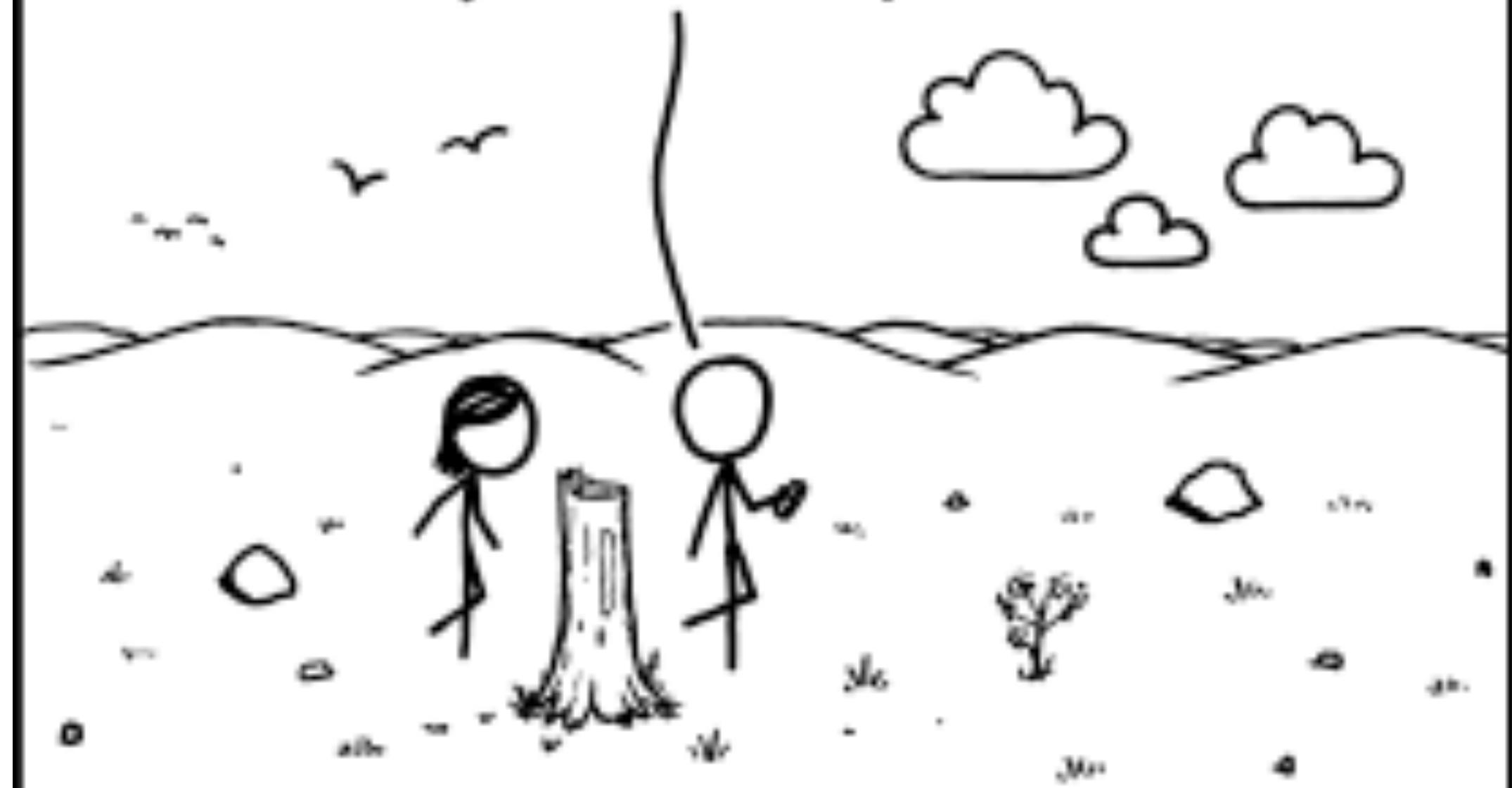


Observational Studies: DAY TWO

Lecture 2a

MY LIFE IS BASICALLY A BIG CONTROLLED TRIAL OF WHETHER I'M MORE LIKELY TO WALK INTO SOMETHING WHILE LOOKING AT A BOOK, MY PHONE, OR THE SKY.



THE WEIRD THING IS THAT THE RATE FOR THE CONTROL GROUP IS SO HIGH.



WALKING IS HARD, OKAY?



Plan for today

- LaLonde: Subclassification
- Regression for obs studies
- Matching
 - Exact matching
 - Approximate matching
 - Matching **and** regression (not matching **or** regression)
- LaLonde: Matching

Cold open: the assumption of "no unmeasured confounding given X" means

Conditional on observed X, treatment assignment is as good as random

Conditional on observed X, there are no unmeasured covariates that predict both...

Conditional on observed X, every unit has some chance of receiving either treatm...

B and C only

A and B only

Cold open: the assumption of "no unmeasured confounding given X" means

Conditional on observed X, treatment assignment is as good as random

0%

Conditional on observed X, there are no unmeasured covariates that predict both the treatment and outcome

0%

Conditional on observed X, every unit has some chance of receiving either treatment or control

0%

B and C only

0%

A and B only

0%

Cold open: the assumption of "no unmeasured confounding given X" means

Conditional on observed X, treatment assignment is as good as random

0%

Conditional on observed X, there are no unmeasured covariates that predict both the treatment and outcome

0%

Conditional on observed X, every unit has some chance of receiving either treatment or control

0%

B and C only

0%

A and B only

0%

Reminder:
Smoking and Mortality

THE EFFECTIVENESS OF ADJUSTMENT BY
SUBCLASSIFICATION IN REMOVING BIAS IN
OBSERVATIONAL STUDIES

W. G. COCHRAN

Harvard University, Cambridge, Mass., U. S. A.

Mortality Rate

Table 5.1: Death rates per 1,000 person-years ([Cochran 1968](#))

Smoking group	Canada	UK	US
Non-smokers	20.2	11.3	13.5
Cigarettes	20.5	14.1	13.5
Cigars/pipes	35.5	20.7	17.4

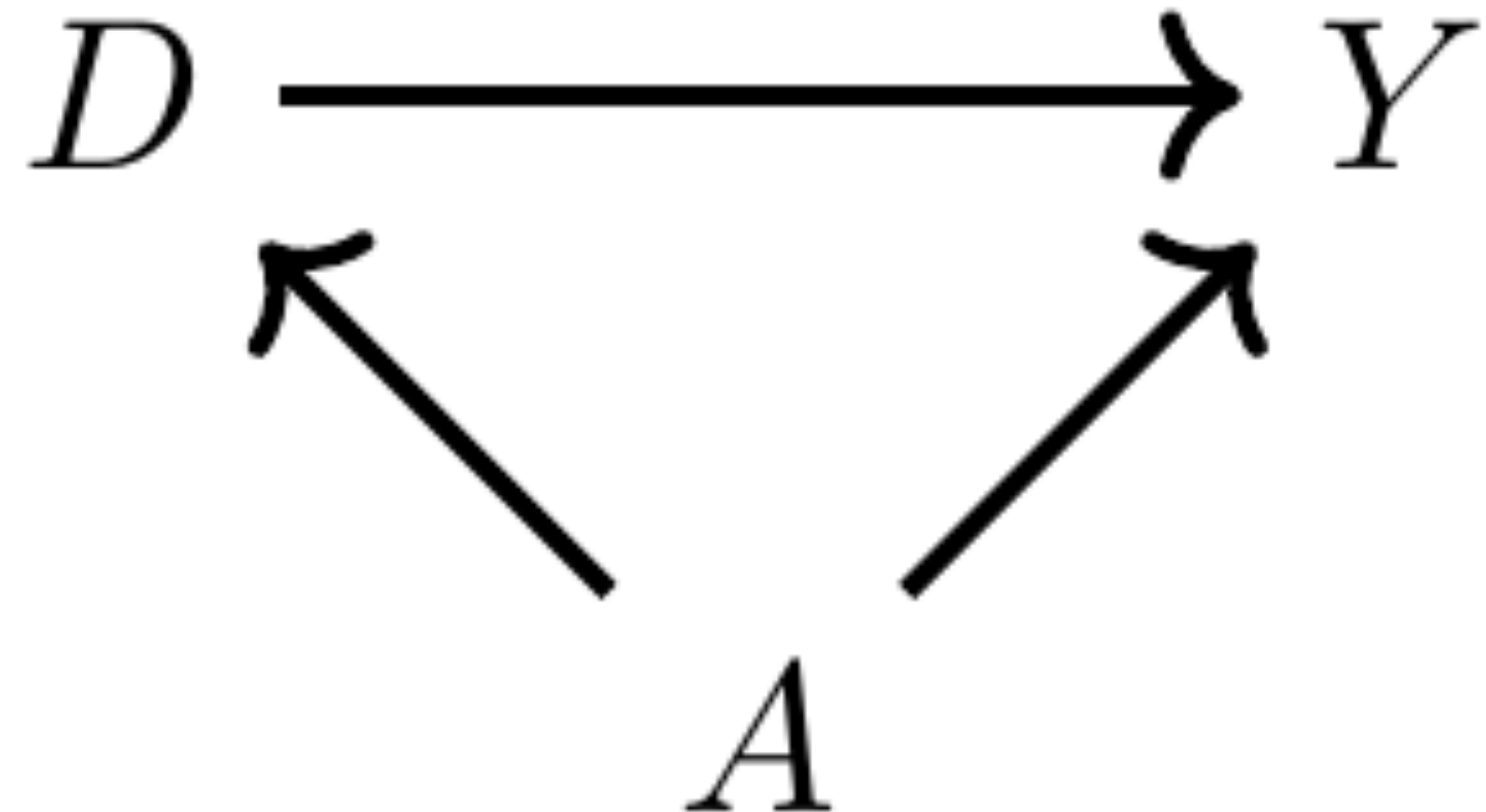
Age

Table 5.2: Mean ages, years ([Cochran 1968](#)).

Smoking group	Canada	British	US
Non-smokers	54.9	49.1	57.0
Cigarettes	50.5	49.8	53.2
Cigars/pipes	65.9	55.7	59.7

SMOKING

DEATH



AGE

Distribution of smoking by age

Table 5.3: Subclassification example.

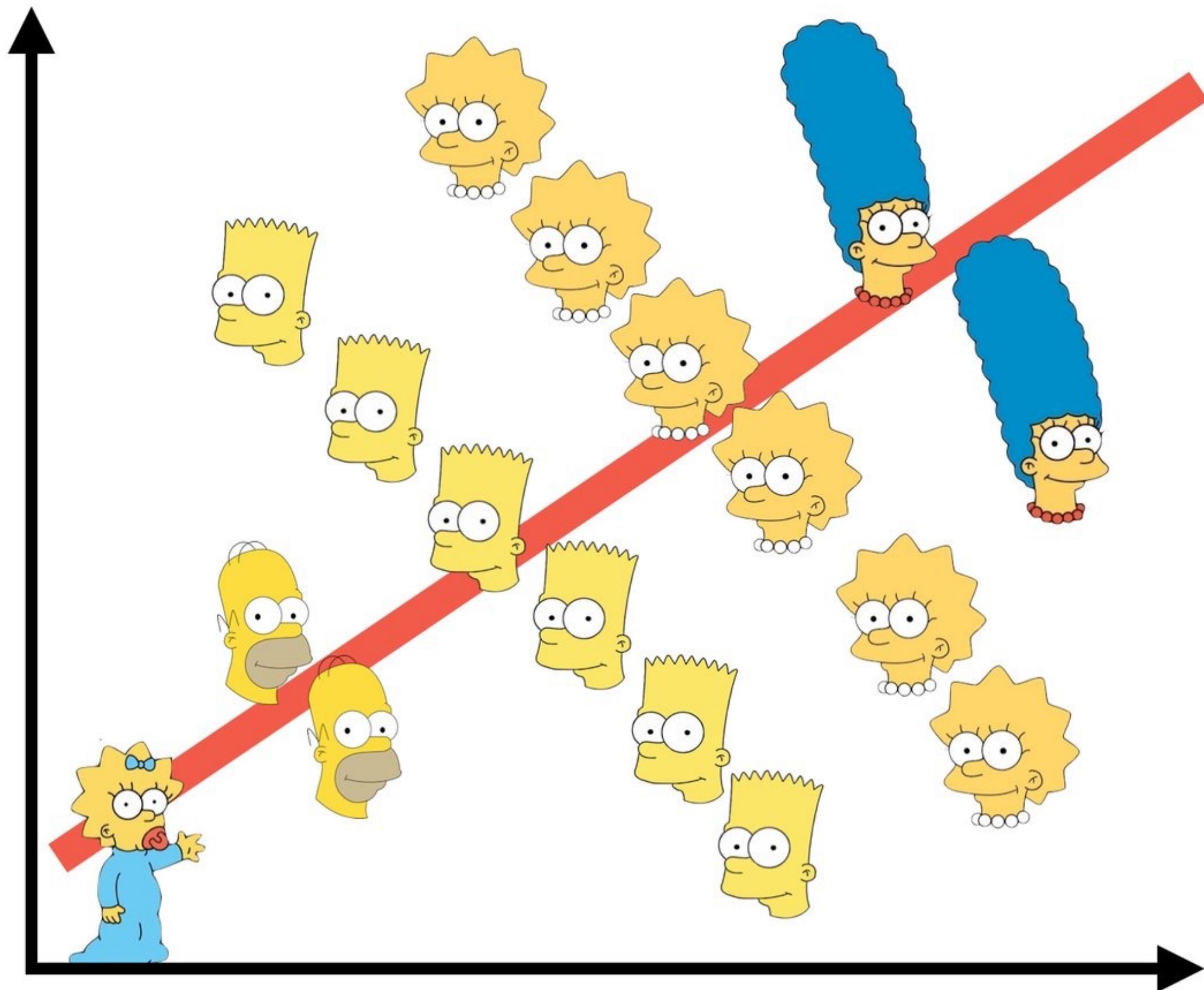
	Death rates	# of Cigarette smokers	# of Pipe or cigar smokers
Age 20-40	20	65	10
Age 41-70	40	25	25
Age ≥ 71	60	10	65
Total		100	100

Age-adjusted mortality

Table 5.4: Adjusted mortality rates using 3 age groups ([Cochran 1968](#)).

