

Econometrics Notes

Jiahui Shui

February 18, 2025

Contents

1	Probability and Statistics Review	2
1.1	Probability Space	2
1.2	Conditional Probability and Independence	2
1.3	Random Variable	2
1.4	Convergence	2
1.5	Law of Large Numbers	4
1.6	Central Limit Theorem	4
1.7	Normal Distribution	4
1.8	Hypothesis Testing	4
1.9	Wald's Test	5
1.10	Introduction to Bayesian Statistics	5
1.11	Some Exercise	5
2	Causal Inference and Potential Outcomes	5
2.1	Potential Outcomes Model	5
2.2	Randomized Controlled Trials	6
2.3	Selection Bias	6
3	Linear Regression	7
3.1	Algebraic Properties	7
3.2	Large Sample Properties	8
3.3	Standard Error	8
3.4	Hypothesis Testing	9
4	Selection on Observables	9
4.1	Overlap	9
4.2	Regression Adjustment	9
4.3	Propensity Score Weighting	11
4.4	Doubly robust	12
5	Instrumental Variables	12

Remark

Incomplete, Preliminary

1 Probability and Statistics Review

1.1 Probability Space

1.2 Conditional Probability and Independence

1.3 Random Variable

1.4 Convergence

Definition 1.1 (Convergence in Probability). Let $\{X_n\}, X$ be random variables defined on $(\Omega, \mathcal{F}, \mathbb{P})$. We say X_n converges to X in probability if, $\forall \varepsilon > 0, \mathbb{P}(|X_n - X| \geq \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$. Or equivalently,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| < \varepsilon) = 1, \quad \forall \varepsilon > 0 \quad (1)$$

We denote it as $X_n \xrightarrow{p} X$.

Definition 1.2 (Almost Surely Convergence). Let $\{X_n\}, X$ be random variables defined on $(\Omega, \mathcal{F}, \mathbb{P})$. We say X_n converges to X almost surely, if

$$\mathbb{P}(\lim_{n \rightarrow \infty} X_n = X) = 1 \quad (2)$$

This can be written as $X_n \xrightarrow{a.s.} X$

It is easy to say that $X_n \xrightarrow{a.s.} X \Rightarrow X_n \xrightarrow{p} X$. But the opposite direction is not true. Counterexample: consider $((0, 1], \mathcal{B}_{(0,1]}, \lambda)$, where λ is Lebesgue measure. Let

$$\xi_n = \mathbf{1}_{(n/2^k - 1, (n+1)/2^k - 1)}, \quad 2^k \leq n < 2^{k+1}$$

Then $\mathbb{P}(|\xi_n| > \varepsilon) \leq 1/2^k \rightarrow 0$ but $\lim \xi_n(\omega)$ does not exist for any $\omega \in (0, 1]$.

Definition 1.3 (Bounded in Probability). $\{X_n\}$ is said to be bounded in probability if $\forall \varepsilon > 0, \exists M > 0$ such that

$$\inf_n \mathbb{P}(|X_n| \leq M) \geq 1 - \varepsilon \quad (3)$$

or equivalently, $\inf_n \mathbb{P}(|X_n| > M) < \varepsilon$

If $X_n \xrightarrow{p} 0$, then we denote it as $X_n = o_p(1)$.

If X_n is bounded in probability, then we denote it as $X_n = O_p(1)$. Moreover:

$$X_n = o_p(a_n) \Leftrightarrow \frac{X_n}{a_n} = o_p(1)$$

and

$$X_n = O_p(a_n) \Leftrightarrow \frac{X_n}{a_n} = O_p(1)$$

Exercise. Prove each of the followings:

- (i) $o_p(1) + o_p(1) = o_p(1)$
- (ii) $o_p(1) + O_p(1) = O_p(1)$
- (iii) $o_p(1)O_p(1) = o_p(1)$
- (iv) $(1 + o_p(1))^{-1} = O_p(1)$
- (v) $o_p(a_n) = a_n o_p(1)$
- (vi) $O_p(a_n) = a_n O_p(1)$
- (vii) $o_p(O_p(1)) = o_p(1)$

Remark. (iii) is an implication of Slutsky theorem. Since $o_p(1)O_p(1) \xrightarrow{d} 0$, then it must converge to 0 in probability by proposition 1.6

A powerful theorem to prove convergence in probability:

Theorem 1.4 (Continuous Mapping Theorem). Suppose that a measurable function g is (a.s.) continuous, then

$$X_n \xrightarrow{p} X_\infty \Rightarrow g(X_n) \xrightarrow{p} g(X_\infty) \quad (4)$$

Moreover, it applies to *a.s.* convergence and convergence in distribution.

Definition 1.5 (Convergence in Distribution). We say $X_n \xrightarrow{d} X$ if the distribution $P_n := \mathbb{P}\{X_n \in \cdot\}$ converges to $P := \mathbb{P}\{X \in \cdot\}$. Or, equivalently

$$F_{X_n}(x) \rightarrow F(x), \quad \text{for any point } x \text{ that } F(x) \text{ is continuous} \quad (5)$$

Another important result is that if for any $t \in \mathbb{R}$, the characteristic function $\phi_{X_n}(t) \rightarrow \phi_X(t)$, then $X_n \xrightarrow{d} X$. We can prove that

$$X_n \xrightarrow{a.s.} X \Rightarrow X_n \xrightarrow{p} X \Rightarrow X_n \xrightarrow{d} X \quad (6)$$

Another important result for convergence in distribution is:

Proposition 1.6. If $X_n \xrightarrow{d} c$ where c is a constant, then $X_n \xrightarrow{p} c$.

Proof. Let $X = c$, then $F_X(x) = \mathbf{1}_{\{x \geq c\}}$. Hence, $\forall \varepsilon > 0$

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}(|X_n - c| < \varepsilon) &= \lim_{n \rightarrow \infty} \mathbb{P}(c - \varepsilon < X_n < c + \varepsilon) \\ &= \lim_{n \rightarrow \infty} (F_{X_n}(c + \varepsilon) - F_{X_n}(c - \varepsilon)) \\ &= 1 - 0 = 1 \end{aligned} \quad (7)$$

Then we know that $X_n \xrightarrow{p} c$. □

Lemma 1.7 (Marginal Convergence and Joint Convergence). • If $X_n \xrightarrow{a.s.} X, Y_n \xrightarrow{a.s.} Y$, then $(X_n, Y_n) \xrightarrow{a.s.} (X, Y)$

- If $X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y$, then $(X_n, Y_n) \xrightarrow{p} (X, Y)$
- If $X_n \xrightarrow{d} X, Y_n \xrightarrow{d} Y$ and X_n, Y_n are independent for all n , X, Y are independent, then $(X_n, Y_n) \xrightarrow{d} (X, Y)$
- If $X_n \xrightarrow{d} X, Y_n \xrightarrow{d} c$, then $(X_n, Y_n) \xrightarrow{d} (X, c)$

Theorem 1.8 (Slutsky's Theorem). Let $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} c$, where c is constant.

