

Empirical Analysis III

The University of Chicago

April 11, 2019

Content: Week 2

- Today we focus on observational data, a single cross-section
 - What can we learn?
-
- How to select controls and argue exogeneity
 - How to perform matching
 - Relation between matching and regression

Selection

Formal definition

- There is **selection** into the treatment state D if

$\underbrace{Y_d|D=d}_{\text{observable}}$ is distributed differently from $\underbrace{Y_d|D=d'}_{\text{unobserved}}$ for $d' \neq d$

- This is not the case under the random assignment assumption
- Expected to occur if agents choose D with knowledge of $\{Y_d\}_{d \in \mathcal{D}}$

Selection is common

- Particularly concerning if you are trained in neoclassical economics
- Agents choose a job training program ($D \in \{0, 1\}$) to max utility
- Utility will incorporate expected future earnings (Y_0, Y_1)
- Agents who choose job training might do so because of low Y_0
- Data typically supports this story ("Ashenfelter's (1978) dip")
- Alternatively, might choose $D = 0$ because of high Y_0

Selection Bias

- Consider the simple treatment/control mean contrast under selection
- This contrast would be the ATE under random assignment
- Decompose the contrast into a causal effect and selection bias:

$$\mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0]$$

$$= \underbrace{(\mathbb{E}[Y_1|D = 1] - \mathbb{E}[Y_0|D = 1])}_{\text{ATT}} + \underbrace{(\mathbb{E}[Y_0|D = 1] - \mathbb{E}[Y_0|D = 0])}_{\text{selection bias}}$$

$$\text{or } = \underbrace{(\mathbb{E}[Y_1|D = 0] - \mathbb{E}[Y_0|D = 0])}_{\text{ATU}} + \underbrace{(\mathbb{E}[Y_1|D = 1] - \mathbb{E}[Y_1|D = 0])}_{\text{selection bias}}$$

- First term is the causal effect for those who were treated/untreated
- Under random assignment would have $\text{ATT} = \text{ATU} = \text{ATE}$
- Second term is how the treated would have been different anyway
Under random assignment this would be 0
- The first expression is more natural if thinking of $D = 0$ as baseline

Selection on Observables

- A simple relaxation of random assignment is **selection on observables**
- Suppose that we observe (Y, D, X) where X are covariates
- The selection on observables assumption is that

$$\{Y_d\}_{d \in \mathcal{D}} \perp\!\!\!\perp D | X$$

- Says: Conditional on X , treatment is as-good-as randomly assigned
- Other terms: **unconfoundedness, ignorable treatment assignment**
- Underlies causal interpretations of linear regression
We will look into this connection more later

Identification under Selection on Observables

Identification argument

- Conditional version of random assignment:

$$F_d(y|x) \equiv \mathbb{P}[Y_d \leq y | X = x]$$

$$= \mathbb{P}[Y_d \leq y | D = d, X = x] = \mathbb{P}[Y \leq y | D = d, X = x]$$

- Second equality requires the **overlap condition**: $\mathbb{P}[D = d | X = x] > 0$
- By averaging over x , one can point identify the marginals

$$F_d(y) \equiv \mathbb{P}[Y_d \leq y] = \mathbb{E}(\mathbb{P}[Y_d \leq y | X]) = \mathbb{E}(\mathbb{P}[Y \leq y | D = d, X])$$

- Treatment effects are then constructed from $D = d_1 \rightarrow d_2$ contrasts

Thought experiment: An RCT given $X = x$

- Fix an $X = x$
- Find treated (d_1) and control agents (d_2) with $X = x$ ("match" on $X = x$)

Identification of Mean Contrasts

- Suppose $D \in \{0, 1\}$ is binary — by far the most common case
- Using essentially the same argument as on the previous page:

$$\text{ATE} \equiv \mathbb{E}[\mathbb{E}[Y_1|X] - \mathbb{E}[Y_0|X]] = \mathbb{E}[\mathbb{E}[Y|D=1, X] - \mathbb{E}[Y|D=0, X]]$$

- Similar expressions for the ATT and ATU have an **important difference**:

$$\text{ATT} = \mathbb{E}[Y|D=1] - \mathbb{E}[\mathbb{E}[Y|D=0, X]|D=1]$$

$$\text{ATU} = \mathbb{E}[\mathbb{E}[Y|D=1, X]|D=0] - \mathbb{E}[Y|D=0]$$

- Helps in estimation since only one **conditional expectation** (more later)
- Note that only **mean independence** is needed for these arguments:

$$\mathbb{E}[Y_d|D=0, X] = \mathbb{E}[Y_d|D=1, X]$$

- Difficult to think of arguments for mean (without full) independence

Conditional independence and choice of controls

When might conditional independence hold?

- You have detailed information about the assignment mechanism
- You have very rich data: personal traits, histories, etc.

Only condition on predetermined observables

- For selection on observables to be plausible, X should be **predetermined**
- Usually this really is a temporal issue (measured before vs. after D)
- Intuition is clear - we want to condition on selection into treatment

Examples

- Suppose we accidentally included Y as part of X
- Then clearly we aren't going to have $(Y_0, Y_1) \perp\!\!\!\perp D | X$. Thus:
- Don't include earnings after the program in X
- Don't include employment after the program in X
- Don't include marital status after the program in X

What happens if your controls are affected by treatment?

Neale and Johnson (1996, JPE)

- What is the role of premarket factors in Black-White wage gap?
- They regress adult earnings (Y) on race ($D = 1$ if white) and some covariates, with and without test scores T
- A key finding is that:
$$\mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0] >> \mathbb{E}[Y|D = 1, T = 1] - \mathbb{E}[Y|D = 0, T = 1]$$

- Potential outcome model for earnings, test scores and race:

$$Y = DY_1 + (1 - D)Y_0$$

$$T = DT_1 + (1 - D)T_0$$

- For simplicity, abstracts from covariates and assume:

$$D \perp\!\!\!\perp (Y_1, Y_0, T_1, T_0)$$

What happens if your controls are affected by treatment?

- Unconditional regression:

$$\mathbb{E}[Y|D = 1] - \mathbb{E}[Y|D = 0] = \mathbb{E}[Y_1 - Y_0]$$

$$\mathbb{E}[T|D = 1] - \mathbb{E}[T|D = 0] = \mathbb{E}[T_1 - T_0]$$

- Conditional on test scores:

$$\begin{aligned} & \mathbb{E}[Y|D = 1, T = 1] - \mathbb{E}[Y|D = 0, T = 1] \\ &= \mathbb{E}[Y_1|D = 1, T_1 = 1] - \mathbb{E}[Y_0|D = 0, T_0 = 1] \\ &= \underbrace{\mathbb{E}[Y_1|T_1 = 1] - \mathbb{E}[Y_0|T_0 = 1]}_{\text{independence}} \\ &= \underbrace{\mathbb{E}[Y_1 - Y_0|T_1 = 1]}_{\text{causal effect conditional on test score}} + \underbrace{(\mathbb{E}[Y_0|T_1 = 1] - \mathbb{E}[Y_0|T_0 = 1])}_{\text{selection bias}} \end{aligned}$$

- Selection bias would be negative if blacks need to be smarter to achieve the same test score, despite their disadvantaged background

Does bias go down if you control for more?

- Suppose that $(Y_0, Y_1) \perp\!\!\!\perp D | X$
- Let X_2 be a subset of X
- Let X_1 be a subset of X_2
- So controlling on X is the most information, and X_1 is the least
- Suppose however that we only have X_2 (hence X_1) in our data

- Perhaps surprisingly, using X_2 can be more biased than using X_1
- That is, adding information (X_2 vs X_1) need not reduce bias
- If $X_2 = X$, then it does, but not more generally
- Point is not well-appreciated but should be concerning
- Means we really need to have the correct X
- Suggests one should be careful with automated model selection
- Insight seems to come from Heckman and Navarro-Lozano (2004)

Propensity Score

Definition

- Binary treatment case, $D \in \{0, 1\}$
- $p(x) \equiv \mathbb{P}[D = 1|X = x]$ is called the **propensity score**
- Let $P \equiv p(X)$ be the random variable $\mathbb{P}[D = 1|X]$

Rosenbaum and Rubin (1983) sufficiency argument

- Result: Selection on observables implies $(Y_0, Y_1) \perp\!\!\!\perp D | p$
- Can you prove it (see problem set)?
- The result implies that we can condition on $p(X)$ instead of X
- Still need overlap, but now with $p(X)$ (scalar) instead of X (vector)

Propensity Score Weighting

- Using p , we can write the ATE as a weighted average of Y
- The following is derived in the supplemental notes:

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

- Average over X to point identify

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

- similar expression can be derived for the ATT and ATU
- can you derive these expressions (see problem set)?

Why Have Different Identification Arguments?

Summarizing identification under selection on observables

- Three different constructive identification result for the ATE
- Match on X , match on P , weight using p
- Each one shows that ATE is point identified
- And they are all derived under the same assumptions

So why have three?

- Identification arguments directly inform the construction of estimators
- Different arguments suggest different estimators
- In general, these different estimators may have different properties:
Efficiency, rates of convergence, finite sample performance...
- See e.g. Imbens (2010, 2015)

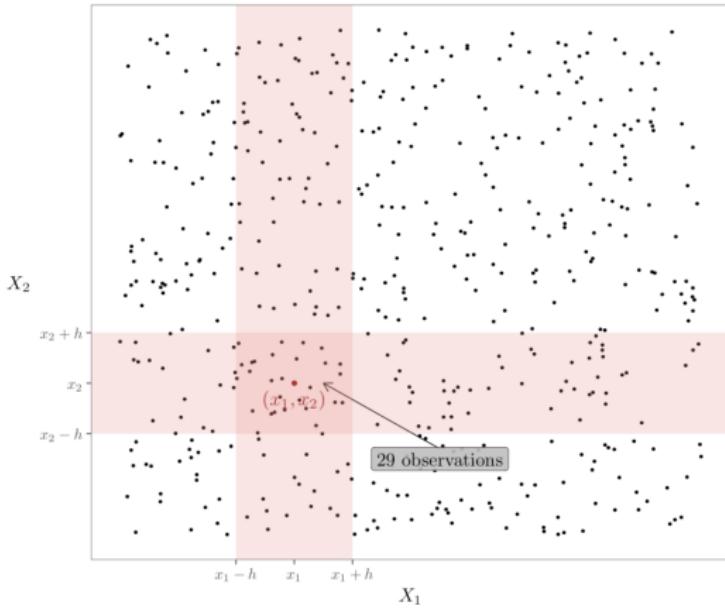
Curse of Dimensionality

- Matching requires the estimation of $\mathbb{E}[Y|D = 0, X]$
- Doing this non-parametrically is typically not feasible:
Curse of Dimensionality
- e.g. $K=30$ binary covariates $\implies \min N = 2^{30+1} = 2147483648$

Overcoming the curse

- Regression analysis
- Match on Propensity score
- Propensity score weighting
- Inexact matching: metric on X (Mahalanobis)

Curse of Dimensionality



- Observations used to estimate $\mathbb{E}[Y|X = x_1, X = x_2]$ with bandwidth h
- Effective number of observations drops by an order of magnitude

Matching vs linear regression

- Consider the following saturated-in- X_i regression model:

$$Y_i = \sum_{k=1}^K d_{ik}\alpha_k + \beta D_i + U_i$$
$$d_{ik} = \mathbb{1}[X_i = x_k]$$

- α_k is the coefficient for $X_k = x_k$
- β is the regression estimand for treatment
- Note: no interactions between X_i and D_i

Regression as pseudo-matching

- Then it can be shown that:

$$\beta = \sum_{k=1}^K \{\mathbb{E}[Y_i | D_i = 1, X_i = x_k] - \mathbb{E}[Y_i | D_i = 0, X_i = x_k]\} w_k$$
$$w_k = \frac{\mathbb{P}[D_i = 1 | X_i = x_k](1 - \mathbb{P}[D_i = 1 | X_i = x_k])\mathbb{P}[X_i = x_k]}{\sum_{k=1}^K \mathbb{P}[D_i = 1 | X_i = x_k](1 - \mathbb{P}[D_i = 1 | X_i = x_k])\mathbb{P}[X_i = x_k]}$$

