

Conditional independence assumption and regressions

From randomisation to regression

- 2 conceptual differences:
 - 1 In Neyman's analysis (finite population), **potential outcomes are fixed** and assignment varies
 - 2 In regression analysis, **realized outcomes and assignment are fixed** but different units, with different **error** (but same treatment status) are sampled
- May seem like some "geeky jargon" details but in many settings, it is very important especially when we leave the experimental ideal.
- Example: Comparing the effect of a policy whose adoption was staggered in different US States: Your PSU are US states, it's a finite population of 51 units !
- Consider a pure randomized control trial, treatment D, outcome Y, individual attributes X.
- Let us consider our analysis sample as a random sample from an infinite population.
- This allows us to think of all variables as random variables with finite moments (e.g. population averages and standard deviation).
- In particular, define $\beta = \mathbb{E}[Y_i(1) - Y_i(0)]$ and $\alpha = \mathbb{E}[Y_i(0)]$.

Conditional independence assumption and regressions

From randomisation to regression

- Consider the population regression:

$$Y = \alpha + \beta D + \eta$$

- η is the **individual error**. In the OLS regression over our sample:

$$Y_i = \alpha + \beta D_i + \varepsilon_i$$

- ε_i is the **residual**. The least squares estimator for β is based on minimizing the sum of squared residuals over α and β

$$\left(\hat{\beta}_{\text{ols}}, \hat{\alpha}_{\text{ols}} \right) = \arg \min_{\beta, \alpha} \sum_{i=1}^N \left(Y_i^{\text{obs}} - \alpha - \beta \cdot D_i \right)^2,$$

- With solutions

$$\hat{\beta}_{\text{ols}} = \frac{\text{Cov}(D_i, Y_i)}{S_N^2} = \frac{\sum_{i=1}^N (D_i - \bar{D}) \cdot (Y_i^{\text{obs}} - \bar{Y}^{\text{obs}})}{\sum_{i=1}^N (D_i - \bar{D})^2} = \bar{Y}_t^{\text{obs}} - \bar{Y}_c^{\text{obs}}$$

- and

$$\hat{\alpha}_{\text{ols}} = \bar{Y}^{\text{obs}} - \hat{\beta}_{\text{ols}} \cdot \bar{D}$$

Conditional independence assumption and regressions

From randomisation to regression

- Note that random assignment of the treatment does **not** imply that the **error term** η is independent of D_i .
- In fact, in general there will be heteroskedasticity, and we need to use the Eicker-Huber White robust standard errors to get valid confidence intervals.
- Mean-independence of the treatment and population error is sufficient.
- The **error term** in the population regression also has a clear interpretation. With simple notation manipulation you can show:

$$\begin{aligned}
 \eta_i &= Y_i(0) - \alpha + D_i (Y_i(1) - Y_i(0) - \beta) \\
 &= (1 - D_i) \cdot \underbrace{(Y_i(0) - \mathbb{E}[Y_i(0)])}_{\text{Individual deviation in } Y_i(0)} + D_i \cdot \underbrace{(Y_i(1) - \mathbb{E}[Y_i(1)])}_{\text{Individual TE Heterogeneity}}
 \end{aligned}$$

Conditional independence assumption and regressions

Reminder

The Frisch Waugh Lovell (FWL) theorem

- Consider a dependent variable Y and two sets of regressors X_1 and X_2 and the linear model

$$Y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon$$

- Frisch and Waugh (1933) then Lovell (2010) prove the following results (Greene 2012, p.73):

Theorem

In the linear least squares regression of vector Y on two sets of variables, X_1 and X_2 , the subvector β_2 is the set of coefficients obtained when the residuals from a regression of Y on X_1 alone are regressed on the set of residuals obtained when each column of X_2 is regressed on X_1 .

- In the Appendix, I give you an illustration of the FWL theorem in action.

[▶ Go to illustration](#)

Conditional independence assumption and regressions

What is conditional independence ?

- So far, we manipulated **conditional expectations** $\mathbb{E}[Y_i|D_i]$ by treatment status and used independence of treatment and potential outcomes ($Y_i(1), Y_i(0) \perp D_i$) to define average treatment effects.
- In some settings, random assignment depend on other factors X (e.g. block randomisation) and the independence hold true **conditional** on the value of these other factors.
- In other settings, it may be plausible that an explanatory variable of interest (e.g. a treatment or policy) is independent of the outcomes conditional on some characteristics.
- This is the **conditional independence assumption**
- You have probably already read papers making causal claim saying things like "all other things being equal", "controlling for X, we find..." or "we matched observation based on the following variables"...
- Sometimes, we see the latin expression *Ceteris paribus*. That's it. That's the conditional independence assumption.

What is conditional independence ?

Few remarks

- You don't need to observe all variables that determine participation/assignment nor all those that determine potential outcomes.
- You need to observe **all** variables that determine **both participation and outcomes**: The **confounders**
- The ignorability of treatment assignment says that if you can't control for confounders, **your statistical model is showing a correlation and not causation**
- if there's a variable that determines participation, but not outcomes, it's not a confounder: that's called an instrument, and you should use it as such (More on that in lecture VI).
- The CIA Is Everywhere (in empirical papers). But unless there are very good reason to believe that there are no latent factor that confound the results (such as a randomized experiment), it is a **strong assumption**

What is conditional independence ?

Careful with the CIA

- What you read in the paper:
"In our preferred model, the set of control variables includes gender, age, country of residence during childhood, marital status, type of residence, wealth and occupation dummies."
- What it means in the data:



Prince Charles

Male
 Born in 1948
 Raised in the UK
 Married Twice
 Lives in a castle
 Wealthy and Famous



Ozzy Osbourne

Male
 Born in 1948
 Raised in the UK
 Married Twice
 Lives in a castle
 Wealthy and Famous

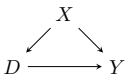
Figure 1: Source: Somewhere on twitter, probably @KhoaVuUm, saved on my phone for this moment

- I mean... Sure they are a fairly good match but whatever the "treatment" condition we would be comparing between the two, I am pretty sure there would be some unobserved factor correlated with treatment and outcomes.

