

Introductory Econometrics

Lecture 26: Treatment Effect Model

Instructor: Ma, Jun

Renmin University of China

December 16, 2021

Introduction

- ▶ In this class, we consider the problem of estimating the causal effect of a binary explanatory variable, which is referred as the treatment effect in the literature. The treatment effect model is different from the linear regression model.
- ▶ In econometrics, the treatment effect model is very often used for evaluating social program/experiment.
- ▶ Example 1: Suppose that a selected set of individuals receive training or education initiated by the government with a view to enhancing their employment prospects. Suppose that the government has collected the earnings data for the individuals who received the training and for the individuals who did not. The main purpose of methods of program evaluations is to quantify and estimate the effect of the training program.

- ▶ Example 2: Suppose that an education program required high schools to agree to assign teachers and students to small (13 to 17 students) or large (22 to 26 students) classes. The government is interested in the effect of class size on student achievement.
- ▶ Such a question can arise in various other situations. A medical experiment studies on the effects of new treatment ask similar questions. One group of patients has received new treatment, and the other group has not.

Potential outcome variables

- ▶ Y_i : outcome variable; $D_i \in \{0, 1\}$: the binary explanatory variable; X_{i1}, \dots, X_{ik} : other observed explanatory variables; ϵ_i : unobserved explanatory factors.
- ▶ The variable D_i is a binary variable taking 1 if the individual has gone through the treatment and 0 otherwise. The treatment here represents the actual treatment. The econometrician usually observes the treatment status for each individual D_i .
- ▶ (X_{i1}, \dots, X_{ik}) represents a vector of various demographic characteristics for individual i . E.g., the variables can be annual income, age, gender, status of marriage, the number of children, education, etc. These represent all the observable characteristics of individual i .
- ▶ Suppose that Y_i is generated by $Y_i = g(D_i, X_{i1}, \dots, X_{ik}, \epsilon_i)$.
- ▶ g is unknown and in the treatment effect model, we do not assume g is linear.

- ▶ The outcome variable $Y_i(1) = g(1, X_{i1}, \dots, X_{ik}, \epsilon_i)$ represents a potential outcome of an individual i in the treatment state (e.g. training is received or studying in a reduced-size class). The variable $Y_i(0) = g(0, X_{i1}, \dots, X_{ik}, \epsilon_i)$ represents a potential outcome of the same individual i in the control state (e.g. training is received or studying in a normal-size class).
- ▶ Thus, each individual has a random vector $(Y_i(1), Y_i(0))$ that represents potential outcomes depending on the state (treatment or control). Certainly, $(Y_i(1), Y_i(0))$ are correlated.
- ▶ The econometrician cannot observe the random vector $(Y_i(1), Y_i(0))$ jointly, because for each individual, only one potential outcome ($Y_i(1)$ or $Y_i(0)$) is realized, depending on whether the individual i has gone through the treatment or not.

The relationship between D_i and $(Y_i(1), Y_i(0))$

- ▶ In a medical experiment, the individual is chosen to be in the treatment group through some randomization device or a lottery. In these cases, $D_i \perp\!\!\!\perp (Y_i(1), Y_i(0))$ (i.e., D_i is independent of $(Y_i(1), Y_i(0))$).
- ▶ For evaluating social experiment/program with observational data, it may not be convincing to assume $D_i \perp\!\!\!\perp (Y_i(1), Y_i(0))$.

Treatment effects

- ▶ The individual treatment effect (ITE) for each individual i is defined as:

$$Y_i(1) - Y_i(0).$$

- ▶ The ITE is the difference between the potential outcomes in two different states for the same person.
- ▶ The ITE is a counterfactual quantity, in the sense that in the actual world, we cannot observe the vector $(Y_i(1), Y_i(0))$.
- ▶ There are mainly two quantities of interest: ATE (average treatment effect)

$$\text{ATE} = E[Y_i(1) - Y_i(0)],$$

and ATT (average treatment effect on the treated)

$$\text{ATT} = E[Y_i(1) - Y_i(0) \mid D_i = 1].$$

- ▶ The average treatment effect on the treated is the treatment effect of the people who have gone through the treatment.

