

The Potential Outcome Framework

Alessandro Caggia

June 2025

Abstract

■ Plan for Research

- What is the **Causal relationship** of interest?
- What is the **Ideal Experiment** to capture it?
- What is the **identification strategy**? that is, how to use what you have to actually approximate 2? eg. you have an exogenous Source of Variation (DiD, RD etc)
- What is your **mode of statistical inference**? what is the population of interest? what is the sample? What assumptions for SE?

■ Causality

- Definition **in terms of counterfactual**: difference between the direct effect of that cause and the state that would have been observed had that specific intervention not taken place (the counterfactual).
- Correlation does not imply causation**: we can define correlation based on the joint distribution of two variables but for causation we need further assumptions (think at OVB).

To define causality, three possible structures:

- A. Potential Outcomes model (Rubin 1974)
- B. Structural model
- C. Granger Sims causality: views causality as a prediction property

1. The Potential Outcome Framework

Causality is tied to an action applied to some units (eg treatment is a job training program).¹

- We primarily consider settings with **2 actions**:
 - $D_i = 1$ if exposed to the treatment,
 - $D_i = 0$ if exposed to no treatment (control).
- For each individual i , there are **two Potential Outcomes**:
 - Y_{1i} : potential outcome if individual i is treated ($D_i = 1$).²
 - Y_{0i} : potential outcome if individual i is not treated ($D_i = 0$).³
- We define the **Individual Causal (or Treatment) Effect** of D as: $Y_{1i} - Y_{0i}$.
- Switching equation:

$$Y_i = Y_{0i} + (Y_{1i} - Y_{0i})D_i.$$

SUTVA

- The potential outcomes for any unit do not vary with the treatments assigned to other units (ie no spillovers).**
- Strong behavioral assumption** that rules out social interactions, peer effects, network effects and many types of general equilibrium (GE) effects and other spillovers.
- Possible workaround: when spillovers are an important part of the analysis, we can choose a more aggregate unit, e.g. one school, or school class instead of one student.

¹Note: careful in distinguishing actions from attributes (e.g. gender).

²This has not to be a single value, as you will see in the following example it could be two values or a distribution of values)

³Remark: 1 vs 0 has nothing to do with time!

Example

- Effect of taking an aspirin on having a headache.
- Actions: $D_i = 1$ if aspirin, $D_i = 0$ if no aspirin
- Unit: individual i
- Outcome:
 - $Y_i = Y_{1i}$: headache ($Y_{1i} = 1$) or no headache ($Y_{1i} = 0$). Only observed if took an aspirin.
 - $Y_i = Y_{0i}$: headache ($Y_{0i} = 1$) or no headache ($Y_{0i} = 0$). Only observed if did not take an aspirin.
- Individual Causal Effect: $Y_{1i} - Y_{0i} \in \{-1, 0, 1\}$

| | Y_{1i} | Y_{0i} | $Y_{1i} - Y_{0i}$ |
|---|--|----------|-------------------|
| 1 | Headache gone only if you take aspirin | 0 | 1 |
| 2 | Headache not gone even if you take aspirin | 1 | 0 |
| 3 | Headache gone even without aspirin | 0 | 0 |
| 4 | Headache gone only if no aspirin taken | 1 | 1 |

1.1. ATE, ATT, ATNT

In practice, we have multiple units ($i = 1\dots N$) and **we want to study aggregate causal effects**. Aggregate the individual causal effects to get:

- ATE(AverageTreatmentEffect):**

$$\text{ATE} = \mathbb{E}(Y_{1i} - Y_{0i}) = \mathbb{E}(Y_{1i}) - \mathbb{E}(Y_{0i})$$

for every individual one takes the difference in potential outcomes.

- ATT(AverageTreatmentEffectontheTreated):**

$$\text{ATT} = \mathbb{E}(Y_{1i} - Y_{0i} | D_i = 1) = \mathbb{E}(Y_{1i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 1)$$

takes the difference in potential outcomes only for treated individuals! (more infra: we will have to reconstruct just the counterfactual outcome when untreated y_0 from the control group)

- ATNT(AverageTreatmentEffectontheNon-Treated):**

$$\text{ATNT} = \mathbb{E}(Y_{1i} - Y_{0i} | D_i = 0) = \mathbb{E}(Y_{1i} | D_i = 0) - \mathbb{E}(Y_{0i} | D_i = 0)$$

This represents the average difference, for non-treated individuals, between their counterfactual outcome under treatment (unobserved) and their observed outcome under no treatment.

- If we could observe both Y_{1i} and Y_{0i} for all units:

$$\text{ATE} = \frac{1}{N} \sum_{i=1}^N (Y_{1i} - Y_{0i}), \quad \text{ATT} = \frac{1}{N_{D=1}} \sum_{i:D_i=1} (Y_{1i} - Y_{0i})$$

- We can also focus on more general functionals of potential outcomes (e.g., medians, percentiles, quantiles).

1.2. Fundamental problem of causal inference

Above we made a crucial, and unrealistic, underlying assumption: **for the same individual, we assumed to have both the outcomes when treated and when not treated**. BUT, in practice, **we only observe one of Y_{1i} or Y_{0i} for each unit**, this is **equivalent to missing data problem**. Hence, we cannot directly estimate the

75 individual Treatment Effect because we only see one potential outcome for each individual, and we cannot use the formulas above for
 76 the ATE and ATT. For example, for a treated individual we only ob-
 77 serve Y_{1i} , whereas Y_{0i} is the unobserved Counterfactual (what would
 78 have been the outcome if not treated).