Conditional independence assumption and regressions

A word on 'bad controls'

Figure 2: Conditioning is "closing the back door"



- A back-door path is any path from D to Y that starts with an arrow pointing **into** D .
- "Backdoor paths" creates **Fork** relationships: $D \leftarrow X \rightarrow Y$. We say X is a **confounder**.
- If one close a backdoor path (by conditioning), then the partial causal effect of D on Y is identified !
- There are two other different path configurations:
 - ① **Chains**: $D \rightarrow M \rightarrow Y$: In that case, M is either a **mediator** or D is an **instrument** for M
 - ② **Colliders**: $D \rightarrow C \leftarrow Y$: In that case, C is a **collider**
- A **Mediator** or mediating variable **transmit** the effect of D to Y through it. Distinction between total, direct and indirect effect.
- A **Collider** is more counter-intuitive ; In general, it's a variable caused by at least two others (arrows colliding...). In the treatment-effect framework, a collider is a variable that is both caused by the treatment and an outcome.

Conditional independence assumption and regressions

A word on 'bad controls'

A notorious collider bias: Survivorship bias

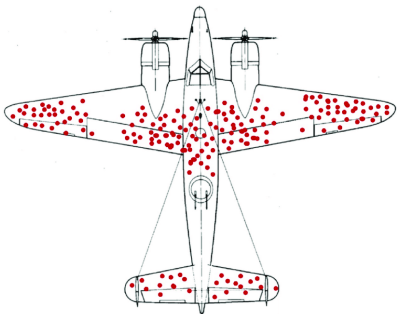


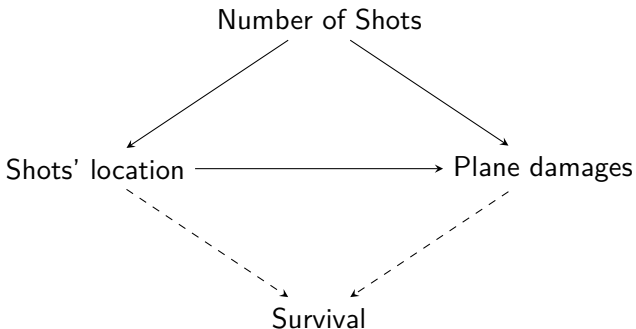
Figure 3: Wald plane as a symbol of survivorship bias

Conditional independence assumption and regressions

A word on 'bad controls'

A notorious collider bias: Survivorship bias

Figure 4: Wald plane in DAG



Conditional independence assumption and regressions

A word on 'bad controls'

How do DAG help us with "bad controls" ?

- From these definition it should be clear that with DAG, you explicitly show which relationships you are modelling and which ones you aren't.
- Clear definition of a causal relationships: identified if there is no backdoor paths left open.
- Also help you think about the role of other variables. **only confounders** should be controlled.
- **General rule:** don't control for post-exposure variables unless you are actually doing a mediation analysis.
- In Pearls word, a confounder is always a parent node, so in case a doubt, map a DAG.

Conditional independence assumption and regressions

A word on 'bad controls'

Collider bias: a high stake example (See the Discussion in Cunningham 2018)

- In a controversial paper, Fryer (2019) Analyse racial bias in police use of force.
- To do that he access a large database of arrest reports and with careful econometric analysis, controlling for many things, he find no evidence of racial bias in these data.
- These conclusions have been heavily criticized for various reasons. One critics from "Mr Selection model" (Pr. J.J. Heckman) is that these police administrative datasets select on officers' post-treatment decisions to detain civilians. Decisions that are potentially **also discriminatory**, thus **omitting all data on encounters not resulting in detentions** and potentially severely understating the extent of racial bias in policing (Knox, Lowe, and Mummolo 2020).
- In other words, the observation itself is conditioned on having been arrested, which is affected by the variable of interest (race) and outcome (if they intend to use force on you they are more likely to arrest you).
- It's fairly easy to notice that using a DAG.

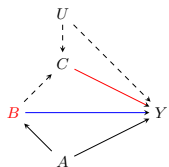
Conditional independence assumption and regressions

A word on 'bad controls'

Collider bias: a high stake example (See the Discussion in Cunningham 2018)

- Denote B for race, Y for use of force and the relationship between the two (in blue) is the one of interest.
- We denote Police controls and arrests C .
- Fryer observe individuals attributes A that affect both the likelihood that a person is from a minority and is brutalized by the police (externalizing behavior, neighborhood, past record),...
- There are unobserved factors U that we call police suspicion (can include many things) that cause both arrests and use of violence.
- From this graph, It's clear that the relationship of interest is not identified, because observations are conditioned on C although $B \rightarrow C$ is not observable.

Figure 5: C is a collider bias



Heterogenous treatment effect

Weird things happening under the hood

- Consider a RCT with block randomization and let P_j denote treatment probabilities across discrete blocks \mathbf{X} (e.g. the share of treated may be higher/lower in some groups).
- Consider the OLS regression with block fixed effects:

$$Y_i = \sum_j \delta_j \mathbb{1}(X_{ij} = 1) + \beta_{OLS} D_i + \varepsilon_i$$

- this regression is **saturated in the covariates**, which means that it is linear in the covariates by construction. It is **not fully saturated** because it doesn't include interactions between treatment and blocks. When is β_{OLS} equal to the ATE or the ATT?

Heterogenous treatment effect

Weird things happening under the hood

- ② "Paradox": $\beta_{OLS} = ATT$ when **very few units are treated** and $\beta_{OLS} = ATU$ when **most units are treated**.

- The very recent (and super clear) paper by Słoczyński 2022 explains what the problem is (you can see why with the proof in this course appendix).
- Suppose you have two strata: X_1 is large with few treated units so $P(X_1)$ is small, and a second X_2 that's small but with a lot of treated units.
- Intuitively, the motivation for using OLS is that the linear projection of Y on D and X is the best predictor of Y given D and X .
- So OLS is best at predicting **actual outcomes**.
- But causal inference is about predicting the **missing outcomes** i.e. the counterfactual values.
- If we wanted to predict "what is", we would put a lot of weights where we have a lot of precision (so, the big strata with few treated).
- That's what OLS does.
- But what we want is to estimate the counterfactual and for that we would need to put more weights where there are lot of treated units because that's where the treatment effect is more precisely estimated.

