

Foundations of Causal Inference

Sebastian Tello-Trillo¹

¹University of Virginia

Spring 2026

Research Designs

Definition

A (causal) research design is a statistical and/or economic statement of how an empirical research paper will estimate a relationship between two (or more variables) that is causal in nature - X causing Y.

- In some cases, the comparisons become so popular they get names:
 - ▶ Mean Comparison
 - ▶ Matching
 - ▶ Instrumental Variables
 - ▶ Regression Discontinuity
 - ▶ Differences-in-Differences
 - ▶ Synthetic Control
- Each research design can recover the causal parameter if the assumptions are validated.

Tools

Definition

Tools are the statistics we use in order to apply the research designs into the data

- Examples:
 - ▶ Averages, Median
 - ▶ OLS
 - ▶ Probit
 - ▶ Logit
 - ▶ Non-parametric Regression

How we will structure learning about designs

- We will mostly follow this structure:
 - ① Understanding the intuition behind the design: What's the comparison? What's the counterfactual?
 - ② What are the assumptions made to obtain causality
 - ★ Using conditional expectation will be helpful in revealing this
 - ③ How to use a tool to apply the design
 - ④ How can we provide evidence for the assumptions
 - ⑤ Robustness and sensitivity checks

Learning about designs

- The key insights about causal design is to think of design that can answer the following questions:
 - ▶ What's the comparison?
 - ▶ How does this design characterize what would have happened, had this “event” not happened?
 - ▶ What's the difference in outcomes between what happened, and what would have happened?
 - ▶ What do we need to assume to characterize that difference as causal?
 - ★ **There will always be assumptions**
- Any new design, fancy term, that claims causality, ask yourself: What's the comparison?

Learning about designs

- Let's practice learning about designs with the simplest design of them all: **Mean Comparison**

Is Batten worth it?

Questions

- What is the simplest comparison you can do in order to determine if going to Batten has an effect on earnings?
- How would you write that comparison using conditional expectation language?
- What other comparisons we can do?

Is Batten worth it?

Skill

Transforming comparisons people propose colloquially into conditional expectation language can help identify (1) the assumptions people are making (that they don't know they are making) and (2) helps you clarify the difference between the comparisons and what they are implying about the causal effect.

Practice with worksheet!

Is Batten worth it?

- Write the comparison in conditional expectation form

$$E(\text{Income}_i^1 | \text{Batten} = 1) - E(\text{Income}_i^0 | \text{Batten} = 0)$$

- Let's this comparison is giving us a number, say 20K. We could then say, Batten increases earnings of students by 20K a year.

Question

Intuitively, why do we think this claim is not quite right? Why is this comparison "not enough"?

Is Batten worth it?

Skill

A great way of finding potential bias is to apply the following thinking: “imagine the true effect is a negative, positive or no effect”, why is it that the finding given is (+,-, or no effect)?. This can help you find a reason of why there is bias.

Is Batten worth it?

- Mean Comparison

$$E(\text{Income}_i^1 | \text{Batten} = 1) - E(\text{Income}_j^0 | \text{Batten} = 0)$$

Question

Can this comparison recover the causal effects?

What assumptions are needed for this design to recover the causal parameter?

Figuring out the assumption

- Using potential outcomes framework is a helpful tool in figuring out each design's assumption
- Let's review some notation

$y_i = \text{Outcome of Person } i \text{ or group } i$

$y_{alex} = \text{Income of alex}$

$y_i^0 = \text{Outcome for } i \text{ when timeline } 0$

$y_i^1 = \text{Outcome for } i \text{ in timeline } 1$

- $y_i^1 | D = 1$ Outcomes for observation that received D, in timeline 1

Figuring out the assumption

- What we really want is $y_i^1 - y_i^0 = \theta_i$. The causal effect of 1 for person i .
- Since we talk about the mean in order to generalize what happens with everyone we then write

$$E(y_i^1) - E(y_i^0) = \theta_i$$

- Conditioning on a group

$$E(y_i^1|D = 1) - E(y_i^0|D = 1) = \theta_i|D = 1$$

Figuring out the assumption

- In the real world we observe if a person has receive a treatment or not

$$E(y_i^1 | D = 1) - E(y_j^0 | D = 0)$$

RMDA I callback

Practice reading this, and ask yourself what's observable and not.

Figuring out the assumption

- So for our example, what we are doing is

$$E(\text{Income}_i^1 | \text{Batten} = 1) - E(\text{Income}_j^0 | \text{Batten} = 0)$$

- And what we want is

$$E(\text{Income}_i^1 | \text{Batten} = 1) - E(\text{Income}_i^0 | \text{Batten} = 1)$$

$$E(y_i^1 | D = 1) - E(y_i^0 | D = 1)$$

Let's do some magic

$$\begin{aligned} & E(y_i^1 | D_i = 1) - E(y_j^0 | D_j = 0) \\ E(y_i^1 | D_i = 1) - E(y_j^0 | D_j = 0) &+ E(y_i^0 | D_i = 1) - E(y_i^0 | D_i = 1) \\ E(y_i^1 | D_i = 1) - E(y_i^0 | D_i = 1) &+ E(y_i^0 | D_i = 1) - E(y_j^0 | D_j = 0) \end{aligned}$$

Causal Effect of Batten Bias

Assumptions

- This means that the simple mean comparison does contain the causal effect, but also contains another term, we call this term bias
- Can we recover the causal parameter from the simple mean comparison?
- YES, if only if the bias term is 0, which means the following:

$$\underbrace{E(y_i^0 | D_i = 1) - E(y_j^0 | D_j = 0)}_{\text{Bias}}$$

$$E(y_i^0 | D_i = 1) = E(y_j^0 | D_j = 0)$$

Bias is 0?

Question

What does this mean? What is this trying to say?

- That the outcomes of the treatment group, the group that went to batten, had they not gone to batten would have been - on average- the same as the outcomes of the group that did not go to Batten.

Bias violations

- There could be many reasons why the bias is not 0, some have names:
 - ▶ selection bias
 - ▶ measurement error
 - ▶ reverse causality
 - ▶ omitted variable bias
 - ▶ etc.

Question

How does RCTs solve this challenge? What are the assumptions in an RCT?

How to apply this research design?

- What tools can we use?
 - ▶ Average
 - ▶ Regression
- How to apply using averages?
- What's the comparison group?
 - ▶ This raises the question of “what's the counterfactual?” What group are we comparing to?
 - ▶ Example: Batten is better relative to what?.
 - ▶ Second example: breastfeeding

How to apply this research design?

- How to apply this method using regression?
- What's the relationship between conditional expectation and regression?

Applying these concepts: COVID-19

- Statement: “I work in a hospital 5 days a week, I’ve seen it all, and I’m tired of this situation, but to be honest I don’t think vaccines are the solution. At the end of the day, if I’m being honest most of the people who end up hospitalized because of COVID-19, have the vaccine. In fact, I think I heard that 75% of the people who are hospitalized were vaccinated.”

Question

What’s the comparison this healthcare worker is making?

Write it in conditional expectation.

What’s the causal question or causal effect they are implying?