

# EFR summary

Applied Microeconomics,

FEM11087

2023-2024



Lectures 1 to 38

Weeks 1 to 7

**Deloitte.**

DeNederlandscheBank

EUROSYSTEEM

## **Details**

**Subject:** applied microeconometrics 2022-2023

**Teacher:** Pilar Garcia-Gomez

**Date of publication:** 18.10.2024

© This summary is intellectual property of the Economic Faculty association Rotterdam (EFR). All rights reserved. The content of this summary is not in any way a substitute for the lectures or any other study material. We cannot be held liable for any missing or wrong information. Erasmus School of Economics is not involved nor affiliated with the publication of this summary. For questions or comments contact [summaries@efr.nl](mailto:summaries@efr.nl)

---

# applied microeconometrics - module 1 - linear regression analysis

## Lecture 1 - Introduction to the linear regression model

**EMPIRICAL ANALYSIS** □ It is a scientific methodology where we use the data to test a theory and to estimate relationships between variables.

First, we have to define our research question. They can come from:

- Existing economic models.
- Via intuitive and less formal reasoning (something that inspires us to be studied from an economic point of view and with a scientific method).

**SIMPLE REGRESSION MODEL** □ It is a model that tells us somehow how a variable defines another one.

We have two variables,  $y$  and  $x$ , and we:

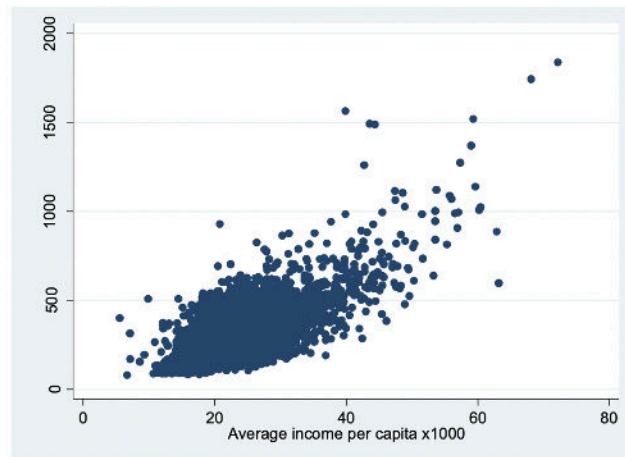
- want to explain  $y$  in terms of  $x$
- want to know how  $y$  varies with changes in  $x$

Example: how does the crime rate ( $y$ ) change with changing the number of police officers ( $x$ ) in a city?

Main example: HOUSE PRICES AND AVERAGE INCOME IN A NEIGHBOURHOOD<sup>1</sup>  
(Garcia-Gomez, 2022)

We could study these variables to know what policies are incrementing.

Figure n. 1<sup>2</sup> (Garcia-Gomez, 2022)



Source: CBS. Information for neighbourhoods in the Netherlands in 2012

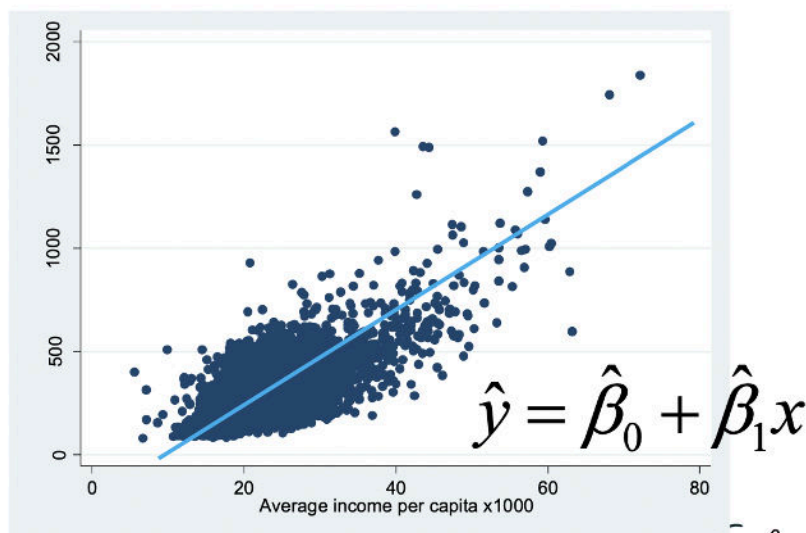
*Ezafun*

We can observe that there is a positive association between wages and house prices (higher the income, higher the house price).

The aim of the linear regression model is to find a line that can summarize all the information given by the scatter plots, to show the predicted value of the average house price as a function of the average income per capita.

The line's formula is:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

Figure n. 2<sup>3</sup> (Garcia-Gomez, 2022)



*Ezafun*

$\hat{\beta}_0$  is the intercept (it is the house price when income = 0, which in this case is not very indicative).

$\hat{\beta}_1$  is the slope (it tells us how house prices change when income does it as well).

## simple regression model

$$y = \beta_0 + \beta_1 x + u$$

- $y$  is the dependent variable.
- $x$  is the independent or explanatory variable.
- $u$  is the error term or disturbance (also referred as the “unobserved”).

$u$  it is all that is unobserved by the researchers that has an impact on the dependent variable (in our example, everything else that affect house prices in a neighbourhood, e.g., the number of amenities present in a neighbourhood).

## ceteris paribus condition

We are interested into knowing how the dependent variable changes when the  $x$  changes, holding all other variables fixed.

If the factors in  $u$  are held fixed  $\square \Delta u = 0$

So  $\square \Delta y = \beta_1 \Delta x$   $\square$  this is the interpretation of a SLOPE in our linear regression model.

## zero conditional mean assumption

- $\square$  The unobserved ( $u$ ) does not change when  $x$  changes in term of expected values.
- $\square E(x) = E(u) = 0$

This makes possible to see the line of the linear regression model in terms of **EXPECTATIONS** (what is the expected value of  $y$  with a given value of  $x$ )

$$E(x) = \beta_0 + \beta_1 x$$

Does the ceteris paribus condition works with our example?

It depends if the zero conditional mean assumption holds.

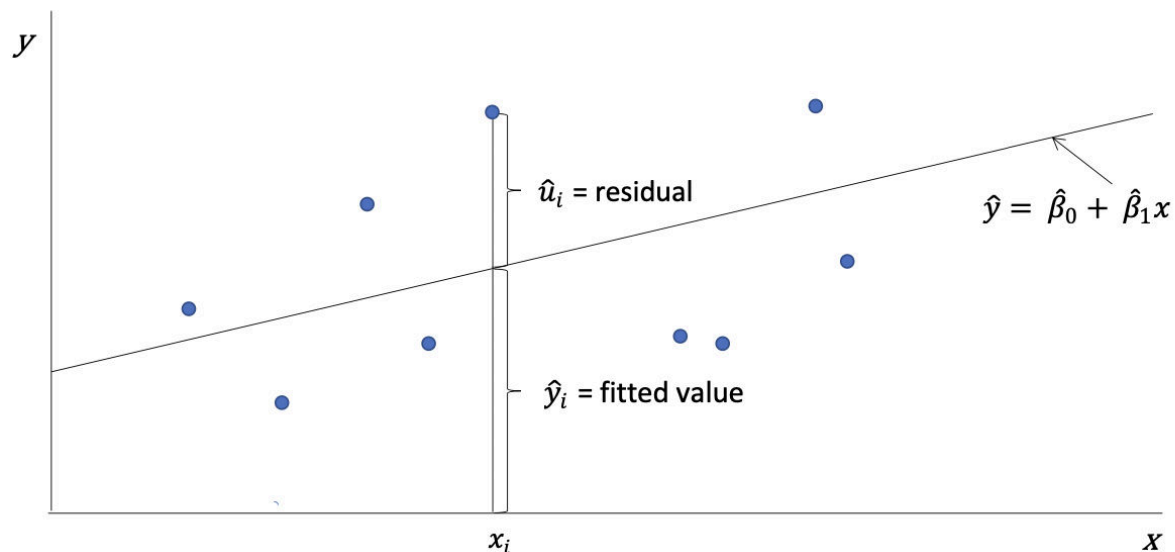
For example:  $u$  contains the quantity and quality of amenities in a neighbourhood.

- If we assume that the number of amenities does not change given the average income in a neighbourhood, then we have  $\Delta u = 0$   $\square$  the zero conditional mean assumption holds.
- If the number of amenities varies depending on the wealth of the neighbourhoods, then the zero conditional mean assumption does not hold up  $\square$  we cannot draw ceteris paribus conclusions.

## Lecture 2 - Estimation and interpretation in the linear regression model

### ORDINARY LEAST SQUARE ESTIMATES

Figure n. 3<sup>4</sup> (Garcia-Gomez, 2022)



We want an estimate of  $\beta_0$  and  $\beta_1$  and we will use the ordinary least square estimates.

1. In the population we expect that  $y$  and  $x$  are related in a linear way.
  - $\square y_i = \beta_0 + \beta_1 x_i + u_i$   $\square$  this is a random sample of the population of interest.
  - We will never know exactly  $\beta_0$  and  $\beta_1$ , but we want a good estimate of those.
  - First, we draw a random sample from the population, and for every individual in the random sample we can plot the value for  $x$  and  $y$  (see the figure n. 3).
  - Given the values for  $x$  and  $y$ , we draw a fitted line, similar as before.

FITTED VALUE ( $\hat{y}_i$ ): is the value that, for a given  $x_i$ , falls on the fitted line.

We see that there is still a difference between the observed value and the value on the fitted line: the difference is called the RESIDUAL  $\hat{u}_i$ .

The fitted value is just a predicted value:  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

The residual is:  $\hat{u}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$  □  $\hat{u}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$

Our aim is to have the residuals as small as possible to have the best possible fit.

2. We then obtain the estimates of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  by minimizing the square of the residuals:

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

This is what the **ordinary least square** estimator does to obtain the estimates. The ordinary least square estimator is the most efficient and unbiased estimator. To calculate these values, we must use STATA.

With STATA we obtain an equation with which we can draw the fitted line.

## MULTIPLE REGRESSION MODEL

It is more difficult to draw ceteris paribus conclusions with this model. For example, house prices are also related to the population density of an area:

$$\text{house price} = \beta_0 + \beta_1 \text{income} + \beta_2 \text{density} + u$$

And before we had:

$$\text{house price} = \beta_0 + \beta_1 \text{income} + u$$

We remember that we can draw ceteris paribus conclusions if the unobserved is not correlated to the explanatory variable. Now we have to see if the population density is correlated to income: richer householders are more likely to live in less populated areas? If yes, the zero conditional mean assumption could not be satisfied □ there would not be any ceteris paribus conclusions and would be better to estimate the simple regression model.

Everything we have seen so far also applies to the multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

Multiple regression analysis allows us to control for many other factors that simultaneously affect the dependent variable.

In addition, controlling for more variables also allows us to have better predictions.

## Lecture 3 – OLS assumptions – unbiasedness

**UNBIASEDNESS:** With the concept of “unbiasedness” we mean that the expected value of our estimated parameter is equal to the population parameter.

$$E(\hat{\beta}_0) = \beta_0$$

$$E(\hat{\beta}_1) = \beta_1$$

Please note that the hat is used for estimated values and anything without the hat is used to indicate population's values.

We need four assumptions for this property:

1. (MLR1) The model must be linear in parameters.
2. (MLR2) We must have random sampling.
3. (MLR3) There cannot be perfect collinearity.
4. (MLR4) The zero conditional mean assumption must be satisfied ( $E(u) = 0$ ).

### 1. LINEARITY IN PARAMETERS ASSUMPTION

We have the population model. In this model we have the variable  $x$  (independent variable) and  $y$  (dependent variable).  $y$  is related to  $x$  and the error  $u$  via the following equation:

$$y = \beta_0 + \beta_1 x + u$$

All the parameters are linearly related.

Important: the assumption is about the linearity in the parameters. So, the linearity is not possible in cases like this:



$$y = \beta_0 + \beta_1 x_1 + \beta_1 \beta_2 x_2 + u$$

$\beta_1 \beta_2$  □ Interaction term between  $\beta_1$  and  $\beta_2$  □ not possible!

However, there can be nonlinearities in the variables:

$$y = \beta_0 + \beta_1 \ln \ln(x_1) + \beta_2 x_2 + \beta_3 x_2^2 + u \quad (\text{Example of quadratic variable})$$

$$\ln \ln(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_2^2 + u \quad (\text{Example of logarithmic variable})$$

In these cases, we need to change our interpretation of the variables, and we are going to see a few more examples later on.

## 2. RANDOM SAMPLING ASSUMPTION

For this assumption to be true, we need to pick a random sample of size  $n$ , following the population model. If, for whatever reason, the sample will not be picked randomly, we will incur in selection bias.

## 3. NO PERFECT COLLINEARITY ASSUMPTION

If we pick a sample, in the following sample (and so in the population):

- Among all the independent variables, none of them should be constant, so that we have variation in all the independent variables. This is important because we use the variation to estimate the effect of our variable  $x$  on our variable  $y$ .
- There cannot be any exact linear relationship between the independent variables considered.

Example of perfect collinearity:

Given the model:

$$\text{houseprice} = \beta_0 + \beta_1 \text{income} + \beta_2 \text{Rotterdam} + \beta_3 \text{density} + \beta_4 \text{percentageyoungs} + \beta_5 \text{percentage elderly} +$$

Avoiding perfect collinearity means that among the independent variables there is not any exact linear relationship.

But if we assume that all and only the elderly people lived in Rotterdam (*Rotterdam = perc\_elderly*), these two variables would suffer the same kind of variation, and therefore the relationship between them would be perfectly linear.

In general, there is perfect collinearity between  $x_1$ ,  $x_2$  and  $x_3$  if  $x_3$  can be expressed as a combination of the other two variables ( $x_3 = ax_1 + bx_2$ ).

We can have two types of collinearity:

- **PERFECT COLLINEARITY:** in these cases, the estimation will not be successful (it could be that some software will not execute the commands given or will give inappropriate results). STATA tries to resolve the problem arbitrarily dropping the faulty variable in order to estimate a model without this problem; the problem is that it could drop one of the variables of most interest in our model. Therefore, we should first clearly and properly define our model in order to avoid any problem of this sort.
- **IMPERFECT COLLINEARITY:** when we estimate the model, the estimation works, but we get imprecise values for the estimates. We must watch out for:
  - $x$ 's with high correlation between them.
  - Imperfect collinearity between variables, for example between  $x_1$  and  $x_2$ . This can happen when we find big F-stat, when  $x_1, x_2$  jointly significant, but small T-statistics.

## 4. ZERO CONDITIONAL MEAN ASSUMPTION

Will be discussed later when talking about exogeneity.

# Lecture 4 – Linear regression analysis: OLS assumptions – inference

## INFERENCE – HYPOTHESIS TESTING

We want to test hypothesis about a parameter, or a group of parameters, in the population. We need not only information about the property of the estimator, so the expected value of the estimator, but also about his distribution (related to the distribution of the errors).

So, the following assumptions are related to the distribution of the errors:

- MLR5 – Homoskedasticity
- MLR6 – Normality

Under assumptions 1-6 OLS estimator is the unbiased estimator with minimum variance.

