

# 150C Causal Inference

## Causal Inference Under Selection on Observables

Jonathan Mummolo

# Today's Random Medical News

from the New England  
Journal of  
Panic-Inducing  
Gobbledygook

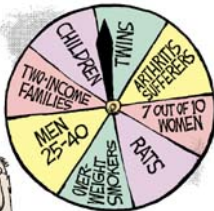
V. M. B. P. M. A. N. IN 1981 ENGLISH 101



CAN CAUSE



IN



Lancet 2001: negative correlation between coronary heart disease mortality and level of vitamin C in bloodstream (controlling for age, gender, blood pressure, diabetes, and smoking)

# Today's Random Medical News

from the New England  
Journal of  
Panic-Inducing  
Gobbledygook

V. M. B. G. M. A. N. IN 1984 12/28/02/02/02/02/02/02



Lancet 2002: no effect of vitamin C on mortality in controlled placebo trial (controlling for nothing)

# Today's Random Medical News

from the New England  
Journal of  
Panic-Inducing  
Gobbledygook

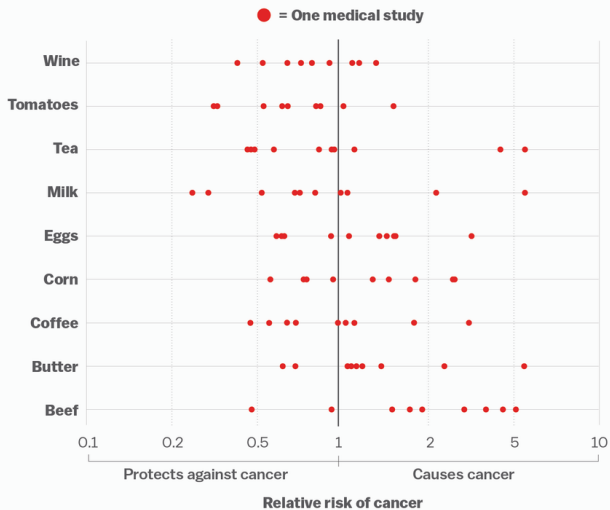
VINCE BRIAN



Lancet 2003: comparing among individuals with the same age, gender, blood pressure, diabetes, and smoking, those with higher vitamin C levels have lower levels of obesity, lower levels of alcohol consumption, are less likely to grow up in working class, etc.

# Observational Medical Trials

## Everything we eat both causes and prevents cancer



SOURCE: Schoenfeld and Ioannidis, *American Journal of Clinical Nutrition*

# Observational Studies

- Randomization forms gold standard for causal inference, because it balances **observed** and **unobserved** confounders
- Cannot always randomize so we do observational studies, where we **adjust** for the **observed covariates** and **hope** that unobservables are balanced
- Better than hoping: **design** observational study to approximate an experiment
  - “The planner of an observational study should always ask himself: How would the study be conducted if it were possible to do it by controlled experimentation” (Cochran 1965)

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

## Treatments, Covariates, Outcomes

- **Randomized Experiment:** Well-defined treatment, clear distinction between covariates and outcomes, control of assignment mechanism
- **Better Observational Study:** Well-defined treatment, clear distinction between covariates and outcomes, precise knowledge of assignment mechanism
  - Can convincingly answer the following question: Why do two units who are identical on measured covariates receive different treatments?
- **Poorer Observational Study:** Hard to say when treatment began or what the treatment really is. Distinction between covariates and outcomes is blurred, so problems that arise in experiments seem to be avoided but are in fact just ignored. No precise knowledge of assignment mechanism.



## How were treatments assigned?

- **Randomized Experiment:** Random assignment
- **Better Observational Study:** Assignment is not random, but circumstances for the study were chosen so that treatment seems haphazard, or at least not obviously related to potential outcomes (sometimes we refer to these as natural or quasi-experiments)
- **Poorer Observational Study:** No attention given to assignment process, units self-select into treatment based on potential outcomes

## **What is the problem with purely cross-sectional data?**

- Difficult to know what is pre or post treatment.
- Many important confounders will be affected by the treatment and including these “bad controls” induces post-treatment bias.
- But if you do not condition on the confounders that are post-treatment, then often only left with a limited set of covariates such as socio-demographics.

## Were treated and controls comparable?

- **Randomized Experiment:** Balance table for observables.
- **Better Observational Study:** Balance table for observables. Ideally sensitivity analysis for unobservables.
- **Poorer Observational Study:** No direct assessment of comparability is presented.

## Eliminating plausible alternatives to treatment effects?

- **Randomized Experiment:** List plausible alternatives and experimental design includes features that shed light on these alternatives (e.g. placebos). Report on potential attrition and non-compliance.
- **Better Observational Study:** List plausible alternatives and study design includes features that shed light on these alternatives (e.g. multiple control groups, longitudinal covariate data, etc.). Requires more work than in experiment since there are usually many more alternatives.
- **Poorer Observational Study:** Alternatives are mentioned in discussion section of the paper.

# Good Observational Studies

**Design features** we can use to handle unobservables:

- Design comparisons so that unobservables are likely to be balanced (e.g. sub-samples, groups where treatment assignment was accidental, etc.)
- Unobservables may differ, but comparisons that are unaffected by differences in time-invariant unobservables
- Instrumental variables, if applied correctly
- Multiple control groups that are known to differ on unobservables
- Sensitivity analysis and bounds

# Seat Belts on Fatality Rates

**Table 1.1** Crashes in FARS 1975–1983 in which the front seat had two occupants, a driver and a passenger, with one belted, the other unbelted, and one died and one survived.

	Driver Passenger	Not Belted Belted	Belted Not Belted
Driver Died	Passenger Survived	189	153
Driver Survived	Passenger Died	111	363

Evans (1986)

## Who Gets a Swiss Passport? A Natural Experiment in Immigrant Discrimination

JENS HAINMUELLER *Massachusetts Institute of Technology*

DOMINIK HANGARTNER *London School of Economics & University of Zurich*

**W**e study discrimination against immigrants using microlevel data from Switzerland, where, until recently, some municipalities used referendums to decide on the citizenship applications of foreign residents. We show that naturalization decisions vary dramatically with immigrants' attributes, which we collect from official applicant descriptions that voters received before each referendum. Country of origin determines naturalization success more than any other applicant characteristic, including language skills, integration status, and economic credentials. The average proportion of "no" votes is about 40% higher for applicants from (the former) Yugoslavia and Turkey compared to observably similar applicants from richer northern and western European countries. Statistical and taste-based discrimination contribute to varying naturalization success; the rewards for economic credentials are higher for applicants from disadvantaged origins, and origin-based discrimination is much stronger in more xenophobic municipalities. Moreover, discrimination against specific immigrant groups responds dynamically to changes in the groups' relative size.

# A known treatment assignment process

C//////, G'//////, italienische Staatsangehörige,  
Gerliswilstrasse 26, 6020 Emmenbrücke



**Geburtsort:** Pietrelcina (I)  
**Geburtsdatum:** 9. Dezember 1939  
**Zivilstand:** geschieden  
**Ausbildung:** Volksschule  
**Bisherige Tätigkeiten:** Mitarbeit auf elterlichem Bauerngut,  
Lingerie-Mitarbeiterin in Hotels  
**Jetzige Tätigkeit:** IV-Rentnerin seit 1997  
**Arbeitgeber:** –  
**Einreise in die Schweiz:** 15. Oktober 1962  
**Zuzug nach Emmen:** 23. September 1970  
**Hobbys:** –  
**Steuern:** Steuerbares Einkommen Fr. 33 900.–  
Steuerbares Vermögen Fr. 28 000.–  
**Kinder:** –  
**Einbürgerungstaxe:** Fr. 123.–  
**Einbürgerungsgebühr:** Fr. 500.–

D'//////, J'//////, ungarischer Staatsangehöriger, Ghürschweg 13,  
6020 Emmenbrücke

