for $Y \sim Z$

$$\text{DML for ATE} \quad Y = Z_I + m(X) + V \quad E(V|X, Z) = 0$$

$$Z = e(X) + V \quad E(V|X) = 0$$

- steps:
1. split sample into K parts $\{J_k\}_{k=1}^K \cdot J_k \cdot j_k$
 2. estimate Z from X using J_k^c
 3. estimate Y from X using J_k^c
 4. estimate ATE in J_k . regress residuals.
 5. aggregate over K folds . repeat 2-3 for $k=1 \dots K$

DML for CATE

$$Y = Z_I(X) + m(X) + V \quad E(V|X, Z) = 0$$

$$Z = e(X) + V \quad E(V|X) = 0$$

$$E(UV|X, Z) = 0$$

$$\hat{\tau}(x) = \underset{I_E}{\operatorname{argmin}} E_n[(\hat{Y} - \tau(x)\hat{V})^2]$$

Sensitivity Analysis

perform sensitivity

analysis to assess how

sensitive the causal
analysis is if unconfound-

two key overlap (testable)

assumption unconfoundedness (untestable) is violated

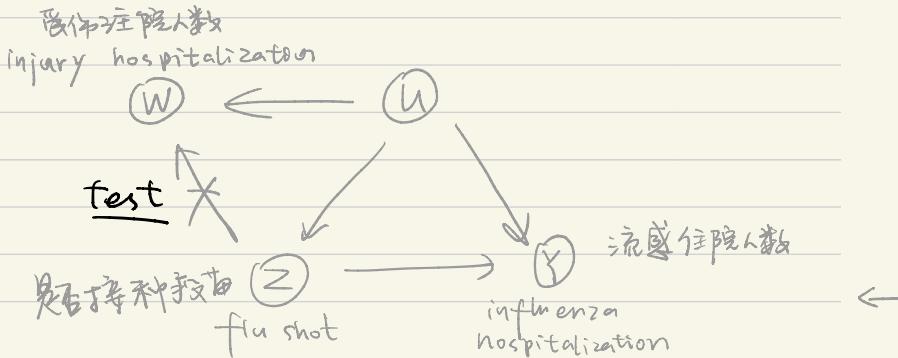
Unconfoundedness and Balance

$$\{Y_{10}, Y_{10}\} \perp\!\!\!\perp Z | X$$

implies: within strata of observed covariates,
potential outcomes corresponding to both treatment
conditions would be balanced between groups

e.g. $P(Y_{10} | Z=1)$ vs. $P(Y_{10} | Z=0)$

In practice, we use balance in covariates
 $P(X | Z=1)$ $P(X | Z=0)$ as a proxy to balance in potential outcomes.



Assess unconfoundedness indirectly

Falsification tests

Method 1. use multiple control groups

indicator $T_i \in \{1, 0, 1\}$

$$\text{suppose } Y_i = \begin{cases} Y_{i(0)} & \text{if } T_i \in \{1, 0\} \\ Y_{i(1)} & \text{if } T_i = 1 \end{cases}$$

using data from alternative control group

extend the unconfoundedness to $\{Y_{i(0)}, Y_{i(1)}\} \perp T_i \mid X_i$

testable $Y_{i(0)} \perp \{T_i = 0\} \mid X_i, T_i \in \{1, 0\}$

$Y_i^{\text{obs}} \perp \{T_i = 0\} \mid X_i, T_i \in \{1, 0\}$

using data from prior periods

Method 2. using lagged outcomes
(observed before the treatment)

Y_{lag} as a proxy for $Y^{(0)}$ $V = (X \mid Y_{\text{lag}})$

test $H_0: E(Y_{\text{lag}, 2=1} - Y_{\text{lag}, 2=0} \mid V=v) = 0 \forall v$

special case

vs. $H_1: \exists v. \text{s.t. } E(Y_{\text{lag}, 2=1} - Y_{\text{lag}, 2=0} \mid V=v) \neq 0$

using alternative placebo outcome supposed not

Method 3. using negative control outcomes affected by treatment (NCO)

Select another outcome variable known not be causally affected by the treatment of interest

$W = Z W_{(1)} + (1-Z) W_{(0)}$ unconfoundedness implies $\{W_{(1)}, W_{(0)}\} \perp Z \mid X$

$\rightarrow E[W_{(1)} - W_{(0)} \mid X] = 0$ e.g. flu shot (Z). influenza hospital (X). injury hospitalization (W)

sensitivity analysis
in causal

assess the bias of causal
effect estimates when the
unconfoundedness is assumed to fail