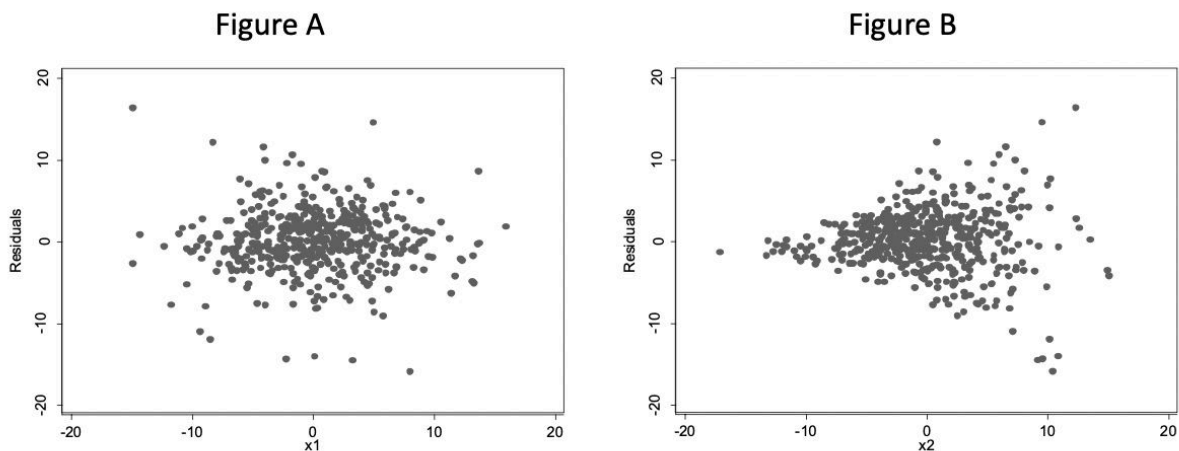
## 5. HOMOSKEDASTICITY ASSUMPTION

The variance of the error term does not change regardless of the values assumed by the independent variables:

$$\text{var}(x_1, x_2) = \text{var}(u) = \sigma^2$$

- For every individual the error term has the same importance, despite of the characteristics.
- The outcome of  $y$  has the same magnitude of uncertainty for every level of  $x$ 's.

Figure n. 4<sup>5</sup> (Garcia-Gomez, 2022)



In this case, the Figure A has homoskedastic residuals, while in Figure B we can see that the variance of the residuals grows with  $x$ .

If this assumption does not hold, we have a condition called **heteroskedasticity**:

$$\text{var}(x_1, x_2) = f(x_1, x_2) \text{ or } f(x_1) \text{ or } f(x_2)$$

When we have heteroskedasticity:

- OLS estimates  $\hat{\beta}$  are still unbiased and inefficient.

- OLS standard errors for the estimators  $se(\hat{\beta})$  are incorrect.

This is a problem for inference, but it is not a big issue since with STATA we can easily adjust the SE and the statistics used for inference.

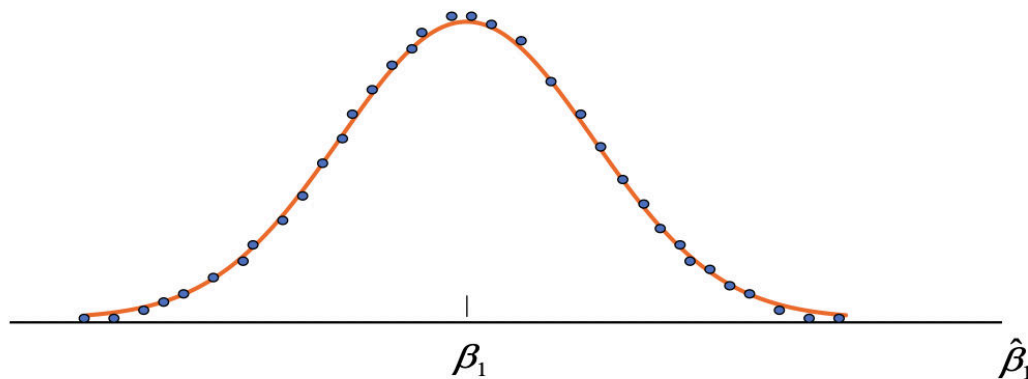
□ ALWAYS use heteroskedasticity-robust standard errors.

## 6. NORMALITY ASSUMPTION

This assumption implies that the population error  $u$  is independent of the explanatory variables  $x_1, x_2, \dots, x_k$  and follows the distribution  $u \sim \text{Normal}(0, \sigma^2)$ .

This implies that if we could draw many samples of size  $n$  and estimate linear regression model by OLS with each of these samples plotting the estimated  $\hat{\beta}_1$  that we get in every case, those follow a normal distribution centred at  $\beta_1$ .

Figure n. 5<sup>6</sup> (Garcia-Gomez, 2022)



With large samples  $\hat{\beta}_1$  always approximately follows a normal distribution centered in  $\beta_1$ , so even if  $u$  does not follow a normal distribution, OLS estimator is asymptotically normally distributed (i.e., approximately normally distributed in large samples).

□ When we have large sample sizes, we can always carry on using the standard tests for hypothesis testing (only IF we have them).

On the other hand, if the sample is small and we have non-normal errors, we need to worry about the normality assumption.

TO SUM UP:

- We need the first 4 assumption (MLR1-MLR4) to obtain unbiased estimates of the population parameters.

- The fifth and sixth assumptions (MLR5–MLR6) are important for inference, but we may face some problems with them. Anyway:
  - Even if assumption MLR5 is not satisfied, standard errors and test can be easily adjusted.
  - When our samples are large, the non-normality of the errors is not an issue.

Said that, now we have all the components for hypothesis testing.

## Lecture 5 – Linear regression analysis: inference I. – single population parameter

We have seen in our main example using STATA that with an increase of €1.000 derives an increase of €16.000 in the average house price of the neighbourhood, *ceteris paribus*.

Is this effect different from 0? 16 should be far enough from 0.

But what about 15? Is it far enough from 16?

Those questions are hypothesis, which we could test comparing how different our estimates coefficient is to the number we want to test (0 or 15). We are going to see how different it is with a statistical tool/testing.

### INFERENCE

First, we are going to test a hypothesis (the null hypothesis  $H_0$ ) for a single population parameter  $H_0: \beta_j = \beta_0$ .

To test this hypothesis (that the population parameter is equal to 0 or 15) we compute a t-statistic:

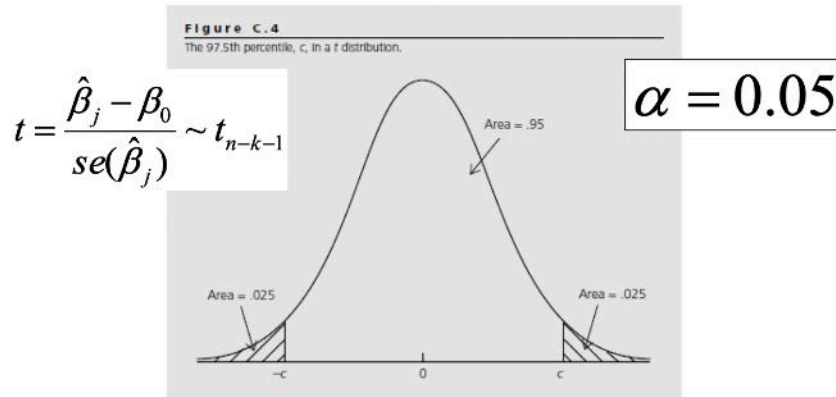
$$t = \frac{\hat{\beta}_j - \beta_0}{se(\hat{\beta}_j)} \sim t_{n-k-1}$$

The t-statistic is computed comparing the estimated population parameters minus the value that we want to test divided by the estimated standard error. The standard error gives us an idea of the level of precision of our estimate, and this is related to the fact that we have a random sample of the population.

Under the null hypothesis we know that this number is distributed following a t-distribution in which  $n$  is the number of observation and  $k$  is the number of explanatory variables.

The distribution looks like the following figure below.

Figure n. 6<sup>7</sup> (Garcia-Gomez, 2022)



The idea is that under the null hypothesis the t-statistic is going to be very close to 0. The more we go away from 0 the more we go towards the tails of the distribution, the less likely it is that our null hypothesis is true.

How further away are we going to accept a number, in order not to reject our null hypothesis?

We need to set a significance level ( $\alpha$ ), which is the tolerance for a type I error: it is the probability of rejecting the null hypothesis given that it is true. The most common values for the significance level  $\alpha$  are 0.10, 0.05 and 0.01. For example, having a significance level of 5% means if the null hypothesis was true then only 5% of the random samples would provide an estimate in the area of the tail of the distribution, 0.025 on the left and 0.025 on the right. If our estimated value falls in these areas, it is very unlikely that our null hypothesis is true, therefore we reject it. We will never be certain, so we base this in this probability that it is less likely that it could happen.

We reject the null hypothesis if the absolute value of the t-statistic is larger than the critical value  $c$  for our significance level (1.96 for a 0.05 significance level):

$$|t| > c$$

We often use the *p-value*: which tells us what is the largest significance level at which we could carry out the test and still fail to reject the null hypothesis<sup>8</sup> (Garcia-Gomez, 2022). We reject the hypothesis if:

$$p - value < \alpha$$

The p-value gives us the probability that we have left on the tails given our t-statistic<sup>9</sup> (Garcia-Gomez, 2022).

## CONFIDENCE INTERVALS

To avoid having only a point estimate, we use confidence intervals, which provide us a range of likely values for the population parameter.

We know that:

$$\frac{\hat{\beta}_j - \beta_0}{se(\hat{\beta}_j)} \sim t_{n-k-1}$$

Looking at the graph above, we can ask ourselves what the value of  $\beta$  that corresponds to  $-c$  and  $c$  is.

We know that testing the null hypothesis for any of those values we could not reject them because they fall within the area in the centre of the distribution (the 95% of it).

We can use this formula, where we add and subtract at the estimated parameter the value  $c$  (the critical value in the distribution) times the standard error we estimated for our  $\beta$ :

$$\hat{\beta}_j \pm c se(\hat{\beta}_j)$$

### Interpretation of the confidence intervals

Important: we are not 95% sure that the real value is in this interval, because, in fact, we are not, but it means that:

- From all the possible samples that we can possibly draw, in 95% of the cases the true value of the coefficient will be inside the interval.  
We just hope that our random sample is one of those containing the real value.
- Using a two-sided hypothesis, it gives us the set of all values that cannot be rejected (in this case at 5%) (the values that are in the centre of the t-distribution between  $-c$  and  $c$ ).
- Again, this is NOT equivalent to the probability of 95% of having the real value inside this exact interval.

## Lecture 6 – Inference II

## Testing multiple restrictions: multiple hypotheses test or joint hypotheses test

- We are interested to see if a group of variables has no effect on the dependent variable. So, we are going to test this.

For example, are the house prices affected by the demographic structure?

$$houseprice = \beta_0 + \beta_1 income + \beta_2 perc_{young} + \beta_3 + perc_{elderly} + u$$

We have multiple explanatory variables:

Our null hypothesis is that  $\beta_2$  is equal to 0 as well as  $\beta_3$ :

- $H_0: \beta_2 = 0, \beta_3 = 0$   $\square$  if true, the demographic structure will not have an effect in houses prices.

We have two restrictions: if  $H_0$  is true, then this group of variables has no effect on house prices after we controlled for the income per capita.

- $H_1: H_0$  is not true  $\square$  AT LEAST ONE of the coefficients is different than 0, or either one of them, but is sufficient that only one of the coefficients is different than 0.

To test these hypotheses, we define an F-statistic in which the null hypothesis  $H_0$  is that the two coefficients are equal to 0; the alternative hypothesis is  $H_1$ .

- $H_0: \beta_1 = 0, \beta_2 = 0;$
- $H_1: \beta_1 \neq 0 \text{ and/or } \beta_2 \neq 0$

$$F = \frac{1}{2} \left( \frac{t_1^2 + t_2^2 - 2\hat{p}_{t_1 t_2} t_1 t_2}{1 - \hat{p}_{t_1 t_2}^2} \right)$$

Note that this formula is robust to heteroskedasticity errors.

- $t_1$  is the t-statistic of test  $\beta_1 = 0$
- $t_2$  is the t-statistic of test  $\beta_2 = 0$
- $\hat{p}$  is the estimator of the correlation between these two t-statistics.

Once completed we know that the F-statistic in large samples follows an F distribution.

$$F - statistic \sim F_{q, \infty}$$



( $q$  is the number of restrictions; infinity refers to the fact that we use a very large sample)

As usual, we reject the null hypothesis if  $F$  is large in statistical terms. If  $F >$  critical value of  $F_{q, \infty}$  we reject the hypothesis.

- At 10% the critical value could be 2.30.
- At 5% the critical value could be 3.00.
- At 1% the critical value could be 4.61.

We test these hypotheses with STATA.

With STATA we see that the  $F$  value is 568,95 (a lot more than 2.30), therefore we reject the null hypothesis.

In alternative we can look at the  $p$ -value, which in this case is smaller than the 5% significance level, and we get to the same conclusion.

We can use the same procedure to test a linear combination of the parameters:

$H_0: \beta_2 = \beta_3$ ;  $H_0: \beta_2 - \beta_3 = 0$  □ test command in STATA after running the regression model.

When  $H_0$  is not rejected, we must say that: WE FAIL TO REJECT  $H_0$  AT THE  $x\%$  SIGNIFICANCE LEVEL

We cannot say that  $H_0$  is accepted at the  $x\%$  significance level, but only that it is not rejected because it can assume different values.

## Lecture 7 – Interpretation and categorical variables

**CATEGORICAL VARIABLES:** variables that contains qualitative information (for example, the variable “*religion*” is a categorical variable).

**BINARY/DUMMY VARIABLES:** categorical variables that can only take two categories.

For example, we could ask ourselves if, *ceteris paribus*, house prices in Rotterdam would be different from the rest of the Netherlands. In order to do so, we can create one of these two variables:

- Variable Rotterdam (it assumes the value 1 if the neighbourhood is in Rotterdam, 0 if in other parts).
- Variable Other (1 if in other parts, 0 if in Rotterdam).

And we can estimate this model:

$$\text{House price} = \beta_0 + \beta_1 \text{income} + \beta_3 \text{Rotterdam} + u$$

Or, alternatively, this other model:

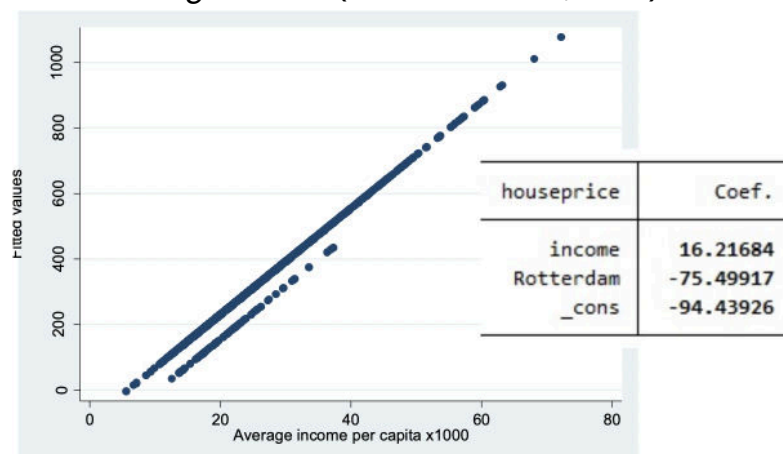
$$\text{House price} = \beta_0 + \beta_1 \text{income} + \beta_4 \text{other} + u$$

The coefficient  $\beta_3$  is representing the effect on the average house price of being in Rotterdam compared to other regions, *ceteris paribus*.

We cannot include both "Rotterdam" and "other" in the same equation because there is perfect collinearity between them, since every time "Rotterdam" is 1 "other" is 0, and vice versa  $\square$  there is no separate variation in one variable compared to the other.

If we add "Rotterdam" and "other" in the same equation, the sum will always be 1, which is the same as the constant.

Figure n. 7<sup>10</sup> (Garcia-Gomez, 2022)



In STATA we can choose a model or the other: the conclusions will be the same, even if we get different coefficients for each equation, and that's because the constant is

giving us different values in each of the models (in the first model it tells us the expected house price for neighbourhoods not in Rotterdam, while in the second it tells us the expected house price for neighbourhoods in Rotterdam).

