# Propensity Score Matching

Rajeev Dehejia

New York University

# Adjustment for observables in observational studies (Conditioning)

1. Subclassification

2. Matching

3. Regression

4. Propensity Score Methods

# Recall the goal and challenge

- We are in the universe of selection on observables:

$$\{Y_{i1}, Y_{i0} \perp\!\!\!\perp T_i\}|X_i.$$

- The problem is that to estimate

$$\tau\Big|_{T_i=1} = E_X[E(Y_{i1}|T_i=1,X_i)] - E_X[E(Y_{i0}|T_i=1,X_i)]$$

when **the treatment and comparison groups are very different**, and when $X$ is high dimensional.

# Why is this hard? Lalonde example

Why Non-Experimental Methods are Difficult: Lalonde Example

| | A. Lalonde's original sample | | | | |
|---|---|---|---|---|---|
| Comp-arison Group | NSW Treatment Earnings Less Comparison Group Earnings, 1978 | | Unrestricted Difference in Differences: Quasi-Differ-ence in Earn-ings Growth: 1975-1978 | | Contr-olling for All Vari-ables[f] |
| | Unad-justed[b] | Ad-justed[c] | Unad-justed[d] | Adjusted[e] | |
| | (1) | (2) | (3) | (4) | (5) |
| NSW | 886 (472) | 798 (472) | 879 (467) | 802 (468) | 820 (468) |
| PSID-1 | -15,578 (913) | -8,067 (990) | -2,380 (680) | -2,119 (746) | -1,844 (762) |
| PSID-2 | -4,020 (781) | -3,482 (935) | -1,364 (729) | -1,694 (878) | -1,876 (885) |
| PSID-3 | 697 (760) | -509 (967) | 629 (757) | -552 (967) | -576 (968) |
| CPS-1 | -8,870 (562) | -4,416 (577) | -1,543 (426) | -1,102 (450) | -987 (452) |
| CPS-2 | -4,195 (533) | -2,341 (620) | -1,649 (459) | -1,129 (551) | -1,149 (551) |
| CPS-3 | -1,008 (539) | -1 (681) | -1,204 (532) | -263 (677) | -234 (675) |

PSID-1: All male household heads under age 55 who did not classify themselves as retired in 1975.
PSID-2: Selects from PSID-1 all men who were not working when surveyed in the spring of 1976.
PSID-3: Selects from PSID-2 all men who were not working in 1975.
CPS-1: All CPS males under age 55.
CPS-2: Selects from CPS-1 all males who were not working when surveyed in March 1976.
CPS-3: Selects from CPS-2 all the unemployed males in 1976 whose income in 1975 was below the poverty level.

# Why is this hard? Lalonde example

Why Non-Experimental Methods are Difficult: Lalonde Example

| Comparison Group | NSW Treatment Earnings Less Comparison Group Earnings, 1978 | | Unrestricted Difference in Differences: Quasi-Differ-ence in Earn-ings Growth: 1975-1978 | | Controlling for All Variables[f] |
|---|---|---|---|---|---|
| | Unadjusted[b] | Adjusted[c] | Un-adjust-ed[d] | Ad-justed[e] | |
| | (1) | (2) | (3) | (4) | (5) |
| NSW | 1,794 (633) | 1,672 (637) | 1,750 (632) | 1,631 (637) | 1,612 (639) |
| PSID-1 | -15,205 (1155) | -7,741 (1175) | -582 (841) | -265 (881) | 186 (901) |
| PSID-2 | -3,647 (960) | -2,810 (1082) | 721 (886) | 298 (1004) | 111 (1032) |
| PSID-3 | 1,070 (900) | 35 (1101) | 1,370 (897) | 243 (1101) | 298 (1105) |
| CPS-1 | -8,498 (712) | -4,417 (714) | -78 (537) | 525 (557) | 709 (560) |
| CPS-2 | -3,822 (671) | -2,208 (746) | -263 (574) | 371 (662) | 305 (666) |
| CPS-3 | -635 (657) | 375 (821) | -91 (641) | 844 (808) | 875 (810) |

B. RE74 Subsample (results do not use RE74)

# Why is this hard? Lalonde example

Why Non-Experimental Methods are Difficult: Lalonde Example

| Comp-arison Group | NSW Treatment Earnings Less Comparison Group Earnings, 1978 | | Unrestricted Difference in Differences: Quasi-Differ-ence in Earn-ings Growth: 1975-1978 | | Contr-olling for All Vari-ables[f] |
|---|---|---|---|---|---|
| | Unad-justed[b] | Ad-justed[c] | Un-adjust-ed[d] | Ad-justed[e] | |
| | (1) | (2) | (3) | (4) | (5) |
| NSW | 1,794 | 1,688 | 1,750 | 1,672 | 1,655 |
| | (633) | (636) | (632) | (638) | (640) |
| PSID-1 | -15,205 | -879 | -582 | 218 | 731 |
| | (1155) | (931) | (841) | (866) | (886) |
| PSID-2 | -3,647 | 94 | 721 | 907 | 683 |
| | (960) | (1042) | (886) | (1004) | (1028) |
| PSID-3 | 1,070 | 821 | 1,370 | 822 | 825 |
| | (900) | (1100) | (897) | (1101) | (1104) |
| CPS-1 | -8,498 | -8 | -78 | 739 | 972 |
| | (712) | (572) | (537) | (547) | (550) |
| CPS-2 | -3,822 | 615 | -263 | 879 | 790 |
| | (671) | (672) | (574) | (654) | (658) |
| CPS-3 | -635 | 1,270 | -91 | 1,326 | 1,326 |
| | (657) | (798) | (641) | (796) | (798) |

**C. RE74 Subsample** (results use RE74)

# Why is this hard? Cochran

Cochran (1965) wrote:

> If the original x-distributions [of comparison groups] diverge widely, none of the methods [e.g., regression adjustment] can be trusted to remove all, or nearly all, the bias. This discussion brings out the importance of finding comparison groups in which the initial difference among the distributions of the distributing variables are small.

*In other words, try to compare apples to apples as much as possible…*

# Why is this hard? Lalonde revisited

*Table 1.   Sample Means of Characteristics for NSW and Comparison Samples*

| | No. of observations | Age | Education | Black | Hispanic | No degree | Married | RE74 (U.S. $) | RE75 (U.S. $) |
|---|---|---|---|---|---|---|---|---|---|
| **NSW/Lalonde:[a]** | | | | | | | | | |
| Treated | 297 | 24.63 | 10.38 | .80 | .09 | .73 | .17 | | 3,066 |
| | | (.32) | (.09) | (.02) | (.01) | (.02) | (.02) | | (236) |
| Control | 425 | 24.45 | 10.19 | .80 | .11 | .81 | .16 | | 3,026 |
| | | (.32) | (.08) | (.02) | (.02) | (.02) | (.02) | | (252) |
| **RE74 subset:[b]** | | | | | | | | | |
| Treated | 185 | 25.81 | 10.35 | .84 | .059 | .71 | .19 | 2,096 | 1,532 |
| | | (.35) | (.10) | (.02) | (.01) | (.02) | (.02) | (237) | (156) |
| Control | 260 | 25.05 | 10.09 | .83 | .1 | .83 | .15 | 2,107 | 1,267 |
| | | (.34) | (.08) | (.02) | (.02) | (.02) | (.02) | (276) | (151) |
| **Comparison groups:[c]** | | | | | | | | | |
| PSID-1 | 2,490 | 34.85 | 12.11 | .25 | .032 | .31 | .87 | 19,429 | 19,063 |
| | | [.78] | [.23] | [.03] | [.01] | [.04] | [.03] | [991] | [1,002] |
| PSID-2 | 253 | 36.10 | 10.77 | .39 | .067 | .49 | .74 | 11,027 | 7,569 |
| | | [1.00] | [.27] | [.04] | [.02] | [.05] | [.04] | [853] | [695] |
| PSID-3 | 128 | 38.25 | 10.30 | .45 | .18 | .51 | .70 | 5,566 | 2,611 |
| | | [1.17] | [.29] | [.05] | [.03] | [.05] | [.05] | (686) | [499] |
| CPS-1 | 15,992 | 33.22 | 12.02 | .07 | .07 | .29 | .71 | 14,016 | 13,650 |
| | | [.81] | [.21] | [.02] | [.02] | [.03] | [.03] | [705] | [682] |
| CPS-2 | 2,369 | 28.25 | 11.24 | .11 | .08 | .45 | .46 | 8,728 | 7,397 |
| | | [.87] | [.19] | [.02] | [.02] | [.04] | [.04] | [667] | [600] |
| CPS-3 | 429 | 28.03 | 10.23 | .21 | .14 | .60 | .51 | 5,619 | 2,467 |
| | | [.87] | [.23] | [.03] | [.03] | [.04] | [.04] | [552] | [288] |

NOTE:   Standard errors are in parentheses. Standard error on difference in means with RE74 subset/treated is given in brackets. Age = age in years; Education = number of years of schooling; Black = 1 if black, 0 otherwise; Hispanic = 1 if Hispanic, 0 otherwise; No degree = 1 if no high school degree, 0 otherwise; Married = 1 if married, 0 otherwise; RE*x* = earnings in calendar year 19*x*.

# Why is this hard? Recall Abadie-Imbens result

A standard technique for estimating the treatment effect is to matching on the covariates *X.*

## Idea

In the above expression we are conditioning on *X.* So if we can somehow match units with the same (similar) *X*, then within these cells we can directly estimate $E(Y_{i1}|T_i=1,X_i)$ and

$E(Y_{i0}|T_i=1,X_i)$.

## Problem

Theorem (Abadie-Imbens): If there are two or more continuous covariates, then the simple matching estimator is asymptotically biased (not $N^{-1/2}$ consistent).

# An obvious, efficient, but impractical solution

Hahn (1998) shows that if we estimate $E(Y_{i1}|T_i=1,X_i)$ and $E(Y_{i0}|T_i=1,X_i)$ non-parametrically, the semi-parametric efficiency bound is achieved.

