

Big Data for Public Policy

Applied Micro Methods I

Sergio Galletta

ETHZ Zurich

09/03/2023

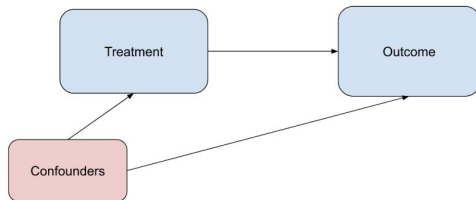
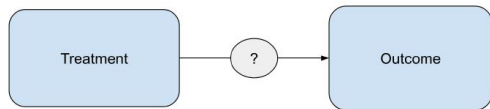
Empirical methods

- ▶ While economists are often motivated by **why** questions, in research, we proceed to address **what if** questions
- ▶ This is because we are typically interested in estimating a **causal effect** (if any) of a “**treatment**” on an “**outcome**”, **but is not easy!!!**

Empirical methods

- ▶ While economists are often motivated by **why** questions, in research, we proceed to address **what if** questions
- ▶ This is because we are typically interested in estimating a **causal effect** (if any) of a “**treatment**” on an “**outcome**”, **but is not easy!!!**
- ▶ Examples:
 - ▶ How does taking this course affect the grade in your master thesis?
 - ▶ This is different from the predictive question: “What is the grade that students taking this course will obtain with their master thesis?”
 - ▶ If Zurich imposed a special tax on Uber drivers, how would that affect the supply of Uber rides?

Empirical methods



A formal framework (Neyman-Rubin Causal model)

The outcome of interest is denoted by $Y_i(D_i)$, the effect that we want to attribute to the treatment. The notation indicates that it may depend on D_i

- ▶ $Y_i(1) = \text{if } D_i = 1$
- ▶ $Y_i(0) = \text{if } D_i = 0$

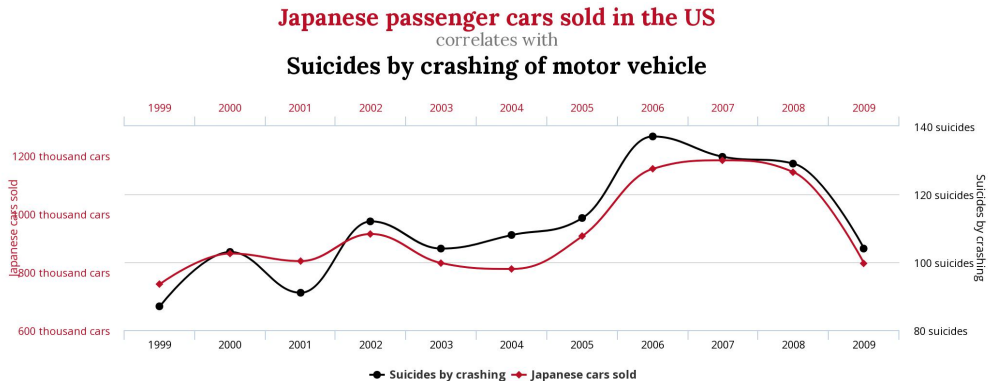
The outcome for each individual i can be written as:

- ▶
$$Y_i(D_i) = D_i Y_i(1) + (1 - D_i) Y_i(0)$$

This is enough for correlation: $\text{Cov}(D_i, Y_i) / \text{var}(D_i)$

But does this imply **causality**?

Correlation does not imply causation



tylervigen.com

The fundamental problem of causal inference

For every individual i , the event $\{D_i = 1 \text{ instead of } D_i = 0\}$ causes the effect $\Delta_i = Y_i(1) - Y_i(0)$

The fundamental problem of causal inference

For every individual i , the event $\{D_i = 1 \text{ instead of } D_i = 0\}$ causes the effect $\Delta_i = Y_i(1) - Y_i(0)$

The “Fundamental Problem of Causal Inference”

- ▶ It is impossible to observe for the same individual i the values $D_i = 1$ and $D_i = 0$ as well as the values $Y_i(1)$ and $Y_i(0)$ and, therefore, it is impossible to observe the effect of D on Y for unit i (Holland, 1986)
- ▶ Another way to express this problem is to say that we cannot infer the effect of a treatment because we do not have the counterfactual evidence

The fundamental problem of causal inference

i	D_i	Y_i	$Y_i(1)$	$Y_i(0)$	$Y_i(1) - Y_i(0)$
1	1	0	0	?	?
2	0	1	?	1	?
3	1	0	0	?	?
4	1	1	1	?	?
5	0	0	?	0	?
6	1	0	0	?	?
7	0	1	?	1	?