We then can look at the predictions. We have 2 lines of dots which tell us how house prices change when income does in a neighbourhood that is in Rotterdam (one line) and for a neighbourhood which is in another region (the other line).

Which one is the highest and which one is the lowest? We could look at what is the expected value as we did in the previous passage  $\square$  this will give us the line with the lowest intercept.

Are they parallel? It seems like they are. It is because we assume that an income variation has the same effect in Rotterdam as well as in any other region  $\square$  the slope is always the same.

Sometimes the categorical variable has more than two categories.

It is important that all the categories contained in the variable could be interpreted in a quantitative way, and we create a dummy variable (the variable with two variables) for each of the categories.

But if we have, for example, four categories, we cannot include all of them in the same equation because of the collinearity. Therefore, we use three categorical variables in one equation and take the one we left behind as the reference category. Once set the reference category, we must interpret the coefficients of the included dummy variables compared to the reference category. When interpreting the intercept, we also need to consider it compared to the excluded variable and remember that the intercept is the value our model assumes when the INCLUDED explanatory variables are equal to zero, but not the reference category.

We interpret the coefficients compared to the reference category because if we changed it, the estimated coefficients would also change.

## Lecture 8 – Model selection in linear regression analysis

What could be the variables to include in a model?

We should begin with a (theoretical) framework:

We could, for example, examine the probability of committing a crime: Gary S. Becker elaborated in 1968 an economic model in order to describe the grade of participation in a crime by an individual.

The individual participation in a crime (our dependent variable  $y$ ) could be described by the following variables: the hourly “wage” obtained through a criminal activity ( $x_1$ ), the hourly wage obtained with a legal job ( $x_2$ ), other income ( $x_3$ ), the probability of getting caught ( $x_4$ ), the probability of being convicted if caught ( $x_5$ ), the expected sentence if convicted ( $x_6$ ), the individual’s age ( $x_7$ ).

After choosing a model framework and having seen the data, it is possible to see which variable are available. If we cannot observe some of them, they will be part of the unobserved.

We must not start with bivariate associations, because it could be that the  $\beta$ , if the zero conditional mean assumption is not satisfied, is biased.

So, we could choose a variable from a bivariate association, but it is going to be biased in most of the cases.

What is the point of the analysis?

There are two possible alternatives:

1. The aim could be having the best predictive model. In this case we can use **goodness of fit** measures.
2. If the aim is to estimate the causal effect of  $x$  on  $y$ , the goodness of fit measures will not be informative: in this case it is needed to control for sufficient confounders to satisfy the zero conditional mean assumption.

First, we analyse the case where our aim is to have the best predictive models. As we said, we can use goodness of fit variables.

These variables are useful to know how well our model fits the data, which means how much of the overall of the dependent variable could be explained by our independent variables.

The standard measures of goodness of fit are the following:

- $R^2$
- Adjusted  $R^2$
- Overall F-test

## R<sup>2</sup>

$$R^2 = \frac{SSE}{SST} = 1 - \frac{SSR}{SST}$$

Where:

$$\text{TOTAL SUM OF SQUARE: } SST = \sum_{i=1}^N (y_i - \bar{y})^2$$

$$\text{EXPLAINED SUM OF SQUARES: } SSE = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2$$

$$\text{SUM OF SQUARED RESIDUALS (UNEXPLAINED VARIATION): } SSR = \sum_{i=1}^N u_i^2$$

R<sup>2</sup> captures how much of the total variation (which is the total sum of the square), is explained by our model.

We can define it in two ways, either as the ratio of the explained sum of the squares divided by the total sum of the squares, or as 1 minus the sum of the square of the residuals (that is the unexplained variation) divided by the total variation.

This number is going to be between 0 and 1 and the closest it is to 1, the highest the share of the variation that our model explains.

Our problem with this measure is that it always increases with the inclusion of variables in the model: every time we add one explanatory variable, even if it explains very little, is going to cause the R-squared to increase.

## ADJUSTED R<sup>2</sup>

Our aim is to have a model which is both simple and complete: we want to control for sufficient x's, but at the same time we should not overdo the model, therefore we should add variables to the point they provide a sufficient contribution in our predictions.

$\bar{R}^2$  (adjusted R<sup>2</sup>) corrects the increase we discussed above by penalising for the number of coefficients: if we add a new variable with a |t-stat| > 1 or, in the case of a set of added variables, F-stat > 1, it increases.

## OVERALL F-TEST

Another useful measure is the **overall F test**: with this measure we can check whether the joint hypothesis that the coefficients of all our variables is equal to zero, compared to the alternative that at least one is different than 0. This is exactly the

same F-test we have seen in the second inference lecture, even if in that case we used it for all the variables.

The unrestricted is the model that includes as many restrictions as explanatory variables in our model. This tells us whether our model explains anything at all, so at least one of the variables is statistically significant.

Overall F test of  $H_0: \beta_1 = \dots = \beta_k = 0$  vs  $H_1: \text{at least one } \beta_j \neq 0 \text{ with } j = 1, \dots, k$

- Unrestricted model (UR):  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$
- Restricted model (R):  $y = \beta_0 + u$  (with  $k$  restrictions)

We can do this for all the variables as well as for sub-groups of variables:

F test:  $H_0: \beta_1 = \dots = \beta_q = 0$  vs  $H_1: \text{at least one } \beta_j \neq 0 \text{ with } j = 1, \dots, q$

- Unrestricted model (UR):  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$
- Restricted model (R):  $y = \beta_0 + \beta_q + 1x_q + 1 + \dots + \beta_k x_k + u$  (with  $q$  restrictions)

This is going to add information if those variables are jointly significant. We can test for this using an F-test in which we test this sub-group of variables. This is also a way of deciding whether we want to keep those variables in the model or not. Because if they are jointly statistically significant there is no need to have those variables in our model.

## Lecture 9 – Exogeneity

Unbiasedness of OLS (recap):

$$E(\hat{\beta}_0) = \beta_0 \text{ and } E(\hat{\beta}_1) = \beta_1$$

To be able to say this, we need the four assumptions:

- Linearity in the parameter (MLR1)
- Random sampling (MLR2)
- No perfect collinearity (MLR3)
- Zero conditional mean assumption, i.e.,  $E(x) = 0$  (MLR4)

The MLR4 is the assumption we need to focus on the most: does it hold in our model?

Let us consider a model with the explanatory variables  $x_1$  and  $x_2$ :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

For the assumption MLR4:  $E(x_1, x_2) = 0$   $Cov(x_1 = 0)$   $Cov(x_2 = 0)$

The error term is independent of  $x_1$  and  $x_2$ , so when they change, the error remains constant.

This assumption does not hold if:

- Is present an incorrect or misspecified functional form: the model is missing, for example, powers of  $x_1$  or  $x_2$ , or, always for example, using  $y$  in level whether we should be using its logarithmic form.
- There is a correlation with other unobserved factors which are part of  $u$ .

Why  $E(x_1, x_2)$  is different from 0 when we miss nonlinearities?

- True model:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_2^2 + u$
- Estimate:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$

Where can we find  $x_2^2$ ?

Everything that has an effect on  $y$  and is not in our model, is contained in the unobserved part of the equation.

So, is  $u$  independent, not correlated with our  $x$ 's?

$x_2^2$  is correlated with  $x_2$  so in that sense the zero conditional mean assumption does not hold.

We proceed in the following order:

1. We need to estimate a model of interest (for example, a model with two explanatory variables  $x_1$  and  $x_2$ )

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

□ we can test for functional form misspecification using the RESET test. In the RESET test we first estimate the model of interest (the one that we think is the true model) and then obtain fitted/predicted values.

2. After the first step, we obtain fitted/predicted values:  $\hat{y}$

3. Re-estimate the model adding powers of  $\hat{y}$  as independent variables.

Example in the Wooldridge textbook:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \delta_1 \hat{y}^2 + \delta_2 \hat{y}^3 + u$ .  
 Here the coefficients of the added powers are equal to zero, so there is no evidence of misspecification.

Test implement in STATA:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \delta_1 \hat{y}^2 + \delta_2 \hat{y}^3 + \delta_3 \hat{y}^4 + u$

4. F-test of joint significance of added powers of:  $\hat{y}$
- o If **insignificant**: there is no evidence of misspecification
  - o if **significant**: there is evidence of misspecification  $\square$  we reject the model.

Powers are nonlinear functions of the  $x$ 's. Significance means that model is missing some important nonlinearities (they were contained in the error term).

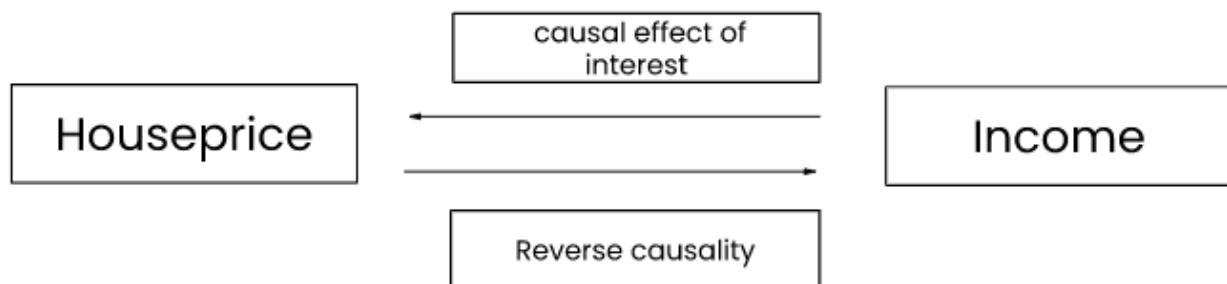
So, if those estimated deltas are statistically significant, it means that the model is missing important non linearities. We should add powers or interaction again, and re-estimate and repeat the process. So, the RESET test does not tell us what to do, but if there is a problem with misspecification in our model. It also does not tell us what to do in case of rejection, we have to try a different specification.

The other case in which the zero conditional mean assumption does not hold up, is when we have correlation with others unobserved factors.

If we take a simple model:

$$\text{Houseprice} = \beta_0 + \beta_1 \text{income} + u$$

We are interested on the fact that the average income has on the average price in a neighbourhood.



However, it could be possible (we are not sure yet) that if a house price in a neighbourhood is higher, the income is higher.

If that is the case, if  $y$  has an effect on  $x$ , we say that there is reverse causality.



If there are other factors such as: amenities in the neighbourhood, social capital, the size of the houses and other factors that have an effect on the average house price and on the income, if we estimate the model by OLS, we are estimating a mixture of all this instead of just the effect of interest.  $\beta_1$  is going to catch not only the effect of interest, but also a mix of all of this.

If that is the case, the fourth assumption does not hold, i.e. income is said to be endogenous. OLS do not estimate causal effect consistently.

An explanatory variable is endogenous when it is correlated with the error term, either because there is reverse causality, or because there are other factors correlated with  $x$  and with  $y$ . and in that case OLS do not estimate causal effect consistently and what we obtain is an estimate of the association of the two variables; for example, in this case,  $\beta_1$  could give us the association between income and the average house price.

There are possible solutions.

The first one, when we do not have reverse causality, is to think about other characteristics that are not included in the model but have an effect both on income and house prices and estimate a more complete model:

$$\text{houseprice} = \beta_0 + \beta_1 \text{income} + \beta_2 \text{rotterdam} + \beta_3 \text{density} + \beta_4 \text{perc\_young} + \beta_5 \text{perc\_elderly} + u$$

But is it enough? Sometimes could be, but often no.

We could try to include more relevant variables but difficult to account for all relevant variables that influence income and house prices.

There is a gold standard for estimating causal effects, which is to randomize the explanatory variable of interest.

For example, in our example, it consists in randomly allocating houses to people to see if there are different effects on house prices and income; people cannot choose where to live. Another idea could be re-allocating income in a similar way.

But in reality, we just cannot randomly tell some people to buy some houses and to others not, and we as well cannot take people's income and randomly distribute it across the population.

There are other possible solutions:

- Instrumental variables

- Panel data

We are going to look at them in future lectures.

## Lecture 10 – Beyond statistical significance

This lecture is relevant not only for the linear regression, but for the entire course as well.

We start by asking the following question: is only significance important for us?

### statistical significance and economical significance

Is important to distinguish the statistical significance and the economical significance, or whether the magnitude of the estimate is really relevant.

If the aim of our research is to test a hypothesis, we point out that a not statistically significant result is also a result.

For example, if we are interested in seeing if a new intervention has an effect, if we find out that we do not reject the null hypothesis that the effect is 0, this is also a result per se: we could be interested to know if we are making bad policies, for example.

Significant does not mean certain: we base our analyses on the p-value and on the confidence intervals □ if we take a random sample of the population, many time of the same size, in the 95% of the cases the true population parameter is going to be there. However, we are not sure if in this specific case our sample is there. The fact that we do not reject the null hypothesis or that we do reject it does not mean that we are accepting a given hypothesis.

*Significant* does not mean neither *relevant* nor *important* nor *substantial*: the size matters.

If we have a huge sample size, we are going to detect very small effects. However, if the effect is really small, we may just conclude that an intervention is not really needed, as the effect is so small that it could also be 0, even if it is statistically significant.

Our ability to find significant results also depend on other characteristics of the data like the sample size, how good are our measures, how much noise do we have. We always have to see how much uncertainty there is around our estimate  this is going to come from the confidence intervals.

Example: We want to estimate the returns to education using OLS<sup>II</sup> (Garcia-Gomez, 2022):

$$wages = 8 + 2.5 * years\_education$$

The wages are equal to €8 + €2.5 per year of education.  
Focusing on the coefficient of years\_education:

If we focus on the coefficient of years\_education:

- P-value = 0.01; p-value = 0.0001

The coefficient is statistically significant in both cases.

But if the P-value = 0.30 the coefficient is statistically insignificant.

When could this happen?

If we start from a very small sample, we get the same point estimate, but the standard errors are going to be large, but then, again, we get an estimate of 2.5. Once we get larger and larger, the standard errors become smaller and the confidence interval as well. In this way we will obtain a more precise estimate of the effect.

If we are interested not in testing the hypothesis, but in whether there are returns to the education or not (so we do not care about testing this null hypothesis), we should be aware that there is a certain amount of uncertainty around our estimate.

For example, we make three regressions and have to make a decision based on the last one. Can we decide to get another year of education?

We get the 2.5 but we should not think only at the statistical significance. Instead, we should think at how much our wage is, and if it is going to increase either in relative or absolute terms thanks to this 2.5 increase. If €2,5 per year seems a good relative increase, it might be worth to add a few years of education at least.

economic significance (or practical significance or relevance)

We have to think if it is a large increase ( for knowing this we have to see the magnitude of the coefficient).

The increase of €2.5 per year seems like a large increase, but what if, estimating a model, our  $\beta$  is going to be equal to 0.01 and the p-value to 0.01? It is not a statistically significant effect. Can we conclude that it is a relevant effect if our wage increases by €0.01 for every additional year of education? The effect is very small and close to 0 independently of this effect not being a statistically significant effect.

So, when we look at our results, we need to go beyond saying whether they are statistically significant or not. It is important to look at the uncertainty around our estimate and the confidence intervals. It is equally important as well to think about the magnitude. We can then look at the estimated coefficient, the absolute effect, but it also is useful to compare it with another magnitude, like average values, to inform about relative effect.

## validity of results

When we think about the validity, the first question is thinking about our study per se, our samples, our analysis, and whether the study provides causal estimates for the population and setting of our study.

Does the assumption to get a causal effect hold? That refers to **internal validity**.

We talk about **external validity** when we want to use the conclusions from our study and then extrapolate them to other populations and settings.

For example, in our example we were looking at the house prices, and the income in neighbourhoods in the Netherlands. Can we extrapolate those conclusions to other countries? This will depend on how the housing market works in the different countries.