## Problem

This involves estimating two non-parametric equations. When $X$ is high-dimensional and has many continuous variables (typical application – 10 variables, 5 of them continuous) this is very difficult computationally.

- E.g., local linear regression like we used in RD would not be practical

# Propensity score methods

Propensity score methods help to control for sample election bias, under the assumption of selection on observables.

Experience suggests that dealing with the problem of sample selection bias is difficult.

- Extrapolation and linear regression regression is unreliable under many circumstances.
- Parametric selection models are not robust.
- Instruments are not always available.
- Higher-order non-parametric regression is difficult.

Propensity score methods offer a solution to this problem, when selection is on observable variables.

# Propensity score methods: General setup

- The treatment may be an explicit or organized program (such as labor training, a medical trial, an educational experiment) or may be a policy variation (such as certain subjects being exposed to a new set of regulations).

- Treatment group: exposed to the variation of phenomenon of interest. Denote treatment-exposed individuals, $i$, in this group with $T_i=1$.

- Control group: not exposed to the treatment; exposed to an alternative, baseline, or status quo program. Denote individuals, $i$, in this group with $T_i=0$.

# Propensity score methods

## Definition

Propensity score is defined as the **selection probability conditional on the confounding variables**: $p(X) = P(T = 1|X)$

## Identification Assumption

1. $(Y_1, Y_0) \perp\!\!\!\perp T | X$ *(selection on observables)*
2. $0 < \Pr(T = 1|X) < 1$ *(common support)*

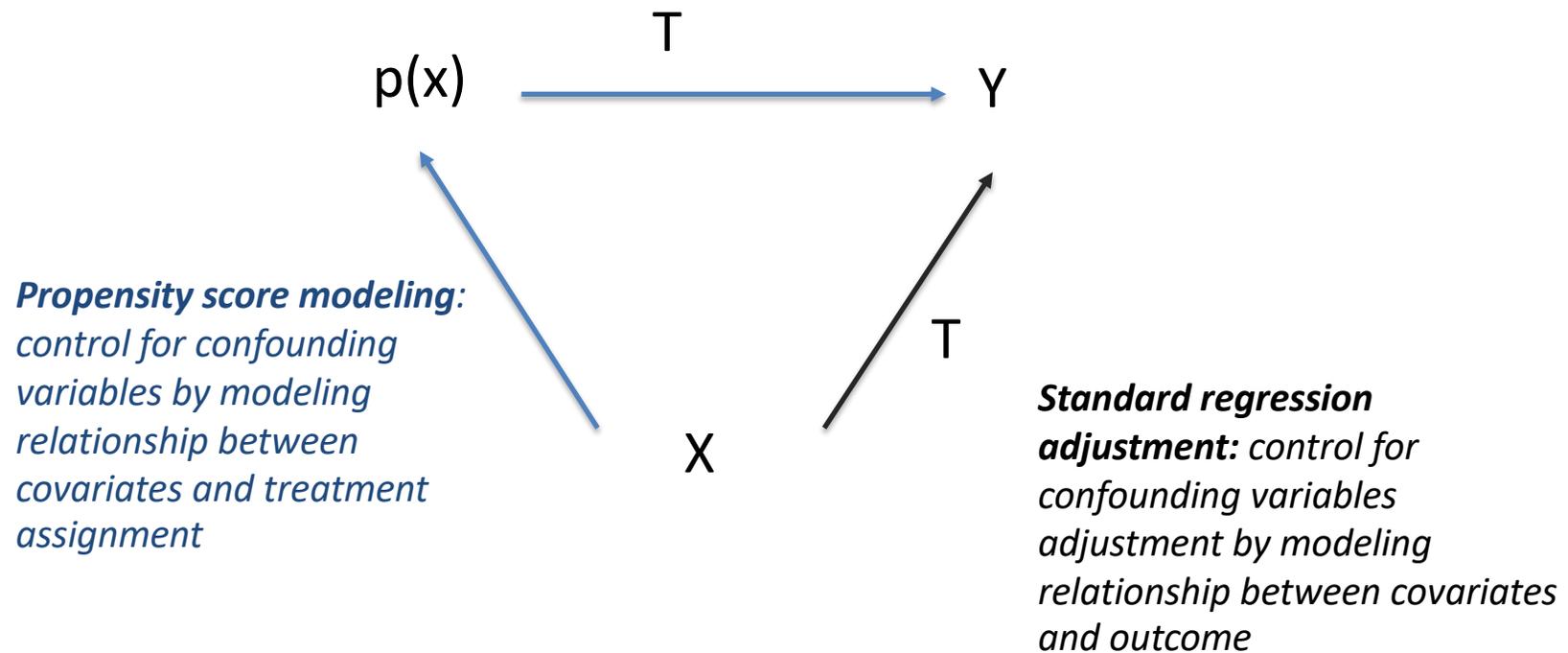## Identification Result (Rosenbaum and Rubin, 1983)

*Selection on observables implies:*

$$(Y_1, Y_0) \perp\!\!\!\perp T | X \Rightarrow (Y_1, Y_0) \perp\!\!\!\perp T | p(X)$$

→ *conditioning on the propensity score is enough to have independence between the treatment indicator and the potential outcomes*

→ *substantial dimension reduction in the matching variables!*

# Propensity Score Matching: Graphical Representation

p(x) $\xrightarrow{\quad T \quad}$ Y

**Propensity score modeling**: control for confounding variables by modeling relationship between covariates and treatment assignment

T

X

**Standard regression adjustment**: control for confounding variables adjustment by modeling relationship between covariates and outcome

# Propensity score methods

Proof.

Assume that $(Y_1, Y_0) \perp\!\!\!\perp T \mid X$. Then:

$$
\begin{aligned}
P(T = 1 \mid Y_1, Y_0, p(X)) &= E[T \mid Y_1, Y_0, p(X)] \\
&= E\{E_X[E(T \mid Y_1, Y_0, X)] \mid Y_1, Y_0, p(X)\} \\
&= E\{E_X[E(T \mid X)] \mid Y_1, Y_0, p(X)\} \\
&= E\{E_X[p(X)] \mid Y_1, Y_0, p(X)\} \\
&= p(X)
\end{aligned}
$$

Using a similar argument, we obtain

$$
\begin{aligned}
P(T = 1 \mid p(X)) &= E[T \mid p(X)] = E[E[T \mid X] \mid p(X)] \\
&= E[p(X) \mid p(X)] = p(X)
\end{aligned}
$$

➔ $P(T = 1 \mid Y_1, Y_0, p(X)) = P(T = 1 \mid p(X))$

➔ $(Y_1, Y_0) \perp\!\!\!\perp T \mid p(X)$

# Matching on the propensity score

Corollary

If $(Y_1, Y_0) \perp\!\!\!\perp T \mid X$, then

$$E[Y_1 - Y_0 \mid p(X)] = E[Y \mid T = 1, p(X)] - E[Y \mid T = 0, p(X)]$$

Suggests a two-step procedure to estimate causal effects under selection on observables:

1. Estimate the propensity score $p(X) = P(T = 1 \mid X)$ (e.g., using logit or probit regression)
2. Do matching or subclassification on the estimated propensity score

Hugely popular method to estimate treatment effects. However, valid method to calculate standard errors not known until (very) recently.

# Propensity score: balancing property

Because the propensity score, $p(X)$ is a function of $X$:

$$\Pr(T = 1 \mid p(X)) = \Pr(T = 1 \mid X)$$
$$= p(X)$$

→ conditional on $p(X)$, the probability that $T = 1$ does not depend on $X$ beyond $p(X)$. *Think back to the exclusion restriction from IV....*

→ $T$ and $X$ are independent conditional on $p(X)$:

$$T \perp\!\!\!\perp X \mid p(X).$$

So we obtain the **balancing property** of the propensity score:

$$P(X \mid T = 1, p(X)) = P(X \mid T = 0, p(X)),$$

*conditional on the propensity score, the distribution of the covariates is the same for treated and non-treated.*

We can use this to check if our estimated propensity score actually produces balance:

$$P\left(X \middle| T = 1, \hat{p}(X)\right) = P\left(X \middle| T = 0, \hat{p}(X)\right)$$

# How to do it: using balancing to estimate $p(X)$

1. Start with a parsimonious logit specification to estimate the score.

2. Sort data according to estimated propensity score (ranking from lowest to highest).

3. Stratify all observations such that estimated propensity scores within a stratum for treated and comparison units are close (no significant difference); e.g., start by dividing observations into strata of equal score range (0-0.2,...,0.8-1).

| Unit ID | Earnings | College Grad | + k descriptive characteristics |
|---|---|---|---|
| 1 | $100 | 0 | |
| 2 | $300 | 1 | *Several fixed observables such as gender, age, parents' education, etc. for each unit* |
| 3 | $150 | 1 | |
| 4 | $50 | 0 | |
| 5 | $75 | 1 | |
| 6 | $225 | 0 | |

```
Install one of a
few stata programs,
such as pscore, to
run the logit reg,
estimate the
propensity scores,
and create a
propensity score
variable by which
you can sort the
data
```

| Unit ID | Earnings | College Grad | *Propensity Score* | *Quartile* |
|---|---|---|---|---|
| 2 | $300 | 1 | *0.98* | *1* |
| 4 | $50 | 0 | *0.91* | *1* |
| 5 | $75 | 1 | *0.61* | *2* |
| 6 | $225 | 0 | *0.59* | *2* |
| 1 | $100 | 0 | *0.40* | *3* |
| 3 | $150 | 1 | *0.39* | *3* |

# How to do it: using balancing to estimate $p(X)$

**Statistical test:** for all covariates, differences in means across treated and comparison units within each stratum are not significantly different from zero.

- ✓ If covariates are balanced between treated and comparison observations for all strata, stop.
- ✓ If covariates are not balanced for some stratum, divide the stratum into finer strata and re-evaluate.
- ✓ If a covariate is not balanced for many strata, modify the logit by adding interaction terms and/or higher-order terms of the covariate and re-evaluate.

