

Part B: Covariate Adjustment. Selection on Observables

Summary of Main Ideas

Kirill Borusyak

ARE 213 Applied Econometrics

UC Berkeley, Fall 2025

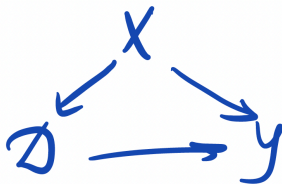
Outline

- 1 The concept of control variables
- 2 Regression adjustment
- 3 Propensity score methods
- 4 Doubly-robust methods

What if treatment is not randomly assigned?

One basic approach for causal inference with observational data is to control for observables

Unconfoundedness:



Violation of unconfoundedness:

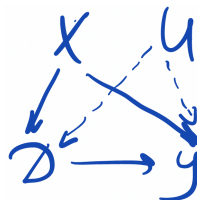


Illustration: Returns to College Degree

i	X_i (urban)	D_i (college)	Y (earnings)	Y_0	Y_1
1	1	1	51	?	51
2	1	1	50	?	50
3	1	1	49	?	49
4	1	0	40	40	?
5	0	1	15	?	15
6	0	0	11	11	?
7	0	0	9	9	?
<i>Treated</i> minus control (not causal!)			21.25		

Identification assumptions

- Unconfoundedness = Ignorability = Conditional independence assumption (**CIA**) = Selection on observables: $(Y_i(0), Y_i(1)) \perp\!\!\!\perp D_i \mid X_i$
 - ▶ Sometimes viewed as *defining* that X_i is a control variable
- **Overlap**: $0 < P(D_i = 1 \mid X_i) < 1$ on the support of X_i
 - ▶ $P(D_i = 1 \mid X_i)$ is called the **propensity score**

Identification under CIA + overlap

- **Conditional average treatment effect:** $CATE(x) \equiv \mathbb{E}[Y_i(1) - Y_i(0) \mid X_i = x]$
 - ▶ Let $h_1(x) = \mathbb{E}[Y \mid D = 1, X = x]$ and $h_0(x) = \mathbb{E}[Y \mid D = 0, X = x]$
 - ▶ CATE can be identified as $CATE(x) = h_1(x) - h_0(x)$ (*prove this!*)
- $ATE = \mathbb{E}[CATE(X_i)] = \mathbb{E}[h_1(X_i) - h_0(X_i)]$, with expectation taken across X_i (in both treatment and control groups)
- And $ATT = \mathbb{E}[CATE(X_i) \mid D_i = 1]$, with expectation taken across X_i in the treated group

Imputation representation

- We also have

$$ATT = \mathbb{E} [Y_i - h_0(X_i) \mid D_i = 1]$$

- ▶ For treated i , we impute $\widehat{Y_i(0)} = h_0(X_i)$ and

- And

$$ATE = \mathbb{E} [(Y_i - h_0(X_i)) D_i + (h_1(X) - Y_i)(1 - D_i)]$$

- ▶ For untreated i , also impute $\widehat{Y_i(1)} = h_1(X_i)$

Matching / Imputation

i	X_i (urban)	D_i (college)	Y (earnings)	\hat{Y}_0	\hat{Y}_1	$\hat{Y}_1 - \hat{Y}_0$
1	1	1	51	40	51	11
2	1	1	50	40	50	10
3	1	1	49	40	49	9
4	1	0	40	40	$\frac{51+50+49}{3}$	10
<i>Urban averages & CATE</i>				40	50	10
5	0	1	15	$\frac{11+9}{2}$	15	5
6	0	0	11	11	15	4
7	0	0	9	9	15	6
<i>Rural averages and CATE</i>				10	15	5
<i>Overall averages and ATE:</i>				27.1	35	7.9

Regression adjustment

Regress Y_i on X_i in the subsample $D_i = 1$; take fitted values as $\hat{h}_1(X_i)$; repeat for $D_i = 0$

i	X_i (urban)	D_i (college)	Y (earnings)	\hat{Y}_0	\hat{Y}_1	$\hat{Y}_1 - \hat{Y}_0$
1	1	1	51	40	50	10
2	1	1	50	40	50	10
3	1	1	49	40	50	10
4	1	0	40	40	50	10
5	0	1	15	10	15	5
6	0	0	11	10	15	5
7	0	0	9	10	15	5
Overall averages and ATE:				27.1	35	7.9
Treatment group averages and ATT:				32.5	41.3	8.8