Causal estimands

- ▶ We could approach the problem by focusing on the **Average Treatment Effect** (ATE) for the entire population

$$E\{\Delta_i\} = E\{Y_i(1) - Y_i(0)\} = E\{Y_i(1)\} - E\{Y_i(0)\}$$

- ▶ Alternatively, one could focus on the **Average Treatment Effect on the Treated** (ATT):

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

Statistical solution

- ▶ ATE, and ATT are not identified as we cannot compute the expectations on the right-hand side (because of the missing data)
- ▶ A comparison of output by observed treatment status gives a biased estimate of the ATT

Statistical solution

- ▶ ATE, and ATT are not identified as we cannot compute the expectations on the right-hand side (because of the missing data)
- ▶ A comparison of output by observed treatment status gives a biased estimate of the ATT

$$\begin{aligned} & E\{Y_i|D_i = 1\} - E\{Y_i|D_i = 0\} \\ &= E\{Y_i(1)|D_i = 1\} - E\{Y_i(0)|D_i = 0\} \\ &= \underbrace{E\{Y_i(1)|D_i = 1\} - E\{Y_i(0)|D_i = 1\}}_{\text{Treatment Effect on Treated}} + \underbrace{E\{Y_i(0)|D_i = 1\} - E\{Y_i(0)|D_i = 0\}}_{\text{"Selection Bias"}} \end{aligned}$$

Statistical solution

- ▶ The observed difference in treatment status adds to this causal effect a term called **selection bias**
- ▶ This selection bias term is the difference in average $Y_i(0)$ between those who were and those who were not treated
- ▶ Alternatively, selection bias is the phenomenon that the distribution of the observed group is not representative to the group we are interested in

Randomized experiments

- ▶ Randomization solves this problem!

Randomized experiments

- ▶ Randomization solves this problem!
- ▶ Random assignment makes D_i independent of potential outcomes:
 $(Y_i(1), Y_i(0)) \perp D_i$
- ▶ Very often we deal with (ignorability assumption) $(Y_i(1), Y_i(0)) \perp D_i | X_i$

Randomized experiments

- ▶ Randomization solves this problem!
- ▶ Random assignment makes D_i independent of potential outcomes:
 $(Y_i(1), Y_i(0)) \perp D_i$
- ▶ Very often we deal with (ignorability assumption) $(Y_i(1), Y_i(0)) \perp D_i | X_i$
- ▶ Consider two random samples C and T, from the population - by construction, these samples are statistically identical to the entire population therefore:

$$E\{Y_i(0)|i \in C\} = E\{Y_i(0)|i \in T\} = E\{Y_i(0)\}$$

and

$$E\{Y_i(1)|i \in C\} = E\{Y_i(1)|i \in T\} = E\{Y_i(1)\}$$

Randomized experiments

- ▶ Randomization allows us to use the control units C as an image of what would happen to the treated unit T in the counterfactual situation of no treatment, and vice-versa
- ▶ Going back to the first equation

$$E\{\Delta_i\} = E\{Y_i(1) - Y_i(0)\} = E\{Y_i(1)|i \in T\} - E\{Y_i(0)|i \in C\}$$

- ▶ However, randomization is rarely a feasible solution for ethical concerns and technical implementation
- ▶ But: always useful benchmark. And increasingly feasible in some settings (also because of increasing awareness among policymakers and researchers)