# How to do it: using balancing to estimate $p(X)$

✓ The difference in the means of the propensity scores in the two groups being compared must be small (e.g., the means must be less than half a standard deviation apart), unless the situation is benign in the sense that: (a) the distributions of the covariates in both groups are nearly symmetric, (b) the distributions of the covariates in both groups have nearly the same variances, and (c) the sample sizes are approximately the same.

✓ The ratio of the variances of the propensity score in the two groups must be close to one (e.g., 1/2 or 2 are far too extreme)

✓ The ratio of the variances of the residuals of the covariates after adjusting for the propensity score must be close to one (e.g., 1/2 or 2 are far too extreme)

# How to do it: controversies with balancing

- Imai, King, and Stuart (2007) argue that traditional hypothesis testing is inappropriate to test balance. And in principle they are right:
  - Balance is a property of the sample at hand, not a population, so not clear that starting hypothesis tests based on sampling apply.
  - Mechanically reducing sample size (matching, throwing observations away even at random) reduces precision, so t-stats come down.
- Their alternative is to look at the matched sample:
  - Compare means, relative to standard deviations.
  - QQ plots (differences of quantiles of the the samples).
  - Useful, but does render the decision of whether balance has been achieved as judgmental.

# How to do it: stratification

1. **Sort** treatment and comparison observations based on the estimated propensity score.

2. **Drop** comparison observations whose estimated $p(X)$ is lower than the $\min(p(X_i))$ or $\max(p(X_i))$ (for $i$ treated) (called *common support or overlap*).

3. **Check:** If there is a range where $p(X_i)$ for treated observations is less than $\min(p(X_i))$ or more likely greater than $\max(p(X_i))$, then the non-overlap problem is fatal – we can't estimate the average treatment effect.

4. **Create strata** (bins, sub-classes) on $p(X_i)$ such that within each sub-class the distribution of $(X_i, p(X_i))$ is balanced across the treated and comparison groups.

5. Within each stratum **compute** $E(Y_i | p(X_i), T_i=1) - E(Y_i | p(X_i), T_i=0)$, and **weight** these by the number of treated observations in each stratum to estimate the average treatment effect for the treated (corresponds to a piece-wise linear functional form).

# Stratification: Dehejia-Wahba Example

**Table 5: Estimating Treatment Effect for NSW Male Group
with PSID-1 Control Group Stratifying on the Score[a]**

| Score Range | Number of Observations | | Average Score[b] | NSW Treatment Earnings Less Comparison group Earnings | |
|---|---|---|---|---|---|
| | NSW | PSID | | Unadjusted | Adjusted[c] |
| 0.0004-0.1 | 9 | 942 | 0.022 | -2,693 (3660) | -4,896 (2830) |
| 0.1-0.2 | 10 | 57 | 0.13 | -2,759 (5462) | 1,078 (5355) |
| 0.2-0.45 | 19 | 39 | 0.34 | -1,171 (1410) | -1,070 (1424) |
| 0.45-0.6 | 20 | 15 | 0.53[*] | -3,044 (2354) | -4,244 (2700) |
| 0.6-0.85 | 35 | 10 | 0.70 | -837 (2603) | -265 (2822) |
| 0.85-1 | 92 | 7 | 0.94 | 4,819 (3412) | 4,918 (3469) |
| Weighted Average | 185 | 1,070 | 0.70 | 1,509 (1823) | 1,647 (1862) |

# Stratification: Dehejia-Wahba example

**Table 6: Estimating Treatment Effect for NSW Male Group with CPS-1 Control Group Stratifying on the Score[a]**

| Score Range | Number of Observations | | Average Score[b] | NSW Treatment Earnings Less Comparison group Earnings | |
|---|---|---|---|---|---|
| | NSW | CPS | | Unadjusted | Adjusted[c] |
| 0.001-0.01 | 10 | 2767 | 0.003 | -1,740 (2596) | -2,242 (2134) |
| 0.01-0.1 | 30 | 935 | 0.078[*] | -139 (1354) | 547 (1175) |
| 0.1-0.3 | 32 | 150 | 0.187 | 500 (1109) | 493 (1038) |
| 0.3-0.5 | 35 | 46 | 0.39 | 200 (1397) | 456 (1462) |
| 0.5-0.6 | 22 | 17 | 0.56 | 6,445 (2604) | 4,895 (3116) |
| 0.6-0.85 | 26 | 12 | 0.7 | 2,364 (2213) | 3,683 (2521) |
| 0.85-1 | 30 | 5 | 0.90 | 3,740 (5772) | 3,299 (5869) |
| Weighted Average | 185 | 3,932 | 0.43 | 1,713 (1115) | 1,774 (1152) |

# Stratification: Dehejia-Wahba example

**Table 9a: Treatment Effect for NSW Male and CPS Control Sample**

**Blocking on Selected Sample Characteristics (s.e.)**

| Sample | All | Black | Non-Black | Ndgree=1 | Ndgree=0 | Educ >=11 | Educ <11 |
|---|---|---|---|---|---|---|---|
| NSW | 1,794 | 2,029 | 803 | 1,154 | 3,192 | 3,085 | 402 |
| | (633) | (706) | (1331) | (696) | (1517) | (1033) | (753) |
| CPS-1: | | | | | | | |
| Unadjusted | -8,498 | -5,870 | -7,578 | -6,936 | -7,750 | -8,040 | -7,541 |
| | (712) | (785) | (1789) | (793) | (1329) | (987) | (967) |
| Adjusted | 738 | 1,487 | 1,087 | 511 | 1,946 | 1,628 | 223 |
| | (547) | (617) | (1311) | (635) | (1004) | (760) | (749) |
| Stratifying on the Score: | | | | | | | |
| -Unadjusted | 1,713 | 1,738 | 1,367 | 1,439 | 2,475 | 2,501 | 1,012 |
| | (1115) | (1191) | (1397) | (1345) | (1420) | (985) | (883) |
| -Adjusted | 1,774 | 1,905 | 1,799 | 1,144 | 2,683 | 2,377 | 738 |
| | (1152) | (1331) | (1847) | (1536) | (1256) | (1913) | (1149) |

# Stratification: Dehejia-Wahba example

**Table 9a: Treatment Effect For NSW Male And CPS Control Sample (cont.)**

**Blocking on Selected Sample Characteristics (s.e.)**

| Sample | All | U74=1 | U74=0 | U75=1 | U75=0 |
|---|---|---|---|---|---|
| NSW | 1,794 | 2,692 | -685 | 1,711 | 1,691 |
|  | (633) | (722) | (1278) | (681) | (1289) |
| CPS-1: |  |  |  |  |  |
| Unadjusted | -8,498 | 1,710 | -10,735 | 1,949 | -9,198 |
|  | (712) | (647) | (1248) | (665) | (1058) |
| Adjusted | 738 | 3,189 | -3,275 | 2,812 | -1,132 |
|  | (547) | (699) | (982) | (728) | (844) |
| Stratifying on the Score: |  |  |  |  |  |
| -Unadjusted | 1,713 | 3,334 | -1,912 | 2,582 | 214 |
|  | (1115) | (1398) | (1085) | (1070) | (1334) |
| -Adjusted | 1,774 | 3,445 | -1,064 | 2,523 | -153 |
|  | (1152) | (1578) | (1340) | (1278) | (1061) |

Note: Adjusted training effect uses least squares regressions of Table 4.

# Stratification: Dehejia-Wahba example

**Table 9b: Treatment Effect for NSW Male and PSID Control Sample**

**Blocking on Selected Sample Characteristics (s.e.)**

| Sample | All | Black | Non-Black | Ndgree=1 | Ndgree=0 | Educ >=11 | Educ <11 |
|---|---|---|---|---|---|---|---|
| NSW | 1,794 | 2,029 | 803 | 1,154 | 3,192 | 3,085 | 402 |
| | (633) | (706) | (1331) | (696) | (1517) | (1033) | (753) |
| PSID-1: | | | | | | | |
| Unadjusted | -15,205 | -9,733 | -15,961 | -9,701 | -16,233 | -16,117 | -10,043 |
| | (1154) | (1002) | (3008) | (1148) | (2172) | (1623) | (1335) |
| | | | | | | | |
| Adjusted | 218 | 1,091 | -632 | 1,695 | 179 | 1,071 | 474 |
| | (866) | (916) | (2078) | (993) | (1569) | (1222) | (1165) |
| Stratifying on the Score: | | | | | | | |
| -Unadjusted | 1,509 | 1,486 | 2,880 | 1,667 | 1,137 | 1,806 | 1,381 |
| | (1823) | (2067) | (2366) | (2298) | (2907) | (2522) | (2163) |
| | | | | | | | |
| -Adjusted | 1,647 | 1,936 | .. | 1,826 | 1,435 | 1,001 | 1,694 |
| | (1862) | (2146) | | (2435) | (3425) | (2725) | (2305) |

Note: Adjusted training effect uses least squares regressions of Table 4.

