

ScPoEconometrics: Advanced Intro and Recap 1

Bluebery Planterose SciencesPo Paris 2023-01-24

Welcome to *ScPoEconometrics: Advanced*!

Today

- 1. Who Am I
- 2. This Course
- 3. Recap 1 of topics from intro course

Next time

- Quiz 1 (before next time)
- Recap 2



Who Am I

- I'm an PhD candidate at the Paris School of Economics. Check out my website!
- I work on tax evasion, climate policies, and macro topics:
 - 1. Acceptability of climate policies: who support/oppose climate policies and why?
 - 2. *Offshore real-estate in Dubai using leaked data*: how large is it, who owns it, and what does it tell us about global offshore real-estate?
 - 3. *Excess Profit Tax*: how to tax excess profits from energy firms that benefited from the war in Ukraine?





Prerequisites

- This course is the *follow-up* to Introduction to Econometrics with R which is taught to 2nd years.
- You are supposed to be familiar with all the econometrics material from the slides of that course and/or chapters 1-9 in our textbook.
- We also assume you have basic **R** working knowledge at the level of the intro course!
 - basic data.frame manipulation with dplyr
 - simple linear models with 1m
 - basic plotting with ggplot2
 - Quiz 1 will try and test for that 😉, so be on top of this chapter



This Course

Grading

- 1. There will be *four quizzes* on Moodle roughly every two weeks => 40%
- 2. There will be *two take home exams / case studies* => 60%
- 3. There will be *no* final exam \cong .



This Course

Grading

- 1. There will be *four quizzes* on Moodle roughly every two weeks => 40%
- 2. There will be *two take home exams / case studies* => 60%
- 3. There will be *no* final exam \cong .

Course Materials

- 1. Book chapter 10 onwards
- 2. The Slides
- 3. The interactive shiny apps
- 4. Quizzes on Moodle



Syllabus

- 1. Intro, Recap 1 (Quiz 1)
- 2. Recap 2 (*Quiz 2*)
- 3. Intro, Difference-in-Differences
- 4. Tools: Rmarkdown and data.table
- 5. Instrumental Variables 1 (Quiz 3)
- 6. Instrumental Variables 2 (*Midterm exam*)

- 7. Panel Data 1
- 8. Panel Data 2 (Quiz 4)
- 9. Discrete Outcomes
- 10. Intro to Machine Learning 1
- 11. Intro to Machine Learning 2
- 12. Recap / Buffer (Final Project)
- 12. Recap / Buffer (Final Project)



Course Organization

| | Lundi 23/1 | Mardi 24/1 | Mercredi 25/1 | Jeudi 26/1 | Vendredi 27/1 | Samedi 28/1 | Dimanche 29/1 |
|-------|--|--|--|--|--|-------------|---------------|
| | | | | | | | |
| 08:00 | | 08 : 00 - 10 : 00 ເສ Taux d'indisponibilité : 100.00% | 08 : 00 - 10 : 00 ໝ Taux d'indisponibilité : 15.00% | 08 : 00 - 10 : 00 ₪ Taux d'indisponibilité : 05.00% | 08 : 00 - 10 : 00 🕱 Taux d'indisponibilité : 10.00% | | |
| 09:00 | | | | | | | |
| 10:00 | 10 • 15 - 12 • 15 🕾 | 10 · 15 - 12 · 15 m | 10 • 15 - 12 • 15 🕅 | 10 • 15 - 12 • 15 🕅 | 10 • 15 - 12 • 15 @ | | |
| 11:00 | Taux d'indisponibilité : 10.00% | Taux d'indisponibilité : 35.00% | Taux d'indisponibilité : 45.00% | Taux d'indisponibilité : 40.00% | Taux d'indisponibilité : 25.00% | | |
| 12:00 | | | | | | | |
| | 12 . 20 - 14 . 20 | 12 . 20 - 14 . 20 | 12 . 20 - 14 . 20 @ | 12 . 20 - 14 . 20 ** | 12 . 20 - 14 . 20 | | |
| 13:00 | Taux d'indisponibilité : 30.00% | Taux d'indisponibilité : 20.00% | Taux d'indisponibilité : 80.00% | Taux d'indisponibilité : 35.00% | Taux d'indisponibilité : 05.00% | | |
| 14:00 | | | | | | | |
| 15:00 | 14 : 45 - 16 : 45 ⊠ Taux d'indisponibilité : 20.00% | 14:45-16:45 🕸 Taux d'indisponibilité : 30.00% | 14 : 45 - 16 : 45 ⊠ Taux d'indisponibilité : 30.00% | 14 : 45 - 16 : 45 🛱 Taux d'indisponibilité : 50.00% | 14:45-16:45 🛱 Taux d'indisponibilité : 20.00% | | |
| 16:00 | | | | | | | |
| 17:00 | 17 : 00 - 19 : 00 @ Taux d'indisponibilité : 30.00% | | 17 : 00 - 19 : 00 ☺ Taux d'indisponibilité : 10.00% | 17 : 00 - 19 : 00 ⊠ Taux d'indisponibilité : 30.00% | 17 : 00 - 19 : 00 ☺ Taux d'indisponibilité : 25.00% | | |
| 18:00 | | | | | | | |
| 19:00 | 19 : 15 - 21 : 15 ଞ Taux d'indisponibilité : | 19 : 15 - 21 : 15 অ Taux d'indisponibilité : | 19 : 15 - 21 : 15 🛱 Taux d'indisponibilité : | | 19 : 15 - 21 : 15 🛱 Taux d'indisponibilité : | | |
| 20:00 | | 15.00% | | | 15.00% | | |
| 21:00 | | | | | | | |
| | | | | | | | |