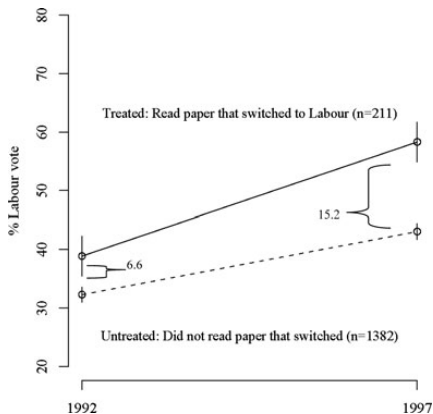


**Geburtsort:** Bucsa (H)  
**Geburtsdatum:** 14. Mai 1936  
**Zivilstand:** geschieden  
**Ausbildung:** Volksschule, Lehre als Mineur und Sprengmeister,  
Zusatzausbildung als Maler  
**Bisherige Tätigkeiten:** Bau-Hilfsarbeiter, selbstständiger Maler  
**Jetzige Tätigkeit:** IV-Rentner seit 1987 (Verkehrsunfall)  
**Arbeitgeber:** –  
**Einreise in die Schweiz:** 17. November 1956  
**Zuzug nach Emmen:** 26. Juni 1991  
**Hobbys:** Fischen, Pilze sammeln, Modellflugzeuge basteln  
**Steuern:** Steuerbares Einkommen Fr. 28 400.–  
Steuerbares Vermögen Fr. 0.–  
**Kinder:** –  
**Einbürgerungstaxe:** Fr. 100.–  
**Einbürgerungsgebühr:** Fr. 500.–



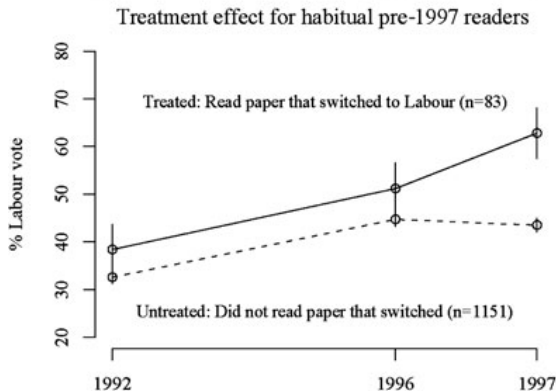
“Official voting leaflets summarizing the applicant characteristics were sent to all citizens usually about two to six weeks before each naturalization referendum. Since we retrieved the voting leaflets from the municipal archives, **we measure the same applicant information from the leaflets that the citizens observed when they voted on the citizenship applications.** Since most voters simply draw on the leaflets to decide on the applicants, this design enables us to greatly minimize potential omitted variable bias and attribute differences in naturalization outcomes to the effects of differences in measured applicant characteristics.”  
-Hainmueller and Hangartner (2013)

# Persuasive Effect of Endorsement Changes on Labour Vote



This figure shows that reading a paper that switched to Labour is associated with an  $(15.2 - 6.6 =) 8.6$  percentage point shift to Labour between the 1992 and 1997 UK elections. Paper readership is measured in the 1996 wave, before the papers switched, or, if no 1996 interview was conducted, in an earlier wave. Confidence intervals show one standard error.

# Persuasive Effect of Endorsement Changes on Labour Vote



Using the hypothetical vote choice question asked in the 1996 wave, this figure shows that the treatment effect only emerges after 1996. Habitual readers are those who read a paper that switched in every wave in which they were interviewed before the 1997 wave.

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Adjustment for Observables in Observational Studies

- Subclassification
- Matching
- Propensity Score Methods
- Regression

# Smoking and Mortality (Cochran (1968))

TABLE 1  
DEATH RATES PER 1,000 PERSON-YEARS

Smoking group	Canada	U.K.	U.S.
Non-smokers	20.2	11.3	13.5
Cigarettes	20.5	14.1	13.5
Cigars/pipes	35.5	20.7	17.4

# Smoking and Mortality (Cochran (1968))

TABLE 2  
MEAN AGES, YEARS

Smoking group	Canada	U.K.	U.S.
Non-smokers	54.9	49.1	57.0
Cigarettes	50.5	49.8	53.2
Cigars/pipes	65.9	55.7	59.7

To control for differences in age, we would like to compare different smoking-habit groups with the same age distribution

One possibility is to use subclassification:

- for each country, divide each group into different age subgroups
- calculate death rates within age subgroups
- average within age subgroup death rates using fixed weights (e.g. number of cigarette smokers)



## Subclassification: Example

	Death Rates Pipe Smokers	# Pipe- Smokers	# Non- Smokers
Age 20 - 50	15	11	29
Age 50 - 70	35	13	9
Age + 70	50	16	2
Total		40	40

What is the average death rate for Pipe Smokers?

## Subclassification: Example

	Death Rates Pipe Smokers	# Pipe- Smokers	# Non- Smokers
Age 20 - 50	15	11	29
Age 50 - 70	35	13	9
Age + 70	50	16	2
Total		40	40

What is the average death rate for Pipe Smokers?

$$15 \cdot (11/40) + 35 \cdot (13/40) + 50 \cdot (16/40) = 35.5$$

## Subclassification: Example

	Death Rates Pipe Smokers	# Pipe- Smokers	# Non- Smokers
Age 20 - 50	15	11	29
Age 50 - 70	35	13	9
Age + 70	50	16	2
Total		40	40

What is the average death rate for Pipe Smokers if they had same age distribution as Non-Smokers?

## Subclassification: Example

	Death Rates Pipe Smokers	# Pipe- Smokers	# Non- Smokers
Age 20 - 50	15	11	29
Age 50 - 70	35	13	9
Age + 70	50	16	2
Total		40	40

What is the average death rate for Pipe Smokers if they had same age distribution as Non-Smokers?

$$15 \cdot (29/40) + 35 \cdot (9/40) + 50 \cdot (2/40) = 21.2$$

# Smoking and Mortality (Cochran (1968))

TABLE 3  
ADJUSTED DEATH RATES USING 3 AGE GROUPS

Smoking group	Canada	U.K.	U.S.
Non-smokers	20.2	11.3	13.5
Cigarettes	28.3	12.8	17.7
Cigars/pipes	21.2	12.0	14.2

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Identification Under Selection on Observables

## Identification Assumption

- 1  $(Y_1, Y_0) \perp\!\!\!\perp D|X$  (selection on observables)
- 2  $0 < \Pr(D = 1|X) < 1$  with probability one (common support)

## Identification Result

Given selection on observables we have

$$\begin{aligned}\mathbf{E}[Y_1 - Y_0|X] &= \mathbf{E}[Y_1 - Y_0|X, D = 1] \\ &= \mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]\end{aligned}$$

Therefore, under the common support condition:

$$\begin{aligned}\tau_{ATE} &= \mathbf{E}[Y_1 - Y_0] = \int \mathbf{E}[Y_1 - Y_0|X] dP(X) \\ &= \int (\mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]) dP(X)\end{aligned}$$

# Identification Under Selection on Observables

## Identification Assumption

- 1  $(Y_1, Y_0) \perp\!\!\!\perp D | X$  (*selection on observables*)
- 2  $0 < \Pr(D = 1 | X) < 1$  with probability one (*common support*)

## Identification Result

*Similarly,*

$$\begin{aligned}\tau_{ATT} &= \mathbf{E}[Y_1 - Y_0 | D = 1] \\ &= \int (\mathbf{E}[Y | X, D = 1] - \mathbf{E}[Y | X, D = 0]) dP(X | D = 1)\end{aligned}$$

*To identify  $\tau_{ATT}$  the selection on observables and common support conditions can be relaxed to:*

- $Y_0 \perp\!\!\!\perp D | X$  (*SOO for Controls*)
- $\Pr(D = 1 | X) < 1$  (*Weak Overlap*)