# Stratification: Dehejia-Wahba example

| | | | Table 9b: Treatment Effect for NSW Male and PSID Control Sample (cont.) | | |
|---|---|---|---|---|---|
| | | | Blocking on Selected Sample Characteristics (s.e.) | | |
| **Sample** | **All** | **U74=1** | **U74=0** | **U75=1** | **U75=0** |
| NSW | 1,794 | 2,692 | -685 | 1,711 | 1,691 |
| | (633) | (722) | (1278) | (681) | (1289) |
| PSID-1: | | | | | |
| Unadjusted | -15,205 | -446 | -17,465 | 501 | -16,364 |
| | (1154) | (1499) | (2022) | (1364) | (1720) |
| Adjusted | 218 | 4,534 | -4,428 | 1,823 | -1527 |
| | (866) | (1702) | (1431) | (1813) | (1241) |
| Stratifying on the Score: | | | | | |
| -Unadjusted | 1,509 | 4,444 | -4,681 | 4,160 | -931 |
| | (1823) | (25000 | (1576) | (2449) | (3518) |
| -Adjusted | 1,647 | 4,408 | -3,285 | 4,935 | -2,854 |
| | (1862) | (2458) | (1939) | (2514) | (4883) |

# How to do it: Propensity score matching

Matching on the propensity score is essentially a weighting scheme, which determines what weights are placed on comparison units when computing the estimated treatment effect:

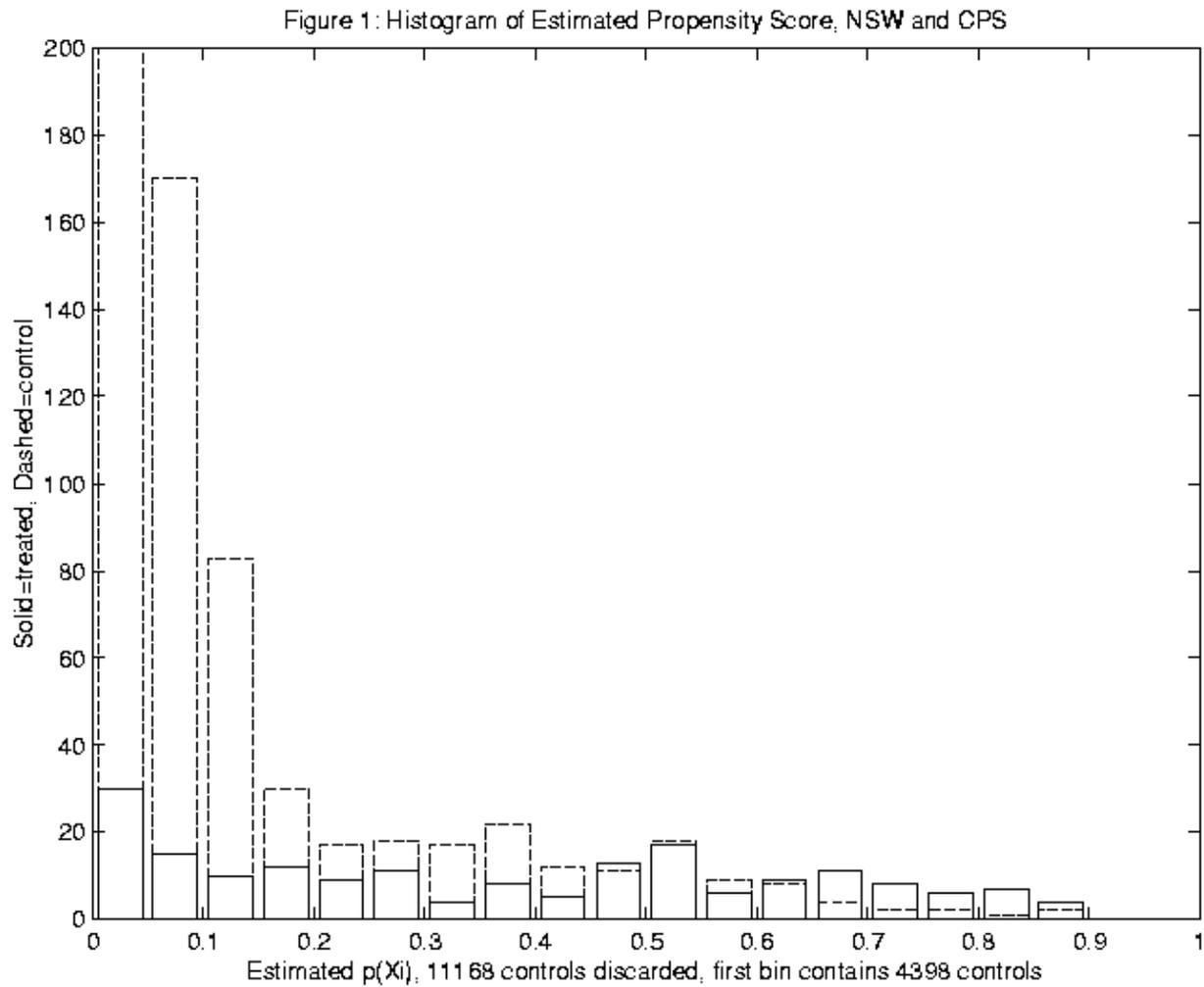$$\hat{\tau}\Big|_{T=1} = \frac{1}{|N|}\sum_{i\in N}\left(Y_i - \frac{1}{|J(i)|}\sum_{j\in J(i)}Y_j\right)$$

where *N* is the treatment group, |*N*| the number of units in the treatment group, *J(i)* is the set of comparison units matched to treatment unit *i*.
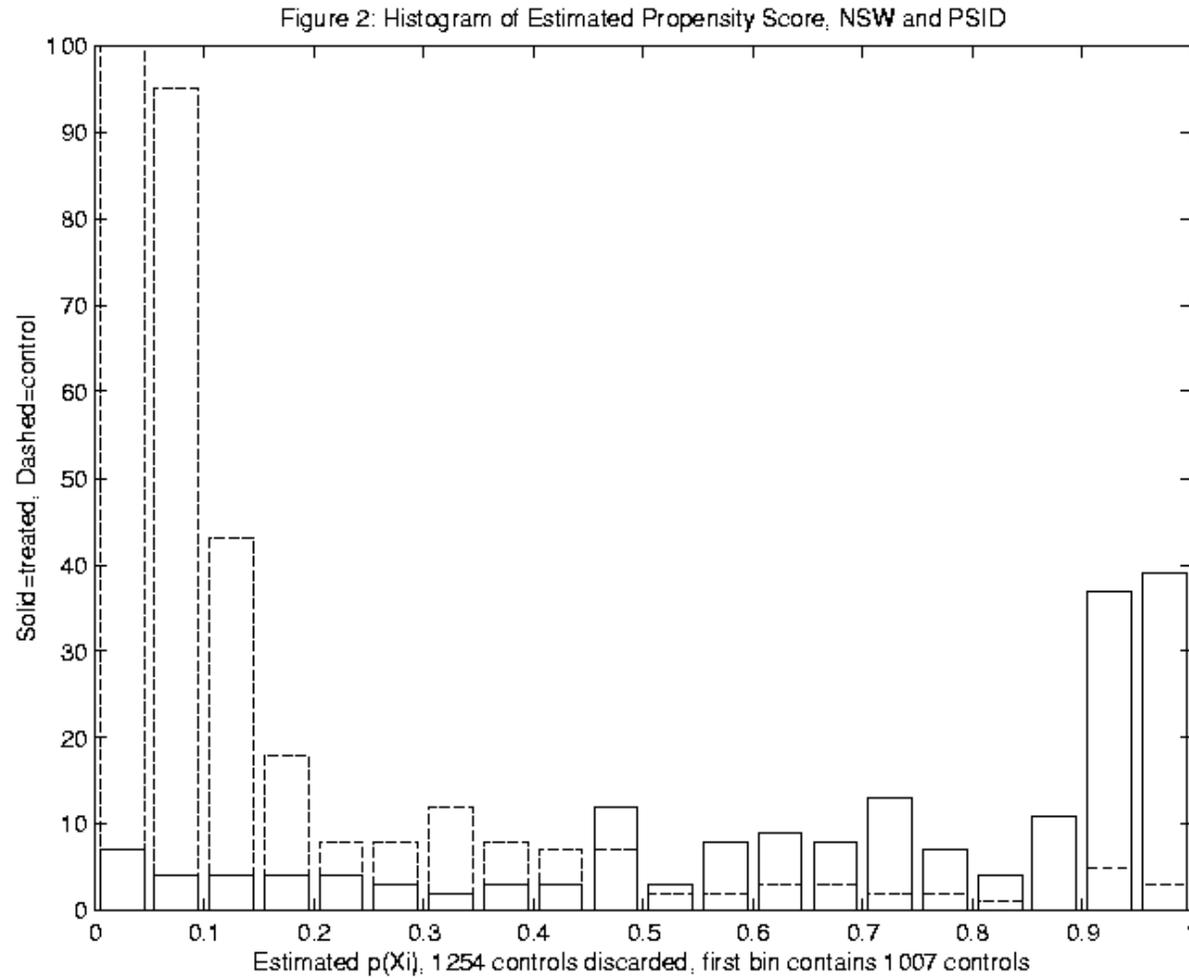
Issues:

