# Randomization & matching

**February 26, 2020**

PMAP 8521: Program Evaluation for Public Service
Andrew Young School of Policy Studies
Spring 2020

Fill out your reading report on iCollege!

# Plan for today

The magic of randomization

The "Gold" Standard

Matching

# The magic of randomization

# Why randomize?

## Fundamental problem of causal inference

$$\delta_i = Y_i^1 - Y_i^0$$

**Individual-level effects are impossible to observe**

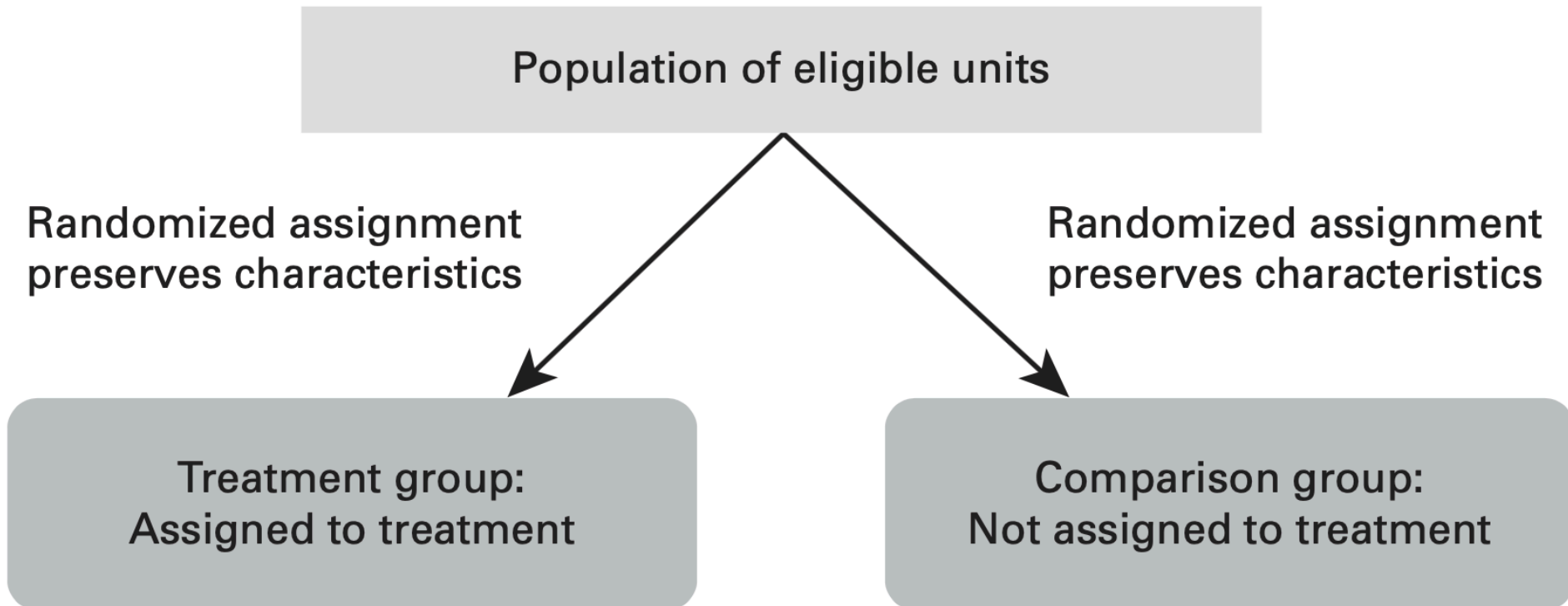# Why randomize?

$$\delta = (\bar{Y}|P = 1) - (\bar{Y}|P = 0)$$

**This only works if subgroups that received/didn't receive treatment look the same**
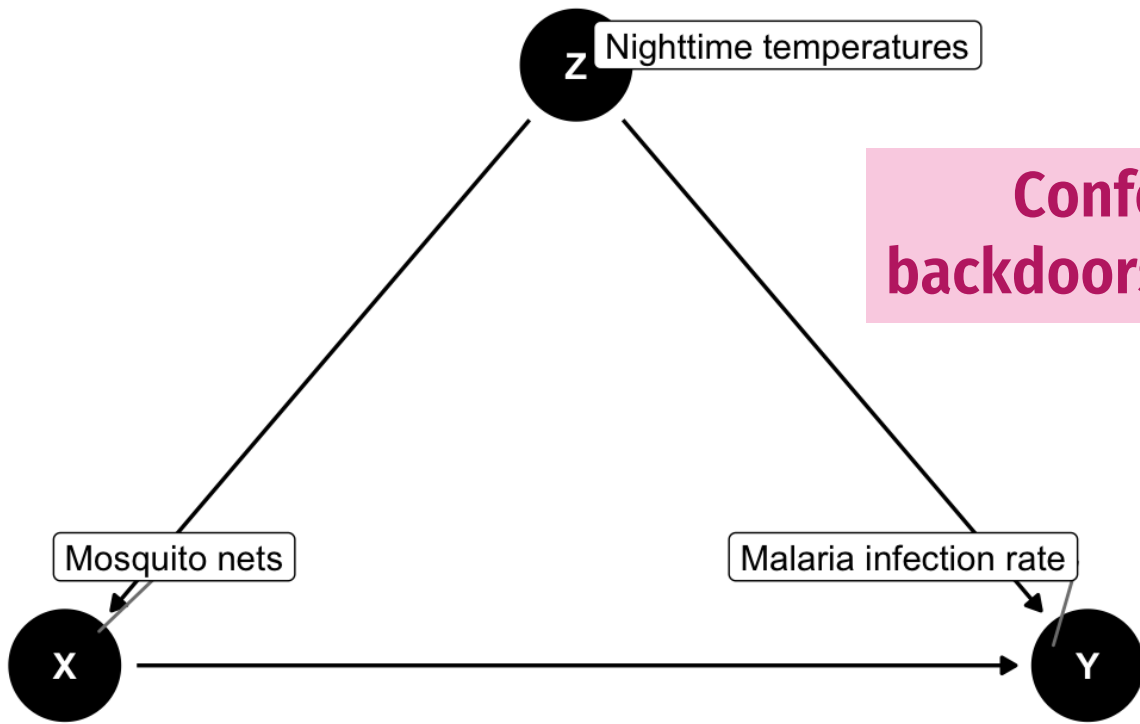
# Why randomize?

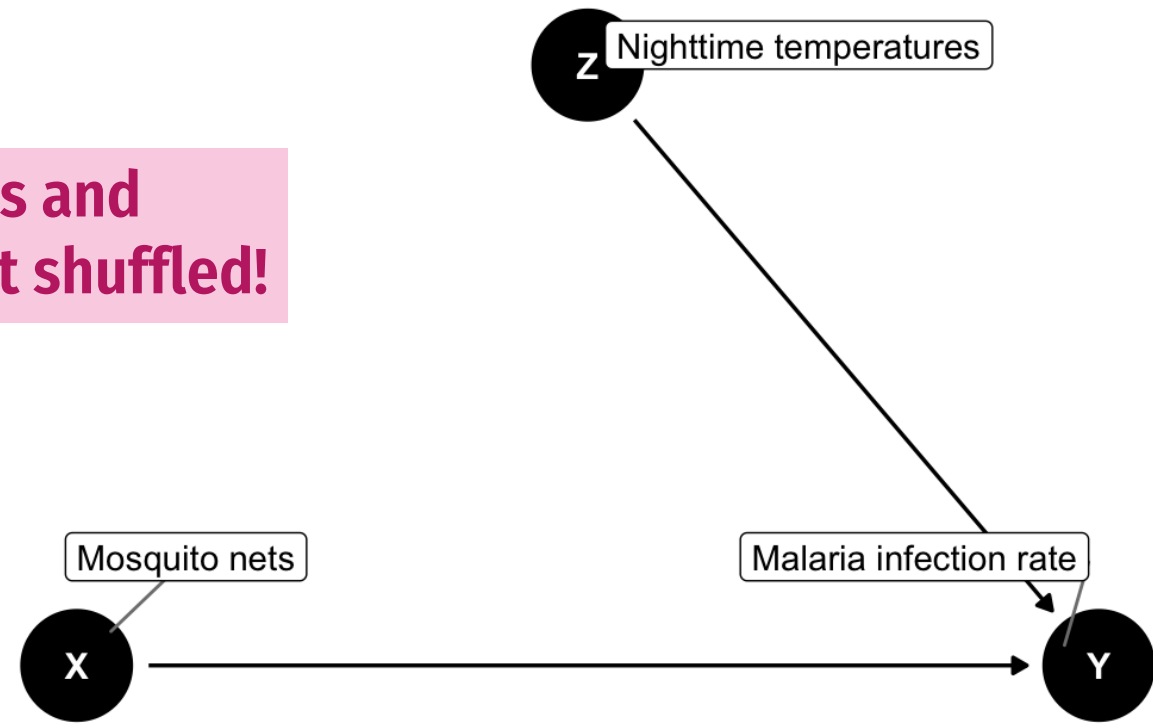With big enough numbers, the magic of randomization helps make comparison groups comparable

Population of eligible units

Randomized assignment preserves characteristics

Randomized assignment preserves characteristics

Treatment group: Assigned to treatment

Comparison group: Not assigned to treatment

# RCTs and DAGs

P(Malaria infection rate | *do*(Mosquito net))

When you *do*() X, remove all arrows into it

Confounders and backdoors all get shuffled!



Observational

Experimental

# How to randomize?



1. Define eligible units

2. Select the evaluation sample

3. Randomize assignment to treatment

Comparison

Treatment

External validity

Internal validity

Ineligible

Eligible

# Random assignment

Coins

Dice

Unbiased lottery

Random numbers + threshold

Atmospheric noise

# How big of a sample?

**R example**

# The "Gold" Standard

# Types of research

Experimental studies vs. observational studies

Which is better?