- ▶ Note that the expectation in the definition of ATE involves the joint distribution of $(Y_i(1), Y_i(0))$, and the expectation in the definition of ATT involves the joint distribution of $(Y_i(1), Y_i(0), D_i)$, which are both unobserved.
- ▶ ATE or ATT can not be estimated accurately merely by collecting a large size of samples.

The observed information

- ▶ The econometrician observes the treatment status D_i and covariates X_i . She also observes the outcome variable:

$$Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0).$$

- ▶ The observed outcome variable Y_i is not the same as the potential outcomes $Y_i(1)$ or $Y_i(0)$. It is a realized outcome for an individual i depending on whether she has received treatment (Y_i is realized to be $Y_i(1)$) or not (Y_i is realized to be $Y_i(0)$).
- ▶ Identification of these parameters is concerned with the following question: can we uniquely determine the value of these parameters once we know the joint distribution of the observable random variables?

Randomized experiments

- ▶ In medical experiments, the treatment is performed using a randomization device. More specifically, for patient i , a lottery is run, and the patient is selected into the treated group with the design probability p , and stays in the control group with the design probability $1 - p$.
- ▶ In these cases, we have $D_i \perp (Y_i(1), Y_i(0), X_{i1}, \dots, X_{ik})$. Randomized experiment assumption requires that knowing whether patient i is treated or not gives one no informational advantage in predicting the potential outcomes of i over another who does not know whether patient i is treated or not.
- ▶ This assumption is still possibly violated in medical studies if only those patients who have higher potential treatment effect are selected into treatment among all the patients in the study on purpose.
- ▶ In this case, observing D_i will give information about the treatment effect ($Y_i(1) - Y_i(0)$) for individual i .

- ▶ We use the following result from probability theory: if $V \perp\!\!\!\perp W$, then for any function f ,

$$E[f(V, W) \mid W = w] = E[f(V, w)]. \quad (1)$$

- ▶ By (1) and the randomized experiment assumption, $D_i \perp\!\!\!\perp (Y_i(1), Y_i(0))$, we have

$$\begin{aligned} \text{ATE} &= E[Y_i(1) - Y_i(0)] \\ &= E[Y_i(1)] - E[Y_i(0)] \\ &= E[D_i Y_i(1) + (1 - D_i) Y_i(0) \mid D_i = 1] \\ &\quad - E[D_i Y_i(1) + (1 - D_i) Y_i(0) \mid D_i = 0] \\ &= E[Y_i \mid D_i = 1] - E[Y_i \mid D_i = 0]. \end{aligned}$$

- By LIE,

$$\begin{aligned} E[Y_i D_i] &= E[E[Y_i D_i | D_i]] \\ &= \Pr[D_i = 1] E[Y_i D_i | D_i = 1] \\ &\quad + \Pr[D_i = 0] E[Y_i D_i | D_i = 0] \\ &= E[D_i] E[Y_i | D_i = 1], \end{aligned}$$

where

$$\begin{aligned} E[Y_i D_i | D_i = 0] &= E[(D_i Y_i(1) + (1 - D_i) Y_i(0)) D_i | D_i = 0] \\ &= 0 \end{aligned}$$

follows from (1).

- Similarly, we have

$$E[Y_i | D_i = 0] = \frac{E[Y_i (1 - D_i)]}{E[1 - D_i]}.$$

- ▶ We can write

$$\text{ATE} = \frac{E[Y_i D_i]}{E[D_i]} - \frac{E[Y_i (1 - D_i)]}{E[1 - D_i]},$$

where the right hand side depends on the joint distribution of the observed random variables.

- ▶ For estimation, we replace the population mean by the sample mean (this is sometimes called the analogue principle):

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

- ▶ We can check its consistency by using LLN and Slutsky's lemma.
- ▶ This randomization assumption is not convincing when the individuals in the social experiments are people who may select into the treatment or not.