Assumption $p(z | Y_{(0)}, Y_{(1)}, X) \neq p(z | X)$
 $\exists H \{Y_{(1)}, Y_{(0)}\} | X$

$p(z | Y_{(0)}, Y_{(1)}, X, V) = p(z | X, V)$
 $\exists H \{Y_{(1)}, Y_{(0)}\} | (X, V)$

Method 1. Rubin

$(\pi, \alpha, \beta_1, \beta_0) \left\{ \begin{array}{l} V \sim \text{Bern}(IV) \\ \text{logit}(p(z=1|V)) = \gamma + \alpha V \\ \text{logit}(p(Y(z)=1|V)) = \beta_z + \beta_2 V \end{array} \right.$

$\begin{matrix} \swarrow & \searrow \\ V & \\ z \rightarrow & P \mid X \end{matrix}$

Method 2. Rosenbaum's bounds (only for testing no causal effect)

$$\frac{1}{T} \leq \frac{\Pr(z=1 | V=1) / \Pr(z=0 | V=1)}{\Pr(z=1 | V=0) / \Pr(z=0 | V=0)} \leq T.$$

Method 3. E-value

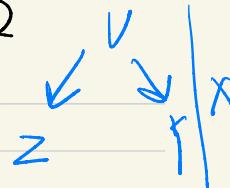
$$Y_{(1)} = Y, \quad z=1$$

$$Y_{(0)} = Y, \quad z=0.$$

BP's work (assume binary outcomes)

observed relative risk

$$\begin{aligned} RR_{z \neq 1|x}^{\text{obs}} &= \frac{p(Y=1 | Z=1, x)}{p(Y=1 | Z=0, x)} \\ &= \frac{\sum_u p(Y=1 | Z=1, x, u) p(u | Z=1, x)}{\sum_u p(Y=1 | Z=0, x, u) p(u | Z=0, x)} \end{aligned}$$



Causal RR

$$\begin{aligned} RR_{z \neq 1|x}^{\text{true}} &= \frac{p(Y_{(u)}=1 | x)}{p(Y_{(0)}=1 | x)} \\ &= \frac{\sum_u p(Y=1 | Z=1, x, u) p(u | x)}{\sum_u p(Y=1 | Z=0, x, u) p(u | x)} \end{aligned}$$

② why not $V \rightarrow Z$

treatment-confounder association ($Z \rightarrow u$)

PR of level $V=u$ (not necessarily binary)

$$RR_{zu|x}^{(u)} = \frac{p(u | Z=1, x)}{p(u | Z=0, x)}$$

measures of confounding

$$RR_{zu|x} = \max_u RR_{zu|x}^{(u)}$$

confounder-outcome association ($u \rightarrow Y$)

$$RR_{uy(z)|x} = \frac{\max_u p(Y(z)=1 | x, u)}{\min_u p(Y(z)=0 | x, u)}$$

$$RR_{uy|x} = \max \{ RR_{uy(0)|x}, RR_{uy(1)|x} \}$$

$$RR_{Z|X}^{obj} = \frac{\Pr(Y=1 | Z=1, X=x)}{\Pr(Y=1 | Z=0, X=x)}$$

其中 $\Pr(Y=1 | Z=1, X=x)$

$$= \Pr(Y=1 | Z=1, X=x, V=1) \Pr(V=1 | Z=1, X=x) +$$

$$\Pr(Y=1 | Z=1, X=x, V=0) \Pr(V=0 | Z=1, X=x)$$

$$= \Pr(Y=1 | X=x, V=1) \Pr(V=1 | Z=1, X=x) +$$

$$\Pr(Y=1 | X=x, V=0) \Pr(V=0 | Z=1, X=x)$$

$$\text{let } \Pr(V=1 | Z=1, X=x) = f_1, x$$

$$\Pr(V=0 | Z=1, X=x) = f_0, x$$

$$= \Pr(Y=1 | X=x, V=1) \cdot f_1, x + (1-f_1, x) \cdot \Pr(Y=1 | X=x, V=0)$$

$$= p_1(Y=1 | V=0, X=x) \left(RR_{V|X} \cdot f_1, x + (1-f_1, x) \right)$$

$$\text{同理 } 1-p_1 = \Pr(Y=1 | V=0, X=x) \left(RR_{V|X} \cdot f_0, x + (1-f_0, x) \right)$$