80 1.3. Solving the fundamental problem of causal inference

- 81 We have a population of N units: $i=1,\dots,N$.
- 82 • **The problem:** we never observe both Y_{1i} and Y_{0i} for the same
 83 individual only one of them is observed depending on treat-
 84 ment status D_i .
- 85 • A common first approach is to compare group means:

$$\begin{aligned} \mathbb{E}(Y_i | D_i = 1) - \mathbb{E}(Y_i | D_i = 0) &= \mathbb{E}(Y_{1i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 0) \\ &= \underbrace{\mathbb{E}(Y_{1i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 1)}_{\text{ATT}} \\ &\quad + \underbrace{\mathbb{E}(Y_{0i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 0)}_{\text{Selection Bias}} \end{aligned}$$

- 87 - To get to the first step we used the switching equation: if $D_i = 1$
 → $Y_i = Y_{1i}$
- 88 - To get to the second step you added and subtracted $\mathbb{E}(Y_{0i} | D_i = 1)$
- 91 • The root of the problem is **Selection Bias**: the observed control
 92 group may not be a valid counterfactual for the treated group.

- 93 – **The fact is that you don't have the outcome the
 94 treated would have had hadn't been treated $\mathbb{E}(Y_{0i} | D_i = 1)$. You use as counterfactual the outcome of
 95 the control**, but if this does not correctly mimic the char-
 96 acteristics of the treated group you have selection bias.
- 97 – For example, treated individuals may have self-selected
 98 into treatment due to the fact they had higher benefits
 99 from treatment ($ATT > ATE$).

101 1.4. Relationship between the ATT and the ATE

What is the relation between ATT and ATE?

$$\begin{aligned} ATE &= \mathbb{E}(Y_{1i} - Y_{0i}) \\ &= \mathbb{E}[\mathbb{E}(Y_{1i} - Y_{0i} | D_i)] \quad (\text{by LIE}) \\ &= \mathbb{P}(D_i=1) \cdot \mathbb{E}(Y_{1i} - Y_{0i} | D_i=1) + \mathbb{P}(D_i=0) \cdot \mathbb{E}(Y_{1i} - Y_{0i} | D_i=0) \\ &= \gamma ATT + (1-\gamma)ATNT \end{aligned}$$

102 Comment: the ATE is a weighted average of the ATT and the
 103 ATNT with the weights (γ) being determined by the fraction of the
 104 population being treated.

105 Useful to express ATT in terms of ATE and treatment effect hetero-
 106 geneity ($ATT - ATNT$):

$$ATT = ATE + (1 - \gamma)(ATT - ATNT)$$

106 Comment: The formula above recomposes the ATT! So the TE on
 107 the Treated is already there, while the TE on the fraction of the non-
 108 Treated is corrected by adding the gap they have with respect to the
 109 non-Treated.

Hence now we can recognize that **our initial simple difference in
 107 means is equal to:**

$$\begin{aligned} \mathbb{E}(Y_i | D_i = 1) - \mathbb{E}(Y_i | D_i = 0) &= ATT + \text{SelectionBias} \\ &= ATE + \underbrace{(1 - \gamma)(ATT - ATNT)}_{\text{Treatment Effect heterogeneity}} + SelB \end{aligned}$$

110 Comment: Even with no selection bias, recovery of ATE requires
 111 equal treatment effects across treated and non-treated: $ATT =$
 112 $ATNT$. So we need both the absence of *a priori* and *a posteriori* dif-
 113 ferences! This is solved by randomization (for the same reasons why
 114 selection bias is solved by randomization), see below!

115 1.5. Random Assignment

To solve this problem, we need to impose the following crucial assumption:

- **Random Assignment:** treatment is independent of potential outcomes.

$$(Y_{1i}, Y_{0i}) \perp\!\!\!\perp D_i$$

- Here the idea is that **the control group has now the exact same composition as the treatment group**. So the baseline outcome y_0 of treated individuals will be the same as the baseline outcome y_0 of treated individuals:

$$\mathbb{E}(Y_{0i} | D_i = 1) = \mathbb{E}(Y_{0i} | D_i = 0)$$

- An (extra) logic step further: think again at the formulas above for the ATE and the ATT.

- Now for the ATE we will need to assign to each individual in the treatment group the counterfactual outcome y_0 extracted from a (similar) individual in the control group. And to each individual in the control group we will assign the counterfactual outcome y_1 extracted from a similar individual in the treatment group.
- for the ATT we will just consider the outcome y_1 of treated individuals and subtract their counterfactual outcome extracted from the outcome y_0 of (similar) individuals in the control group.
- for the ATNT we will just consider the outcome y_0 of control individuals and subtract it to the counterfactual outcome y_1 extracted from the outcome y_1 of (similar) individuals in the treatment group.

- **Random Assignment solves the selection problem** by guaranteeing:

$$D_i \perp (Y_{0i}, Y_{1i}) \Rightarrow \mathbb{E}(Y_{0i} | D_i = 1) = \mathbb{E}(Y_{0i} | D_i = 0)$$

So, the naive difference in means equals the true treatment effect:

$$\begin{aligned} \mathbb{E}(Y_i | D_i = 1) - \mathbb{E}(Y_i | D_i = 0) &= \mathbb{E}(Y_{1i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 1) \\ &= \mathbb{E}(Y_{1i} - Y_{0i} | D_i = 1) = ATT \end{aligned}$$

Moreover, given independence:

$$ATT = \mathbb{E}(Y_{1i} - Y_{0i} | D_i = 1) = \mathbb{E}(Y_{1i} - Y_{0i}) = ATE$$

Example

- $D_i = 1$ if individual has a bank account. Y_i : total savings.
- If we compare individuals with and without a bank account:

$$\mathbb{E}(Y_i | D_i = 1) - \mathbb{E}(Y_i | D_i = 0).$$

- We will not get the ATT (the causal effect of interest) if

$$\mathbb{E}(Y_{0i} | D_i = 1) - \mathbb{E}(Y_{0i} | D_i = 0) > 0.$$

- **Selection bias:** those with a bank account ($D_i = 1$) have a priori higher savings than those who do not, and would have had more savings even without a bank account (Y_{0i}) than those without one ($D_i = 0$).
- As $\mathbb{E}(Y_{0i} | D_i = 1)$ is unobservable we need assumptions or methods (e.g. matching, IV, RCT) to ensure:

$$\mathbb{E}(Y_{0i} | D_i = 1) = \mathbb{E}(Y_{0i} | D_i = 0).$$

More on Random assignment

- The key requirement is that Treated and Control units must be comparable. Hence, understanding the **assignment mechanism** i.e., the rule that determines treatment status is fundamental.