Smoking group	Canada	UK	US
Non-smokers	20.2	11.3	13.5
Cigarettes	29.5	14.8	21.2
Cigars/pipes	19.8	11.0	13.7

Simpson's Paradox



THERE IS NO MAGIC

We're assuming away the problem!

(once we know **educ**)

The Obligatory LaLonde Study

Evaluating the Econometric Evaluations of Training Programs with Experimental Data

By ROBERT J. LALONDE*

This paper compares the effect on trainee earnings of an employment program that was run as a field experiment where participants were randomly assigned to treatment and control groups with the estimates that would have been produced by an econometrician. This comparison shows that many of the econometric procedures do not replicate the experimentally determined results, and it suggests that researchers should be aware of the potential for specification errors in other nonexperimental evaluations.

TABLE 5—EARNINGS COMPARISONS AND ESTIMATED TRAINING EFFECTS FOR THE NSW MALE PARTICIPANTS USING COMPARISON GROUPS FROM THE *PSID* AND THE *CPS-SSA*^{a,b}

Name of Comparison Group ^d	Comparison Group Earnings Growth 1975–78 (1)	NSW Treatment Earnings Less Comparison Group Earnings				Difference in Differences: Difference in Earnings Growth 1975–78 Treatments Less Comparisons		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		Without Age (6)	With Age (7)	Unadjusted (8)	Adjusted ^c (9)	
		Unadjusted (2)	Adjusted ^c (3)	Unadjusted (4)	Adjusted ^c (5)					
Controls	\$2,063 (325)	\$39 (383)	\$–21 (378)	\$886 (476)	\$798 (472)	\$847 (560)	\$856 (558)	\$897 (467)	\$802 (467)	\$662 (506)
<i>PSID</i> -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)
<i>PSID</i> -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)
<i>PSID</i> -3	(\$3,322 (780))	(\$455 (539))	\$455 (704)	\$697 (760)	–\$509 (967)	\$242 (884)	–\$1,325 (1078)	\$629 (757)	–\$552 (967)	\$397 (1103)
<i>CPS-SSA</i> -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)
<i>CPS-SSA</i> -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)
<i>CPS-SSA</i> -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)

LaLonde: Covariate Imbalance

Variable	NSW		Full Samples	
	Treated (1)	Control (2)	CPS-1 (3)	CPS-3 (4)
Age	25.82	25.05	33.23	28.03
Years of schooling	10.35	10.09	12.03	10.24
Black	0.84	0.83	0.07	0.20
Hispanic	0.06	0.11	0.07	0.14
Dropout	0.71	0.83	0.30	0.60
Married	0.19	0.15	0.71	0.51
1974 earnings	2,096	2,107	14,017	5,619
1975 earnings	1,532	1,267	13,651	2,466
Number of Obs.	185	260	15,992	429

(From MHE Table 3.3.2)

Covariate balance for race/ethnicity

Race/Eth. Category	N	Prop. Treated
White	299	0.06
Hispanic	72	0.15
Black	243	0.64

	Prop. Black	Prop. Hispanic
Treated	0.84	0.06
Control	0.20	0.14

Subclassify by race/ethnicity

Naive Difference in Means: -\$635

[i.e., program has an *adverse* effect]

Subclassify by race/ethnicity

Naive Difference in Means: -\$635

[i.e., program has an *adverse* effect]

Race/Eth. Category	Diff-in-Means	N	Prop. Treated
White	+\$103	299	0.06
Hispanic	+\$19	72	0.15
Black	+\$1,283	243	0.64

Subclassify by race/ethnicity

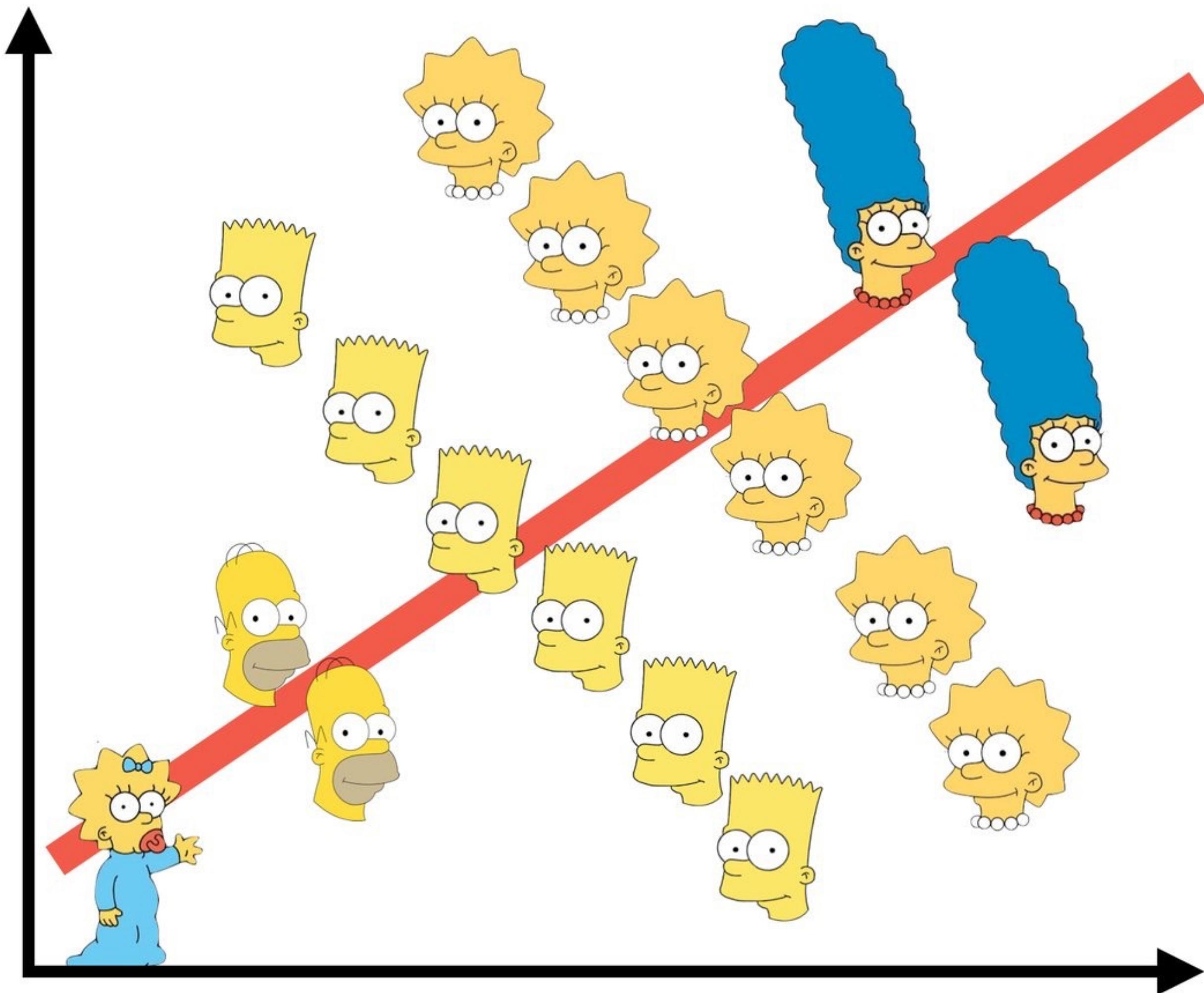
Naive Difference in Means: -\$635

[i.e., program has an *adverse* effect]

Race/Eth. Category	Diff-in-Means	N	Prop. Treated
White	+\$103	299	0.06
Hispanic	+\$19	72	0.15
Black	+\$1,283	243	0.64

Re-Weighted Difference in Means: +\$560

Simpson's Paradox



Regression

Regression with discrete X
 $Y \sim \text{treat} + \text{race_eth}$

Call:

```
lm_robust(formula = re78 ~ treat + black + hispan, data = laonode)
```

Standard error type: HC2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	7570.1	435.5	17.3824	2.840e-55	6714.9	8425	610
treat	896.2	742.1	1.2077	2.276e-01	-561.1	2354	610
black	-2468.5	700.9	-3.5217	4.608e-04	-3845.0	-1092	610
hispan	-600.4	935.9	-0.6415	5.214e-01	-2438.3	1238	610

Multiple R-squared: 0.01702 , Adjusted R-squared: 0.01218

F-statistic: 4.243 on 3 and 610 DF, p-value: 0.005568

[Differs from stratified est of +\$560. Why?]

Regression

Regression with discrete X
 $Y \sim \text{treat} + \text{race_eth}$

Call:

```
lm_robust(formula = re78 ~ treat + black + hispan, data = laonode)
```

Standard error type: HC2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	700.9	435.5	17.3824	2.840e-55	6714.9	8425	610
treat	896.2	742.1	1.2077	2.276e-01	-561.1	2354	610
black	-159.7	700.9	-3.5217	4.608e-04	-3845.0	-1092	610
hispan	-600.4	935.9	-0.6415	5.214e-01	-2438.3	1238	610

Multiple R-squared: 0.01702 , Adjusted R-squared: 0.01218

F-statistic: 4.243 on 3 and 610 DF, p-value: 0.005568

[Differs from stratified est of +\$560. Why?]

Subclassify by more?

Race/Eth. Category	HS Degree	ATE	N	Prop. Treated
White	No	+\$1,868	163	0.06
Hispanic	No	+\$385	55	0.16
Black	No	+\$885	169	0.67
White	Yes	-\$1,903	136	0.07
Hispanic	Yes	-\$1,677	17	0.12
Black	Yes	+\$2,700	74	0.58

Subclassify by more?

Race/Eth. Category	HS Degree	ATE	N	Prop. Treated
White	No	+\$1,868	163	0.06
Hispanic	No	+\$385	55	0.16
Black	No	+\$885	169	0.67
White	Yes	-\$1,903	136	0.07
Hispanic	Yes	-\$1,677	17	0.12
Black	Yes	+\$2,700	74	0.58

Re-Weighted Difference in Means: +\$632

Subclassify by more?

Race/Eth. Category	HS Degree	ATE	N	Prop. Treated
White	No	+\$1,868	163	0.06
Hispanic	No	+\$385	55	0.16
Black	No	+\$885	169	0.67
White	Yes	-\$1,903	136	0.07
Hispanic	Yes	-\$1,677	17	0.12
Black	Yes	+\$2,700	74	0.58

Re-Weighted Difference in Means: +\$632

Regression (take 2)

Regression with discrete X ?
 $Y \sim \text{treat} + \text{race_eth} + \text{hs_degree}$

Call:

```
lm_robust(formula = re78 ~ treat + black + hispan + hs_degree,  
          data = lalonde)
```

Standard error type: HC2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	6633.0	491.2	13.5047	1.527e-36	5668.4	7597.6	609
treat	1003.2	745.2	1.3461	1.788e-01	-460.4	2466.8	609
black	-2223.1	705.1	-3.1531	1.695e-03	-3607.8	-838.5	609
hispan	-162.7	947.0	-0.1718	8.637e-01	-2022.5	1697.2	609
hs_degree	2046.2	650.8	3.1442	1.746e-03	768.1	3324.2	609

Multiple R-squared: 0.03395 , Adjusted R-squared: 0.0276
F-statistic: 5.625 on 4 and 609 DF, p-value: 0.0001887

Is this the
"right" model?
Interactions?

[Differs from stratified est of +\$632. Why?]