- weights depend on the conditional variance of treatment status

Linear regression and matching

- Differ in weighting of $\mathbb{E}[Y_i | D_i = 1, X] - \mathbb{E}[Y_i | D_i = 0, X]$
- Same if treatment effects do not vary with X
- Biased with selection on unobservable factors of non-treatment level

Matching and ATT

- Suppose X takes on values x_1, \dots, x_K
- Then the matching estimator $\hat{\beta}^{ATT}$ can be written as:

$$\sum_{k=1}^K \{ \mathbb{E}[Y_i | D_i = 1, X_i = x_k] - \mathbb{E}[Y_i | D_i = 0, X_i = x_k] \} \mathbb{P}[X_i = x_k | D_i = 1]$$

- Where $\mathbb{P}[X_i = x_k | D_i = 1]$ is the probability mass function for X_i given $D_i = 1$
- That is, partition the treated and control sample in K cells by X
 - ▶ Calculate the mean outcome difference in each cell
 - ▶ Take a weighted average of the mean differences
 - ▶ Using the fraction of treated observations in each cell as the weights

Matching and ATE

- Suppose X takes on values x_1, \dots, x_K
- Then the matching estimator $\hat{\beta}^{ATE}$ can be written as:

$$\sum_{k=1}^K \{ \mathbb{E}[Y_i | D_i = 1, X_i = x_k] - \mathbb{E}[Y_i | D_i = 0, X_i = x_k] \} \mathbb{P}[X_i = x_k]$$

- Where $\mathbb{P}[X_i = x_k]$ is the probability mass function for X_i in the whole population
- That is, partition the treated and control sample in K cells by X
 - ▶ Calculate the mean outcome difference in each cell
 - ▶ Take a weighted average of the mean differences
 - ▶ Using the fraction of the total sample in each cell as the weights

Numerical example

- Suppose that X is binary and $\mathbb{P}[X = 1] = 0.5$, $\mathbb{P}[D_i = 1|X = 0] = 0.9$ and $\mathbb{P}[D_i = 1|X = 1] = 0.5$
- Recall Bayes' rule:

$$\mathbb{P}[X_i = x_k | D_i = 1] = \mathbb{P}[D_i = 1 | X_i = x_k] \mathbb{P}[X_i = x_k] / \mathbb{P}[D_i = 1]$$

- Then, we get:

$$\begin{aligned}\beta^{ATT} &= 0.64(\mathbb{E}[Y_i | D_i = 1, X_i = 0] - \mathbb{E}[Y_i | D_i = 0, X_i = 0]) \\ &\quad + 0.36(\mathbb{E}[Y_i | D_i = 1, X_i = 1] - \mathbb{E}[Y_i | D_i = 0, X_i = 1])\end{aligned}$$

$$\begin{aligned}\beta^{ATE} &= 0.5(\mathbb{E}[Y_i | D_i = 1, X_i = 0] - \mathbb{E}[Y_i | D_i = 0, X_i = 0]) \\ &\quad + 0.5(\mathbb{E}[Y_i | D_i = 1, X_i = 1] - \mathbb{E}[Y_i | D_i = 0, X_i = 1])\end{aligned}$$

$$\begin{aligned}\beta^{OLS} &= 0.26(\mathbb{E}[Y_i | D_i = 1, X_i = 0] - \mathbb{E}[Y_i | D_i = 0, X_i = 0]) \\ &\quad + 0.74(\mathbb{E}[Y_i | D_i = 1, X_i = 1] - \mathbb{E}[Y_i | D_i = 0, X_i = 1])\end{aligned}$$

How Propensity Score matching works

- Where before we had troubles estimating

$$\mathbb{E}[Y|D = 0, X]$$

- Now we can estimate

$$\mathbb{E}[Y|D = 0, p(X)]$$

- This is just a non-parametric regression of Y on 1 dimensional $p(X)$

- There is a little catch: we need to estimate $p(X)$
- This reintroduces the curse of dimensionality in the estimation of $p(X)$

How Propensity Score matching works

- Estimate $p(X)$
 - ▶ Choice of X : all the attributes X influencing both D and Y_0, Y_1
guidance: economic theory, a priori considerations, institutional set-up, previous literature
plausibility of matching depends on the choice of X
sensitivity of results to the set of regressors used
 - ▶ Estimation: e.g. Probit or Logit
 - ▶ Note: we are not interested in the behavioral interpretation, all we want:

$$X \perp\!\!\!\perp D | p(X)$$

- Match treated to controls
- Check Common Support
- Check balancing of X 's on CS, if it holds → Calculate treatment effect
- Reiterate until balancing condition holds

Matching Estimators Overview

- Pair to each treated individual i some group of “comparable” non-treated individuals
- Associate to the outcome Y of treated individual, a matched outcome \hat{Y} given by the (weighted) outcomes of the “neighbors” in the comparison group:

$$\hat{Y}_0 = \sum_{j \in C^0} w_j Y_j$$

- C^0 - set of neighbours of treated individual in the $D=0$ group
- w_j - weight of non-treated j in forming a comparison with treated i , where $\sum_{j \in C^0} w_j = 1$
- Weights differ for each treated individual

Classes of Matching Estimators

- Traditional matching estimators
 - ▶ one-to-one (nearest neighbour) matching
- Simple smoothed matching estimators
 - ▶ K-nearest neighbours
- Weighted smoothed matching estimators
 - ▶ kernel-based matching
 - ▶ local linear regression-based matching
- Trade-off between matching quality and variance

One-to-one matching

- To each treated unit, match only one non-treated unit:

$$C^0 = \{j : |p(X_i) - p(X_j)| = \min_{k \in \{D=0\}} |p(X_i) - p(X_k)|\}$$

$$\forall m \in \{D = 0\} : w_m = \mathbb{1}[m = j]$$

With replacement

- Many treated units may be matched to the same non-treated unit
- Less bias, more variance

Without replacement

- Once a non-treated unit has been matched, it cannot be used again
- More bias, less variance

Simple smoothed matching

K-nearest neighbours

- $C^0 = \{ \text{the } K \text{ units in } D = 0 \text{ with } p(X_k) \text{ closest to } p(X_i) \}$

$$w_m = \begin{cases} 1/K & m \in C^0 \\ 0 & o/w \end{cases}$$

- With or without replacement

Radius matching

- $C^0 = \{ \text{all } j : |p(X_i) - p(X_j)| < \delta \}$

$$w_m = \begin{cases} 1/|C^0| & m \in C^0 \\ 0 & o/w \end{cases}$$

Kernel matching

- The matched outcome for treated i is a kernel - weighted average of the non - treated outcomes, where the weights are in proportion to the non - treated closeness to i :

$$\hat{Y}_0 = \frac{\sum_{j \in \{D=0\}} K\left(\frac{p(X_i) - p(X_j)}{h}\right) Y_j}{\sum_{j \in \{D=0\}} K\left(\frac{p(X_i) - p(X_j)}{h}\right)} = \sum_{j \in \{D=0\}} w_j Y_j$$

- Non-treated j 's outcome Y_j is weighted by

$$w_j = \frac{K\left(\frac{p(X_i) - p(X_j)}{h}\right)}{\sum_{j \in \{D=0\}} K\left(\frac{p(X_i) - p(X_j)}{h}\right)}$$

Non-parametric background

- For each treated i , estimate $\hat{Y}_0 \equiv [Y|D = 0, p(X) = p(X_i)]$ non-parametrically:
 - Fit a constant estimated on a local neighbourhood of $p(X_i)$
 - Apply a weighting scheme with weights K_j , for $p = p(X)$:

$$\min_{\theta_0} \sum_{j \in C^0} (Y_j - \theta_0)^2 K\left(\frac{p(X_i) - p(X_j)}{h}\right)$$

$\implies \hat{Y}_0$ is approximated locally by a constant:

$$\hat{Y}_0 = \hat{\theta}_0 = \sum_j \frac{K_j}{\sum_j K_j} Y_j$$

Non-parametric background

Kernel choice

- Gaussian

$$K(u) \propto \exp\left(-\frac{u^2}{2}\right)$$

uses all the non-treated units

- Epanechnikov

$$K(u) \propto \begin{cases} (1 - u^2) & |u| < 1 \\ 0 & o/w \end{cases}$$

has a moving window within the $D = 0$ group

Bandwidth h selection

- Bias-variance trade-off
- $h \uparrow$, more “tolerant” in terms of closeness of matches

Local linear regression matching

- For each treated i , estimate $\hat{Y}_0 \equiv \mathbb{E}[Y|D = 0, p(X) = p(X_i)]$ non-parametrically:
 - ▶ Fit a line estimated on a local neighbourhood of $p(X_i)$
 - ▶ Apply a weighting scheme with weights:

$$\min_{\theta_0, \theta_1} \sum_{j \in C^0} (Y_j - \theta_0 - \theta_1(p(X_i) - p(X_j)))^2 K\left(\frac{p(X_i) - p(X_j)}{h}\right)$$

Mahalanobis-metric matching

Alternative to propensity score matching:

- Combine the X 's into a distance measure and then match on the resulting scalar:

$$d(i,j) = (X_i - X_j)^T V^{-1} (X_i - X_j)$$

with V the pooled within-sample covariance matrix of X

- Properties:
 - Unit free
 - Assigns weight to each w in proportion to the inverse of the variance of w
 - reduces differences in X within matched pairs in all directions
 - Uses $d(i,j)$ instead of $|p(X_i) - p(X_j)|$

Matching in practice using LaLonde data

- NSW was a randomized experiment that assigned people to training positions
- mid 1970s:
 - ▶ AFDC women, ex-drug addicts, ex-criminal offenders, high school drop-outs of both sexes
 - ▶ guaranteed a job for 9 to 18 months
- Baseline earnings and demographics (1975)
- Post treatment earnings demographics (1979)
- LaLonde (AER,1986) compared experimental with non-experimental estimates
- For now, we will just use this data to see how matching can be done in practice
- But get to know paper and data – we will return to this again several times

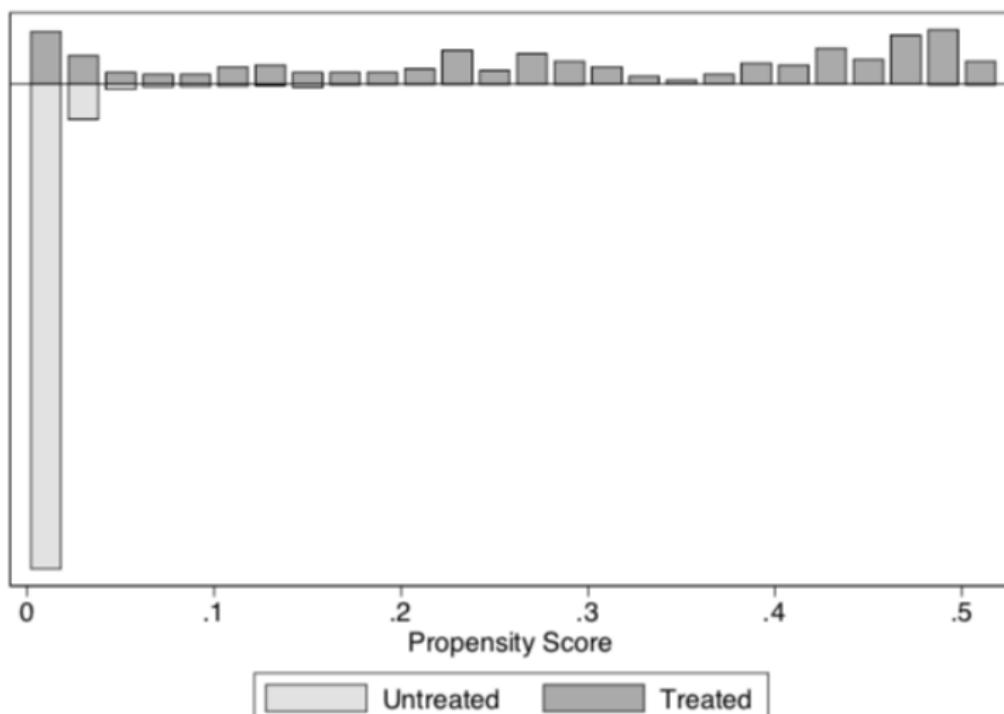
Estimating the Propensity Score

- Consider treated units from NSW sample - 297 obs., and control units from CPS sample - 15992 obs.

```
. probit treatment age education black hispanic married nodegree earnings75  
  
Probit regression  
Number of obs      =      16289  
LR chi2(7)        =     1421.68  
Prob > chi2       =     0.0000  
Pseudo R2         =     0.4791  
  
Log likelihood = -772.77819  
  
-----  
treatment |      Coef.    Std. Err.      z     P>|z| [95% Conf. Interval]  
-----+-----  
age |   -.0106229   .0043183    -2.46    0.014   -.0190866   -.0021591  
education |   .0220374   .0194534     1.13    0.257   -.0160905   .0601653  
black |   1.918753   .0834945    22.98    0.000    1.755107    2.0824  
hispanic |   .8917917   .116267     7.67    0.000    .6639125   1.119671  
married |   -.6246106   .091131    -6.85    0.000   -.803224   -.4459972  
nodegree |   .5700337   .1043422     5.46    0.000    .3655266   .7745407  
earnings75 |   -.0000518   6.25e-06    -8.29    0.000   -.000064   -.0000395  
.cons |   -2.513416   .3092365    -8.13    0.000   -3.119508   -1.907324  
-----
```