? Whether or not to match with **replacement**

? How many **comparison units** to match to each treated unit

? Which **matching method** to choose

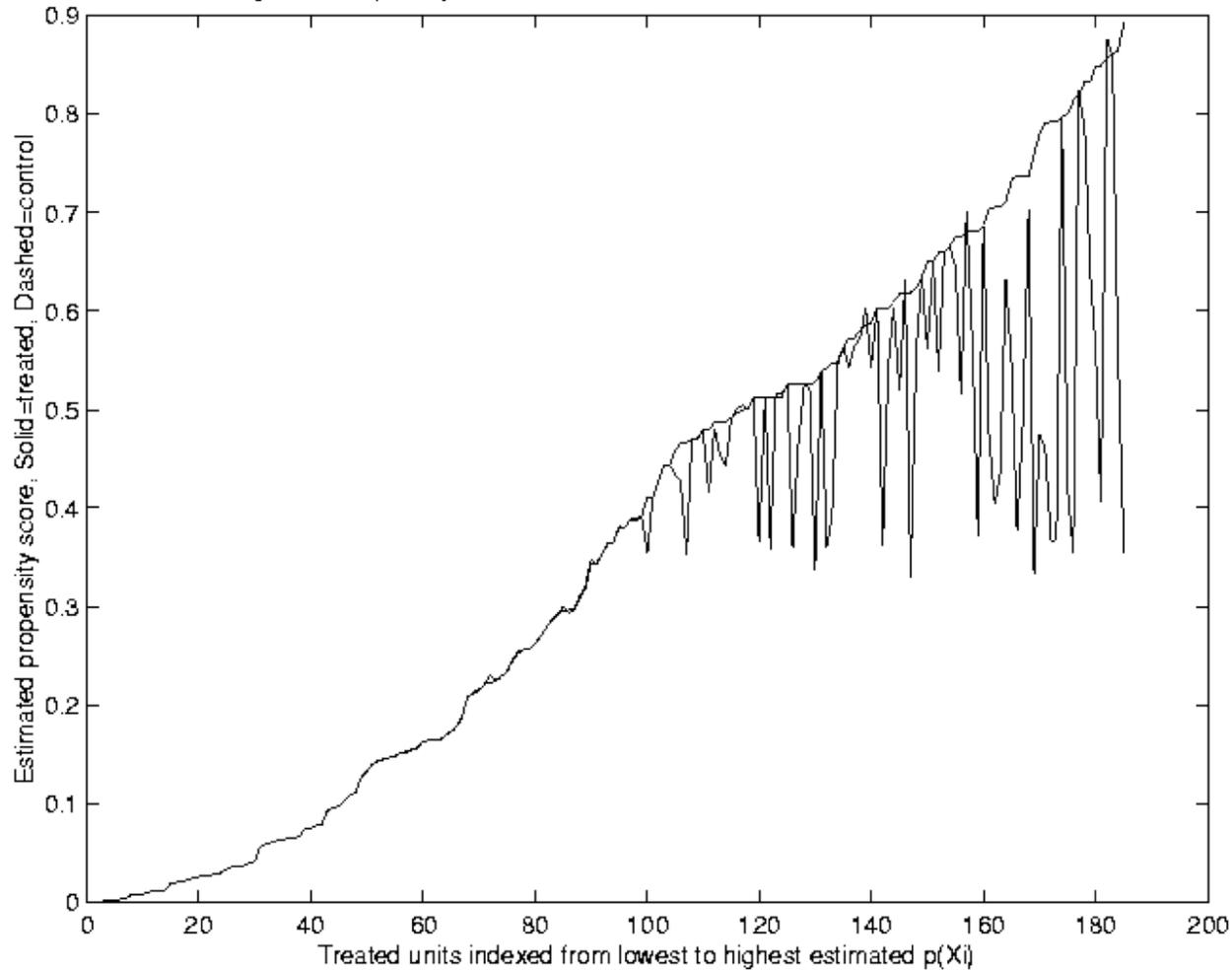# Propensity score matching: Dehejia-Wahba example



Figure 1: Histogram of Estimated Propensity Score, NSW and CPS

# Propensity score matching: Dehejia-Wahba example



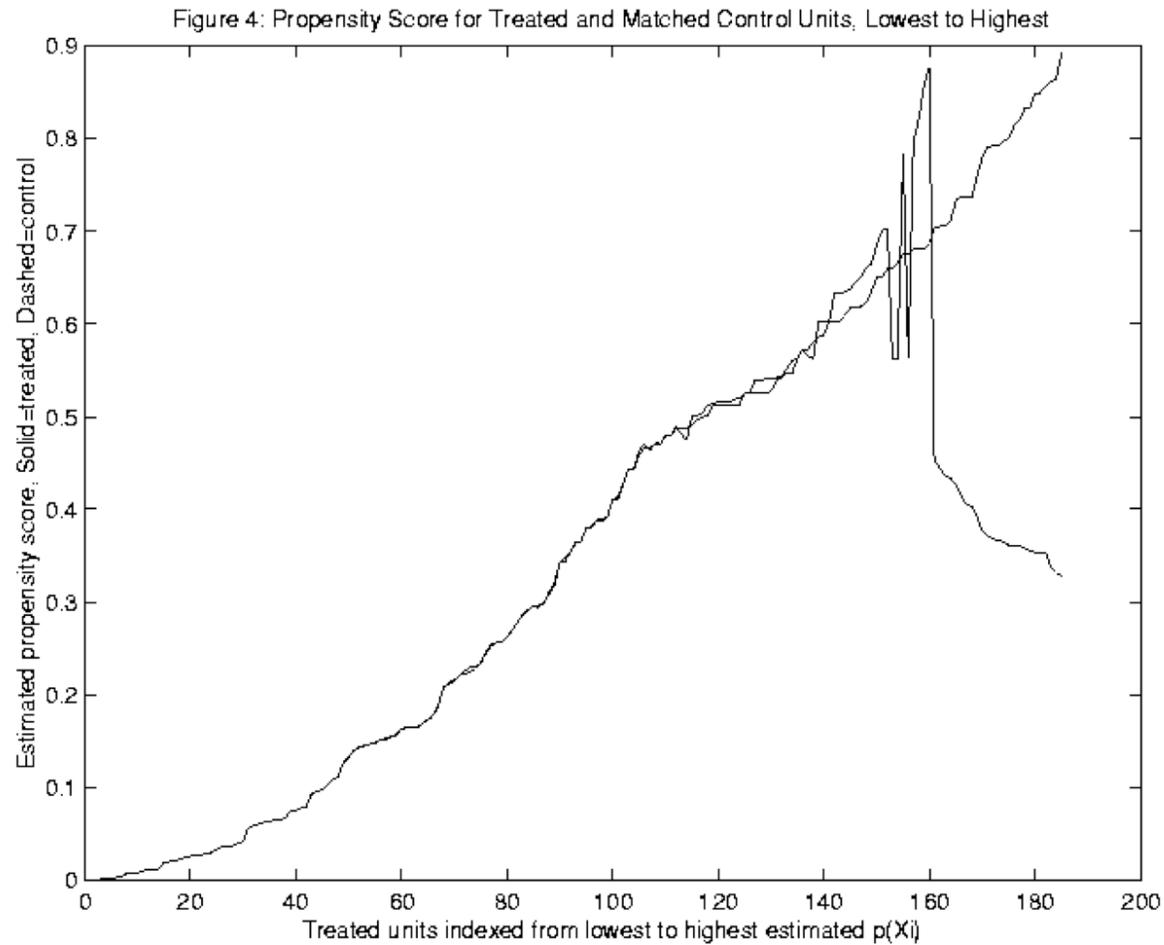Figure 2: Histogram of Estimated Propensity Score, NSW and PSID

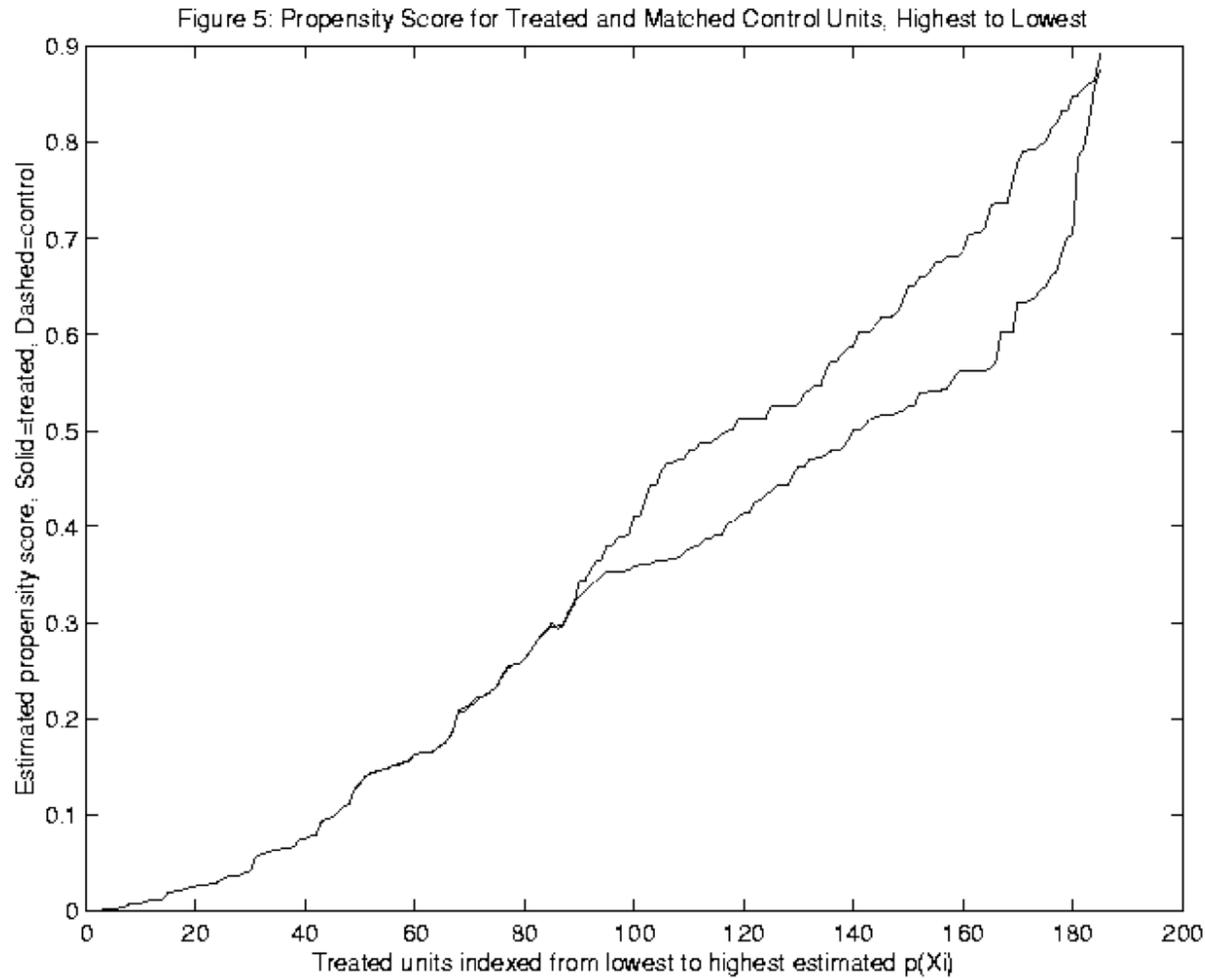# Propensity score matching: Dehejia-Wahba example