Regression (take 2)

Regression with discrete X ?
 $Y \sim \text{treat} + \text{race_eth} + \text{hs_degree}$

Call:

```
lm_robust(formula = re78 ~ treat + black + hispan + hs_degree,  
          data = lalonde)
```

Standard error type: HC2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	CI Lower	CI Upper	DF
(Intercept)	1003.2	491.2	13.5047	1.527e-36	5668.4	7597.6	609
treat	1003.2	745.2	1.3461	1.788e-01	-460.4	2466.8	609
black	-162.7	705.1	-3.1531	1.695e-03	-3607.8	-838.5	609
hispan	-162.7	947.0	-0.1718	8.637e-01	-2022.5	1697.2	609
hs_degree	2046.2	650.8	3.1442	1.746e-03	768.1	3324.2	609

Multiple R-squared: 0.03395 , Adjusted R-squared: 0.0276
F-statistic: 5.625 on 4 and 609 DF, p-value: 0.0001887

Is this the
"right" model?
Interactions?

[Differs from stratified est of +\$632. Why?]

OLS? OLS!

Two perspectives on regression

Regression models make it all too easy to substitute technique for work.... Regression models often seem to be used to compensate for problems in measurement, data collection, and study design. By the time the models are deployed, the scientific position is nearly hopeless.

- Freedman (1991)

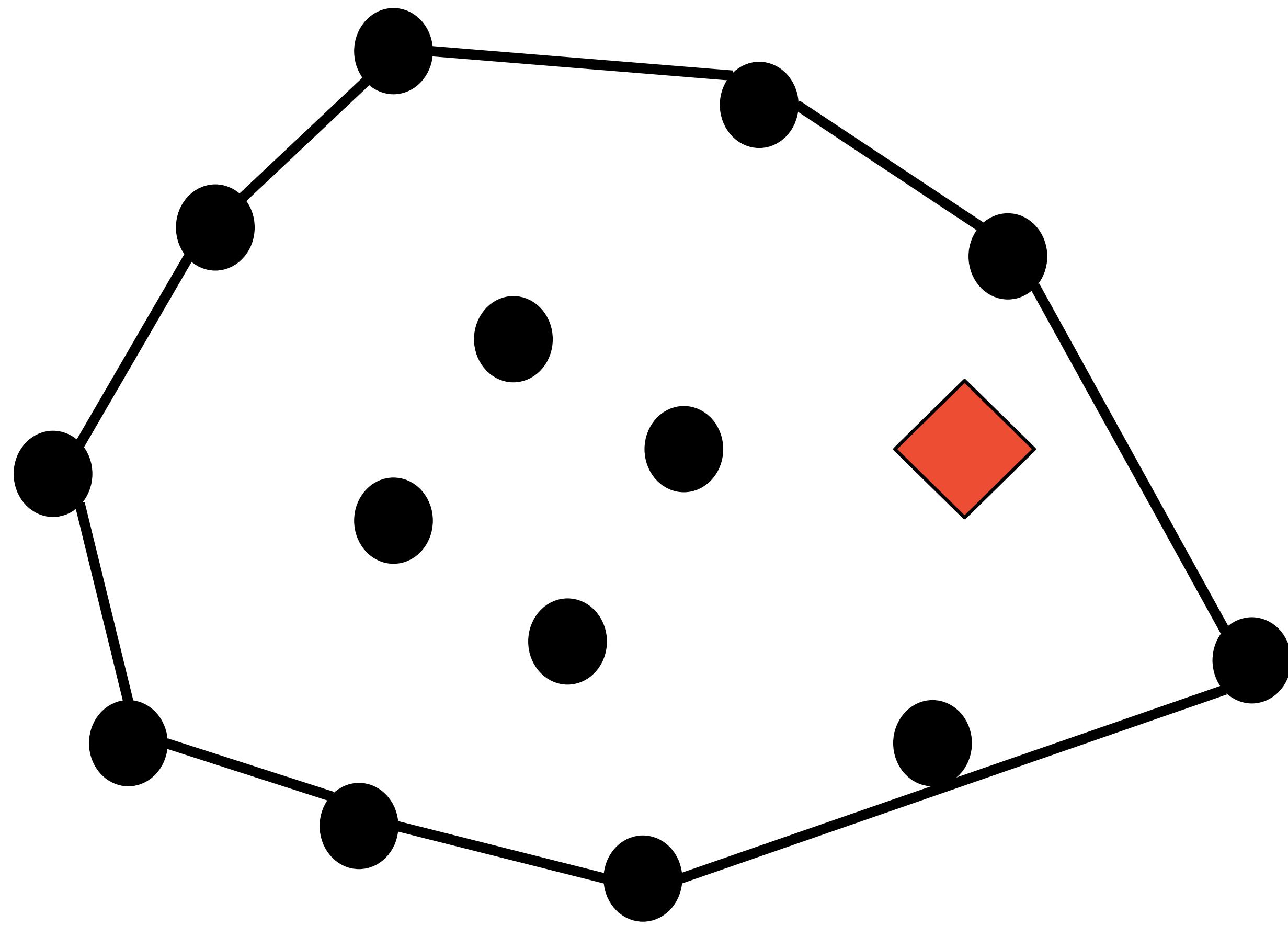
Nevertheless, we believe regression should be the starting point for most empirical projects...

- Angrist & Pischke (2012)

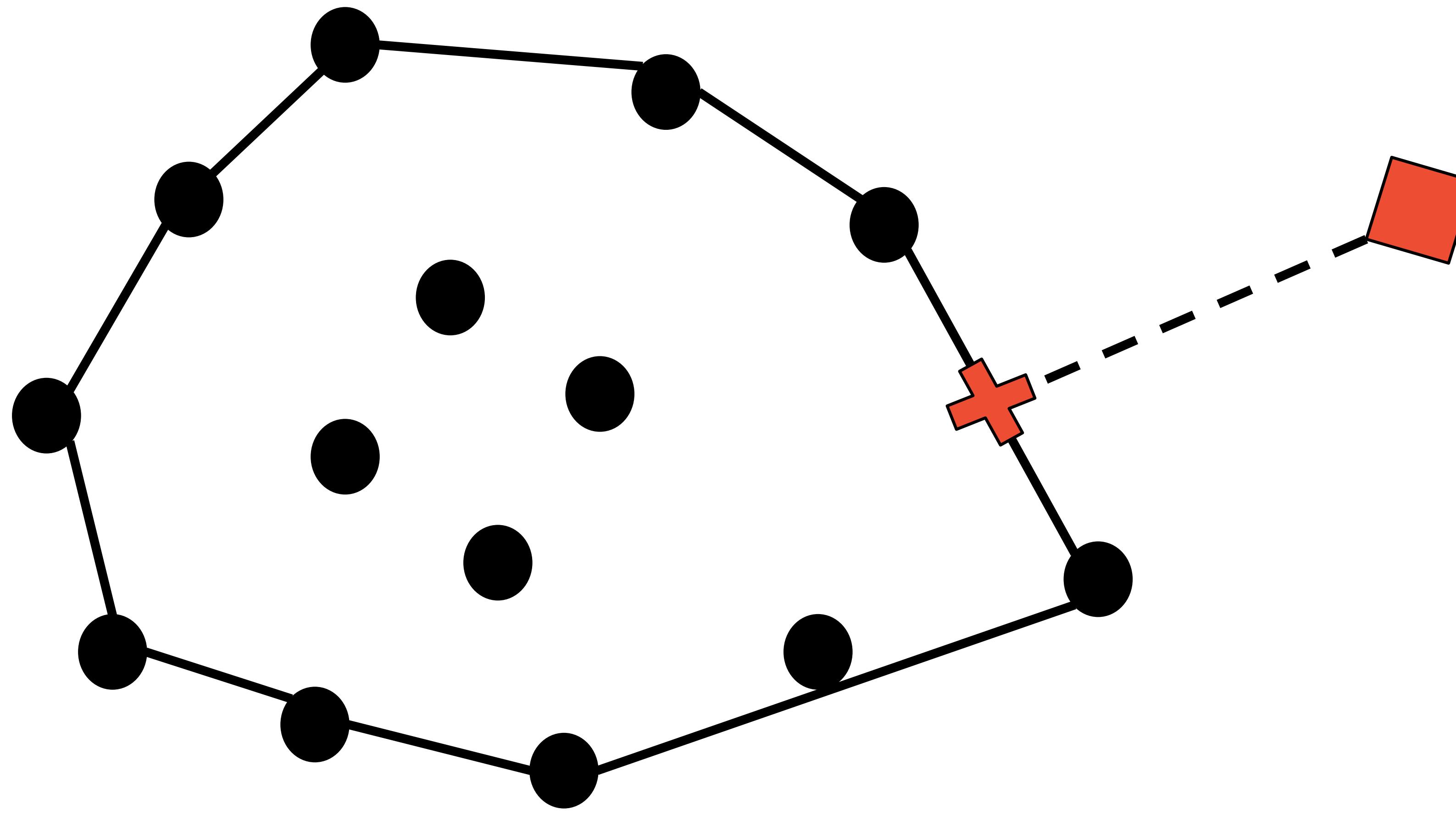
Issue #1: Mis-specification

- Estimate: $Y \sim 1 + D + X_1 + X_2$ but truth is $Y \sim 1 + D + X_1 + X_2 + X_1 \times X_2$
- Classic **omitted variable bias**
 - Remember: $\{Y_i(0), Y_i(1)\} \perp\!\!\!\perp D_i \mid X_i$ doesn't specify functional form
- With binary/discrete X , mis-specification alone usually isn't a huge problem
 - [Maybe ML can help? We'll come back to this]