Heterogenous treatment effect

Comments and solutions

- If the probabilities do not vary much, this is of little importance. But if treatment probabilities vary a lot, the OLS results will be far from the true ATE.
- If we know the conditional treatment probabilities (also called **propensity scores** (by design), we can re-weight the observations (Imbens 2004):
 - To estimate the ATE, weight treated observations by $w_1 = \frac{1}{Pr(D_i|X_i=x)}$ and controls by $w_0 = \frac{1}{(1-Pr(D_i|X_i=x))}$;
 - To estimate ATT, weight treated observations by $w_1 = \frac{(Pr(D_i|X_i=x))}{(1-Pr(D_i|X_i=x))}$ and controls get unit weights ;
 - To estimate ATU, weight treated observations by $w_1 = \frac{(1-Pr(D_i|X_i=x))}{(Pr(D_i|X_i=x))}$ and controls get unit weights ;
- The idea that you can correct for non-random sampling by weighting by the reciprocal of the probability of selection dates back to Horvitz and Thompson (1952).
- When we don't know the probability, we need consistent estimators.

Conditional independence assumption and regressions

Common support

- When the conditional independence assumption holds, we actually need an extra assumption/condition to estimate the effects

$$0 < Pr(D_i = 1 | \mathbf{X} = \mathbf{x}) < 1 \quad \forall \mathbf{x} \quad (4)$$

- “At all x 's, there must be both treatment and control observations”*
- Implies $f(X|D = 1)$ overlaps with $f(X|D = 0)$
- If unmet, restrict sample to observations with overlap
- This is important because OLS will project over the support of X and if one group has no support on some values, then we rely on extrapolation.
- Always check that your observations are balanced over the support of the X . The distribution of the X impacts the weights OLS give to different observations.
- In practice, check for outliers, observe densities, scatter plots etc. Observation without common support should usually be dropped (Lechner and Strittmatter 2017).

Common support

Covariates with RCT: How to ?

1 Special case: fully saturated regressions

- $Y_i = \sum_x \mathbb{1}(X_i = x)\delta_x + \beta D_i + \sum_x \tau_x D_i \times \mathbb{1}(X_i = x) + \varepsilon_i$
- The regression fits the CEF perfectly (whatever the distribution of Y) because the true CEF is linear in parameters.
- Thus, the OLS estimate of β is an unbiased estimate of the ATE (Athey and Imbens 2017).
- The coefficients τ estimate treatment effect heterogeneity interpreted as deviation from the ATE.

2 General case where X may be continuous:

- **Pooled regression:** The regression of Y , D and X yields consistent estimate of the ATE provided D and X are uncorrelated, which follows under random assignment (Negi and Wooldridge 2021).
- It estimates the ATE if we assume constant treatment effect or a random sample of a large population (where heterogeneity is in the error term).
- **(Lin 2013) regressions:** Estimate $Y_i = \alpha + \beta D_i + \dot{X}_i \gamma + D_i \times \dot{X}_i \tau + \varepsilon_i$
- Where $\dot{X}_i = X_i - \bar{X}_i$ is the deviation from the population average (in practice, use sample mean).
- The demeaning of the covariates ensures that the coefficient on D is the treatment effect.

Common support

More advanced methods

- **Inverse propensity score weighting:** If we can access the conditional treatment probability (or propensity score) to construct weights, we can use weighted regression methods based on the inverse propensity score (see e.g. Imbens (2004)).
- **doubly robust methods:** We can use IPW in regressions with covariates. This inverse-probability-weighting regression adjustment (IPWRA) has the doubly robust property: we only need either the propensity score model to be true or the outcome model to be true to recover unbiased estimates. (See the recent paper by (Słoczyński, Uysal, and Wooldridge 2022))
- **Matching:** more on that in lecture 11.
- **Machine learning:** not this year, but if you are curious : [Susan Athey and Guido W Imbens. 2019. "Machine Learning Methods Economists Should Know About." *Annual Review of Economics* 11:62.](#) Important contribution when the set of potential covariates is large (Belloni, Chernozhukov, and Hansen 2014; Chernozhukov, Hansen, and Spindler 2015; Chernozhukov et al. 2017), or to identify treatment effect heterogeneity (Imai and Ratkovic 2013; Athey and Imbens 2016; Demirer et al. 2017; Wager and Athey 2018; Athey, Tibshirani, and Wager 2019)

Outline

- ① Introduction
- ② Conditional independence assumption and regressions
- ③ Case study: The STAR experiment by Krueger (1999)
 - Main specification and results
 - Problem 1: Dealing with attrition
 - Problem 2: Students changed class after random assignments
- ④ Inference: a history of variance
- ⑤ Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Case study: The STAR experiment by Krueger (1999)

Does class size or teacher/student ratio improve student achievement ?

The effect of small class size on IQ in Kindergarten

Figure 7: OLS estimates of the effect of class size on average percentile of Stanford Binet IQ test from Krueger (1999)[p. 512]

Explanatory variable	OLS: actual class size			
	(1)	(2)	(3)	(4)
A. Kindergarten				
Small class	4.82 (2.19)	5.37 (1.26)	5.36 (1.21)	5.37 (1.19)
Regular/aide class	.12 (2.23)	.29 (1.13)	.53 (1.09)	.31 (1.07)
White/Asian (1 = yes)	—	—	8.35 (1.35)	8.44 (1.36)
Girl (1 = yes)	—	—	4.48 (.63)	4.39 (.63)
Free lunch (1 = yes)	—	—	-13.15 (.77)	-13.07 (.77)
White teacher	—	—	—	- .57 (2.10)
Teacher experience	—	—	—	.26 (.10)
Master's degree	—	—	—	- .51 (1.06)
School fixed effects	No	Yes	Yes	Yes
R ²	.01	.25	.31	.31

Case study: The STAR experiment by Krueger (1999)

Does class size or teacher/student ratio improve student achievement ?