Pscore reweighting (for ATE)

i	X_i (urban)	D_i (college)	Y (earnings)	\hat{Y}_0	\hat{Y}_1	$\hat{Y}_1 - \hat{Y}_0$
1	1	1	51	0	$51 \cdot \frac{1}{3/4}$	$51 \cdot \frac{1}{3/4}$
2	1	1	50	0	$50 \cdot \frac{1}{3/4}$	$50 \cdot \frac{1}{3/4}$
3	1	1	49	0	$49 \cdot \frac{1}{3/4}$	$49 \cdot \frac{1}{3/4}$
4	1	0	40	$40 \cdot \frac{1}{1/4}$	0	$-40 \cdot \frac{1}{1/4}$
5	0	1	15	0	$15 \cdot \frac{1}{1/3}$	$15 \cdot \frac{1}{1/3}$
6	0	0	11	$11 \cdot \frac{1}{2/3}$	0	$-11 \cdot \frac{1}{2/3}$
7	0	0	9	$9 \cdot \frac{1}{2/3}$	0	$-9 \cdot \frac{1}{2/3}$
<i>Overall averages and ATE:</i>				27.1	35	7.9

Continuous covariates

- With continuous or high-dimensional covariates, can't do any of this
 - ▶ Approximate matching (*skipped*)
 - ▶ Regression adjustment
 - ▶ Propensity score methods
 - ▶ Doubly-robust methods; double machine learning
 - ▶ *Note:* Different methods are no longer equivalent
- Also, do simpler strategies, like regressing Y_i on D_i and X_i , recover anything useful?

Outline

- 1 The concept of control variables
- 2 Regression adjustment
- 3 Propensity score methods
- 4 Doubly-robust methods

Regression adjustment

- Recall that, under CIA, we have for $d = 0, 1$

$$\mathbb{E}[Y(d) \mid X] = \mathbb{E}[Y(d) \mid D = d, X] = \mathbb{E}[Y \mid D = d, X] \equiv h_d(X)$$

- If we estimate $h_0(\cdot)$ and $h_1(\cdot)$, we get

$$\widehat{ATE} = \frac{1}{N} \sum_i \left(\hat{h}_1(X_i) - \hat{h}_0(X_i) \right), \quad \widehat{ATT} = \frac{1}{N_1} \sum_i \left(\hat{h}_1(X_i) - \hat{h}_0(X_i) \right) D_i$$

- or via imputation:

$$\widehat{ATT} = \frac{1}{N_1} \sum_i \left(Y_i - \hat{h}_0(X_i) \right) D_i$$

$$\widehat{ATE} = \frac{1}{N} \sum_i \left\{ \left(Y_i - \hat{h}_0(X_i) \right) D_i + \left(\hat{h}_1(X_i) - Y_i \right) (1 - D_i) \right\}$$

How to estimate $h_0(\cdot), h_1(\cdot)$?

- $h_d(\cdot)$ is the CEF of Y_i given X_i for $D_i = d$
- Can use nonparametric regression, e.g. local linear regression
 - ▶ For each x , estimate $h_d(x)$ by an intercept from a regression of Y_i on $(X_i - x)$, keeping observations in the neighborhood of x (and with $D_i = d$) only
- If $\mathbb{E}[Y(d) | X] = \gamma'_d X$ is linear in X (e.g. X is saturated): Oaxaca-Blinder estimator
 - ▶ Run linear regressions of Y on X within treated/control groups separately
 - ▶ Or a single fully-interacted regression (for X_i including an intercept)

$$Y_i = \gamma_0 + \gamma'_1 X_i + \beta D_i + \tau' X_i D_i + \text{error}_i, \quad \widehat{ATE} = \hat{\beta} + \hat{\tau}' \bar{X}$$

- ▶ Or its convenient reformulation

$$Y_i = \gamma_0 + \gamma'_1 X_i + \beta D_i + \tau' (X_i - \bar{X}) D_i + \text{error}_i, \quad \widehat{ATE} = \hat{\beta}$$

- ▶ Note: interactions are helpful even if you are not interested in effect heterogeneity

Uninteracted regression

What about regression of Y_i on D_i and X_i without interactions?