Figure 3: Propensity Score for Treated and Matched Control Units, Random

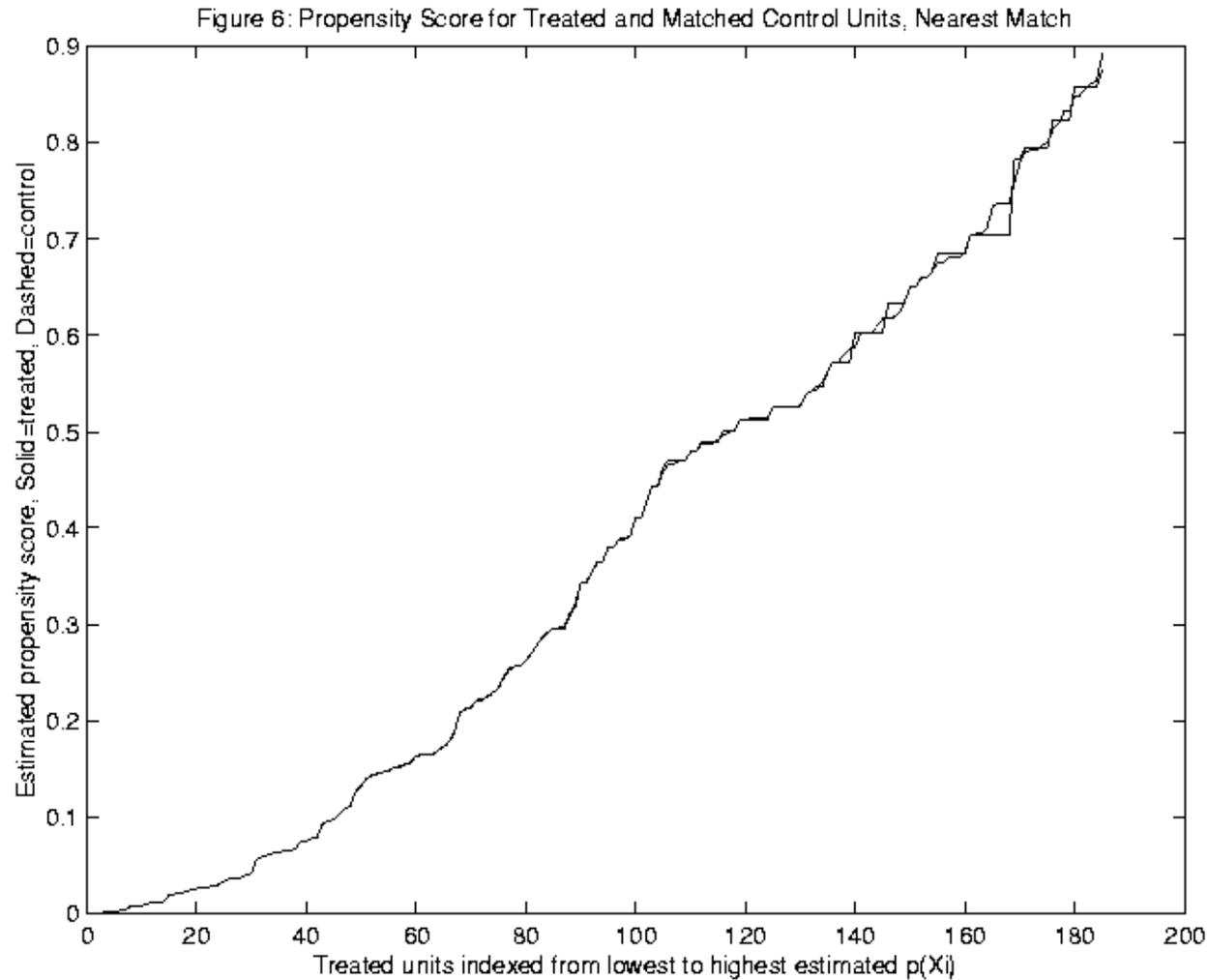# Propensity score matching: Dehejia-Wahba example



Figure 4: Propensity Score for Treated and Matched Control Units, Lowest to Highest

# Propensity score matching: Dehejia-Wahba example



Figure 5: Propensity Score for Treated and Matched Control Units, Highest to Lowest

# Propensity score matching: Dehejia-Wahba example



Figure 6: Propensity Score for Treated and Matched Control Units, Nearest Match

# Propensity score matching: Dehejia-Wahba example

Sample Means of Characteristics for Matched Control Samples

| Matched Samples | No. of Obs. | Age | School | Black | Hisp | Nodegree | Married | RE74 US$ | RE75 US$ |
|---|---|---|---|---|---|---|---|---|---|
| **NSW** | 185 | 25.81 | 10.35 | 0.84 | 0.06 | 0.71 | 0.19 | 2,096 | 1,532 |
| **MPSID-1** | 56 | 26.39 | 10.62 | 0.86 | 0.02 | 0.55 | 0.15 | 1,794 | 1,126 |
|  |  | [2.56] | [0.63] | [0.13] | [0.06] | [0.13] | [0.12] | [1406] | [1146] |
| **MPSID-2** | 49 | 25.32 | 11.10 | 0.89 | 0.02 | 0.57 | 0.19 | 1,599 | 2,225 |
|  |  | [2.63] | [0.83] | [0.14] | [0.08] | [0.16] | [0.16] | [1905] | [1228] |
| **MPSID-3** | 30 | 26.86 | 10.96 | 0.91 | 0.01 | 0.52 | 0.25 | 1,386 | 1,863 |
|  |  | [2.97] | [0.84] | [0.13] | [0.08] | [0.16] | [0.16] | [1680] | [1494] |
| **MCPS-1** | 119 | 26.91 | 10.52 | 0.86 | 0.04 | 0.64 | 0.19 | 2,110 | 1,396 |
|  |  | [1.25] | [0.32] | [0.06] | [0.04] | [0.07] | [0.06] | [841] | [563] |
| **MCPS-2** | 87 | 26.21 | 10.21 | 0.85 | 0.04 | 0.68 | 0.20 | 1,758 | 1,204 |
|  |  | [1.43] | [0.37] | [0.08] | [0.05] | [0.09] | 0.08 | [896] | [661] |
| **MCPS-3** | 63 | 25.94 | 10.69 | 0.87 | 0.06 | 0.53 | 0.13 | 2,709 | 1,587 |
|  |  | [1.68] | [0.48] | [0.09] | [0.06] | [0.10] | [0.09] | [1285] | [760] |

# Issues with matching

! When matching with replacement, standard non-parametric bootstrap standard errors are not valid

! Options are to use new asymptotic standard errors or to use the *m* of *n* bootstrap or subsampling

! Also note that bootstrap *does* work for matching without replacement

# Weighting on the propensity score

## Proposition

*If $Y_1, Y_0 \perp\!\!\!\perp T | X$, then*

$$\alpha_{ATE} = E[Y_1 - Y_0] = E\left[Y \cdot \frac{T - p(X)}{p(X) \cdot (1 - p(X))}\right]$$

$$\alpha_{ATET} = E[Y_1 - Y_0 | T = 1] = \frac{1}{\Pr(T = 1)} \cdot E\left[Y \cdot \frac{T - p(X)}{1 - p(X)}\right]$$

## Proof.

$$E\left[Y \frac{T - p(X)}{p(X)(1 - p(X))} \Big| X\right] = E\left[\frac{Y}{p(X)} \Big| X, T = 1\right] p(X) + E\left[\frac{-Y}{1 - p(X)} \Big| X, T = 0\right](1 - p(X))$$

$$= E[Y | X, T = 1] - E[Y | X, T = 0]$$

*The results follow from integration over $P(X)$ and $P(X | T = 1)$*

# Weighting on the propensity score
## *Why this works pt. 1*

**Population:** 50 men and 50 women.

**Outcome of interest:** All men have earnings of $10 and all women have earnings of $12

**Sampling:** For some reason you randomly sample men at probability 1/10 (from the male sample or 1/20 from the full sample) and women at probability 1/2 (from the female sample or ¼ from the full sample).

**Sample:** Your sample consists of 5 men (earnings 10) and 25 women (earnings 12).

**Population mean $\mu$:** You know the true population mean is ½ 10 + ½ 12 = 11.

**Sample mean $\bar{x}$:** A naïve (wrong) mean would be (10 x 5 + 25 x 12) / 30 = 11.83, because women are overrepresented in the sample.

**Adjustment:** How can you recover the true mean from the sample? Weight by the probability of sampling:

$$(5*10/(1/10) + 25*12 /(1/2))/(5/10+25/2) =11$$

# Weighting on the propensity score
## *Why this works pt. 2*

More generally, consider a population of $N$ units, made up of $k$ types (e.g., covariate values) of $n_k$ units each. Sample of $k=1...K$ units each of category with probality $p_k = p(n_k) = 1/n_k$.

Expectation of the sample mean as the sample gets bigger:

$$\frac{E\left(\sum \frac{x_i}{p_i}\right)}{\left(\sum \frac{1}{p_i}\right)} \rightarrow \frac{E\left(\sum n_k x_k\right)}{\left(\sum n_k\right)} \rightarrow \frac{E\left(\sum n_k x_k\right)}{N}$$

**What this means:** weighting each observation by its probability of selection adjusts its weight so it is in proportion to the overall popuation.

# Weighting on the propensity score
## *Why this works pt. 3*

$$\alpha_{ATET} = E\left[Y_1 - Y_0 \middle| T = 1\right] = \frac{1}{\Pr(T=1)} \cdot E\left[Y \cdot \frac{T - p(X)}{1 - p(X)}\right]$$

If T=1, weight $= \dfrac{1}{\Pr(T=1)}$.

If T=0, weight $= \dfrac{1}{\Pr(T=1)} \dfrac{p(X)}{1-p(X)} = \dfrac{1}{\Pr(T=1)} \dfrac{1}{(1-p(X))/p(X)}$

For the treated group (T=1), each observation gets equal weight and all observations together are reweighted by probability of being in the treated group.