Recap 1

Let's get cracking! 🦾



8 / 42

Models and notation

We write our (simple) population model

$$y_i=eta_0+eta_1x_i+u_i$$

and our sample-based estimated regression model as

$$y_i = {\hat eta}_0 + {\hat eta}_1 x_i + e_i$$

An estimated regression model produces estimates for each observation:

$${\hat y}_i = {\hat eta}_0 + {\hat eta}_1 x_i$$

which gives us the *best-fit* line through our dataset.



(A lot of this set slides - in particular: pictures! - have been taken from Ed Rubin's outstanding material. Thanks Ed A)

Task 1: Run Simple OLS (4 minutes)

- 1. Load data here. in dta format. (Hint: use haven::read_dta("filename") to read this
 format.)
- 2. Obtain common summary statistics for the variables classize, avgmath and avgverb. Hint: use the skimr package.
- 3. Estimate the linear model

 $\mathrm{avgmath}_i = eta_0 + \mathrm{classize}_i x_i + u_i$



Task 1: Solution

1. Load the data

grades = haven::read_dta(file ="https://www.dropbox.com/s/wwp2cs9f0dubmhr/grade5.dta?dl=1")

2. Describe the dataset:

library(dplyr)
grades %>%
 select(classize,avgmath,avgverb) %>%
 skimr::skim()

3. Run OLS to estimate the relationship between class size and student achievement?

summary(lm(formula = avgmath ~ classize, data = grades))





Population



 $y_i=2.53+0.57x_i+u_i$ $y_i=eta_0+eta_1x_i+u_i$



Sample 1: 30 random individuals



Sample 1: 30 random individuals



Population relationship $y_i = 2.53 + 0.57x_i + u_i$

Sample relationship $\hat{y}_i = 2.36 + 0.61 x_i$



Sample 2: 30 random individuals



Population relationship $y_i = 2.53 + 0.57 x_i + u_i$

Sample relationship $\hat{y}_i = 2.79 + 0.56 x_i$



Sample 3: 30 random individuals

Population relationship $y_i = 2.53 + 0.57x_i + u_i$

Sample relationship $\hat{y}_i = 3.21 + 0.45 x_i$

Let's repeat this **10,000 times**.

(This exercise is called a (Monte Carlo) simulation.)



Question: Why do we care about *population vs. sample*?



- On **average**, our regression lines match the population line very nicely.
- However, **individual lines** (samples) can really miss the mark.
- Differences between individual samples and the population lead to **uncertainty** for the econometrician.

Question: Why do we care about *population vs. sample*?

Question: Why do we care about *population vs. sample*?

Answer: Uncertainty matters.

- Every random sample of data is different.
- Our (OLS) estimators are computed from those samples of data.
- If there is sampling variation, there is variation in our estimates.

Question: Why do we care about *population vs. sample*?

Answer: Uncertainty matters.

- Every random sample of data is different.
- Our (OLS) estimators are computed from those samples of data.
- If there is sampling variation, there is variation in our estimates.

- OLS inference depends on certain assumptions.
- If violated, our estimates will be biased or imprecise.
- Or both. 😧

Linear regression

The estimator

We can estimate a regression line in R ($lm(y \sim x, my_data$)) and stata (reg y x). But where do these estimates come from?

A few slides back:

$${\hat y}_i = {\hat eta}_0 + {\hat eta}_1 x_i$$

which gives us the *best-fit* line through our dataset.

But what do we mean by "best-fit line"?