# How the Illinois Wellness Program Affected …



Randomized controlled trial · Observational study

Participation in running events

Number of gym visits

Estimate

Ends employment

Hospital spending

Total medical spending

Half as much · No effect · Twice as much

Source: What Do Workplace Wellness Programs Do? Evidence from the Illinois Workplace Wellness Study

Google

rct "gold standard"

About 636,000 results (0.67 seconds)

## Randomised controlled trials—the gold standard for effectiveness research

Eduardo Hariton, MD, MBA[1] and Joseph J. Locascio, PhD[2]

‣ Author information ‣ Copyright and License information Disclaimer

## Randomized Assignment of Treatment

When a program is assigned at random—that is, using a lottery—over a large eligible population, we can generate a robust estimate of the counterfactual. *Randomized assignment* of treatment is considered the gold standard of impact evaluation. It uses a random process, or chance, to decide who is granted access to the program and who is not.[1] Under randomized assignment, every eligible unit (for example, an individual, household, business.

**RCTs are great!**

**Super impractical to do all the time though!**

**Business**

# 3 share Nobel Prize in economics for 'experimental approach' to solving poverty

Esther Duflo, who at 46 is the award's youngest winner, shares the hor... fellow MIT economist Abhijit Banerjee and Harvard's Michael Kremer

Pioneers in fight against poverty win 2019 Nobel economics prize
THE PRIZE IN ECONOMIC SCIENCES 2019
KUNGL. VETENSKAPS AKADEMIEN

J-PAL
ABDUL LATIF JAMEEL POVERTY ACTION LAB

**Massachusetts Institute of Technology (MIT)** ✔ @MIT · 5h

Professors Esther Duflo and Abhijit Banerjee, co-directors of MIT's @JPAL, receive congratulations on the big news this morning. They share in the #NobelPrize in economic sciences "for their experimental approach to alleviating global poverty."

Photo: Bryce Vickmark

💬 12     ↻ 112     ♥ 510

**Grad School Imposter** @darinself · 6h

Siri, can you sum up the issues of gender and Economics in one headline??

> **Rohini Mohan** ✔ @rohini_mohan · 7h
>
> Oh COME ON @EconomicTimes!
>
> Business News › News › Politics and Nation › Indian-American MIT Prof Abhijit Banerjee and wife wins Nobel in Economics
>
> | Benchmarks › | NSE Loser-Large Cap › | SPONSORED FU |
> |---|---|---|
> | Sensex ● CLOSED | Infosys | Axis Long Term Plan-Growth |
> | 38,214.47 ↑ 87.39 | 786.10 ↓ -28.70 | ★★★★★ |
>
> # Indian-American MIT Prof Abhijit Banerjee and wife wins Nobel in Economics
>
> *Banerjee, born in 1961 in Mumbai, bagged the award for his "experimental approach to alleviating global poverty".*
>
> PTI | Updated: Oct 14, 2019, 04.18 PM IST
>
> Save
>
> f  🐦  in  ◐  MORE  💬  1 Comments
>
> A+  🖨  ✉  🔖
>
> BCCL
>
> STOCKHOLM: Indian-American Abhijit

# "Gold standard"

## "Gold standard" implies that all causal inferences will be valid if you do the experiment right

We don't care if studies are experimental or not

We care if our causal inferences are valid

RCTs are a helpful baseline/rubric for other methods

# RCTs and validity

Randomization fixes a ton of internal validity issues

## Selection
Treatment and control groups are comparable; people don't self-select

## Trends
Maturation, secular trends, seasonality, regression to the mean all generally average out

# RCTs and validity

RCTs don't fix attrition!

Worst threat to internal validity in RCTs

If attrition is correlated with treatment, that's bad

People might drop out because of the treatment, or because they got/didn't get the control group

# Addressing attrition

## Recruit as effectively as possible

You don't just want weird/WEIRD participants

## Get people on board

Get participants invested in the experiment

## Collect as much baseline information as possible

Check for randomization of attrition

# RCTs and validity

## Randomization failures

### Check baseline pre-data

## Noncompliance

Some people assigned to treatment won't take it; some people assigned to control will take it

Intent-to-treat (ITT) vs. Treatment-on-the treated (TTE)

# Other limitations

RCTs don't magically fix construct validity and statistical conclusion validity

RCTs definitely don't magically fix external validity

# The Nobel Prize in economics goes to three groundbreaking antipoverty researchers

In the last 20 years, development economics has been transformed. These researchers are the reason why.

By Kelsey Piper | Oct 14, 2019, 3:30pm EDT

## Empiricism and development economics

The transformation of development economics into an intensely empirical field that leans heavily on randomized controlled trials hasn't been uncontroversial, and many of **the responses** to the Nobel Prize announcement acknowledge that controversy.

Critics have **complained that** randomization feels much more scientific than other approaches but doesn't necessarily answer our questions any more definitively. **Others worry** that the focus on small-scale questions — Do wristbands increase vaccination rates? Do textbooks improve school performance? — might distract us from addressing larger, structural contributors to poverty.

# When to randomly assign

Demand for treatment exceeds supply

Treatment will be phased in over time

Treatment is in equipoise

Local culture open to randomization

When you're a nondemocratic monopolist

When people won't know (and it's ethical!)

When lotteries are going to happen anyway

# When to not randomly assign

When you need immediate results

When it's unethical or illegal

When it's something that happened in the past

When it involves universal ongoing phenomena

# Matching

|  |  | Private | | | Public | | |  |
| Applicant group | Student | Ivy | Leafy | Smart | All State | Tall State | Altered State | 1996 earnings |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| A | 1 |  | Reject | Admit |  | Admit |  | 110,000 |
|  | 2 |  | Reject | Admit |  | Admit |  | 100,000 |
|  | 3 |  | Reject | Admit |  | Admit |  | 110,000 |
| B | 4 | Admit |  |  | Admit |  | Admit | 60,000 |
|  | 5 | Admit |  |  | Admit |  | Admit | 30,000 |
| C | 6 |  | Admit |  |  |  |  | 115,000 |
|  | 7 |  | Admit |  |  |  |  | 75,000 |
| D | 8 | Reject |  |  | Admit | Admit |  | 90,000 |
|  | 9 | Reject |  |  | Admit | Admit |  | 60,000 |

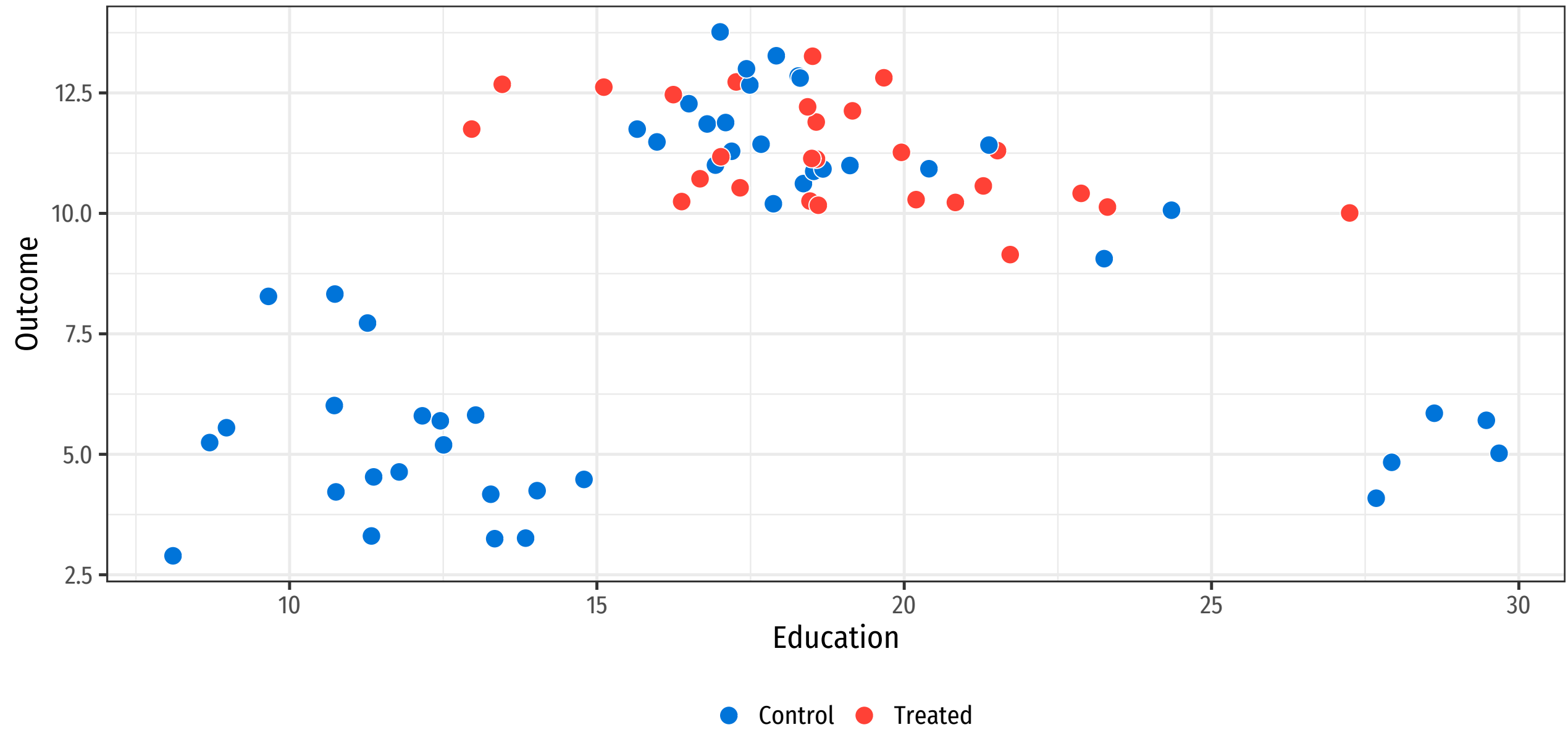*Note:* Enrollment decisions are highlighted in gray.