Drop $X = \infty$ for simplicity (conduct analysis)

within strata of covariates or propensity scores

Goal of sensitivity analysis:

recover RR_{2Y}^{true} from RR_{2Y}^{obs} and (RR_{2U}, RR_{UY})

$$\rightarrow RR_{2Y}^{\text{true}} \geq RR_{2Y}^{\text{obs}} / \frac{RR_{2U} \times RR_{UY}}{RR_{2U} + RR_{UY} - 1}$$

stronger $\rightarrow \max(RR_{2U}, RR_{UY})$

$$\geq \left\{ RR_{2Y}^{\text{obs}} + \sqrt{RR_{2Y}^{\text{obs}}(RR_{2Y}^{\text{obs}} - RR_{2Y}^{\text{true}})} \right\} / RR_{2Y}^{\text{true}}$$

to explain away the observed RR ($RR_{2Y}^{\text{true}} = 1$)

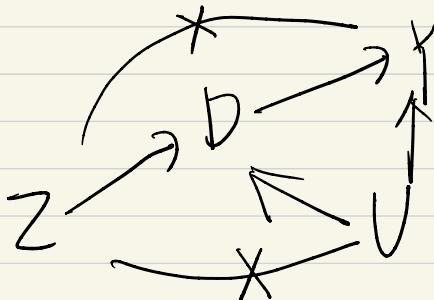
$$\rightarrow \max(RR_{2U}, RR_{UY}) \geq RR + \underbrace{\sqrt{RR(RR-1)}}_{E\text{-value}}$$

Instrument Variables

Handle unmeasured
confounding

Main idea

1. find an IV that influences treatment assignment but is independent of unmeasured confounders and has no direct effect on the outcome (except through its effect on treatment).
(natural experiment random)
2. use this variable to extract variation in the treatment that is free of the unmeasured confounders
3. use this confounder-free variation in the treatment to estimate the causal effect of the treatment



$$\hat{\beta}_{1,IV} = \frac{\widehat{\text{cov}}(Y_i, Z_i)}{\widehat{\text{cov}}(D_i, Z_i)} = \frac{\frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})(Z_i - \bar{Z})}{\frac{1}{N} \sum_{i=1}^N (D_i - \bar{D})(Z_i - \bar{Z})}$$

When Z binary (Wald estimator)

$$\hat{\beta}_{1,IV} = \frac{\bar{Y}_1 - \bar{Y}_0}{\bar{D}_1 - \bar{D}_0}$$

Interpretation I (indirect least squares)

$$Y_i = \pi_{10} + \pi_{11} \cdot Z_i + \varepsilon_{1i} (+ \pi_{12}' X_i)$$

$$D_i = \pi_{20} + \pi_{21} \cdot Z_i + \varepsilon_{2i} (+ \pi_{22}' X_i)$$

$$\hat{\beta}_{1,IV} = \frac{\hat{\pi}_{11}}{\hat{\pi}_{21}}$$

Interpretation II (2SLS)

$$\text{Stage 1. } \hat{D}_i = \hat{\pi}_{20} + \hat{\pi}_{21} \cdot Z_i (+ \hat{\pi}_{22}' X_i)$$

$$\text{Stage 2. } Y_i = \beta_0 + \beta_1 \hat{D}_i + \eta_i (+ \beta_2' X_i)$$

$$\hat{\beta}_{1,2SLS}$$

Discuss 2V using potential outcome notation

2V

focusing on the case of binary treatment and binary

IV $Z=0,1$; treatment $D=0,1$; outcome Y ; covariates X

Potential outcomes: $D(z) \cdot Y(z, d)$

Motivating context: randomized experiments with noncompliance ($Z_i \neq D_i$ for some units)

ITT approach:

$$T^{ITT} = E[Y_{i(1)} - Y_{i(0)}]$$

$$\hat{T}^{ITT} = \sum Y_i Z_i / \sum Z_i - \sum Y_i (1-Z_i) / \sum (1-Z_i)$$

estimates the causal effect of the assignment
on outcome, but not the effect of treatment received

Effectiveness: "Does the vaccine help people?"

Efficacy: "Does the vaccine work?"