# Identification Under Selection on Observables

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	$E[Y_1 X = 0, D = 1]$	$E[Y_0 X = 0, D = 1]$	1	0
2			1	0
3	$E[Y_1 X = 0, D = 0]$	$E[Y_0 X = 0, D = 0]$	0	0
4			0	0
5	$E[Y_1 X = 1, D = 1]$	$E[Y_0 X = 1, D = 1]$	1	1
6			1	1
7	$E[Y_1 X = 1, D = 0]$	$E[Y_0 X = 1, D = 0]$	0	1
8			0	1

# Identification Under Selection on Observables

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	$E[Y_1 X = 0, D = 1]$	$E[Y_0 X = 0, D = 1] =$	1	0
2		$E[Y_0 X = 0, D = 0]$	1	0
3	$E[Y_1 X = 0, D = 0]$	$E[Y_0 X = 0, D = 0]$	0	0
4			0	0
5	$E[Y_1 X = 1, D = 1]$	$E[Y_0 X = 1, D = 1] =$	1	1
6		$E[Y_0 X = 1, D = 0]$	1	1
7	$E[Y_1 X = 1, D = 0]$	$E[Y_0 X = 1, D = 0]$	0	1
8			0	1

$(Y_1, Y_0) \perp\!\!\!\perp D|X$  implies that we conditioned on all confounders. The treatment is randomly assigned within each stratum of  $X$ :

$$\begin{aligned}
 E[Y_0|X = 0, D = 1] &= E[Y_0|X = 0, D = 0] \text{ and} \\
 E[Y_0|X = 1, D = 1] &= E[Y_0|X = 1, D = 0]
 \end{aligned}$$

# Identification Under Selection on Observables

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	$E[Y_1 X = 0, D = 1]$	$E[Y_0 X = 0, D = 1] =$	1	0
2		$E[Y_0 X = 0, D = 0]$	1	0
3	$E[Y_1 X = 0, D = 0] =$	$E[Y_0 X = 0, D = 0]$	0	0
4	$E[Y_1 X = 0, D = 1]$		0	0
5	$E[Y_1 X = 1, D = 1]$	$E[Y_0 X = 1, D = 1] =$	1	1
6		$E[Y_0 X = 1, D = 0]$	1	1
7	$E[Y_1 X = 1, D = 0] =$	$E[Y_0 X = 1, D = 0]$	0	1
8	$E[Y_1 X = 1, D = 1]$		0	1

$(Y_1, Y_0) \perp\!\!\!\perp D | X$  also implies

$$E[Y_1|X = 0, D = 1] = E[Y_1|X = 0, D = 0] \text{ and}$$

$$E[Y_1|X = 1, D = 1] = E[Y_1|X = 1, D = 0]$$

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables**
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Subclassification Estimator

## Identification Result

$$\tau_{ATE} = \int (\mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]) dP(X)$$

$$\tau_{ATT} = \int (\mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]) dP(X|D = 1)$$

Assume  $X$  takes on  $K$  different cells  $\{X^1, \dots, X^k, \dots, X^K\}$ . Then the analogy principle suggests estimators:

# Subclassification Estimator

## Identification Result

$$\tau_{ATE} = \int (\mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]) dP(X)$$

$$\tau_{ATT} = \int (\mathbf{E}[Y|X, D = 1] - \mathbf{E}[Y|X, D = 0]) dP(X|D = 1)$$

Assume  $X$  takes on  $K$  different cells  $\{X^1, \dots, X^k, \dots, X^K\}$ . Then the analogy principle suggests estimators:

$$\hat{\tau}_{ATE} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N^k}{N}\right); \quad \hat{\tau}_{ATT} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N_1^k}{N_1}\right)$$

- $N^k$  is # of obs. and  $N_1^k$  is # of treated obs. in cell  $k$
- $\bar{Y}_1^k$  is mean outcome for the treated in cell  $k$
- $\bar{Y}_0^k$  is mean outcome for the untreated in cell  $k$

## Subclassification by Age ( $K = 2$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old	28	24	4	3	10
Young	22	16	6	7	10
Total				10	20

What is  $\hat{\tau}_{ATE} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N^k}{N}\right)$ ?

## Subclassification by Age ( $K = 2$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old	28	24	4	3	10
Young	22	16	6	7	10
Total				10	20

What is  $\hat{\tau}_{ATE} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N^k}{N}\right)$ ?

$$\hat{\tau}_{ATE} = 4 \cdot (10/20) + 6 \cdot (10/20) = 5$$



## Subclassification by Age ( $K = 2$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old	28	24	4	3	10
Young	22	16	6	7	10
Total				10	20

What is  $\hat{\tau}_{ATT} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N_1^k}{N_1}\right)$ ?

## Subclassification by Age ( $K = 2$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old	28	24	4	3	10
Young	22	16	6	7	10
Total				10	20

What is  $\hat{\tau}_{ATT} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N_1^k}{N_1}\right)$ ?

$$\hat{\tau}_{ATT} = 4 \cdot (3/10) + 6 \cdot (7/10) = 5.4$$

## Subclassification by Age and Gender ( $K = 4$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old, Male	28	22	4	3	7
Old, Female		24		0	3
Young, Male	21	16	5	3	4
Young, Female	23	17	6	4	6
Total				10	20

What is  $\hat{\tau}_{ATE} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N^k}{N}\right)$ ?

## Subclassification by Age and Gender ( $K = 4$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old, Male	28	22	4	3	7
Old, Female		24		0	3
Young, Male	21	16	5	3	4
Young, Female	23	17	6	4	6
Total				10	20

What is  $\hat{\tau}_{ATE} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N^k}{N}\right)$ ?

Not identified!

## Subclassification by Age and Gender ( $K = 4$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old, Male	28	22	4	3	7
Old, Female		24		0	3
Young, Male	21	16	5	3	4
Young, Female	23	17	6	4	6
Total				10	20

What is  $\hat{\tau}_{ATT} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N_1^k}{N_1}\right)$ ?

## Subclassification by Age and Gender ( $K = 4$ )

$X_k$	Death Rate Smokers	Death Rate Non-Smokers	Diff.	# Smokers	# Obs.
Old, Male	28	22	4	3	7
Old, Female		24		0	3
Young, Male	21	16	5	3	4
Young, Female	23	17	6	4	6
Total				10	20

What is  $\hat{\tau}_{ATT} = \sum_{k=1}^K (\bar{Y}_1^k - \bar{Y}_0^k) \cdot \left(\frac{N_1^k}{N_1}\right)$ ?

$$\hat{\tau}_{ATT} = 4 \cdot (3/10) + 5 \cdot (3/10) + 6 \cdot (4/10) = 5.1$$

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables**
  - Subclassification
  - Matching**
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Matching

When  $X$  is continuous we can estimate  $\tau_{ATT}$  by “imputing” the missing potential outcome of each treated unit using the observed outcome from the “closest” control unit:

$$\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$$

where  $Y_{j(i)}$  is the outcome of an untreated observation such that  $X_{j(i)}$  is the **closest** value to  $X_i$  among the untreated observations.



# Matching

When  $X$  is continuous we can estimate  $\tau_{ATT}$  by “imputing” the missing potential outcome of each treated unit using the observed outcome from the “closest” control unit:

$$\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$$

where  $Y_{j(i)}$  is the outcome of an untreated observation such that  $X_{j(i)}$  is the **closest** value to  $X_i$  among the untreated observations.

We can also use the average for  $M$  closest matches:

$$\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} \left\{ Y_i - \left( \frac{1}{M} \sum_{m=1}^M Y_{j_m(i)} \right) \right\}$$

Works well when we can find good matches for each treated unit

## Matching: Example with a Single $X$

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	?	1	3
2	1	?	1	1
3	0	?	1	10
4		0	0	2
5		9	0	3
6		1	0	-2
7		1	0	-4

What is  $\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$ ?

## Matching: Example with a Single $X$

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	?	1	3
2	1	?	1	1
3	0	?	1	10
4		0	0	2
5		9	0	3
6		1	0	-2
7		1	0	-4

What is  $\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$ ?