Comparison with the linear regression

- ▶ It seems that D_i is nothing but a dummy variable. Can we run a regression of Y_i on D_i and X_{i1}, \dots, X_{ik} to estimate the ATE? Can the parameter of interest, the ATE, be formulated as a coefficient in a regression model.
- ▶ One possible assumption is that

$$Y_i = g(D_i, X_{i1}, \dots, X_{ik}, \epsilon_i) = \gamma_0 + \gamma_1 D_i + \sum_{j=1}^k \beta_j X_{ij} + \epsilon_i.$$

In this case, the ITE $Y_i(1) - Y_i(0) = \gamma_1$ is constant. This is very unrealistic. We investigate alternative model assumptions.

- ▶ We first consider the following model assumption

$$Y_i(0) = \mu_0 + U_i(0)$$

$$Y_i(1) = \mu_1 + U_i(1),$$

where μ_0 and μ_1 are constants common across individuals and assumed to be nonstochastic and $(U_i(0), U_i(1))$ are stochastic components.

- ▶ We denote $\mathbf{X}_i = (X_{i1}, \dots, X_{ik})^\top$ for the vector of observed covariates.
- ▶ We assume $E[U_i(0) | \mathbf{X}_i] = E[U_i(1) | \mathbf{X}_i]$, which implies

$$E[Y_i(1) - Y_i(0) | \mathbf{X}_i] = \mu_1 - \mu_0,$$

i.e., the ITE is mean independent of \mathbf{X}_i but it can be random.
And by LIE,

$$ATE = E[Y_i(1) - Y_i(0)] = \mu_1 - \mu_0.$$

- ▶ We assume $E[Y_i(1) | D_i, \mathbf{X}_i] = E[Y_i(1) | \mathbf{X}_i]$ and $E[Y_i(0) | D_i, \mathbf{X}_i] = E[Y_i(0) | \mathbf{X}_i]$, i.e., the conditional mean independence of potential outcomes with treatment status, conditional on demographic status \mathbf{X}_i .
- ▶ When we focus on a sub-population of individuals with specific demographic status \mathbf{X}_i , $Y_i(1)$ and $Y_i(0)$ are both mean independent of D_i .

- ▶ Let us write

$$\begin{aligned} E[Y_i | D_i, \mathbf{X}_i] &= D_i E[Y_i(1) | D_i, \mathbf{X}_i] + (1 - D_i) E[Y_i(0) | D_i, \mathbf{X}_i] \\ &= D_i E[Y_i(1) - Y_i(0) | D_i, \mathbf{X}_i] + E[Y_i(0) | D_i, \mathbf{X}_i] \\ &= D_i E[Y_i(1) - Y_i(0) | \mathbf{X}_i] + E[Y_i(0) | \mathbf{X}_i], \end{aligned}$$

where the last equality follows from the conditional mean independence assumption.

- ▶ By the assumption $E[U_i(0) | \mathbf{X}_i] = E[U_i(1) | \mathbf{X}_i]$, we have

$$\begin{aligned} D_i E[Y_i(1) - Y_i(0) | \mathbf{X}_i] + E[Y_i(0) | \mathbf{X}_i] \\ &= D_i (\mu_1 - \mu_0) + E[Y_i(0) | \mathbf{X}_i] \\ &= \mu_0 + D_i (\mu_1 - \mu_0) + g(X_{i1}, \dots, X_{ik}), \end{aligned}$$

where we denote $g(X_{i1}, \dots, X_{ik}) = E[U_i(0) | \mathbf{X}_i]$.