- Suppose causal effects are homogeneous, $Y_i = \beta D_i + Y_i(0)$
- And $\mathbb{E}[Y(0) \mid X] = \gamma'X$ is linear,

$$Y_i = \beta D_i + \gamma'X_i + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i \mid D_i, X_i] = 0$$

where $h_0(X) = \gamma'X$ and $h_1(X) = \beta + \gamma'X$

⇒ Regression that linearly controls for X_i identifies the causal effect

- But Mostly Harmless Econometrics advocates for uninteracted regressions even when the effects are heterogeneous. Why?

Uninteracted regression with heterogeneous effects

- Assume the propensity score $p(X_i) \equiv \mathbb{E}[D_i | X_i] = P(D_i = 1 | X_i)$ is linear in X_i
 - ▶ Angrist (1998) focused on saturated controls $X_i \implies$ trivially satisfied

- Then

$$\beta_{OLS} = \frac{\mathbb{E}[CATE(X_i) \cdot \omega(X_i)]}{\mathbb{E}[\omega(X_i)]}, \quad \omega(X_i) = \text{Var}[D_i | X_i] = p(X_i)(1 - p(X_i))$$

- ▶ Groups with $p(X_i) \approx 1/2$ get the most weight (relative to their size)
 - ▶ Groups where overlap is limited ($p(X_i) \approx 0$ or $p(X_i) \approx 1$) get little weight
 - ▶ $\beta_{OLS} = ATE$ if $CATE(X_i)$ is constant, $\omega(X_i)$ is constant, or they are uncorrelated with each other
- *Note:* linearity of $\mathbb{E}[Y_i(0) | X_i]$ is not needed

Variance weighting: Proof

- By linearity of the pscore, partialling out X_i from D_i yields residuals $\tilde{D}_i = D_i - \mathbb{E}[D_i | X_i]$
- By Frisch-Waugh-Lovell, $\beta_{OLS} = \mathbb{E}[\tilde{D}_i Y_i] / \mathbb{E}[\tilde{D}_i D_i]$ (where $\mathbb{E}[\tilde{D}_i D_i] = \text{Var}[\tilde{D}_i]$)
- Using CIA and $\mathbb{E}[\tilde{D}_i | X_i] = 0$,

$$\begin{aligned}\mathbb{E}[\tilde{D}_i Y_i] &= \mathbb{E}\left[\mathbb{E}\left[\tilde{D}_i (Y_i(0) + (Y_i(1) - Y_i(0)) D_i) \mid X_i\right]\right] \\ &= \mathbb{E}\left[CATE(X_i) \cdot \mathbb{E}[\tilde{D}_i D_i \mid X_i]\right] = \mathbb{E}[CATE(X_i) \cdot \text{Var}[D_i \mid X_i]]\end{aligned}$$

- Analogously, $\mathbb{E}[\tilde{D}_i D_i] = \mathbb{E}[\text{Var}[D_i \mid X_i]]$

Outline

- 1 The concept of control variables
- 2 Regression adjustment
- 3 Propensity score methods**
- 4 Doubly-robust methods

Propensity score theorems (Rosenbaum and Rubin, 1983)

- Consider binary D . Recall $p(X) \equiv P(D = 1 \mid X)$
- **Proposition 1:** $D \perp\!\!\!\perp X \mid p(X)$
 - ▶ i.e. propensity score balances X between treated and control groups
 - ▶ Proof: $P(D = 1 \mid X, p(X)) = P(D = 1 \mid X) = p(X)$
- **Proposition 2:** $D \perp\!\!\!\perp (Y_0, Y_1) \mid X \implies D \perp\!\!\!\perp (Y_0, Y_1) \mid p(X)$
 - ▶ i.e., under CIA, controlling for scalar $p(X)$ is enough
 - ★ A stronger version of the OVB idea
 - ▶ Proof: $P(D = 1 \mid p(X), Y_0, Y_1) = \mathbb{E} [\mathbb{E} [D \mid p(X), X, Y_0, Y_1] \mid p(X), Y_0, Y_1] = p(X)$, doesn't depend on (Y_0, Y_1)