- $X_n + Y_n \xrightarrow{d} X + c$
- $X_n Y_n \xrightarrow{d} cX$
- $X_n / Y_n \xrightarrow{d} X/c$ if $c \neq 0$

Exercise. Let $\{X_n\}$ be independent with $X_n \sim \text{Gamma}(\alpha_n, \beta_n)$. $\alpha_n \rightarrow \alpha, \beta_n \rightarrow \beta$ for some positive real number α, β . Now, let $\hat{\beta}_n$ be a consistent estimator for β . Prove that $X_n / \hat{\beta}_n \xrightarrow{d} \text{Gamma}(\alpha, 1)$

Theorem 1.9 (Delta Method). First order expansion: Suppose that g is differentiable at \mathbf{c} , for any sequence $0 < a_n \rightarrow \infty$, we have

$$a_n(\mathbf{X}_n - \mathbf{c}) \xrightarrow{d} \mathbf{X} \Rightarrow a_n[g(\mathbf{X}_n) - g(\mathbf{c})] \xrightarrow{d} [\nabla g(\mathbf{c})]^\top \mathbf{X} \quad (8)$$

If $\nabla g(\mathbf{c}) = 0$, then we have similar expansion: [To be completed](#)

Example. Suppose that $\{X_i\}$ i.i.d with mean μ and variance σ^2 . Consider the following estimator:

$$\hat{\theta} = \bar{x}^2 \quad (9)$$

Let $\theta := \mu^2$. (1) Find the asymptotic distribution of $\sqrt{n}(\hat{\theta} - \theta)$ provided $\mu \neq 0$. (2) If $\mu = 0$, find the convergence rate of $\hat{\theta}$ and its limit distribution.

Theorem 1.10 (Prohorov's Theorem). If $X_n \xrightarrow{d} X$, then $X_n = O_p(1)$

Proof. $\forall \varepsilon > 0$, we can choose M_0 sufficiently large such that

$$\mathbb{P}(|X| > M_0) < \varepsilon \quad (10)$$

Then, since $\mathbb{P}(|X_n| > M_0) \rightarrow \mathbb{P}(|X| > M_0)$, then we can choose n_0 such that for all $n \geq n_0$, $\mathbb{P}(|X_n| > M_0) < \varepsilon$. Now, we can select M_1 such that

$$\mathbb{P}(|X_i| > M_1) < \varepsilon, \quad \forall i = 1, \dots, n_0 - 1 \quad (11)$$

Then let $M = \max(M_0, M_1)$ we have $\mathbb{P}(|X_n| > M) < \varepsilon$ for all n . □

A natural question is: will bounded in probability imply convergence in distribution? The answer is **No**. Consider $X_n = 2 + 1/n$ for even n and $X_n = 1 + 1/(n+1)$ for odd n . Then the sequence (X_{2k}) converges in distribution to $Y = 2$. And (X_{2k-1}) converges in distribution to $W = 1$. Since $Y \neq W$ then the sequence does not converge in distribution. Since all X_n lie in the interval $[1, 5/2]$, then we can easily show that $X_n = O_p(1)$.

1.5 Law of Large Numbers

Theorem 1.11 (WLLN, Khintchin). If $\{X_n\}$ are i.i.d with $\mathbb{E}[X_1] = \mu < \infty$, then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu \quad (12)$$

Theorem 1.12 (SLLN, Kolmogorov). If $\{X_n\}$ are i.i.d with $\mathbb{E}[X_1] = \mu < \infty$, then

$$\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mu \quad (13)$$

1.6 Central Limit Theorem

Theorem 1.13 (Levy CLT). Suppose that $\{X_n\}$ i.i.d with mean μ and variance σ^2 , then

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2) \quad (14)$$

where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.

1.7 Normal Distribution

Consider multivariate normal distribution: $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$. The moment generating function is

$$M_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}[e^{\mathbf{t}'\mathbf{X}}] = e^{\mathbf{t}'\boldsymbol{\mu} + \frac{1}{2}\mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}}, \quad \mathbf{t} \in \mathbb{R}^d \quad (15)$$

For any $\mathbf{A} \in \mathbb{R}^{m \times d}$, we have $\mathbf{AX} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$. The probability density function of \mathbf{X} is given by

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \det(\boldsymbol{\Sigma})^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}, \quad \mathbf{x} \in \mathbb{R}^d \quad (16)$$

- If $X_1, \dots, X_n \sim N(0, 1)$ i.i.d, then $X_1^2 + \dots + X_n^2 \sim \chi_n^2$.
- If $X \sim N(0, 1)$ and $Q \sim \chi_n^2$ are independent, then $\frac{X}{\sqrt{Q/n}} \sim t_n$.
- If $Q_1 \sim \chi_m^2, Q_2 \sim \chi_n^2$ are independent, then $\frac{Q_1/m}{Q_2/n} \sim F_{m,n}$.

1.8 Hypothesis Testing

Consider $H_0 : \theta \in \Theta$, this is called the null hypothesis. The alternative hypothesis is $H_1 : \theta \in \Theta_1$, where $\Theta_1 = \Theta \setminus \Theta_0$. We have to decide between H_0 and H_1 . Let R denotes the reject region. A Test might have two types of mistake.

- **Type I Error:** Reject H_0 when $\theta \in \Theta_0$. $\mathbb{P}_{\theta}(\mathbf{X} \in R)$ for $\theta \in \Theta_0$
- **Type II Error:** Accept H_0 when $\theta \in \Theta_1$. $\mathbb{P}_{\theta}(\mathbf{X} \in R^c)$ for $\theta \in \Theta_1$

In this notes, we will use φ to denote *power function* for a hypothesis test, i.e.

$$\beta(\theta) = \mathbb{P}_{\theta}(\mathbf{X} \in R) \quad (17)$$

When $\theta \in \Theta_0$, then $\beta(\theta) = \mathbb{P}(\text{Type I Error})$. If $\theta \in \Theta_1$, then $\beta(\theta) = 1 - \mathbb{P}(\text{Type II Error})$

Definition 1.14. For $\alpha \in [0, 1]$, a test with power function $\beta(\theta)$ is a *size α test* if

$$\sup_{\theta \in \Theta_0} \beta(\theta) = \alpha,$$

is a *level α test* if

$$\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$$

Definition 1.15. A test is *unbiased* if $\beta(\theta') \geq \beta(\theta'')$ for all $\theta' \in \Theta_1$ and $\theta'' \in \Theta_0$. A test is *consistent* if

$$\lim_{n \rightarrow \infty} \beta(\theta) = 1, \quad \forall \theta \in \Theta_1$$

Definition 1.16. Let \mathcal{C} be a class of tests. A test in class \mathcal{C} with power function $\beta(\theta)$ is a *uniformly most powerful* (UMP) class \mathcal{C} test if $\beta(\theta) \geq g(\theta)$ for every $\theta \in \Theta_1$ and every $g(\theta)$ that is a power function of a test in class \mathcal{C} . Generally we take \mathcal{C} as the class of all level α test.