2V Approach to Noncompliance

D is a post-assignment. $b(z)$, $z=0,1$

$\gamma(z)$ for $z=0,1$ (omit the double index $\gamma(z_{rd})$ for simplicity)

Observed data: Z_i , $D_i = b(Z_i)$. $\gamma_i = \gamma(Z_i)$

Central idea: (i) random assignment Z is an I^V
 (ii) divide units into latent subgroups based on compliance behavior

Compliance type: $S_i = (D_{i(0)}, D_{i(1)})$

		$D_{i(0)}$	
	0	0	1
$D_{i(1)}$	0	NT	b
	1	c	AT

The observed cells of 2.D are mixture

Z	D	S
0	0	[C, NT]
0	1	[AT, b]
1	0	[NT, b]
1	1	[c, AT]

Define ΔTT for each compliance type ($s=c, n, a, d$)

$$T_S^{ITT} = E[\gamma_i(1) - \gamma_i(0) \mid S_i = s]$$

Global ΔTT

$$T_Y^{ITT} = \pi_c T_C^{ITT} + \pi_n T_n^{ITT} + \pi_a T_a^{ITT} + \pi_d T_d^{ITT}$$

Assumptions

equivalent to DV
estimand when the
random assignment be
viewed as DV

A1. SVTR

A2. $\text{cov}(Z_i, D_i) > 0$

A3. IV is randomized

A4. ER

$$(T_n^{ITT} = T_a^{ITT} = 0)$$

A5. monotonicity

$$(\pi_d = 0)$$

$$\therefore T_Y^{ITT} = \pi_c T_C^{ITT}$$

(global ΔTT may be viewed as a conservative estimate of treatment effect
 $T^{ITT} = T_c^{ITT}$)

Complier Average Causal Effect (CACE)

$$T^{CACE} = T_C^{ITT} = E[\gamma_i(1) - \gamma_i(0) \mid S_i = c]$$

$$= \frac{T_Y^{ITT}}{\pi_c} = \frac{T_Y^{ITT}}{T_D^{ITT}} = \frac{E(Y|Z=1) - E(Y|Z=0)}{E(D|Z=1) - E(D|Z=0)}$$

Principal Stratification

Post-treatment confounding

$$Z \rightarrow D \rightarrow Y$$

Examples:

Non compliance in randomised experiments

$$\tau^{CACE} = E[Y_{i(1)} - Y_{i(0)} \mid D_i(0) = 0, D_i(1) = 1]$$

Selection Bias in cluster randomised trials

$$\tau^{\text{overall}} = E[Y_{i(1)} - Y_{i(0)}]$$

$$\tau^{\text{recruited}} = E[Y_{i(1)} - Y_{i(0)} \mid D_i = 1]$$

Censoring by death

$$\tau^{SACE} = E[Y_{i(1)} - Y_{i(0)} \mid D_i(0) = 1, D_i(1) = 1]$$

Survival Average Causal Effect

The individual principle strata memberships are + observed \rightarrow need additional assumption for identifying ps.

Assumption 1:

$$\{Y_{i(0)}, Y_{i(1)}, D_{i(0)}, D_{i(1)}\} \perp Z_i | X_i$$

The observed (Z_i) consist of mixtures of principal strata.

Estimation of PS inherently involves latent mixture models

Estimation	}	moment-based
mixture model approach		

$$\prod_i \Pr(Y_{i(0)}, Y_{i(1)}, D_{i(0)}, D_{i(1)}, Z_i, X_i; \theta) \quad (?)$$

$$= \prod_i \Pr(Z_i | Y_{i(0)}, Y_{i(1)}, S_i, X_i; \theta) \Pr(Y_{i(0)}, Y_{i(1)} | S_i, X_i; \theta)$$

$$\propto \prod_i \Pr(Y_{i(0)} | S_i, X_i; \theta)^{z_i} \Pr(Y_{i(1)} | S_i, X_i; \theta)^{1-z_i} \Pr(S_i | X_i; \theta)$$

Regression Discontinuity Designs

RDD. quasi-experimental design

The treatment status changes discontinuously according to some underlying pre-treatment variable (so-called forcing / running variable)

Basic idea: in these studies, comparing units with similar values of this variable, but different levels of treatment would lead to causal effect of the treatment at the threshold.