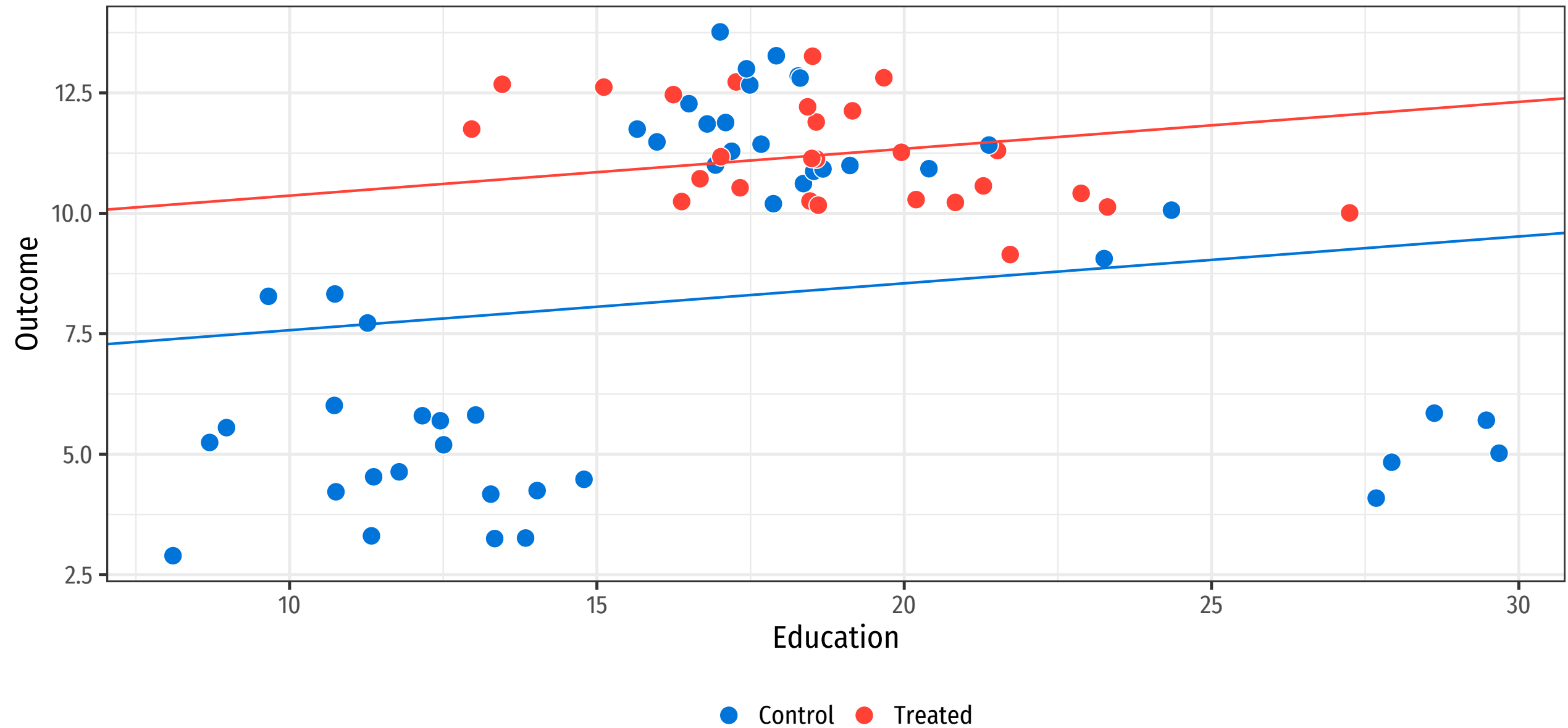# Why match?

Reduce model dependence

Imbalance → model dependence → researcher discretion → bias
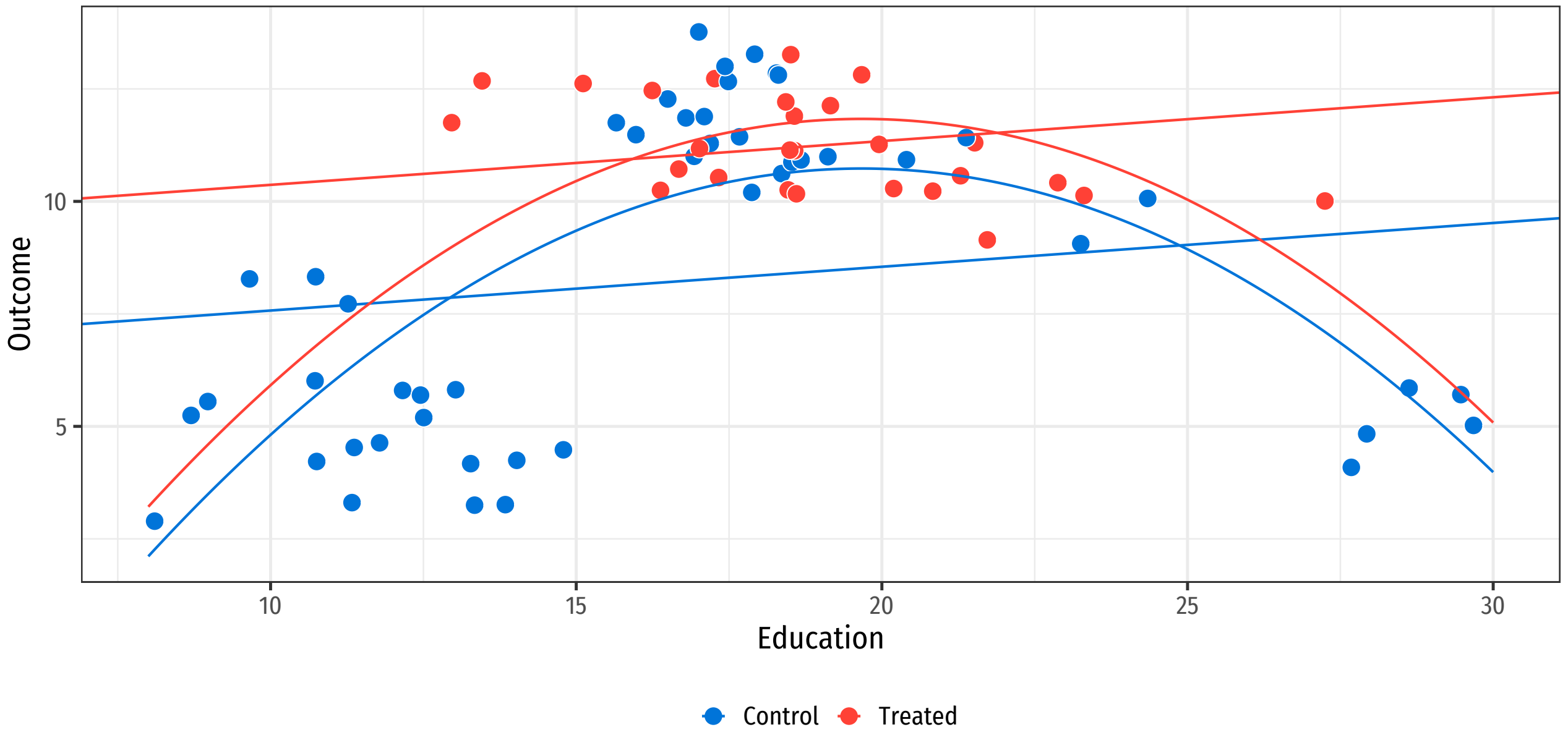
Compare apples to apples
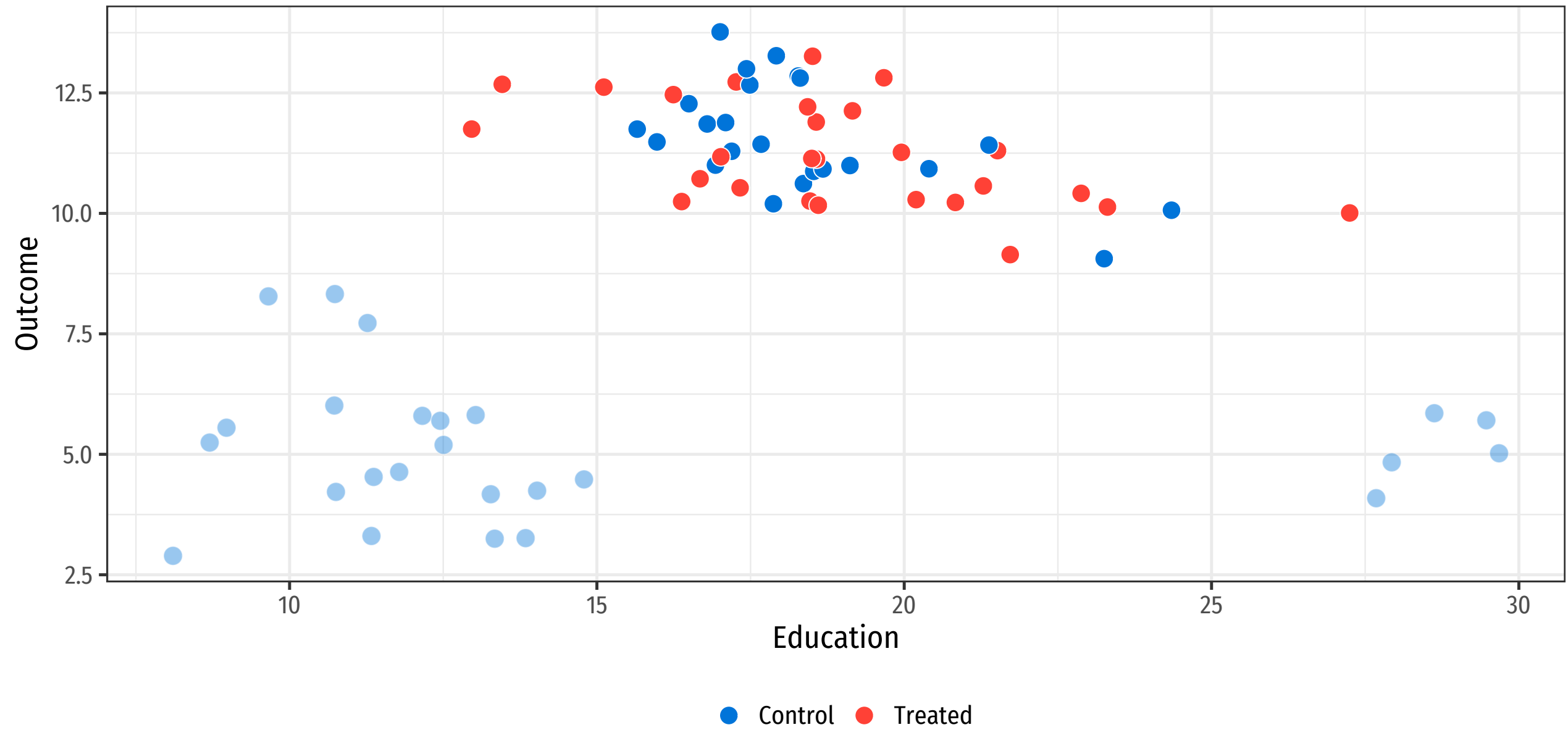
It's a way to adjust for backdoors!