Causality without experiments

- ▶ The **research design, identification strategy, or empirical strategy** is the approach used with observational data (i.e., data not generated by a randomized trial) to approximate a randomized experiment.
- ▶ Standard methods
 - ▶ Linear Regression
 - ▶ Difference-in-differences
 - ▶ Event Studies, Synthetic Control + Synthetic DiD
 - ▶ Instrumental Variables
 - ▶ Regression Discontinuity

Introduction to Regression

- ▶ How does schooling affect income?
- ▶ Assume a linear model

$$Y_i = \alpha + s_i\beta + \epsilon_i$$

- ▶ Y_i is wage as a function of s_i , years of education
- ▶ β is the slope parameter summarizing how wages vary with schooling.
- ▶ α , the “intercept” or “constant”, gives the expected income with no schooling ($s_i = 0$)
- ▶ ϵ_i includes all other factors affecting income besides schooling, including randomness

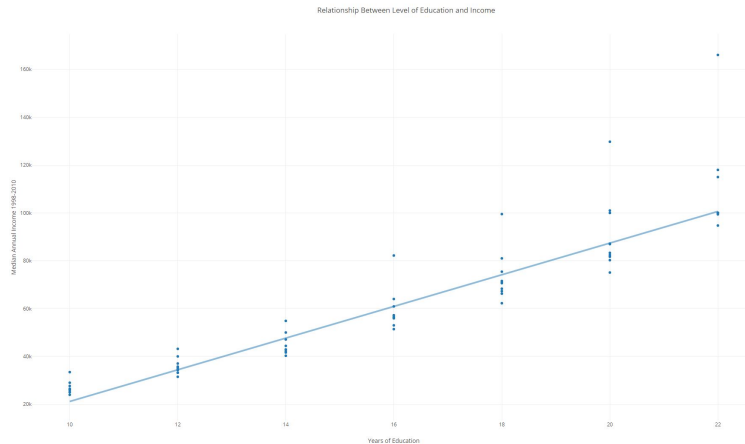
OLS Estimator

$$Y_i = \alpha + s_i\beta + \epsilon_i$$

- ▶ The Ordinary Least Squares (OLS) Estimator is the workhorse of applied microeconometrics
- ▶ It consists in minimizing the sum of squared residuals
- ▶ The OLS estimator is given by

$$\hat{\beta} \approx \frac{\text{Cov}(Y, s)}{\text{Var}(s)}$$

OLS Estimator



- ▶ $\hat{\beta}$ is the slope of the regression – how Y responds to a change of 1 in s
- ▶ Is this relationship causal?
- ▶ No, we need specific assumptions

Unbiased Estimates

- ▶ The **OLS exogeneity assumption** is $\text{Cov}(s, \epsilon) = 0$
- ▶ This should mean that the treatment is uncorrelated with the error term, i.e., no confounders
- ▶ When **conditional independence** is not satisfied, we say that “ s is endogenous”:
 - ▶ That is, an explanatory variable s_i is said to be endogenous if it is correlated with unobserved factors (confounders) that are also correlated with the outcome variable.
- ▶ Since the error term ϵ_i includes all unobserved factors affecting the outcome, we can define endogeneity when there is correlation between an explanatory variable and the error term $\text{Cov}(s, \epsilon) \neq 0$

Omitted variable bias

- ▶ Assume individuals who choose to get more education likely differ from those who don't: maybe they have a higher innate ability, enjoy schooling, and are good at it
- ▶ The true model could be

$$Y_i = \alpha + s_i\beta + \gamma a_i + \eta_i$$

- ▶ $\epsilon_i = a_i + \eta_i$ and we cannot measure a_i , while η_i is random error
- ▶ This would make $\hat{\beta}$ biased estimate of β

$$\hat{\beta} = \beta + \underbrace{\gamma \frac{\text{Cov}(s_i, a_i)}{\text{Var}(s_i)}}_{\text{Omitted Variable Bias}} + \underbrace{\frac{\text{Cov}(s_i, \eta_i)}{\text{Var}(s_i)}}_{=0 \text{ by assumption}}$$

Omitted variable bias

$$\underbrace{\gamma \frac{\text{Cov}(s_i, a_i)}{\text{Var}(s_i)}}_{\text{Omitted Variable Bias}}$$