We can also think that those results were using data from 2012; can we extrapolate those conclusions to nowadays?

We can then think on how the housing market has changed. Are there reasons to think that the association between the household income and the household prices has become stronger or less strong given the dynamics over the past years?

When we have our results and it is the moment to write the conclusions, what evidence should be of our interest?

We should look for the absolute and relative effects about the importance, about what those number mean, whether they are **statistically significant or not**, the **amount of uncertainty around that estimate**, the presence of **potential biases**

(maybe our sample was not completely random, for example, and this type of discussion regards the **internal validity** of our studies, or if we extrapolate the conclusions we also have to think about the external validity.

There could be then other dimensions like, for example, the heterogeneity of effects, so whether the effects for different groups of the population are different. For example, if we think about the returns to education, we expect the returns to education to be different for men and women, for people with different ethnic backgrounds, people from different socioeconomic groups. If this is the case, we can also analyse this heterogeneity of the effects.

# applied microeconometrics – module 2 – endogeneity and instrumental variables estimation

## Lecture 11 (2.1) – Introduction

### causality and ceteris paribus

The economist's goal is to infer a ceteris paribus relationship or to make sure that one variable has a causal effect on another variable. Another goal is knowing what is behind this causal effect.

We must be careful when we think we have found an association, because it may be tempting and we could think we had found one even if we actually don't.

#### **Example**<sup>12</sup> (Garcia-Gomez, 2022)

Imagine an alien observing human behaviour. He sees that some people go to the hospitals and others do not. He compares the outcome of those two groups of humans and conclude that people that go to hospitals are more likely to die. So, the alien, trying to improve human's well-being decides to close all the hospitals to save human lives.

He was inferring from an association, a causal effect, but was missing some very key information: for example, humans that go to hospitals are more likely to die beforehand because their health was worse.

We have to focus on this distinction, between associations and causal estimates □ when do OLS provide causal estimates?

□ We introduce a new instrument, the **instrumental variables regression**, that allows us to estimate a causal effect in certain conditions under some assumptions. We will discuss when those assumptions hold and if we have the right data to do so.

For the following lectures, our example will focus on what is the effect of retirement on depression<sup>13</sup> (Garcia-Gomez, 2022).

We see that in many countries, governments are increasing the retirement age due to the pressure on public budgets. This can alleviate how much government spend on old age pensions. But we could also wonder whether these types of policies have a spillover to other budgets.

### **Could these choices have an adverse effect on people's health?**

If people have to work, if there is a negative effect on our health because we have to work longer, we are going to have a decreased productivity, or it could be that we end up retiring, leaving the labour force through other pathways like disability benefits.

We can also think that retiring has a bad effect on our health because we are more likely to lose our cognitive ability because we practice less, some of us gets bored, and then this could have a negative effect on mental health. If that is the case, increasing the retirement age can also have a positive effect on health □ it all comes up to it being an empirical question.

Our data come from a large European health survey, the Survey of Health, Ageing and Retirement in Europe (wave 1, 2004)<sup>14</sup> (A. Börsch-Supan, coordinator, 2004).

The information come from 11 countries including individuals and their spouses that are older than 50.

The dependent variable of interest is a depression scale, the EuroD (from 0 to 12). Higher values of the depression scale are associated with a worse mental health.

Our main explanatory variable of interest is the variable "retired" (it takes value 1 if the individual is retired and 0 if not).

We will also control for other explanatory variables like age (in years), marital status and the education level, whether the individual has a low, medium, or high educational attainment.

If we estimate the regression model on STATA, we estimate an OLS with robust standard errors and then we ask ourselves what the effect of retirement on depression is. We see that the estimated effect, or association (we do not know yet what is), is positive, which means that the mental health of the individuals worsens when they retire compared to not being retired, and the magnitude is 0.7 points □ a retired individual has a mental health score of 0.7 points higher compared to a non-retired individual, *ceteris paribus*.

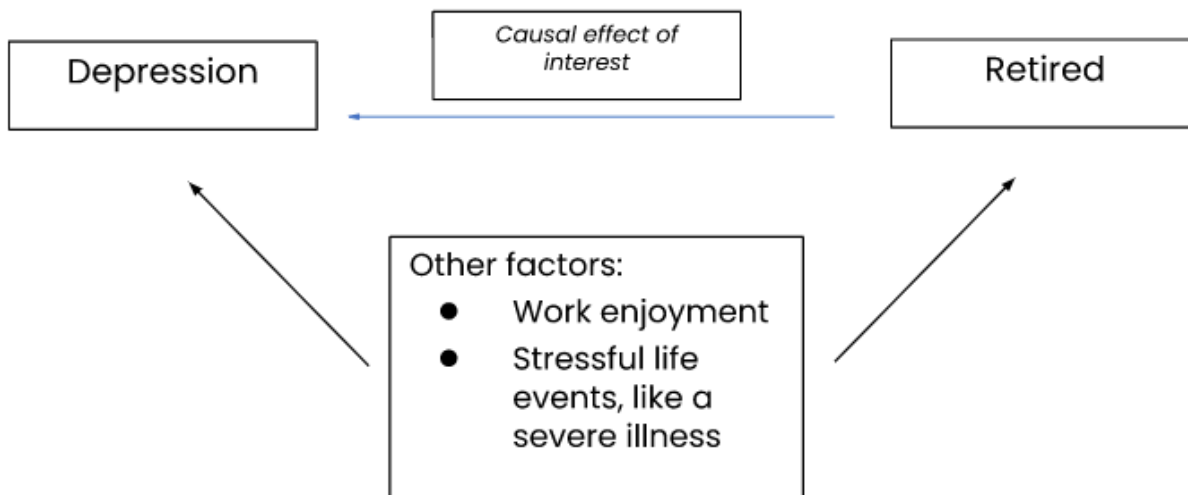
Is this a true causal effect or an association? For answering this question, we need to remember the zero conditional mean assumption: it will not be a causal effect if the zero conditional mean assumption does not hold. This happens when we have an incorrect or misspecified functional form; this is something that we can test with the RESET test, or when we have a correlation with other unobserved factors that are part of  $u$ . In this case our zero conditional mean assumption will not hold if there are unobserved characteristics that determine mental health and are correlated with retirement status, age, educational attainment, or marital status. If those unobserved characteristics are not correlated, the zero conditional mean assumption will hold, and if they are correlated, it will not.

So, we have to discuss whether there are correlations with other unobserved factors that are part of  $u$ ; we will see that this will happen if we have omitted variables bias, collider bias and reverse causality.

## Lecture 12 (2.2) – Omitted variable bias

The following graph shows how the following causal relationship works:

$$y = \beta_0 + \beta_1 \text{retired} + \beta_2 \text{married} + \beta_3 \text{age} + \beta_4 \text{education} + u$$



What can be present in the unobserved term?

Other factors such as the individual's level of enjoyment of its work, and meaningful stressful event of its life, like a serious illness such as cancer.

Enjoying work influences only mental health, or does it influence the probability to retire too? If the answer to this question is affirmative, then this variable is going to be correlated with retirement, and it is also going to be a determinant of depression.

Another relevant factor in this example are stressful life events: are they also likely to have an effect on whether you retire or not?

Some yes and other no.

A friend having a bad accident does not have an effect on the probability that we retire, then it is not a factor of interest for us, even if it is a stressful life event. On the other hand, there is evidence that having a severe health condition (like a cancer) increases the probability of retirement.

If this also affects our mental health (but not through retirement) then this is also going to have an effect on depression. Omitting any of those variables creates an **omitted variables bias**: we have an omitted variables if those variables are correlated with  $x$  and are also the determinant of  $y$ .

## omitted variables bias

If we have an omitted variables bias, it could go in both directions. It could be:

- An **upward bias**
- A **downward bias**



There are going to be situations where we will not have this additional factor, otherwise, we would put this in the model. We could argue about this bias being positive or negative:

- To say that a bias is positive is the same as saying we have an upward bias: in these cases, the estimated coefficient (the value we obtain from our regression model) is larger than the population value (that we do not know: we get an estimate, and we can argue that the estimate is larger than the population value).
- To say that a bias is negative is the same as saying we have a downward bias: in these cases, the estimated coefficient is smaller than the population value.

Whether we get a positive or a negative bias, depends on the sign of the correlation of the omitted variable  $x_2$  with  $x_1$  (endogenous variable), the variable that in our model is correlated with the unobserved characteristics. It also depends on the sign of the correlation of the omitted variable and  $y$ , the dependent variable, as well.

	$Corr(x_1, x_2) > 0$	$Corr(x_1, x_2) < 0$
$\beta_2 > 0$	Positive bias	Negative bias
$\beta_2 < 0$	Negative bias	Positive bias

<sup>15</sup> (J. M. Wooldridge, 2014)

When the sign of those two correlations (between  $x_1$  and  $x_2$ ) is positive, and when  $\beta_2 > 0$  (which means that the omitted variable has a positive effect on our dependent variable), we have a positive bias.

We also have a positive bias when those two are negative. In any of the other two combinations we have a negative bias. For example, we get a negative bias when the correlation between the endogenous variable and the omitted variable is negative, so a correlation between  $x_1$  and  $x_2$  is negative, and when the correlation between the omitted variable and the dependent variable is positive. So, the coefficient of  $\beta_2$  would be positive if we could have this variable in our model.

Let us think about this in our example.

What would be the bias of retired if we do not control for education? We assume for the moment that there are no other things that could be contained in  $u$ .

The true population regression is the following, in which we control for retired marital status, age and education:

$$y = \beta_0 + \beta_1 \text{retired} + \beta_2 \text{married} + \beta_3 \text{age} + \beta_4 \text{education} + u$$

But if in our data the variable "education" is missing, the model that we can estimate is a model in which we can control to estimate the effect of retirement on mental health.

But we can only control for marital status and for age:

$$y = \beta_0 + \beta_1 \text{retired} + \beta_2 \text{married} + \beta_3 \text{age} + u$$

Endogenous variable  $x_1$ : retired

Omitted variable  $x_2$ : education

	$Corr(x_1, x_2) > 0$	$Corr(x_1, x_2) < 0$
$\beta_2 > 0$	Positive bias	Negative bias
$\beta_2 < 0$	Negative bias	Positive bias

(Wooldridge, 2014)

We must look at what is the correlation between  $x_1$  and  $x_2$  and what would be the effect of education on mental health (we remember that higher the score for mental health the worse it is).

What is the expected correlation between retired and education?

Do we expect highly educated individuals to retire earlier or later? We could expect this correlation to be negative: we expect that highly educated individuals to have less physically challenging jobs, and this could allow them to retire later.

But somebody else could have a different expectation and suppose the correlation to be positive instead of negative.

Anyway, for us the correlation is negative. Then the bias would be negative or positive? For answering this we should reflect on  $\beta_2$ , so we need to think about what the expected correlation between education and (bad) mental health is.

Do highly educated individuals have worse or better mental health?

For us the correlation is negative: highly educated individuals have a better mental health  $\square \beta_2 < 0$

So, if  $\beta_2$  is negative and the  $Corr(x_1, x_2)$  is negative as well, according to our expectations, we would expect a positive bias. With different expectations, we could have a different bias.

It is recommended to reflect on our expectations and argue why we expect each of those two relationships to be positive and negative.

In the data we have, we could run an estimate on the model with education and also on the model without education.

We can then compare the coefficient of retired in the two models. The estimated coefficient in the second model (the one where we exclude the variable education), it is larger than the expected coefficient in the model that includes the education variable.

So, if the model with education is a good representation of the population model, we see that the coefficient of retired is smaller than when we exclude education.

$\square$  the model without education is a model that suffers from omitted variable bias, and then we see that we could have an upward/positive bias, because the number is higher.

It seems that probably our expectations were correct.

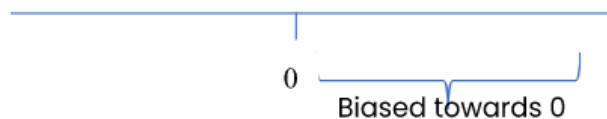
If we look at the coefficients, we cannot directly say anything about what is the correlation between education and retire. We can do that once we have access to the data. But we can see that the coefficient for edmed and edhigh are negative, so people with medium and high education have a mental health score that is lower than individuals with low education  $\square$  there is a negative correlation between education and bad mental health.

We can have upward/positive bias either if the estimated coefficient is positive or negative  $\square$  the estimated coefficient is larger than the population parameter and it does not really matter if it is on the positive or negative range of values. The same happens when we have downward/negative bias.

- Upward (or positive) bias:



- Downward (or negative) bias:



**Biased towards 0:** the estimated coefficient  $E(\hat{\beta}_1)$  is closer to 0 than the population parameter  $\beta_1$ .

It is important because if we know that our estimate is biased towards 0, we can argue this is a lower bound of the estimated effect.

If we are in the range of an upward/positive bias, our estimate is going to be biased towards 0 if it is negative.

If we are in the range of a downward/negative bias, our estimate is going to be biased towards 0 if it is positive.

Should we always add more controls?

If we can control for omitted variables, if we are able to observe that they affect both  $x$  and  $y$ , we should add them to our model to be able to get this causal effect, a ceteris paribus conclusion.

We could think to add omitted variables that affect only  $y$ . By doing this, we will increase the goodness of fit, but adding those variables to the model is not needed for ceteris paribus conclusions.

Adding irrelevant variables (not correlated with  $y$ ) is not a problem to get ceteris paribus conclusions, but it will have a cost in terms of efficiency: we will have a less efficient estimator  $\square$  it is better to leave them out than to include them.

Sometimes adding an additional control introduce a collider bias  $\square$  it must not be done!

## Lecture 13 (2.3) – Collider bias

The collider variable is a variable that will introduce bias if we control for it<sup>16</sup>.

It is different from what we have seen in the previous lecture, when we were worried about variables that we did not have.

The collider is a variable in our dataset that we may think it is positive to include as part of our explanatory variables, but if we include it, it will make things worse.

We will see some examples using **Directed Acyclical Graphs (DAG)**, similar to the graphs previously used.

In these graphs the arrows indicate that there is a causal effect  $x \rightarrow y$

" $\square$ " is a causal effect

In these graphs we also include all the common causes affecting the two variables.

The first example of a collider bias is the selection bias, which is a specific type of collider bias.

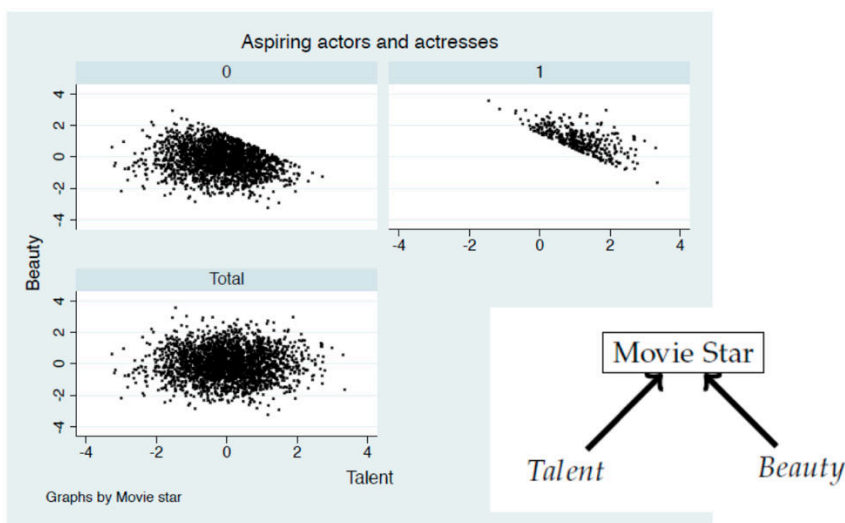
### THE SELECTION BIAS

What is the importance of talent and beauty into being a movie star?

Is there a masked relation between talent and beauty in this industry?