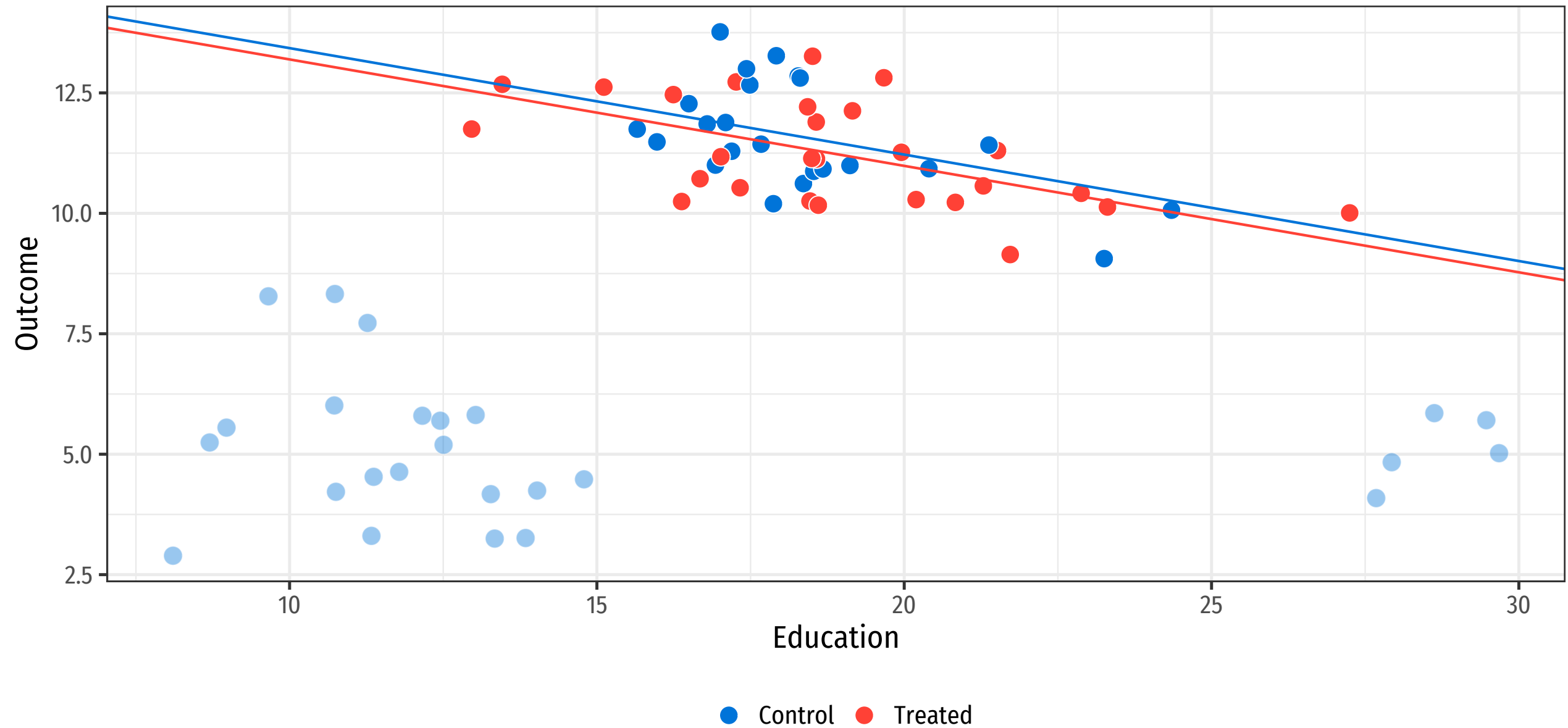$$\text{Outcome} = \beta_0 + \beta_1 \text{Education} + \beta_2 \text{Treatment}$$

$$\text{Outcome} = \beta_0 + \beta_1\text{Education} + \beta_2\text{Education}^2 + \beta_3\text{Treatment}$$