Match and plugin in

## Matching: Example with a Single $X$

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	9	1	3
2	1	0	1	1
3	0	9	1	10
4		0	0	2
5		9	0	3
6		1	0	-2
7		1	0	-4

What is  $\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$ ?

## Matching: Example with a Single $X$

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	9	1	3
2	1	0	1	1
3	0	9	1	10
4		0	0	2
5		9	0	3
6		1	0	-2
7		1	0	-4

What is  $\hat{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)})$ ?

$$\hat{\tau}_{ATT} = 1/3 \cdot (6 - 9) + 1/3 \cdot (1 - 0) + 1/3 \cdot (0 - 9) = -3.7$$

## Matching Distance Metric

“Closeness” is often defined by a **distance metric**. Let  $X_i = (X_{i1}, X_{i2}, \dots, X_{ik})'$  and  $X_j = (X_{j1}, X_{j2}, \dots, X_{jk})'$  be the covariate vectors for  $i$  and  $j$ .

A commonly used distance is the **Mahalanobis distance**:

$$MD(X_i, X_j) = \sqrt{(X_i - X_j)' \Sigma^{-1} (X_i - X_j)}$$

where  $\Sigma$  is the Variance-Covariance-Matrix so the distance metric is scale-invariant and takes into account the correlations. For an exact match  $MD(X_i, X_j) = 0$ .

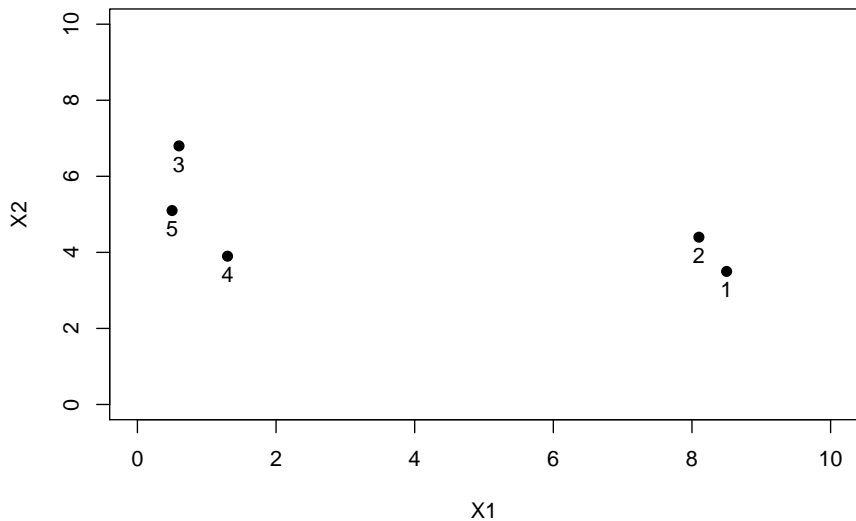
Other distance metrics can be used, for example we might use the normalized Euclidean distance, etc.

# Euclidean Distance Metric

R Code

```
> X
      X1 X2
[1,] 8.5 3.5
[2,] 8.1 4.4
[3,] 0.6 6.8
[4,] 1.3 3.9
[5,] 0.5 5.1
>
> Xdist <- dist(X,diag = T,upper=T)
> round(Xdist,1)
      1  2  3  4  5
1 0.0 1.0 8.6 7.2 8.2
2 1.0 0.0 7.9 6.8 7.6
3 8.6 7.9 0.0 3.0 1.7
4 7.2 6.8 3.0 0.0 1.4
5 8.2 7.6 1.7 1.4 0.0
```

# Euclidean Distance Metric





# Useful Matching Functions

The workhorse model is the `Match()` function in the `Matching` package:

```
Match(Y = NULL, Tr, X, Z = X, V = rep(1, length(Y)),  
      estimand = "ATT", M = 1, BiasAdjust = FALSE, exact = NULL,  
      caliper = NULL, replace = TRUE, ties = TRUE,  
      CommonSupport = FALSE, Weight = 1, Weight.matrix = NULL,  
      weights = NULL, Var.calc = 0, sample = FALSE, restrict = NULL,  
      match.out = NULL, distance.tolerance = 1e-05,  
      tolerance = sqrt(.Machine$double.eps), version = "standard")
```

Default distance metric (`Weight=1`) is normalized Euclidean distance

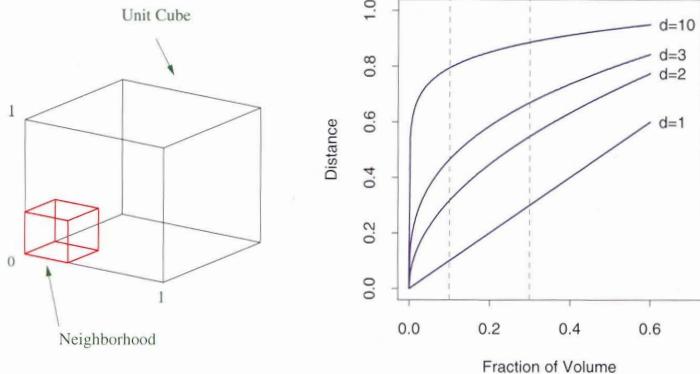
- `MatchBalance(formu)` for balance checking

# Local Methods and the Curse of Dimensionality

**Big** Problem:

# Local Methods and the Curse of Dimensionality

**Big Problem:** in a mathematical space, the volume increases **exponentially** when adding extra dimensions.



**FIGURE 2.6.** *The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction  $r$  of the volume of the data, for different dimensions  $p$ . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.*

## Matching with Bias Correction

Matching estimators may behave badly if  $X$  contains multiple continuous variables.

Need to adjust matching estimators in the following way:

$$\tilde{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} (Y_i - Y_{j(i)}) - (\hat{\mu}_0(X_i) - \hat{\mu}_0(X_{j(i)})),$$

where  $\mu_0(x) = E[Y|X = x, D = 0]$  is the population regression function under the control condition and  $\hat{\mu}_0$  is an estimate of  $\mu_0$ .

$X_i - X_{j(i)}$  is often referred to as the **matching discrepancy**.

These “bias-corrected” matching estimators behave well even if  $\mu_0$  is estimated using a simple linear regression (ie.  $\mu_0(x) = \beta_0 + \beta_1 x$ ) (Abadie and Imbens, 2005)

# Matching with Bias Correction

Each treated observation contributes

$$\mu_0(X_i) - \mu_0(X_{j(i)})$$

to the bias.

Bias-corrected matching:

$$\tilde{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} \left( (Y_i - Y_{j(i)}) - (\hat{\mu}_0(X_i) - \hat{\mu}_0(X_{j(i)})) \right)$$

The large sample distribution of this estimator (for the case of matching with replacement) is (basically) standard normal.  $\mu_0$  is usually estimated using a simple linear regression (ie.  $\mu_0(x) = \beta_0 + \beta_1 x$ ).

In R: `Match(Y,Tr, X,BiasAdjust = TRUE)`

## Bias Adjustment with Matched Data

unit	Potential Outcome under Treatment	Potential Outcome under Control		
$i$	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	?	1	3
2	1	?	1	1
3	0	?	1	10
4		0	0	2
5		9	0	3
6		1	0	8

What is  $\tilde{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} \left( (Y_i - Y_{j(i)}) - (\hat{\mu}_0(X_i) - \hat{\mu}_0(X_{j(i)})) \right)$ ?

# Bias Adjustment with Matched Data

unit	Potential Outcome under Treatment	Potential Outcome under Control		
i	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	9	1	3
2	1	0	1	1
3	0	1	1	10
4		0	0	2
5		9	0	3
6		1	0	8

What is  $\tilde{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} \left( (Y_i - Y_{j(i)}) - (\hat{\mu}_0(X_i) - \hat{\mu}_0(X_{j(i)})) \right)$ ?

Estimate  $\hat{\mu}_0(x) = \beta_0 + \beta_1 x = 5 - .4x$ .