- In most observational settings, the assignment mechanism is unknown or uncontrolled. An important exception is the **Randomized Control Trial (RCT)**, which ensures a well-behaved assignment via randomization. RCTs satisfy the following properties:

- **Unconfounded Assignment:** the treatment is independent of potential outcomes.

- 150 – **Probabilistic Assignment:** every unit has a positive
151 probability of receiving either treatment status (no deter-
152 ministic assignment).
153 – **Individualistic Assignment:** the treatment probability
154 of one unit does not depend on the characteristics or out-
155 comes of other units.

156 2. Statistical Inference

157 Now that you have run your experiment, is time to run statistical
158 inference!

159 2.1. Neyman Inference

- 160 • We use statistical inference to calculate standard errors, confi-
161 dence intervals, and conduct hypothesis testing.
- 162 • Standard approach: **Neyman's repeated sampling frame-**
163 **work**, where:
 - 164 – Sampling is assumed to be drawn from an infinite **super-**
165 **population**.
 - 166 – **Uncertainty** (variance) arises from this sampling process
167 (no uncertainty if we had the full population). The concept
168 is that the **true parameter is fixed**, and the **uncertainty**
169 we face is driven by the fact that we are analysing a
170 **subsample from a population**.
 - 171 • To derive standard errors, we rely on the **CLT** and con-
172 struct confidence intervals using the normal distribution
173 (for n large enough the sampling distribution of each sample mean
174 (mean of T vs mean of C) converges to a normal distribution!)

175 2.2. Randomization Inference

- 176 • Idea: Different perspective, uncertainty comes from the fact
177 that we only observe one assignment realization, while many
178 others were possible.
- 179 • Randomization inference becomes possible under the **sharp**
180 **null hypothesis**:

$$H_0 : Y_{1i} = Y_{0i} \quad \text{for all } i$$

181 That is, treatment has **no effect on any unit**. Hence, under this
182 hypothesis, each unit's outcome is invariant to the treatment
183 status:

- 184 • **Simulate** counterfactual scenarios by:

- 185 1. **Plot their distribution of the observed outcomes**
186 Y_1, \dots, Y_N in the particular treatment realization you had.
- 187 2. **Build the counterfactual**
 - 188 – 2a. **Re-randomize the treatment assignments** D_i
189 according to the known experimental design (e.g., ran-
190 domization with fixed number treated)
 - 191 – 2b. **Recompute the test statistic** (e.g., difference in
192 means) under each simulated assignment;
- 193 3. **Building the randomization distribution** of the statis-
194 tic under the null;
- 195 4. **Comparing** the observed test statistic (in 1.) to this null
196 distribution (in 3.) to compute the p -value.

| Unit i | Observed D_i | Observed Y_i | D_i^* | Y_i^* | D_i^{**} | Y_i^{**} |
|--|----------------|----------------|--------------------|----------------|--------------------|------------|
| 1 | 1 | 4 | 0 | 4 | 0 | 4 |
| 2 | 1 | 6 | 1 | 6 | 0 | 6 |
| 3 | 0 | 2 | 1 | 2 | 1 | 2 |
| 4 | 0 | 5 | 0 | 5 | 1 | 5 |
| 5 | 1 | 3 | 1 | 3 | 1 | 3 |
| 6 | 0 | 1 | 0 | 1 | 0 | 1 |
| Means (Observed D_i) | | Treated: 1,2,5 | $\bar{Y}_T = 4.33$ | Control: 3,4,6 | $\bar{Y}_C = 2.67$ | |
| | | | | ATE = +1.67 | | |
| Means (D_i^*) | | Treated: 2,3,5 | $\bar{Y}_T = 3.67$ | Control: 1,4,6 | $\bar{Y}_C = 3.33$ | |
| | | | | ATE = +0.33 | | |
| Means (D_i^{**}) | | Treated: 3,4,5 | $\bar{Y}_T = 3.33$ | Control: 1,2,6 | $\bar{Y}_C = 3.67$ | |
| | | | | ATE = -0.33 | | |

Remark: what changes from one column to the other is not the outcomes but the composition of the treated and control groups.

- 195 • **Key Insight:** Under the sharp null H_0 , outcomes are invariant to treatment, so we
196 can treat them as known constants and generate the entire distribution of the test
197 statistic using only variation in D_i . So the following is possible because we assume
198 we have the same counterfactual outcomes and we are just redrawing treatment.
199 • This makes randomization inference a **nonparametric** (not assuming any para-
200 metric distribution for the outcomes (e.g., normality, linearity)) method for hy-
201 pothesis testing, justified by our sharp null hypothesis.
202 • In **small samples**, randomization inference is often preferred for valid exact p -
203 values, or at least should be reported alongside parametric results to show robust-
204 ness.