Figure n. 1

(Cunningham, 2020 [https://www.scunning.com/causalinference\\_norap.pdf](https://www.scunning.com/causalinference_norap.pdf))<sup>17</sup>



Let us study every aspiring actor and actress.

We can measure for each of them their beauty on a scale and what is their talent on a scale centred to zero.

We can plot all these data in a scatter plot.

Is there any relationship between these two variables? Looking at the third image in the figure n. 1 seems that there is no relationship. The ones who make it into the industry are those who are on top of the distribution of beauty and talent.

We can also observe the distribution among those who do not make it and those who make it, respectively the first and second images of figure n.1.

Focusing on the ones that do make it, it seems that beauty and talent are negatively associated: the ones who are more talented are less beauty and the other way around.

But if we do not split the distribution and select our sample, we would not find any association.

This selection creates this spurious correlation which is not really in the overall data □ this is what we call a selection bias

Now, let us return to the retirement and depression example. We assume that the DAG represents the true effect, and all the variables are dummies.



There is no arrow between retirement and depression □ There is no real effect between the two variables

The third variable, obesity, is caused by both retirement and depression.

For example, when people retire, they change eating habits, physical activity habits, etc. This will have an effect on obesity. We can make the same reasoning for depressed people.

- If retirement = 1  it increases the chances of becoming obese (because people move less)
- If retirement = 0  people are less likely to be obese
- If depression = 1  people are more likely to be obese
- If depression = 0  people are less likely to be obese

What happens when we control for obesity?

What is the relationship between retirement and depression when obesity is equal to 0 and when is equal to 1?

This is what we do when controlling for the effects of other variables.

If we run a regression between retirement and depression, there would not be any association.

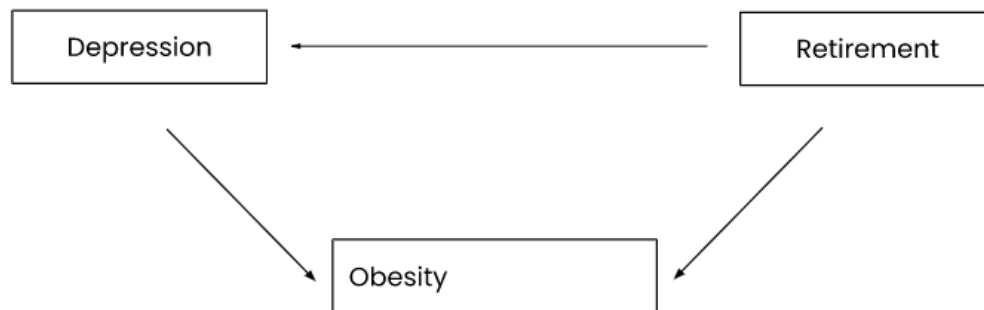
When obesity assumes value 1, it is likely that the person is retired because if it is, it is more likely that he is obese, and it is likely that he also is depressed, because both retirement and depression cause obesity.

When obesity assumes value 0, we have less retired people and less depressed people.

Once we control for obesity it seems that there is a relationship between retirement and depression. If we estimate now this OLS, we will likely get a positive estimated coefficient between those two variables.

- Controlling for obesity creates a bias
- Obesity is a collider because it creates a collider bias; this is a type of bias that exists because we introduce a variable in our model that is caused by both the explanatory and the dependent variables.

If the DAG would have been:



The reasoning would have been the same: we could get an estimate that it is not the true causal effect when we control for obesity because it is introducing the collider bias.

We could get, in this DAG, a causal effect by just including the variables “retirement” and “depression”.

So, in this example, it is better not to include obesity in the model, because is a collider and introduces the collider bias in our estimates.

How do we select then which variables to include in our model?

There is not any existing formula that does it for us. We have to reflect on the theoretical model or the previous evidence and if there is any correlation between the two variables already known in the scientific world.

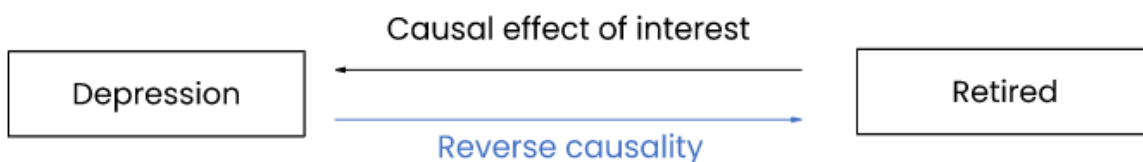
Then we have to reason about what is the expected relation between the variables. In this regard a DAG can help us because it illustrates how the different variables are related to each other.

We will need to make assumptions, as it is required every time we estimate a model and use estimators: it is all based on assumptions. We should also be able to defend those assumptions.

## Lecture 14 (2.4) – Reverse causality

We currently studying the effects of retirement on mental health, and we also are controlling for marital status, age, and education.

$$y = \beta_0 + \beta_1 \text{retired} + \beta_2 \text{married} + \beta_3 \text{age} + \beta_4 \text{education} + u$$



The lower arrow indicates reverse causality.

To better understand: if our mental health worsens, we are more likely to retire.

In this case, OLS will underestimate or overestimate the effect if reverse causality is present.



□  $\beta_1$  will be biased □ the expected value will be different from the population parameter.

If depression has a negative effect on retirement, the causal effect will be underestimated.

If depression has a positive effect on retirement, the causal effect will be overestimated.

In this specific case we expect depression to have a positive effect on retirement, therefore the explanatory value of the estimated parameter is going to be larger than the true population parameter.

How can reverse causality cause a bias?

When we have an omitted variable bias, it is easier to see the reason of the correlation between the error term and the explanatory variable.

With reverse causality it is more complicated, but easy to see using a couple of equations.

$$depression = \beta_0 + \beta_1 retired + \beta_2 married + \beta_3 age + \beta_4 education + u$$

The reverse causality implies that we can estimate a model like the following:

$$retired = \gamma_0 + \gamma_1 depression + \gamma_2 married + \gamma_3 age + \gamma_4 education + v$$

Retired is the dependent variable and is described by the mental health, marital status, age, and education.

The zero conditional mean assumption means that in the first equation there is no correlation between retired and all the other variables and the error term. It means that we have to check if the covariance between  $u$  and retired is equal to zero or not.

So  $cov(u, retired)$  needs to be equal to 0.

But we can also write  $cov(u, retired)$  as following:

$$cov(u, \gamma_0 + \gamma_1 depression + \gamma_2 married + \gamma_3 age + \gamma_4 education + v) = 0$$

But it is impossible that  $cov(u, \gamma_0 + \gamma_1 depression + \gamma_2 married + \gamma_3 age + \gamma_4 education + v)$  would be equal to 0 because we have the  $u$ , which represents all the things that explain depression other than retirement, marital status, age and education. So, this is going to be correlated with depression and this covariance cannot be equal to zero.

In this case the zero conditional mean assumption does not hold.

What should we do now?

The model where the explaining variable of interest is endogenous (without any possibilities of having this fixed) could still be very useful. Having causal effect is always preferable, but in many cases we will only be able to get an association.

So, we must recognize that we cannot measure causal effect but simply associations. These associations often are a combination of the effects of third factors, reverse causality and the real causal effect.

This will have an effect on how we interpret our results. We have seen in our example that the zero conditional mean assumption does not hold, so we will not say that retiring will increase mental health for sure, but we need to express our interpretation in terms of an association.

In our example, we can say that on average the mental health status of retired individuals is higher than the one of non-retired individuals by 0.7 points, after controlling for the other variables.

We cannot draw *ceteris paribus* conclusions as we cannot be sure that all the other things do not change at same time, and we mention the variables we control for as they are included in our model.

## Lecture 15 (2.5) – Estimation and interpretation of instrumental variables

### instrumental variables

One of the tools we can use to estimate the causal effect is an instrumental variable, if we have the appropriate data and the required assumptions are satisfied.

With an instrumental variable we are able to isolate exogenous variation. If we think about the variation in any given variable, it consists of two parts:

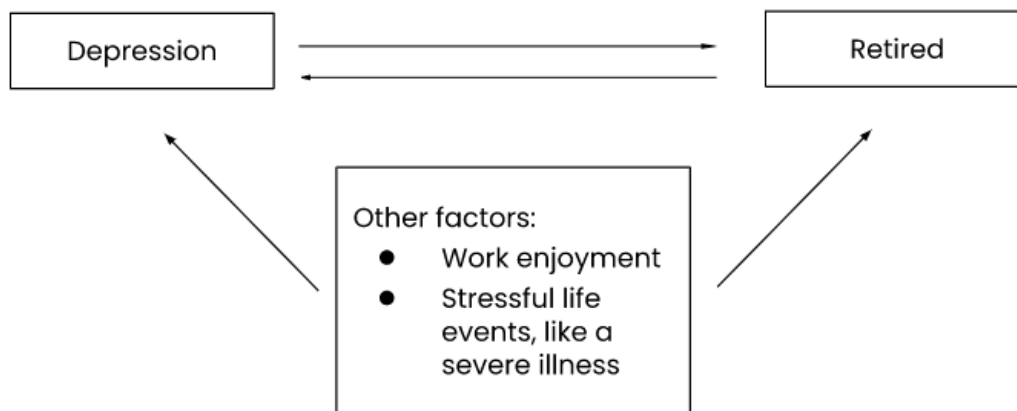
1. **Endogenous variation:** the part of the variation that is correlated with the unobserved.

2. **Exogenous variation:** the part of the variation that is independent of the unobserved.

If we think about all the reasons people retire, there are reasons that are endogenous to the individual, so they are related to what we cannot observe. Examples of this could be, if that person is enjoying his/her work, whether that person has a health impairment, etc. All these variables are related to mental health, and it also explain why people do retire. Maybe there are other reasons that motivate people to retire, related to retirement, but not related to the mental health of the individual. These reasons are what we have identified with exogenous variations.

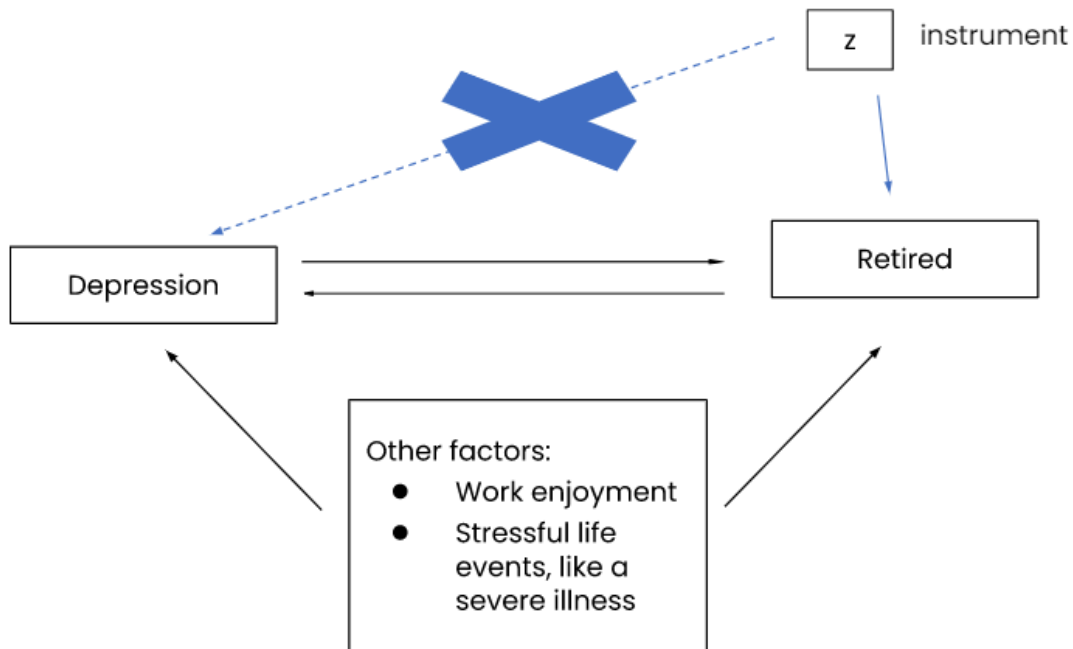
So, we have to isolate the exogenous variation to focus on the causes of why people retire and then use them to estimate a causal effect.

If we look at the following graph:

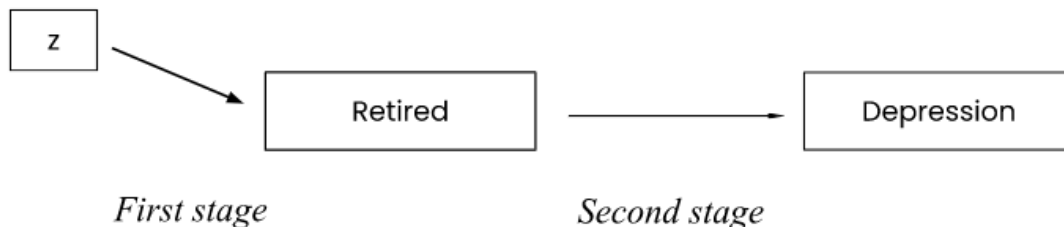


Mental health can also have an effect on the probability that someone retires. Then we have other factors that have an effect both on retired and depression.

The idea behind this instrumental variable is that this variable is able to capture some information that explains why some people retire (this variable is our instrument). This variable has not to be correlated to depression, it has only to influence retirement. This does not have an effect also on the unobserved characteristics.



The instrumental variable estimator looks first at the effect of this instrument on the probability that someone retires. This is the first stage. Once we get these predicted reasons on why people retire for exogenous reasons, then we use this exogenous variation to look at the effect of retirement on mental health. This is the second stage.



We will first look at the first arrow in the figure and then at the second. What could be a variable that we could put instead of  $z$ ? What could explain that people retire but it is uncorrelated with depression? In many countries, once people reach the early retirement age or the normal retirement age they do retire at that specific age. If we plot for any country the probability that someone retires, there will be a spike at those specific ages.

We can use this variation whether someone is over their normal retirement age, for example, to predict whether someone is retiring or not, and then use this exogenous variation to look at the effect of retirement on depression.

To do this we can use a **two stage least squares estimator**.

In the first stage:

$$Retired = \pi_0 + \pi_1 edmed + \pi_2 edhigh + \pi_3 married + \pi_4 age + \pi_5 full + v_2$$

We estimate retired as a function of education, marital status and age, with an additional variable, the full, which is the age at which you are entitled to *full* retirement benefits in the European countries.

We have variation in this variable because it varies through different countries, and even within some countries there are different retirement age for men and women (at least at the year of the data).

From the first stage we get the predicted retirement:  $\hat{retired}$

We then insert this predicted value in the second stage:

$$y = \beta_0 + \beta_1 \hat{retired} + \beta_2 edmed + \beta_3 edhigh + \beta_4 married + \beta_5 age + error$$

In the second stage we estimate by OLS. This model has the same explanatory variables as before, but instead of having the observed retired we have the predicted retired.

Given that in the first stage we control for educational status, educational attainment, the marital status and age, the variation that is picking from retired comes from the fact that someone has reached their normal retirement age, the *full* age.

We can estimate by OLS the first stage and then the second stage. If we were to do this, the standard errors in the second stage could not be correct because we need to consider that retired is not the observed retired, but only a prediction.

Stata does everything together to correct the standard errors in one go, so we do not have much to worry about. We only need to tell Stata which variable is our endogenous variable, and which is the instrument. We can also ask for heteroskedasticity robust standard errors as we did in OLS.

In Stata we first see the first stage regression where we see that we have this variable full. We will not interpret this output because our dependent variable "retired" is a 0-1

variable and we have not seen yet how to interpret such an outcome (we will do this when we talk about binary models).

What we cannot really say is that we see that full is statistically significant in the first stage, and that those who go over this age are more likely to be retired. This is how we can interpret the sign of this coefficient.