P-score methods: Steps

1. Obtain/estimate $p(X)$
 - ▶ Known in complex RCTs
 - ▶ Parametric estimation, e.g. logit of D on X
 - ▶ Non-parametric regression of D on X
2. Verify balance
 - ▶ Within bins of $\hat{p}(X)$ compare X among treated and controls
 - ▶ If balance fails (with sufficiently many bins), make the p-score model richer
3. Assess overlap
 - ▶ Compare p-score distributions in treated & control groups
4. **Adjust for pscore differences** between treated and control groups
 - ▶ Regression, matching, blocking, reweighting

P-score adjustment methods: Regression

- With constant effects, enough to control linearly

$$Y_i = \beta D_i + \gamma p(X_i) + \text{error}$$

- ▶ *Exercise:* Why?
 - ▶ *Exercise:* if $p(X)$ is estimated from a linear probability model, both ways are numerically the same as linearly controlling for X
- With heterogeneous effects, this yields the variance-weighted average of effects (*Exercise:* Why?)

P-score adjustment methods: Blocking/Matching

- **Blocking** (stratifying) = exact matching on binned pscore
 - ▶ Split data into bins of $p(X_i)$
 - ▶ Estimate difference-in-means within bins
 - ▶ Average across bins weighting by $\#$ obs. (ATE) or $\#$ treated obs. (ATT)
- Approximate **matching**: For each treated obs., compare the outcome with the untreated one with the closest $p(X_i)$; and the other way round
 - ▶ Discard observations with the pscore outside the range for the other group, such that the nearest match is very far

P-score adjustment methods: Reweighting (IPW)

- In the bin with $p(X_i) = \pi$ we have fraction π of observed $Y_i(1)$ (for treated) and fraction $1 - \pi$ of comparable but missing $Y_i(1)$ (for controls)
- Horvitz-Thompson (1952) “**inverse probability weighting**” (IPW): reweighting by $1/\pi$ makes the sample of $Y_i(1)$ representative

$$\mathbb{E} \left[\frac{YD}{p(X)} \right] = \mathbb{E} \left[\frac{Y_1 D}{p(X)} \right] = \mathbb{E} \left[\mathbb{E} \left[\frac{Y_1 D}{p(X)} \mid X \right] \right] = \mathbb{E} [Y_1]$$

- Similarly for Y_0 : $\mathbb{E} \left[\frac{Y(1-D)}{1-p(X)} \right] = \mathbb{E} [Y_0]$. Thus, under CIA+overlap:

$$ATE = \mathbb{E} \left[\left(\frac{D}{p(X)} - \frac{1-D}{1-p(X)} \right) Y \right] = \mathbb{E} \left[\frac{D - p(X)}{p(X)(1-p(X))} \cdot Y \right]$$

- *Exercise:* derive the reweighting expression for ATT

P-score adjustment methods: Reweighting (2)

- Plug-in (Horvitz-Thompson) estimator: $\widehat{ATE}_{HT} = \frac{1}{N} \sum_i \left(\frac{D_i}{\hat{p}(X_i)} - \frac{1-D_i}{1-\hat{p}(X_i)} \right) Y_i$
- Issue: weights on treated $\left(\frac{1}{N} \frac{1}{\hat{p}(X_i)} \right)$ and control $\left(\frac{1}{N} \frac{1}{1-\hat{p}(X_i)} \right)$ do not exactly sum to 1
 - ▶ So adding a constant to Y_i changes the estimator
- Normalizing the weights improves performance: the Hajek estimator

$$\widehat{ATE}_{Hajek} = \frac{\sum_i \frac{D_i Y_i}{\hat{p}(X_i)}}{\sum_i \frac{D_i}{\hat{p}(X_i)}} - \frac{\sum_i \frac{(1-D_i) Y_i}{1-\hat{p}(X_i)}}{\sum_i \frac{1-D_i}{1-\hat{p}(X_i)}}$$

Outline

- 1 The concept of control variables
- 2 Regression adjustment
- 3 Propensity score methods
- 4 Doubly-robust methods