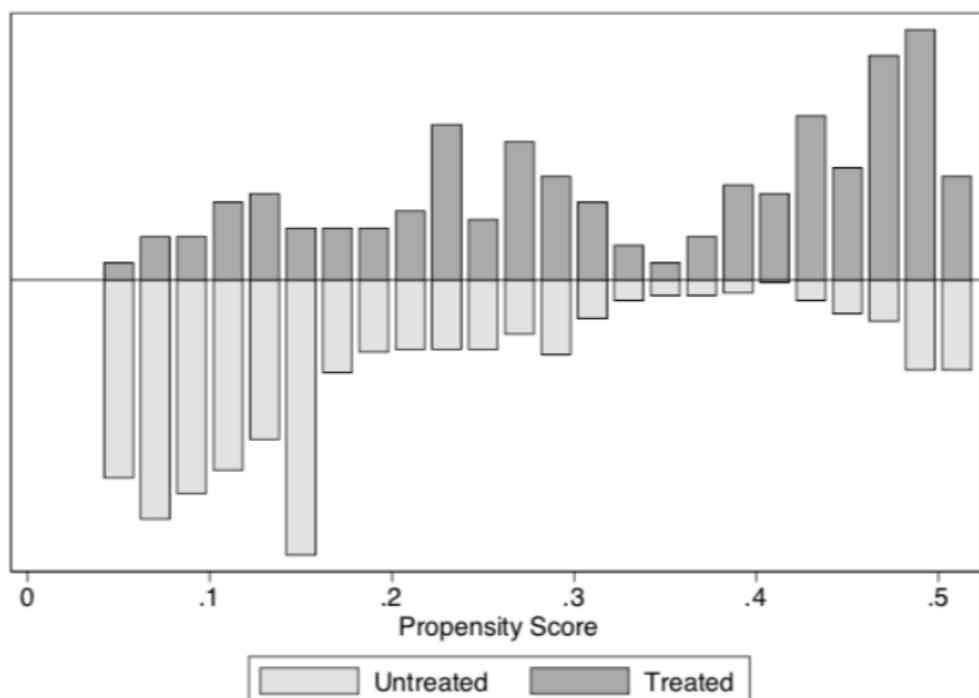
Estimating the Propensity Score

- Checking common support: vertical axis - density



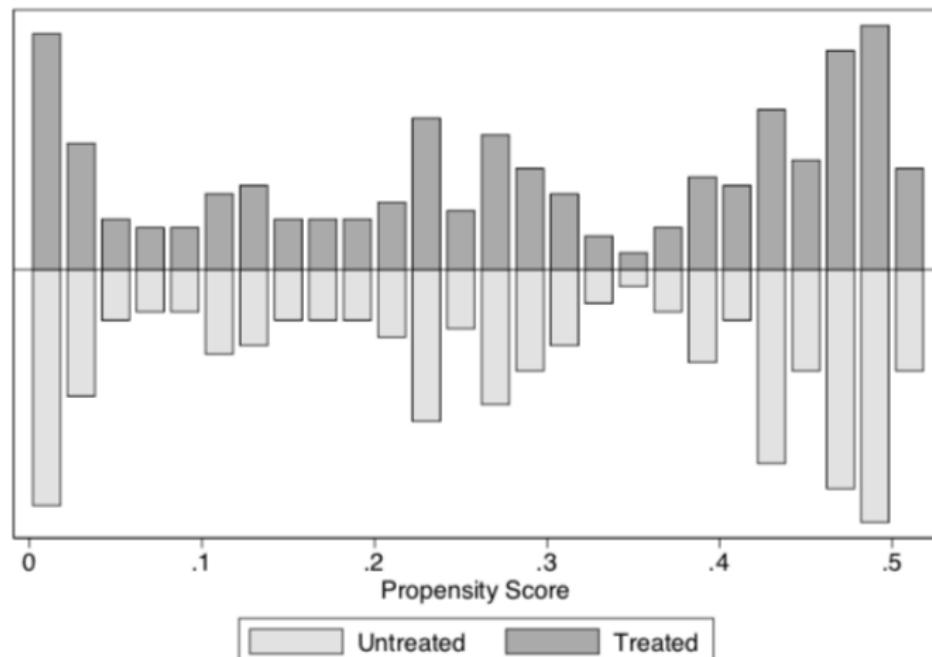
Estimating the Propensity Score

- Checking common support: trimming



One-to-One Matching with Replacement

- Match the treated units with closest (in $p(X)$ metrics) untreated units:
193 out of 15992 untreated units left



Assessing Matching Quality

- Check balancing of observables:

$$D \perp\!\!\!\perp X | p(X)$$

- For each variable compute standardized % bias before and after matching:

$$B_{\text{before}}(X) = 100 \frac{\bar{X}_1 - \bar{X}_0}{\sqrt{(V_1(X) + V_0(X))/2}}$$

$$B_{\text{after}}(X) = 100 \frac{\bar{X}_{1M} - \bar{X}_{0M}}{\sqrt{(V_{1M}(X) + V_{0M}(X))/2}}$$

- Achieved % reduction in |bias|
- t-tests for equality of means before and after matching

Assessing Matching Quality

```
. pstest age education black hispanic married nodegree earnings75
```

Variable	Sample	Mean		%reduct	t-test		
		Treated	Control	%bias	bias	t	p> t
age	Unmatched	24.626	33.225	-94.2		-13.37	0.000
	Matched	24.626	24.923	-3.2	96.6	-0.43	0.666
education	Unmatched	10.38	12.028	-68.5		-9.85	0.000
	Matched	10.38	10.36	0.8	98.8	0.11	0.914
black	Unmatched	.80135	.07354	215.6		47.04	0.000
	Matched	.80135	.83838	-11.0	94.9	-1.17	0.241
hispanic	Unmatched	.09428	.07204	8.1		1.47	0.143
	Matched	.09428	.07071	8.5	-6.0	1.04	0.297
married	Unmatched	.16835	.71173	-130.7		-20.54	0.000
	Matched	.16835	.17172	-0.8	99.4	-0.11	0.913
nodegree	Unmatched	.73064	.29584	96.5		16.27	0.000
	Matched	.73064	.68013	11.2	88.4	1.35	0.178
earnings75	Unmatched	3066.1	14017	-144.2		-19.67	0.000
	Matched	3066.1	3362.9	-3.9	97.3	-0.76	0.446

Assessing Matching Quality

Overall measures

- Summary indicators of the distribution of $|bias|$ before and after matching
- Pseudo R^2 from Probit of treatment on covariates before matching and on matched samples
- P-values of the likelihood-ratio test of joint insignificance of covariates before and after matching

Common support

- % of treated lost?
- Compare to those off-the-support: covariates, Y

Assessing Matching Quality

- Probit on matched sample: note the decrease in pseudo R^2 comparing to the unmatched sample

```
. probit treatment age education black hispanic married nodegree earnings75 [iw=_w]

Probit regression                                         Number of obs     =      490
                                                               LR chi2(7)      =     6.61
                                                               Prob > chi2    =    0.4707
Log likelihood = -408.42502                           Pseudo R2       =    0.0080

-----  
treatment |      Coef.    Std. Err.      z     P>|z|    [95% Conf. Interval]  
-----+-----  
age |   .0006298   .0067915    0.09    0.926    -.0126814    .013941  
education |   .0476399   .0324732    1.47    0.142    -.0160065    .1112863  
black |  -.1498752   .1765545   -0.85    0.396    -.4959156    .1961653  
hispanic |   .0943199   .2484845    0.38    0.704    -.3927008    .5813405  
married |  -.0069582   .1480393   -0.05    0.963    -.2971098    .2831935  
nodegree |   .3281476   .161253     2.03    0.042    .0120976    .6441975  
earnings75 |  -9.52e-06   .0000111   -0.86    0.392    -.0000313    .0000123  
.cons |  -.5939974   .516367   -1.15    0.250    -.1606058    .4180634  
-----  
. testparm *
```



```
chi2( 7) =      6.56  
Prob > chi2 =    0.4756
```

Assessing Matching Quality

Iterate the following until balancing is achieved:

- Estimate propensity score
 - ▶ Specification, Probit/Logit, probability/index/odds ratio
- Match
 - ▶ Metric: $X, \hat{p}(X), \{X, \hat{p}(X)\}$
 - ▶ Type of matching
 - ▶ Smoothing parameters
- Common support
- Assessment of matching quality

Matching Estimate

```
. psmatch2 treatment age education black hispanic married  
> nodegree earnings75, out(earnings78)
```

Variable	Sample	Treated	Controls	Difference	S.E.	T-stat
-----+-----+-----+-----+-----+-----+-----						
earnings78	Unmatched	5976.35202	13650.8035	-7674.4515	540.709664	-14.19
	ATT	5976.35202	4349.83269	1626.51934	677.626206	2.40

Note: S.E. does not take into account that the propensity score is estimated.

	psmatch2:
psmatch2:	Common
Treatment	support
assignment	On suppor Total
-----+-----+-----+-----+-----+-----+-----	
Untreated	15,992 15,992
Treated	297 297
-----+-----+-----+-----+-----+-----+-----	
Total	16,289 16,289

Reducing imbalances

- Although matching should go a long way towards balancing the X between the treated and comparisons
- We can perform a final regression correction for remaining imbalances

$$\min_{\alpha, \beta, \delta} \sum_{i=1}^N w_i (Y_i - \alpha - \underbrace{\delta D_i}_{\text{parameter of interest}} - X'_i \beta)^2$$

- Where w_i 's are the matching weights

Post-matching regression adjustment

```
. reg earnings78 treatment age education black hispanic married nodedgree ///
earnings75 [iw=_weight]
```

Source	SS	df	MS	Number of obs	=	594
Model	3.4317e+09	8	428964551	F(8, 585)	=	12.88
Residual	1.9480e+10	585	33299154.1	Prob > F	=	0.0000
Total	2.2912e+10	593	38636967.3	R-squared	=	0.1498

earnings78	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
treatment	1686.398	476.1792	3.54	0.000	751.1692	2621.627
age	68.52453	31.14949	2.20	0.028	7.346074	129.703
education	447.4218	148.3188	3.02	0.003	156.1195	738.724
black	-640.7179	812.5813	-0.79	0.431	-2236.65	955.2141
hispanic	-115.6389	1135.11	-0.10	0.919	-2345.027	2113.749
married	1703.38	679.2145	2.51	0.012	369.3842	3037.376
nodedgree	830.5569	741.3742	1.12	0.263	-625.5223	2286.636
earnings75	.3565398	.0514142	6.93	0.000	.255561	.4575186
_cons	-3504.447	2366.024	-1.48	0.139	-8151.382	1142.488

Examples of studies

Let's look at some studies assuming selection on obs.