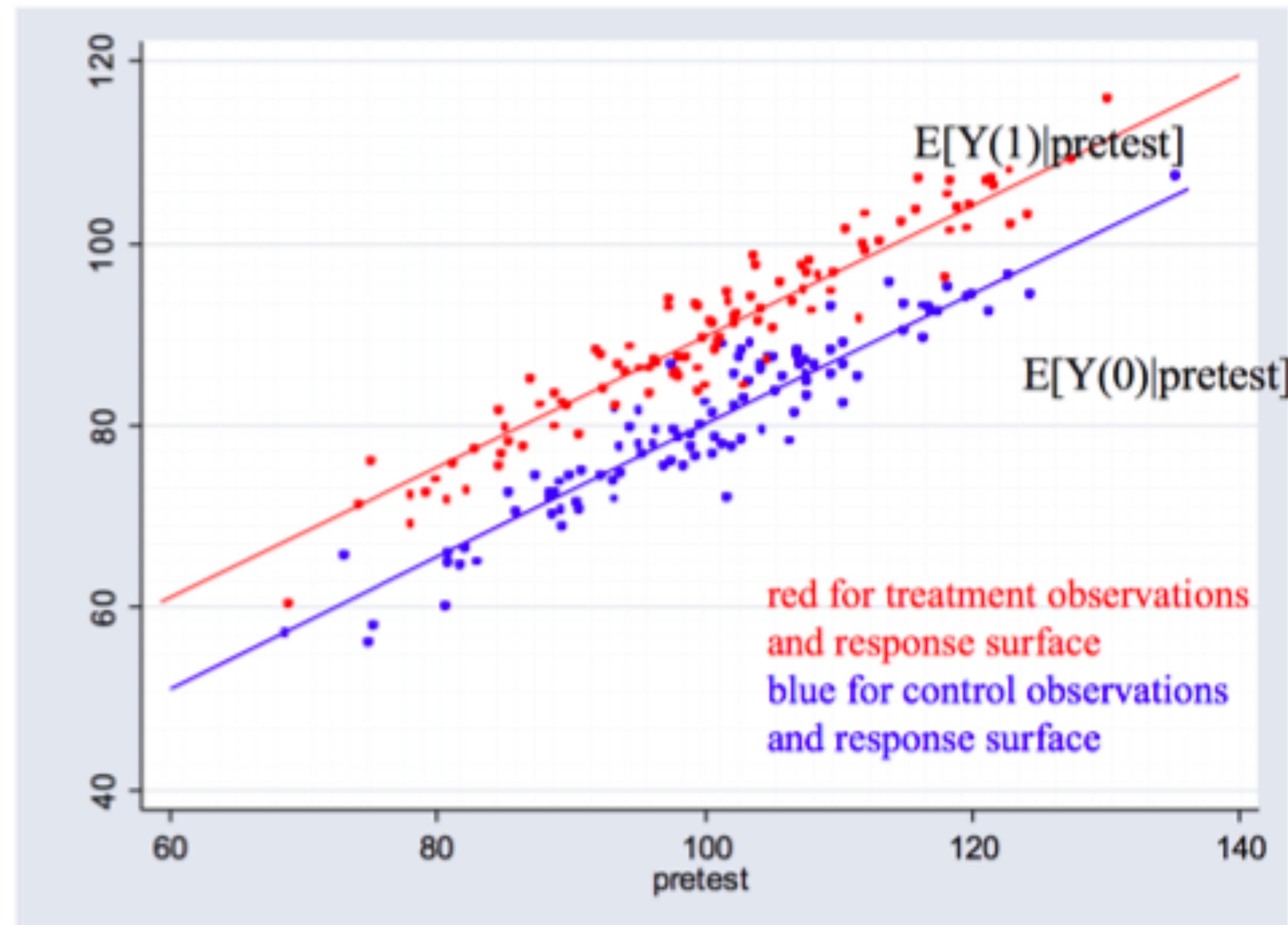
Issue #2: Extrapolation



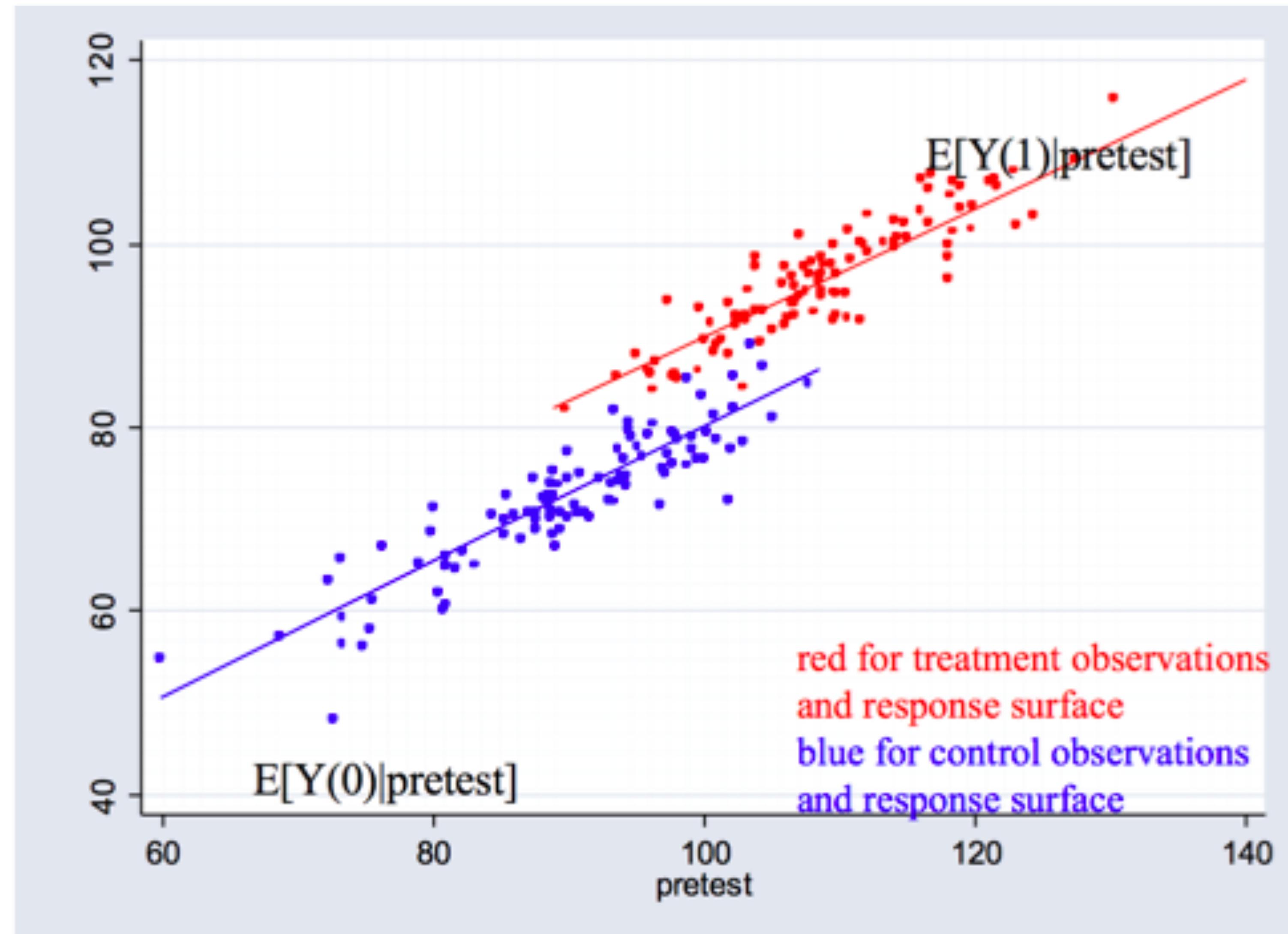
Issue #2: Extrapolation



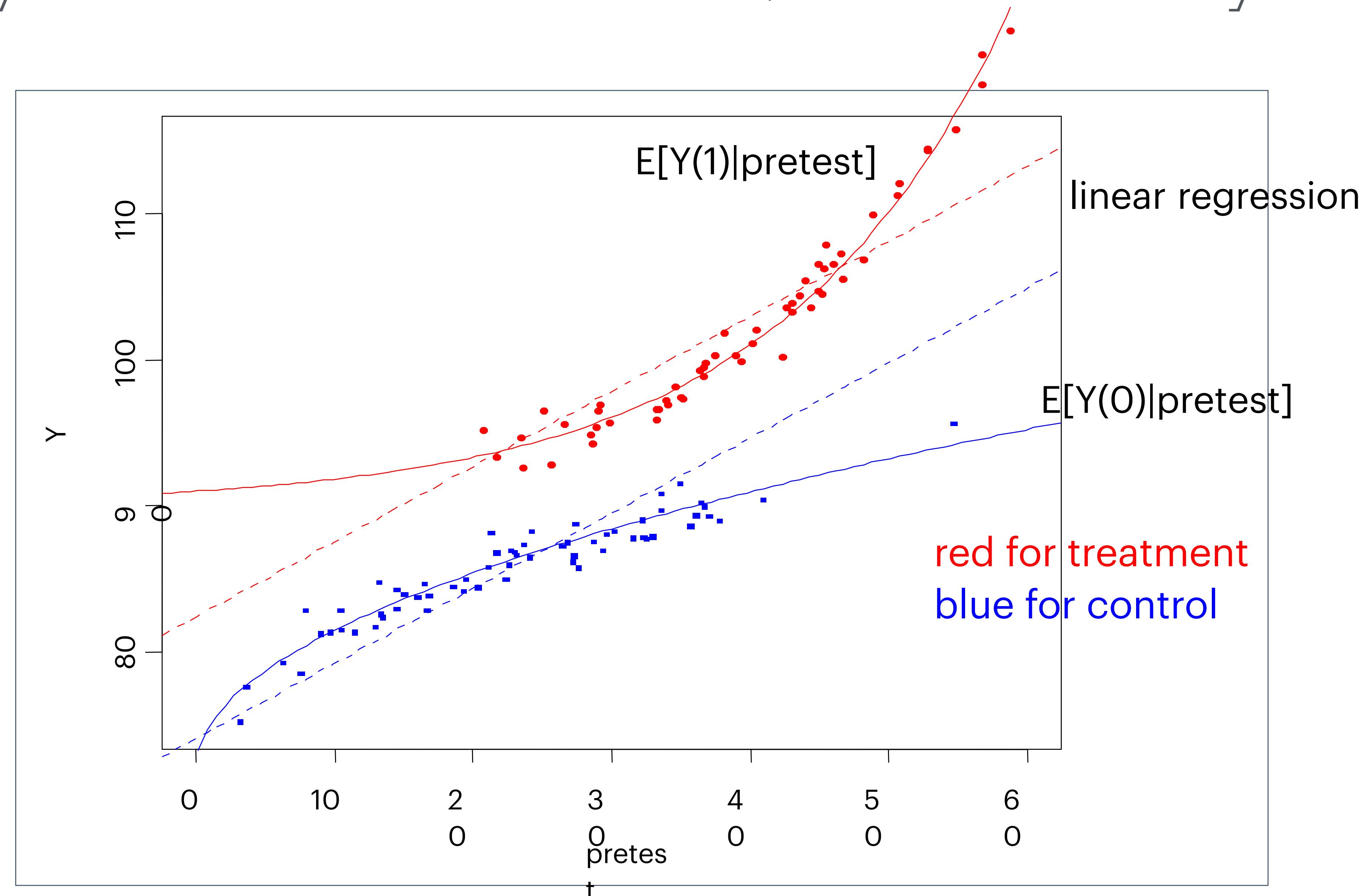
RCT: Regression only for precision



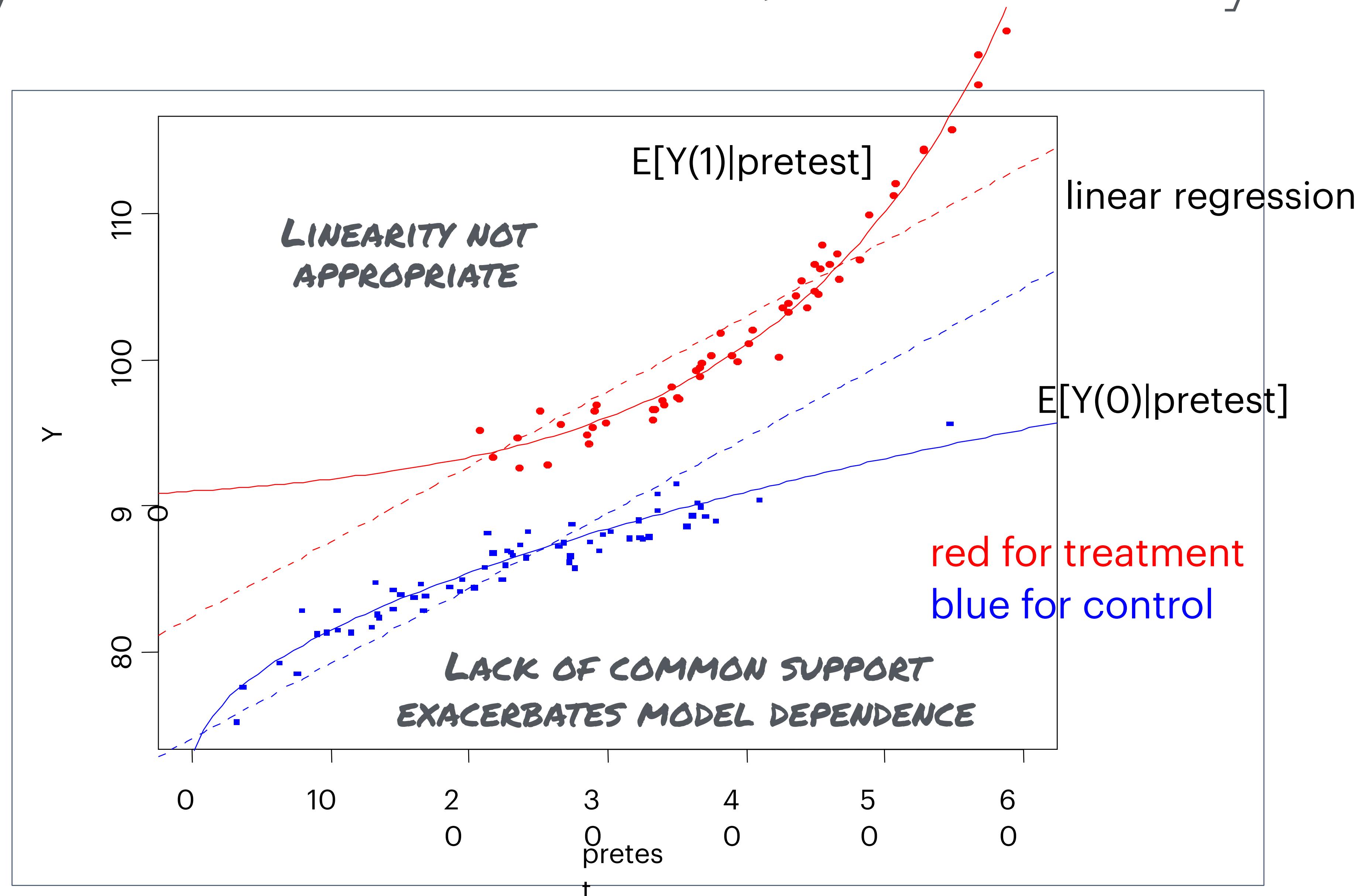
Obs study: One confounder + linearity



Obs study: One confounder, no linearity



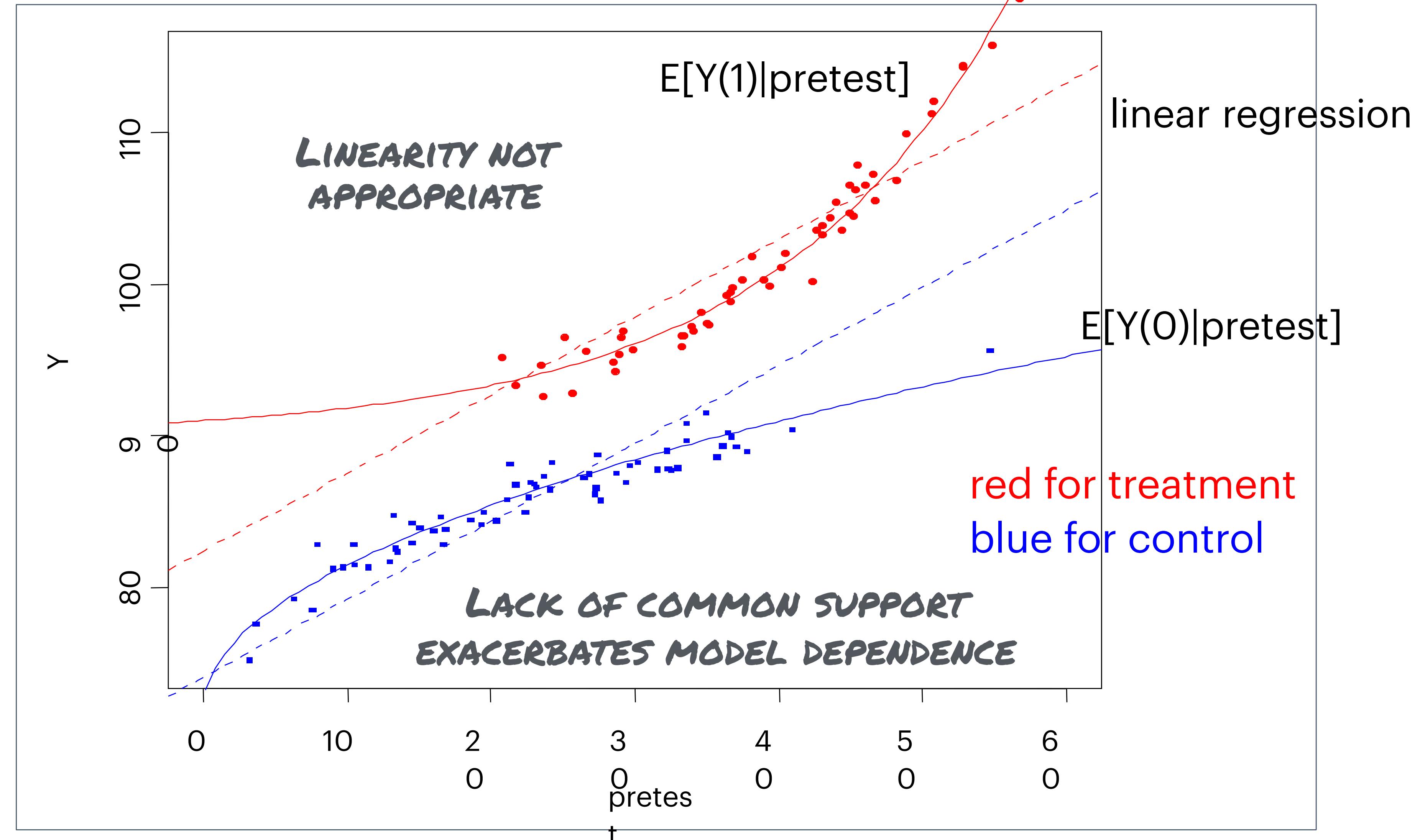
Obs study: One confounder, no linearity



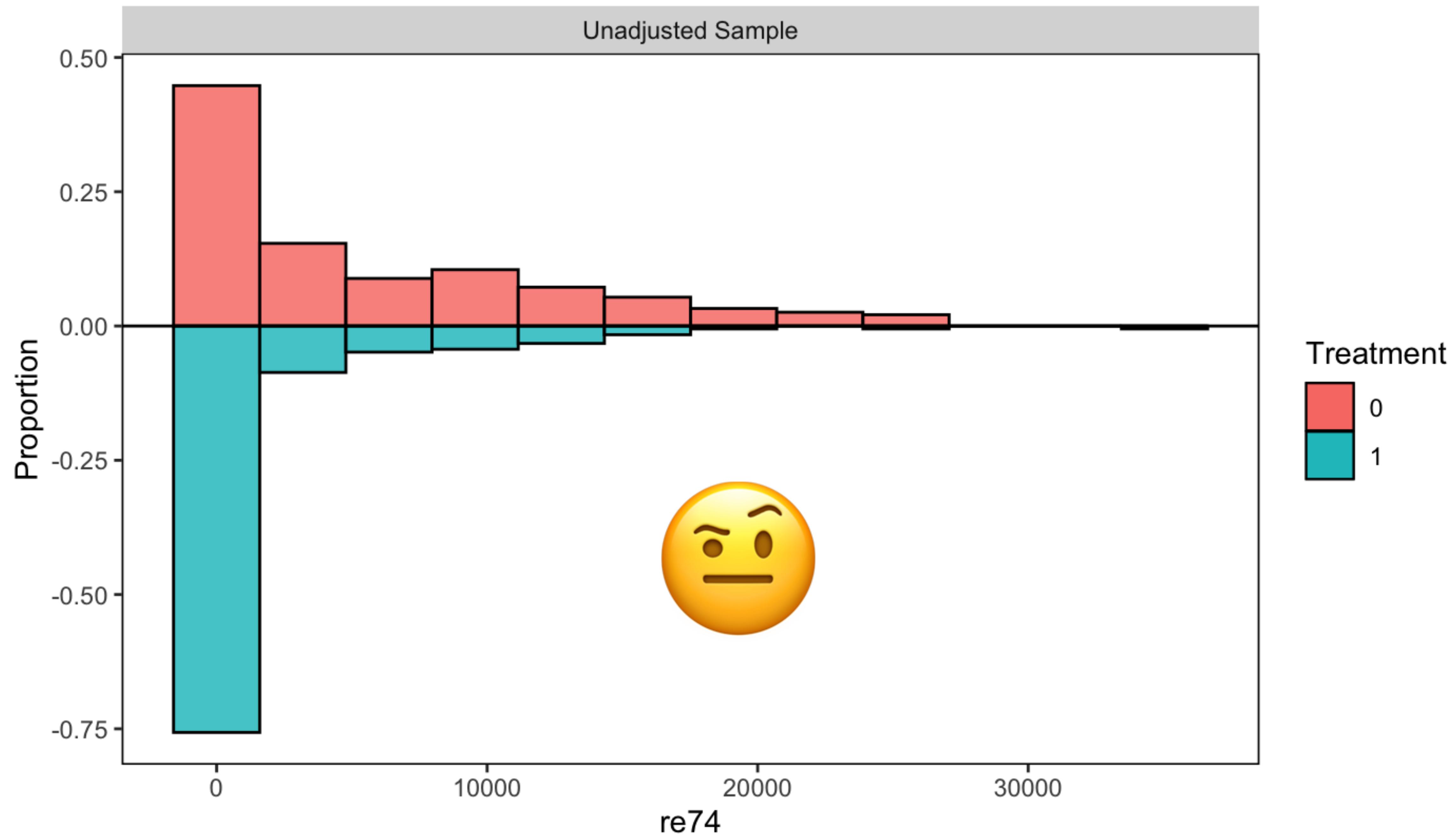
Obs study: One confounder, no linearity



NO RED FLAGS!



Distributional Balance for "re74"

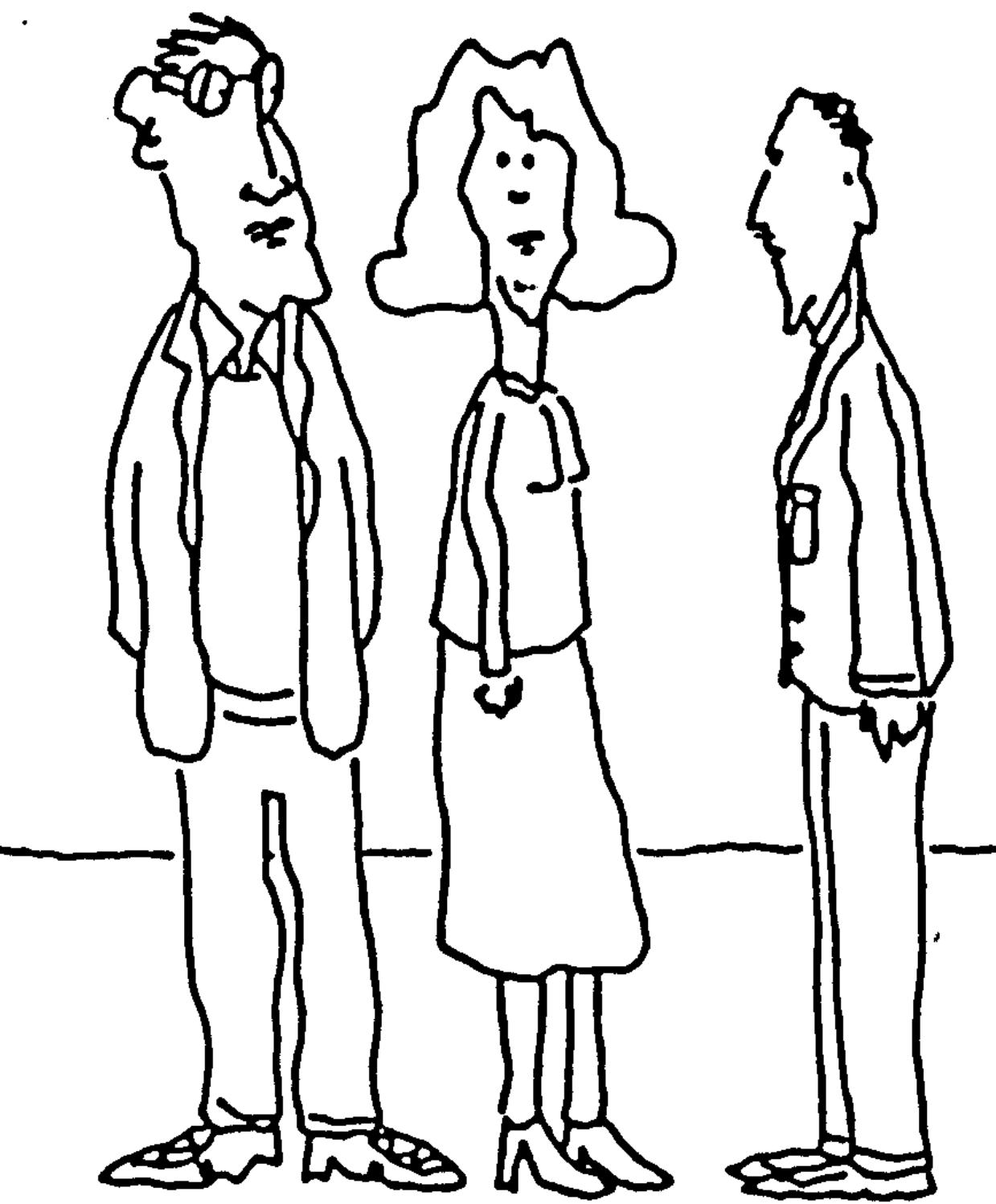


LaLonde: Regression adjustment 😞

Specification	Full Samples		
	NSW	CPS-1	CPS-3
	(1)	(2)	(3)
Raw Difference	1,794 (633)	-8,498 (712)	-635 (657)
Demographic controls	1,670 (639)	-3,437 (710)	771 (837)
1975 Earnings	1,750 (632)	-78 (537)	-91 (641)
Demographics, 1975 Earnings	1,636 (638)	623 (558)	1,010 (822)
Demographics, 1974 and 1975 Earnings	1,676 (639)	794 (548)	1,369 (809)

(From MHE Table 3.3.3. Demographics are age, years of schooling, dummies for Black, Hispanic, high school dropout, and married)

Matching



CONTROL GROUP



OUT OF CONTROL GROUP

“Planners of observational studies should always ask themselves: How would the study be conducted if it were possible to do it by controlled experimentation?”

- Cochran (1965)

design observational study to **approximate experiments**

[i.e., to make them less terrible; aka. “target trials” or “trial emulation”]

Matching ↵ Blocking with $n = 2$

- **Blocking:** find (possibly many) units with same X
 - What if X takes many values? (e.g., age)
 - What if we have many X s?
- **Matching:** for each treated unit, find (at least one) control unit with similar X
 - Approximate randomized trial with many blocks of size $n = 2$
 - Start with **exact matching** if possible; usually have **inexact matching** in practice
 - Don't (necessarily) care about the matches themselves ↵ want good balance overall
- Many choices, no right way to do matching (can be frustrating!)

Warmup: 1-1 matching with replacement

Warmup: 1-1 matching with replacement

- *Repeat:* for every treated unit
 - Find the control unit with the closest value of X
[if there are 2+ good matches, choose one at random]

Warmup: 1-1 matching with replacement

- *Repeat:* for every treated unit
 - Find the control unit with the closest value of X
[if there are 2+ good matches, choose one at random]
- That's it... 