The effect of small class size on IQ in first grade

Figure 8: OLS estimates of the effect of class size on average percentile of Stanford Binet IQ test from Krueger (1999)[p. 512]

	B. First grade			
Small class	8.57 (1.97)	8.43 (1.21)	7.91 (1.17)	7.40 (1.18)
Regular/aide class	3.44 (2.05)	2.22 (1.00)	2.23 (0.98)	1.78 (0.98)
White/Asian (1 = yes)	—	—	6.97 (1.18)	6.97 (1.19)
Girl (1 = yes)	—	—	3.80 (.56)	3.85 (.56)
Free lunch (1 = yes)	—	—	-13.49 (.87)	-13.61 (.87)
White teacher	—	—	—	-4.28 (1.96)
Male teacher	—	—	—	11.82 (3.33)
Teacher experience	—	—	—	.05 (0.06)
Master's degree	—	—	—	.48 (1.07)
School fixed effects	No	Yes	Yes	Yes
R^2	.02	.24	.30	.30

Case study: The STAR experiment by Krueger (1999)

Does class size or teacher/student ratio improve student achievement ?

Consistent estimate with imputed missing outcomes

Figure 9: Comparison of raw estimates with specifications with imputed missing outcomes Krueger (1999)[p. 512]

TABLE VI
EXPLORATION OF EFFECT OF ATTRITION DEPENDENT VARIABLE: AVERAGE
PERCENTILE SCORE ON SAT

Grade	Actual test data		Actual and imputed test data	
	Coefficient on small class dum.	Sample size	Coefficient on small class dum.	Sample size
K	5.32 (.76)	5900	5.32 (.76)	5900
1	6.95 (.74)	6632	6.30 (.68)	8328
2	5.59 (.76)	6282	5.64 (.65)	9773
3	5.58 (.79)	6339	5.49 (.63)	10919

Estimates of reduced-form models are presented. Each regression includes the following explanatory variables: a dummy variable indicating initial assignment to a small class; a dummy variable indicating initial assignment to a regular/aid class, unrestricted school effects; a dummy variable for student gender; and a dummy variable for student race. The reported coefficient on small class dummy is relative to regular classes. Standard errors are in parentheses.

Case study: The STAR experiment by Krueger (1999)

Does class size or teacher/student ratio improve student achievement ?

What's important

- Large experiment with high compliance rate and fairly clean design shows that reducing class size in early grades improve cognitive development
- Teacher assistant don't seem to work that well thus testing the hypothesis that it's less about teacher/student ratio but maybe the learning environment ?
- This paper was very influential and opened the path of a 10-year academic debate on the effects of class size.
- At that time, the consensus was that class size has little to no effect on students' achievement (See Krueger, Hanushek, and Rice (2002))
- The puzzle with teaching assistants also nurtured a broad strand of literature (See the work of Peter Blatchford for instance)

The problems with clusters

What is usually meant when one talks about clusters

- Most econometrics textbooks⁴ approaches the clustering issue as something close to omitted variable bias where, the initial model:

$$Y_{ic} = \alpha + \mathbf{X}'\beta + \mu_{ic}$$

actually hides the fact that the error term μ_{ic} has a group structure s.t.:

$$\mu_{ic} = v_c + \varepsilon_{ic}$$

- And thus, estimating the model without accounting for that yields biased standard errors because $\mathbb{E}[\mu_{ic}\mu_{jc}] = \rho\sigma_\mu^2 > 0$
- This presentation, although pedagogical, reinforces the confusion between fixed effect and clustering.

$$(Y_{ic} - \bar{Y}_c) = (\mathbf{X}_{ic} - \bar{\mathbf{X}}_c)' \beta + \mu_{ic} - \bar{\mu}_c$$

4. For instance Cameron and Trivedi 2005; Angrist and Pischke 2008; Wooldridge 2010; Wooldridge 2012

Conventional wisdom about standard errors

When to cluster according to Colin Cameron and Miller 2015

Until recently, the conventional wisdom was summed up as follows:

“There are settings where one may not need to use cluster-robust standard errors. We outline several though note that in all these cases it is always possible to still obtain cluster-robust standard errors and contrast them to default standard errors. If there is an appreciable difference, then use cluster robust standard errors”. (p.334)

This is actually wrong according to Abadie et al. (2022)

Clustering: illustration through simulations

Simulation with $\rho = .5$ intra-cluster correlation

- We consider again the simple bivariate regression

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

- We generate a sample of 1000 observations allocated across 50 clusters (e.g. municipalities) with $\rho = .5$ correlation within cluster.
- Like before, we generate a true $\beta=0$
- We estimate this equation using OLS and collect the estimate of $\hat{\beta}$
- We do that 10 000 times with a new random sample and plot the distribution of the estimated $\hat{\beta}$

Clustering: illustration through simulations

Distribution of estimations with clustering

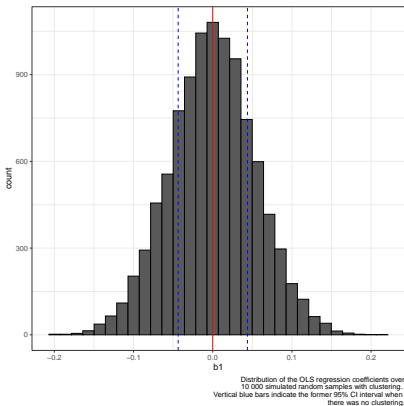


Figure 15: Density of the β estimates with clustering

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Context and motivation

Bruno Crépon et al. 2013. “Do Labor Market Policies Have Displacement Effects? Evidence from a Clustered Randomized Experiment.” *The Quarterly Journal of Economics* 128 (2): 531–580

- In France, unemployment rate is 17.5 % for age 15 30 against 9.2 % in the whole population
- Higher education has traditionally been somewhat protective
 - In France, unemployment rate is 9.4 % for college graduates vs 21.4 % for the others
 - However, even educated youth may experience unemployment and long term unemployment
 - 20 to 30 % of young high school/college graduates have been unemployed for more than 6 months, and around 10 % have been unemployed for more than 12 months

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Context and motivation

- One common policy response is to provide hard to place jobseekers with **reinforced counseling scheme**
- Provide assistance with writing resume, searching for job offers and answering to them, preparing for interviews
- Reinforced counseling programs are **costly** as they mean more frequent meetings with the caseworker
- Intensive support caseworkers have about 30 unemployed in their caseload, instead of 120 in the normal situation
- One strong orientation of the public employment policy was to **use services of private operators** instead of Pôle emploi
- End of the monopoly of the Employment Agency is a key component of the Employment policy in France
- Work through contracts with placement agencies

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

A clever experimental design

- A “super control group” eligible unemployed workers in 0 assignment areas
- Comparing those assigned to **control groups** and those assigned to the **super control group** identify **displacement effects**. Why ?
- Without displacement effects, exit rates of unemployment should not vary by treatment intensity.
- Comparing those assigned to **treatment groups** and those assigned to the **super control group** identify the average effect on the treated.