- ▶ Therefore, we have

$$E[Y_i | D_i, \mathbf{X}_i] = \mu_0 + (\mu_1 - \mu_0) D_i + g(X_{i1}, \dots, X_{ik}).$$

- ▶ Define

$$V_i = Y_i - E[Y_i | D_i, \mathbf{X}_i]$$

and now we have the following regression model:

$$Y_i = \mu_0 + (\mu_1 - \mu_0) D_i + g(X_{i1}, \dots, X_{ik}) + V_i.$$

- ▶ We have $E[V_i | D_i, \mathbf{X}_i] = 0$ by definition.
- ▶ We assume g is linear in X_{i1}, \dots, X_{ik} :

$$g(X_{i1}, \dots, X_{ik}) = \sum_{j=1}^k \beta_j X_{ij},$$

and then

$$Y_i = \mu_0 + (\mu_1 - \mu_0) D_i + \sum_{j=1}^k \beta_j X_{ij} + V_i.$$

- ▶ A multiple linear regression of Y_i on D_i and X_{i1}, \dots, X_{ik} consistently estimates $ATE = (\mu_1 - \mu_0)$.

- ▶ We assume $E[U_i(0) | \mathbf{X}_i] = E[U_i(1) | \mathbf{X}_i]$, which implies

$$E[Y_i(1) - Y_i(0) | \mathbf{X}_i] = \mu_1 - \mu_0.$$

- ▶ This assumption implies that the conditional average treatment effect given \mathbf{X}_i does not depend on \mathbf{X}_i , the characteristics of individual i .
- ▶ This assumption can be unrealistic. E.g., Average treatment of the class-size is the same between students from high-income family and students from low-income family.

Unconfoundedness assumption

- ▶ Unconfoundedness is the key assumption of the basic treatment effect model.
- ▶ Unconfoundedness assumption: $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i \mid X_i$, i.e., $(Y_i(1), Y_i(0))$ and D_i are conditionally independent given X_i .
- ▶ Unconfoundedness can be thought of as an assumption that the decision to take the treatment is purely random for individuals with similar values of the covariates.
- ▶ Suppose that we have three random vectors V , W and X , where (V, W) is a continuous random vector. Then we say V and W are conditionally independent given X , if for all possible values of v , w and x ,

$$f_{(V,W)|X}(v, w \mid x) = f_{V|X}(v \mid x) f_{W|X}(w \mid x).$$

- Unconfoundedness is satisfied if (Y_i, D_i) are generated by the model

$$\begin{aligned} Y_i &= g(D_i, X_{i1}, \dots, X_{ik}, \epsilon_i) \\ D_i &= m(X_{i1}, \dots, X_{ik}, \eta_i) \end{aligned}$$

and $\epsilon_i \perp\!\!\!\perp \eta_i \mid X_{i1}, \dots, X_{ik}$.

More on conditional independence

- ▶ When V and W are conditionally independent given X , one can easily see that for any function φ ,

$$E[\varphi(V) | W, X] = E[\varphi(V) | X].$$

I.e., once we observe X , knowledge of W does not give us any further advantage in predicting the value of $\varphi(V)$.

- ▶ We notice that

$$\begin{aligned} f_{(V,W)|X}(v, w | x) &= \frac{f_{(V,W,X)}(v, w, x)}{f_X(x)} \\ &= \frac{f_{(V,W,X)}(v, w, x)}{f_{(W,X)}(w, x)} \frac{f_{(W,X)}(w, x)}{f_X(x)} \\ &= f_{V|(W,X)}(v | w, x) f_{W|X}(w, x). \end{aligned}$$

- ▶ Therefore, we have $f_{V|X}(v|x) = f_{V|(W,X)}(v|w,x)$, if (V, W) are conditionally independent given X . Hence,

$$\begin{aligned} E[\varphi(V) | W = w, X = x] &= \int \varphi(v) f_{V|(W,X)}(v|w,x) dv \\ &= \int \varphi(v) f_{V|X}(v|x) dv \\ &= E[\varphi(V) | X = x]. \end{aligned}$$

- ▶ Therefore, the unconfoundedness assumption $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i | X_i$ implies the conditional mean independence assumption:

$$\begin{aligned} E[Y_i(1) | D_i, X_i] &= E[Y_i(1) | X_i] \\ E[Y_i(0) | D_i, X_i] &= E[Y_i(0) | X_i]. \end{aligned}$$