Theorem 1.17 (Neyman-Pearson Lemma). Consider testing $H_0 : \theta = \theta_0$ versus $H_1 : \theta = \theta_1$, where the pdf or pmf corresponding to θ_i is $f(x|\theta_i), i = 0, 1$. Then

$$R = \{x : f(x|\theta_1) > k f(x|\theta_0)\}$$

for some $k \geq 0$ and $\alpha = \mathbb{P}_{\theta_0}(X \in R)$ is a UMP level α test.

Definition 1.18 (*p*-value). A *p*-value $p(X)$ is a test statistic satisfying $0 \leq p(x) \leq 1$ for every sample point x . Small values of $p(X)$ give evidence that H_1 is true. A *p*-value is valid if, for every $\theta \in \Theta_0$ and every $0 \leq \alpha \leq 1$,

$$\mathbb{P}_{\theta}(p(X) \leq \alpha) \leq \alpha \quad (18)$$

If $p(X)$ is a valid *p*-value, then it is easy to construct a level α test based on this statistics. We reject H_0 if and only if $p(X) \leq \alpha$.

Theorem 1.19. Suppose that $T(X)$ is a test statistic such that large values of T give evidence that H_1 is true. For each sample point x , define

$$p(x) = \sup_{\theta \in \Theta_0} \mathbb{P}_{\theta}(T(X) \geq T(x)) \quad (19)$$

Then $p(X)$ is a valid *p*-value.

Now, suppose that β is a parameter, and $\hat{\beta}$ is an estimator of β . Moreover, we assume that $\hat{\beta}$ is consistent and asymptotically normal, i.e.

$$\hat{\beta} \xrightarrow{P} \beta, \quad \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma^2)$$

Also, suppose that $\hat{\sigma}^2 \xrightarrow{P} \sigma^2$ is an estimator for asymptotical variance. Now we want to test the hypothesis: $H_0 : \beta = c$, where c is a constant. To this end, we can employ *t*-test:

$$T = \frac{\hat{\beta} - c}{\text{se}(\hat{\beta})} = \frac{\hat{\beta} - c}{\hat{\sigma}/\sqrt{n}} \quad (20)$$

Then under the null, T is asymptotically normal since by Slutsky's theorem, we have

$$\frac{\hat{\beta} - c}{\hat{\sigma}/\sqrt{n}} = \frac{1}{\hat{\sigma}} \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, 1) \quad (21)$$

Then let the rejection region be

$$R := \{|T| \geq z_{\alpha/2}\}, \quad z_{\alpha/2} := \Phi^{-1}(1 - \alpha/2)$$

where $\Phi(x)$ is standard normal cdf. Also, we can construct confidence interval by test inversion:

$$\text{CI}_{1-\alpha} = [\hat{\beta} - z_{\alpha/2} \text{se}(\hat{\beta}), \hat{\beta} + z_{\alpha/2} \text{se}(\hat{\beta})] \quad (22)$$

Exercise. To be completed.

1.9 Wald's Test

1.10 Introduction to Bayesian Statistics

1.11 Some Exercise

2 Causal Inference and Potential Outcomes

2.1 Potential Outcomes Model

We denote x_i the treatment status, i.e. $x_i = 1$ implies the individual is treated and $x_i = 0$ indicates the individual is not treated. Also, the out come $y_i(1)$ is "the outcome of the i^{th} individual had she received the treatment. Therefore, for each i , we have three random variables, $(x_i, y_i(1), y_i(0))$.

Definition 2.1. Individual treatment effect (ITE):

$$\tau_i = y_i(1) - y_i(0)$$

Average treatment effect (ATE):

$$\tau_{\text{ATE}} = \mathbb{E}[\tau_i] = \mathbb{E}[y_i(1) - y_i(0)]$$

Treatment effect on the treated (ATT):

$$\tau_{\text{ATT}} = \mathbb{E}[\tau_i | x_i = 1] = \mathbb{E}[y_i(1) - y_i(0) | x_i = 1]$$

Treatment effect on the untreated (ATU):

$$\tau_{\text{ATU}} = \mathbb{E}[\tau_i | x_i = 0] = \mathbb{E}[y_i(1) - y_i(0) | x_i = 0]$$

We can always write y_i as

$$y_i = y_i(1) \mathbb{1}_{\{x_i=1\}} + y_i(0) \mathbb{1}_{\{x_i=0\}} = x_i y_i(1) + (1 - x_i) y_i(0) \quad (23)$$

In words, we will never be able to observe the two potential outcomes simultaneously for any individual.

2.2 Randomized Controlled Trials

Definition 2.2 (Missing completely at random; independent treatment). The potential outcomes are said to be missing completely at random, if

$$x_i \perp\!\!\!\perp (y_i(1), y_i(0)) \quad (24)$$

It means the treatment status is independent of the potential outcomes. And therefore two individuals with different treatment status should not differ systematically. This is a very strong assumption. But there is a special case that we believe this assumption holds: randomized controlled trials (RCT).

In an RCT, a sample of units are randomized into two groups, the treatment group and the control group. Then individuals who are in the treatment group will be exposed to the treatment, while those in the control group are not. Due to the randomization of treatment, it is plausible to believe that treatment assignment is independent of the potential outcomes, and hence this assumption holds.

How this assumption will help to identify the treatment effects? Consider the ATE:

$$\tau_{ATE} = \mathbb{E}[y_i(1) - y_i(0)] = \mathbb{E}[y_i(1)] - \mathbb{E}[y_i(0)] \quad (25)$$

Since the treatment is independent of the potential outcomes, then

$$\mathbb{E}[y_i(1)] = \mathbb{E}[y_i(1)|x_i = 1] = \mathbb{E}[y_i|x_i = 1] \quad (26)$$

Similarly, $\mathbb{E}[y_i(0)] = \mathbb{E}[y_i|x_i = 0]$. Hence

Proposition 2.3. Under the independent treatment assumption, the ATE is identified by

$$\tau_{ATE} = \mathbb{E}[y_i(1) - y_i(0)] = \mathbb{E}[y_i|x_i = 1] - \mathbb{E}[y_i|x_i = 0] \quad (27)$$

Now we re-write the potential outcomes into an expectation component and an error term:

$$y_i(1) = \mathbb{E}[y_i(1)] + \underbrace{y_i(1) - \mathbb{E}[y_i(1)]}_{u_i(1)}, \quad y_i(0) = \mathbb{E}[y_i(0)] + \underbrace{y_i(0) - \mathbb{E}[y_i(0)]}_{u_i(0)} \quad (28)$$

Then

$$\begin{aligned} y_i &= x_i y_i(1) + (1 - x_i) y_i(0) \\ &= x_i (\mathbb{E}[y_i(1)] + u_i(1)) + (1 - x_i) (\mathbb{E}[y_i(0)] + u_i(0)) \\ &= \mathbb{E}[y_i(0)] + x_i (\mathbb{E}[y_i(1)] - \mathbb{E}[y_i(0)]) + x_i u_i(1) + (1 - x_i) u_i(0) \\ &= \beta_0 + \beta_1 x_i + u_i \end{aligned} \quad (29)$$