# Bias Adjustment with Matched Data

unit	Potential Outcome under Treatment	Potential Outcome under Control		
i	$Y_{1i}$	$Y_{0i}$	$D_i$	$X_i$
1	6	9	1	3
2	1	0	1	1
3	0	1	1	10
4		0	0	2
5		9	0	3
6		1	0	8

What is  $\tilde{\tau}_{ATT} = \frac{1}{N_1} \sum_{D_i=1} \left( (Y_i - Y_{j(i)}) - (\hat{\mu}_0(X_i) - \hat{\mu}_0(X_{j(i)})) \right)$ ?

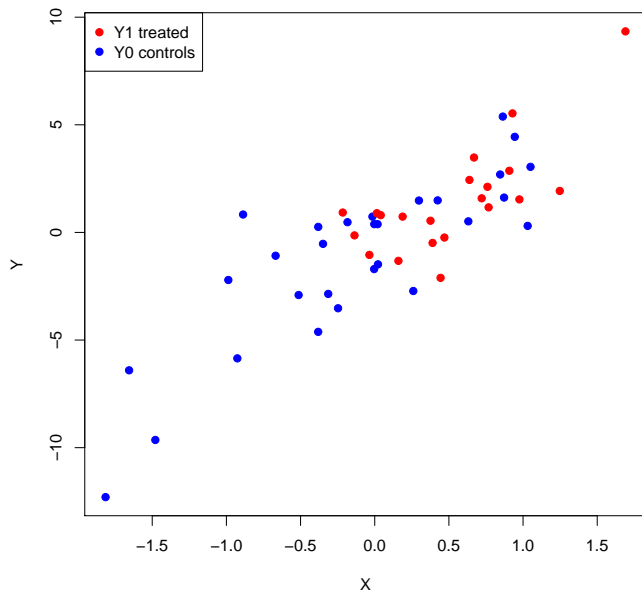
Estimate  $\hat{\mu}_0(x) = \beta_0 + \beta_1 x = 5 - .4x$ . Now plug in:

$$\begin{aligned}\hat{\tau}_{ATT} &= 1/3\{((6 - 9) - (\hat{\mu}_0(3) - \hat{\mu}_0(3))) \\ &+ ((1 - 0) - (\hat{\mu}_0(1) - \hat{\mu}_0(2))) \\ &+ ((0 - 1) - (\hat{\mu}_0(10) - \hat{\mu}_0(8)))\} \\ &= -0.86\end{aligned}$$

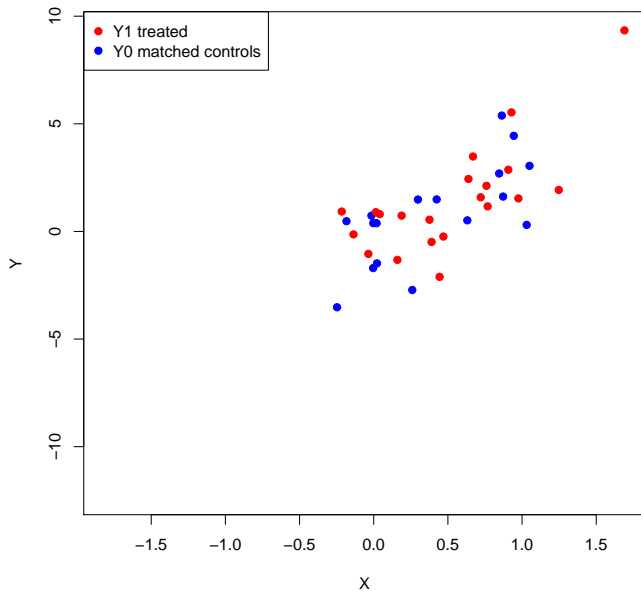
Unadjusted:  $1/3((6 - 9) + (1 - 0) + (0 - 1)) = -1$



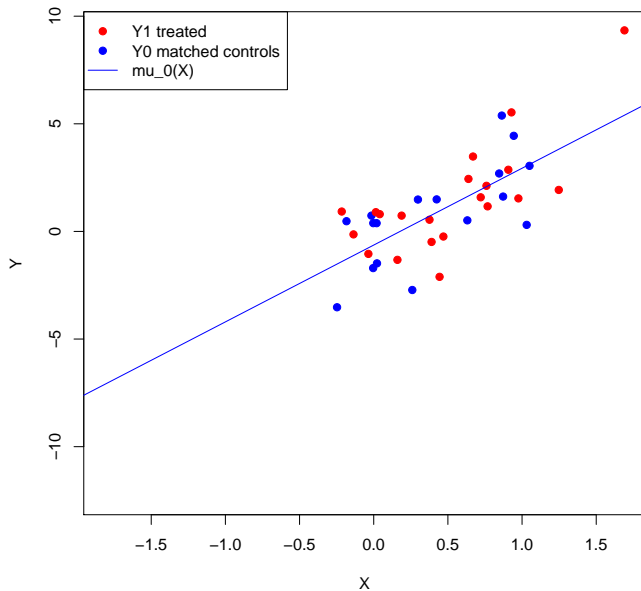
# Before Matching



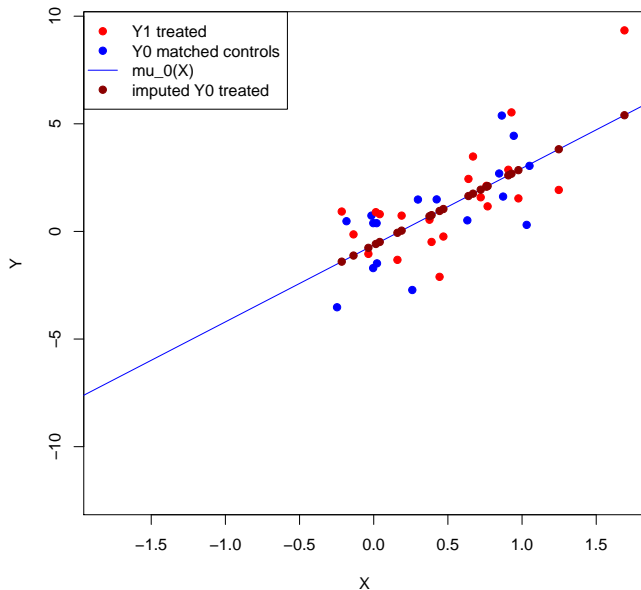
# After Matching



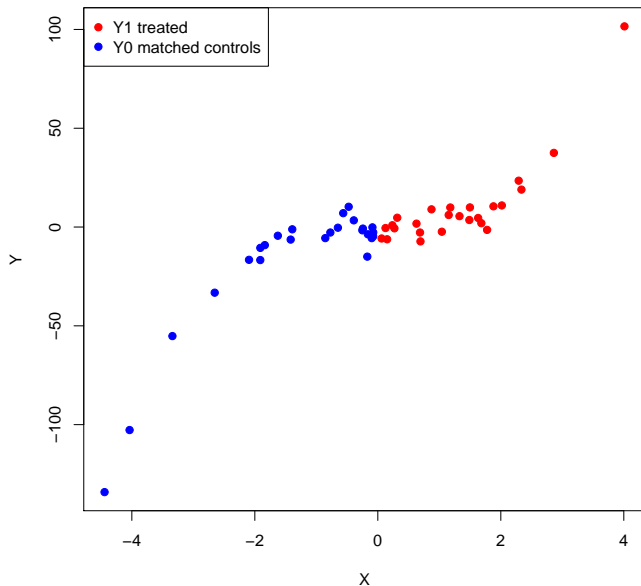
# After Matching: Imputation Function



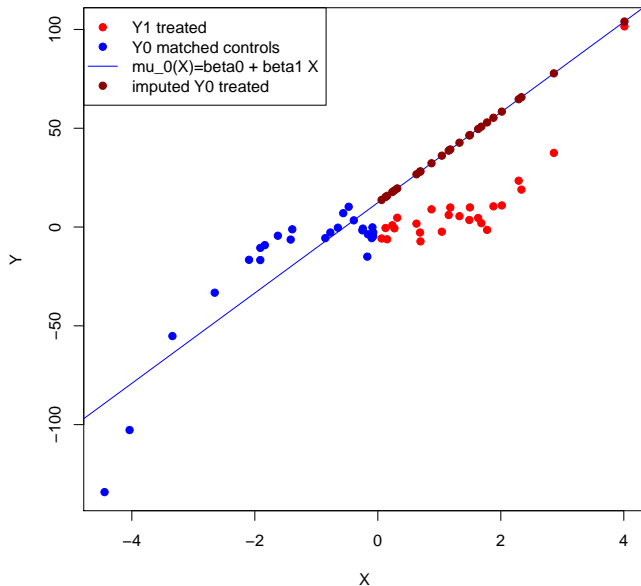
# After Matching: Imputation of missing $Y_0$