Case study: displacement effects of job search assistance (Crépon Et. Al. 2013)

Estimation

Figure 20: Main regressions in Crépon et al. (2013)

We estimate a fully unconstrained reduced form model, and test whether the effect of being assigned to treatment or to control varies by assignment probability. The specification we consider is the following:

$$\begin{aligned}
 y_{ic} &= \beta_{25}Z_{ic}P_{25c} + \beta_{50}Z_{ic}P_{50c} + \beta_{75}Z_{ic}P_{75c} + \beta_{100}Z_{ic}P_{100c} \\
 &+ \delta_{25}P_{25c} + \delta_{50}P_{50c} + \delta_{75}P_{75c} \\
 &+ X_{ic}\gamma_4 + u_{ic}
 \end{aligned} \tag{6}$$

where Z_{ic} is the assignment to treatment variable and P_{xc} is a dummy variable at the area level indicating an assignment rate of $x\%$. ZP_{25} is thus a dummy for being assigned to treatment in a labor market with a rate of 25% assignment. As before, control variables are individual characteristics (gender, education, etc.) and the set of 47 dummy variables for city quintuplets (our randomization strata). Standard errors are clustered at the local area level. The parameter β_x measures the effect of being assigned to treatment in an area where $x\%$ of the eligible population was assigned to treatment, compared to being unassigned in an area of the same type (or, for β_{100} , compared to the super-control). Coefficient δ_x measures the effect of being assigned to the control group in an area where $x\%$ of the eligible population was assigned to treatment, compared to being in the super-control group in which no one was assigned to treatment. Note

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Estimation

Figure 21: Main table in Crépon et al. (2013)

	Labor market outcome: Long term fixed contract			
	All workers (1)	Not employed		
		All (2)	Men (3)	Women (4)
Assigned to treatment in 25% areas	0.016 (0.012)	0.021 (0.014)	0.037 (0.027)	0.015 (0.016)
Assigned to treatment in 50% areas	0.009 (0.012)	0.013 (0.013)	0.021 (0.021)	0.008 (0.020)
Assigned to treatment in 75% areas	-0.015 (0.016)	0.007 (0.019)	0.001** (0.030)	-0.016 (0.021)
Assigned to treatment in 100% areas	0.010 (0.009)	0.025** (0.010)	0.021 (0.014)	0.028** (0.014)
25% areas	-0.002 (0.010)	-0.015 (0.011)	-0.041** (0.019)	-0.001 (0.013)
50% areas	-0.002 (0.010)	-0.014 (0.013)	-0.026 (0.018)	-0.005 (0.017)
75% areas	0.016 (0.016)	-0.006 (0.020)	-0.055** (0.027)	0.014 (0.024)
Control Mean	0.199	0.167	0.150	0.178
F-test for equality of all assigned to treatment coefficients to zero	0.34	0.05	0.07	0.22
F-test for equality of all areas coefficients to zero	0.72	0.48	0.04	0.92
F-test for equality of all areas coefficients	0.52	0.90	0.59	0.77
Number of observations	21431	11806	4387	7419

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Estimation

- Noting significant among all workers, so we focus among those who were unemployed at the time of randomisation.
- The main ITT is in the 4th line. For the unemployment, receiving the treatment instead of not increases fixed-term employment by 2.5pp, slightly more for women.
- Coefficients by treatment intensity and their relative displacement effects are almost perfectly symmetric suggesting that, indeed, job search assistance **impose negative externalities to the untreated**, especially when there are more treated units.
- The F tests at the bottom of the table test the joint hypotheses of null effects or that the effects are equal.
- However, there is little power due to the clustering setting, imperfect compliance and the number of hypotheses to test.
- After that, the regressions **impose some structure** on the parameters and move away from the ITT. Parameters are still causal but their interpretation is not straightforward (OLS weighting and all that).

Case study: displacement effects of job search assistance (Crepon Et. Al. 2013)

Alternative model

Figure 24: Alternative regression equation Crépon et al. (2013)

Table 5: Reduced form: Impact of the program, accounting for externalities						
By job type: share of job seekers who are eligible for program						
	Not employed			Not employed, above third quartile		
	All	Men	Women	All	Men	Women
	(1)	(2)	(3)	(4)	(5)	(6)
Panel A: Long term fixed contract						
Assigned to program (β)	0.023*** (0.008)	0.043*** (0.013)	0.013 (0.010)	0.040** (0.016)	0.072** (0.029)	0.021 (0.022)
In a Program area (δ)	-0.013 (0.009)	-0.036*** (0.013)	-0.001 (0.012)	-0.040* (0.021)	-0.086** (0.035)	-0.013 (0.027)
Net effect of program assignment ($\beta+\delta$)	0.010 (0.008)	0.007 (0.011)	0.012 (0.011)	0.000 (0.019)	-0.014 (0.031)	0.008 (0.024)
Control Mean	0.16	0.131	0.177	0.19	0.161	0.204

Wrap-up

Fairly advanced econometric stuff today

- The formal link between RCT and regression analysis
- Case study : typically example of what could be in a test
- Inference: what it means, and the impact of clustering
- Advanced design: A double-nested randomisation
- Most of the things we'll see from now on try to mimic settings that are akin to RCT. This was the core, this you **have to master**.