Outcome $= \beta_0 + \beta_1 \text{Education} + \beta_2 \text{Treatment}$

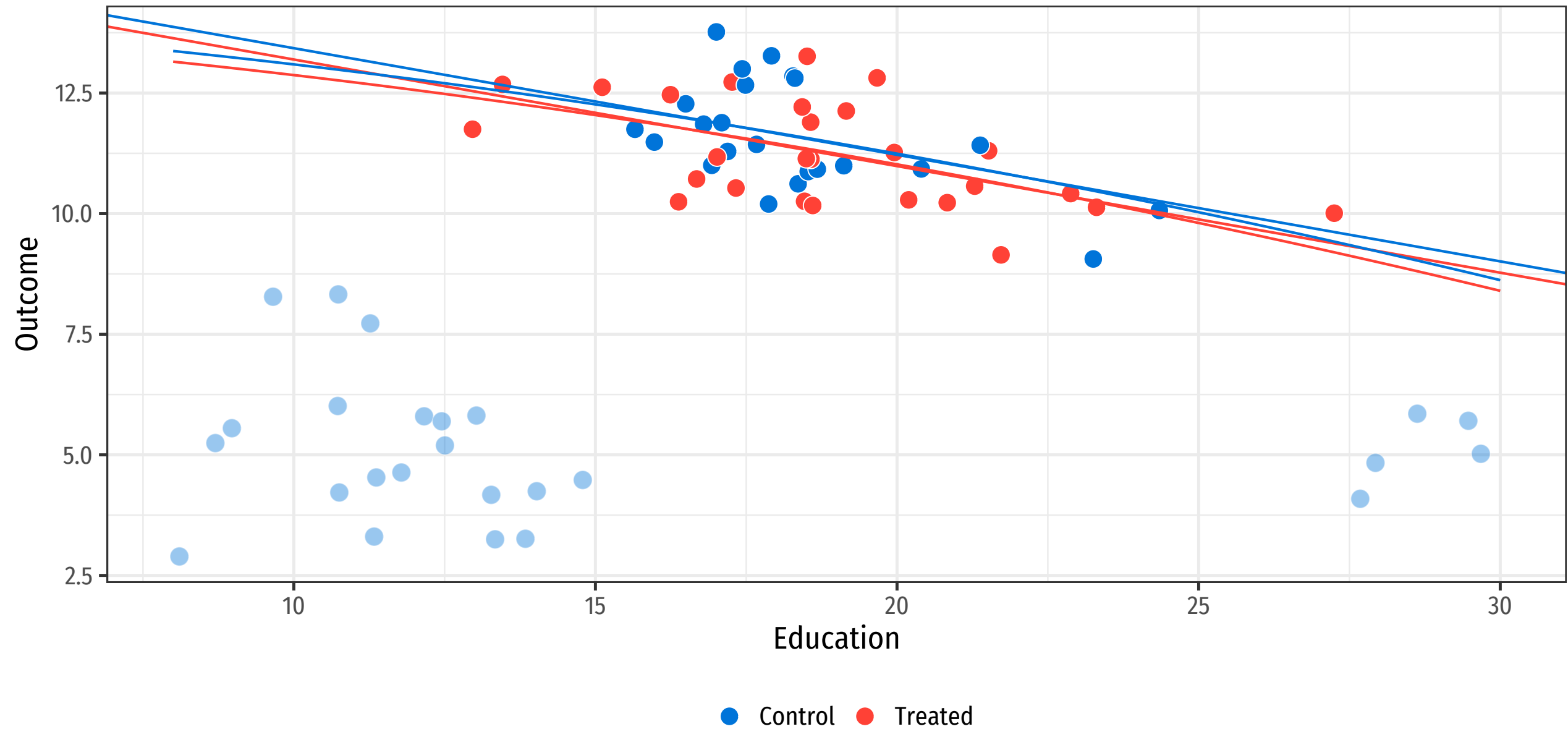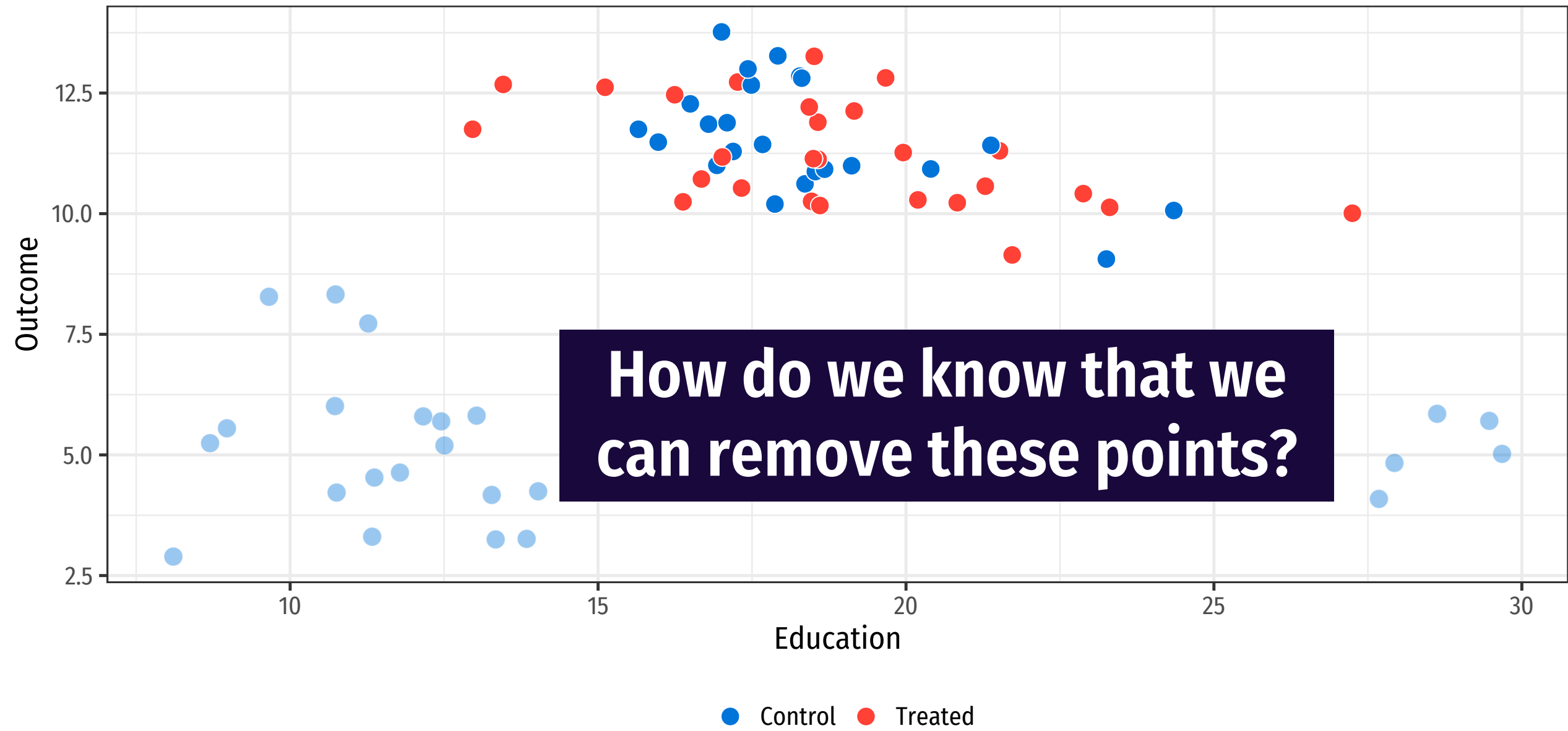$$\text{Outcome} = \beta_0 + \beta_1 \text{Education} + \beta_2 \text{Education}^2 + \beta_3 \text{Treatment}$$

How do we know that we can remove these points?
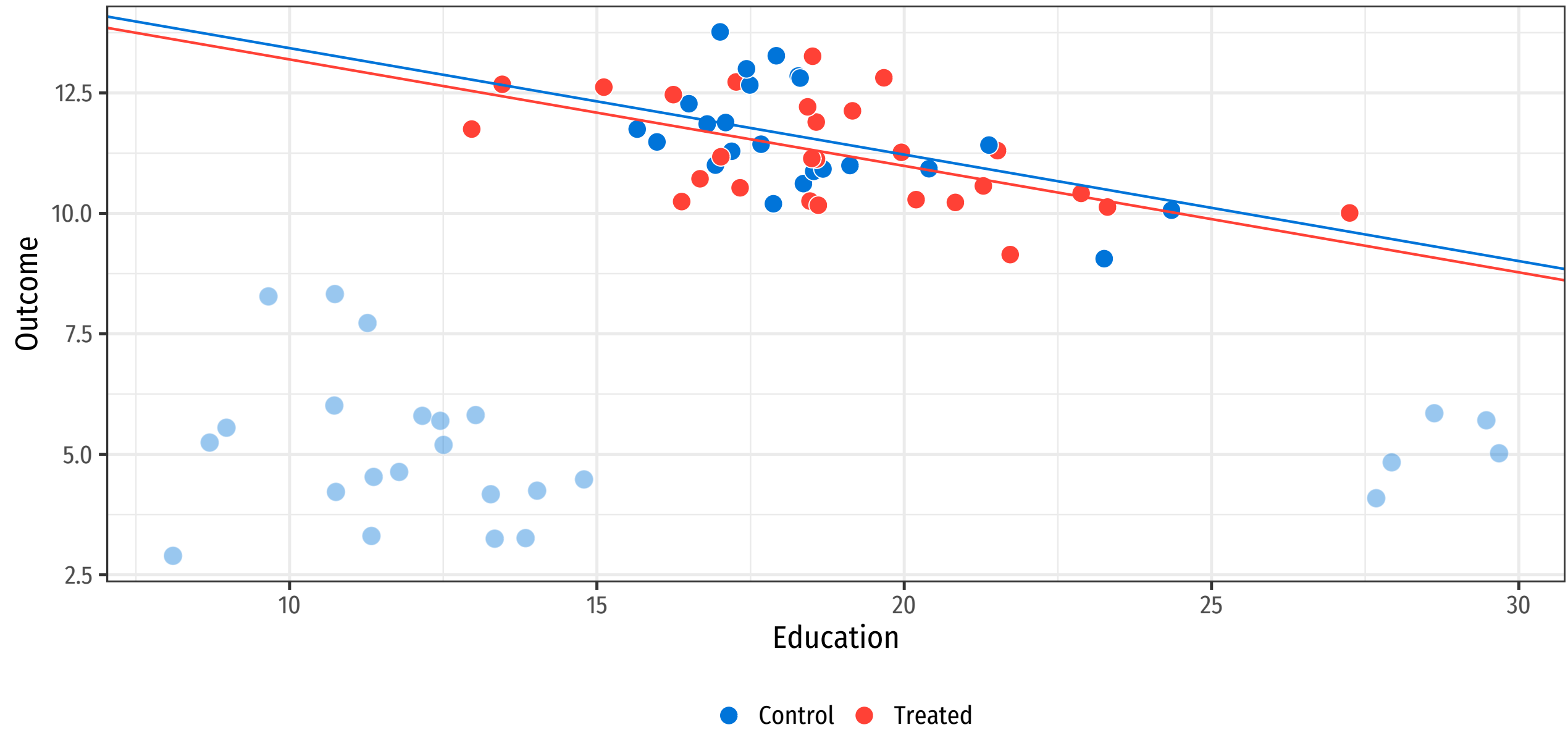
# General process for matching

## 1. Preprocess data

Do something to guess or model the assignment to treatment

Use what you know about the DAG to inform this!

## 2. Estimation

Use the new trimmed/preprocessed data to build a model, calculate difference in means, etc.

# Different methods

Nearest neighbor matching (NN)

Mahalanobis distance / Euclidean distance

Coarsened exact matching (CEM)

~~Propensity score matching (PSM)~~

Inverse probability weighting (IPW)

# Nearest neighbor matching

Find control observations that are very close/similar to treatment observations based on confounders

Lots of mathy ways to measure distance

Mahalanobis and Euclidean distance are most common

# There's a 70% chance of recession in the next six months, new study from MIT and State Street finds

PUBLISHED WED, FEB 5 2020·12:20 PM EST | UPDATED WED, FEB 5 2020·4:13 PM EST

**Pippa Stevens**
**@PIPPASTEVENS13**

SHARE

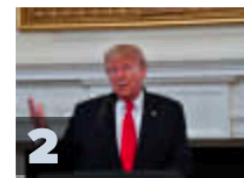**That's just Mahalanobis matching!**

**KEY POINTS**

- A new study from the MIT Sloan School of Management and State Street Associate says there's a 70% chance that a recession will occur in the next six months.

- The researches used a scientific approach initially developed to measure human skulls to determine how the relationship of four factors compares to prior recessions.

- The index currently stands at 76%. Looking at data back to 1916, the researchers found that once the index topped 70%, the likelihood of a recession rose to 70%.