Being the "best"

Question: What do we mean by *best-fit line*?

Answers:

• In general (econometrics), *best-fit line* means the line that minimizes the sum of squared errors (SSE):

$$ext{SSE} = \sum_{i=1}^n e_i^2$$
 where $e_i = y_i - {\hat y}_i$,

- Ordinary least squares (OLS) minimizes the sum of the squared errors.
- Based upon a set of (mostly palatable) assumptions, OLS
 - Is unbiased (and consistent)
 - Is the *best* (minimum variance) linear unbiased estimator (BLUE)

Let's consider the dataset we previously generated.



For any line
$$\left(\hat{y} = {\hat eta}_0 + {\hat eta}_1 x
ight)$$



For any line $\left(\hat{y}=\hat{eta}_{0}+\hat{eta}_{1}x
ight)$, we can calculate errors: $e_{i}=y_{i}-\hat{y}_{i}$



For any line $\left(\hat{y}=\hat{eta}_{0}+\hat{eta}_{1}x
ight)$, we can calculate errors: $e_{i}=y_{i}-\hat{y}_{i}$



For any line $\left(\hat{y}=\hat{eta}_0+\hat{eta}_1x
ight)$, we can calculate errors: $e_i=y_i-\hat{y}_i$



SSE squares the errors $\left(\sum e_i^2\right)$: bigger errors get bigger penalties.



The OLS estimate is the combination of $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize SSE.



ScPoApps::launchApp("reg_simple")



Formally

In simple linear regression, the OLS estimator comes from choosing the $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize the sum of squared errors (SSE), *i.e.*,

 $\min_{\hat{\beta}_0,\,\hat{\beta}_1} \mathrm{SSE}$

OLS

Formally

In simple linear regression, the OLS estimator comes from choosing the $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize the sum of squared errors (SSE), *i.e.*,

$$\min_{\hat{eta}_0,\,\hat{eta}_1}\mathrm{SSE}$$

but we already know $ext{SSE} = \sum_i e_i^2$. Now use the definitions of e_i and \hat{y} .

$$e_i^2 = (y_i - {\hat y}_i)^2 = \left(y_i - {\hat eta}_0 - {\hat eta}_1 x_i
ight)^2 \ = y_i^2 - 2y_i {\hat eta}_0 - 2y_i {\hat eta}_1 x_i + {\hat eta}_0^2 + 2{\hat eta}_0 {\hat eta}_1 x_i + {\hat eta}_1^2 x_i^2$$

OLS

Formally

In simple linear regression, the OLS estimator comes from choosing the $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimize the sum of squared errors (SSE), *i.e.*,

$$\min_{\hat{eta}_0,\,\hat{eta}_1} \mathrm{SSE}$$

but we already know $ext{SSE} = \sum_i e_i^2$. Now use the definitions of e_i and \hat{y} .

$$e_i^2 = (y_i - {\hat y}_i)^2 = \left(y_i - {\hat eta}_0 - {\hat eta}_1 x_i
ight)^2 \ = y_i^2 - 2y_i {\hat eta}_0 - 2y_i {\hat eta}_1 x_i + {\hat eta}_0^2 + 2{\hat eta}_0 {\hat eta}_1 x_i + {\hat eta}_1^2 x_i^2$$

Recall: Minimizing a multivariate function requires (1) first derivatives equal zero (the 1storder conditions) and (2) second-order conditions (concavity).



Interactively

ScPoApps::launchApp("SSR_cone")



10

х





Interactively

We skipped the maths.

We now have the OLS estimators for the slope

$${\hat eta}_1 = rac{\sum_i (x_i - \overline{x})(y_i - \overline{y})}{\sum_i (x_i - \overline{x})^2}$$

and the intercept

$${\hat eta}_0 = \overline{y} - {\hat eta}_1 \overline{x}$$

Remember that *those* two formulae are amongst the very few ones from the intro course that you should know by heart! 💗



Interactively

We skipped the maths.

We now have the OLS estimators for the slope

$${\hat eta}_1 = rac{\sum_i (x_i - \overline{x})(y_i - \overline{y})}{\sum_i (x_i - \overline{x})^2}$$

and the intercept

$${\hat eta}_0 = \overline{y} - {\hat eta}_1 \overline{x}$$

Remember that *those* two formulae are amongst the very few ones from the intro course that you should know by heart! 💗

We now turn to the assumptions and (implied) properties of OLS.

Question: What properties might we care about for an estimator?

Question: What properties might we care about for an estimator?

Tangent: Let's review statistical properties first.