Next week: Difference in differences

- **To read: mandatory:** Card and Krueger 1993
- **To read: mandatory:** Bertrand, Duflo, and Mullainathan 2004
- **Very good DiD paper:** The impact of Glyphosate on children birth outcomes in Brazil [Mateus Dias, Rudi Rocha, and Rodrigo R Soares. 2023. "Down the River: Glyphosate Use in Agriculture and Birth Outcomes of Surrounding Populations." *The Review of Economic Studies* \(February 6, 2023\): rdad011](#)

Outline

- 8 The Frisch-Waugh-Lovell theorem in action
- 9 Refresher on conditional independence
- 10 Estimations with heterogeneous treatment effect

The Frisch-Waugh-Lovell theorem in action

Basic illustration: life expectancy and GDP per capita

- Use GAPMINDER package and data to illustrate regressions with the good old relationship between life expectancy and GDP per capita.

```
head(gapminder)
```

```
## # A tibble: 6 x 6
##   country      continent  year  lifeExp      pop  gdpPercap
##   <fct>        <fct>    <int> <dbl>    <int>    <dbl>
## 1 Afghanistan Asia      1952  28.8  8425333  779.
## 2 Afghanistan Asia      1957  30.3  9240934  821.
## 3 Afghanistan Asia      1962  32.0 10267083  853.
## 4 Afghanistan Asia      1967  34.0 11537966  836.
## 5 Afghanistan Asia      1972  36.1 13079460  740.
## 6 Afghanistan Asia      1977  38.4 14880372  786.
```

The Frisch-Waugh-Lovell theorem in action

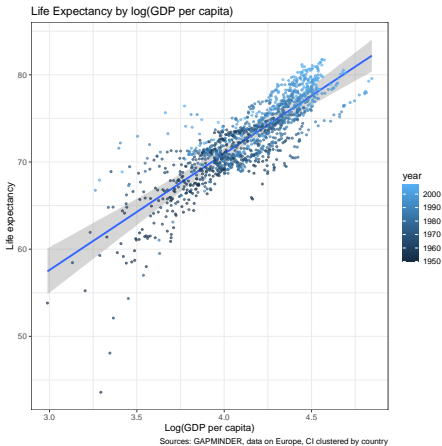


Figure 25: Linear regression of life expectancy on log GDP per capita

The Frisch-Waugh-Lovell theorem in action

What did we do ? What can we do ?

- In the previous slide we estimated the regression:

$$Y_{it} = \alpha + \beta D_{it} + \epsilon_{it}$$

- We want to account for systematic differences between countries and common evolution over years. One way to do that is to define 2 extra sets of dummies $C_i = \mathbb{1}(\text{Country} = i)$ and $T_t = \mathbb{1}(\text{year} = t)$ that we put in a matrix \mathbf{X} .
- Consider the regression:

$$Y_{it} = \alpha + \beta D_{it} + \mathbf{X}'\delta + \epsilon_{it}$$

- By the FWL theorem, estimating this regression gives the same estimate for $\hat{\beta}$ as estimating subsequently:
 - 1 $Y_{it} = \alpha_0 + \mathbf{X}'\rho + \mu_{it}$ removing the time and country fixed effects in the outcome
 - 2 $D_{it} = \alpha_1 + \mathbf{X}'\eta + v_{it}$ removing the time and country fixed effects in GDP per capita
 - 3 $\mu_{it} = \alpha_2 + \beta v_{it} + \epsilon_{it}$ the residualized outcome on the residualized "treatment"
- Let's see that

The Frisch-Waugh-Lovell theorem in action

Frisch-Waugh-Lovell in action in R !

```
FW_S1 <- lm_robust(lifeExp ~ factor(year) + factor(country), data = mygapminder,
  cluster = country)

# Retrieve the residual from this first regression, call it res_Life
mygapminder$res_Life <- mygapminder$lifeExp - FW_S1$fitted.values
# Then regress the log of GDP per capita over the same year dummies
FW_S2 <- lm_robust(log(gdpPercap) ~ factor(year) + factor(country), data = mygapminder,
  cluster = country)
# get the residual, call them res_gdp
mygapminder$res_gdp <- log(mygapminder$gdpPercap) - FW_S2$fitted.values

# Now we regress the first residual on the second:
FW <- lm_robust(res_Life ~ res_gdp, data = mygapminder, cluster = country)
# Compare with the regression with controls
Controls <- lm_robust(res_Life ~ res_gdp + factor(year) + factor(country), data = mygapminder,
  cluster = country)
# plot the residual over years to see there's no correlation anymore, but still
# remaining variation with gdp in particular, that we can color code to add a
# dimension ;-)
restime <- ggplot(mygapminder) + geom_point(aes(y = res_Life, x = year, color = log(gdpPercap))) +
  geom_smooth(aes(y = res_Life, x = year), method = "lm_robust", method.args = list(clusters = mygapminder,
  scale_color_gradient(low = "blue", high = "orange"))
restime
```

The Frisch-Waugh-Lovell theorem in action

We removed time and country fixed variation

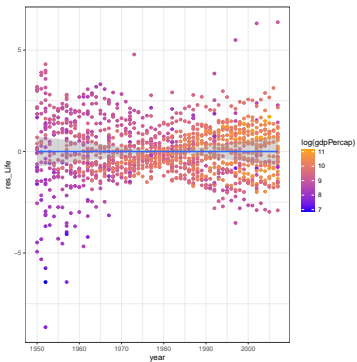
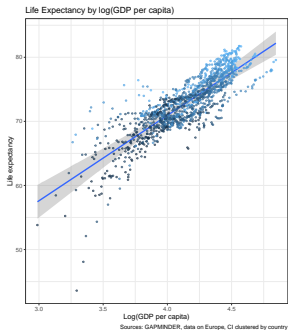
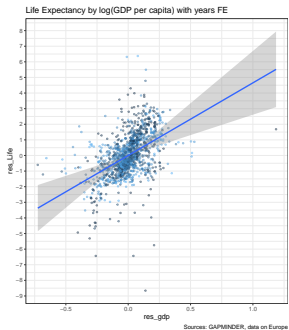


Figure 26: Remaining variation after controlling for time and country fixed effects

The Frisch-Waugh-Lovell theorem in action



(a) No control



(b) After controlling for time and country fixed effects

Figure 27: Evolution of the correlation between GDP and life expectancy after controlling for time and country fixed effects

The Frisch-Waugh-Lovell theorem in action

How to interpret OLS models with log transformations

Table 2: Summary of Functional Forms Involving Logarithms

Model	Dependent Variable	Independent Variable	Interpretation of β_1
Level-level	y	x	$\Delta y = \beta_1 \Delta x$
Level-log	y	$\log(x)$	$\Delta y = (\beta_1/100) \% \Delta x$
Log-level	$\log(y)$	x	$\% \Delta y = (100\beta_1) \Delta x$
Log-log	$\log(y)$	$\log(x)$	$\% \Delta y = \beta_1 \% \Delta x$