Then under the independent treatment assumption, we have the regression expression, where $\beta_0 = \mathbb{E}[y_i(0)]$, $\beta_1 = \tau_{ATE}$ and $\mathbb{E}[u_i|x_i] = 0$. The OLS estimators are

$$\hat{\beta}_0 = \frac{1}{n_0} \sum_{i=1}^n y_i \mathbb{1}_{\{x_i=0\}}, \quad \hat{\beta}_1 = \frac{1}{n_1} \sum_{i=1}^n y_i \mathbb{1}_{\{x_i=1\}} - \frac{1}{n_0} \sum_{i=1}^n y_i \mathbb{1}_{\{x_i=0\}} \quad (30)$$

where $n_1 = \sum_{i=1}^n x_i$ is the size of the treatment group, and $n_0 = \sum_{i=1}^n (1 - x_i)$ is the size of the control group. Also we can re-write them as

$$\hat{\beta}_0 = \frac{n}{n_0} \frac{1}{n} \sum_{i=1}^n y_i (1 - x_i) = \frac{n}{n_0} \frac{1}{n} \sum_{i=1}^n y_i(0) (1 - x_i)$$

Note that by LLN

$$\frac{n}{n_0} \xrightarrow{p} \frac{1}{\mathbb{P}(x_i = 0)}$$

and

$$\frac{1}{n} \sum_{i=1}^n y_i(0) (1 - x_i) \xrightarrow{p} \mathbb{E}[y_i(0) (1 - x_i)] = \mathbb{E}[y_i(0)] \mathbb{E}[1 - x_i] = \mathbb{E}[y_i(0)] \mathbb{P}(x_i = 0)$$

Hence

$$\hat{\beta}_0 \xrightarrow{p} \beta_0 = \mathbb{E}[y_i(0)]$$

Also

$$\frac{n}{n_1} \frac{1}{n} \sum_{i=1}^n y_i \mathbb{1}_{\{x_i=1\}} \xrightarrow{p} \frac{1}{\mathbb{P}(x_i = 1)} \mathbb{E}[y_i(1)] \mathbb{P}(x_i = 1) = \mathbb{E}[y_i(1)]$$

Hence

$$\hat{\beta}_1 \xrightarrow{p} \beta_1 = \mathbb{E}[y_i(1)] - \mathbb{E}[y_i(0)] = \tau_{ATE}$$

2.3 Selection Bias

What if we don't assume the independent assumption? Then $\hat{\beta}_1$ is consistent for

$$\hat{\beta}_1 \xrightarrow{p} \mathbb{E}[y_i(1)|x_i = 1] - \mathbb{E}[y_i(0)|x_i = 0]$$

Without the independent assumption, we can not pull out ATE from the conditional expectation, but we have

$$\mathbb{E}[y_i(1)|x_i = 1] - \mathbb{E}[y_i(0)|x_i = 1] = \underbrace{\mathbb{E}[y_i(1)|x_i = 1] - \mathbb{E}[y_i(0)|x_i = 1]}_{\tau_{ATT}} + \underbrace{\mathbb{E}[y_i(0)|x_i = 1] - \mathbb{E}[y_i(0)|x_i = 0]}_{\text{selection bias}}$$

3 Linear Regression

Definition 3.1 (Missing at random; conditional independence; selection on observables). The potential outcomes are said to be missing at random (or, alternatively, the treatment assignment satisfies the selection on observables assumption), if

$$x_i \perp\!\!\!\perp (y_i(1), y_i(0)) | w_i \quad (31)$$

where w_i is some characteristics of each individual that can be observed by researcher.

Then we have

$$y_i(1) = \underbrace{\mathbb{E}[y_i(1)|w_i]}_{g_1(w_i)} + u_i(1), \quad y_i(0) = \underbrace{\mathbb{E}[y_i(0)|w_i]}_{g_0(w_i)} + u_i(0) \quad (32)$$

where $u_i(1) = y_i(1) - g_1(w_i)$ and $u_i(0) = y_i(0) - g_0(w_i)$. Then

$$y_i = g_0(w_i) + x_i(g_1(w_i) - g_0(w_i)) + \underbrace{x_i u_i(1) + (1 - x_i) u_i(0)}_{u_i} \quad (33)$$

We have to make another very strong assumption:

$$g_1(w_i) = \mathbb{E}[y_i(1)] + w_i^\top \delta, \quad g_0(w_i) = \mathbb{E}[y_i(0)] + w_i^\top \delta \quad (34)$$

Then the outcome variable is

$$y_i = \mathbb{E}[y_i(0)] + w_i^\top \delta + x_i (\mathbb{E}[y_i(1)] + w_i^\top \delta - \mathbb{E}[y_i(0)] - w_i^\top \delta) + u_i = \mathbb{E}[y_i(0)] + x_i \tau_{ATE} + w_i^\top \delta + u_i \quad (35)$$

3.1 Algebraic Properties

Model:

$$y_i = x_i^\top \beta + u_i \quad (36)$$

we assume that $\mathbb{E}[u_i] = 0$ and $\mathbb{E}[u_i x_i] = 0$. The moment condition is

$$\frac{1}{n} \sum_{i=1}^n x_i (y_i - x_i^\top \hat{\beta}) = 0 \quad (37)$$

Standard algebra leads to

$$\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) \quad (38)$$

Matrix form: $\hat{\beta} = (X^\top X)^{-1} X^\top y$.

Definition 3.2 (Sum of squares, R-squared). Let $\hat{y}_i = x_i^\top \hat{\beta}$. Define

$$\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2, \quad \text{ESS} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \quad \text{SSR} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

We have $\text{TSS} = \text{ESS} + \text{SSR}$. A goodness-of-fit measure is the R -squared, which is defined as

$$R^2 = \frac{\text{ESS}}{\text{TSS}} \quad (39)$$

We often define projection matrix

$$P_x = X(X'X)^{-1}X'y, \quad \hat{y} = P_x y \quad (40)$$

It has following properties: (1) P_x is symmetric. (2) P_x is idempotent, i.e. $P_x P_x = P_x$. (3) $P_x X = X$.

Similarly, we can define $M_x = I - P_x$, is called the annihilator or the residual maker. Then $\hat{u} = M_x y$. It has following properties: (1) M_x is symmetric (2) M_x is idempotent, i.e. $M_x M_x = M_x$ (3) $M_x \hat{u} = \hat{u}$ (4) $M_x P_x = P_x M_x = 0$.

Also, for any vector $a \in \mathbb{R}^n$, $\|P_x a\| \leq \|a\|$, $\|M_x a\| \leq \|a\|$.