For the comparison group (T=0), each observation is weighted (up) by the probability of being in the comparison group relative to the treatment group – and again the whole is reweighted to match the probability of being the treated group.

# Weighting on the propensity score
## *Why this works pt. 4*

$$\alpha_{ATET} = E\left[Y_1 - Y_0 \mid T = 1\right] = \frac{1}{\Pr(T=1)} \cdot E\left[Y \cdot \frac{T - p(X)}{1 - p(X)}\right]$$

$$= \frac{1}{\Pr(T=1)} \cdot \Pr(T=1) E\left[Y \cdot \frac{T - p(X)}{1 - p(X)} \mid T = 1\right] + \frac{1}{\Pr(T=1)} \cdot (1 - \Pr(T=1)) E\left[Y \cdot \frac{T - p(X)}{1 - p(X)} \mid T = 0\right]$$

$$= E\left[Y_1 \mid T = 1\right] - \frac{1}{\Pr(T=1)} \cdot (1 - \Pr(T=1)) E\left[Y \cdot \frac{p(X)}{1 - p(X)} \mid T = 0\right]$$

$$= E\left[Y_1 \mid T = 1\right] - \frac{1}{\Pr(T=1)} \cdot (1 - \Pr(T=1)) E\left[E_X\left[Y_0 \frac{p(X)}{1 - p(X)} \mid X, T = 0\right] \mid T = 0\right]$$

$$= E\left[Y_1 \mid T = 1\right] - \frac{1}{\Pr(T=1)} \cdot (1 - \Pr(T=1)) E\left[\frac{p(T=1)}{1 - p(T=1)} E_X\left[Y_0 \mid X, T = 1\right] \mid T = 0\right]$$

$$= E\left[Y_1 \mid T = 1\right] - \frac{1}{\Pr(T=1)} \cdot (1 - \Pr(T=1)) \frac{p(T=1)}{1 - p(T=1)} E\left[Y_0 \mid X, T = 1\right]$$

$$= E\left[Y_1 \mid T = 1\right] - E\left[Y_0 \mid T = 1\right]$$

**What's the intuition?**

✓ **For $Y_1$** (potential outcomes for treated group) it's a simple reweighting argument

✓ **For $Y_0$** (potential outcomes for comparison group) first you reweight to match the full comparison population, and then reweight the comparison to match the treated population.

# Weighting on the propensity score

$$\alpha_{ATE} = E[Y_1 - Y_0] = E\left[Y \cdot \frac{T - p(X)}{p(X) \cdot (1 - p(X))}\right]$$

$$\alpha_{ATET} = E[Y_1 - Y_0 | T = 1] = \frac{1}{Pr(T = 1)} \cdot E\left[Y \cdot \frac{T - p(X)}{1 - p(X)}\right]$$

The analogy principle suggests a two-step estimator:

1. Estimate the propensity score: $\hat{p}(X)$
2. Use estimated score to produce analogous estimators:

$$\hat{\alpha}_{ATE} = \frac{1}{N} \sum_{i=1}^{N} Y_i \cdot \frac{T_i - \hat{p}(X_i)}{\hat{p}(X_i) \cdot (1 - \hat{p}(X_i))}$$

$$\hat{\alpha}_{ATET} = \frac{1}{N_1} \sum_{i=1}^{N} Y_i \cdot \frac{T_i - \hat{p}(X_i)}{(1 - \hat{p}(X_i))}$$

# Weighting on the propensity score

$$\hat{\alpha}_{ATE} = \frac{1}{N} \sum_{i=1}^{N} Y_i \cdot \frac{T_i - \hat{p}(X_i)}{\hat{p}(X_i) \cdot (1 - \hat{p}(X_i))}$$

$$\hat{\alpha}_{ATET} = \frac{1}{N_1} \sum_{i=1}^{N} Y_i \cdot \frac{T_i - \hat{p}(X_i)}{(1 - \hat{p}(X_i))}$$

Standard errors:

- ✓ We need to adjust the s.e.'s for first-step estimation of $p(X)$
- ✓ Parametric first-step: Newey & McFadden (1994)
- ✓ Non-parametric first-step: Newey (1994)
- ✓ Or bootstrap the entire two-step procedure!

# Weighting on the propensity score

**Table 4: Estimated Training Effects for the NSW Male Participants Using Comparison Groups from PSID and CPS-SSA**

| | NSW EARNINGS LESS COMPARISON GROUP EARNINGS | | NSW TREATMENT EARNINGS LESS COMPARISON GROUP EARNINGS, CONDITIONAL ON THE ESTIMATED PROPENSITY SCORE | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | LINEAR WITH SCORE[g] | STRATIFYING ON THE SCORE | | | MATCHING ON THE SCORE | | WEIGHTING WITH THE SCORE[f] |
| | (1) | (2) | (3) | (4) | (5) | Obs.[h] (6) | (7) | (8) | (9) |
| | Unadjusted | Adjusted[a] | | Unadjusted | Adjusted[a] | | Unadjusted | Adjusted[e] | |
| NSW | 1,794 (633) | 1,672 (638) | | | | | | | |
| PSID-1[b] | -15,205 (1154) | 218 (866) | 542 (1197) | 1,509 (1823) | 1,647 (1862) | 1,255 | 2,190 (761) | 1,690 (1020) | 1,129 (2927) |
| PSID-2[c] | -3,647 (959) | 907 (1004) | 434 (1193) | 1,648 (2010) | 2,538 (2063) | 389 | 870 (977) | 826 (962) | 1,951 (1178) |
| PSID-3[c] | 1,069 (899) | 822 (1101) | 862 (1334) | 1,829 (2250) | 2,308 (2468) | 247 | 1,534 (1223) | 1,740 (1063) | 1,618 (1231) |
| CPS-1[d] | -8,498 (712) | 738 (547) | 893 (642) | 1,713 (1115) | 1,774 (1152) | 4,117 | 1,253 (988) | 1,174 (798) | 1,485 (3148) |
| CPS-2[d] | -3,822 (670) | 879 (654) | 399 (765) | 1,358 (1432) | 1,378 (1582) | 1493 | 1,445 (962) | 1,589 (946) | 862 (4059) |
| CPS-3[d] | -635 (657) | 1,326 (796) | 526 (891) | 1,335 (1765) | 1,023 (1956) | 514 | -466 (951) | -372 (945) | 379 (2308) |

# Weighting on the propensity score

- Traditional concern is that misspecification of estimated propensity score can lead to extreme weights.

- Buso, DiNardo, and McCrary argue that weighting performs well in small samples.

- Argue for normalizing the weight (standard option in most regression packages) as a way to deal with this.

- Can also be applied to diff-in-diffs estimators (Abadie 2005).

# Check sensitivity to $p(X)$ specification

**Table 10: Sensitivity of Estimated Training Effects for the NSW Male Participants to Specification of the Propensity Score**

| LOGIT BY SAMPLE | NSW EARNINGS LESS COMPARISON GROUP EARNINGS | | NSW COMPARISON TREATMENT EARNINGS LESS COMPARISON GROUP EARNINGS CONDITIONAL ON THE ESTIMATED PROPENSITY SCORE | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | LINEAR WITH SCORE[d] | STRATIFYING ON THE SCORE | | | MATCHING ON THE SCORE | | WEIGHTING WITH THE SCORE[c] |
| | | | | | | No. Obs.[e] | | | |
| | (1) Unadjusted | (2) Adjusted[a] | (3) | (4) Unadjusted | (5) Adjusted[a] | (6) | (7) Unadjusted | (8) Adjusted[b] | (9) |
| NSW | 1,794 (633) | 1,674 (638) | | | | 445 | | | |
| **PSID-1:** | | | | | | | | | |
| 1 | -15,205 (1154) | 218 (866) | 542 (1197) | 1,509 (1823) | 1,647 (1862) | 1,255 | 2,190 (761) | 1,690 (1020) | 1,129 (2927) |
| 2 | -15,205 (1154) | 105 (863) | -225 (1217) | 1,348 (1558) | 2,128 (1699) | 1,465 | 871 (988) | 795 (937) | 2,017 (2673) |
| 3 | -15,205 (1154) | 105 (863) | 463 (1080) | 1,044 (1087) | 136 (1226) | 1,373 | 2,124 (869) | 2,338 (842) | 2,125 (1570) |
| **CPS-1:** | | | | | | | | | |
| 4 | -8,498 (712) | 738 (547) | 893 (642) | 1,713 (1115) | 1,774 (1152) | 4,117 | 1,253 (988) | 1,174 (798) | 1,485 (3148) |
| 5 | -8,498 (712) | 684 (546) | 1,103 (614) | 1,485 (653) | 1,636 (682) | 6,365 | 1,179 (821) | 1,258 (897) | 1,414 (2221) |
| 6 | -8,498 (712) | 684 (546) | 1,147 (582) | 1,456 (595) | 1,728 (610) | 6,017 | | | 1,236 (1824) |

# Check sensitive to including RE74

**Table 11: Sensitivity of Estimated Training Effects for the NSW Male Participants When Dropping Pre-Treatment Earnings in 1974**

| | NSW EARNINGS LESS COMPARISON GROUP EARNINGS | | NSW COMPARISON TREATMENT EARNINGS LESS COMPARISON GROUP EARNINGS CONDITIONAL ON THE ESTIMATED PROPENSITY SCORE | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | LINEAR WITH SCORE[i] | STRATIFYING ON THE SCORE | | | MATCHING ON THE SCORE | | WEIGHTING WITH THE SCORE[h] |
| | | | | | | Obs[j] | | | |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) | (8) | (9) |
| | Unadjusted | Adjusted[a] | | Unadjusted | Adjusted[a] | | Unadjusted | Adjusted[g] | |
| NSW | 1,794 (633) | 1,631 (637) | | | | | | | |
| PSID-1[b] | -15,205 (1154) | -265 (880) | -1,085 (1110) | -47 (1442) | 635 (1455) | 1284 | 547 (1419) | 1,597 (1308) | -251 (3342) |
| PSID-2[c] | -3,647 (959) | 297 (1004) | -544 (1149) | 223 (1619) | -592 (1615) | 356 | 647 (1232) | 1,298 (1033) | 1,067 (1915) |
| PSID-3[d] | 1,069 (899) | 243 (1100) | 1,366 (1246) | 586 (1551) | 440 (1620) | 252 | 1,558 (1059) | 2,093 (940) | -355 (1478) |
| CPS-1[e] | -8,498 (712) | 525 (557) | 949 (632) | 1,294 (765) | 1,207 (880) | 4,558 | 1,827 (581) | 1,969 (808) | 1,198 (2748) |
| CPS-2[f] | -3,822 (670) | 371 (662) | 656 (706) | 1,475 (1054) | 1,661 (1101) | 1,222 | 924 (777) | 857 (813) | 2,238 (1528) |
| CPS-3[f] | -635 (657) | 844 (807) | 988 (860) | 1,044 (1417) | 1,129 (1509) | 504 | 1,325 (928) | 1,074 (942) | 1,241 (1494) |

# Regression estimators

Can we use regression estimators?

