

Overview

1. Why causality?
2. The fundamental problem of causal inference
3. Problems
4. Solutions
5. The gold-standard (RCT, A/B testing)

Why causality ?

- Every decision we take relies on causal relationships :
 - **Individuals** :
 - If I go vegan, I'll reduce my ecological footprint.
 - If I drink this tequila shot, I'll dance better.
 - **Companies** :
 - Home-office reduces productivity.
 - Spamming users with YouTube Premium Ads will increase the number of subscribers.
 - **Policy makers** :
 - Replacing nuclear power plants with renewables will help to reach the Paris Agreement.
 - Lockdowns will reduce the spread of the covid-19.
- Failing to properly assess causality might lead to costly mistakes.

The fundamental problem of causal inference



- <https://www.youtube.com/watch?v=0zvrgiPkVcs&t=58s>

There is only one Africa - if we try to understand if the aid we are giving is helping -> must understand was impact the aid has , Africa might be better without Aid or with more Aid?

The fundamental problem of causal inference



- <https://www.youtube.com/watch?v=0zvrgiPkVcs&t=58s>
- "How do we know what would have happened without the aid? We have no idea. We don't know what the **counterfactual** is. There is only one Africa."
- It's impossible to observe the outcome with and without treatment for the same entity at the same point in time.

The fundamental problem of causal inference

Rubin Causal Model

- Treatment : $D_i = 1$ if treated, 0 otherwise
D = Dummy = Treated group

- Potential outcome :

- Y_{i0} : Individual i outcome without treatment
- Y_{i1} : Individual i outcome with treatment

$$\Rightarrow Y_i = Y_{i0} \text{ if } D_i = 0 \text{ or } Y_i = Y_{i1} \text{ if } D_i = 1$$
$$\Rightarrow Y_i = Y_{i0} + (Y_{i1} - Y_{i0})D_i$$

can't necessarily observe both
outcomes like eg: Africa but I must
understand -> So what to do?

The fundamental problem of causal inference

Rubin Causal Model

- Treatment : $D_i = 1$ if treated, 0 otherwise
- Potential outcome :
 - Y_{i0} : Individual i outcome without treatment
 - Y_{i1} : Individual i outcome with treatment
 - $\Rightarrow Y_i = Y_{i0}$ if $D_i = 0$ or $Y_i = Y_{i1}$ if $D_i = 1$
 - $\Rightarrow Y_i = Y_{i0} + (Y_{i1} - Y_{i0})D_i$
- Causal effect : $Y_{i1} - Y_{i0}$

The fundamental problem of causal inference

Rubin Causal Model

- Instead of $Y_{i1} - Y_{i0}$ we can measure $E(Y_i^A | D_i = 1) - E(Y_i^B | D_i = 0)$
 A = expected average outcome of treated group
 B = expected average outcome of treated group
 -> look at difference between A & B

$$= E(Y_{i1} | D_i = 1) - E(Y_{i0} | D_i = 0) + \underbrace{E(Y_{i0} | D_i = 0) - E(Y_{i0} | D_i = 1)}_{\text{will never observe this}}$$

The fundamental problem of causal inference

Rubin Causal Model

- Instead of $Y_{i1} - Y_{i0}$ we can measure $E(Y_i|D_i = 1) - E(Y_i|D_i = 0)$

$$E(Y_i|D_i = 1) - E(Y_i|D_i = 0) =$$

$$E(Y_{i1}|D_i = 1) - E(Y_{i0}|D_i = 1) + [E(Y_{i0}|D_i = 1) - E(Y_{i0}|D_i = 0)]$$

- This "naive" comparison includes :
 - Average Treatment of the Treated (ATT)
 - "Selection" Bias



The fundamental problem of causal inference

Rubin Causal Model

- There is a bias if :
- $E(Y_{i0}|D_i = 1) \neq E(Y_{i0}|D_i = 0)$
- The average outcome for the treated and untreated would be different without treatment

The gold-standard : RCT, A/B testing

- A Randomized Control Trial is a controlled experiment (in opposition to a natural experiment) where you randomly allocate the treatment between groups.

⇒ By **randomly allocating the treatment** to the different groups, you can **solve the selection bias**.

if randomly allocated -> they are perfectly comparable , include conditions in simple allocations if the sample is small for example - groups must combine man and woman ..



RCT

Limitations

1. **Not always possible** (e.g. gender) or ethical (e.g. weapons) to manipulate the treatment.

2. **External validity vs. internal validity :**
 - Example with seating position and learning
 - <https://www.sciencedirect.com/science/article/pii/S0959475217305716>
 - <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0236131>

3. **Simple random allocation works with large samples or if the groups are relatively homogenous :**
 - Large sample : by the law of large numbers, on average the groups will be similar.
 - Homogeneity : In a lab experiment, you have inbred strains of rats (almost identical genetically).

⇒ Cluster/Stratified sampling method.

eg: sitting in the further back from the teacher -> there is no causality to the grade -> should check with random allocation

Limitations

4. **Blinding is not always possible !** If possible, the researchers and the subjects do not know who receives the treatment. At least, we should always do our best to prevent the subjects from knowing if they are in the treatment or control group.