- ▶ We can also show: if $V \perp\!\!\!\perp W | X$,

$$E[\eta(V, W) | X, W = w] = E[\eta(V, w) | X]. \quad (2)$$

The unconfoundedness and randomization assumptions

- ▶ It can be shown that the randomization assumption $(Y_i(1), Y_i(0), \mathbf{X}_i) \perp\!\!\!\perp D_i$ implies the unconfoundedness assumption $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i \mid \mathbf{X}_i$.
- ▶ The randomized experiment assumption does not allow X_{i1}, \dots, X_{ik} to be correlated with D_i ,
- ▶ The unconfounded condition allows D_i to be affected by X_{i1}, \dots, X_{ik} , while the randomized experiment assumption does not.

Identification of ATE

- By LIE, we have

$$\begin{aligned}\text{ATE} &= E[Y_i(1) - Y_i(0)] \\ &= E[E[Y_i(1) | \mathbf{X}_i]] - E[E[Y_i(0) | \mathbf{X}_i]],\end{aligned}\quad (3)$$

and

$$\begin{aligned}E[Y_i D_i | \mathbf{X}_i] &= E[E[Y_i D_i | \mathbf{X}_i, D_i] | \mathbf{X}_i] \\ &= \Pr[D_i = 1 | \mathbf{X}_i] E[Y_i D_i | \mathbf{X}_i, D_i = 1] \\ &\quad + \Pr[D_i = 0 | \mathbf{X}_i] E[Y_i D_i | \mathbf{X}_i, D_i = 0].\end{aligned}$$

- By the unconfoundedness assumption: $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i \mid \mathbf{X}_i$, the result (2) and the relation $Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0)$, we have

$$\begin{aligned} & \mathbb{E}[Y_i D_i \mid \mathbf{X}_i, D_i = 1] \\ &= \mathbb{E}[(D_i Y_i(1) + (1 - D_i) Y_i(0)) D_i \mid \mathbf{X}_i, D_i = 1] = \mathbb{E}[Y_i(1) \mid \mathbf{X}_i] \end{aligned}$$

and

$$\mathbb{E}[Y_i D_i \mid \mathbf{X}_i, D_i = 0] = 0.$$

- Therefore, we have

$$\mathbb{E}[Y_i D_i \mid \mathbf{X}_i] = \Pr[D_i = 1 \mid \mathbf{X}_i] \mathbb{E}[Y_i(1) \mid \mathbf{X}_i] \quad (4)$$

and similarly,

$$\mathbb{E}[Y_i(1 - D_i) \mid \mathbf{X}_i] = \Pr[D_i = 0 \mid \mathbf{X}_i] \mathbb{E}[Y_i(0) \mid \mathbf{X}_i]. \quad (5)$$

- ▶ Now (3), (4), (5) and LIE imply

$$\begin{aligned} \text{ATE} &= \mathbb{E} \left[\frac{\mathbb{E}[Y_i D_i \mid \mathbf{X}_i]}{\Pr[D_i = 1 \mid \mathbf{X}_i]} \right] - \mathbb{E} \left[\frac{\mathbb{E}[Y_i (1 - D_i) \mid \mathbf{X}_i]}{\Pr[D_i = 0 \mid \mathbf{X}_i]} \right] \\ &= \mathbb{E} \left[\frac{Y_i D_i}{\Pr[D_i = 1 \mid \mathbf{X}_i]} - \frac{Y_i (1 - D_i)}{\Pr[D_i = 0 \mid \mathbf{X}_i]} \right]. \end{aligned}$$

Now the right hand side depends only on the joint distribution of observed random variables.

- ▶ Denote

$$p(\mathbf{x}) = \Pr[D_i = 1 \mid \mathbf{X}_i = \mathbf{x}].$$

This function is called propensity score. It is the probability of the event that the individual belongs to the treatment group, given that the observed characteristics are $\mathbf{x} \in \mathbb{R}^k$.