Moving to the second stage we have all the variables and the variable retired. Comparing its coefficient to the number we had before, the coefficient is larger and what it, still, suggest, is that those that retire have a worse mental health status compared to those that do not retire. We can interpret this coefficient in the same way as we did before: the mental health index increases by 2.2 points (on a scale 0-12) when an individual retires compared to being non retired, *ceteris paribus*.

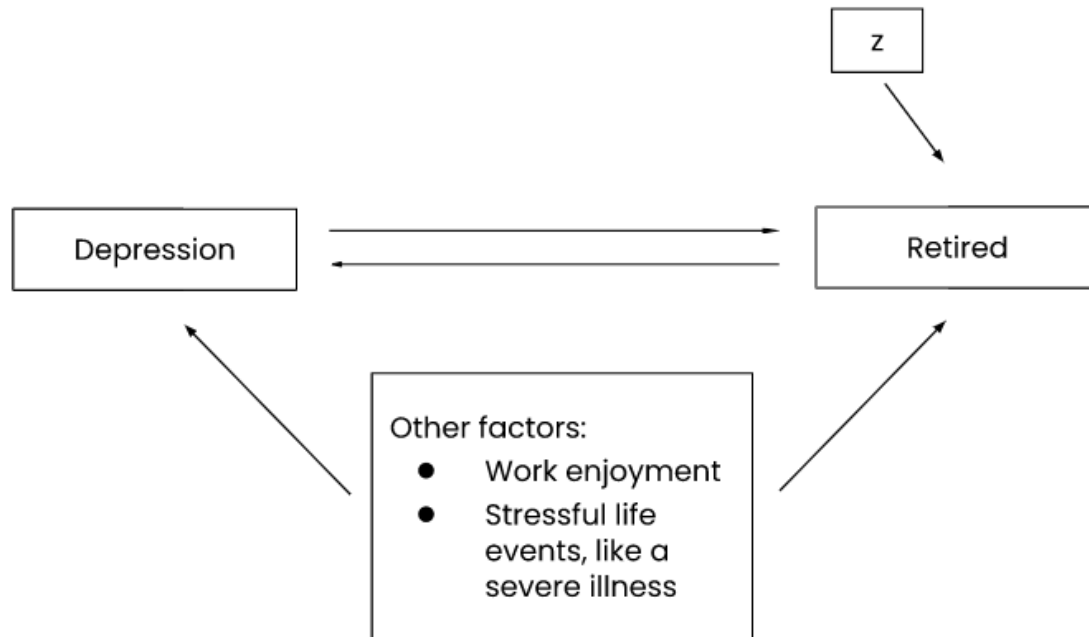
We can see that this effect is statistically significant. If we look at the confidence interval, we see that now it has broadened. One of the things that happen with instrumental variables is that because we only use a small part of all the variation, the estimate is going to be less precise, and it is important to see how much less precise this estimate becomes.

We can use the same method if we have more than one endogenous variable in the model. We can use instrumental variable regression in the same way (we only have to estimate the first stage for each of the endogenous variables and we would need more than one good instrument; we should have one for each of the endogenous variables). In fact, we always need the number of instruments to be larger or equal to the number of endogenous variables. But it is hard to usually find this exogenous source of variation, so it will be nearly impossible to be able to control for more than one endogenous variable in each model. But if, somehow, we can manage this, once we have the second stage, we interpret the coefficients in the same way we used to do when we had an OLS.

## **Lecture 16 (2.6) – Instrumental variables assumptions**

Our model of interest is a model in which we want to estimate the effect of retirement on depression, and we want to use an exogenous variation that it is not related to depression causing retirement, or the influence of other omitted variables.

In order to do that we need to find an instrumental variable that is correlated with the endogenous variable.



The assumptions are the following:

- The instrumental variable must be correlated with the endogenous variable:

$$Cov(\text{retired}, z) \neq 0$$

In this case we say that it is a **relevant instrument**.

We need not only a correlation between the two variables, but a strong correlation. They need to explain a large part of the retirement decision: we need a **strong instrument**.

- The instrument must not be correlated with any other (unobserved) determinants of depression:

$$Cov(u, z) = 0$$

In this case the instrument is exogenous. So, our instrument can only affect the dependent variable through the endogenous variable, after controlling for all the other variables in our model. After this we have a **valid (exogenous) instrument**.

## relevance assumption

The relevance assumption is that the covariance between the endogenous variable and the instrument is different from zero:

$$\text{cov}(\text{endogenous } x, z) \neq 0$$

This is something that we can conclude from looking at the first stage regression, where we predict our endogenous variable, using the instrument, so we can see whether the instrument is statistically significant or not. If it is not, in the second stage we will get very large instrumental variable standard errors. It is important that we have a strong correlation.

If we look at the first stage on Stata at the coefficient of full, we see that it has a t-statistic of 9.81 and that the p-value is very small too  $\square$  we can say that this instrument is highly significant.

The more variation in the variable  $x$  our instrument explains, the more variation will be available in the regression. Then we will have more variation to be used to look at the effect of retirement on depression. There is a little problem when we explain very little (so when our instrument is not strong), but we have a great instrument. In that case, the IV are no longer reliable.

We can only use instrumental variables if we have a strong instrument.

How do we check for weak instruments?

There is a rule of thumb that we can use when we have one endogenous variable: we must look at the F-test of the first stage, and it has to be larger than 10; if it is smaller, it suggests that the instrument is weak.

## validity assumption (or exogeneity assumption)

This assumption refers to the correlation between our instrument  $z$  and the unobserved component. Any assumption that we make about the relationship between the unobserved components and any other variable is an assumption that we have to argue, and we will not be able to fully test it.

There are some partial tests, but they will not be covered here because none of those tests will be able to tell us if our instrument is valid or not, and they are often misused in the literature. Some people claim that passing this test means that the instrument



is valid, but then we could find other reasons for which this assumption could not be satisfied.

We have to use economic theories, our own reasoning and previous evidence to argue about the validity assumption.

In this case, for example, to justify why using the normal retirement age is exogenous or not, we can use our knowledge of the institutional setting and we know that the way in which normal retirement ages have been defined is independent of the mental health of the population.

## Lecture 17 (2.7) – Instrumental variables vs OLS

To choose between IV and OLS there are the following steps to follow:

- We have to use IV only if we have a valid and strong instrument.  
Before even running any IV regression models we must think about the assumptions of our instrument. Does it satisfy them? Is it valid? Is it not correlated with the unobserved term? Is it strong?  
If we have a valid but weak instrument, we will not get a reliable IV estimate.
- We need to have large sample size.  
When we have a valid and relevant instrument, IV is consistent, but under endogeneity IV is biased. So, we must rely on this consistency property that only applies to large sample size. If it is small, we are no better off using IV instead of simple OLS.
- We choose IV if our explanatory variable is endogenous; if it is exogenous, we do not gain anything by using IV, but only lose information because IV are inefficient and do not use all the variation that we have in our explanatory variable. We throw away information if we want to gain in terms of getting a causal estimate, even though it is not a gain if we do not get a more reliable estimate.  
Our IV is going to be better the more highly correlated our instrument is with  $x$ . So, then we are going to get a smaller variance of the IV estimate. The stronger the instrument, the less we will lose in term of efficiency in terms of IV.

We have already discussed the first and the second conditions, but how do we determine if an  $x$  is endogenous or not? We do not observe the error term, and we do

not observe the unobserved characteristics, so we cannot just check the correlation between the unobserved and  $x$  to determine whether is endogenous or not.

We can use our instrumental variable estimator. To inform about this endogeneity. We can use the same logic behind IV to test for the exogeneity of retired:

1. We first estimate the first stage equation:

$$Retired = \pi_0 + \pi_1 edmed + \pi_2 edhigh + \pi_3 married + \pi_4 age + \pi_5 full + v_2$$

2. From the equation we get the residuals:  $\hat{v}_2$ . What are they picking up? The variation in retired that is not explained by the educational attainment, marital status, age, and our instrument. This is the potentially endogenous variation.
3. We then can estimate our main equation with OLS, adding residuals of step 2:

$$y = \beta_0 + \beta_1 retired + \beta_2 edmed + \beta_3 edhigh + \beta_4 married + \beta_5 age + \delta_1 \hat{v}_2 + error$$

We have our main equation and the residuals that are picking the potentially endogenous variation of retired. If this variation is correlated with our dependent variable, then  $\delta_1$  is going to be statistically significant.

4. We test the significant of  $\delta_1$ : if it is insignificant there is no evidence of endogeneity and if it is significant there is evidence of endogeneity. It can be shown that if we do this procedure the  $\beta_1$  we could get for retired, is also similar to the one that we get if we use the IV.

Stata can do all the steps above all at once. The null hypothesis is that the variables are exogenous. We get to different statistics, similar in 99.9% of the cases. In any of those cases the p-value is very small, so we reject the null hypothesis. Is retired endogenous? If we reject the null hypothesis that retired is exogenous, so this evidence that suggest that retired is endogenous.

If retired was endogenous this would be the third case in which we would prefer IV. We have been able to argue about the two assumptions in the previous lectures: we have a large sample size, and our explanatory variable is endogenous: this is a case in which we would prefer IV.

Many times, we will not have a valid and strong instrument. The golden rule for getting a causal effect is randomization, but this is not always possible neither,

because it is very expensive or because there are ethical issues (imagine we want to look at the effects of going to prison on the probability of committing future crime). But there ethical about putting people in prison. We could also see the effects of providing hospital care versus not, or even retirement; we cannot tell people to retire at random.

Another solution we can apply is the use of panel data methods.

## Lecture 18 (2.8) – Potential outcomes and DAGS

We will look at two useful frameworks to analyse causal effects:

- DAGs
- Potential Outcomes Framework

### DAG (Directed Acyclic Graph)

DAGs are the graphical representation of a chain of causal effects<sup>18</sup> (Garcia-Gomez, 2022) in which the nodes are the representation of random variables. The causal effect between two variables is represented by the arrows, which also indicate the direction of causality. If an arrow goes from variable A to variable B this arrow indicates the direction in which we assume that the causal effects work out. If there are no arrows, we assume that there are no causal effects between the two variables.

The causal effect between the independent variable  $x$  and the dependent variable  $y$  can be direct but also indirect: in this latter case we say that the causal effect is mediated by a third variable ( $x \rightarrow D \rightarrow y$  which means that the variable  $x$  influences the variable  $D$ , which has then an effect on the variable  $y$ ). if we think about the causal effect of income on health, there are reasons why an additional euro per se could improve our health, but there are many more reasons on why income may indirectly have an effect on health. For example, higher income makes it possible to have a better quality, or it can grant us the access to health insurance, etc. All these effects are not direct but are mediated by third variables.

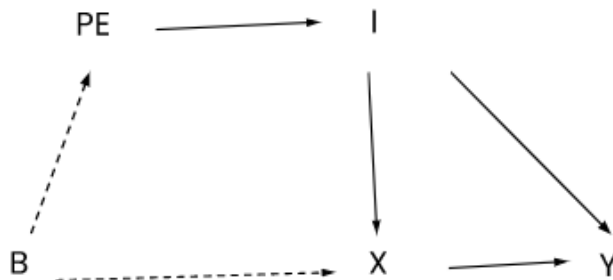
DAGs are useful to graphically represent the relationships variables and to acknowledge the assumptions and the data available in our specific setting. If  $x$  and  $y$  are connected by a solid arrow, we know that we have this information that we observe. If instead we have a dotted line, we know that it is an information we do not observe. Looking at these relationships should help us think whether this information that we do not observe is relevant or not to estimate a causal effect.

We are interested in knowing the returns on education, so how much does a schooling brings in terms of educational attainment.

We have the following variables:

- $X$  = educational attainment
- $Y$  = earnings
- $I$  = family income
- $PE$  = parental education
- $B$  = family background

### Case 1



How much an additional year of education translates into higher earnings?

In our theoretical framework, also family income plays a role: higher family income gives us easier access to certain positions, and that increases our earnings. A higher family income also affects our educational attainments because it allows us to attend better schools.

We also know that some family background characteristics influence both our educational attainment as well as parental education. Parental education has an effect as well on family income.

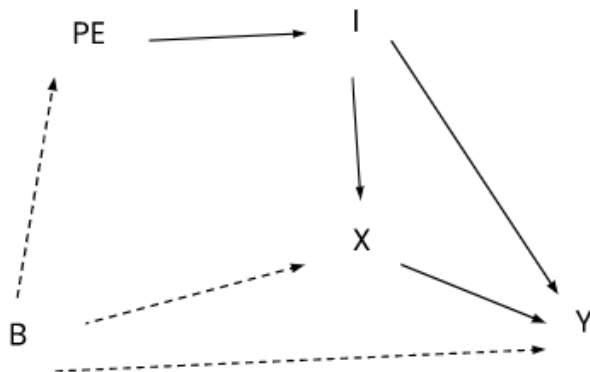
When we have the dotted line, we know that the family background information may be observable for individuals and their families, but not to us.

If we have this information, can we estimate the causal effect on educational attainment and earnings?

That is possible because when we observe for X and Y for family income, there are no other unobserved. This means that there is no other channel in which parental education or family background has an effect on the earnings. If this was the setting, it would be possible to estimate the causal effect.

If family background also influences earnings:

Case 2



If family background also has an effect on earnings, by controlling for X and I could not be sufficient to estimate a causal effect of educational attainment on earnings, because we could have this omitted variable bias as family background also has an effect on Y and X.

The DAG is useful to clearly illustrate the assumptions that we are making, the data we are able to observe and help us see whether we estimate a causal effect or if we are missing because of the effect of these additional variables. In other examples we could also have a collider. In any case it makes it impossible to get a causal effect. We get an unbiased estimate of the causal effect, and we can talk about association or partial associations.

## Potential Outcomes Framework (POF)

It is not a graphical interpretation, but a tool for when we are interested in the effect of an intervention or a variable on an outcome Y.

To keep it simple, let us think that our intervention is a binary.

For example, as we were looking throughout this topic of looking at the effects of retirement on depression, our treatment T could be retired.

If we think about the returns on education, for example we could think about our indicator as getting a college degree versus not. Treatment is a general term: it could be a medical term used in the medical field, but it can also be part of a training program, it can be retired, it can be getting college education □ it is the effect of the variable we are interested in.

If we think for every individual, we have two potential outcomes. Returning to the example of the returns to education, 1 could be the earning that that person could get if goes to college and 0 the earnings that that person does not get because he did not go to college.

Thinking about the effects of retirement on depression, for a given person in a given year, we are interested in the mental status of that person if retires and the mental status of that person if does not retire. Those are our two potential outcomes.

$$Y_i = \{Y_i(1) \text{ if } T_i = 1 \ Y_i(0) \text{ if } T_i = 0$$

We can then define the treatment or our causal effect as the differences in between those potential outcomes. So, the difference in our mental health, if you retire compared to not retiring. The difference in our health if we go to the hospital compared to not going. The difference of the earning if we go to college minus the earning that we could get if we do not go.

$$\Delta = Y_i(1) - Y_i(0) = (T = 1) - (Y|T = 0)$$

This is how our data could look like:

Individual	Treated	$Y(1)$	$Y(0)$	Causal effect
1	0	$Y_1(1)$	$Y_1(0)$	$Y_1(1) - Y_1(0)$
2	0	$Y_2(1)$	$Y_2(0)$	$Y_2(1) - Y_2(0)$
3	0	$Y_3(1)$	$Y_3(0)$	$Y_3(1) - Y_3(0)$
4	1	$Y_4(1)$	$Y_4(0)$	$Y_4(1) - Y_4(0)$
5	1	$Y_5(1)$	$Y_5(0)$	$Y_5(1) - Y_5(0)$
6	1	$Y_6(1)$	$Y_6(0)$	$Y_6(1) - Y_6(0)$

<sup>19</sup> (Garcia-Gomez, 2022)

$$ATE = E[Y_{j(1)} - Y_{j(0)}] = \frac{1}{N} \sum_{j=1}^N [Y_j(1) - Y_j(0)]$$

For a given individual, that individual could be treated □ we would have the outcome if that individual could be treated and the outcome if he could not be treated. The same if the individual retires or not: the causal effect is the difference.