- ▶ LaLonde (1986) and the subsequent papers. You should know the arguments of this (heated) debate
- ▶ Fagereng et al. (2019). The authors try to use knowledge of treatment assignment to argue selection on observables
- ▶ Angrist (1998). Much cited paper for regression vs matching discussion.
But why would selection on obs. hold here?

What are we looking for?

- ▶ What is the argument for selection on observables?
- ▶ How is it implemented?

Selection on observables and information set

Behavioral restrictions

- ▶ Recall that selection on observables imply that (Y_1, Y_0) or Y_0 do not determine (or correlate with) D conditional on X
- ▶ Rules out selection into the program based on unobserved (by the analyst) outcomes (or factors correlate with the outcomes)

Key question:

- ▶ What is the information set of the agent making participation decision (taking or assigning treatment)?
- ▶ What is the information set of the analyst?
- ▶ Is there asymmetry in information? If so, who knows more?

Selection on obs. requires the analyst to know (or assume) the information the agent uses when deciding treatment participation

The Many Ways to Skin a Statistical Cat

What estimator should I choose?

- Many different ways to implement selection on observables. The literature is a bit overwhelming (at least to me)
- Unfortunate truth about statistics is that optimality results are hard
- Even so, some choices may be better than others
Theory plus judgment and trial and error may take you a long way
(If you don't want to get your hands dirty, you can always do theory ...)

Standard errors

- We haven't even touched on standard error estimates
- Many of the estimators are quite complicated with multiple steps
- Bootstrap can be invalid for the non-smooth ones (e.g. matching)
- Why? It's a bit complicated but (my) intuition is that matching is like an extreme order statistic – bootstrap doesn't work then

Discrete Nonparametric Estimation

Binning estimator

- Suppose $X \in \{x_1, \dots, x_K\}$ is discrete with K small relative to N
- This covers cases where X has multiple discrete components e.g.
 $X = \{(male, white), (female, white), (male, nonwhite), (female, nonwhite)\}$
- Then a nonparametric **binning estimator** is very natural: e.g.

$$\hat{\mu}_d(x) = \frac{1}{N_{d,x}} \sum_{i:D_i=d, X_i=x}^N Y_i \text{ where } N_{d,x} \equiv \sum_{i=1}^N \mathbb{1}[D_i = d, X_i = x]$$

Limitations

- Only works if X is completely discrete, otherwise $N_{d,x} = 0$ or 1!
- If X is continuous, could bin it into discrete chunks
→ This can be a bit arbitrary
- Poor finite sample performance if **small bins** (K large relative to N)

Nonparametric Smoothing Regression

Smoothing

- Classical approach to continuous X is to use $X \approx x$ to **smooth**
- Relies on μ_d being a sufficiently smooth function of x
- Old, extremely well-developed, and enormous literature
- Two main approaches considered in economics are very intuitive

Kernel regression

- Take a sample mean of Y over $D_i = d, X_i \in [x - h, x + h]$
- $h > 0$ is a **bandwidth** parameter that needs to be chosen
- Bias goes down as $h \rightarrow 0$, variance goes up

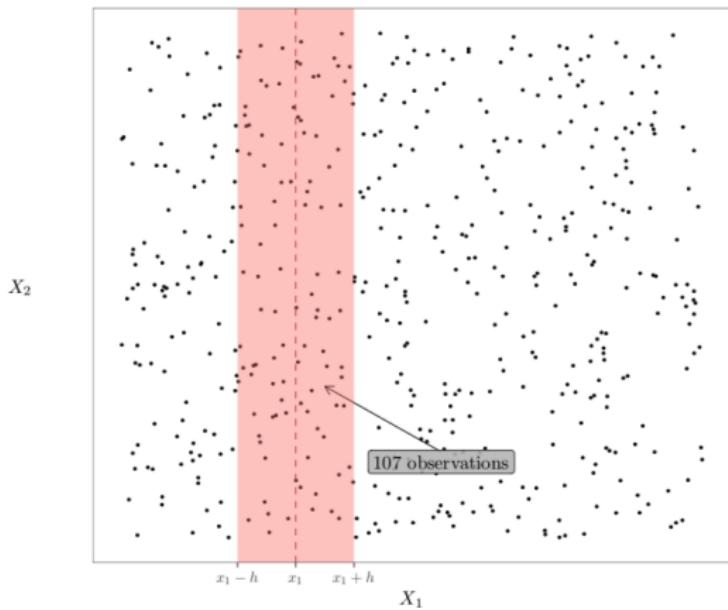
Series/sieve approximations

- Write $\mu_d(x) = \sum_{k=1}^K \theta_k b_k(x)$ for some **basis functions** b_k
- For example, regressing Y on $1, X, X^2, X^3, \dots, X^K$
- Bias goes down as $K \rightarrow \infty$, variance goes up

The Curse of Dimensionality

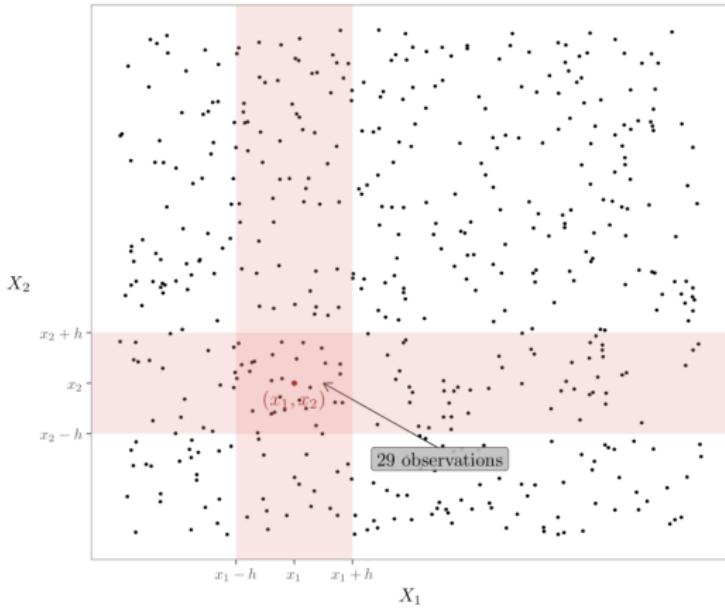
- Smoothing regression is held back by **the curse of dimensionality**
 - Estimator quality rapidly deteriorates with the dimension of X
 - Affects both kernel and sieve methods (indeed, all rate-optimal methods)
-
- Formally, the **rate of convergence** of the estimators goes down
 - Implication is that each new dimension requires **much** more data
 - Kernel with $X \in [0, 1]$ uniformly has $2h \times N$ observations within $\pm h$ of x
 - Kernel with $X \in [0, 1]^2$ uniformly has only $4h^2 \times N$ observations
 - Given h^* need **much larger** N to maintain the same “effective” N
-
- Each $\hat{\mu}_d$ will have a slow rate of convergence, but \widehat{ATE} can still be \sqrt{N}
→ Intuition is that by averaging we are again using all N observations
 - This can be misleading “asymptopia” - practical value is questionable
 - Finite sample performance of \widehat{ATE} will still be poor in high dimensions

Curse of Dimensionality



- Kernel regression with 500 draws from a bivariate uniform $[0, 1]$
- Observations used to estimate $\mathbb{E}[Y|X = x_1]$ with bandwidth h

Curse of Dimensionality



- Observations used to estimate $\mathbb{E}[Y|X = x_1, X = x_2]$ with bandwidth h
- Effective number of observations drops by an order of magnitude

Smoothing with the propensity score

- Recall that selection on observables implies $(Y_0, Y_1) \perp\!\!\!\perp D | p(X)$
- So instead of estimating $\mu_d(x)$, we could estimate

$$\nu_d(p) \equiv \mathbb{E}[Y | D = d, P = p]$$

- Appears to break the curse of dimensionality ...

Estimating p

- An immediate practical problem here is p needs to be estimated by \hat{p}
- However now we see that dimension reduction was an illusion
→ Same curse of dimensionality in nonparametrically estimating p !
- Still, parametric p is arguably better than parametric μ_d
→ p is just a matter of fit, whereas μ_d is a counterfactual object
- In practice, usually see a logit (maybe probit) estimator of p

Subclassification (Blocking)

- One could use kernel or sieve methods for $\nu_d(p)$
- Instead, more common to use blocking on the propensity score
- Just a particular type of nonparametric smoothing regression

Blocking

- Divide $[0, 1]$ into $\{b_0, b_1, \dots, b_J\}$ with $b_0 = 0, b_J = 1$
- Define $B_j = 1$ if $p(X) \in (b_{j-1}, b_j)$ as membership in block j
- If $b_j - b_{j-1}$ is small then roughly random assignment within block
- Estimate $\widehat{ATE}_j = Y_{1,j} - Y_{0,j}$ per block, i.e. conditional on $B_j = 1$
- Then average \widehat{ATE}_j by block size into \widehat{ATE}
- Key question is how to construct the blocks
- Imbens (2015) suggests an algorithm based on testing $D \perp\!\!\!\perp X | \{B_j\}_{j=1}^J$
 $\rightarrow D \perp\!\!\!\perp X | p$ implied by selection on observables - so check within blocks

Blocking with Linear Regression

Combining two approaches

- Imbens (2015) suggests combining blocking with linear regression
- First construct the blocks
- Then within each block, run a linear regression Y on $1, D, X$
- Estimate regression-adjusted ATE for each block, then average up

Why?

- Intuitively, this could potentially reduce both bias and variance
- The variance part is clear if accounting for X reduces variation in Y
- The bias part is less clear (i.e. not necessarily true) ...
 - Recall that linear regression extrapolates if $\bar{X}_1 \neq \bar{X}_0$
 - However within each block $\bar{X}_1 \approx \bar{X}_0$ - little extrapolation
- Adjusting for X reduces remaining differences within blocks
 - But presumably the remaining differences should be small anyway?

Job Training: Background

The National Supported Work (NSW) Demonstration

- Publicly and privately funded training program in the 1975-1977
- Provided 6 - 18 months work experience to disadvantaged individuals
- Applicants were randomly assigned to treatment or control
- Typically focus on post December 1975 \implies 1975 earnings predet.