Estimation of ATE

- ▶ Let $\hat{p}(\mathbf{x})$ be an estimator of the propensity score, then we can estimate the ATE:

$$\widehat{\text{ATE}} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{Y_i D_i}{\hat{p}(\mathbf{X}_i)} - \frac{Y_i (1 - D_i)}{1 - \hat{p}(\mathbf{X}_i)} \right\}.$$

- ▶ It is straightforward to construct $\hat{p}(\mathbf{x})$ if \mathbf{X}_i is discrete:

$$\hat{p}(\mathbf{x}) = \frac{\sum_{i=1}^n 1(D_i = 1, \mathbf{X}_i = \mathbf{x})}{\sum_{i=1}^n 1(\mathbf{X}_i = \mathbf{x})}.$$

- ▶ If \mathbf{X}_i is continuous, we specify a parametric model for the propensity score:

$$\Pr[D_i = 1 \mid \mathbf{X}_i] = \Phi(\beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik})$$

as what we did for the Probit model. This gives a parametric model for the propensity score. $(\beta_0, \dots, \beta_k)$ can be estimated by MLE (denoted by $(\hat{\beta}_0, \dots, \hat{\beta}_k)$).

- ▶ The estimated propensity score is

$$\hat{p}(\mathbf{X}_i) = \Phi(\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_k X_{ik}).$$

- ▶ This estimator is known to be consistent and asymptotically normally distributed, if our propensity score model is correct.
- ▶ This approach has the drawback that if our model for the propensity score is wrong, the ATE estimator is inconsistent.
- ▶ Actually, $p(\mathbf{x}) = E[D_i | \mathbf{X}_i = \mathbf{x}]$ can be estimated without specifying a parametric model for it.

k-NN estimator

- ▶ The k -nearest neighbor (k -NN) estimator is the simplest nonparametric estimator of $p(\mathbf{x})$.
- ▶ Fix \mathbf{x}_0 and suppose that we want to estimate $p(\mathbf{x}_0)$ at this point. Assume that p is a smooth function, which means that its graph does not change too much.
- ▶ $p(\mathbf{x})$ should be close to $p(\mathbf{x}_0)$ when \mathbf{x} is close enough to \mathbf{x}_0 . $p(\mathbf{X}_i)$ would be close to $p(\mathbf{x}_0)$ for observations \mathbf{X}_i close to \mathbf{x}_0 .
- ▶ We simply average these $p(\mathbf{X}_i)$ for observations \mathbf{X}_i close to \mathbf{x}_0 . We do not observe $p(\mathbf{X}_i)$ but use D_i instead.
- ▶ Let

$$d_i(\mathbf{x}_0) = \|\mathbf{X}_i - \mathbf{x}_0\| = \sqrt{(\mathbf{X}_i - \mathbf{x}_0)^\top (\mathbf{X}_i - \mathbf{x}_0)}$$

denote the distance of \mathbf{X}_i to \mathbf{x}_0 .

- ▶ After computing the distance for all n observations in the sample, we sort them in the increasing order

$$d_{(1)}(\mathbf{x}_0) \leq d_{(2)}(\mathbf{x}_0) \leq \cdots \leq d_{(n)}(\mathbf{x}_0).$$

- ▶ Let $N_k(\mathbf{x}_0)$ denote the identities of the k -nearest neighbors of \mathbf{x}_0 :

$$N_k(\mathbf{x}_0) = \{i : d_i(\mathbf{x}_0) \leq d_{(k)}(\mathbf{x}_0)\}.$$

- ▶ The k -NN nonparametric estimator of $p(\mathbf{x}_0)$ is

$$\hat{p}_{kNN}(\mathbf{x}_0) = \frac{1}{k} \sum_{i \in N_k(\mathbf{x}_0)} D_i.$$

- ▶ The k -NN estimator is simply an average of the values of D_i across the k closest observations in terms of \mathbf{X}_i .
- ▶ There is a data-driven procedure to select k in practical applications.
- ▶ The nonparametric ATE estimator using $\hat{p}_{kNN}(\mathbf{X}_i)$ is consistent and asymptotically normal. It does not require a parametric model for the propensity score.