If we want to estimate the average causal effect or the average treatment effect, we could just get the average over all these individual causal effects. So, we have the causal effect for individual one, two, and so on.

The fundamental problem to causal inference and to get a causal effect, is that we can only observe one of these two outcomes: we cannot observe the same person in two different states: at the same moment, I am either retired or not. It is not possible to look at both the states of the world, but we need to think about in which situation we are.

We want to know what the outcome of those non-observed could be, the counterfactual outcomes, which are the outcomes that could be observed if the person was in a different state (we have these gaps in our dataset). We only observe the outcomes of an individual given his or her current condition.

The goal of the analysis is to try to generate those counterfactuals.

	Outcome <i>without</i> Treatment	Outcome <i>with</i> Treatment
Control	$E[Y_i(0) T_i = 0]$	$E[Y_i(0) T_i = 1]$
Treatment	$E[Y_i(1) T_i = 0]$	$E[Y_i(1) T_i = 1]$

<sup>20</sup> (Garcia-Gomez, 2022)

If we think about going to the hospital/doctor example, in the data we see the output of those that did not go to the hospital (the control group), and the output of those who did go to the hospital (the treatment group). We need to be able to estimate what could be the expected output of those that did not go to the hospital if they had gone, and the other way around.

We need to ask ourselves whether our empirical method give us a valid counterfactual. If we think about the effect of going to the hospital, we can just compare the mortality outcomes of people who did go to the hospital and those who did not go to the hospital.

$$E[T = 1] - E[Y_i|T = 0]$$

The effect of going to the hospital is obtained comparing those two. If those that go to the hospital, or those that retire, or those who go to college, are different compared to those that do not receive the treatment, do not retire, do not go to college, then we cannot just use the information from this counterfactual group to estimate our causal effect.

When there are other observed or unobserved characteristics that influence both the treatment and  $y$ , in other words, our explanatory variable of interest and  $y$ , our method will not give us a valid counterfactual.

If those reasons why it is only due to other observed characteristics that confound this relation, then we can just add these explanatory variables to our model, and then OLS with this additional control explanatory variables will give us the causal effect, the average treatment effect, this right counterfactual.

However, if we also have unobserved characteristics that explain why people go to the hospital and at the same time why they die, or why people retire and at the same time their mental health, OLS provides bias estimates  $\neq$  we need to use a different estimator.

This is another way of thinking about what it means to get a causal effect. We should always check how the assumptions of the different methods are fitting in our research. At the same time DAG will help us visualize the assumptions we are making in our specific contribution.

# applied microeconomics - Module 3 - Introduction to empirical methods: panel data

## Lecture 19 (3.1) – Panel data: introduction

Which estimation technique we should use when working with panel data?

To answer this question, we need to know:

- What panel data is
- The notations of panel data
- The distinction in the notation from cross-sectional analysis



- Unbiasedness and efficiency
- Different estimation techniques, including:
  - OLS (and the conditions under OLS will be the best estimation technique with panel data)
- The within effects estimators, which refer to fixed effects, least squared dummy variables and first differences; we also want to know which assumption they require and how these estimators solve potential OLS issues
- Random effects and correlated random effects
- How to choose across the estimators
- Definition and implication of attrition for panel data

At the end of this module, we will also introduce difference in differences, which is an approach that can solve some of the issues that within estimators might not be able to.

## challenges of panel data

- We will introduce a set of new estimators, each with their own working and assumptions
- The choice of which estimator to use depends carefully on the topics of bias and statistical efficiency
- With panel data we now have the distinction between time-variant and time-invariant covariates

After this module we should be able to:

- Describe a panel data set and distinguish it from a cross-sectional
- Understand the assumptions required for each estimation technique and under which of those they are unbiased and efficient
- Knowing how to select the best a most appropriate estimation technique
- Evaluate and understand the importance of attrition in different panel data sets

## Lecture 20 (3.2) – What is panel data?

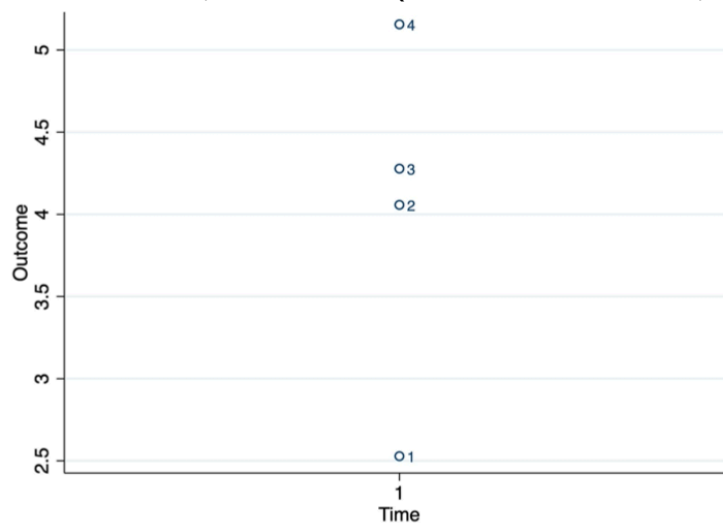
In research, data could have different structures:

- **Cross-sectional data:** the cross-sectional data are those data sets that have collected data on one observed sample of units at only one point of time.

- **Repeated cross-sectional data:** in this case the data collected are the of same type of the cross-sectional data, but on different observed samples and at different points of time. We are not following the same person, but we are collecting the same data.
- **Longitudinal/panel data:** the panel data is a data set that collects the same data on the same observed sample at different points of time □ we will collect the same data, the same variables on the same people and units over time.

In the following figure we can visualize cross-sectional data:

Figure 1 module 3, lecture 3.2<sup>21</sup> (Carlos Riumalló Herl, 2022)



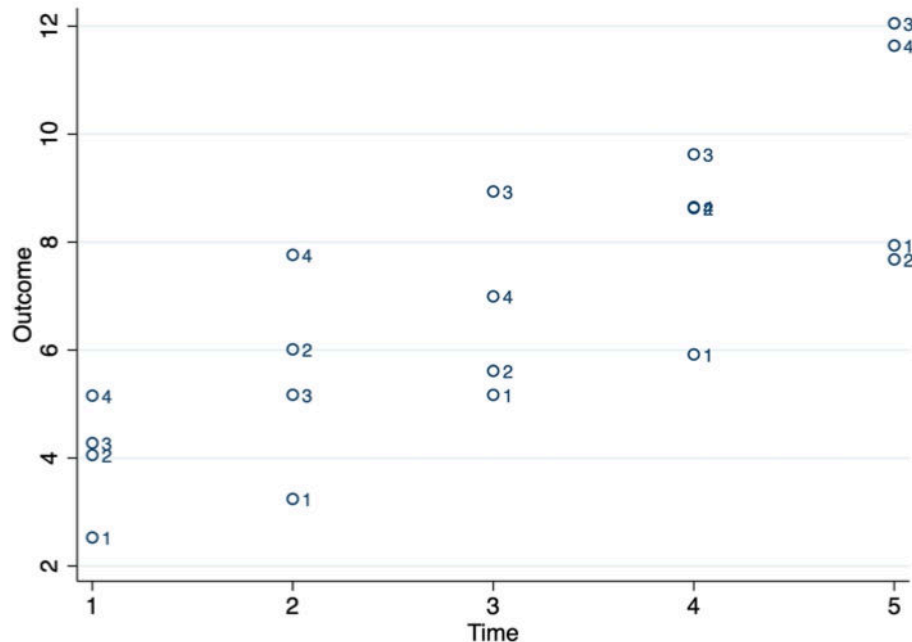
We can think of a cross sectional data as a snapshot on time.

Figure 2 module 3, lecture 3.2<sup>22</sup> (Carlos Riumalló Herl, 2022)



In the case of repeated cross section data (figure above), we are collecting the same data, so the same variables and the same information, but on different samples and in different time periods □ this could be imagined as two snapshots at two different people in two different points of time.

Figure 3 module 3, lecture 3.2<sup>23</sup> (Carlos Riumalló Herl, 2022)



With panel data we collect the same information on the same people over time.

## panel data characteristics

A panel data contains repeated observations of the same unit across time.

The data is now characterized by having two dimensions:

- **Individual unit dimension** □ indicates which units we are observing (the unit could be people, countries, stocks, schools, etc.) □  $i = 1, 2, \dots, N$
- **Time dimension** □ it gives us information on the period of time during which the data are collected. It could be years, month, seconds, etc. □  $t = 1, \dots, T$

Now we can define panel data in two ways:

- **Balanced panel dataset:** all units are observed in all time periods □ we observe all people all the time.

- **Unbalanced panel dataset:** in this (more frequent) case the number of observations vary across individuals □ whatever the reason, we are not able to follow each individual for the entire period of time.

Some of the examples of panel data includes:

- Longitudinal surveys of Ageing (HRS, ELSA, SHARE, etc.) □ in this case the data have been collected on people every two years.
- Administrative tax records □ some countries (including Denmark and the Netherlands) collect annual tax data and make their information available for research □ each person is followed every year.
- Hospital characteristics □ in this case units are not people, and in general units do not necessarily need to be people □ in this case hospitals are observed over a period of a year; we could observe their spendings, the number of operations they do and so on.
- Financial stocks □ we could follow what is happening at a stock in every point of time □ this is a panel data set.

Figure n. 4<sup>24</sup> (Carlos Riumalló Herl, 2022)

	Person	Time	Age	Male	Outcome
1	1	1	47	1	2.089202
2	1	2	48	1	4.522359
3	1	3	49	1	6.552687
4	1	4	50	1	5.527746
5	1	5	51	1	7.380176
6	2	1	47	0	3.993593
7	2	2	48	0	5.705793
8	2	3	49	0	7.461751
9	2	4	50	0	7.695805
10	2	5	51	0	8.791395
11	3	1	48	1	4.401231
12	3	2	49	1	7.169417
13	3	3	50	1	6.500992
14	3	4	51	1	10.39803
15	3	5	52	1	8.463883
16	4	1	47	0	5.891078
17	4	2	48	0	6.662832
18	4	3	49	0	7.206686
19	4	4	50	0	9.569715
20	4	5	51	0	9.225389

When looking into the data □ each observation of our dataset represents one unit time □ it represents a particular outcome of one unit at a particular time set

Looking at the figure n. 4, we can see that the first observation represents the observation for person 1 at time 1, while the second observation corresponds to person 1 at time 2, and so on.

So now we have multiple observations per unit, each one in different points of time.

When we are working with panel data, we always have two fundamental variables:

- The **identifying unit variable**, in our example the variable "person" □ indicates which observation belongs to the same unit □ in our example the variable person contains the number of the ID for the unit we observe (1) □ the first five observations belong to person one because they are identified as such.
- The **time variable**, in our example the variable "time" □ with this variable we can see which observation belongs to which time frame: in this example we have observations that go from time one to time five.

The fundamental characteristic of a panel data set is that we now have a series of observations that belong to the same unit. For each of the person in the panel data of figure 4, we have five observations in five different time points. In this case the variable person allows us to identify which observation belongs to each person.

With panel data sets we also have another two types of variables:

- **Time-variant variables**
- **Time-invariant variables**

In our example, the variable *age* is a time-variant variable, as the value of age changes for each person at each point of time, while the variable *male* is a time-invariant variable, as it contains the same value over the whole timeframe in which we are following people up.

□ The concept of time-variant and time-invariant might not be absolute, but might depend on the data set we are using.

For example<sup>25</sup> (Carlos Riumalló Herl, 2022), if we are using a survey that collects data on older people and we have information on things like education, those might not change for older people, but it could for young people. In the end, the concept of time-variant and time-invariant might depend on the dataset that we are using.

panel data structure in practice

Whenever we are looking to a variable, there will be variations in its values.

While in the case of cross-section data we already had variations, in the case of panel data sets we need to distinguish the variation in two parts:

- How much of the total variation is due to variation within units
- How much of the total variation is due to variation between units

As for the within variation, we look at a single individual and we are going to look at how much the value of a certain variable changes for that individual with regards to the unit's mean value for that variable.

□ The concept of within variation is relative to the unit.

In the case of between variation, we want to know how much variation of the variable exists between units. In this case we will calculate this by comparing the mean value across different units.

This analysis of the variation will tell us how much variation is occurring within individuals (how much that variable changes over time within units), or whether most of the variation exists between units.

We can ask ourselves: do people change a lot over time in reference to that certain variable, or is most of the difference due to differences between people for that variable?

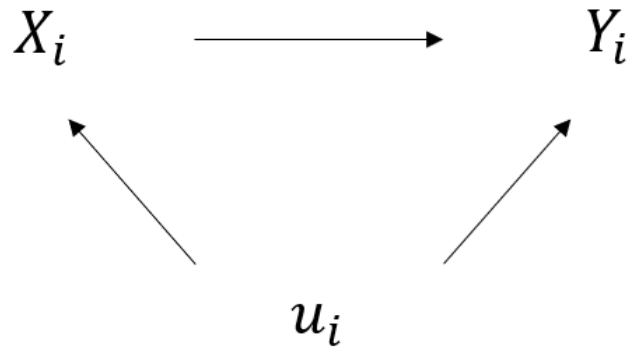
Do the units per se change over time for a given variable or the difference is given by differences between the individuals?

In the applied part it will be seen how to describe and interpret these information on STATA.

## Lecture 21 (3.3) – Notations

With panel data we also change the way we visualize some aspects of the models and the definitions.

This is a DAG for a cross-sectional data:



$X_i$  has an effect on  $Y$ ; we have then an unobserved error term that might affect both our covariate of interest,  $X$  and  $Y$ .

We can use the classical OLS notation:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

We can write  $y$  as a linear function of  $\beta_0$  and  $\beta_1 X_i$  and the error term  $u_i$ .

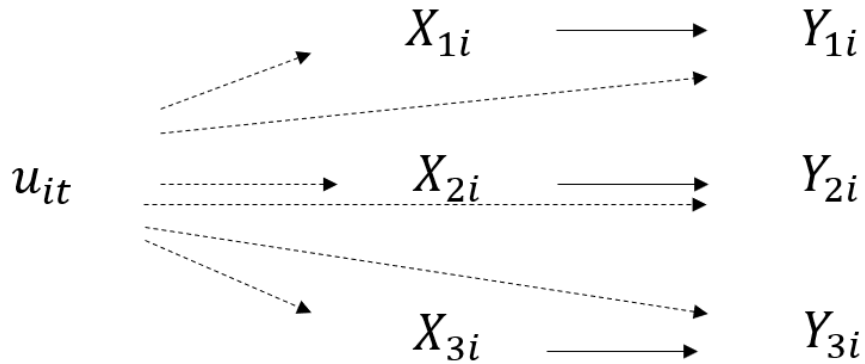
We can add more covariates, there will be only one outcome, defined for a unit  $i$  and in this case there is no time dimension.

In this case using OLS gives us an unbiased estimate on the parameter  $\beta_1$  if the expected value of the error term conditional on  $x$  is equal to 0:

□  $\hat{\beta}_1$  is unbiased if  $E(u_i | X_i) = 0$

With panel data we have another dimension: time.

How can this influence the DAGs and the type of notations we use?



For each individual  $i$  now we have multiple observations. In this case we have the effect of  $X$  on  $Y$  at three different points of time, all for the same individual,  $i$ .