**TRENDING NOW**
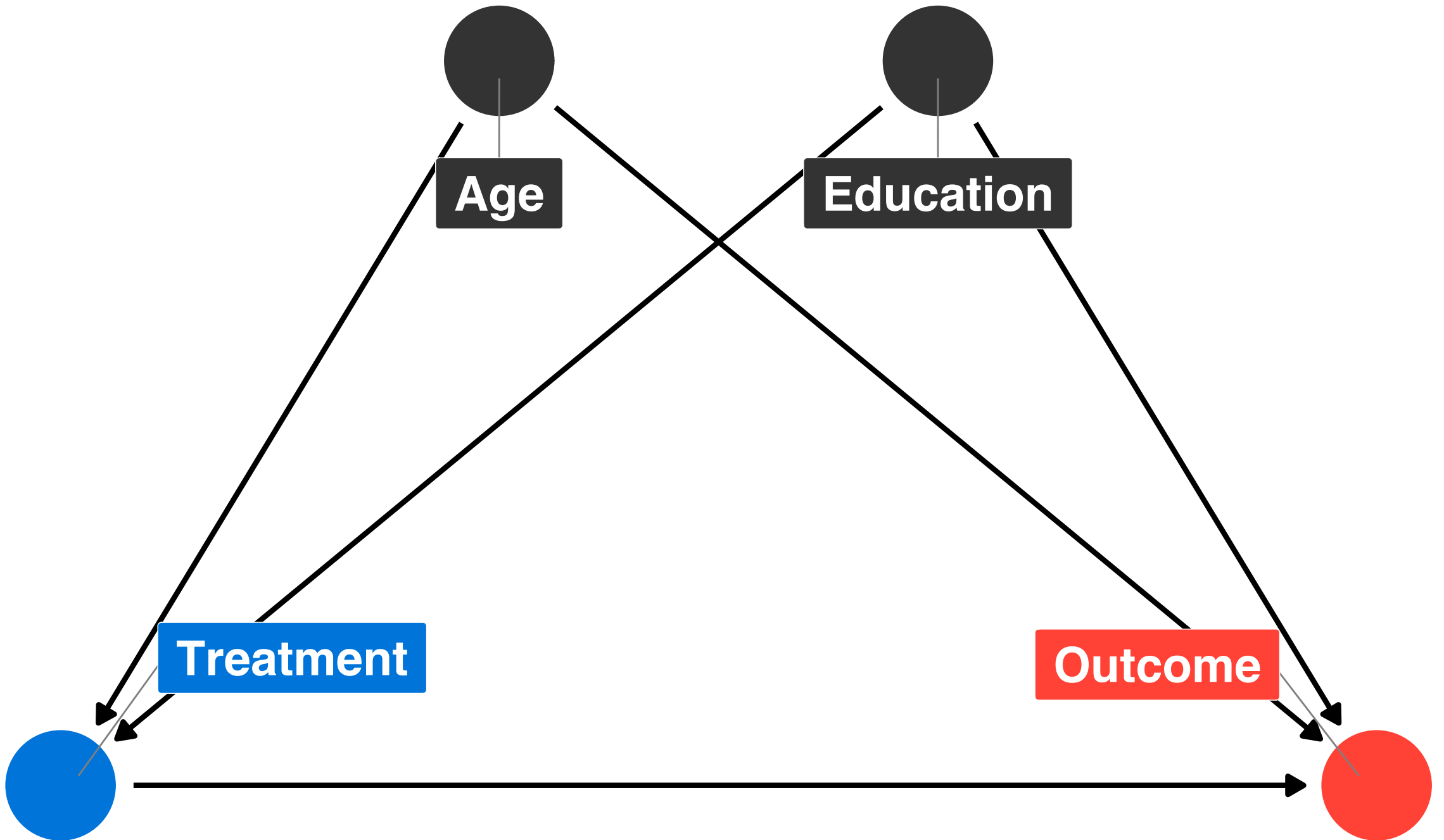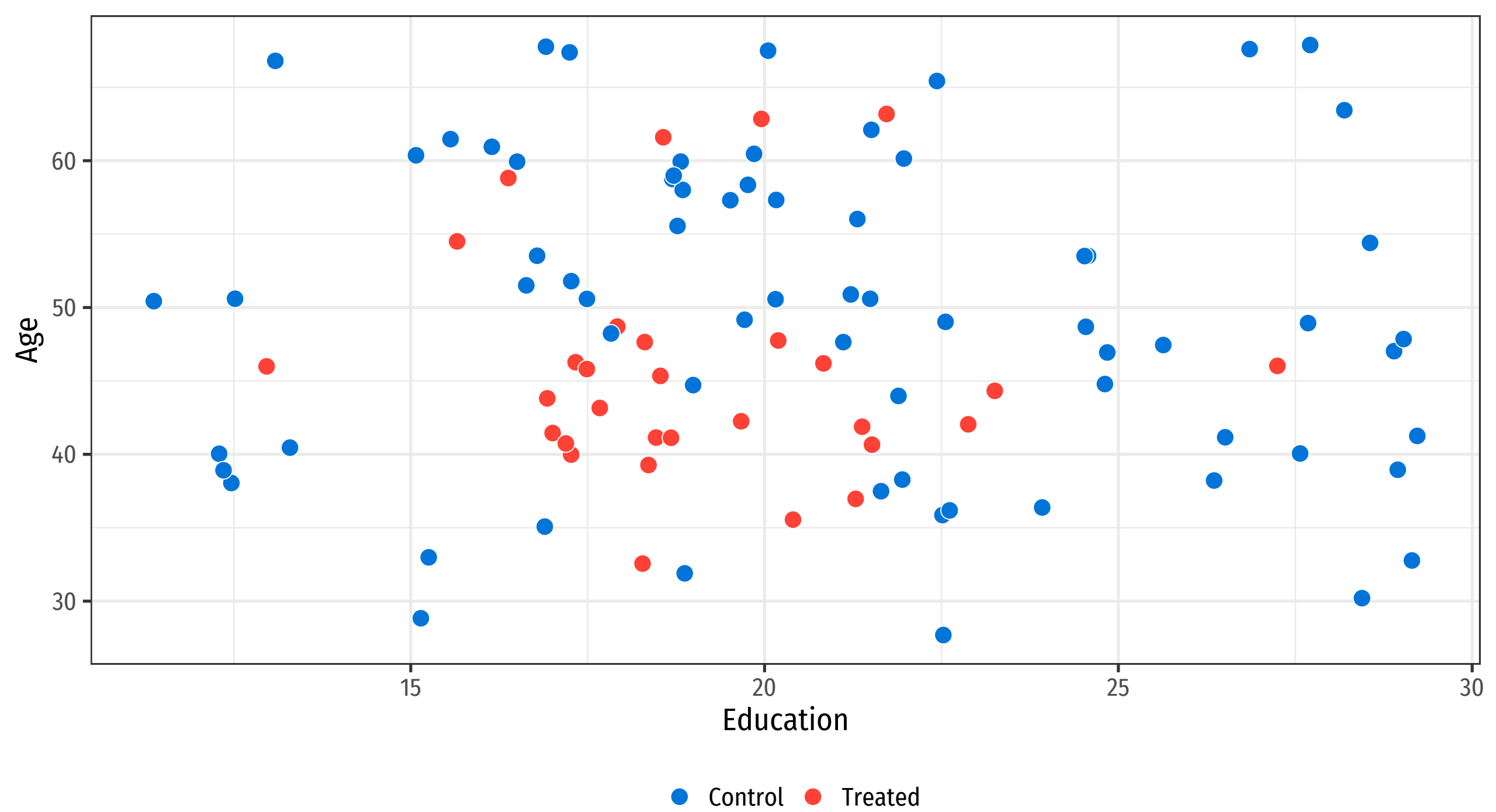
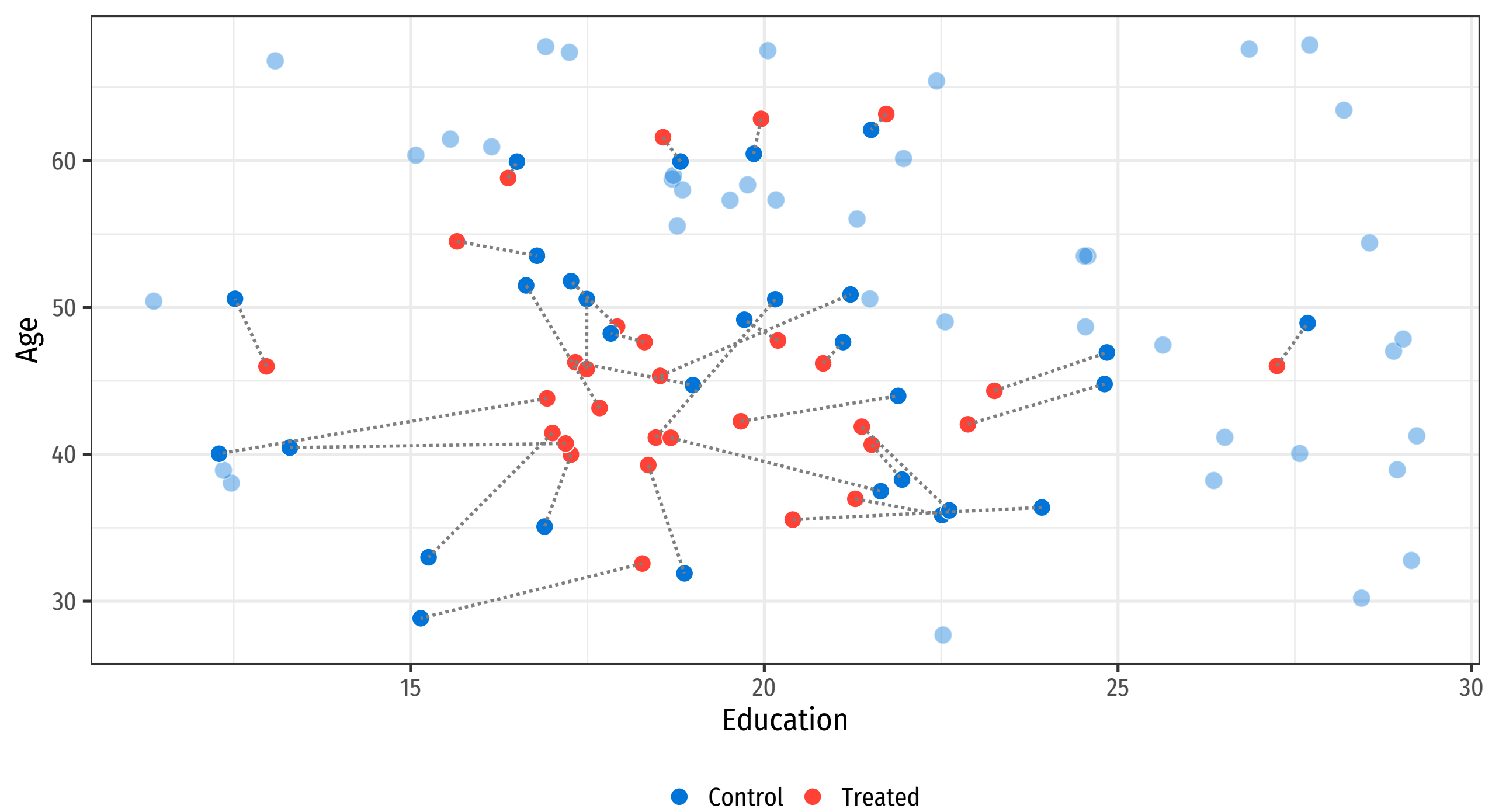1. Coror Brazi travel 'irrele