Trainees			Non-Trainees		
unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900
2	34	10200	2	50	31000
3	29	14400	3	30	21000
4	25	20800	4	27	9300
5	29	6100	5	54	41100
6	23	28600	6	48	29800
7	33	21900	7	39	42000
8	27	28800	8	28	8800
9	31	20300	9	24	25500
10	26	28100	10	33	15500
11	25	9400	11	26	400
12	27	14300	12	31	26600
13	29	12500	13	26	16500
14	24	19700	14	34	24200
15	25	10100	15	25	23300
16	43	10700	16	24	9700
17	28	11500	17	29	6200
18	27	10700	18	35	30200
19	28	16300	19	32	17800
Average:	28.5	16426	20	23	9500
			21	32	25900
			Average:	33	20724

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900			
2	34	10200	2	50	31000			
3	29	14400	3	30	21000			
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900			
2	34	10200	2	50	31000			
3	29	14400	3	30	21000			
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900			
2	34	10200	2	50	31000			
3	29	14400	3	30	21000			
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000			
3	29	14400	3	30	21000			
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000			
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300			
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100			
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800			
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000			
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800			
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500			
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500			
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	11,13	26	8450
11	25	9400	11	26	400			
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600			
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500			
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200			
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	9,16	24	17700
15	25	10100	15	25	23300			
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700			
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700	1	43	20900
17	28	11500	17	29	6200			
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

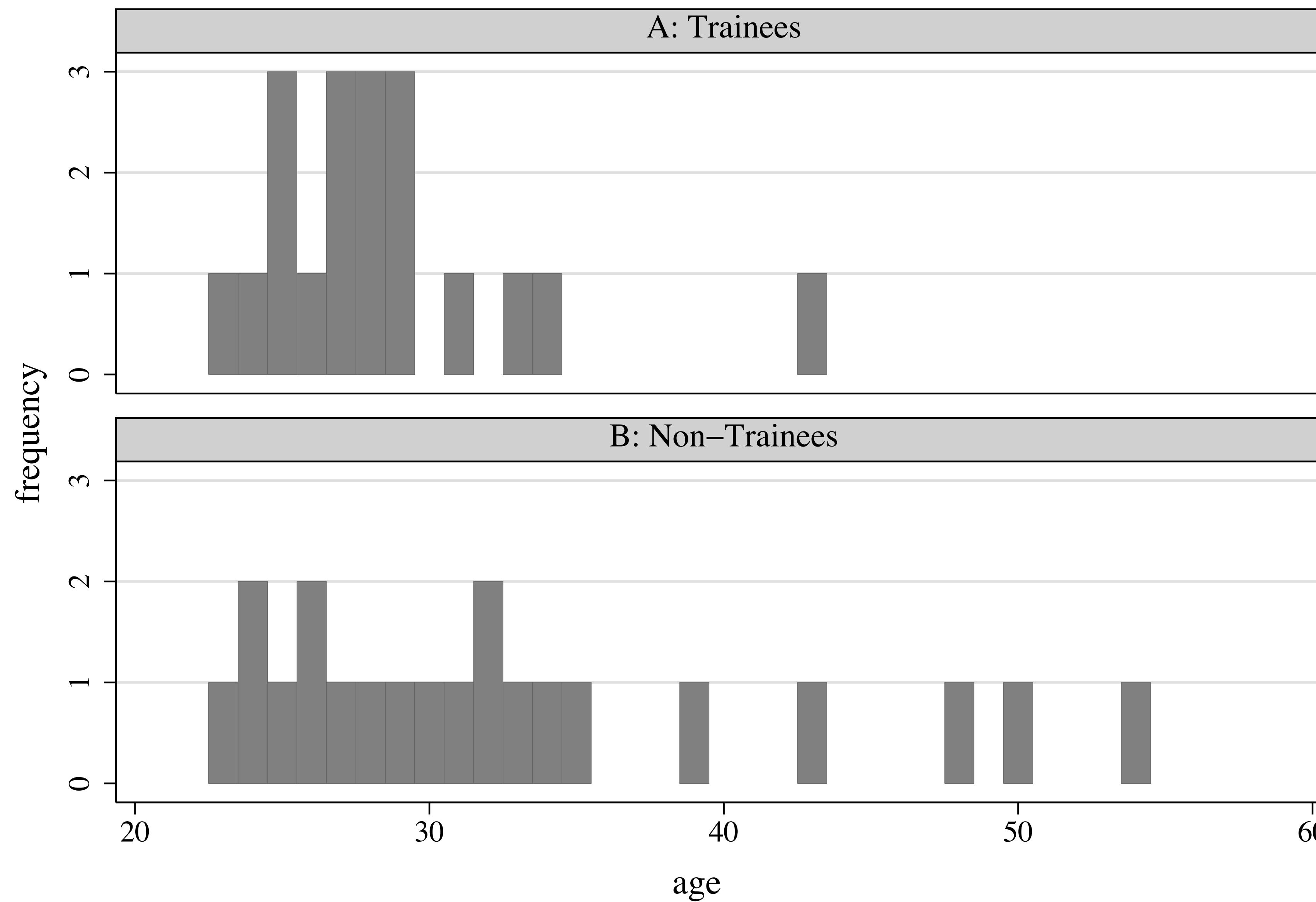
Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700	1	43	20900
17	28	11500	17	29	6200	8	28	8800
18	27	10700	18	35	30200			
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700	1	43	20900
17	28	11500	17	29	6200	8	28	8800
18	27	10700	18	35	30200	4	27	9300
19	28	16300	19	32	17800			
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

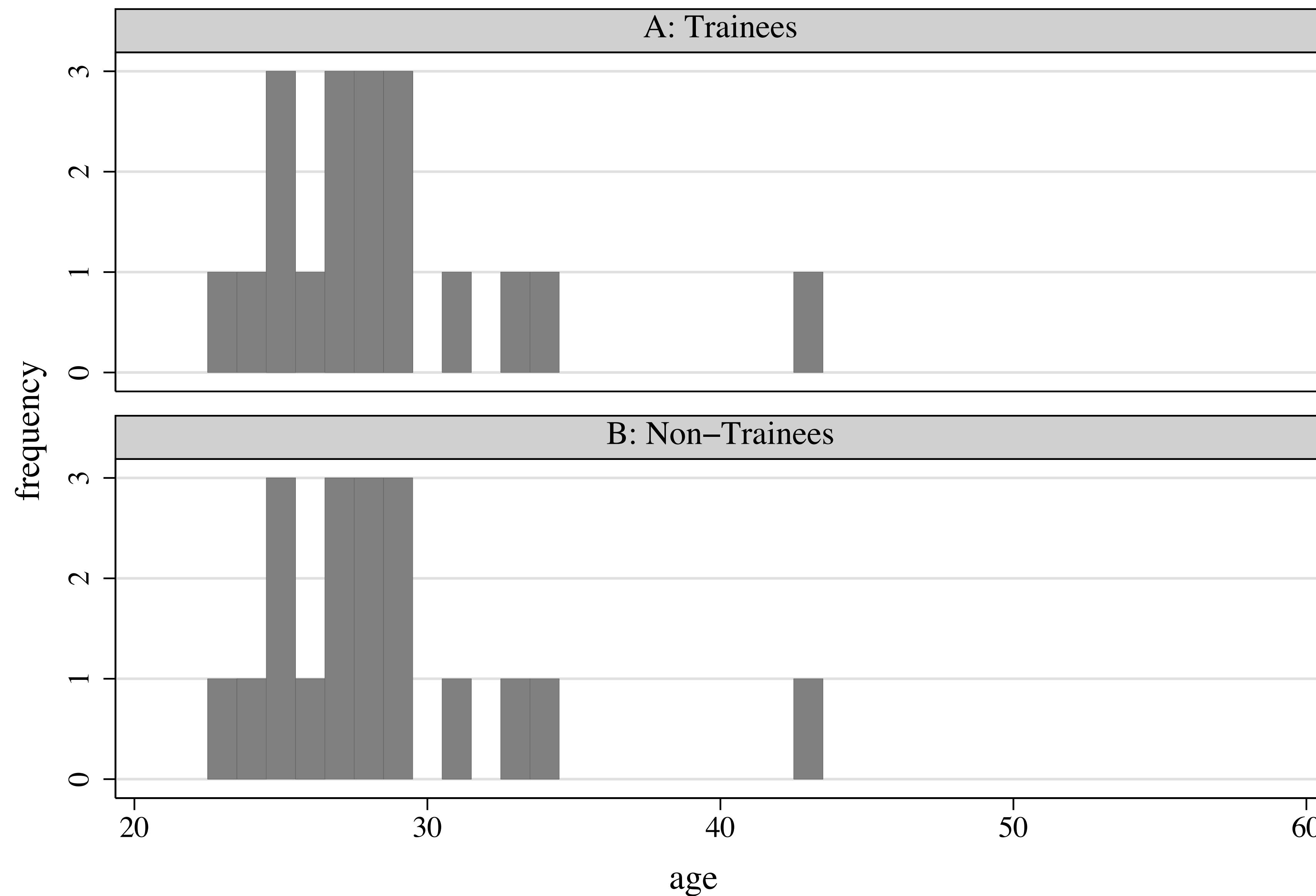
Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700	1	43	20900
17	28	11500	17	29	6200	8	28	8800
18	27	10700	18	35	30200	4	27	9300
19	28	16300	19	32	17800	8	28	8800
Average:	28.5	16426	20	23	9500	Average:		
			21	32	25900			
			Average:	33	20724			

Trainees			Non-Trainees			Matched Sample		
unit	age	earnings	unit	age	earnings	unit	age	earnings
1	28	17700	1	43	20900	8	28	8800
2	34	10200	2	50	31000	14	34	24200
3	29	14400	3	30	21000	17	29	6200
4	25	20800	4	27	9300	15	25	23300
5	29	6100	5	54	41100	17	29	6200
6	23	28600	6	48	29800	20	23	9500
7	33	21900	7	39	42000	10	33	15500
8	27	28800	8	28	8800	4	27	9300
9	31	20300	9	24	25500	12	31	26600
10	26	28100	10	33	15500	13	26	16500
11	25	9400	11	26	400	15	25	23300
12	27	14300	12	31	26600	4	27	9300
13	29	12500	13	26	16500	17	29	6200
14	24	19700	14	34	24200	16	24	9700
15	25	10100	15	25	23300	15	25	23300
16	43	10700	16	24	9700	1	43	20900
17	28	11500	17	29	6200	8	28	8800
18	27	10700	18	35	30200	4	27	9300
19	28	16300	19	32	17800	8	28	8800
Average:	28.5	16426	20	23	9500	Average:	28.5	13982
			21	32	25900	Average:	33	20724

Age Distribution: Before Matching



Age Distribution: After Matching



Treatment effect estimates

- Difference in average earnings between trainees and non-trainees:
- *Before matching*: $\$16,426 - \$20,724 = \textbf{-\$4, 298}$
- *After matching*: $\$16,426 - \$13, 982 = \textbf{-\$2, 444}$
- [use diff-in-means here; could also use regression after matching]

What's the estimand?

- Find matches for each **treated unit**
 - i.e., trying to get X to look like distribution for treated group
- ↗ **Average Treatment Effect on the Treated (ATT)**

$$\begin{aligned} \text{ATT} &= \mathbb{E}[Y(1) - Y(0) \mid D = 1] \\ &= \mathbb{E}[Y(1) \mid D = 1] - \mathbb{E}[Y(0) \mid D = 1] \\ &\quad \bar{Y}_{\text{treated}} \qquad \qquad \qquad \bar{Y}_{\text{matched controls}} \end{aligned}$$

- Still need: no unmeasured confounding given X , overlap

Matching Assumptions

Selection on Observables Assumption:



Prince Charles

Male
Born in 1948
Raised in the UK
Married Twice
Lives in a castle
Wealthy and Famous



Ozzy Osbourne

Male
Born in 1948
Raised in the UK
Married Twice
Lives in a castle
Wealthy and Famous

THANKS, JAKE!

Matching

9

69

Matching complexities

- Example above **1:1 matching with replacement**
- Even here, many decisions to make!
 - What if there are multiple possible matches? [avg them? pick one?]
 - What if there are multiple treated units with the same X ?
- Biggest question: *what if we can't find a match??*

Matching when we can't find an exact match

- What to do when a treated unit doesn't have an exact match on X in the control group?
- **Option 1:** Discard unit
 - Then perfectly matched sample — but **we've changed the estimand!**
 - [More on this when we get to propensity scores]
- **Option 2:** Accept a “close” match
 - Have to decide what is “close enough”? Introduce imbalance in $X \rightsquigarrow$ possible imbalance in $Y(0)$
- With more X s, harder to find an exact match [*curse of dimensionality*]

Distance measures for matching



How do we measure “closeness” between units?

- ▶ **Euclidean distance:**

$$\text{dist}(i,j) = \sqrt{\sum_k (X_{ki} - X_{kj})^2}$$

- ▶ **Normalized Euclidean distance:**

$$\text{dist}(i,j) = \sqrt{\sum_k \frac{(X_{ki} - X_{kj})^2}{\text{var}[X_k]}}$$

- ▶ **Mahalanobis distance:**

$$\text{dist}(i,j) = \sqrt{(\mathbf{X}_i - \mathbf{X}_j)' \widehat{\Sigma}_X^{-1} (\mathbf{X}_i - \mathbf{X}_j)}$$

(Propensity scores coming up...)