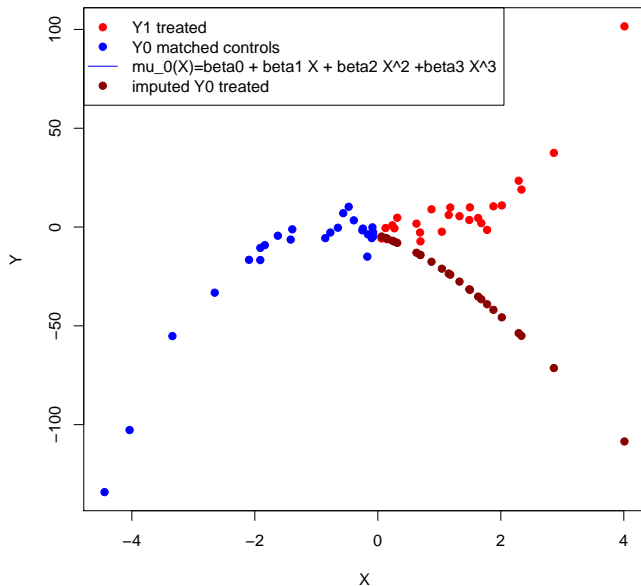
# After Matching: No Overlap in $Y_0$



# After Matching: Imputation of missing $Y_0$



# After Matching: Imputation of missing $Y_0$



# Choices when Matching

- With or Without Replacement?



# Choices when Matching

- With or Without Replacement?
- How many matches?

# Choices when Matching

- With or Without Replacement?
- How many matches?
- Which Matching Algorithm?
  - Genetic Matching
  - Kernel Matching
  - Full Matching
  - Coarsened Exact Matching
  - Matching as Pre-processing
  - Propensity Score Matching

# Choices when Matching

- With or Without Replacement?
- How many matches?
- Which Matching Algorithm?
  - Genetic Matching
  - Kernel Matching
  - Full Matching
  - Coarsened Exact Matching
  - Matching as Pre-processing
  - Propensity Score Matching
- Use whatever gives you the best balance! Checking balance is important to get a sense for how much extrapolation is needed
  - Should check balance on interactions and higher moments
- With insufficient overlap, all adjustment methods are problematic because we have to heavily rely on a model to impute missing potential outcomes.

# Balance Checks

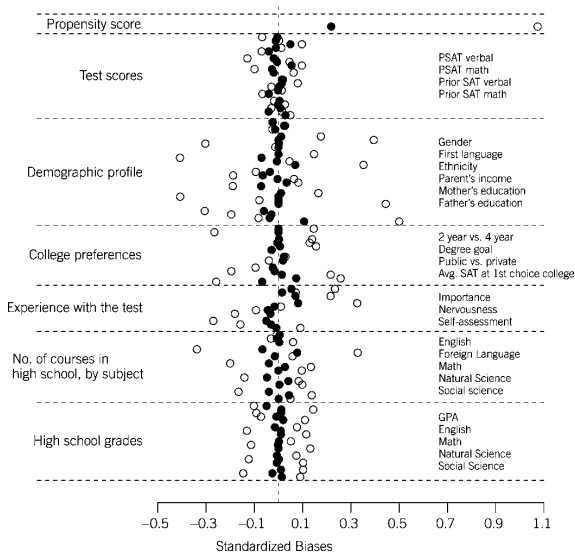


Figure 3. Standardized Biases Without Stratification or Matching, Open Circles, and Under the Optimal [.5, 2] Full Match, Shaded Circles.

# Balance Checks - Lyall (2010)

**TABLE 2. Balance Summary Statistics and Tests: Russian and Chechen Sweeps**

Pretreatment Covariates	Mean Treated	Mean Control	Mean Difference	Std. Bias	Rank Sum Test	K-S Test
<i>Demographics</i>						
Population	8.657	8.606	0.049	0.033	0.708	0.454
Tariqa	0.076	0.048	0.028	0.104	0.331	—
Poverty	1.917	1.931	-0.016	-0.024	0.792	1.000
<i>Spatial</i>						
Elevation	5.078	5.233	-0.155	-0.135	0.140	0.228
Isolation	1.007	1.070	-0.063	-0.096	0.343	0.851
Groznyy	0.131	0.138	-0.007	-0.018	0.864	—
<i>War Dynamics</i>						
TAC	0.241	0.282	-0.041	-0.095	0.424	—
Garrison	0.379	0.414	-0.035	-0.072	0.549	—
Rebel	0.510	0.441	0.070	0.139	0.240	—
<i>Selection</i>						
Presweep violence	3.083	3.117	-0.034	0.009	0.454	0.292
Large-scale theft	0.034	0.055	-0.021	-0.115	0.395	—
Killing	0.117	0.090	0.027	0.084	0.443	—
<i>Violence Inflicted</i>						
Total abuse	0.970	0.833	0.137	0.124	0.131	0.454
Prior sweeps	1.729	1.812	-0.090	-0.089	0.394	0.367
<i>Other</i>						
Month	7.428	6.986	0.442	0.130	0.260	0.292
Year	2004.159	2004.110	0.049	0.043	0.889	1.000

Note: 145 matched pairs. Matching with replacement.

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables**
  - Subclassification
  - Matching
  - Propensity Scores**
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Identification with Propensity Scores

## Definition

Propensity score is defined as the selection probability conditional on the confounding variables:  $\pi(X) = \Pr(D = 1|X)$

## Identification Assumption

- 1  $(Y_1, Y_0) \perp\!\!\!\perp D | X$  (*selection on observables*)
- 2  $0 < \Pr(D = 1|X) < 1$  with probability one (*common support*)

## Identification Result

*Under selection on observables we have  $(Y_1, Y_0) \perp\!\!\!\perp D | \pi(X)$ , ie. conditioning on the propensity score is enough to have independence between the treatment indicator and potential outcomes. Implies substantial dimension reduction.*

# Matching on the Propensity Score

## Corollary

If  $(Y_1, Y_0) \perp\!\!\!\perp D | X$ , then

$$\mathbf{E}[Y|D = 1, \pi(X) = \pi_0] - \mathbf{E}[Y|D = 0, \pi(X) = \pi_0] = \mathbf{E}[Y_1 - Y_0 | \pi(X) = \pi_0]$$

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

- 1 Estimate the propensity score  $\pi(X) = P(D = 1|X)$  (e.g. using logit/probit regression, machine learning methods, etc)



# Matching on the Propensity Score

## Corollary

If  $(Y_1, Y_0) \perp\!\!\!\perp D | X$ , then

$$\mathbf{E}[Y|D = 1, \pi(X) = \pi_0] - \mathbf{E}[Y|D = 0, \pi(X) = \pi_0] = \mathbf{E}[Y_1 - Y_0 | \pi(X) = \pi_0]$$

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

- 1 Estimate the propensity score  $\pi(X) = P(D = 1|X)$  (e.g. using logit/probit regression, machine learning methods, etc)
- 2 Match or subclassify on propensity score.