205 3. Regression Analysis of Experiments

- 206 • We can estimate treatment effects either:
207 – Non-parametrically: by comparing conditional means (up
208 to now);
209 – Parametrically: using a regression framework.
- 210 • Define the observed outcome as:

$$Y_i = Y_{0i} + (Y_{1i} - Y_{0i})D_i$$

- 211 • Assume a constant treatment effect: $\rho = Y_{1i} - Y_{0i}$ and define:

$$\alpha = \mathbb{E}(Y_{0i}), \quad \eta_i = Y_{0i} - \mathbb{E}(Y_{0i})$$

212 take the expected value and define η_i in a way to center the re-
213 gression residuals around 0 leading to the regression model:

$$Y_i = \alpha + \rho D_i + \eta_i$$

- 214 • From this regression, the expected difference in outcomes be-
215 tween treated and control is:

$$\mathbb{E}(Y_i | D_i = 1) - \mathbb{E}(Y_i | D_i = 0) = \rho + \mathbb{E}(\eta_i | D_i = 1) - \mathbb{E}(\eta_i | D_i = 0)$$

- 216 • This shows: **Selection Bias = correlation between D_i and η_i** . Another (microeconomic) way to frame the problem of endogeneity! Imagine: OVB creates Selection Bias
217 • only if the assumption above holds we can interpret ρ as the causal effect of D on Y . $\rho = ATT = ATE_y$
218 • In correctly implemented randomized experiments, this condi-
219 tion holds **by design**, so OLS recovers the causal effect.

220 3.1. Do you need controls in your RCT regression?

- 221 • **In theory:** Randomization ensures independence be-
222 tween treatment assignment and any covariate. So:

- 223 – The treatment should be independent of any control vari-
224 able, thus it is not necessary to include controls in order to
225 get an unbiased estimate for the treatment effect (think of
226 condition for omitted variable bias: treatment is now inde-
227 pendent from any other covariate!);

- 228 • **In practice:** You must check whether randomization actu-
229 ally balanced the covariates. To verify:

- 230 1. **Conduct balance tests:** compare pre-treatment covariate
231 means in treatment vs. control groups. *Remark: expect
232 differences to be statistically significant at 5% level in 5% of
233 cases (1 out of 20 variables) due to false positives.*

- 234 2. **Run a joint F-test** to check if covariates jointly predict
235 treatment.

- 236 3. **Check whether including controls changes the coefficient on D .** It shouldn't under successful randomization.

- 237 4. **Robustness Rule:** You should always include controls
238 for unbalanced variables to absorb bias and improve
239 credibility.

- 240 • **Why include controls if there is no imbalance?**

- 241 – If covariates help predict Y , they reduce residual variance:

$$\text{If } \sigma_{Y|D,X}^2 < \sigma_{Y|D}^2, \text{ then precision increases.}$$

- 241 – But in general after the first few covariates, gains in pre-
 242 cision diminish as those covariates are weakly related to
 243 Y .

244 **Remark: Never include controls that might be affected by**
 245 **the treatment.** They would capture part of the effect of interest
 246 and including them would bias the results. Eg effect of being in
 247 the back of the class on grades. You are in the back of the class,
 248 you feel less pressure by prof, you use more the pc, you get lower
 249 grade. If you control for pc you stop the transition channel (will
 250 anlaysis individuals at the same level of pc use!).

251 3.2. Interacting Covariates with the Treatment Dummy

- 252 • In addition to including covariates linearly, we can interact
 253 them with the treatment dummy D_i .
 254 • This is useful when we expect the association between covari-
 255 ates and the outcome to vary by treatment status (i.e., heteroge-
 256 neous covariate effects).
 257 • These interactions allow the covariate slopes to differ across
 258 treatment and control, potentially improving the precision of
 259 the treatment effect estimate as well.
 260 • A useful specification is:

$$Y_i = \alpha + \rho D_i + X_i \beta + D_i(X_i - \bar{X})\gamma + \eta_i$$

261 where:

- 262 – X_i is a vector of covariates,
- 263 – $(X_i - \bar{X})$ is mean-centered (demeaned) to ensure ρ is inter-
 264 pretable as the average treatment effect (ATE)¹,
- 265 – $D_i(X_i - \bar{X})$ allows covariate effects to vary with treatment.

266 **1 Remark:** If we use the raw (non-demeaned) interaction $D_i X_i$,
 267 the coefficient $\hat{\rho}$ reflects the treatment effect only for the special
 268 case when $X_i = 0$. In this case, to recover the ATE we would
 269 need to adjust $\hat{\rho}$ by adding the average effect of the interaction
 270 term: $\text{ATE} = \hat{\rho} + \hat{\gamma} \cdot \bar{X}$. By demeaning X_i , the average of the
 271 interaction term $D_i(X_i - \bar{X})$ is zero by construction, so $\hat{\rho}$ directly
 estimates the ATE without further correction.

272 3.3. Randomization Techniques

273 Randomization is a critical component of experimental design. It
 274 ensures that treatment assignment is independent of potential out-
 275 comes and covariates.

276 3.3.1. How do we actually randomize?

- 277 • Randomization can be conducted **privately** (e.g., via random
 278 number generators) **or publicly** (e.g., via lotteries). Trade-offs,
 279 e.g., lotteries, are better to create trust, but they let people know
 280 they are treated.
 281 • Several techniques are available to improve covariate balance:

- 282 1. **Pure Randomization (Single Draw):** Assign units ran-
 283 domly. By chance, imbalances may occur (e.g., 1 in 10 con-
 284 variates may differ at the 10% significance level).
- 285 2. **Stratification (Blocking):** Divide units into strata based
 286 on key covariates (e.g., baseline outcomes, gender, location),
 287 then randomize within each stratum. This is to in-
 288 crease the likelihood covariates are balanced (10 females-
 289 90 males, randomizing globally could have 5-45. Random-
 290 ize within stratum: sample 5 on 10 and then 45 on 90).
 291 Issue: if you create too many subcategories you don't have
 292 enough T and C for each subgroup.
- 293 3. **Pair-wise Matching:** Form pairs based on similarity in
 294 covariates; assign one unit in each pair to treatment and
 295 the other to control.
- 296 4. **Re-randomization methods:** Repeat randomization
 297 many times and select the one with the best covariate bal-
 298 ance:

- 299 – Redraw if any covariate shows imbalance (e.g., $p <$
 300 0.05).
 301 – Choose the draw minimizing (between) the worst t-
 302 statistic from regressors (within) on treatment.