Diagnostics: Covariate Balance

- ▶ **Key idea:** A “good match” is a match that yields good covariate balance
- ▶ Important to assess covariate balance for every match
 - Differences in means are most obvious
 - But also interested in variances, overall shape
 - People often use p -values, F -stat, etc.
- ▶ If poor balance, then match again!

(Stylized) Matching Process

A LITTLE CONVOLUTED...

- ▶ Check overall balance and common support.
- ▶ Match.
- ▶ Check overall balance.
- ▶ Re-match with adjustment.
- ▶ Finalize balance.
- ▶ Estimate treatment effect.

TYPICALLY VIA REGRESSION
("BIAS CORRECTION")

Back to LaLonde

We construct a matched sample in which we exactly match on *years of education.* What can we say about covariate imbalance on *age*?

The matched sample will be exactly balanced on age.

The matched sample will be approximately balanced on age.

The matched sample will be heavily imbalanced on age.

We can't say anything about covariate imbalance on age before looking at the data.

We construct a matched sample in which we exactly match on *years of education.* What can we say about covariate imbalance on *age*?

The matched sample will be exactly balanced on age.

0%

The matched sample will be approximately balanced on age.

0%

The matched sample will be heavily imbalanced on age.

0%

We can't say anything about covariate imbalance on age before looking at the data.

0%

We construct a matched sample in which we exactly match on *years of education.* What can we say about covariate imbalance on *age*?

The matched sample will be exactly balanced on age.

0%

The matched sample will be approximately balanced on age.

0%

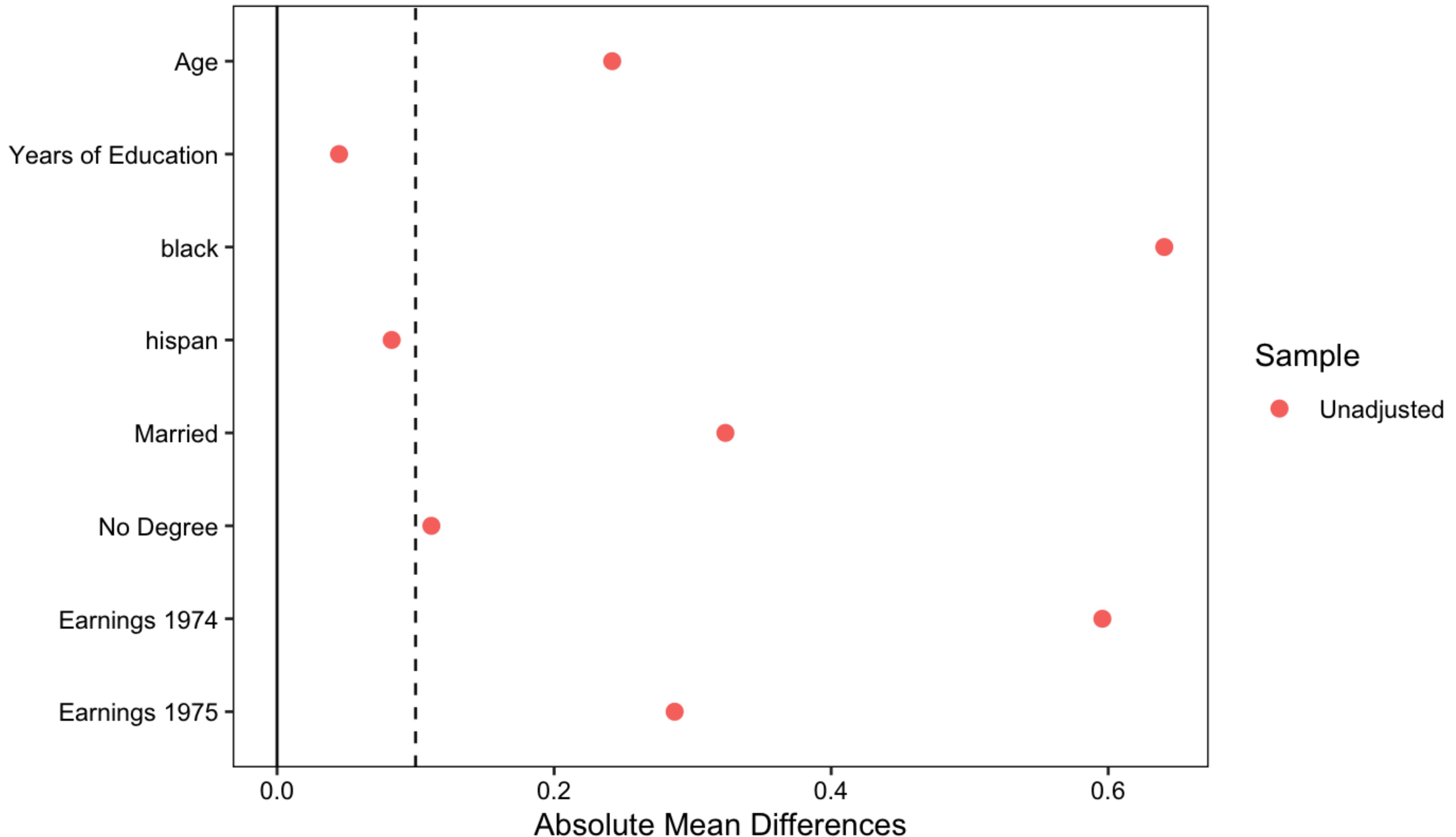
The matched sample will be heavily imbalanced on age.

0%

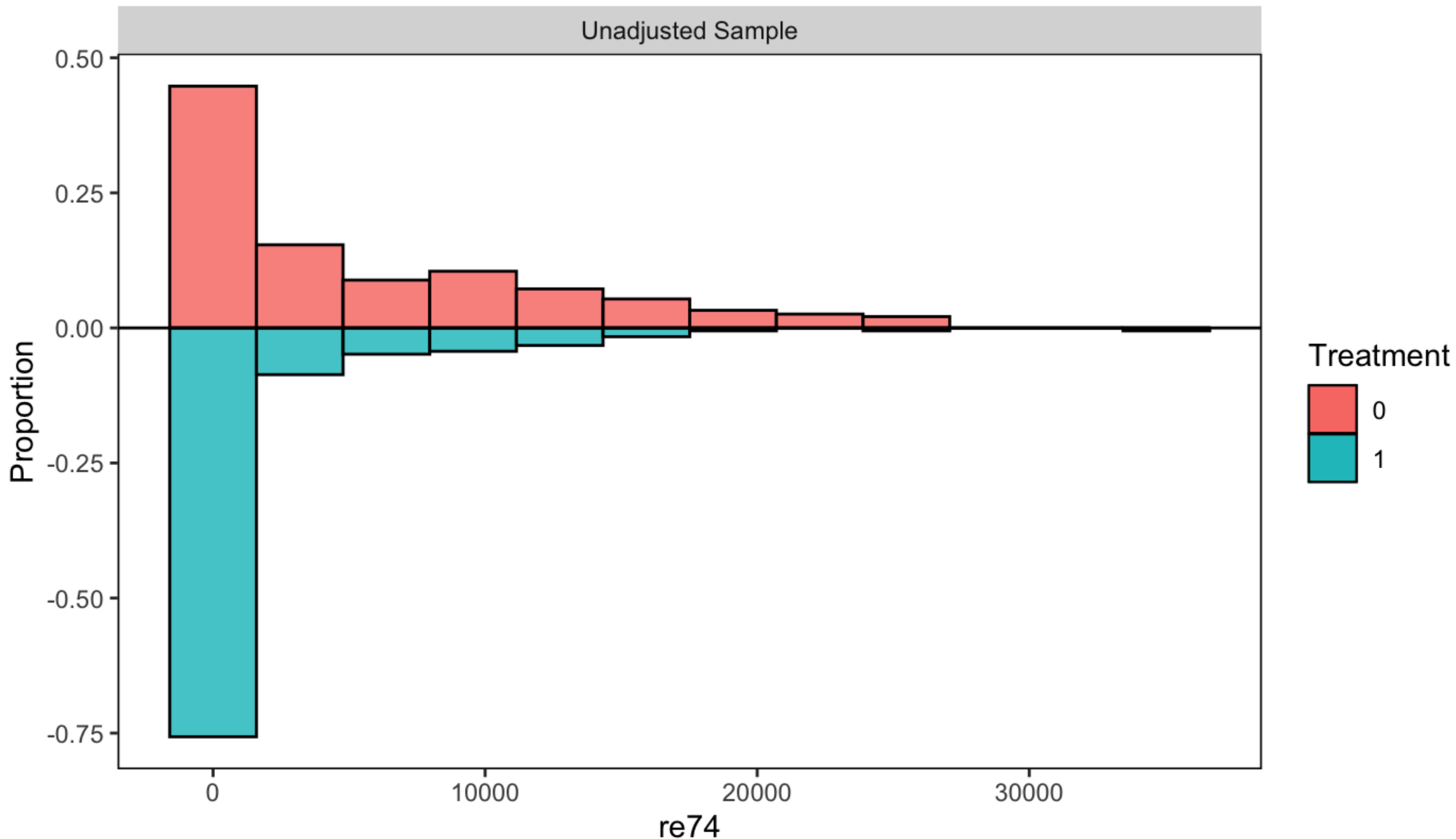
We can't say anything about covariate imbalance on age before looking at the data.

0%

Covariate Balance



Distributional Balance for "re74"



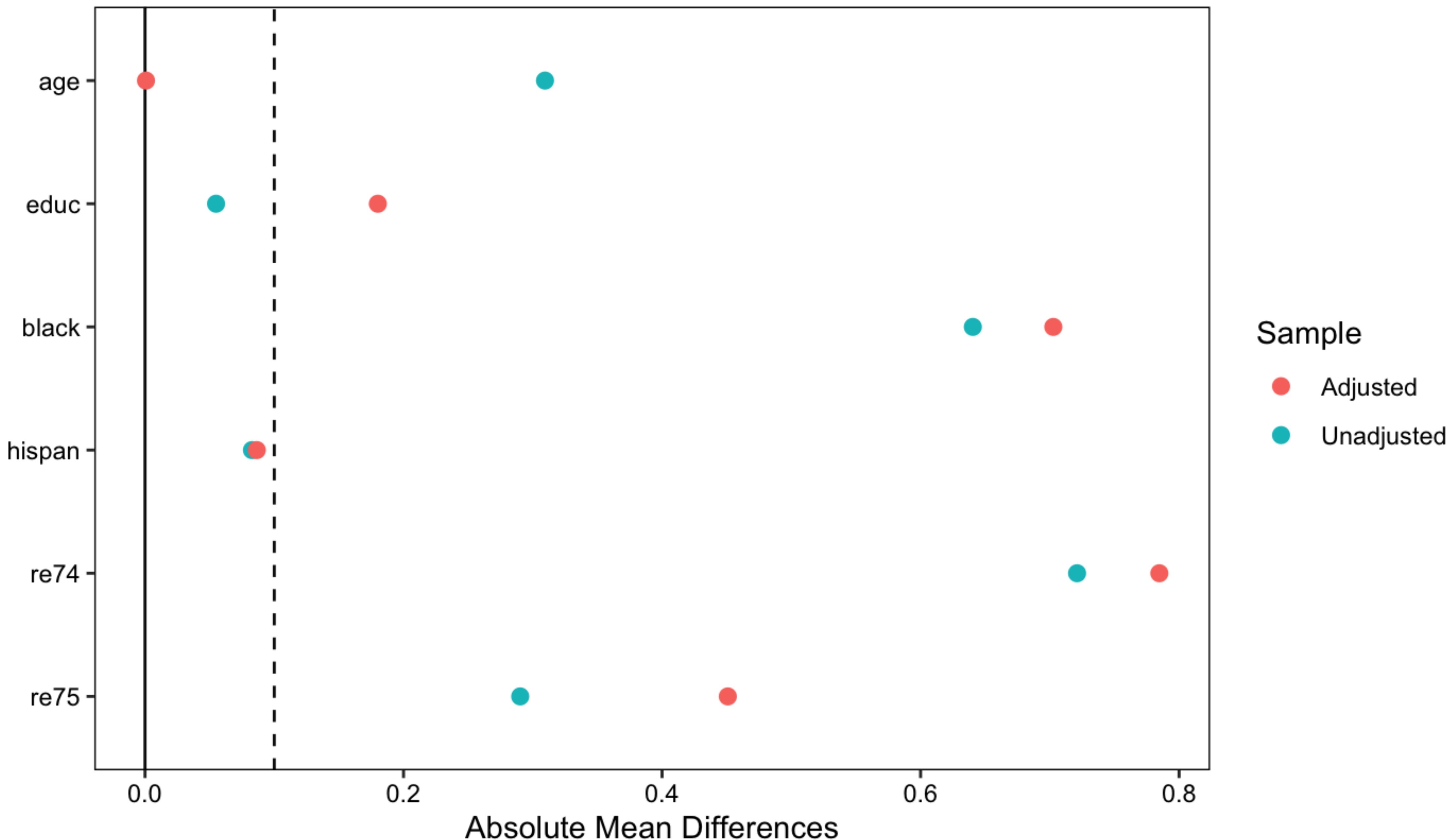
Let's Match!