## Estimating the Propensity Score

- Given selection on observables we have  $(Y_1, Y_0) \perp\!\!\!\perp D \mid \pi(X)$  which implies the balancing property of the propensity score:

$$\Pr(X \mid D = 1, \pi(X)) = \Pr(X \mid D = 0, \pi(X))$$

## Estimating the Propensity Score

- Given selection on observables we have  $(Y_1, Y_0) \perp\!\!\!\perp D | \pi(X)$  which implies the balancing property of the propensity score:

$$\Pr(X|D = 1, \pi(X)) = \Pr(X|D = 0, \pi(X))$$

- We can use this to check if our estimated propensity score actually produces balance:  $P(X|D = 1, \hat{\pi}(X)) = P(X|D = 0, \hat{\pi}(X))$

# Estimating the Propensity Score

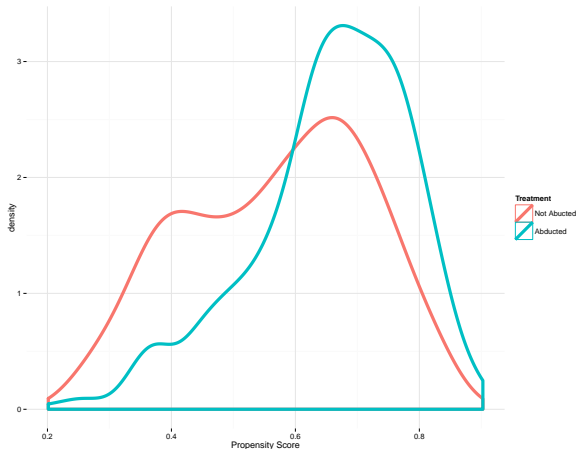
- Given selection on observables we have  $(Y_1, Y_0) \perp\!\!\!\perp D | \pi(X)$  which implies the balancing property of the propensity score:

$$\Pr(X|D = 1, \pi(X)) = \Pr(X|D = 0, \pi(X))$$

- We can use this to check if our estimated propensity score actually produces balance:  $P(X|D = 1, \hat{\pi}(X)) = P(X|D = 0, \hat{\pi}(X))$
- To properly model the assignment mechanism, we need to include important confounders correlated with treatment and outcome
- Need to find the correct functional form, miss-specified propensity scores can lead to bias. Any methods can be used (probit, logit, etc.)
- Estimate  $\mapsto$  Check Balance  $\mapsto$  Re-estimate  $\mapsto$  Check Balance

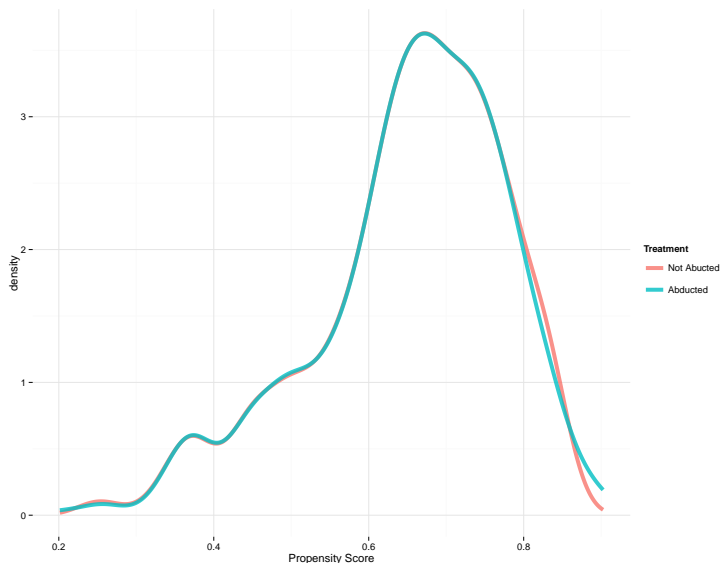
# Example: Blattman (2010)

```
pscore.fmla <- as.formula(paste("abd~",paste(names(covar),collapse="+")))
abd <- data$abd
pscore_model <- glm(pscore.fmla, data = data,
  family = binomial(link = logit))
pscore <- predict(pscore_model, type = "response")
```



# Blattman (2010): Match on the Propensity Score

```
match.pscore <- Match(Tr=abd, X=pscore, M=1, estimand="ATT")
```



## Blattman (2010): Check Balance

```
match.pscore <-  
+ MatchBalance(abd ~ age, data=data, match.out = match.pscore)
```

```
***** (V1) age *****
```

	Before Matching	After Matching
mean treatment.....	21.366	21.366
mean control.....	20.151	20.515
std mean diff.....	24.242	16.976
var ratio (Tr/Co).....	1.0428	0.98412
T-test p-value.....	0.0012663	0.0034409
KS Bootstrap p-value..	0.016	0.034
KS Naive p-value.....	0.024912	0.070191
KS Statistic.....	0.11227	0.077899

# Blattman (2010): Mahalanobis Distance Matching

```
match.mah <- Match(Tr=abd, X=covar, M=1, estimand="ATT", Weight = 3)
MatchBalance(abd ~ age, data=data, match.out = match.mah)
```

```
***** (V1) age *****
```

	Before Matching	After Matching
mean treatment.....	21.366	21.366
mean control.....	20.151	21.154
std mean diff.....	24.242	4.2314
var ratio (Tr/Co).....	1.0428	1.0336
T-test p-value.....	0.0012663	3.0386e-05
KS Bootstrap p-value..	0.008	0.798
KS Naive p-value.....	0.024912	0.94687
KS Statistic.....	0.11227	0.034261



# Blattman (2010): Genetic Matching

```
genout <- GenMatch(Tr=abd,X=covar,BalanceMatrix=covar,estimand="ATT",
                  pop.size=1000)
match.gen <- Match(Tr=abd, X=covar,M=1,estimand="ATT",Weight.matrix=genout)
gen.bal <- MatchBalance(abd~age,match.out=match.gen,data=covar)
```

```
***** (V1) age *****
```

	Before Matching	After Matching
mean treatment.....	21.366	21.366
mean control.....	20.151	21.225
std mean diff.....	24.242	2.8065
var ratio (Tr/Co).....	1.0428	1.1337
T-test p-value.....	0.0012663	0.21628
KS Bootstrap p-value..	0.008	0.454
KS Naive p-value.....	0.024912	0.68567
KS Statistic.....	0.11227	0.046512

# Outline

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables**
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression**
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Identification under Selection on Observables: Regression

Consider the linear regression of  $Y_i = \beta_0 + \tau D_i + X_i' \beta + \epsilon_i$ .

Given **selection on observables**, there are mainly three identification scenarios:

# Identification under Selection on Observables: Regression

Consider the linear regression of  $Y_i = \beta_0 + \tau D_i + X_i' \beta + \epsilon_i$ .

Given **selection on observables**, there are mainly three identification scenarios:

- 1 Constant treatment effects and outcomes are linear in  $X$ 
  - $\tau$  will provide unbiased and consistent estimates of ATE.
- 2 Constant treatment effects and unknown functional form
  - $\tau$  will provide well-defined linear approximation to the average causal response function  $\mathbf{E}[Y|D = 1, X] - \mathbf{E}[Y|D = 0, X]$ . Approximation may be very poor if  $\mathbf{E}[Y|D, X]$  is misspecified and then  $\tau$  may be biased for the ATE.

# Identification under Selection on Observables: Regression

Consider the linear regression of  $Y_i = \beta_0 + \tau D_i + X_i' \beta + \epsilon_i$ .

Given **selection on observables**, there are mainly three identification scenarios:

- 1 Constant treatment effects and outcomes are linear in  $X$ 
  - $\tau$  will provide unbiased and consistent estimates of ATE.
- 2 Constant treatment effects and unknown functional form
  - $\tau$  will provide well-defined linear approximation to the average causal response function  $\mathbf{E}[Y|D = 1, X] - \mathbf{E}[Y|D = 0, X]$ . Approximation may be very poor if  $\mathbf{E}[Y|D, X]$  is misspecified and then  $\tau$  may be biased for the ATE.
- 3 Heterogeneous treatment effects ( $\tau$  differs for different values of  $X$ )
  - If outcomes are linear in  $X$ ,  $\tau$  is unbiased and consistent estimator for conditional-variance-weighted average of the underlying causal effects. This average is often different from the ATE.