The LaLonde (1986) critique

- Want to know impact of being offered training ($D = 1$) on applicants
- Applying/being offered job training is usually not randomly assigned
- In the NSW it is (for those who apply) $\implies \mathbb{E}[Y_1 - Y_0 | D = 1]$ identified
- Suppose we didn't have this - then we need to infer $\mathbb{E}[Y_0 | D = 1]$
- Could try this with (e.g.) PSID/CPS using a variety of methods
- How close do we get to the experimental benchmark?
- **Answer:** Not very close \implies extremely influential paper
 \rightarrow Findings were challenged (e.g. Heckman & Hotz, 1987)

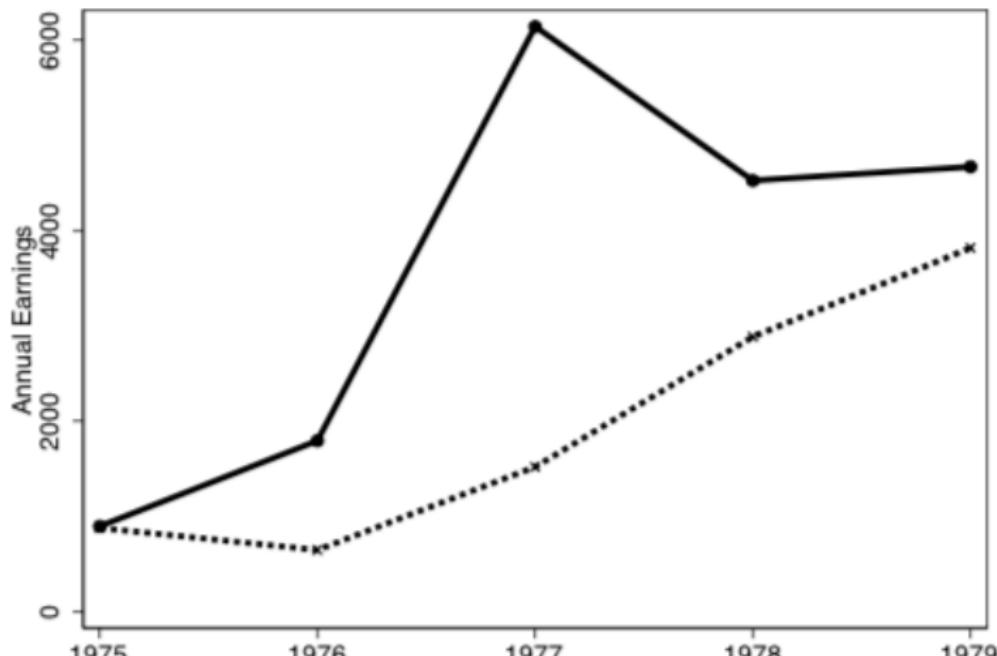
Balancing

TABLE 1—THE SAMPLE MEANS AND STANDARD DEVIATIONS OF PRE-TRAINING EARNINGS AND OTHER CHARACTERISTICS FOR THE NSW AFDC AND MALE PARTICIPANTS

Variable	Full National Supported Work Sample			
	AFDC Participants		Male Participants	
	Treatments	Controls	Treatments	Controls
Age	33.37 (7.43)	33.63 (7.18)	24.49 (6.58)	23.99 (6.54)
Years of School	10.30 (1.92)	10.27 (2.00)	10.17 (1.75)	10.17 (1.76)
Proportion High School Dropouts	.70 (.46)	.69 (.46)	.79 (.41)	.80 (.40)
Proportion Married	.02 (.15)	.04 (.20)	.14 (.35)	.13 (.35)
Proportion Black	.84 (.37)	.82 (.39)	.76 (.43)	.75 (.43)
Proportion Hispanic	.12 (.32)	.13 (.33)	.12 (.33)	.14 (.35)
Real Earnings 1 year Before Training	\$393 (1,203)	\$395 (1,149)	1472 (2656)	1558 (2961)
Real Earnings 2 years Before Training	\$854 (2,087)	\$894 (2,240)	2860 (4729)	3030 (5293)
Hours Worked 1 year Before Training	90 (251)	92 (253)	278 (466)	274 (458)
Hours Worked 2 years Before Training	186 (434)	188 (450)	458 (654)	469 (689)
Month of Assignment (Jan. 78 = 0)	-12.26 (4.30)	-12.30 (4.23)	-16.08 (5.97)	-15.91 (5.89)
Number of Observations	800	802	2083	2193

Note: The numbers shown in parentheses are the standard deviations and those in the square brackets are the standard errors.

Annual Earnings - Females



Source: LaLonde (1986), Table 2

Choice of comparison group - Females

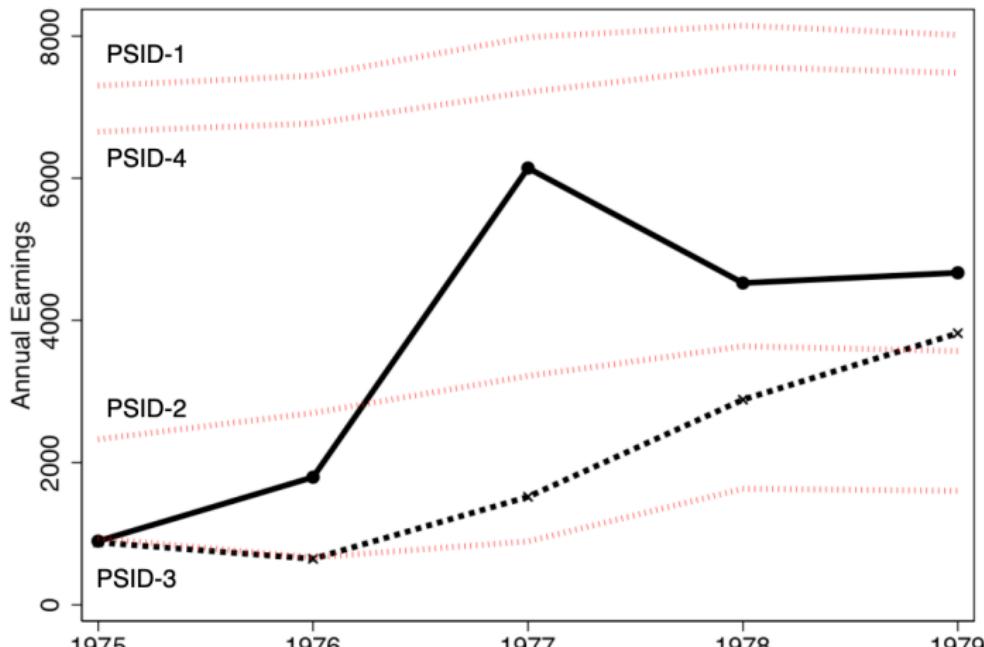
PSID

1. All female household heads aged 20-55
2. (1) + received AFDC in 1975
3. (2) + not working in 1976
4. (1) + no children below 5

CPS SSA

1. Matched sample
2. (1) + received AFDC in 1975
3. (1) + not working in 1976
4. (2) + not working in 1976

Annual Earnings - Females + PSID Comparisons

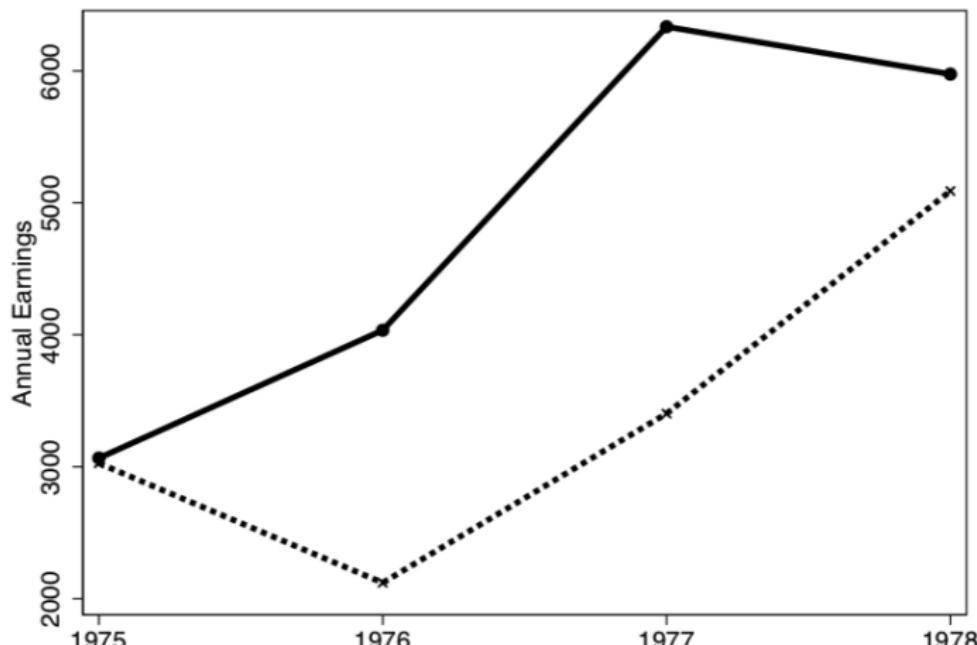


Source: LaLonde (1986), Table 2

Treatment effects Females

Name of Comparison Group ^d	Comparison Group Earnings Growth 1975–79	NSW Treatment Earnings Less Comparison Group Earnings				Difference in Differences: Difference in Earnings		Unrestricted Difference in Differences: Quasi Difference in Earnings Growth 1975–79		Controlling for All Observed Variables and Pre-Training Earnings	
		Pre-Training Year, 1975		Post-Training Year, 1979		Growth 1975–79 Treatments Less Comparisons		Unadjusted With Age		Without AFDC With AFDC	
		Unadjusted ^e	Adjusted ^e	Unadjusted ^e	Adjusted ^e	Without Age	With Age	Unadjusted ^e	Adjusted ^e	Without AFDC	With AFDC
Controls	2,942 (220)	-17 (122)	-22 (122)	851 (307)	861 (306)	833 (323)	883 (323)	843 (308)	864 (306)	854 (312)	-
<i>PSID-1</i>	713 (210)	-6,443 (326)	-4,882 (336)	-3,357 (403)	-2,143 (425)	3,097 (317)	2,657 (333)	1746 (357)	1,354 (380)	1664 (409)	2,097 (491)
<i>PSID-2</i>	1,242 (314)	-1,467 (216)	-1,515 (224)	1,090 (468)	870 (484)	2,568 (473)	2,392 (481)	1,764 (472)	1,535 (487)	1,826 (537)	-
<i>PSID-3</i>	665 (351)	-77 (202)	-100 (208)	3,057 (532)	2,915 (543)	3,145 (557)	3,020 (563)	3,070 (531)	2,930 (543)	2,919 (592)	-
<i>PSID-4</i>	928 (311)	-5,694 (306)	-4,976 (323)	-2,822 (460)	-2,268 (491)	2,883 (417)	2,655 (434)	1,184 (483)	950 (503)	1,406 (542)	2,146 (652)
<i>CPS-SSA-1</i>	233 (64)	-6,928 (272)	-5,813 (309)	-3,363 (320)	-2,650 (365)	3,578 (280)	3,501 (282)	1,214 (272)	1,127 (309)	536 (349)	1,041 (503)
<i>CPS-SSA-2</i>	1,595 (360)	-2,888 (204)	-2,332 (256)	-683 (428)	-240 (536)	2,215 (438)	2,068 (446)	447 (468)	620 (554)	665 (651)	-
<i>CPS-SSA-3</i>	1,207 (166)	-3,715 (226)	-3,150 (325)	-1,122 (311)	-812 (452)	2,603 (307)	2,615 (328)	814 (305)	784 (429)	-99 (481)	1,246 (720)
<i>CPS-SSA-4</i>	1,684 (524)	-1,189 (249)	-780 (283)	926 (630)	756 (716)	2,126 (654)	1,833 (663)	1,222 (637)	952 (717)	827 (814)	-

Annual Earnings - Males



Source: LaLonde (1986), Table 3

Choice of comparison group - Males

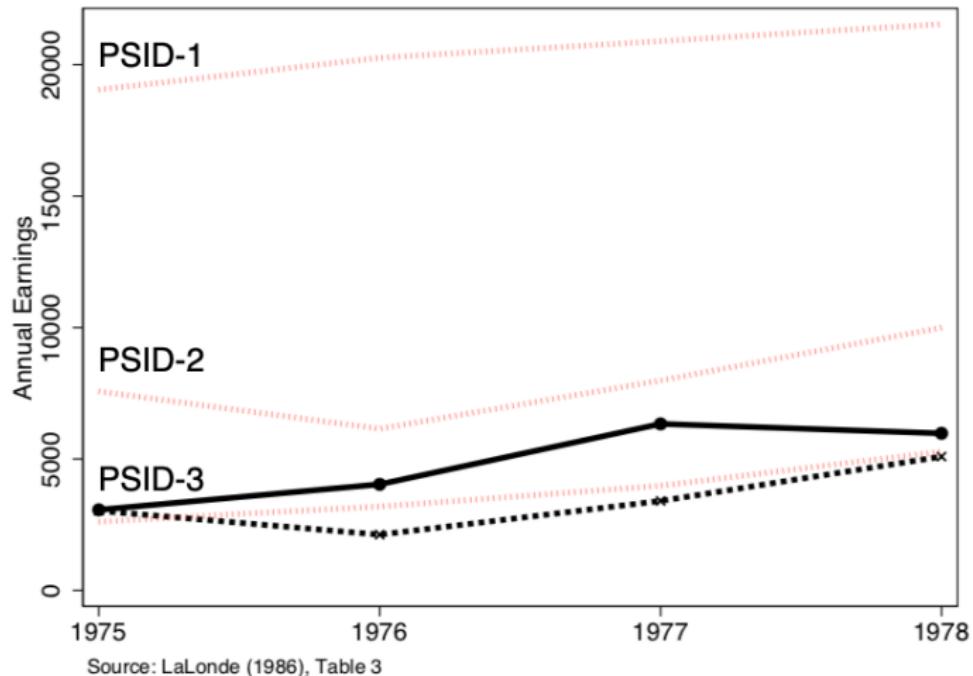
PSID

1. All male household heads not older than 55
2. (1) + not working in 1976
3. (2) + not working in 1975 or 1976

CPS SSA

1. Matched sample not older than 55
2. (1) + not working in 1976
3. (2) + not working in 1976 and income in 1975 below poverty level