Theorem 3.3 (Frisch-Waugh-Lovell). Suppose that $x_{i1} \in \mathbb{R}^{d_1}$, $x_{i2} \in \mathbb{R}^{d_2}$ for every i . And

$$y_i = x_{i1}^\top \beta_1 + x_{i2}^\top \beta_2 + u_i \quad (41)$$

The estimated coefficient of x_{i2} in the regression of y_i on x_{i1} and x_{i2} is given by

$$\hat{\beta}_2 = (X_2^\top M_{X_1} X_2)^{-1} (X_2^\top M_{X_1} y) = (\check{X}_2^\top \check{X}_2)^{-1} (\check{X}_2^\top \check{y}) \quad (42)$$

where \check{X}_2 and \check{y} are the residuals obtained from regression x_{i2} and y_i on x_{i1} , respectively.

Proof. Consider

$$y = X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 + \hat{u}$$

Then pre-multiply M_{X_1} at BHS,

$$M_{X_1} y = 0 + M_{X_1} X_2 \hat{\beta}_2 + M_{X_1} \hat{u} \Rightarrow M_{X_1} y - M_{X_1} X_2 \hat{\beta}_2 = \hat{u}$$

Then, multiply X_2^\top at BHS

$$X_2^\top M_{X_1} y - X_2^\top M_{X_1} X_2 \hat{\beta}_2 = X_2^\top \hat{u} = 0 \Rightarrow \hat{\beta}_2 = (X_2^\top M_{X_1} X_2)^{-1} (X_2^\top M_{X_1} y)$$

□

A simple, but useful application of this is, consider $x_i = (z_i, 1)$. Then run the regression, what will we get?

3.2 Large Sample Properties

Note that

$$\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right) = \beta + \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i u_i \right) \quad (43)$$

We first make an assumption on the data generating process.

- (i) x_i has a finite second moments, $\mathbb{E}[|x_i|^2] < \infty$. And the matrix $\mathbb{E}[x_i x_i^\top]$ is invertible.
- (ii) The error term u_i is mean-zero, $\mathbb{E}[u_i] = 0$. And it has finite variance, also uncorrelate with x : $\mathbb{E}[x_i u_i] = 0$

Proposition 3.4. Under those assumptions:

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \xrightarrow{p} \mathbb{E}[x_i x_i^\top], \quad \mathbb{P} \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \text{ is non-singular} \right) \rightarrow 1 \quad (44)$$

Also,

$$\frac{1}{n} \sum_{i=1}^n x_i u_i \xrightarrow{p} 0$$

Hence

$$\hat{\beta} \xrightarrow{p} \beta, \quad \hat{\beta} = \beta + o_p(1) \quad (45)$$

If we further assume x_i and u_i has finite fourth moments, and $\mathbb{E}[x_i x_i^\top u_i^2]$ is non-singular, then

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, V), \quad V = (\mathbb{E}[x_i x_i])^{-1} (\mathbb{E}[x_i x_i^\top u_i^2]) (\mathbb{E}[x_i x_i])^{-1} \quad (46)$$

3.3 Standard Error

The challenge is to estimate $\mathbb{E}[x_i x_i^\top u_i^2]$, since it involves unknown u_i . Now consider

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \hat{u}_i^2, \quad \hat{u}_i = y_i - x_i^\top \beta$$

We can not apply LLN here since \hat{u}_i are not i.i.d. Hence we have to decompose it first.

$$\hat{u}_i^2 - u_i^2 = (\hat{u}_i - u_i)(\hat{u}_i + u_i), \quad \hat{u}_i - u_i = x_i^\top (\beta - \hat{\beta}) \Rightarrow \hat{u}_i^2 = u_i^2 + x_i^\top (\beta - \hat{\beta})(y_i - x_i^\top \hat{\beta} + u_i)$$

Hence

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \hat{u}_i^2 = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top u_i^2 + \frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) x_i^\top (\beta - \hat{\beta}) - \frac{1}{n} \sum_{i=1}^n x_i x_i^\top x_i^\top \hat{\beta} x_i^\top (\beta - \hat{\beta}) \quad (47)$$

Then we know that

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^\top u_i^2 \xrightarrow{p} \mathbb{E}[x_i x_i^\top u_i^2], \quad \left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) x_i^\top (\beta - \hat{\beta}) \right\| \leq \|\beta - \hat{\beta}\| \left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) \right\| \quad (48)$$

Since $\beta - \hat{\beta} \xrightarrow{p} 0$, and $\frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) \xrightarrow{p} \mathbb{E}[x_i x_i^\top x_i^\top (y_i + u_i)]$. Hence, by $o_p(1)O_p(1) = o_p(1)$, we know that

$$\left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) x_i^\top (\beta - \hat{\beta}) \right\| \xrightarrow{p} 0 \quad (49)$$

Moreover, impose $\beta - \hat{\beta} = O_p(\frac{1}{\sqrt{n}})$, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top (y_i + u_i) x_i^\top (\beta - \hat{\beta}) \right\| = O_p(\frac{1}{\sqrt{n}}) = o_p(1) \quad (50)$$

Similarly,

$$\left\| \frac{1}{n} \sum_{i=1}^n x_i x_i^\top x_i^\top \hat{\beta} x_i^\top (\beta - \hat{\beta}) \right\| \leq \|\beta - \hat{\beta}\| \|\hat{\beta}\| \frac{1}{n} \sum_{i=1}^n \|x_i\|^4 = O_p(\frac{1}{\sqrt{n}}) = o_p(1)$$

Hence

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \hat{u}_i^2 \xrightarrow{p} \mathbb{E}[x_i x_i^\top u_i^2]$$

Then

$$\hat{V} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left[\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \hat{u}_i^2 \right] \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \xrightarrow{p} V \quad (51)$$

Sometimes we also write

$$\hat{\beta} \overset{a}{\sim} N(\beta, \frac{\hat{V}}{n}) \quad (52)$$

The standard error (or, more precisely, the variance estimator) we discussed above is widely known as the Huber-Eicker-White standard error. We also call it HC0, where HC stands for “heteroskedasticity consistent. `robust` in `Stata` corresponds to the original HC0 standard error. Compared to HC0, HC1 incorporates a degrees of freedom adjustment of $n/(n-d)$.

$$\begin{aligned} \hat{V}_{HC1} &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left[\frac{1}{n-d} \sum_{i=1}^n x_i x_i^\top \hat{u}_i^2 \right] \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} = \frac{n}{n-d} \hat{V} \\ \hat{V}_{HC2} &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left[\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\hat{u}_i^2}{1-p_{ii}} \right] \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1}, \quad p_{ii} = x_i^\top (X'X)^{-1} x_i, \\ \hat{V}_{HC3} &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left[\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\hat{u}_i^2}{(1-p_{ii})^2} \right] \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \\ \hat{V}_{HC4} &= \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1} \left[\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \frac{\hat{u}_i^2}{(1-p_{ii})^{\delta_i}} \right] \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right)^{-1}, \quad \delta_i = \min \left\{ 4, \frac{np_{ii}}{d} \right\} \end{aligned} \quad (53)$$

There are a lot of exercises to be done here

Exercise. Prove the Gauss-Markov Theorem

3.4 Hypothesis Testing

Consider general problem

$$H_0 : R\beta = r \quad (54)$$

Note that $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, V)$. Then under the null

$$\sqrt{n}(R\hat{\beta} - R\beta) = \sqrt{n}(R\hat{\beta} - r) \xrightarrow{d} N(0, RV R')$$

We consider the Wald's statistic for this test:

$$n(R\hat{\beta} - r)'[RV R']^{-1}(R\hat{\beta} - r) \xrightarrow{d} \chi_k^2$$

4 Selection on Observables

4.1 Overlap

Assumption: Let $\{(y_i, x_i, w_i) : 1 \leq i \leq n\}$ be a random sample of size n , where y_i is the outcome variable of interest, x_i is a binary indicator of treatment status, and w_i are some additional observable characteristics.