Omitted Variable Bias

		Correlation of omitted variable with explanatory variable	
		$\text{Cov}[s, a] > 0$	$\text{Cov}[s, a] < 0$
Correlation of omitted variable with outcome	$\gamma > 0$	$\hat{\beta} > \beta$	$\hat{\beta} < \beta$
	$\gamma < 0$	$\hat{\beta} < \beta$	$\hat{\beta} > \beta$

Statistical Significance

- ▶ The value for β provides a prediction for the effect of the explanatory variable on the outcome
- ▶ But if this prediction is very noisy, then it might not be useful for policy analysis.
- ▶ To do causal inference, we have to determine whether the effect is statistically significant.
- ▶ This is generally achieved by computing a standard error for each coefficient and then using the standard error to compute confidence intervals and a p-value for the hypothesis that $\beta \neq 0$

Residuals and Standard Errors

- ▶ The residuals or errors from an OLS regression (\neq st. errors) are defined as

$$\begin{aligned}\tilde{\epsilon}_i &= Y_i - \hat{Y}_i \\ &= Y_i - \hat{\alpha} - \hat{\beta}s_i\end{aligned}$$

- ▶ The standard error (SE) for the OLS estimate $\hat{\beta}$ is

$$\hat{\sigma}_{\beta} = \sqrt{\frac{1}{n} \sum_{i=1}^n \tilde{\epsilon}_i^2},$$

the square root of the average of the squared residuals.

- ▶ SE provides information about the estimate's precision: a lower standard error is a more precise estimate.

t-statistics, p-values and confidence intervals

- ▶ A rule of thumb for statistical significance is to compute the **t-statistic**:

$$t = \frac{\hat{\beta}}{\hat{\sigma}_{\beta}}$$

- ▶ $t > 2 \rightarrow$ statistically significant positive effect, $t < -2 \rightarrow$ statistically significant negative effect
- ▶ A high t (in absolute value) is associated with a small **p-value** (e.g., $t = \pm 1.96 \rightarrow p = .05$)
- ▶ 95% **confidence intervals** indicate (roughly) that the coefficient is 95% likely to reside within that interval

$$CI_{0.95}^{\beta} = [\hat{\beta} - 1.96 \times \hat{\sigma}_{\beta}, \hat{\beta} + 1.96 \times \hat{\sigma}_{\beta}]$$

Causal inference toolkit

- ▶ Difference-in-differences
- ▶ Event Studies, Synthetic Control + Synthetic DiD
- ▶ Instrumental Variables
- ▶ Regression Discontinuity

Difference-in-differences

- ▶ Assume we have n units (i) and T time periods (t)
- ▶ Consider a binary policy D_{it} , and we are interested in estimating its effect on outcomes Y_{it}
- ▶ The inherent problem is that D_{it} is *not* necessarily randomly assigned
- ▶ Diff-in-Diff works under the assumption that **in the absence of the treatment, the Y_{it} across units evolve in parallel – their γ_t are identical.**

$$Y_{it}(D_{it}) = \alpha_i + \gamma_t + \tau_i D_{it}$$
$$\text{s.t. } Y_{it}(1) - Y_{it}(0) = \tau_i$$

- ▶ Absent the policy, units may have different *levels* (α_i) but their changes would evolve in parallel

Difference-in-differences (using linear regression)

- ▶ A simple linear regression will identify $E(\tau_i | D_{i1} = 1)$ with two time periods:

$$Y_{it} = \alpha_i + \gamma_t + D_{it}\beta + \epsilon_{it} \quad (1)$$

- ▶ This setup is sometimes referred to as the Two-way Fixed Effects estimator (TWFE)
- ▶ Note: we could have also estimated τ directly:

$$\hat{\tau} = \underbrace{E(Y_{i1} - Y_{i0} | D_i = 1)}_{\Delta \bar{Y}_1} - \underbrace{E(Y_{i1} - Y_{i0} | D_i = 0)}_{\Delta \bar{Y}_0}$$

- ▶ Intuitively, we generate a counterfactual for the treatment using the changes in the untreated units: $E(Y_{i1} - Y_{i0} | D_i = 0)$
- ▶ Necessary: two time periods! What if we have more?