Annual Earnings - Males + PSID Comparisons



Source: LaLonde (1986), Table 3

Treatment effects Males

Name of Comparison Group ^a	Comparison Group Earnings Growth 1975–78 (1)	NSW Treatment Earnings Less Comparison Group Earnings				Difference in Differences: Difference in Earnings		Unrestricted Difference in Differences: Quasi Difference in Earnings Growth 1975–78		Controlling for All Observed Variables and Pre-Training Earnings (10)	
		Pre-Training Year, 1975		Post-Training Year, 1978		Growth 1975–78 Treatments Less Comparisons					
		Unadjusted (2)	Adjusted ^c (3)	Unadjusted (4)	Adjusted ^c (5)	Without Age (6)	With Age (7)	Unadjusted (8)	Adjusted ^c (9)		
Controls	\$2,063 (325)	\$39 (383)	\$–21 (378)	\$886 (476)	\$798 (472)	\$847 (560)	\$856 (558)	\$897 (467)	\$802 (467)	\$662 (506)	
PSID-1	\$2,043 (237)	–\$15,997 (795)	–\$7,624 (851)	–\$15,578 (913)	–\$8,067 (990)	\$425 (650)	–\$749 (692)	–\$2,380 (680)	–\$2,119 (746)	–\$1,228 (896)	
PSID-2	\$6,071 (637)	–\$4,503 (608)	–\$3,669 (757)	–\$4,020 (781)	–\$3,482 (935)	\$484 (738)	–\$650 (850)	–\$1,364 (729)	–\$1,694 (878)	–\$792 (1024)	
PSID-3	(\$3,322 (780))	(\$455 (539))	(\$455 (704))	(\$697 (760))	(\$509 (967))	\$242 (884)	–\$1,325 (1078)	\$629 (757)	–\$552 (967)	\$397 (1103)	
CPS-SSA-1	\$1,196 (61)	–\$10,585 (539)	–\$4,654 (509)	–\$8,870 (562)	–\$4,416 (557)	\$1,714 (452)	\$195 (441)	–\$1,543 (426)	–\$1,102 (450)	–\$805 (484)	
CPS-SSA-2	\$2,684 (229)	–\$4,321 (450)	–\$1,824 (535)	–\$4,095 (537)	–\$1,675 (672)	\$226 (539)	–\$488 (530)	–\$1,850 (497)	–\$782 (621)	–\$319 (761)	
CPS-SSA-3	\$4,548 (409)	\$337 (343)	\$878 (447)	–\$1,300 (590)	\$224 (766)	–\$1,637 (631)	–\$1,388 (655)	–\$1,396 (582)	\$17 (761)	\$1,466 (984)	

Background

- LaLonde (1986) used a variety of methods available at the time
- The huge array of fancy selection on observables methods came later
- DW repeat the exercise using some of these methods
- They find results that suggest selection on observables works well
→ "Works well" here means matches the experimental benchmark
- See also Dehejia and Wahba (2002), which is a bit clearer

Impact

- This has also been a hugely influential paper
- Often cited in support of propensity score methods
- Findings have again been controversial (more later)

Data

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 (.32)	10.38 (.09)	.80 (.02)	.09 (.01)	.73 (.02)	.17 (.02)	3,066 (236)	
Control	425	24.45 (.32)	10.19 (.08)	.80 (.02)	.11 (.02)	.81 (.02)	.16 (.02)	3,026 (252)	
RE74 subset: ^b									
Treated	185	25.81 (.35)	10.35 (.10)	.84 (.02)	.059 (.01)	.71 (.02)	.19 (.02)	2,096 (237)	1,532 (156)
Control	260	25.05 (.34)	10.09 (.08)	.83 (.02)	.1 (.02)	.83 (.02)	.15 (.02)	2,107 (276)	1,267 (151)
Comparison groups: ^c									
PSID-1	2,490	34.85 [.78]	12.11 [.23]	.25 [.03]	.032 [.01]	.31 [.04]	.87 [.03]	19,429 [991]	19,063 [1,002]
PSID-2	253	36.10 [1.00]	10.77 [.27]	.39 [.04]	.067 [.02]	.49 [.05]	.74 [.04]	11,027 [853]	7,569 [695]
PSID-3	128	38.25 [1.17]	10.30 [.29]	.45 [.05]	.18 [.03]	.51 [.05]	.70 [.05]	5,566 [686]	2,611 [499]
CPS-1	15,992	33.22 [.81]	12.02 [.21]	.07 [.02]	.07 [.02]	.29 [.03]	.71 [.03]	14,016 [705]	13,650 [682]
CPS-2	2,369	28.25 [.87]	11.24 [.19]	.11 [.02]	.08 [.02]	.45 [.04]	.46 [.04]	8,728 [667]	7,397 [600]
CPS-3	429	28.03 [.87]	10.23 [.23]	.21 [.03]	.14 [.03]	.60 [.04]	.51 [.04]	5,619 [552]	2,467 [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; RE74 = earnings in calendar year 1974.

^a NSW sample as constructed by Lalonde (1986).

^b The subset of the Lalonde sample for which RE74 is available.

^c Definition of comparison groups (Lalonde 1986):

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.

PSID1-3 and CPS-1 are identical to those used by Lalonde. CPS2-3 are similar to those used by Lalonde, but Lalonde's original subset could not be recreated.

- NSW treatment and control are quite similar (balance)
- NSW and observational samples (PSID/CPS) are quite different

Propensity Score Overlap

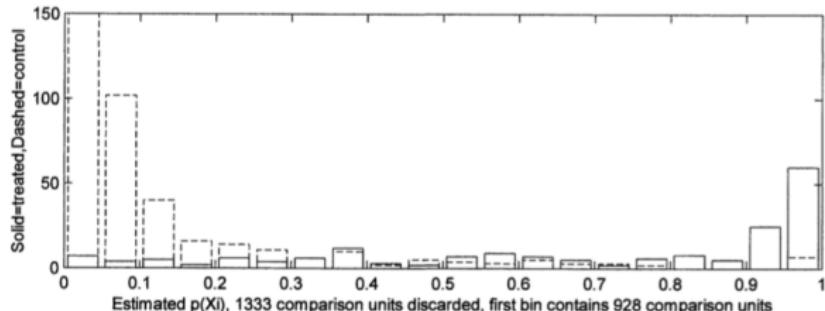


Figure 1. Histogram of the Estimated Propensity Score for NSW Treated Units and PSID Comparison Units. The 1,333 PSID units whose estimated propensity score is less than the minimum estimated propensity score for the treatment group are discarded. The first bin contains 928 PSID units. There is minimal overlap between the two groups. Three bins (.8-.85, .85-.9, and .9-.95) contain no comparison units. There are 97 treated units with an estimated propensity score greater than .8 and only 7 comparison units.

- Here this is a regression of treatment on X
- Assess both the overlap and fit of the propensity score
- Most comparison group units have very small propensity scores
- Treated units are more evenly spread out, but more close to 1

Results

Table 3. Estimated Training Effects for the NSW Male Participants Using Comparison Groups From PSID and CPS

	NSW earnings less comparison group earnings		NSW treatment earnings less comparison group earnings, conditional on the estimated propensity score					
			Quadratic in score ^b	Stratifying on the score			Matching on the score	
	(1) Unadjusted	(2) Adjusted ^a		(4) Unadjusted	(5) Adjusted	(6) Observations ^c	(7) Unadjusted	(8) Adjusted ^d
NSW	1,794 (633)	1,672 (638)						
PSID-1 ^e	-15,205 (1,154)	731 (886)	294 (1,389)	1,608 (1,571)	1,494 (1,581)	1,255	1,691 (2,209)	1,473 (809)
PSID-2 ^f	-3,647 (959)	683 (1,028)	496 (1,193)	2,220 (1,768)	2,235 (1,793)	389	1,455 (2,303)	1,480 (808)
PSID-3 ^g	1,069 (899)	825 (1,104)	647 (1,383)	2,321 (1,994)	1,870 (2,002)	247	2,120 (2,335)	1,549 (826)
CPS-1 ^h	-8,498 (712)	972 (550)	1,117 (747)	1,713 (1,115)	1,774 (1,152)	4,117	1,582 (1,069)	1,616 (751)
CPS-2 ⁱ	-3,822 (670)	790 (658)	505 (847)	1,543 (1,461)	1,622 (1,346)	1,493	1,788 (1,205)	1,563 (753)
CPS-3 ^j	-635 (657)	1,326 (798)	556 (951)	1,252 (1,617)	2,219 (2,082)	514	587 (1,496)	662 (776)

^a Least squares regression: RE78 on a constant, a treatment indicator, age, age², education, no degree, black, Hispanic, RE74, RE75.

^b Least squares regression of RE78 on a quadratic on the estimated propensity score and a treatment indicator, for observations used under stratification; see note (g).

^c Number of observations refers to the actual number of comparison and treatment units used for (3)-(5); namely, all treatment units and those comparison units whose estimated propensity score is greater than the minimum, and less than the maximum, estimated propensity score for the treatment group.

^d Weighted least squares: treatment observations weighted as 1, and control observations weighted by the number of times they are matched to a treatment observation [same covariates as (a)]. Propensity scores are estimated using the logistic model, with specifications as follows:

^e PSID-1: Prob (7 = 1) = F(age, age², education, education², married, no degree, black, Hispanic, RE74, RE75, RE74², RE75², u74*black).

^f PSID-2 and PSID-3: Prob (7 = 1) = F(age, age², education, education², no degree, married, black, Hispanic, RE74, RE74², RE75, RE75², u74, u75).

^g CPS-1, CPS-2, and CPS-3: Prob (7 = 1) = F(age, age², education, education², no degree, married, black, Hispanic, RE74, RE75, u74, u75, education*RE74, age²).

- Focus on the first two rows - experimental ATT is about 1800
- They prefer (4)&(5)(blocking), and (7)&(8) (matching)
- Simple regression in column (2) does less well

How much should we trust DW's findings?