(i) The potential outcomes are missing at random

$$x_i \perp\!\!\!\perp (y_i(1), y_i(0)) | w_i \quad (55)$$

(ii) The conditional probability of receiving treatment satisfies

$$0 < \mathbb{P}(x_i = 1 | w_i) < 1 \quad (56)$$

Here we do not assume the linearity on conditional expectation. Instead, define

$$g_1(w) = \mathbb{E}[y_i(1) | w_i = w], \quad g_0(w) = \mathbb{E}[y_i(0) | w_i = w]$$

It is easy to show, under the assumption (i), we must have

$$g_1(w) = \mathbb{E}[y_i | x_i = 1, w_i = w], \quad g_0(w) = \mathbb{E}[y_i | x_i = 0, w_i = w]$$

Assume that w_i is discrete, then $0 < \mathbb{P}(x_i = 1 | w_i = w)$ is equivalent to

$$\mathbb{P}(w_i = w | x_i = 1), \mathbb{P}(w_i = w | x_i = 0) > 0 \quad \text{and} \quad 0 < \mathbb{P}(x_i = 1) < 1$$

If w_i is continuous, then $0 < \mathbb{P}(x_i = 1 | w_i = w) < 1$ is equivalent to positive conditional densities of w_i given $x_i = 1$. Also, by Bayes theorem,

$$\mathbb{P}(x_i = 1 | w_i = w) = \frac{f_{w|x_i=1}(w)\mathbb{P}(x_i = 1)}{f_w(w)} = \frac{f_{w|x=1}(w)\mathbb{P}(x_i = 1)}{f_{w|x=1}(w)\mathbb{P}(x_i = 1) + f_{w|x=0}(w)\mathbb{P}(x_i = 0)}$$

Thus, the condition (ii) implies that $f_{w|x=0} > 0, f_{w|x=1} > 0$.

4.2 Regression Adjustment

Under the assumption in this chapter, we can identify the conditional expectation functions:

$$\begin{aligned} g_1(w) &= \mathbb{E}[y_i(1) | w_i = w] = \mathbb{E}[y_i | x_i = 1, w_i = w] \\ g_0(w) &= \mathbb{E}[y_i(0) | w_i = w] = \mathbb{E}[y_i | x_i = 0, w_i = w] \end{aligned} \quad (57)$$

The next question is how we can recover the different treatment effect parameters/estimands. We will assume that w_i is continuous. The conditional densities are denoted $f_{w|x=1}$ and $f_{w|x=0}$ respectively. Then

$$\tau_{ATE} = \mathbb{E}[y_i(1) - y_i(0)] = \mathbb{E}[g_1(w) - g_0(w)] = \int [g_1(w) - g_0(w)] f_w(w) dw$$

$g_1(w) - g_0(w)$ is sometimes called the conditional treatment effect (CTE). We will use $\tau(w)$ to denote CTE. To estimate, assume we already have consistent estimators for the two conditional expectation functions, which we will denote $\hat{g}_1(w)$ and $\hat{g}_0(w)$ correspondingly. Then

$$\hat{\tau}_{ATE} = \frac{1}{n} \sum_{i=1}^n (\hat{g}_1(w_i) - \hat{g}_0(w_i))$$

We consider a special case where the two conditional expectation functions are estimated using linear regressions:

$$\begin{aligned} \hat{g}_1(w_i) &= w_i^\top \hat{\beta}_1, \quad \hat{\beta}_1 = \operatorname{argmin}_b \sum_{i=1}^n x_i (y_i - w_i^\top b)^2 \\ \hat{g}_0(w_i) &= w_i^\top \hat{\beta}_0, \quad \hat{\beta}_0 = \operatorname{argmin}_b \sum_{i=1}^n (1 - x_i) (y_i - w_i^\top b)^2 \end{aligned}$$

where $\hat{\beta}_1$ and $\hat{\beta}_0$ are interpreted regression y_i on w_i using treated and untreated sample correspondingly, i.e.

$$\hat{\beta}_1 = \left(\frac{1}{n} \sum_{i=1}^n w_i^\top w_i x_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n w_i y_i x_i \right), \quad \hat{\beta}_0 = \left(\frac{1}{n} \sum_{i=1}^n w_i^\top w_i (1 - x_i) \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n w_i y_i (1 - x_i) \right) \quad (58)$$

Then the ATE estimator can be written as

$$\hat{\tau}_{ATE} = \bar{w}^\top (\hat{\beta}_1 - \hat{\beta}_0), \quad \bar{w} = \frac{1}{n} \sum_{i=1}^n w_i$$

Then one can decompose it as

$$\hat{\tau}_{ATE} = \underbrace{\mathbb{E}[w_i]^\top (\beta_1 - \beta_0)}_{\text{True ATE}} + \mathbb{E}[w_i]^\top (\hat{\beta}_1 - \beta_1) - \mathbb{E}[w_i]^\top (\hat{\beta}_0 - \beta_0) + (\bar{w} - \mathbb{E}[w_i])^\top (\hat{\beta}_1 - \hat{\beta}_0)$$

Hence it is a consistent estimator. Now consider the asymptotic distribution of $\hat{\tau}_{ATE}$:

$$\sqrt{n}(\hat{\tau}_{ATE} - \tau_{ATE}) = \mathbb{E}[w_i]^\top \sqrt{n}(\hat{\beta}_1 - \beta_1) - \mathbb{E}[w_i]^\top \sqrt{n}(\hat{\beta}_0 - \beta_0) + \sqrt{n}(\bar{w} - \mathbb{E}[w_i])^\top (\hat{\beta}_1 - \hat{\beta}_0)$$

Recall that

$$\frac{1}{n} \sum_{i=1}^n w_i w_i^\top x_i \xrightarrow{P} \mathbb{P}(x_i = 1) \mathbb{E}[w_i w_i^\top | x_i = 1] \quad (59)$$