1. Least squares as an approximation

2. Least squares as a weighting scheme

3. Estimators of average treatment effects based on nonparametric regression

# 1. Least squares as an approximation

$(Y_1, Y_0) \perp\!\!\!\perp T \mid X$ implies that the conditional expectation $E[Y \mid T, X]$ can be interpreted as a conditional causal response function:

$$E[Y \mid T = 1, X] \;=\; E[Y_1 \mid T = 1, X] \;=\; E[Y_1 \mid X],$$

$$E[Y \mid T = 0, X] \;=\; E[Y_0 \mid T = 0, X] \;=\; E[Y_0 \mid X]$$

So $E[Y \mid T, X]$ provides average potential responses with and without the treatment.

The functional form of $E[Y \mid T, X]$ is typically unknown, but Ordinary Least Squares provides a well-defined approximation.

# 2. Least squares as a weighting scheme

Suppose that the covariates take on a finite number of values $K$. Then, from the subclassification section we know that:

$$\alpha_{ATE} = \sum_{k=1}^{K} \left( E\left[Y \middle| T = 1, X = x^k\right] - E\left[Y \middle| T = 0, X = x^k\right]\right) \Pr\left(X = x^k\right)$$

Now, we suppose that you run a regression that is **saturated** in $X$: a regression including one dummy variable, $Z^k$, for each possible value of $X$:

$$Y = \alpha T + \sum_{k=1}^{K} Z^k \beta_k + u,$$

where

$$Z^k = \begin{cases} 1 & \text{if } X = x^k \\ 0 & \text{if } X \neq x^k \end{cases}$$

# 2. Least squares as a weighting scheme

It can be shown that the coefficient $\hat{\alpha}_{\text{OLS}}$ of the saturated regression converges to:

$$\alpha_{\text{OLS}} = \sum_{k=1}^{K} \left( E\left[Y \middle| T = 1, X = x^k\right] - E\left[Y \middle| T = 0, X = x^k\right] \right) w_k$$

where

$$w_k = \frac{\text{var}\left(T \middle| X = x^k\right) \text{Pr}\left(X = x^k\right)}{\sum_{k=1}^{K} \text{var}\left(T \middle| X = x^k\right) \text{Pr}\left(X = x^k\right)}.$$

- Strata $k$ with a higher

$$\text{var}(T \mid X = x^k) = \text{Pr}(T \mid X = X^k)(1 - \text{Pr}(T \mid X = X^k))$$

(i.e. propensity scores close to 0.5) receive higher weight. Strata with propensity score close to 0 or 1 receive lower weight.

- OLS down-weights strata where the average causal effects are less precisely estimated.

# 3. Estimators based on nonparametric regression

$$\alpha_{ATE} = \int \left( E[Y|X, T=1] - E[Y|X, T=0] \right) dp(X)$$

$$\alpha_{ATET} = \int \left( E[Y|X, T=1] - E[Y|X, T=0] \right) dp(X|T=1)$$

This suggests the following estimators:

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

$$\hat{\alpha}_{ATET} = \frac{1}{N_1} \sum_{Ti=1}^{N} \left( \hat{\mu}_1(X_i) - \hat{\mu}_0(X_i) \right),$$

where $\hat{\mu}_1(\ )$ and $\hat{\mu}_0(\ )$ are nonparametric estimators of $E[Y|X, T=1]$ and $E[Y|X, T=0]$, respectively.

But estimating these regressions nonparametrically is difficult if the dimension of $X$ is large (curse of dimensionality again!).

# 3. Estimators based on nonparametric regression

It is enough, however, to adjust for the propensity score. So we can proceed in two steps:

1. Estimate the propensity score, $\hat{p}(X)$.

2. Calculate:

$$\hat{\alpha}_{ATE} = \frac{1}{N} \sum_{i=1}^{N} \left( \hat{V}_1 \left( \hat{p}(X_i) \right) - \hat{V}_0 \left( \hat{p}(X_i) \right) \right)$$

$$\hat{\alpha}_{ATET} = \frac{1}{N_1} \sum_{i=1}^{N} \left( \hat{V}_1 \left( \hat{p}(X_i) \right) - \hat{V}_0 \left( \hat{p}(X_i) \right) \right),$$

where $\hat{v}_1( \ )$ and $\hat{v}_0( \ )$ are nonparametric estimators of $E[Y|p(X), T = 1]$ and $E[Y|p(X), T = 0]$, respectively.

# Issues: what about efficiency?

**Question:** How can matching be efficient, if we are throwing away information?

**Answer:** It's not.

**Solution:**

1. **Kernel matching** (Heckman, Ichimura, and Todd): Match each treatment unit to all the comparison unit using a kernel weight, which uses a smoother set of weights on the comparison units (comparison units that are very far away get very low, rather than no, weight).

2. **Weighting by the propensity score**: Hirano, Imbens, and Ridder prove that weighting by the estimated propensity score is efficient. (Problem: efficiency pre-supposes a non-parametric series estimator of the propensity score. How is that any easier than a non-parametric series estimate of the outcome equation?)

# Assessing unconfoundedness: multiple control groups

Suppose we have a three-valued indicator $T_i \in \{0, -1, 1\}$ for the groups (e.g., ineligibles, eligible nonparticipants and participants), with the treatment indicator equal to

$W_i = 1\{T_i = 1\}$ so that

$$Y_i = \begin{cases} Y_i(0) & \text{if } T_i \in \{-1, 0\} \\ Y_i(1) & \text{if } T_i = 1. \end{cases}$$

Suppose we extend the unconfoundedness assumption to independence of the potential outcomes and the three-valued group indicator given covariates

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

# Assessing unconfoundedness: multiple control groups

Now a testable implication is

$$Y_i(0) \perp\!\!\!\perp 1\{T_i = 0\} \big| X_i, T_i \in \{-1, 0\}.$$

and thus

$$Y_i \perp\!\!\!\perp 1\{T_i = 0\} \big| X_i, T_i \in \{-1, 0\}.$$

An implication of this independence condition is being tested by the tests discussed above – should be no difference at Y in across the two non-treatment groups. Whether this test has much bearing on the unconfoundedness assumption depends on whether the extension of the assumption is plausible given unconfoundedness itself.

# Basic Procedure

✓ **Define "closeness":** the distance measure that will be used to determine whether one individual is a good match for another. (e.g. propensity score ranges, Mahalanobis distance, exact matching...)

✓ Implement a **matching method** (e.g. nearest neighbor, with or without replacement...)

✓ **Assess the quality of the matching**: typically looking at balance on covariates (e.g. graph density of propensity scores by treatment, control, matched, unmatched; test for statistically significant differences in mean variance between groups)

✓ **Analyze the data**: run your regression (reweight to get ATE or ATET; you should still control for covariates within subclassifications!)

# Conclusions

(1) Propensity score methods succeed in reducing the dimensionality problem of dealing with observables.

(2) Practically speaking, when we have a large comparison group, and we are looking for the best matches to the treatment group, this method helps us to zoom in on the right subset of comparison units.

(3) By reducing dimensionality, it becomes very easy to allow for a flexible functional form.

(4) Economists are guilty of using functional forms that are "excessively" linear and expecting thee functional forms to successfully extrapolate between very different groups.

(5) This method only works if we observe all (enough) of the variables that determine the selection process.

(6) In training programs, earnings seem to be crucial.

# Current Research and Discussion

- For health policy folks: Brown DW, DeSantis SM, Greene TJ, et al. [A novel approach for propensity score matching and stratification for multiple treatments: Application to an electronic health record-derived study](). Stat Med. 2020

- For international monetary policy folks: "[The effectiveness of macroprudential policies and capital controls against volatile capital inflows]()" from the Bank for International Settlements (June 2020)

- For machine learning and/or labor econ folks: Goller, Daniel & Lechner, Michael & Moczall, Andreas & Wolff, Joachim, 2019. "[Does the Estimation of the Propensity Score by Machine Learning Improve Matching Estimation? The Case of Germany's Programmes for Long Term Unemployed]()," Labour Economics. August 2020

# Additional Resources

- #EconTwitter! Prof. Nick Huntington-Klein has [another animated graphic](#) for matching

- Some [accessible info](#) on conducting PSM in R (if of interest, not needed for this course)

- Very [high-level overview](#) of the PSM concept from the Urban Institute

- Prof. Elizabeth Stuart's (Johns Hopkins School of Public Health) [presentation on PSM](#) for the 2011 Society of Prevention Research conference (workshop slides)

# During Recitation…

✓ Discuss main skills and questions related to PS4

✓ Discussion of progress on replication exercise

✓ Chat about how propensity score matching is evolving and is being used in applied research in industry (example from energy regulation)

✓ Review of readings if time allows