303 Lost internal validity: your sample is no longer a random
 304 draw from the population \Rightarrow No longer Neyman inference,
 305 need adjustments.

306 3.3.2. Which randomization method is better?

- 307 • Athey and Imbens (2016) recommend **stratified randomiza-**
 308 **tion** as superior, especially when covariates are predictive of
 309 outcomes (i.e., we really need balance around those covariates).
 310 • Stratification should be done *ex-ante* to improve balance, rather
 311 than relying on regression adjustment *ex-post*.
 312 • If:
 - 313 – Covariates are weak predictors of outcomes.
 - 314 – Sample size is large (> 300)

315 Then balancing on them does not matter much (quite obvious)

316 3.3.3. Randomization recommendations (Bruhn & McKenzie, 2009):

- 317 • The method of randomization affects which covariates to con-
 318 trol for later: there is no true variability in sample draws; the
 319 sample is synthetic and we need to fix se. To do this, we add
 320 controls to resort to within-group comparisons.
 - 321 – *Stratification*: include strata dummies in regressions (not
 322 for bias, actually we stratify just to give balance).
 - 323 – *Pair-wise matching*: include matching variables in reg.
 - 324 – *Re-randomization*: include all variables for which balance
 325 was checked to decide whether to re-randomize.

326 3.3.4. Clustered Randomized Experiments

- 327 • When spillovers are possible (e.g., peer effects), it is better to
 328 randomize at a more aggregate level, i.e. the **cluster level**.
 329 • Instead of assigning treatment randomly to units within a clus-
 330 ter, we assign treatment to entire clusters. All units within a
 331 cluster receive the same treatment.
 332 • Then, econometric analysis can be done:
 - 333 – At the cluster level (preferred for inference).
 - 334 – At the individual level, using cluster-robust standard er-
 335 rors.

336 3.4. Power Computation

337 How large a sample do we need in order to detect a meaningful effect
 338 with high probability? That is, how much will I be able to learn from
 339 this experiment? Power computation! Suppose we aim to estimate
 340 the effect of a randomly assigned treatment D_i on an outcome Y_i
 341 using:

$$Y_i = \alpha + \delta D_i + \eta_i,$$

342 Power computation

- 343 • **Power:** If the true effect is δ , what is the probability of getting
 344 an estimate big enough to reject $\delta = 0$ (reject when false)?
 345 • To achieve a power of $1 - \beta$ at a confidence level α , we need the
 346 true effect to be such as (see image below):

$$\delta / SE(\hat{\delta}) > (t_\beta + t_{\alpha/2})$$

- 347 • Then, the minimum detectable effect (MDE) is:

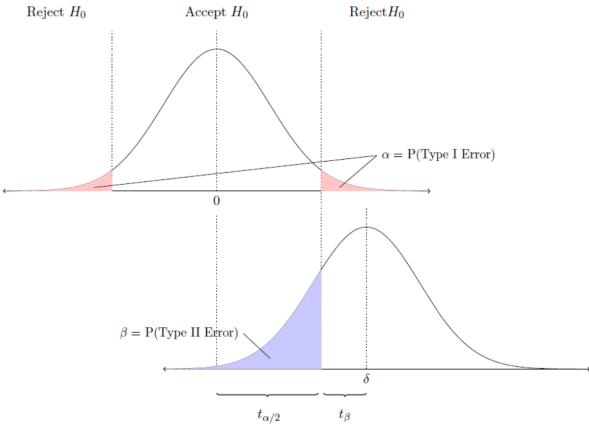
$$\delta_{MDE} = (t_\beta + t_{\alpha/2}) \cdot SE(\hat{\delta}) = (t_\beta + t_{\alpha/2}) \cdot \sqrt{\frac{\sigma^2}{P(1-P)N}}.$$

348 This tells us that we will be able to detect a smaller effect size if:

- 349 – Power increases with sample size N ,

- 349 – Power is maximized when $P = 0.5$ (equal-sized treatment
350 and control),
351 – Power increases when error variance σ^2 is smaller. High
352 error variance means unexplained noise in the data.

353 **Remark:** why $\delta/SE(\delta) > (t_\beta + t_{\alpha/2})$



- 354 • Upper panel assumes H_0 is true (ATE = 0).
355 – To control the Type I error at a rate of α , we reject H_0 only
356 if the absolute value of the observed t-statistic is equal to
357 or larger than the critical value $t_{\alpha/2}$.
358 • Bottom panel assumes H_0 is false (ATE = δ).
359 – The power to detect an effect of size δ is the fraction of the
360 area under the distribution that falls to the right of $t_{\alpha/2}$.
361 – In that region we correctly reject the null hypothesis.
362 – We limit the Type II error to some $\beta > \alpha$ (and then power
363 is $1 - \beta$), such β will be characterized by some t_β .
364 – the δ plotted in the bottom panel is exactly the MDE! The
365 MDE is the smallest value for which we can correctly reject
366 H_0 with probability $1 - \beta$ at a significance level α .