To start: **1:1 nearest neighbor matching without replacement**

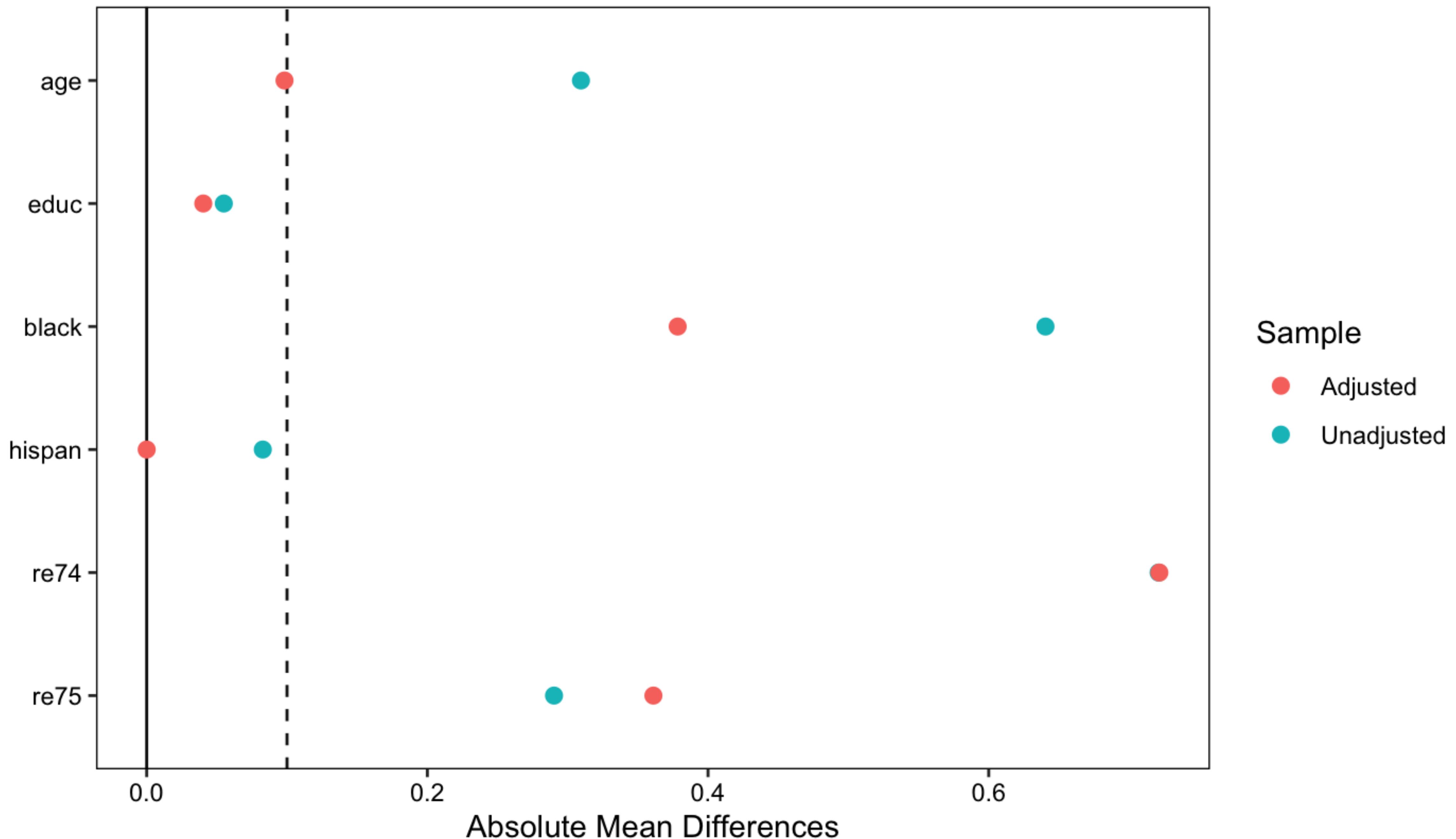
Practical questions:

- What should we use to compute “distance”?
- Should we “exact” match on anything?
(i.e., only match treated units with a HS degree to control units with a HS degree)

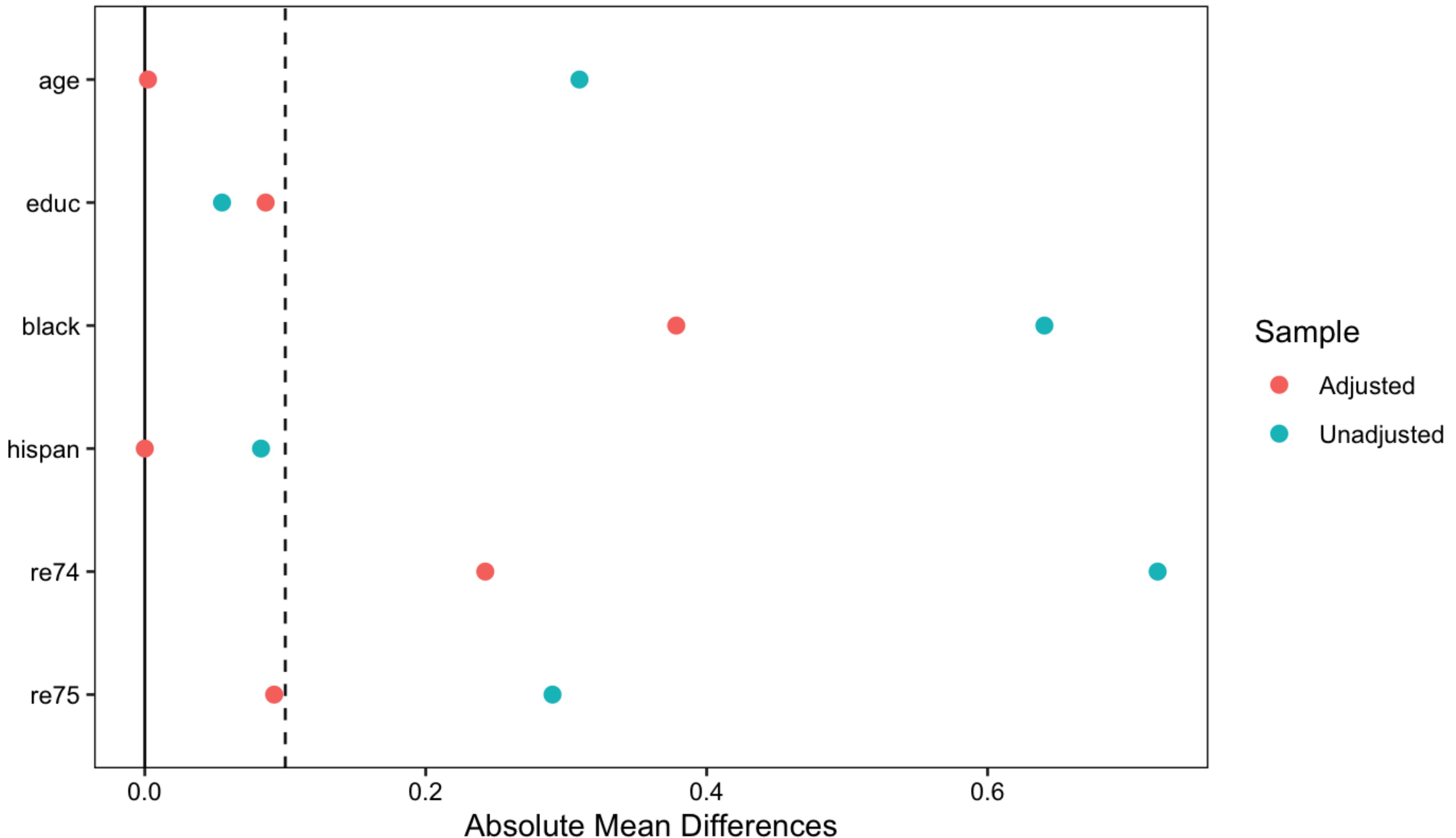
Cov balance matching on: age



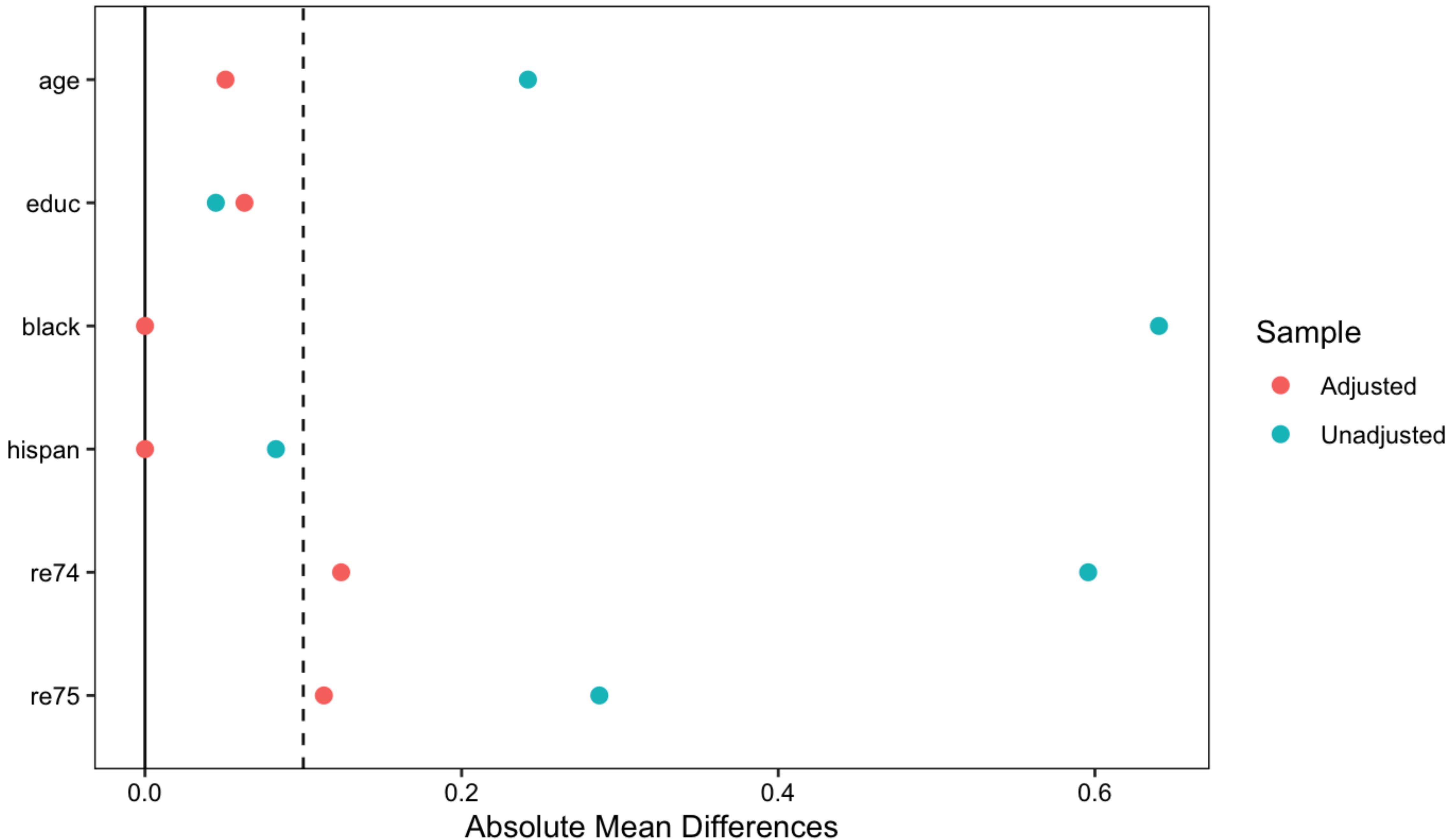
Cov balance matching on: age, educ, race/eth



Cov balance matching on: age, educ, race/eth, prior earnings



Cov balance with exact match on race/eth



Dropping units

When we only do 1:1 nearest neighbor matching:

- Treated units: 185 (out of 185)
- Control units: 185 (out of 429)

When we also exact match on race/ethnicity:

- Treated units: 116 (out of 185)
- Control units: 116 (out of 429)

Why?

	Control	Treat
Non-Black	342	29
Black	87	156

Target trial

Our target trial is a matched pair RCT where we match on:

- Race/ethnicity
- Age
- Education
- Prior earnings

Many other possible target trials

We like our match. Now what?

Option 1: simple difference-in-means on matched sample

→ **+\$997** (*robust se = \$867*)

Option 2: regression-adjusted estimate on matched sample
(aka. “bias correction”)

→ **+\$998** (*robust se = \$874*)

Next up, the propensity score...