Heckman et al. studies

- Dehejia and Wahba's (1999, 2002) finding of low bias from applying propensity score matching to LaLonde's (1986) data is arguably surprising
- Contrasts with results from Heckman et al. (1997a) and Heckman et al. (1996, 1998a)) using the experimental data from the JTPA Study
- They conclude that in order for matching estimators to have low bias, it is important that the data include
 - Variables affecting program part. and labor market outcomes
 - Nonexperimental comparison group be drawn from the same local labor markets as the participants
 - Comparable data for participants and nonparticipants (e.g earnings)

All these conditions fail to hold in the NSW data analyzed by LaLonde (1986) and Dehejia and Wahba (1999, 2002)

Critique and Lessons

Specific critique

- Smith and Todd (2005) show that DW results are not robust
- They argue changing sample from LaLonde (1986) is important
→ Sample selection choices as well as specification
- Also point out that income is measured differently in NSW/CPS
→ Self-reported vs. matched SS earnings data in March CPS
- Indeed, this is a critique of the whole LaLonde (1986) exercise
→ These points and others: Heckman, Ichimura and Todd (1997)

My reading of the literature

- The preferred nonexperimental evaluation strategy in a given context depends critically on the available data and on the institutions governing selection into the program
- Understanding the sources of bias may help a lot: i) lack of support, ii) quality of match, iii) selection on unobservables

Heckman et al: Understanding the biases

TABLE 2

*Decomposition of difference in post-programme mean earnings
 Bootstrapped standard errors shown in parentheses†
 Percentage of mean difference attributable to components in square brackets
 Earnings measured in average monthly dollars*

Experimental Controls and eligible nonparticipants (ENPs)†						Selection bias** (\hat{B}_{Sp}) as a % of treatment impact
	Mean difference \hat{B}	Non-overlap* \hat{B}_1	Density weighting \hat{B}_2	Selection bias \hat{B}_3	Average bias (\hat{B}_{Sp})	
Adult males (std. err.) [%]	-342 (47)	218 (38) [-64%]	-584 (41) [170%]	23 (33) [-7%]	38 (63)	87%
Adult females (std. err.) [%]	33 (26)	80 (13) [242%]	-78 (18) [-235%]	31 (26) [94%]	38 (33)	129%
Male youth (std. err.) [%]	20 (57)	142 (28) [704%]	-131 (35) [-650%]	9 (42) [46%]	14 (64)	23%
Female youth (std. err.) [%]	42 (36)	74 (17) [177%]	-67 (26) [-161%]	35 (28) [84%]	49 (42)	7239%

Heckman et al: Understanding the biases

Experimental Controls and SIPP Eligibles††

	Mean difference \hat{B}	Non-overlap* \hat{B}_1	Density weighting \hat{B}_2	Selection bias \hat{B}_3	Average bias (\hat{B}_{Sp})	Selection bias** (\hat{B}_{Sp}) as a % of treatment impact
Adult males (std. err.) [%]	-145 (56)	151 (30) [-104%]	-417 (44) [287%]	121 (33) [-83%]	192 (57)	440%
Adult females (std. err.) [%]	47 (23)	97 (19) [206%]	-172 (16) [-367%]	122 (15) [260%]	198 (26)	676%
Male youth (std. err.) [%]	-188 (106)	65 (108) [-35%]	-263 (53) (139%)	9 (25) [-5%]	21 (90)	36%
Female youth (std. err.) [%]	-88 (38)	83 (22) [-94%]	-168 (27) [191%]	-3 (10) [3%]	-13 (58)	1969%

† They are based on 50 replications of the data with 100% sampling.

† The best predictor Control-ENP participation models for all the demographic groups include indicator variables for site, age, race, education, marital status, children less than 6 and labour force transitions. In addition to these variables, the adult male model also includes an indicator for vocational training history, the number of household members, earnings in the month of random assignment or eligibility determination (RA or EL) and number of jobs held in 18 months before RA or EL. The adult female model includes an indicator for recent schooling, earnings in the month of RA or EL and number of labour force transitions in the 24 months prior to RA or EL. The male youth model includes average earnings in the 6 months and 12 months prior to RA or EL and average positive earnings in the 6 months before RA or EL. The female youth model includes earnings in the 12 months before RA or EL.

†† The best predictor Control-SIPP model includes indicators for age, race, education, marital status, children age less than 6, labour force transition patterns and levels of earnings in the preceding year. The data used are SIPP eligibles and experimental JTPA controls.

* A 2% trimming rule was used for adult males and females and a 5% trimming rule for youth was used in determining the overlapping support region (see Appendix C for a description of how the support is determined). The proportion of Controls and ENPs falling in the overlap region (S_p) are: 60% and 96% of adult males, 82% and 96% of adult females, 67% and 92% of male youth and 71% and 93% of female youth. The proportion of SIPPs and Controls falling in the overlap region are: 63% and 96% of adult males, 61% and 96% of adult females, 41% and 90% of male youth and 20% and 89% of female youth. A 0.06 fixed bandwidth and a biweight kernel, defined in Appendix A, were used for the nonparametric estimates.

** The final column displays the ratio of the absolute value of (\hat{B}_{S_p}) to the absolute value of experimental impact estimate.

Background and motivation

- Many studies document strong intergenerational associations in wealth
- Natural question: Why do wealthy parents have wealthy children?

Several possible explanations

- Selection: Parents genetically pass on abilities and preferences, creating intergen. links in income, savings behavior or financial risk taking
- Causation: Children's accumulation of wealth depends on the actions of their parents
- Several possible causal channels, including direct wealth transfers, parental investment, learning of attitudes and traits
- Distinguishing between selection and causality is key to understand how economic conditions or government policies may shape the persistence of wealth inequality across generations

Research question and design

What they do

- Fagereng et al. investigate the role of family background in determining children's wealth accumulation and investor behavior as adults
- Analysis is made possible by linking Korean-born children who were adopted at infancy by Norwegian parents
- Mechanism by which these Korean-Norwegian adoptees were assigned to adoptive families is known and effectively random

How they argue selection on observables

- First step: Submission of an application to review by case examiners
Parents were not given the opportunity to specify gender, family background or anything else about their future adoptee
- Next step: Approved files sent to the adoption agency, Holt Korea
Adoptees assigned to the approved adoptive families in the order the applications arrived
- Thus, as good as random assignment conditional on time of adoption

Testing for Quasi-Random Assignment

Table 2. Testing for quasi-random assignment of Korean-Norwegian adoptees

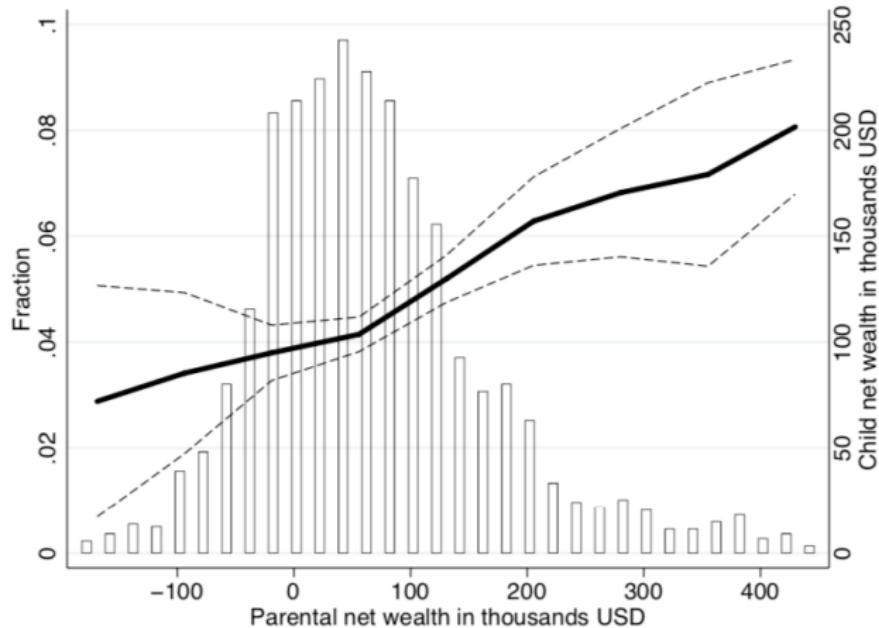
Regressors	Dependent variable:			
	Age at adoption		Gender	
	Specification:			
	Bivar. reg.	Multivar. reg.	Bivar. reg.	Multivar. reg.
Parent net wealth	-0.002 (0.003)	-0.002 (0.0037)	0.005 (0.004)	0.004 (0.004)
Mother's years of schooling	0.002 (0.002)	0.003 (0.003)	0.002 (0.003)	0.001 (0.004)
Father's years of schooling	0.001 (0.002)	-0.000 (0.002)	0.002 (0.003)	-0.000 (0.004)
(Log) parent income at birth	0.001 (0.035)	0.007 (0.038)	0.059 (0.0488)	0.037 (0.054)
Median (log) income in childhood municipality	-0.046 (0.034)	-0.047 (0.035)	0.051 (0.0459)	0.036 (0.047)
Dependent mean	0.78	0.78	0.75	0.75
F-stat, joint significance of regressors		0.737		0.648
[p-value]		[0.596]		[0.663]

Notes: The table contains estimates from regressions of a pre-determined characteristic of the adoptee (age at adoption or indicator for female) on family background variables such as parental net wealth, education (in years) of the mother and father, the log of parents income and the log the median income in parents' municipality of residence, all measured at the time of birth of the child. In columns 1 and 3, we run separate regressions for each of the family background variables (conditional on a full set of indicators for adoption years of the children). In columns 2 and 4, we run multivariate regressions with all the family characteristics (conditional on a full set of indicators for adoption years of the children). The estimation sample consists of 2,254 Korean-Norwegian adoptees.



Intergenerational Links in Wealth

Figure 4. Association between adoptee's net wealth and adoptive parents' net wealth



Notes: This figure is based on the baseline sample consisting of 2,254 Korean adoptees and their parents. The histogram shows the density of parental wealth (the left y-axis). The solid line shows estimates from a local linear regression of net wealth of the adoptee as an adult (measured as an average of 2012-2014) on the net wealth of her adoptive parents (measured as an average of 1994-1996), conditional on full set of indicators for year of adoption
► (UChicago)