2. Trump furiou marke coror

# Prasanta Chandra Mahalanobis
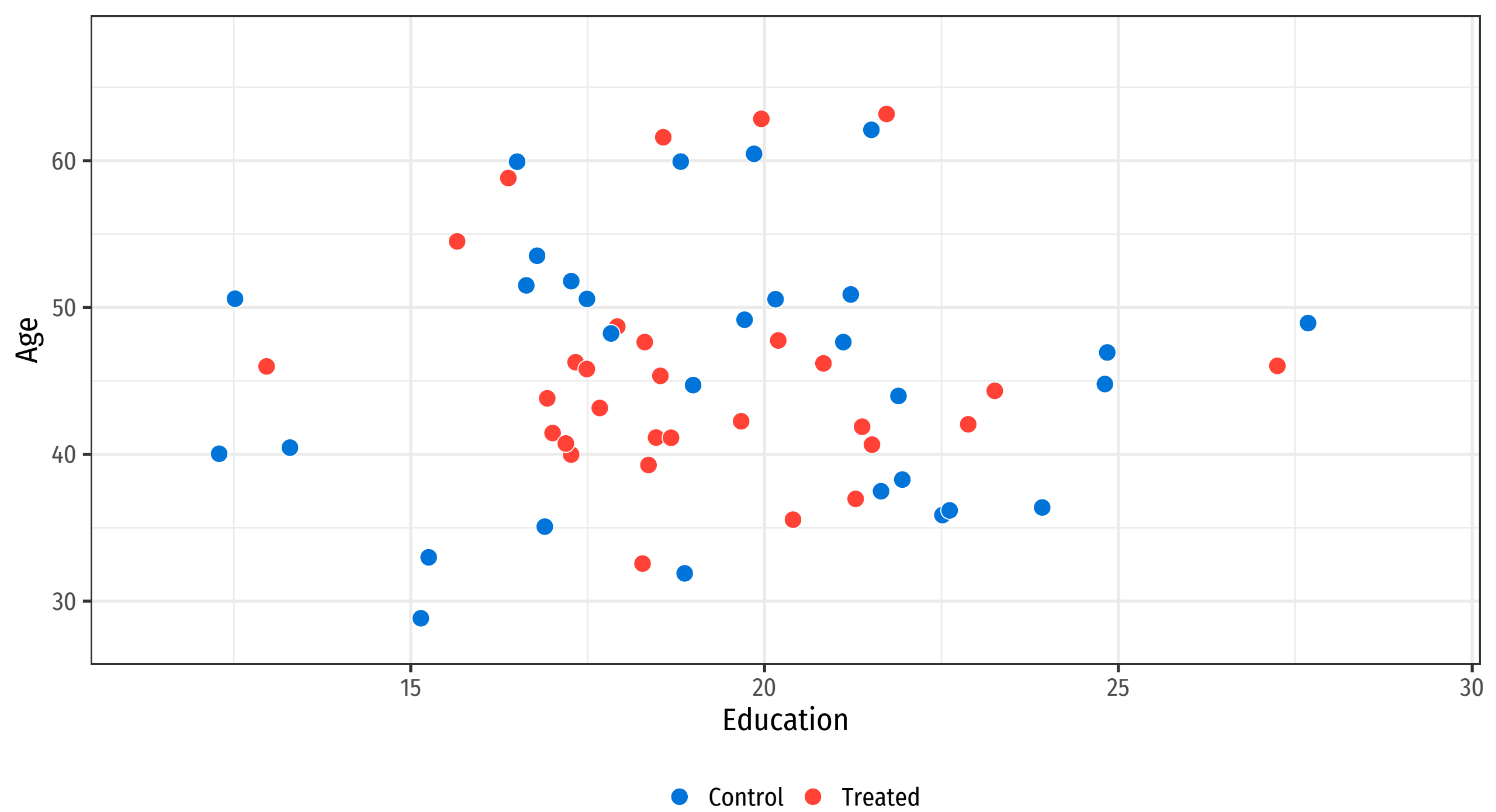
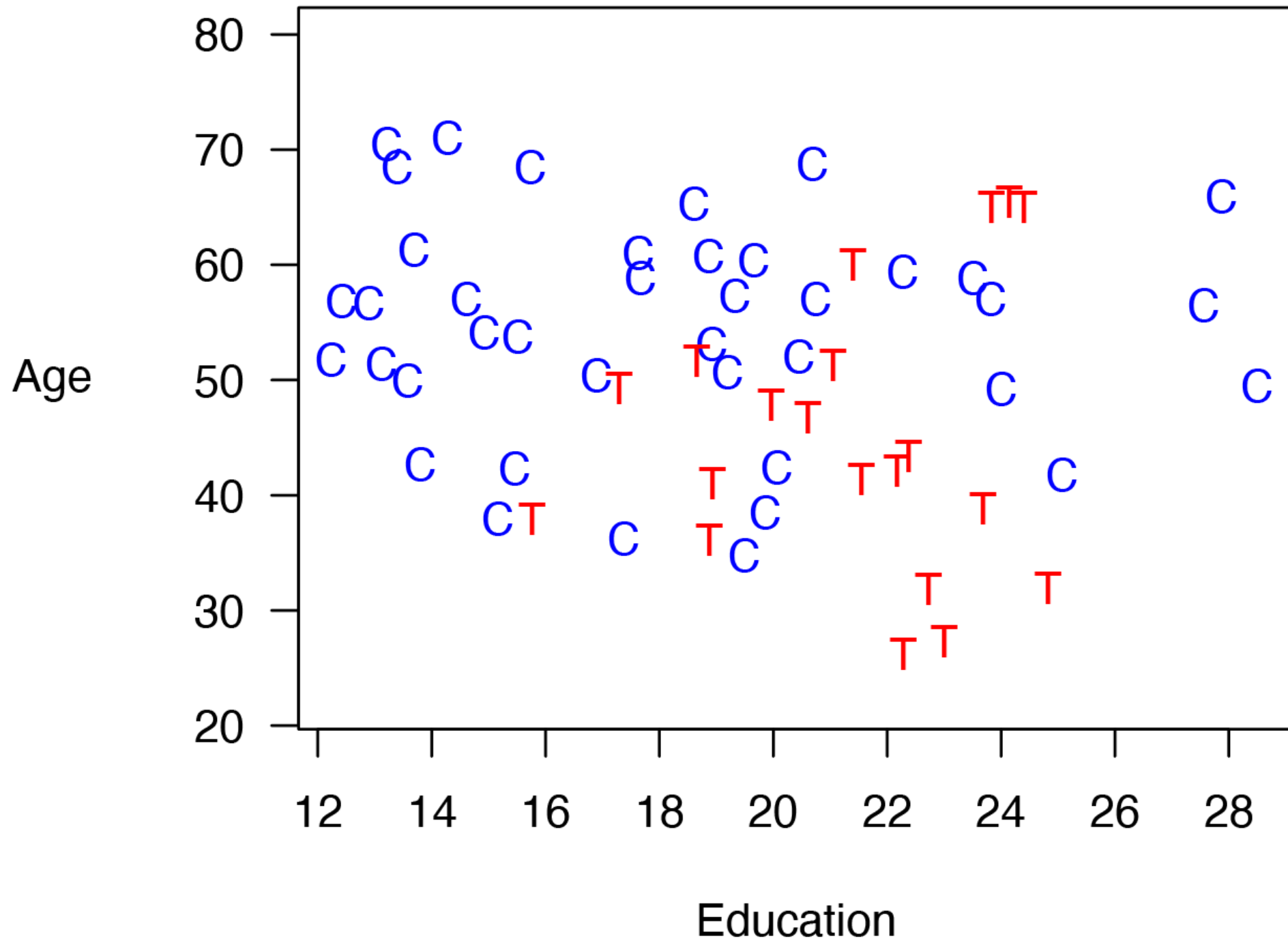Tried to prove brain size differences between castes; low-key eugenicist

# Coarsened exact matching

Use rules to partition data into clusters

Treatment should be random within clusters

Unconfoundedness again!

Some clusters will be more/less important

# Potential problems with matching

Nearest neighbor matching and CEM can be greedy!