Sources:(Wooldridge 2012, p.72)

The Frisch-Waugh-Lovell theorem in action

By FWL the estimated semi-elasticities are the same

Table 3: Illustration of the Frisch-Waugh-Lowell theorem

	Simple Correlation	Controls	Residualized
Log(GDP)	5.788*** (0.509)	4.653*** (1.067)	
Log(GDP), residualized			4.653*** (1.041)
Constant	X	X	X
Year+country Fixed effect		X	
Num.Obs.	1302	1302	1302
R2	0.720	0.927	0.244
R2 Adj.	0.720	0.922	0.243
AIC	5964.9	4395.5	4213.5
BIC	5980.4	4881.6	4229.0
RMSE	2.39	1.22	1.22
Std.Errors	by: country	by: country	by: country

† It's a *level-log* model in base 10, so we interpret the coefficient on *gdp* by saying "when GDP increases by 1 point, the life expectancy increases by $\hat{\beta}/100$ years"

[▶ Back to CIA](#)

Outline

- 8 The Frisch-Waugh-Lovell theorem in action
- 9 Refresher on conditional independence
 - Independence
 - Conditional independence
- 10 Estimations with heterogeneous treatment effect

Refresher on conditional independence

Independence

- This is a good time for a quick refresher on independence. Two random variables are independent if and only if: $f_{X,Y}(x,y) = f_X(x)f_Y(y)$.
- For discrete random variables: $P(X = x, Y = y) = P(X = x)P(Y = y)$
- In terms of events: $P(A \cap B) = P(A)P(B)$. These definitions are not that intuitive but: What is the conditional probability if two events are **independent**?

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$$

- So the probability of A given that B occurs is just $P(A)$. In words, B happening does not affect $P(A)$ (and vice versa)
- Better: knowing one doesn't tell you anything about the other event chances of happening

Refresher on conditional independence

Conditional independence

- Conditional independence is an important concept and closely related to regression models and the conditional independence assumption
- Events A and B are conditionally independent if
$$P(A \cap B | Z) = P(A | Z)P(B | Z)$$
- More useful: If A and B are conditional independent given Z , then
$$P(A | B, Z) = P(A | Z)$$
- In words, knowing B doesn't tells us anything about $P(A)$ once we know Z

Outline

- 8 The Frisch-Waugh-Lovell theorem in action
- 9 Refresher on conditional independence
- 10 Estimations with heterogeneous treatment effect**
 - Derivating the ATT and ATU
 - What the saturated-in-X regression gives

Estimations with heterogeneous treatment effect

Derivating the ATT and ATU

- First, it's useful to remember that the ATE is a weighted sum of the conditional ATEs: $\beta_X = \mathbb{E}[Y_i(1) - Y_i(0)|X_i = x]$. We compute the overall ATE as $\beta = \sum_x \beta_X Pr[X_i = x]$.
- Note that, by ignorability,

$$\begin{aligned}\beta_X &= \mathbb{E}[Y_i(1)|X_i = x, D_i = 1] - \mathbb{E}[Y_i(0)|X_i = x, D_i = 0] \\ &= \mathbb{E}[Y_i(1) - Y_i(0)|X_i = x]\end{aligned}$$

- There is a similar derivation for the ATT:

$$\begin{aligned}\beta_{ATT} &= \mathbb{E}[Y_i(1) - Y_i(0)|D_i = 1] \\ &= \mathbb{E}\left[\mathbb{E}[Y_i(1) - Y_i(0)|X_i, D_i = 1]|D_i = 1\right] \text{ by the LIE} \\ &= \mathbb{E}\left[(\mathbb{E}[Y_i(1)|X_i, D_i = 1] - \mathbb{E}[Y_i(0)|X_i, D_i = 1])|D_i = 1\right] \\ &= \mathbb{E}\left[(\mathbb{E}[Y_i(1)|X_i] - \mathbb{E}[Y_i(0)|X_i]) \mid D_i = 1\right] \text{ by ignorability} \\ &= \mathbb{E}[\beta_X \mid D_i = 1] \\ &= \sum_x \beta_X Pr[X_i = x \mid D_i = 1]\end{aligned}$$

Estimations with heterogeneous treatment effect

Derivating the ATT and ATU

- Remember Bayes's formula:

$\Pr[X_i = x | D_i = 1] \Pr[D_i = 1] = \Pr[D_i = 1 | X_i = x] \Pr[X_i = x]$, so we can rewrite the ATT as a propensity score-weighted function of the CATEs (with a normalizing factor):

$$\beta_{ATT} = \frac{\sum_x \beta_X \cdot \Pr[D_i = 1 | X_i = x] \Pr[X_i = x]}{\sum_x \Pr[D_i = 1 | X_i = x] \Pr[X_i = x]}$$

- In words, the Average treatment effect on the treated is the weighted average of block-specific ATE, with weights equal to the conditional treatment probability in the block times the probability of being in this block.
- So the weight is estimated by the share of treated in a block times the share of the sample in the block.
- We make the same derivation with the ATU.

Estimations with heterogeneous treatment effect

What the saturated-in- X regression gives

- Consider the saturated-in- X regression:

$$Y_i = \sum_j \delta_j \mathbb{1}(X_{ij} = 1) + \beta_{OLS} D_i + \varepsilon_i$$

- And the auxiliary regression of treatment over block indicators:

$$D_i = \sum_j \pi_j \mathbb{1}(X_{ij} = 1) + v_i$$

- This equation is fully saturated and thus estimates $\mathbb{E}[D_i | \mathbf{X}_i]$ and thus $v_i = D_i - \mathbb{E}[D_i | \mathbf{X}_i]$.
- By the FWL theorem, β_{OLS} is equivalent to the regression of Y_i on the residual of the previous auxiliary regression:

$$\begin{aligned} \beta_{OLS} &= \frac{\text{Cov}(Y_i, v_i)}{\mathbb{V}[v_i]} \\ &= \frac{\mathbb{E}\left[Y_i \cdot (D_i - \mathbb{E}[D_i | \mathbf{X}_i])\right]}{\mathbb{E}\left[(D_i - \mathbb{E}[D_i | \mathbf{X}_i])^2\right]} \end{aligned} \quad (5)$$

