

## Lecture 4: Propensity Score, Part 1

$$y_i = \theta T_i + X_i' \beta + \varepsilon_i$$

$T_i$  is a binary (0-1) treatment with homogenous treatment effects:

$$T_i = \begin{cases} 1 & \text{if treated} \\ 0 & \text{not treated} \end{cases}$$

↳ There's nothing in between! No being "half-way" treated. You are either treated or not. 0 or 1.

↳  $\theta_i = \theta \quad \forall i$  (although we'll consider heterogeneous treatment effects later).

$i$  has two potential outcomes:

$y_{0i}$  if  $T_i = 0$

$y_{1i}$  if  $T_i = 1$

$$\theta_i = y_{1i} - y_{0i} \Rightarrow E(y_{1i} - y_{0i}) = E(\theta_i) = \theta$$

$(\bar{y}_1 - \bar{y}_0)$

Problem:

For individual  $i$ , observe either  $y_{1i}$  or  $y_{0i}$ , but not both! Unobserved counterfactual.

What's the "gold standard" solution?

- Random assignment of  $T_i$

$$\Rightarrow \{y_{1i}, y_{0i}\} \perp\!\!\!\perp T_i$$

- Control group ( $y_{0i}$ ) identifies correct counterfactual as  $N \rightarrow \infty$

$$\bar{y}_1 - \bar{y}_0 \xrightarrow{P} \theta \text{ as } N \rightarrow \infty$$

- Indirect test of random assignment:

$$\bar{X}_1 \approx \bar{X}_0 \quad \forall x_{ik}$$

Linear model:

$$y_i = \theta T_i + X_i' \beta + \varepsilon_i \quad \text{where we assume } E(\varepsilon_i | T_i) = 0$$

Comparing mean of  $y_i$  conditional on  $T_i$ :

$$E(y_i | T_i=1) - E(y_i | T_i=0) = E(y_{1i} | T_i=1) - E(y_{0i} | T_i=0)$$

$$= E(y_{1i} | T_i=1) - \underbrace{E(y_{0i} | T_i=1)} + \underbrace{E(y_{0i} | T_i=1) - E(y_{0i} | T_i=0)}$$

$$= \underbrace{E(y_{1i} - y_{0i} | T_i=1)} + \underbrace{[E(y_{0i} | T_i=1) - E(y_{0i} | T_i=0)]}$$

$$= ATE$$

= 0 if  $T_i$  randomly assigned

$= 0$  + selection bias if  $T_i$  not randomly assigned.

Example: Birthweight and Smoking.

$$E(y_{1i} - y_{0i} | T_i = 1) + [E(y_{0i} | T_i = 1) - E(y_{0i} | T_i = 0)]$$

$$= -210 \text{ g} + [2950 \text{ g} - 3000 \text{ g}]$$

$$= -210 \text{ g} + [-50 \text{ g}]$$

$$= -260 \text{ g}$$

The counterfactual  
BW for women  
who smoke if  
they didn't  
smoke: 2950g

Why lower?

Because they're  
drinking, not  
going to  
prenatal  
visits, etc.

Their BW is  
lower than  
for women  
who don't  
smoke (3000g)

- Random assignment conditional on observables:

$$\{y_{1i}, y_{0i}\} \perp\!\!\!\perp T_i \mid X_i$$

↳  $T_i$  is independent of potential outcomes conditional on  $X_i$ .

$$= E(y_{1i} - y_{0i} \mid T_i = 1, X_i) + \underbrace{[E(y_{0i} \mid T_i = 1, X_i) - E(y_{0i} \mid T_i = 0, X_i)]}$$

$$= \theta$$

$= 0$  if only source of bias before was due to  $X_i$  and now we've removed this bias.

$\neq 0$  if omitted variables or misspecification.  
of  $\underbrace{g(X_i)}$ .

$g(X_i) = X_i' \beta$  is only one (linear) possibility.  
What if not the right one?

### - Matching: Univariate Case

- For each treatment observation, match control case with "identical"  $X_i$ .

(Problem if  $X_i, T_i$  are collinear  $\Rightarrow$  impossible to match).

- Using the matched pairs, run a regression controlling for "pair identifier" fixed effects.

```

1  cd "C:\Users\lfeler1\Documents\Applied Econometrics Course\Notes\Weeks 4-6"
2  clear
3  use smoking2
4
5  **This is the regression we would run if smoking was randomly assigned**
6  reg dbirwt tobacco
7
8
9  **Check if smoking randomly assigned**
10 sort tobacco
11 ttest dimage, by(tobacco)
12 ttest dmeduc, by(tobacco)
13 ttest dmar, by(tobacco)
14 ttest nprevist, by(tobacco)
15 ttest alcohol, by(tobacco)
16 ttest anemia, by(tobacco)
17 ttest mblack, by(tobacco)
18     *So no, it does not look like smoking is randomly assigned!*
19
20
21 **Assume smoking is randomly assigned conditional on observables**
22 reg dbirwt tobacco dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete ///
23 phyper alcohol drink foreignb plural deadkids mblack motherr mhispan fblack fotherr
24 fhispan first
25 estat imtest, white
26     *So reject homo in favor of heteroskedasticity
27
28 reg dbirwt tobacco dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete ///
29 phyper alcohol drink foreignb plural deadkids mblack motherr mhispan fblack fotherr
30 fhispan first, robust
31 estat ovtest
32     *This regresses y on x y-hat^2 y-hat^3 y-hat^4 and jointly tests that the coeffs on
33 y-hat^2 y-hat^3 y-hat^4 are zero.
34     *We cannot reject that the model as no omitted variables.
35
36 *Let's try matching on... education.*
37 sort dmeduc
38 areg dbirwt tobacco, absorb(dmeduc) robust
39
40
41
42
43
44
45

