

Randomization

Silje Synnøve Lyder Hermansen

03-12-2019

Where are we? And what are we at?

We've completed the first part of the course: Congrats!

- ▶ Our focus has been on *describing* data: GLMs
- ▶ Now, we'll focus on *research design*: causal inference

The goal of the social sciences

Why do we run regressions?

We run regressions to learn about the world, which means...

- ▶ To describe data → observe the world
 - ▶ ... but how do we know if it's not an illusion?
- ▶ To make causal claims → manipulate the world
 - ▶ ... in the social sciences, that's not always possible

⇒ *We design studies to approximate manipulation*

We want to make causal claims

Two (compatible) approaches.

- ▶ **Logic of inference:** (King, Keohane and Verba, 1994)
 - ▶ We can only imperfectly observe the world
 - ▶ ... but we can theorize (causal mechanism)
 - ▶ ... and test hypotheses (observable implications)

⇒ *A closer connection between theory and statistics (e.g. EITM).*

We want to make causal claims

Two (compatible) approaches.

- ▶ **Logic of inference:** (King, Keohane and Verba, 1994)
- ▶ **Potential outcomes** (Donald Rubin)

What is causation?

A sequence of events in which – if the first didn't happen – the second wouldn't occur either.

- ▶ We can manipulate the first event → what happens then?
- ▶ Can we infer what *would have* happened if we did not manipulate?

⇒ *Potential outcomes*

We want to make causal claims

Two (compatible) approaches.

- ▶ **Logic of inference:** (King, Keohane and Verba, 1994)
- ▶ **Potential outcomes** (Donald Rubin)
 - ▶ *causal effect*: difference between what is and could have been

⇒ *a set of methods designed for causal inference with observational data*

The conundrum

The true causal effect

What is causal effect?

Imagine two versions of me.

- ▶ I have a headache and I take an aspirine ($Y_{1,Silje}$).
- ▶ I have a headache but receive no treatment ($Y_{0,Silje}$).

⇒ *the causal effect is $Y_1 - Y_0$*

True causal effect $Y_{1, \text{silje}}$  $Y_{0, \text{silje}}$

A causal effect is the difference between two potential outcomes

- ▶ ... but – at best – I can only observe one outcome.

**True causal effect is
NOT POSSIBLE
to observe**

$Y_{1, \text{silje}}$



~~$Y_{0, \text{silje}}$~~

A causal effect is the difference between two potential outcomes

- ▶ ... but – at best – I can only observe one outcome.

⇒ *We have to compare two different individuals*

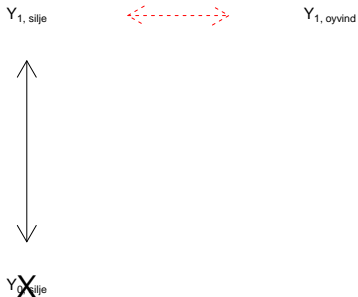
Plan B

Plan B: Can we compare across cases?

Let's compare my headache now with Øyvind's current headache
 $(Y_{1,Silje} - Y_{1,Oyvind})$

Let's compare my headache now with Øyvind's current headache
($Y_{1,Silje} - Y_{1,Oyvind}$)

Can we compare two individuals
post treatment?

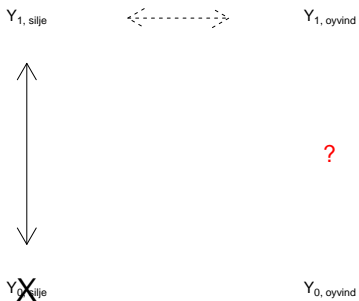


Let's compare my headache now with Øyvind's current headache
 $(Y_{1,Silje} - Y_{1,Oyvind})$

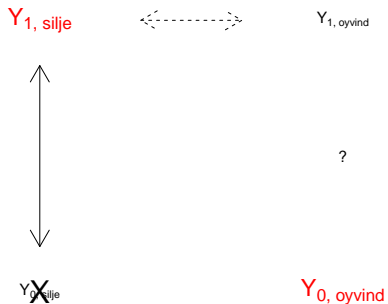
- ▶ ... but did he even have a headache before?