Estimations with heterogeneous treatment effect

What the saturated-in-X regression gives

- Remember, estimating the regression of Y on X and D is the same as estimating Y on $\mathbb{E}[Y_i | \mathbf{X}_i, D_i]$, so in the big expectation in the numerator, we can substitute Y_i by $\mathbb{E}[Y_i | \mathbf{X}_i, D_i]$:

$$\beta_{OLS} = \frac{\mathbb{E}[Y_i | D_i, \mathbf{X}_i] \cdot \mathbb{E}[(D_i - \mathbb{E}[D_i | \mathbf{X}_i])]}{\mathbb{E}[(D_i - \mathbb{E}[D_i | \mathbf{X}_i])^2]} \quad (6)$$

- We can expand the CEF $\mathbb{E}[Y_i | \mathbf{X}_i, D_i]$ further to get:

$$\mathbb{E}[Y_i | \mathbf{X}_i, D_i] = \mathbb{E}[Y_i | D_i = 0, \mathbf{X}_i] + \beta_X D_i$$

- We then plug this expression in the numerator of the previous equation

$$\begin{aligned} \mathbb{E}[Y_i | D_i, \mathbf{X}_i] \cdot \mathbb{E}[(D_i - \mathbb{E}[D_i | \mathbf{X}_i])] &= \mathbb{E} \left[(D_i - \mathbb{E}[D_i | \mathbf{X}_i]) \mathbb{E}[Y_i | D_i = 0, \mathbf{X}_i] \right] \\ &\quad + \mathbb{E} \left[D_i (D_i - \mathbb{E}[D_i | \mathbf{X}_i]) \beta_X \right] \end{aligned}$$

Estimations with heterogeneous treatment effect

What the saturated-in- X regression gives

- Because $\mathbb{E}[Y_i|D_i = 0, \mathbf{X}_i]$ is a function of \mathbf{X}_i , it is uncorrelated with $(D_i - \mathbb{E}[D_i|\mathbf{X}_i])$, so the first hand term on the right-hand side is zero ;

$$\mathbb{E}\left[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])\mathbb{E}[Y_i|D_i = 0, \mathbf{X}_i]\right] = 0$$

- For the same reason, D_i is uncorrelated with $(D_i - \mathbb{E}[D_i|\mathbf{X}_i])$ so the numerator actually becomes:

$$\begin{aligned}\mathbb{E}[Y_i D_i, \mathbf{X}_i] \cdot \mathbb{E}[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])] &= \mathbb{E}\left[D_i(D_i - \mathbb{E}[D_i|\mathbf{X}_i])\beta_X\right] \\ &= \mathbb{E}\left[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])^2 \beta_X\right]\end{aligned}$$

- At this point, we have shown:

$$\beta_{OLS} = \frac{\mathbb{E}\left[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])^2 \beta_X\right]}{\mathbb{E}\left[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])^2\right]} = \underbrace{\frac{\mathbb{E}\left[\mathbb{E}[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])^2 | \mathbf{X}_i] \beta_X\right]}{\mathbb{E}\left[\mathbb{E}[(D_i - \mathbb{E}[D_i|\mathbf{X}_i])^2 | \mathbf{X}_i]\right]}}_{\text{Using the LIE again}} \quad (7)$$

Estimations with heterogeneous treatment effect

What the saturated-in- X regression gives

- There are components in this expression we can make sense of, in particular:

$$\mathbb{E}[(D_i - \mathbb{E}[D_i | \mathbf{X}_i])^2 | \mathbf{X}_i] = \sigma_D^2(\mathbf{X}_i)$$

- Is the conditional variance of the treatment given \mathbf{X} . So:

$$\beta_{OLS} = \frac{\mathbb{E}[\sigma_D^2(\mathbf{X}_i)\beta_X]}{\mathbb{E}[\sigma_D^2(\mathbf{X}_i)]} \quad (8)$$

- This establishes that the regression model of Y on block fixed effect and a treatment produces a treatment-variance weighted average of block specific average treatment effects β_X .
- Finally, because D_i is binary

$$\sigma_D^2(\mathbf{X}_i) = \text{Pr}(D_i | \mathbf{X}_i) \cdot (1 - \text{Pr}(D_i | \mathbf{X}_i))$$

- So,

$$\begin{aligned} \beta_{OLS} &= \frac{\mathbb{E}[\text{Pr}(D_i | \mathbf{X}_i) \cdot (1 - \text{Pr}(D_i | \mathbf{X}_i)) \beta_X]}{\mathbb{E}[\text{Pr}(D_i | \mathbf{X}_i) \cdot (1 - \text{Pr}(D_i | \mathbf{X}_i))]} \\ &= \frac{\sum_x \beta_x (\text{Pr}(D_i | \mathbf{X}_i = x) \cdot (1 - \text{Pr}(D_i | \mathbf{X}_i = x))) \text{Pr}(\mathbf{X}_i = x)}{\sum_x (\text{Pr}(D_i | \mathbf{X}_i = x) \cdot (1 - \text{Pr}(D_i | \mathbf{X}_i = x))) \text{Pr}(\mathbf{X}_i = x)} \end{aligned} \quad (9)$$

Estimations with heterogeneous treatment effect

What the saturated-in- X regression gives

- This shows that the regression estimand weights each covariate-specific treatment effect by
$$(Pr(D_i|X_i = x) \cdot (1 - Pr(D_i|X_i = x))) Pr(X_i = x)$$
- The OLS regression put more weights where the treatment variance is highest, not where the treatment probability is highest.
- Comparing the weights for ATT and that of OLS, we clearly see that the weights are "polluted" by 1 - the probability of treatment and if we looked at the weights for the ATU that would be the pollution of the treatment probability.
- Therefore, unless treatment effect is constant or the conditional treatment probabilities are constant and equal to .5, the OLS estimand is not the ATE, nor the ATT, nor the ATU, but a conditional treatment-variance weighted parameter of the ATE.

▶ [Back to Heterogeneous treatment effects](#)