Refresher: Density functions

Recall that we use **probability density functions** (PDFs) to describe the probability a **continuous random variable** takes on a range of values. (The total area = 1.)

These PDFs characterize probability distributions, and the most common/famous/popular distributions get names (*e.g.*, normal, *t*, Gamma).

Here is the definition of a *PDF* f_X for a *continuous* RV X:

$$\Pr[a \leq X \leq b] \equiv \int_a^b f_X(x) dx$$

Refresher: Density functions

The probability a standard normal random variable takes on a value between -2 and 0: ${
m P}(-2 \leq X \leq 0) = 0.48$



Refresher: Density functions

The probability a standard normal random variable takes on a value between -1.96 and 1.96: $P(-1.96 \le X \le 1.96) = 0.95$



Refresher: Density functions

The probability a standard normal random variable takes on a value beyond 2: ${
m P}(X>2)=0.023$



Imagine we are trying to estimate an unknown parameter β , and we know the distributions of three competing estimators. Which one would we want? How would we decide?



Question: What properties might we care about for an estimator?

Question: What properties might we care about for an estimator?

Answer one: Bias.

On average (after *many* samples), does the estimator tend toward the correct value?

More formally: Does the mean of estimator's distribution equal the parameter it estimates?

$$extsf{Bias}_{eta} \Big(\hat{eta} \Big) = oldsymbol{E} \Big[\hat{eta} \Big] - eta$$

Answer one: Bias.

Unbiased estimator:
$$\boldsymbol{E}\left[\hat{\boldsymbol{\beta}}\right] = \boldsymbol{\beta}$$



Answer one: Bias.



Answer two: Variance.

The central tendencies (means) of competing distributions are not the only things that matter. We also care about the **variance** of an estimator.

$$\mathrm{Var}\!\left(\hat{eta}
ight) = oldsymbol{E}\!\left[\left(\hat{eta} - oldsymbol{E}\!\left[\hat{eta}
ight]
ight)^2
ight]$$

Lower variance estimators mean we get estimates closer to the mean in each sample.

Answer two: Variance.



Answer one: Bias.

Answer two: Variance.

Subtlety: The bias-variance tradeoff.

Should we be willing to take a bit of bias to reduce the variance?

In econometrics, we generally stick with unbiased (or consistent) estimators. But other disciplines (especially computer science) think a bit more about this tradeoff.

The bias-variance tradeoff.

ß



Properties

As you might have guessed by now,

- OLS is **unbiased**.
- OLS has the **minimum variance** of all unbiased linear estimators.

Properties

But... these (very nice) properties depend upon a set of assumptions:

- 1. The population relationship is linear in parameters with an additive disturbance.
- 2. Our X variable is **exogenous**, *i.e.*, E[u|X] = 0.
- 3. The X variable has variation. And if there are multiple explanatory variables, they are not perfectly collinear.
- 4. The population disturbances u_i are independently and identically distributed as normal random variables with mean zero ($\boldsymbol{E}[u] = 0$) and variance σ^2 (*i.e.*, $\boldsymbol{E}[u^2] = \sigma^2$). Independently distributed and mean zero jointly imply $\boldsymbol{E}[u_i u_j] = 0$ for any $i \neq j$.

Assumptions

Different assumptions guarantee different properties:

- Assumptions (1), (2), and (3) make OLS unbiased.
- Assumption (4) gives us an unbiased estimator for the variance of our OLS estimator.

We will discuss solutions to **violations of these assumptions**. See also our discussion in the book

- Non-linear relationships in our parameters/disturbances (or misspecification).
- Disturbances that are not identically distributed and/or not independent.
- Violations of exogeneity (especially omitted-variable bias).

Conditional expectation

For many applications, our most important assumption is **exogeneity**, *i.e.*,

E[u|X]=0

but what does it actually mean?

Conditional expectation

For many applications, our most important assumption is **exogeneity**, *i.e.*,

E[u|X]=0

but what does it actually mean?

One way to think about this definition:

For any value of X, the mean of the residuals must be zero.

- E.g., E[u|X=1] = 0 and E[u|X=100] = 0
- E.g., $E[u|X_2 = ext{Female}] = 0$ and $E[u|X_2 = ext{Male}] = 0$
- Notice: E[u|X] = 0 is more restrictive than E[u] = 0

Graphically...

Valid exogeneity, *i.e.*, E[u|X]=0



40 / 42

Invalid exogeneity, i.e., E[u|X] eq 0







- bluebery.planterose@sciencespo.fr
- � Original Slides from Florian Oswald
- 🗞 Book
- O @ScPoEcon