```

## - Propensity Score

- What is it? An index (one variable) constructed out of all the  $X_k$ 's, and possibly their squares, cubics, and interactions.
- Why? Reduces multi-dimensional  $X_k$  into one dimension! Makes matching possible.

### ↳ Propensity Score Theorem:

If  $\{y_{1i}, y_{0i}\} \perp\!\!\!\perp T_i | X_i$ , then  
 $\{y_{1i}, y_{0i}\} \perp\!\!\!\perp T_i | \underbrace{p(X_i)}_{p_i}$ .

↳ where  $p_i \equiv \Pr(T_i=1 | X_i) = E(T_i | X_i) \equiv p(X_i)$   
Probability of treatment conditional on  $X_i$ .

- Now just control or match for single index  $p(X_i)$  rather than all  $X_k$ 's.



~> How?

- (1) Estimate propensity score,  $\hat{p}(x_i)$  such that it balances  $X_k$ 's.
- (2) Estimate  $\theta \equiv ATE$  by controlling for  $\hat{p}(x_i)$ .

---

Step (1): Estimate propensity score,  $\hat{p}(x_i)$ .

$$- \Pr(T_i=1|X_i) = \frac{e^{h(X_i)}}{1 + e^{h(X_i)}}$$

- $h(X_i)$  contains linear and possibly higher order terms and interactions.  
~> include enough terms so that Treatments and Controls with similar  $\hat{p}(x_i)$  have similar  $X_k$ 's.

"Algorithm" for estimating  $p(x_i)$

- (1) start with parsimonious logit  $\rightarrow$  estimate  $\hat{p}(x_i)$
- (2) stratify data into 5 blocks of  $\hat{p}(x_i)$
- (3) test  $\bar{X}_1 = \bar{X}_0$  for all  $K$  within each block, using  $t$ -test of significant differences in sample means.
  - (a) if  $X_k$ 's "balanced" in each block, STOP.
  - (b) if  $X_k$ 's not balanced in some blocks, divide block into 2 blocks ⑧



and reevaluate.

- (c) if  $X_k$ 's not balanced in all blocks, add interaction and/or polynomial of  $X_k$ 's to logit and reevaluate.

Goal: Balance  $X_k$ 's in Treatment and Control groups in each block.

→ overlap in  $\hat{p}(x_i)$  for  $T_i = 1$  and  $T_j = 0$  implies overlap in  $X_i$ 's.

Stopping Rule: Stop when fail to reject  $\bar{X}_{1k} = \bar{X}_{0k}$  for over 90% of  $t$ -tests within a block.

```

1  cd "C:\Users\lfeler1\Documents\Applied Econometrics Course\Notes\Weeks 4-6"
2  clear
3  use smoking2
4  set seed 1000
   sample 20
5
6
7  pscore tobacco  dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete phyper
   alcohol ///
8  drink foreignb plural deadkids mblack motherr mhispan fblack fotherr fhispan first, ///
9  logit pscore(phat1) blockid(block1) numblo(5) detail
10
11
12
13  pscore tobacco  dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete phyper
   alcohol ///
14  drink foreignb plural deadkids mblack motherr mhispan fblack fotherr fhispan first, ///
15  logit pscore(phat2) blockid(block2) numblo(5)
16
17
18  gen dmeduc2=dmeduc^2
19  gen dimageXdmeduc=dimage*dmeduc
20
21  pscore tobacco  dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete phyper
   alcohol ///
22  drink foreignb plural deadkids mblack motherr mhispan fblack fotherr fhispan first ///
23  dmeduc2 dimageXdmeduc, ///
24  logit pscore(phat3) blockid(block3) numblo(15)
25
26
27  gen dmeducXdmeduc=dmeduc*dmeduc
28
29  pscore tobacco  dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete phyper
   alcohol ///
30  drink foreignb plural deadkids mblack motherr mhispan fblack fotherr fhispan first ///
31  dmeduc2 dimageXdmeduc dmeducXdmeduc, ///
32  logit pscore(phat4) blockid(block4) numblo(15)
33
34
35  gen mblackXfblack=mblack*fblack
36
37  pscore tobacco  dimage dmeduc dmar ddivord nprevist dfage dfeduc anemia diabete phyper
   alcohol ///
38  drink foreignb plural deadkids mblack motherr mhispan fblack fotherr fhispan first ///
39  dmeduc2 dimageXdmeduc dmeducXdmeduc mblackXfblack, ///
40  logit pscore(phat5) blockid(block5) numblo(20)
41
42  ****This last one just made things worse, so go back****
43
44  graph box phat4, by(tobacco)
45
46

```