3.4.1. Interpretation of Minimum Detectable Effect (MDE)

- 368 • If we aim at power $1 - \beta = 0.8$ with confidence level $\alpha = 0.05$
• Then, the minimum detectable effect is:

$$\delta_{MDE} = 2.8 \times \frac{\sigma}{\sqrt{P(1-P)N}}$$

- 369 • If $P = 0.5$, and $N = 1000$, we can detect an effect of 0.18σ , i.e.,
370 0.18 standard deviations of the outcome.
371 • Is 0.18 standard deviations a meaningful effect size? It depends
372 on the ratio between the mean and standard deviation of the
373 outcome!
374 • Example 1: If the ratio is 1 (e.g., mean and sd of employment
375 rate are both 0.5), the effect is 18% of the mean (or 9 percentage
376 points), which seems large.
377 • Example 2: If the ratio is 2 (e.g., mean = 100, sd = 50), then
378 $0.18\sigma = 9$, which is only 9% of the mean. This may not seem
379 like a large effect.

3.4.2. Required Sample Size:

381 we want to know what is the required sample size to detect a pre-
382 specified effect δ_0 , we rearrange the MDE formula:

$$N > (t_\beta + t_{\alpha/2})^2 \cdot \frac{\sigma^2}{\delta_0^2 P(1-P)}.$$

383 Take $(t_\beta + t_{\alpha/2})^2 = 2.8^2$, $P = 0.5$, $1 - \beta = 0.8$, if the effect is of 2
384 dollars with variance $\sigma^2 = 100$:

$$N > (2.8)^2 \cdot \frac{100}{4 \cdot 0.5 \cdot 0.5} = 784.$$

385 **Interpretation:** We need at least 784 observations to detect a 2
386 dollar effect with 80% power at the 5% significance level.

4. Internal Validity and Its Threats

- **Internal Validity:** ability of a study to estimate causal effects within the study sample.
- **External Validity:** ability to generalize results beyond the study sample.
- To maximize both, randomization should ideally be performed in two stages:
 1. Take a random sample from the population of interest (maximize external validity).
 2. Randomly allocate part of the sample to the treatment group and part to the control group (guarantees internal validity).

4.1. Main Threats to Internal Validity

An experiment eliminates the most important threats to internal validity in non-experimental settings (e.g. omitted variable bias). Yet, randomization is not a panacea.

- **Attrition:** Some individuals drop out of the sample after assignment.
- **Partial Compliance:** Some treated individuals do not receive treatment, or some controls receive it.
- **Externalities:** Spillovers across units (e.g., treated individuals affect control individuals).
- **Experimental Effects:** Subjects may change behavior due to being observed (Hawthorne/Henry effects).
- **Data Mining:** Testing many hypotheses without proper correction can lead to spurious findings.

These are also threats to non-experimental settings, except maybe experimental effects.

4.1.1. Attrition and Its Implications

A. Intuition

- If **random (unrelated to treatment)**: only reduces statistical power.
- If **correlated with treatment or outcomes** (e.g., low-benefit individuals drop out), introduces **attrition bias**.

B. See it in the model

- R_{1i}, R_{0i} : potential response indicators if the same individual is treated or controlled.
- switching eq: $R_i = R_{0i} + D_i(R_{1i} - R_{0i})$
- Y_i is only observed if $R_i = 1$. Meaning that, clearly, I can observe outcomes only of those who respond.
- Suppose $D_i \perp (Y_{0i}, Y_{1i}, R_{0i}, R_{1i})$. This means treatment is randomly assigned and independent of both potential outcomes and potential response behavior.
- However, the above could actually be violated. Outcomes and response can still be correlated: the issue arises *ex post*, after treatment is assigned. That is, the treatment itself may not affect average response rates, but eg:

1. Individuals assigned to the control group may choose not to respond; and
2. The probability of responding may depend on potential outcomes (e.g., people with lower Y_{0i} are less likely to respond if assigned to control).

Even though treatment is randomly assigned, the **composition of respondents across groups may differ**, leading to biased comparisons unless response is also independent of outcomes.

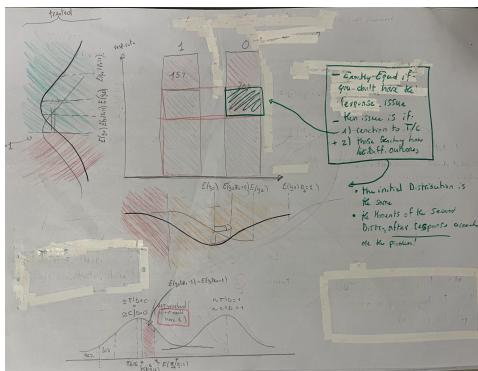
- 442 • What you observe: $\Delta = \mathbb{E}[Y_i | D_i = 1, R_i = 1] - \mathbb{E}[Y_i | D_i = 0, R_i = 1]$. Then by the switching equation:

$$\Delta = \underbrace{\mathbb{E}[Y_{1i} | R_{1i} = 1]}_{\text{Treated who responds}} - \underbrace{\mathbb{E}[Y_{0i} | R_{0i} = 1]}_{\text{ATT on respondents when treated}}$$

$$= \underbrace{\mathbb{E}[Y_{1i} | R_{1i} = 1] - \mathbb{E}[Y_{0i} | R_{1i} = 1]}_{\text{Differential response across treat arms bias}} + (\underbrace{\mathbb{E}[Y_{0i} | R_{1i} = 1] - \mathbb{E}[Y_{0i} | R_{0i} = 1]}_{})$$

444 where second equation adds and subtract that term to give economic meaning. Note we have assumed the experimtn has
445 been successfull, so we have no selection bias. **Two problems:**
446

- 447 1. $\underbrace{\mathbb{E}[Y_{1i} - Y_{0i}]}_{\text{Overall Outcome}} \neq \underbrace{\mathbb{E}[Y_{1i} - Y_{0i} | R_{1i} = 1]}_{\text{outcome just for respondents}}$ threatens External
448 validity: not generalizable to the whole population, as
449 these are only those who responded.
- 450 2. Treatment induces different units to respond: $\mathbb{E}[Y_{0i} | R_{1i} = 1] \neq \mathbb{E}[Y_{0i} | R_{0i} = 1]$, Internal validity: I may be estimating
451 the wrong treatment effect if the baseline outcome of treated respondent is different from the baseline outcome of control respondents. So in the Control you have eg
452 less respondents and those who are responding less are individuals with lower outcomes (would create downwards
453 bias).