Formulation

$$\Pr(Z=1 | S_0^+) \neq \Pr(Z=1 | S_0^-)$$

$$\left\{ \begin{array}{l} \text{Sharp RD } \Pr(Z=1 | S_0^+) = 1 \quad \Pr(Z=1 | S_0^-) = 0 \\ \text{Fuzzy RD } \Pr(Z=1 | S_0^+) > \Pr(Z=1 | S_0^-) \end{array} \right.$$

Two frameworks to prove { continuity-based
local randomization

Sharp RD continuity approach

S_i = the running variable ; s_0 is the cutoff

Z_i : binary treatment indicator $Z_i = I\{S_i \leq s_0\}$

Assumption 1 (continuity of conditional Regression functions)

$E[Y_{(0)} | S=s]$ and $E[Y_{(1)} | S=s]$ are continuous in s

$$\Rightarrow E[Y_{(0)} | S=s_0] = \lim_{s \rightarrow s_0^-} E[Y_{(0)} | S=s]$$

$$= \lim_{s \rightarrow s_0^-} E[Y_{(0)} | Z=0, S=s]$$

$$= \lim_{s \rightarrow s_0^-} E[Y | S=s]$$

$$E[Y_{(1)} | S=s_0] = \lim_{s \rightarrow s_0^+} E[Y | S=s]$$

Estimand:

$$T_{RD}^{true} = \lim_{s \rightarrow s_0^-} E[Y | S=s] - \lim_{s \rightarrow s_0^+} E[Y | S=s]$$

the causal effect at the threshold

Estimation:

$$\hat{T}_{RD}^{true} = \hat{f}_+(s_0) - \hat{f}_-(s_0)$$

→ Sharp RD Remarks

Sharp RD completely violates the overlap assumption

assumption: there is zero overlap in the treatment probability at the threshold

continuity can be indirectly tested. e.g. distributions of all pre-treatment variables should be very similar on either side of the cutoff

→ Bandwidth Selection

In local linear regression, choice of h boils down to a bias-variance tradeoff.

smaller h : closer to threshold so less bias.

smaller sample size so more variance

Fuzzy RD continuity approach

S_i = the running variable ; s_0 is the cutoff

Z_i : binary assignment indicator $Z_i = I\{S_i \leq s_0\}$

D_i : treatment received

Estimation :

$$\hat{T}_{FPRD} = \frac{\lim_{s \rightarrow s_0^-} E[Y|s=s] - \lim_{s \rightarrow s_0^+} E[Y|s=s]}{\lim_{s \rightarrow s_0^-} E[D|s=s] - \lim_{s \rightarrow s_0^+} E[D|s=s]}$$

Assumption 1 (continuity of conditional Regression functions)

$$E[Y(0)|s=s_0] = \lim_{s \rightarrow s_0^-} E[Y|s=s]$$

$$E[Y(1)|s=s_0] = \lim_{s \rightarrow s_0^+} E[Y|s=s]$$

$$\left. \begin{array}{l} \\ \end{array} \right\} \begin{array}{l} T_{FPRD} \\ = T_{EXCE}(s_0) \end{array}$$

Assumption 2 (monotonicity)

$D_i(s)$ is monotone in s at $s=s_0$

Estimation :

$$\hat{T}_{FPRD} = \frac{\hat{f}_+(s_0) - \hat{f}_-(s_0)}{\hat{g}_+(s_0) - \hat{g}_-(s_0)}$$

Sharp RD Local Randomization Approach

the running variable as a random variable

Assumption 1. Local overlap : There exists a subset of units, V_{S_0} , in the random sample (or population) in the study such that for each $i \in V_{S_0}$, $\Pr(S_i \leq S_0) > \varepsilon$ and $\Pr(S_i > S_0) > \varepsilon$ for some sufficient large $\varepsilon > 0$.

Assumption 2. Local SRD-SVTVA : For each $i \in V_{S_0}$,

$D_i' = \mathbb{I}\{S_i' \leq S_0\}$ $D_i'' = \mathbb{I}\{S_i'' \leq S_0\}$ with possibly $S_i' \neq S_i''$

if $D_i' = D_i''$, that is, if either $S_i' \leq S_0$ and $S_i'' \leq S_0$

or $S_i' > S_0$ and $S_i'' > S_0$, then $\gamma_i(D') = \gamma_i(D'')$

Estimand :-

$$\tau_{S_0} = E[\gamma_i(y) - \gamma_i(b)] \mid i \in V_{S_0}$$

Assumption 3. Local randomization : within V_{S_0}

$$\Pr(S_i \mid \gamma_i(y), \gamma_i(b), X_i) = \Pr(S_i)$$

$$\tau_{S_0} = E[\gamma_i \mid b_i = 1, i \in V_{S_0}] - E[\gamma_i \mid b_i = 0, i \in V_{S_0}]$$