Multiple time periods in basic setup

- ▶ Let's consider a policy that occurs all at t_0 (e.g. single timing rolled out to treated units)
- ▶ More time periods helps in several ways:
 1. If we have multiple periods *before* the policy implementation, we can partially test the underlying assumptions
 2. If we have multiple periods *after* the policy implementation, we can examine the timing of the effect
- ▶ How do we implement this?

$$Y_{it} = \alpha_i + \gamma_t + \sum_{t=1, t \neq t_0}^T \delta_t D_{it} + \epsilon_{it},$$

- ▶ One of the coefficients is fundamentally unidentified because of α_i
- ▶ All coefficients measure the effect *relative* to period t_0 .

Recent warning

- ▶ This literature has had a certain amount of upheaval over the past 5-6 years
- ▶ Tension: provide context for how people currently and historically have studied diff-in-diff
- ▶ For example, Roth and Sant'Anna (2021) or Rambachan and Roth (2020)

Example: Facebook and Mental Health

American Economic Review 2022, 112(11): 3660–3693
<https://doi.org/10.1257/aer.20211218>

Social Media and Mental Health[†]

By LUCA BRAGHERI, RO'EE LEVY, AND ALEXEY MAKARIN*

We provide quasi-experimental estimates of the impact of social media on mental health by leveraging a unique natural experiment: the staggered introduction of Facebook across US colleges. Our analysis couples data on student mental health around the years of Facebook's expansion with a generalized difference-in-differences empirical strategy. We find that the rollout of Facebook at a college had a negative impact on student mental health. It also increased the likelihood with which students reported experiencing impairments to academic performance due to poor mental health. Additional evidence on mechanisms suggests the results are due to Facebook fostering unfavorable social comparisons. (JEL D91, I12, I23, L82)

Example: Facebook and Mental Health

$$Y_{icgt} = \alpha_g + \delta_t + \beta \times \text{Facebook}_{gt} + \mathbf{X}_i \times \gamma + \mathbf{X}_c \times \lambda + \epsilon_{icgt}$$

TABLE 1—BASELINE RESULTS: INDEX OF POOR MENTAL HEALTH

	Index of poor mental health			
	(1)	(2)	(3)	(4)
Post-Facebook introduction	0.137 (0.040)	0.124 (0.022)	0.085 (0.033)	0.077 (0.032)
Observations	374,805	359,827	359,827	359,827
Survey-wave fixed effects	✓	✓	✓	✓
Facebook-expansion-group fixed effects	✓	✓		
Controls		✓	✓	✓
College fixed effects			✓	✓
FB-expansion-group linear time trends				✓

Example: Facebook and Mental Health

$$Y_{igt} = \alpha_g + \delta_t + \beta_k \times \sum_{k=-8}^5 D_{k(gt)} + \epsilon_{icgt}$$

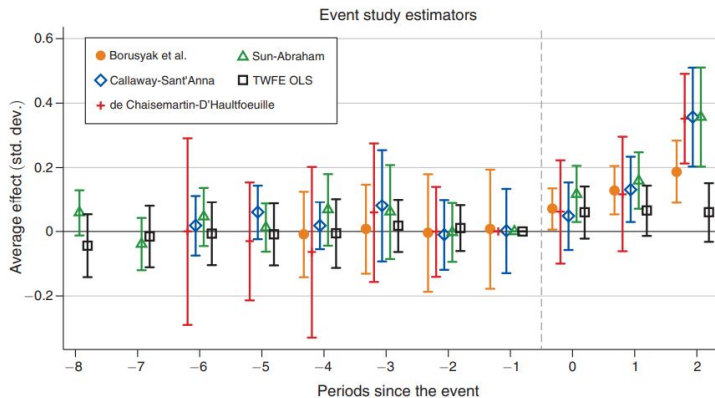


FIGURE 2. EFFECTS OF FACEBOOK ON THE INDEX OF POOR MENTAL HEALTH BASED ON DISTANCE TO/FROM FACEBOOK INTRODUCTION