C. Reasons for attrition

- 458 • Question not clear for $T = 0$, they donot understand, what if
459 those who don't understand have different outocmes? (non-natives?)
- 460 • $t = 1$ has more up to date contact (updated the email)
- 461 • unhappy to be in the control group!

464 *note always need systematically different outcomes (or is it just an
465 issue of power)?

D. Suggestions

- 467 1. **Present the attrition rate** in control vs treatment.
- 468 2. **Compare baseline characteristics** of attritors vs non-attritors.
- 469 3. **Run a regression:**

$$\text{Attritor}_i = \alpha + \beta D_i + \gamma X_i + \delta(D_i \cdot X_i) + \varepsilon_i$$

470 and test significance:

- 471 • Coefficient on treatment evidence of differential attrition
472 between treatment arms.
- 473 • Interaction Tests whether the relationship between covariates
474 and attrition differs by treatment status. That is: do certain
475 subgroups drop out more in treatment than in control?

477 4. If attrition is related to treatment, apply:

- 478 • Parametric methods: model attrition using covariates
479 (control for).

- 480 • Non-parametric methods: **bounding approaches**. However note: Bounds are going to be uninformative if attrition is very large.

E. Bounding Methods: Manski Bounds (Binary Outcomes: eg employed 483 484 vs unemployed 0)

- 485 • facile: la riposta è 1 (employed) o 0 (unemployed). ti
486 hanno risposto solo l'80% ed hai una media di 0.4. LB: assumi
487 il restante 20% abbia tutta 0. UB: oppure tutto 1.
- 488 • focus on the effects on the treated for the moment
- 489 • Your expected valye is, reasonably, a weighted average of the outcome
490 of the reposndents and the outcome of non reponsnts.
491 $\mathbb{E}[R | D = 1]$ is the share of respondents among the treated:

$$\mathbb{E}[Y | D = 1] = \mathbb{E}[R | D = 1] \mathbb{E}[Y | D = 1, R = 1] \quad (1)$$

- 492 • $\mathbb{E}[Y | D = 1, R = 0]$ is **unobserved**, but bounded: $0 \leq \mathbb{E}[Y | D = 1, R = 0] \leq 1$
- 493 • Thus, bound $\mathbb{E}[Y | D = 1]$. So trivial: see the expression [1]
494 above, the dummy takes 0,1. So worst case scenario all the un-
495 observed share is multiplied by an expectaiton of y equal to 0
496 (LB), best case sceanrio all the unobserved share is multiplied
497 by an expectation of y equal to 1.
- 498 • Provides ATE bounds. ATE no longer point identified but set
499 identified.

Bounding Methods: Lee Bounds (for Continuous or Unbounded Outcomes)

- 501 • Suppose treatment increases the probability of response (sup-
502 pose response rate is 10% higher in the treated group). Then
503 the **treated sample may include people who would have
504 not responded if assigned to control**. So the treated sample
505 includes people for which you do not have a valid counterfac-
506 tual in the control group.
- 507 • Lee bounds(non-binary): trim the treatment group to only
508 those who would have responded in the control group.
- 509 • Two assumptions:

- 510 – Extra 10% are those with **highest** 5% and **lowest** 5% out-
511 comes: drop top (bottom) 5% from treated. ATT:

- 512 • We take them out from the treated mean to set-identify $\mathbb{E}[Y_1 | R_0 = 1]$

$$\mathbb{E}[Y_1 | R_0 = 1] \in [\mathbb{E}[Y | D = 1, R = 1, Y \leq p_{95}], \mathbb{E}[Y | D = 1, R = 1, Y \geq p_{5}]$$

- 513 • Look at the CI above, the more you trim (the more the differ-
514 ence in response rate or the larger the % you trim) the larger
515 the CI (the first CI element is the lower bound, because it is the
516 expected value with the top 5% out, the second element is the
517 upper bound, because you are removing the worst 5%). So the
518 idea is that the lower bound gets lower and the upper bound
519 gets larger.

520 Remark: another option could have been to predict the proba-
521 bility of response on covariates, etc.

4.1.2. Partial compliance

Definition

- 523 • When only a fraction of the individuals who are offered the treat-
524 ment take it up.
- 525 • Or, when some members of the control group manage to get the
526 treatment.
- 527 • Default case in "encouragement" designs, individuals randomly
528 encouraged to take treatment. Randomization affects the proba-
529 bility of getting treatment, not the treatment itself.

532 **Solutions**

- 533 • Focus on the *Intention-To-Treat* (ITT) effect:
 - 534 – Compare all treated *assigned* to treatment to all controls
535 *assigned* to no-treatment.
 - 536 – Effect of being offered treatment.
- 537 • Try to estimate the ATE instead of the ITT?
 - 538 – We will see that we can use the dummy for being offered
539 treatment as an Instrumental Variable (IV) for treatment,
540 but retrieves LATE.
 - 541 • We can use bound analysis (partial identification) to provide
542 bounds for the ATE (assumptions of TE on Compliers and NT:
543 e.g., you estimate proportions from actual take-up rates and
544 then you impose assumptions on the outcomes e.g., Manski
545 bounds).