Solution: Don't throw everything away

# Propensity scores

**Predict the probability of assignment to treatment using a model**

**Logistic regression, probit regression, machine learning**

$$\log \frac{p_{\text{Treatment}}}{1 - p_{\text{Treatment}}} = \beta_0 + \beta_1 \text{Education} + \beta_2 \text{Age}$$

$$\log \frac{p_{\text{Manual}}}{1 - p_{\text{Manual}}} = \beta_0 + \beta_1 \text{MPG}$$

```
model_transmission <- glm(am ~ mpg, data = mtcars, family = binomial(link = "logit"))
```



```
> tidy(model_transmission)
# A tibble: 2 x 5
  term          estimate std.error statistic p.value
  <chr>            <dbl>     <dbl>     <dbl>   <dbl>
1 (Intercept)      -6.60      2.35     -2.81 0.00498
2 mpg               0.307     0.115     2.67 0.00751
```

```
> tidy(model_transmission, exponentiate = TRUE)
# A tibble: 2 x 5
  term          estimate std.error statistic p.value
  <chr>            <dbl>     <dbl>     <dbl>   <dbl>
1 (Intercept)    0.00136      2.35     -2.81 0.00498
2 mpg            1.36         0.115     2.67 0.00751
```

# Plug all the values of MPG into the model and find the predicted probability

```
augment(model_transmission, data = mtcars, type.predict ="response")
```

```
# A tibble: 32 x 3
     mpg    am propensity
   <dbl> <dbl>      <dbl>
 1  21       1      0.461
 2  21       1      0.461
 3  22.8     1      0.598
 4  21.4     0      0.492
 5  18.7     0      0.297
 6  18.1     0      0.260
 7  14.3     0      0.0986
 8  24.4     0      0.708
 9  22.8     0      0.598
10  19.2     0      0.330
# … with 22 more rows
```

**Highly unlikely to be manual**

**Highly likely to be manual (1)**

# Propensity score matching

Super popular method

There are mathy reasons why it's not great for matching

Propensity scores are fine!
Using them for matching isn't!

# Why Propensity Scores Should Not Be Used for Matching

## Gary King[1] and Richard Nielsen[2]

[1] Institute for Quantitative Social Science, Harvard University, 1737 Cambridge Street, Cambridge, MA 02138, USA.
Email: king@harvard.edu, URL: http://GaryKing.org

[2] Department of Political Science, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA 02139,
USA. Email: rnielsen@mit.edu, URL: http://www.mit.edu/~rnielsen

## Abstract

We show that propensity score matching (PSM), an enormously popular method of preprocessing data for causal inference, often accomplishes the opposite of its intended goal—thus increasing imbalance, inefficiency, model dependence, and bias. The weakness of PSM comes from its attempts to approximate a completely randomized experiment, rather than, as with other matching methods, a more efficient fully blocked randomized experiment. PSM is thus uniquely blind to the often large portion of imbalance that can be eliminated by approximating full blocking with other matching methods. Moreover, in data balanced enough to approximate complete randomization, either to begin with or after pruning some observations, PSM approximates random matching which, we show, increases imbalance even relative to the original data. Although these results suggest researchers replace PSM with one of the other available matching methods, propensity scores have other productive uses.

*Keywords:* matching, propensity score matching, coarsened exact matching, Mahalanobis distance matching, model dependence

https://www.youtube.com/watch?v=rBv39pK1iEs

# Weighting in general

Make some observations more important than others

| | Young | Middle | Old |
|---|---|---|---|
| Population | 30% | 40% | 30% |
| Sample | 60% | 30% | 10% |

# Weighting in general

Make some observations more important than others

| | Young | Middle | Old |
|---|---|---|---|
| Population | 30% | 40% | 30% |
| Sample | 60% | 30% | 10% |
| Weight | 30 / 60 = 0.5 | 40 / 30 = 1.333 | 30 / 10 = 3 |

Multiply weights by average values (or use in regression) to adjust for importance

# Inverse probability weighting

## Use propensity scores to weight observations by how "weird" they are

Observations with high probability of treatment who don't get it (and vice versa) have higher weight
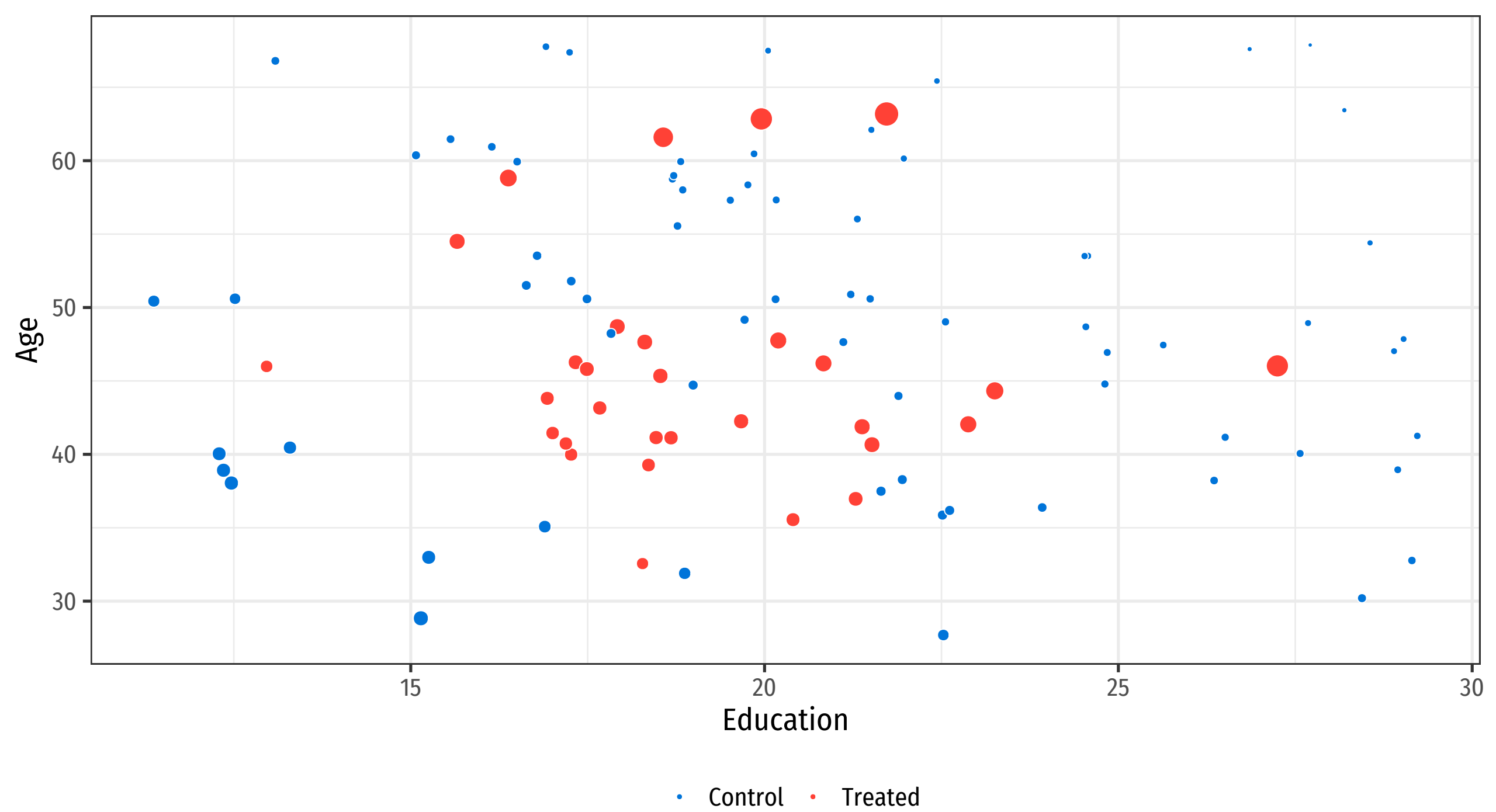
$$\frac{\text{Treatment}}{\text{Propensity}} + \frac{1 - \text{Treatment}}{1 - \text{Propensity}}$$

```
augment(model_transmission, data = mtcars,
        type.predict = "response") %>%
  select(mpg, am, propensity = .fitted) %>%
  mutate(ip_weight = (am / propensity) +
                     ((1 - am) / (1 - propensity)))
```

```
# A tibble: 32 x 4
     mpg    am propensity ip_weight
   <dbl> <dbl>      <dbl>     <dbl>
 1  21       1      0.461      2.17
 2  21       1      0.461      2.17
 3  22.8     1      0.598      1.67
 4  21.4     0      0.492      1.97
 5  18.7     0      0.297      1.42
 6  18.1     0      0.260      1.35
 7  14.3     0      0.0986     1.11
 8  24.4     0      0.708      3.43
 9  22.8     0      0.598      2.49
10  19.2     0      0.330      1.49
# … with 22 more rows
```

**Unlikely to be manual and isn't**

**Highly likely to be manual but isn't. Weird!**

# Other weights

**This gets you the ATE**

$$\frac{\text{Treatment}}{\text{Propensity}} + \frac{1 - \text{Treatment}}{1 - \text{Propensity}}$$

**Other versions of weights**
(Z = treatment;
e = propensity score)

$$w_{ATE} = \frac{Z_i}{e_i} + \frac{1 - Z_i}{1 - e_i}$$

$$w_{ATT} = \frac{e_i Z_i}{e_i} + \frac{e_i(1 - Z_i)}{1 - e_i}$$

$$w_{ATC} = \frac{(1 - e_i)Z_i}{e_i} + \frac{(1 - e_i)(1 - Z_i)}{1 - e_i}$$

$$w_{ATM} = \frac{\min\{e_i, 1 - e_i\}}{Z_i e_i + (1 - Z_i)(1 - e_i)}$$

$$w_{AT0} = (1 - e_i)Z_i + e_i(1 - Z_i)$$

R example