Idea of double robustness

- Regression adjustment methods require estimates of $h_0(X)$, $h_1(X)$
 - ▶ No need for $p(X) = \mathbb{E}[D | X]$
- Propensity score methods are the opposite
- Doubly robust methods take estimates of both $h_d(X)$ and $p(X)$ as inputs
 - ▶ But validity requires only one of those estimates to be correct

Automatic double robustness

Some methods already possess some double robustness

- Under constant effects, regression $Y_i = \beta D_i + \gamma' X_i + \text{error}_i$ is causal if:
 - ▶ $h_0(X)$ is linear in X **or** $p(X)$ is linear in X
- Kline (2011): the Oaxaca-Blinder estimator for ATT is consistent if:
 - ▶ $h_0(X), h_1(X)$ are linear in X **or** $\frac{p(X)}{1-p(X)}$ is linear in X

Augmented IPW (AIPW)

- But double-robustness can be achieved for arbitrary estimators of $h_d(\cdot)$ and $p(\cdot)$
- Augmented IPW (**AIPW**) idea:

$$\begin{aligned}\mathbb{E}[Y_{1i}] &= \mathbb{E}[h_1(X_i)] = \mathbb{E}\left[\frac{D_i}{p(X_i)} Y_i\right] = \mathbb{E}\left[h_1(X_i) + \frac{D_i}{p(X_i)} (Y_i - h_1(X_i))\right] \\ &= \mathbb{E}\left[\tilde{h}_1(X_i) + \frac{D_i}{\tilde{p}(X_i)} (Y_i - \tilde{h}_1(X_i))\right] \quad \text{if } p(\cdot) = \tilde{p}(\cdot) \text{ or } h_1(\cdot) = \tilde{h}_1(\cdot)\end{aligned}$$

- ▶ If the model of $h_1(\cdot)$ is correct, IPW adjustment doesn't change the estimand
- ▶ If the model of $p(\cdot)$ is correct, the adjustment fixes mistakes in $h_1(\cdot)$

Augmented IPW (2)

- Combining with a similar expression for Y_{0i} ,

$$ATE = \mathbb{E} \left[h_1(X_i) + \frac{D_i}{p(X_i)} (Y_i - h_1(X_i)) - h_0(X_i) - \frac{1 - D_i}{1 - p(X_i)} (Y_i - h_0(X_i)) \right]$$

if $(h_0(\cdot), h_1(\cdot))$ **or** $p(\cdot)$ are correctly specified

- ▶ The sample analog based on preliminary estimates of h_0, h_1, p yields the estimator

Double robustness and Neyman orthogonality

- The usefulness of double robustness is not obvious
- But a closely related concept is very helpful when X_i are high-dimensional:
 - ▶ Estimation of **nuisance functions** $h_0(\cdot), h_1(\cdot), p(\cdot)$, e.g. via machine learning methods, is inevitably imprecise
 - ▶ So we want to use moments that have low sensitivity to such mistakes around the true values — **Neyman-orthogonal** moments
 - ★ *Optional exercise*: verify that the “derivative” of the IAPW moment w.r.t. nuisance functions is zero at the true values
 - ▶ Covariate adjustment by Double Machine Learning (Chernozhukov et al., 2018) uses this idea

References I

- ANGRIST, J. (1998): "Estimating the labor market impact of voluntary military service using social security data on military applicants," *Econometrica*, 66, 249–288.
- CHERNOZHUKOV, V., D. CHETVERIKOV, M. DEMIRER, E. DUFLO, C. HANSEN, W. NEWHEY, AND J. ROBINS (2018): "Double/debiased machine learning for treatment and structural parameters," *Econometrics Journal*, 21, C1–C68.
- HORVITZ, D. G. AND D. J. THOMPSON (1952): "A generalization of sampling without replacement from a finite universe," *Journal of the American statistical Association*, 47, 663–685.
- KLINE, P. (2011): "Oaxaca-Blinder as a reweighting estimator," *American Economic Review*, 101, 532–537.
- ROSENBAUM, P. R. AND D. B. RUBIN (1983): "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika*, 70, 41–55.