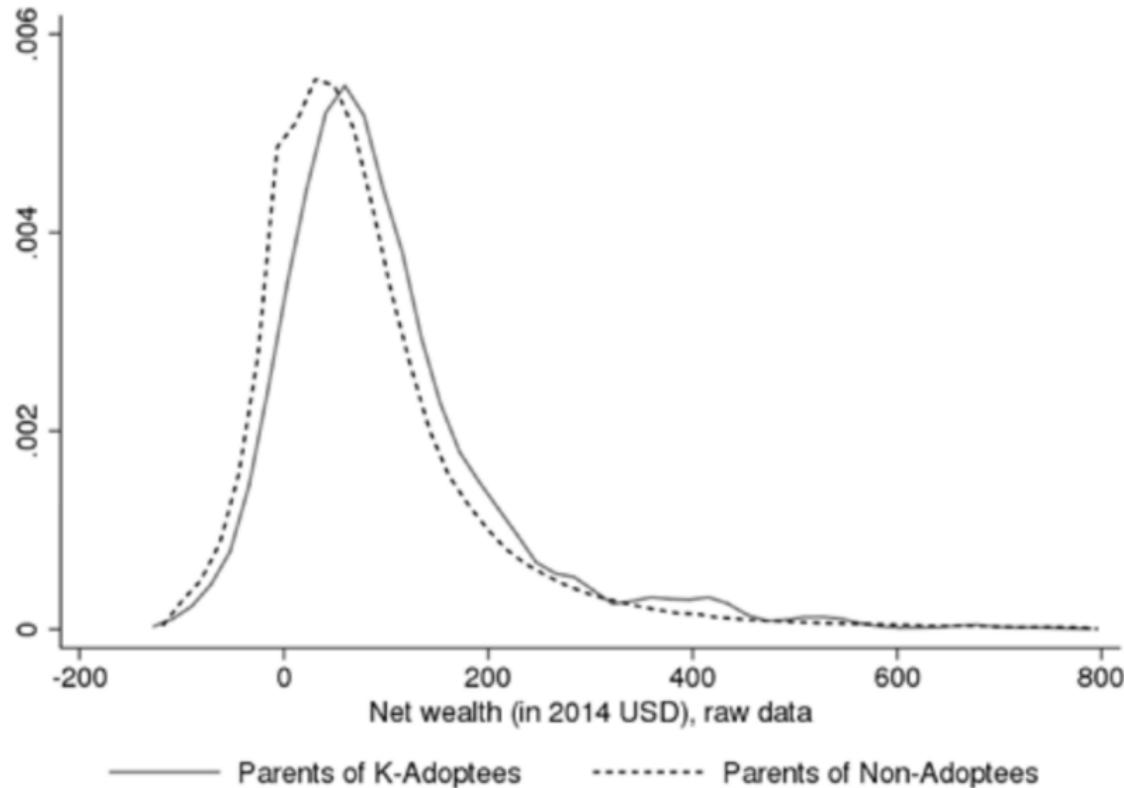
Intergenerational Links in Wealth

Table 3. Intergenerational links in wealth

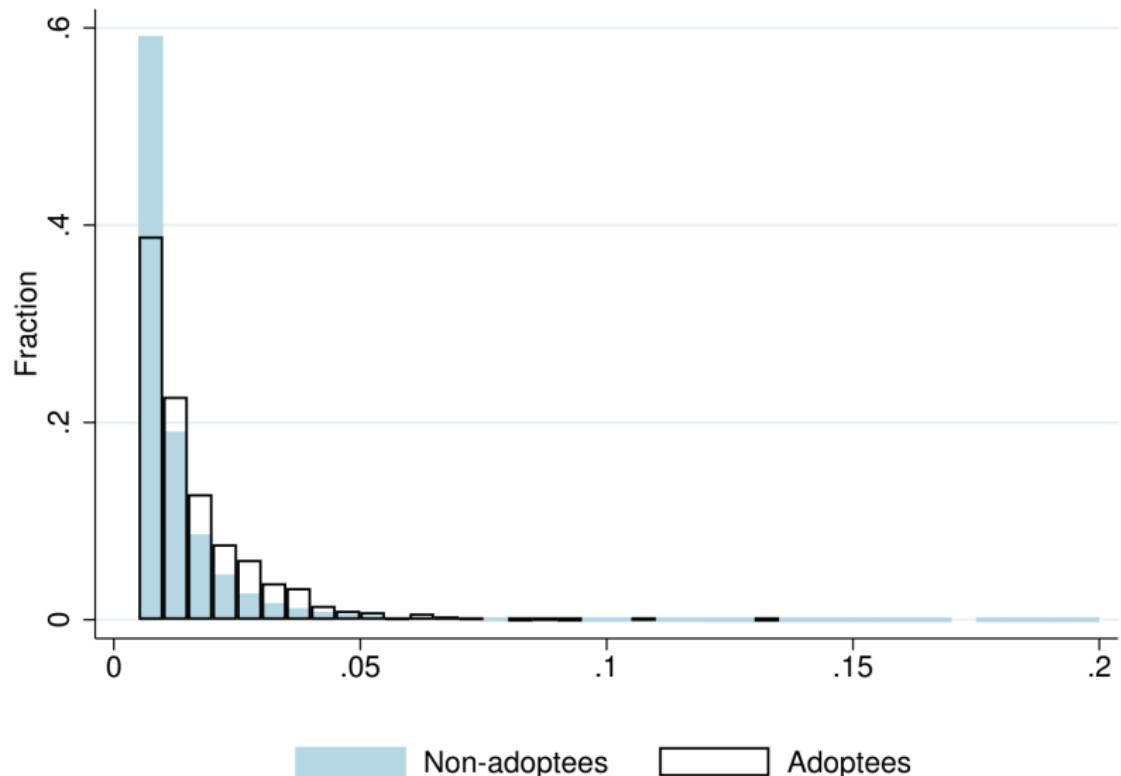
	Korean-Norwegian adoptees			Non-adoptees		
	(1)	(2)	(3)	(4)	(5)	(6)
Child-parent net wealth relation	0.225*** (0.041)	0.225*** (0.041)	0.204*** (0.042)	0.575*** (0.011)	0.547*** (0.011)	0.546*** (0.026)
Adoption year indicators	Yes	Yes	Yes			
Birth year ind. of child & parents	Yes	Yes	Yes	Yes	Yes	Yes
Gender		Yes	Yes	Yes	Yes	Yes
Adoption age (in days)		Yes	Yes			
Family characteristics			Yes		Yes	Yes
Matched sample (prop. score)						Yes
Observations	2,254			1,206,650	1,174,330	

Notes: The Korean-Norwegian adoptees are born in South Korea between 1965 and 1986, and adopted at infancy by Norwegian parents. Family characteristics include education (in years) of the mother and father, the number of siblings, the (log of) parents income and the (log of) the median income in parents' municipality of residence.

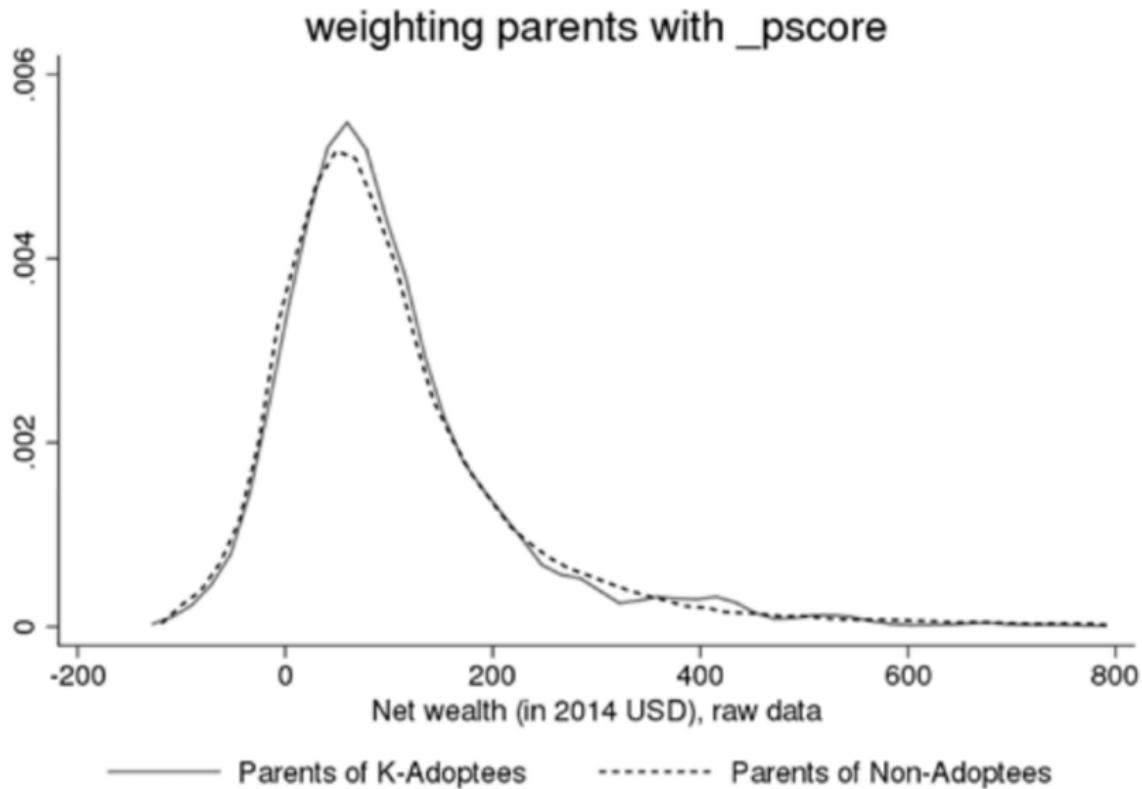
Distribution of Parental Wealth



Propensity Score



Weighted Distribution of Parental Wealth



Intergenerational Links in Wealth

Table 3. Intergenerational links in wealth

	Korean-Norwegian adoptees			Non-adoptees		
	(1)	(2)	(3)	(4)	(5)	(6)
Child-parent net wealth relation	0.225*** (0.041)	0.225*** (0.041)	0.204*** (0.042)	0.575*** (0.011)	0.547*** (0.011)	0.546*** (0.026)
Adoption year indicators	Yes	Yes	Yes			
Birth year ind. of child & parents	Yes	Yes	Yes	Yes	Yes	Yes
Gender		Yes	Yes	Yes	Yes	Yes
Adoption age (in days)		Yes	Yes			
Family characteristics			Yes		Yes	Yes
Matched sample (prop. score)						Yes
Observations	2,254			1,206,650	1,174,330	

Notes: The Korean-Norwegian adoptees are born in South Korea between 1965 and 1986, and adopted at infancy by Norwegian parents. Family characteristics include education (in years) of the mother and father, the number of siblings, the (log of) parents income and the (log of) the median income in parents' municipality of residence.

Criticisms of Assuming Selection on Observables

Inherent unobservables

- Selection on observables can often be difficult to argue
→ **Inherent unobservables:** preferences, private info, expectations, ...
- Boils down to arguing i) observationally identical people behave differently due to A, and ii) A is like a coin flip because of B

Controlling for more isn't an (convincing) answer

- Often argued that large X makes selection on observables more likely
→ Strictly speaking, this is not necessarily true-we saw this earlier
- Even if it were, still raises an uncomfortable friction with overlap
→ If we could perfectly explain D with X then $\mathbb{P}[D = 1|X] = 0$ or 1

Fancy methods for choosing observables will not solve this

- Selection on observables is seeing a resurgence with "machine learning"
- Fancier methods, but the identifying assumption is still the same
- Estimation and bias/variance trade-off is rarely the first-order issue here

Angrist (1998)

A bad use of selection on observables

- Well-known economic application of selection on observables:
 - Y is a labor market outcome (employment/earnings)
 - $D \in \{0, 1\}$ is veteran status (participation in the military)
 - X are socioeconomic variables (race, year, schooling, AFQT, age)
- Assumption is given X , military participation as-if randomly assigned
- **Observationally similar people randomly join the military?!**
- Ignores first-order issues such as outside employment options
 - Also unobservable screening factors (fitness, interpersonal skills)
- These are unobserved and inherently unobservable

Allowing for selection on unobservables

- Most applied microeconomists seem to share this skepticism
- This motivates the other methods that we will discuss in the course
- All use different arguments to allow for **selection on unobservables**

Parametric identification: A cautionary tale

How to deal with selection on unobservables

- One possibility is to find an instrument (next week!)
- Another possibility is to make functional form / parametric assumptions about the relationship between the observed data and the unobservables
- In practice, it is often necessary to do both. Not necessarily anything wrong with that.
- But be wary; with strong enough assumptions you can draw inferences about anything (but learn nothing)

Heckman (2008); first draft 1980

- The Effect of Prayer on God's Attitude Toward Mankind
- Data on intensity of prayer from the National Opinion Research Center's (NORC) survey on religious attitudes
- God's attitude is unobserved to the analyst but why would that stop him?

Data and assumption

- $Y \in [0, 1]$ - God's attitude, this is an unobserved variable.
- $X \in [0, 1]$ - the intensity of prayer in the population.
- The population density of prayer is summarized by a univariate density $f(X)$ which has been estimated by Father Greeley (1972).
- Accept on faith that the conditional density of X given Y is of the form

$$g(X|Y) = a(Y) \exp(XY) \quad (1)$$

where $a(Y)$ is an unknown, continuous, positive, and differentiable function.

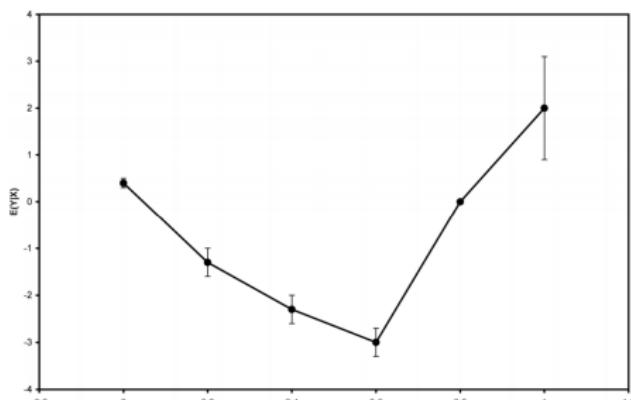
Estimating equation and results

Heckman (2008)

- Singh (1977) demonstrates that:

$$\mathbb{E}[Y|X = x] = \frac{f'(x)}{f(x)}$$

- From the population distribution of prayer we can estimate the population regression function of God's attitude as a function of prayer.



What's the lesson?

- If one is willing to make a parametric assumption, it is not necessary to observe a variable in order to compute its conditional expectation with respect to another variable
- Admittedly an extreme example, but it makes an important point:
 - Think carefully about (the implications of) the assumptions you make
 - Some important questions might not be possible to rigorously study or credibly answer
 - How would you define, identify and estimate parameters answering questions such as: What's the effect of inequality on democracy? Why are some countries rich? What's the meaning of life?