Hence

$$\sqrt{n}(\hat{\beta}_1 - \beta_1) = \mathbb{E}[w_i w_i^\top | x_i = 1]^{-1} \frac{1}{\mathbb{P}(x_i = 1)} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n x_i w_i (y_i - w_i^\top \beta_1) \right) + o_p(1)$$

and

$$\sqrt{n}(\hat{\beta}_0 - \beta_0) = \mathbb{E}[w_i w_i^\top | x_i = 0]^{-1} \frac{1}{\mathbb{P}(x_i = 0)} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n (1 - x_i) w_i (y_i - w_i^\top \beta_0) \right) + o_p(1)$$

Define

$$\begin{aligned} \psi_{i,1} &= \mathbb{E}[w_i]^\top \mathbb{E}[w_i w_i^\top | x_i = 1]^{-1} \frac{1}{\mathbb{P}(x_i = 1)} x_i w_i (y_i - w_i^\top \beta_1) \\ \psi_{i,2} &= -\mathbb{E}[w_i]^\top \mathbb{E}[w_i w_i^\top | x_i = 0]^{-1} \frac{1}{\mathbb{P}(x_i = 0)} (1 - x_i) w_i (y_i - w_i^\top \beta_0) \end{aligned} \quad (60)$$

The remaining term is

$$\sqrt{n}(\bar{w} - \mathbb{E}[w_i])^\top (\hat{\beta}_1 - \hat{\beta}_0) = \sqrt{n}(\bar{w} - \mathbb{E}[w_i])^\top (\beta_1 - \beta_0 + \hat{\beta}_1 - \hat{\beta}_0 + \hat{\beta}_0 - \beta_0)$$

Since $\hat{\beta}_1 - \beta_1 = o_p(1)$, $\hat{\beta}_0 - \beta_0 = o_p(1)$, then

$$\sqrt{n}(\bar{w} - \mathbb{E}[w_i])^\top (\hat{\beta}_1 - \hat{\beta}_0) = (\beta_1 - \beta_0)^\top \frac{1}{\sqrt{n}} \sum_{i=1}^n (w_i - \mathbb{E}[w_i])$$

Define $\psi_{i,3} = (\beta_1 - \beta_0)^\top (w_i - \mathbb{E}[w_i])$, then

$$\sqrt{n}(\hat{\tau}_{ATE} - \tau_{ATE}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\psi_{i,1} + \psi_{i,2} + \psi_{i,3}) \xrightarrow{d} N(0, \sigma^2), \quad \sigma^2 = \operatorname{Var}(\psi_{i,1} + \psi_{i,2} + \psi_{i,3})$$

To estimate the asymptotic variance, we can employ the sample analogue

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\hat{\psi}_{i,1} + \hat{\psi}_{i,2} + \hat{\psi}_{i,3})^2$$

where

$$\begin{aligned}\hat{\psi}_{i,1} &= \bar{w}^\top \left(\frac{1}{n_1} \sum_{i=1}^n x_i w_i w_i^\top \right)^{-1} \frac{n}{n_1} x_i w_i (y_i - w_i^\top \hat{\beta}_1) \\ \hat{\psi}_{i,2} &= -\bar{w}^\top \left(\frac{1}{n_0} \sum_{i=1}^n (1-x_i) w_i w_i^\top \right)^{-1} \frac{n}{n_0} (1-x_i) w_i (y_i - w_i^\top \hat{\beta}_0) \\ \hat{\psi}_{i,3} &= (\hat{\beta}_1 - \hat{\beta}_0)^\top (w_i - \bar{w})\end{aligned}$$

Next we study the identification and estimation of the average treatment effect on the treated (ATT). We again rewrite the ATT as a weighted average

$$\tau_{ATE} = \int [g_1(w) - g_0(w)] f_{w|x=1}(w) dw = \mathbb{E}[y_i | x_i = 1] - \int g_0(w) f_{w|x=1}(w) dw$$

Hence

$$\hat{\tau}_{ATT} = \frac{1}{n_1} \sum_{i=1}^n x_i (y_i - \hat{g}_0(w_i)) = \frac{1}{n_1} \sum_{i=1}^n x_i (y_i - w_i^\top \hat{\beta}_0) = \bar{y}_1 - \bar{w}_1^\top \hat{\beta}_0$$

where

$$\bar{y}_1 = \frac{1}{n_1} \sum_{i=1}^n y_i x_i, \quad \bar{w}_1 = \frac{1}{n_1} \sum_{i=1}^n w_i x_i$$

4.3 Propensity Score Weighting

The propensity score is defined as

$$e(w) = \mathbb{P}(x_i = 1 | w_i = w)$$

Then we know that (note that $x_i y_i = x_i y_i(1)$)

$$\mathbb{E}[x_i y_i | w_i] = \mathbb{E}[x_i | w_i] \mathbb{E}[y_i | w_i] = e(w_i) g_1(w_i)$$

Hence

$$\mathbb{E} \left[\frac{x_i y_i}{e_i} \right] = \mathbb{E} \left[\mathbb{E} \left(\frac{x_i y_i(1)}{e_i} \middle| w_i \right) \right] = \mathbb{E}[\mathbb{E}(y_i(1) | w_i)] = \mathbb{E}[y_i(1)]$$

Similarly, we can derive

$$\mathbb{E} \left[\frac{(1-x_i) y_i}{1-e_i} \right] = \mathbb{E}[y_i(0)]$$

Then

$$\begin{aligned}\tau_{ATE} &= \mathbb{E} \left[\frac{x_i y_i}{e_i} - \frac{(1-x_i) y_i}{1-e_i} \right] \\ \tau_{ATT} &= \mathbb{E} \left[\frac{x_i y_i}{e} - \frac{e_i (1-x_i) y_i}{e(1-e_i)} \right] \\ \tau_{ATU} &= \mathbb{E} \left[\frac{(1-e_i) x_i y_i}{(1-e) e_i} - \frac{(1-x_i) y_i}{1-e} \right]\end{aligned}$$

where $e = \mathbb{E}[e_i] = \mathbb{P}(x_i = 1)$ is the unconditional probability of receiving treatment. For simplicity, we will assume that the parameter of interest is $\beta = \mathbb{E}[y_i(1)]$, then

$$\hat{\beta} = \frac{1}{n} \sum_{i=1}^n \frac{x_i y_i}{\hat{e}_i}, \quad \hat{e}_i = w_i^\top \hat{\gamma}, \quad \hat{\gamma} = \left(\frac{1}{n} w_i w_i^\top \right)^{-1} \left(\frac{1}{n} w_i x_i \right)$$

Given the linear regression specification, we are implicitly assuming that the true propensity score is linear, that is $x_i = e_i + v_i = w_i^\top \gamma + v_i$, where $\mathbb{E}[v_i | w_i] = 0$. We only consider this linear propensity score model to simplify the asymptotic analysis. In practice, it is more common to use nonlinear models such as the logit or probit, which we will discuss later