## Identification Assumption

- 1 *Constant treatment effect:  $\tau = Y_{1i} - Y_{0i}$  for all  $i$*
- 2 *Control outcome is linear in  $X$ :  $Y_{0i} = \beta_0 + X_i'\beta + \epsilon_i$  with  $\epsilon_i \perp\!\!\!\perp X_i$  (no omitted variables and linearly separable confounding)*

## Identification Result

*Then  $\tau_{ATE} = \mathbf{E}[Y_1 - Y_0]$  is identified by a regression of the observed outcome on the covariates and the treatment indicator*

$$Y_i = \beta_0 + \tau D_i + X_i'\beta + \epsilon_i$$

# Regression with Heterogeneous Effects

What is regression estimating when we allow for heterogeneity?

# Regression with Heterogeneous Effects

What is regression estimating when we allow for heterogeneity?

Suppose that we wanted to estimate  $\tau_{OLS}$  using a **fully saturated** regression model:

$$Y_i = \sum_x B_{xi} \beta_x + \tau_{OLS} D_i + e_i$$

where  $B_{xi}$  is a dummy variable for unique combination of  $X_i$ .

Because this regression is fully saturated, it is linear in the covariates (i.e. linearity assumption holds by construction).



# Heterogenous Treatment Effects

With two  $X$  strata there are two stratum-specific average causal effects that are averaged to obtain the ATE or ATT.

Subclassification weights the stratum-effects by the marginal distribution of  $X$ , i.e. weights are proportional to the share of units in each stratum:

$$\begin{aligned}\tau_{ATE} &= (\mathbf{E}[Y|D = 1, X = 1] - \mathbf{E}[Y|D = 0, X = 1])\Pr(X = 1) \\ &+ (\mathbf{E}[Y|D = 1, X = 2] - \mathbf{E}[Y|D = 0, X = 2])\Pr(X = 2)\end{aligned}$$

Regression weights by the marginal distribution of  $X$  and the conditional variance of  $\text{Var}[D|X]$  in each stratum:

$$\begin{aligned}\tau_{OLS} &= (\mathbf{E}[Y|D = 1, X = 1] - \mathbf{E}[Y|D = 0, X = 1]) \frac{\text{Var}[D|X = 1]\Pr(X = 1)}{\sum_X \text{Var}[D|X = x] \Pr(X = x)} \\ &+ (\mathbf{E}[Y|D = 1, X = 2] - \mathbf{E}[Y|D = 0, X = 2]) \frac{\text{Var}[D|X = 2]\Pr(X = 2)}{\sum_X \text{Var}[D|X = x] \Pr(X = x)}\end{aligned}$$

- So strata with a higher  $\text{Var}[D|X]$  receive higher weight. These are the strata with propensity scores close to .5
- Strata with propensity score close to 0 or 1 receive lower weight
- OLS is a minimum-variance estimator of  $\tau_{OLS}$  so it downweights strata where the average causal effects are less precisely estimated.

# Heterogenous Treatment Effects

With two  $X$  strata there are two stratum-specific average causal effects that are averaged to obtain the ATE or ATT.

Subclassification weights the stratum-effects by the marginal distribution of  $X$ , i.e. weights are proportional to the share of units in each stratum:

$$\begin{aligned}\tau_{ATE} &= (\mathbf{E}[Y|D = 1, X = 1] - \mathbf{E}[Y|D = 0, X = 1])\Pr(X = 1) \\ &+ (\mathbf{E}[Y|D = 1, X = 2] - \mathbf{E}[Y|D = 0, X = 2])\Pr(X = 2)\end{aligned}$$

Regression weights by the marginal distribution of  $X$  and the conditional variance of  $\text{Var}[D|X]$  in each stratum:

$$\begin{aligned}\tau_{OLS} &= (\mathbf{E}[Y|D = 1, X = 1] - \mathbf{E}[Y|D = 0, X = 1]) \frac{\text{Var}[D|X = 1]\Pr(X = 1)}{\sum_X \text{Var}[D|X = x] \Pr(X = x)} \\ &+ (\mathbf{E}[Y|D = 1, X = 2] - \mathbf{E}[Y|D = 0, X = 2]) \frac{\text{Var}[D|X = 2]\Pr(X = 2)}{\sum_X \text{Var}[D|X = x] \Pr(X = x)}\end{aligned}$$

- Whenever both weighting components are misaligned (e.g. the PS is close to 0 or 1 for relatively large strata) then  $\tau_{OLS}$  diverges from  $\tau_{ATE}$  or  $\tau_{ATT}$ .
- We still need overlap in the data (treated/untreated units in all strata of  $X$ )! Otherwise OLS will interpolate/extrapolate  $\rightarrow$  model-dependent results

# Conclusion: Regression

Is regression evil? 😞

# Conclusion: Regression

Is regression evil? 😊

- Its ease sometimes results in lack of thinking. So only a little. 😊
- For descriptive inference, very useful!
  - Good tool for characterizing the conditional expectation function (CEF)
  - But other less restrictive tools are also available for that task (machine learning)

# Conclusion: Regression

Is regression evil? 😊

- Its ease sometimes results in lack of thinking. So only a little. 😊
- For descriptive inference, very useful!
  - Good tool for characterizing the conditional expectation function (CEF)
  - But other less restrictive tools are also available for that task (machine learning)
- For causal analysis, always need to ask yourself if *linearly* separable confounding is plausible.
  - A regression is causal when the CEF it approximates is causal.
  - Still need to check common support!
  - Results will be highly sensitive if the treated and controls are far apart (e.g. standardized difference above .2)
- Think about what your **estimand** is: because of variance weighting, coefficient from your regression may not capture ATE if effects are heterogeneous

- 1 What Makes a Good Observational Study?
- 2 Removing Bias by Conditioning
- 3 Identification Under Selection on Observables
  - Subclassification
  - Matching
  - Propensity Scores
  - Regression
- 4 When Do Observational Studies Recover Experimental Benchmarks?

# Dehejia and Wabha (1999) 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					
	(1) Unadjusted	(2) Adjusted <sup>a</sup>	Quadratic in score <sup>b</sup> (3)	Stratifying on the score			Matching on the score	
				(4) Unadjusted	(5) Adjusted	(6) Observations <sup>c</sup>	(7) Unadjusted	(8) Adjusted <sup>d</sup>
NSW	1,794 (633)	1,672 (638)						
PSID-1 <sup>e</sup>	-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 <sup>f</sup>	-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 <sup>f</sup>	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 <sup>g</sup>	-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 <sup>g</sup>	-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 <sup>g</sup>	-635 (657)	1,326 (798)	556 (951)	1,252 (1,617)	2,219 (2,082)	514	587 (1,496)	662 (776)

<sup>a</sup> Least squares regression: RE78 on a constant, a treatment indicator, age, age<sup>2</sup>, education, no degree, black, Hispanic, RE74, RE75.

<sup>b</sup> 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).

<sup>c</sup> 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.

<sup>d</sup> 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:

<sup>e</sup> PSID-1: Prob ( $T_i = 1$ ) = F(age, age<sup>2</sup>, education, education<sup>2</sup>, married, no degree, black, Hispanic, RE74, RE75, RE74<sup>2</sup>, RE75<sup>2</sup>, u74 \* black).

<sup>f</sup> PSID-2 and PSID-3: Prob ( $T_i = 1$ ) = F(age, age<sup>2</sup>, education, education<sup>2</sup>, no degree, married, black, Hispanic, RE74, RE74<sup>2</sup>, RE75, RE75<sup>2</sup>, u74, u75).

<sup>g</sup> CPS-1, CPS-2, and CPS-3: Prob ( $T_i = 1$ ) = F(age, age<sup>2</sup>, education, education<sup>2</sup>, no degree, married, black, Hispanic, RE74, RE75, u74, u75, education \* RE74, age<sup>3</sup>).