Two types of statistical errors

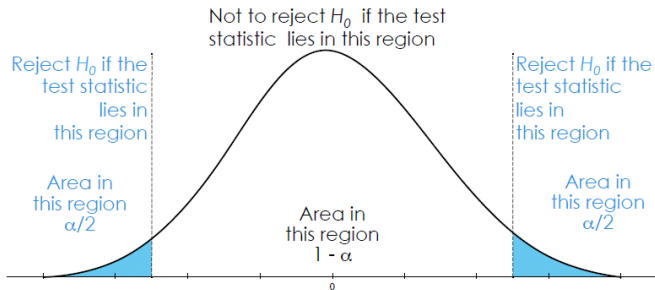
		True state of nature	
		Effect present	Effect absent
Conclusion of statistical analysis	Effect present (reject H_0)	Correct	Type I error (α)
	No effect (accept H_0)	Type II error (β)	Correct

Goal is to reject the null hypotheses

Type I error : α

- In statistics you “never” get an answer with a 100% certitude
- α is the probability of wrongly rejecting the null hypothesis
- Defined by the researcher (most frequently 5%)
- α will define the rejection rule of H_0
 - if $\text{p-value} \leq \alpha \Rightarrow \text{reject } H_0$
 - if $\text{p-value} > \alpha \Rightarrow \text{not reject } H_0$

Example



Weakness of the p-value

- **Statement on p-values (ASA Statement (2016))**

- Indicates how incompatible the data are with a specified statistical model
- Scientific conclusions should not be based only on the p-value
- Proper inference requires full reporting and transparency
- p-value states nothing about the magnitude

Need a rational understanding why the results I am getting are reasonable

- My view on this :

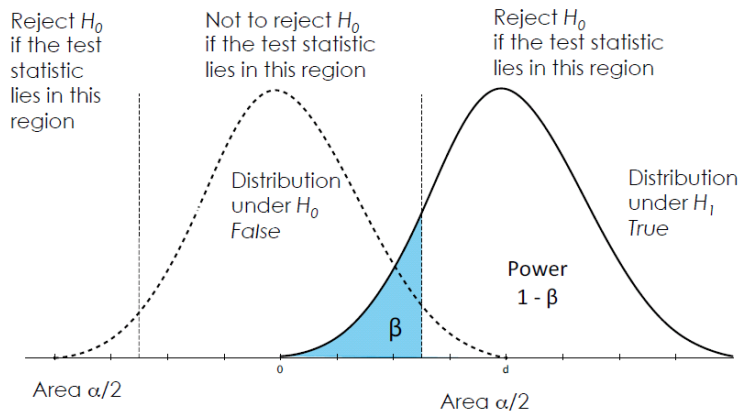
- The statistical significance should be used as a necessary condition to interpret the magnitude, nothing more.
- (The worst idea ever is to get rid of the stars in tables. Rather show the p-values without stars.)
- Read : *Why and how to use forest plots efficiently?*

<https://medium.com/towards-data-science/>

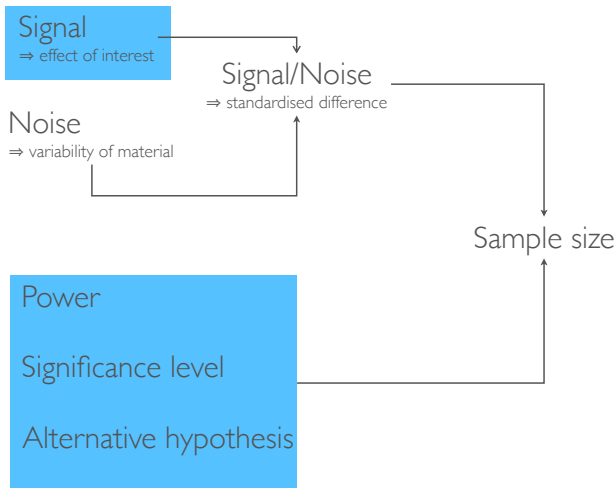
unhappy-with-statistical-significance-p-value-here-is-a-simple

The p-value says nothing about the magnitude of the effect eg- shampoo promising new hairgrowth -> I care about the magnitude

Example



Sample size is influenced by five variables



Power and sample size

- Relationship between six variables
 - The effect size of practical interest
 - The standard deviation
 - The significance level
 - The desired power of the experiment
 - The sample size
 - The alternative hypothesis (i.e. one or two-sided test)

3. The significance level (α)

- $\alpha \nearrow \Rightarrow \text{Power} \nearrow$
- However, other things being equal, specifying a low chance of a false-positive result will increase the chance of a false-negative result
- By convention, fixed to 5%

4. Power of the experiment γ

- Power $\nearrow \Rightarrow$ sample size \nearrow
- $\gamma = 1 - \beta$ (β = Allowed probability of wrongly not rejecting the null hypothesis H_0) The aim should be to have powerful experiments that have a high chance of detecting an effect if it exists (i.e. low β error)
- **Usually set between 80% and 90%**

⚠ To be considered in the interpretation of “negative results” : One can only conclude that differences in effects between two treatments were certainly absent if the study would have had enough POWER to detect them. FALSE! post-hoc test doesn't work!

6. The alternative hypothesis (H_1)

- The usual null hypothesis \Rightarrow no differences among treatment means
 $\Rightarrow H_1 =$ there is an effect \Rightarrow two-sided significance test
- If the alternative is that means differ in a particular direction \Rightarrow one-sided test

Sample size - Principle

- Specify the smallest true difference between the treatments that would be of practical relevance
- **Choose the smallest sample size allowing to test your hypothesis efficiently (power)**

What do we need to compute the sample size

- Type of the outcome variable (continuous, categorical, proportion etc.) and number of comparisons
 - ⇒ Choice of statistical test
- Direction of the test : one or two-sided
- Standard deviation (or variance)
- Effect size

⇒ Use Python or G-power (let's see both)