546 **How does Partial Compliance affect power?**

- 547 • Assume that only a fraction c of the treatment group receives
548 the treatment, and a fraction d of the control group receives it.
- 549 • Then:

$$\delta_{MDE} = 2.8 \cdot \frac{1}{c-d} \cdot \frac{\sigma}{\sqrt{P(1-P)N}}$$
- 550 • With $P = 0.5$, and $N = 1000$, with $c - d = 1$, we can detect an
551 effect of 0.18σ .
- 552 • If $c - d = 0.5$, we can detect an effect of only 0.36σ .
- 553 • To ensure the same MDE, the sample size under 50% partial
554 compliance would be almost 4 times larger than the sample under
555 **full compliance (no linear relationship!)**. Why? recall
556 you have to put N out of the sq root.

557 **4.1.3. Externalities / spillovers**

Individuals in the control group might be affected by the treatment.
Examples: vaccines, imitation effects.

SUTVA is violated: the potential outcome of each individual depends on the vector of allocations to treatment and comparison groups.

- 562 • One solution: randomize at a level that allows capturing externalities (e.g., school level instead of individual level).
- 563 • Smart design used leveraging the violation of the SUTVA: compare adoption of fertilizer among the friends of treated farmers to that of the friends of control farmers.
- 564 • Actually could be interesting for research: to study peer effects, randomly assign individuals to different peer groups (e.g., Moving to Opportunity experiment in the US).

570 **4.1.4. Experimental Effects (intrinsic to experiments)**

The evaluation itself may lead individuals to change their behavior: experimenter effects. Example: farmers working harder in the experimental plot with randomly allocated fertilizer than in the rest of their farm.

- 575 • If the behavioral change is in the treatment group, we call it Hawthorne Effect (e.g., observed workers exhibit higher productivity).
- 576 • If in the control group: John Henry Effect. (Legendary American Steel driver, when he heard his output was being compared with that of a steam drill, he worked harder to outperform the machine; he died in the process).
- 577 • funny to try PLACEBO TESTS!
- 578 • how to fight these effects? justify saying they are not very strong, OR INDUCE them and bound their strength (eg through placebo!).

4.1.5. Data Mining: P-hacking when Studying Heterogeneity of Impact

- 586 • Underlying Issue: As for balancing: if you estimate many treatment effects (and interactions) you may have significance by chance (false positives)... and then you report just the significant!
- 588 • If we estimate the treatment effects for many different subgroups we need to **correct p-values for multiple testing**.
- 590 • Extra Good Practices:
 - 592 – write **pre-analysys plan**
 - 594 – Stratify before randomizing: be sure there are no imbalances across characteristics in T and C group. Pre-randomization stratification ensures balance in key subgroups, thus Increasing power to detect true heterogeneity.
 - 596 – Avoid complex modeling to recover unbiased ATE. MAKE ANALYSS SIMPLER AND MORE CREDIBLE Recall: One argument in favor of randomization is that everyone will get the same result and the researcher is bounded by the ex-ante design (no need to choose method, variables, specification).
 - 598 – study heterogeneity with machine learning causal forest. Non parametric (no linear form, normality etc): 1) flexible (no functional form assumed for treatemtn efetc), 2) algoritmic discovery (he does the job, no manipulation), data drievn learning

4.1.6. External Validity

- 611 • Concerns about external validity are more relevant when there is significant **treatment effect heterogeneity** (cross-country analysis).
- 612 • **Generalization / Extrapolation** of results may not be reliable if:
 - 614 – **Issue of scalability**
 - 616 – It is difficult to **extrapolate results to other populations** (e.g., will poor farmers in Uganda respond like those in South America?).
 - 618 – One approach to improve external validity is to **accumulate evidence across different settings and locations**.
 - 620 – Another approach is to **account for differences in the distribution of characteristics** across settings that may drive treatment effect heterogeneity. Effect: So you see which characteristics are driving the treatment effect heterogeneity, then you can try to successfully predict the ATE on a training sample of which you know the values for those characteristics + the ATE (works)
 - 622 – **Behavioral theory** can guide whether and how results should be extrapolated to new contexts.
 - 624 – In general, randomized evaluations cannot capture **general equilibrium effects**. These effects may be observed if the unit of observation is large enough (e.g., country, state, or city).

5. Natural Experiments

5.1. Idea

- 636 • The exogenous source of variation comes from a randomly assigned variable that can be used as an instrument for treatment.
- 637 • Classic examples: Vietnam draft lottery to study military service effects, Mariel boatlift for immigration impacts, changes in divorce laws, etc.
- 638 • These settings mimic random assignment, enabling causal inference without explicit experimental control.
- 639 • **Threats to internal validity** increase:
 - 641 – policy endogeneity (to an area (unobervabels)-> affects treatemtn and outcome),

- 647 – omitted variable bias (Suppose we estimate the effect of
648 a health insurance reform but do not control for baseline
649 regional health infrastructure),
650 – outcome trends (Did with wrong pretrends),
651 – simultaneity (A government might increase subsidies for
652 schools in response to declining test scores),
653 – sample selection (if only a non-random subset of the pop-
654 ulation is observed or affected by the treatment, estimates
655 will not be representative or causal.).

- 656 • Similar challenges as RCTs regarding **external validity**, with
657 one key difference: natural experiments typically do not require
658 informed consent.
659 • Rosenzweig and Wolpin (2000) emphasize *natural natural experiments* those that rely solely on natural exogenous events,
660 not policy changes or institutional assignments.

662 **5.2. Debate**

- 663 • Natural experiments and quasi-experimental methods are cen-
664 tral to the so-called **credibility revolution** in microeconomics.
665 • Critics argue that **overenthusiasm** for these methods can lead
666 to oversimplified claims e.g., reliance on single-equation mod-
667 els with robust standard errors without deeper understanding.