To establish asymptotic properties, we will assume the true propensity score is bounded away from 0 and 1, that is, there exists some $\delta > 0$ such that $\delta \leq e_i \leq 1 - \delta$. Then

$$\hat{\beta} - \beta = \underbrace{\frac{1}{n} \sum_{i=1}^n \frac{x_i y_i}{e_i} - \beta}_{\xrightarrow{p} 0} + \frac{1}{n} \sum_{i=1}^n \frac{x_i y_i}{e_i \hat{e}_i} w_i^\top (\gamma - \hat{\gamma})$$

Then we further assume that $|w_i| \leq M$, and consider two events:

$$A = \{|\hat{\gamma} - \gamma| \leq \frac{\delta}{2M}\}$$

and A^c . Note that on A we must have $\hat{e}_i \geq \delta/2$. Then

$$\left| \frac{1}{n} \sum_{i=1}^n \frac{x_i y_i}{e_i \hat{e}_i} w_i^\top \right| \leq \frac{1}{n} \sum_{i=1}^n \frac{2|x_i y_i w_i|}{\delta e_i} + \frac{1}{n} \sum_{i=1}^n \frac{|x_i y_i w_i|}{e_i \hat{e}_i} \mathbf{1}_{A^c}$$

As long as the expectation $\mathbb{E}[|x_i y_i w_i|/e_i]$ is finite, then the first term would be $O_p(1)$. Note that the second term is $o_p(1)$.

(Or there is a better way to do this, illustrated in the asymptotic) Also

$$\sqrt{n}(\hat{\beta} - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{x_i y_i}{e_i} - \beta \right) + \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \left(\frac{1}{\hat{e}_i} - \frac{1}{e_i} \right)$$

Note that

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \left(\frac{1}{\hat{e}_i} - \frac{1}{e_i} \right) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \frac{1}{e_i \hat{e}_i} w_i^\top (\gamma - \hat{\gamma}) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \frac{1}{e_i^2} w_i^\top (\gamma - \hat{\gamma}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \left(\frac{1}{e_i \hat{e}_i} - \frac{1}{e_i^2} \right) w_i^\top (\gamma - \hat{\gamma}) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \frac{1}{e_i^2} w_i^\top (\gamma - \hat{\gamma}) - \frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \frac{1}{e_i^2 \hat{e}_i} w_i^\top (\gamma - \hat{\gamma}) w_i^\top (\gamma - \hat{\gamma}) \\ &= -\frac{1}{\sqrt{n}} \sum_{i=1}^n x_i y_i \frac{1}{e_i^2} w_i^\top (\hat{\gamma} - \gamma) + o_p(1) \end{aligned}$$

By the theory of linear regression, we know that

$$\sqrt{n}(\hat{\gamma} - \gamma) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbb{E}[w_i w_i^\top]^{-1} w_i v_i + o_p(1)$$

Hence

$$\sqrt{n}(\hat{\beta} - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(\frac{x_i y_i}{e_i} - \beta - \mathbb{E} \left[\frac{x_i y_i w_i^\top}{e_i^2} \right] \mathbb{E}[w_i w_i^\top]^{-1} w_i v_i \right) + o_p(1)$$

Asymptotic normality then follows from the central limit theorem.

4.4 Doubly robust

Regression adjustment gives the following estimator:

$$\hat{\mathbb{E}}[y_i(1)] = \frac{1}{n} \sum_{i=1}^n \hat{g}_1(\mathbf{w}_i)$$

Propensity score weighting method will give the estimator:

$$\hat{\mathbb{E}}[y_i(1)] = \frac{1}{n} \sum_{i=1}^n \frac{x_i y_i}{e(\mathbf{w}_i)}$$

Then we consider the combination of those two:

$$\hat{\mathbb{E}}[y_i(1)] = \frac{1}{n} \sum_{i=1}^n \left[\frac{x_i(y_i - \hat{g}_1(\mathbf{w}_i))}{\hat{e}(\mathbf{w}_i)} + \hat{g}_1(\mathbf{w}_i) \right] = \frac{1}{n} \sum_{i=1}^n \left[\frac{x_i y_i}{\hat{e}(\mathbf{w}_i)} + \frac{\hat{e}(\mathbf{w}_i) - x_i}{\hat{e}(\mathbf{w}_i)} \hat{g}_1(\mathbf{w}_i) \right]$$

Case 1: If \hat{g}_1 is correctly specified but the propensity score is mis-specified, i.e. $\hat{g}_1 \xrightarrow{p} g_1$ but $\hat{e} \xrightarrow{p} \tilde{e} \neq e$. Then

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left[\frac{x_i(y_i - \hat{g}_1(\mathbf{w}_i))}{\hat{e}(\mathbf{w}_i)} + \hat{g}_1(\mathbf{w}_i) \right] &\xrightarrow{p} \mathbb{E} \left[\frac{x_i(y_i - g_1(\mathbf{w}_i))}{\tilde{e}(\mathbf{w}_i)} + g_1(\mathbf{w}_i) \right] + o_p(1) \\ &= \mathbb{E} \left[\frac{x_i(y_i - g_1(\mathbf{w}_i))}{\tilde{e}(\mathbf{w}_i)} \right] + \mathbb{E}[y_i(1)] + o_p(1) \end{aligned}$$

Note that the first equality requires more assumptions. Also, the first term is zero, i.e.

$$\mathbb{E} \left[\frac{x_i(y_i - g_1(\mathbf{w}_i))}{\tilde{e}(\mathbf{w}_i)} \right] = \mathbb{E} \left[\mathbb{E} \left[\frac{x_i(y_i - g_1(\mathbf{w}_i))}{\tilde{e}(\mathbf{w}_i)} \middle| \mathbf{w}_i \right] \right] = 0$$

Hence it is consistent for ATE. On the other hand, if $\hat{g}_1 \xrightarrow{p} \tilde{g}_1 \neq g_1$ but $\hat{e} \rightarrow e$, then

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left[\frac{x_i y_i}{\hat{e}(\mathbf{w}_i)} + \frac{\hat{e}(\mathbf{w}_i) - x_i}{\hat{e}(\mathbf{w}_i)} \hat{g}_1(\mathbf{w}_i) \right] &\xrightarrow{p} \mathbb{E} \left[\frac{x_i y_i}{e(\mathbf{w}_i)} + \frac{e(\mathbf{w}_i) - x_i}{e(\mathbf{w}_i)} \tilde{g}_1(\mathbf{w}_i) \right] + o_p(1) \\ &= \mathbb{E} \left[\frac{x_i y_i}{e(\mathbf{w}_i)} \right] + \mathbb{E} \left[\frac{e(\mathbf{w}_i) - x_i}{e(\mathbf{w}_i)} \tilde{g}_1(\mathbf{w}_i) \right] + o_p(1) \end{aligned}$$

Also, applying law of iterated expectation will give us

$$\mathbb{E} \left[\frac{e(\mathbf{w}_i) - x_i}{e(\mathbf{w}_i)} \tilde{g}_1(\mathbf{w}_i) \right] = 0$$

Hence it is doubly robust.

5 Instrumental Variables