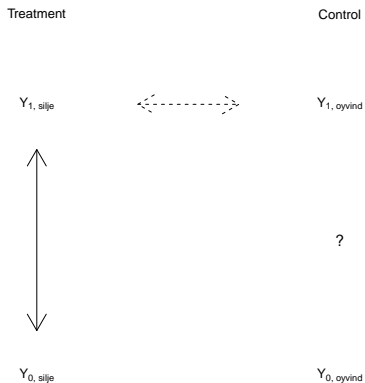
Is there a selection bias?

How did Øyvind's case
look untreated?



How did Øyvind's case look untreated?



What do we compare?

Where's the selection bias?

The solution

We have to observe Øyvind's untreated headache ($Y_{0,Oyvind}$) and compare with treated me ($Y_{1,Silje}$)

$$\begin{aligned}
 Y_{Silje} - Y_{Oyvind} &= Y_{1,Silje} - Y_{0,Oyvind} \\
 &= Y_{1,Silje} - Y_{0,Silje} + Y_{0,Silje} - Y_{0,Oyvind}
 \end{aligned} \tag{1}$$

- ▶ **Causal effect:** $Y_{1,Silje} - Y_{0,Silje}$
- ▶ **Selection bias:** $Y_{0,Silje} - Y_{0,Oyvind}$

How to do it?

How to do it?

We use statistics

We cannot observe two potential outcomes, but we can rely on the law of large numbers (LLN).

- ▶ We use **average** causal effect

Average causal effect = Differences in means - Selection bias

Differences in means

- ▶ We create a **dummy** for treated vs. untreated observations:

$$D_i = \begin{cases} 1 & \Leftrightarrow \textit{treated} \\ 0 & \Leftrightarrow \textit{untreated} \end{cases} \quad (2)$$

- ▶ We calculate the **differences in means**

$$= \text{Avg}_n[Y_i | D_i = 1] - \text{Avg}_n[Y_i | D_i = 0] \quad (3)$$

Differences in means

- ▶ We create a **dummy** for treated vs. untreated observations:

$$D_i = \begin{cases} 1 & \Leftrightarrow \text{treated} \\ 0 & \Leftrightarrow \text{untreated} \end{cases} \quad (4)$$

- ▶ We calculate the **differences in means**

$$\begin{aligned} &= \text{Avg}_n[Y_i | D_i = 1] - \text{Avg}_n[Y_i | D_i = 0] \\ &= \text{Avg}_n[Y_{1,i} | D_i = 1] - \text{Avg}_n[Y_{0,i} | D_i = 0] \end{aligned} \quad (5)$$

Basic assumption

We have to assume that the treatment has the same effect across all units

- ▶ then we can compare across units
- ▶ contrast that with the effect of β in OLS vs GLM

Selection bias

Now we have to get rid of the selection bias!

- ▶ **A priori** selecting units without bias: randomization
- ▶ **A posteriori** assessing the bias and extract it: Rubin's contribution

Why not just compare?

Consider the fate of young mothers

[https://www.thelancet.com/journals/lancet/article/PIIS0140-6736\(16\)31411-8/fulltext](https://www.thelancet.com/journals/lancet/article/PIIS0140-6736(16)31411-8/fulltext)

The gold standard

Randomization

Randomization is the gold standard. This requires

- ▶ manipulation → experiments
- ▶ a sufficient number of units (LLN) → statistical power

⇒ *Randomization eliminates bias*

Checking on observables

Even when we randomize, we check for signs of selection bias

- ▶ we cannot observe the bias
- ▶ but we can check the balance of possible correlates (of bias)

⇒ *Here comes the social science theories back in!*

Checking on observables

Even when we randomize, we check for signs of selection bias

⇒ *We verify the balance of pre-treatment variables*

The post hoc fixes