In this scenario we also have the error term, which has the dimensions unit and time as well. The error term not only influence  $X$  and  $Y$  in one particular time period, but it can influence the variables for the entire period of observation.

We can think of a DAG for panel data as the cross-sectional diagram multiplied for all the time periods we have.

This is a simplified version of the DAG, because we can also consider past  $x$ s: for example,  $X_{1i}$  having an impact on  $Y_{2i}$  and vice versa.

□ Now we will have multiple  $X$ 's that will depend on the time period we observe, and now we have an error term that does not only influence  $X$  and  $Y$  of one time period, but of all time periods.

## panel data notation: decomposing the error

The outcome now is unit and time specific. What is its notation?

$$Y_{it} = \beta_0 + \beta_1 X_{it} + u_{it}$$

This is a simplified model: we can always add more covariates and eventually past covariates (this is called a "**lag**"). Now the error term, that varies both across unit and time, can be decomposed in two parts:

$$Y_{it} = \beta_0 + \beta_1 X_{it} + u_{it} = \beta_0 + \beta_1 X_{it} + \alpha_i + \varepsilon_{it}$$

Now the error term is distinguished into:

- $\alpha_i$  □ in the literature, this component is referred as the **unit heterogeneity**, or **individual fixed effect** or **unit fixed effect**. It is the part of the error term that is the same for all observations of one unit □ is a persistent component and the same for all observations of a unit. It is also unique to each unit and can distinguish the units between each other. It does not necessarily have to be different than zero, but if it exists it is time-invariant in the person □ this is the part of the error term that is time-invariant within units.
- $\varepsilon_{it}$  □ it is referred as the **idiosyncratic error**, and it is the time-variant component of the error term. It can be different for each unit and time period. We are referring to a part of the error term that is unique for each unit and time observation and for it we will create this idiosyncratic variation. The



idiosyncratic variation could be equal to zero, but does not have to, and in this case it will distinguish which types of estimation we can use and which types we cannot use.

**Example**<sup>26</sup> (Carlos Riumalló Herl, 2022): Let us study again what is the effect of education on income: with panel data, income (the outcome for and individual  $i$  at time  $t$ ) would be a linear function of the parameters  $\beta_0$  and  $\beta_1$  and the education of the individual  $i$  at time  $t$ .

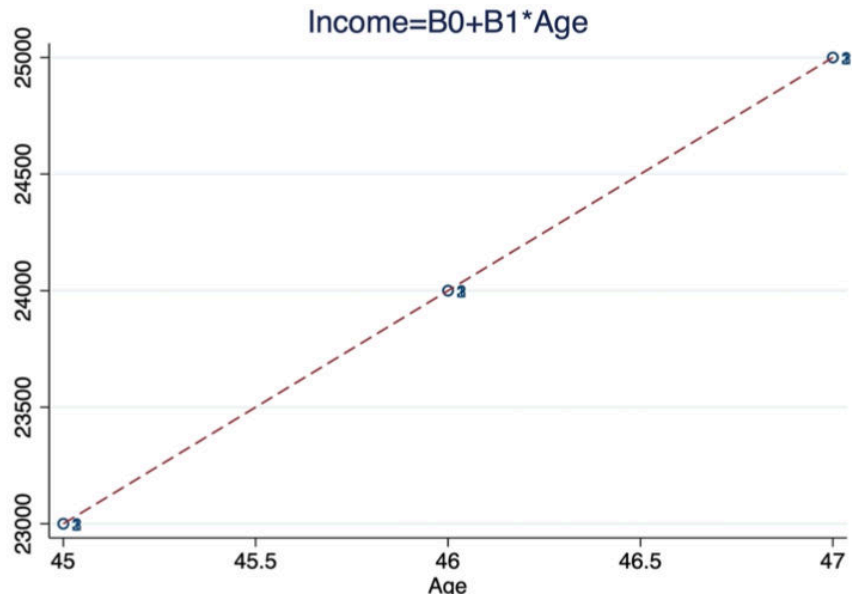
The data set for this example is the annual tax register, that allows us to follow individuals on an annual basis. From the register we have information on the education level as well as what type of income they receive.

We have the error term  $u_{it}$  (the error term for individual  $i$  at time  $t$ ) that can be decomposed in the time-invariant component (the individual heterogeneity  $\alpha_i$ ) and in the time-variant component (the idiosyncratic shock  $\varepsilon_{it}$ )

When we think about the time-invariant component of the error we think about all the possible unobserved factors that might influence education and income that are time-invariant. For example, in this case we could think about genetic material (as our DNA is time-invariant and it could also determine our capacity to become more educated or having an higher income). The genetic material makes each individual different and is time-invariant: it is an individual heterogeneity  $\alpha_i$  it makes each individual different in a time-invariant way.

We can then think about time-variant components that can influence our model and that could be shocks on a daily basis. Skills development could be one of these components: it could change for each individual at each time point, and it could be based on things like mental health or capacity to work  $\varepsilon_{it}$  these could be idiosyncratic shocks.

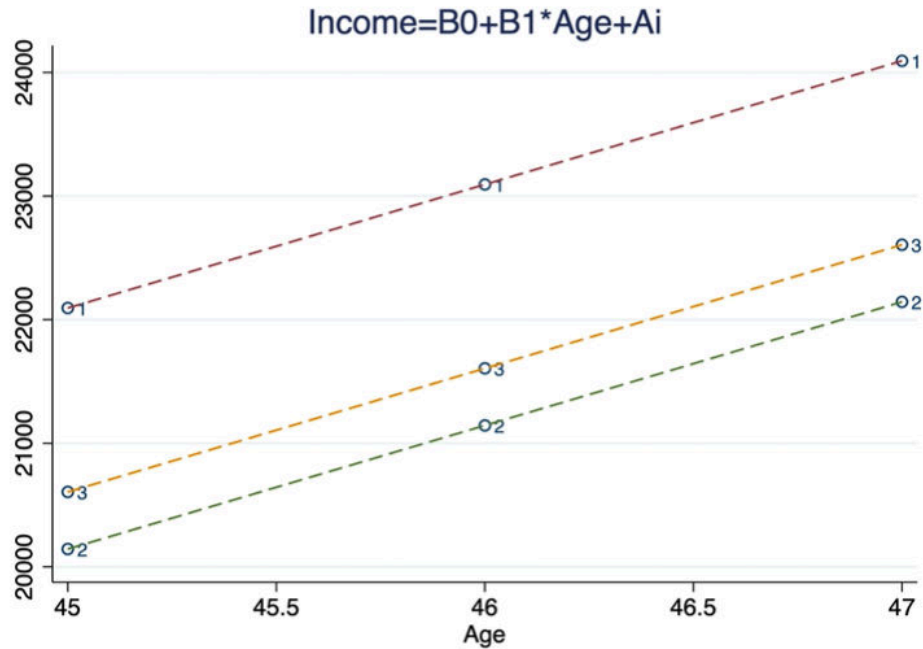
Figure n. 5<sup>27</sup> (Carlos Riumalló Herl, 2022)



In this case we want to see if there is any correlation between income and age. For this example have been collected data for three people. If we look at the example without thinking about the errors, we see that income is a linear function of age. At age 45 individuals have an income of around 23000, at age 46 of 24000 and at 47 of 25000.

If we add the individual heterogeneity, so if we add a time-invariant component of the error term that make each of these individuals different, we now have a scenario where the intercept for each individual will be different and will be given by each individual heterogeneity. In this case, we have that age and income have the same linear relationship □ the distinction now is that the individual heterogeneity will make the intercept for each individual different □ this is what is called a “between difference” (how different people are between each other).

Figure n. 6<sup>28</sup> (Carlos Riumalló Herl, 2022)

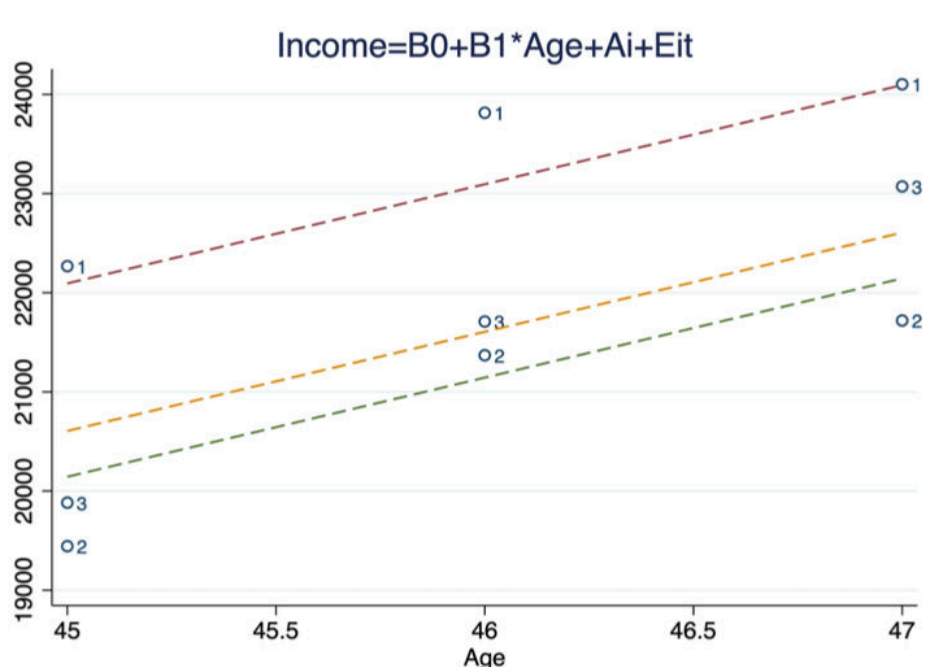


In this setting we have that the top line is for individual 1, the second for individual 3 and the third for individual 2.

Now we add the idiosyncratic shock, an error term that varies for each individual  $i$  at time  $t$  □ it is a different value for each unique observation

In this setting we add  $\varepsilon_{it}$  in the formula:

Figure n. 7<sup>29</sup> (Carlos Riumalló Herl, 2022)



We can now see how the actual values get different or get spread out from the lines. Now our dots are both made up of the individual heterogeneity as well as the idiosyncratic shock.

□ Whenever we are talking about the individual heterogeneity, we are referring to the time-invariant part of the error term that distinguishes units between each other and the idiosyncratic shock, that distinguishes and is unique to each point  $i$  and  $t$  and will help distinguish the within variation.

We can go further with the notation already seen:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit}^T + \sum_{c=C}^{C+N} \beta_c X_{cit} + \alpha_i + \varepsilon_{it}$$

When writing the model we can distinguish time-invariant and time-variant variables. In the case of the notation above, we have added two terms: The sum from  $c = 1$  to  $C$  of the  $X^T$  (which is the time-variant variable) and the sum from  $c = C$  to  $C + N$  of the  $X$  (the time-invariant variable).

□ This will be important as some of the estimation techniques we will analyze will exploit some of these variations.

In many cases we will just write down the variables we will be using in the model without making it clear or distinguishing them as time-invariant or time-variant. But in the case of this video this equation is meant to show us that we can distinguish in the model between the two types and decompose the error term. We also have to remember that some characteristics can be time-invariant in our data, depending on the population and sample that we are using.

## Lecture 22 (3.4) – Unbiasedness and efficiency

The decision to use one estimation technique over another is often based on whether an estimation technique is unbiased and efficient.

For this reason we are going to analyze the unbiasedness and efficiency in the context of panel data.

An **estimator** is a method or a technique that researchers use to estimate certain parameters with a given data.

When we have to choose estimators, we have to consider two properties of interest:

- We have to see if they are unbiased  $\square$  they are if they will give us a parameter equal to the true population parameter on expectation.
- We have to see if they are efficient  $\square$  the estimators are efficient if they give us precise estimates.

Unbiasedness often relates to the value of the parameter estimated, while efficiency relates to the standard errors  $\square$  it influences how confident or not we are on rejecting the null hypothesis that the parameter is different than zero.

In the case of panel data, unbiasedness builds upon what we have already seen for cross-sectional data and refers to the lack of correlation between the error term and the variable of interest. But as we saw in the previous lecture, in the case of panel data we can decompose that error term into two parts (into the individual heterogeneity and into the idiosyncratic shock). Therefore, in the case of panel data, an estimator will be unbiased if both components of the error term are uncorrelated with the variable of interest:

- $\alpha_i$  must not be correlated with  $x_{it}$   $\square E[X, \alpha] = 0$

- $\varepsilon_{it}$  must not be correlated with  $x_{it}$   $\square E[X, \varepsilon] = 0$

In some estimation techniques in the field of panel data some assumptions can be relaxed, but in the end a panel data estimate will be unbiased if and only if the individual heterogeneity component and the idiosyncratic shock are uncorrelated with a variable of interest ( $X$  or  $Y$ ).

Failure to do so, unless we can relax those assumption in specific estimation techniques will lead to a biased estimate.

In the case of efficiency, we have to consider the structure of the error term and how the error terms are correlated over time and between each other (and over people, with other covariates).

One example that harms efficiency in OLS from cross-sectional data is heteroskedasticity  $\square$  if the error term is correlated somehow with our variable of interest, we will have heteroskedastic error terms and we will have to adjust for it, correct it, and finally have precision.

This is not different from the issues we could find in panel data. In the case of panel data, an estimation technique will be inefficient if there is a structure in the error terms that we are unable to account for correctly.

If we have inefficient estimators, this will produce incorrect standard errors and we will need to solve this problem if we do not want to reach incorrect conclusions.

In the case of panel data, when thinking about the structure of the error terms, we also need to account for the fact that we have two dimensions  $\square$  not only we have to consider whether an error term is uncorrelated with the error term of another individual, but we also have to understand and account if that error term is correlated within units.

The question that should come to mind when thinking about data structure, is whether the error terms are correlated or not across time. We will see some estimation techniques that assume that that is the case, but often is unlikely that the error terms in panel data are uncorrelated with each other (so that  $Corr(\mu_t, \mu_s) = 0$ ).

The reason for this is simple: if we go to the construction of that error term, if we look at how the error term is constructed at time  $t$  and at time  $t - 1$ , as shown in the following equation, we can see that both of the error terms at time  $t$ , and at time  $t - 1$ , are constructed and based on the individual heterogeneity:

$$Corr(\mu_{it}, \mu_{it-1}) \neq 0 \{ \mu_{it} = \alpha_i + \varepsilon_{it} \quad \mu_{it-1} = \alpha_i + \varepsilon_{it-1}$$

□ both error terms have a component within them that is the same over time. That makes it such that then the error terms are likely going to be correlated over time. This is a part of the data structure, of the error structure in panel data, that we need to account for future models that we will see. This will determine whether a model might be efficient or not.

In general, in panel data analysis, efficiency will depend on how the error terms will be structure. But still, it is likely that the error terms will be correlated with each other across time: we need to take care of this.

## Lecture 23 (3.5) – Pooled OLS

We've already seen that OLS is one of the best estimation techniques. Under which circumstances can we use (pooled) OLS, and what are the challenges when using it?

Pooled OLS is the estimation of parameters using OLS on data that combines multiple observations of units in a sample. It is like using OLS but using pool data (multiple observation for the same unit).

In this scenario, we estimate the following equation:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit}^T + \sum_{c=C}^{C+N} \beta_c X_{cit} + u_{it}$$

The outcome  $Y$  for an individual  $i$  in a time  $t$  is a linear function of our parameters, the time-invariant variables and time variant-variables and an error term.

The difference between this situation and the one with cross-sectional data, is that now our data points have two dimensions □ we have data for an individual  $i$  at a time  $t$ .

In general, pooled OLS is the same estimation done for cross-sectional data □ this estimation will be unbiased under the same assumptions and will have the same properties as a normal OLS estimation for cross-sectional data.

A difference between pooled OLS and OLS estimation is that a pooled OLS exploits all data variations including the between and within.

Also the normal OLS exploits all data variations, but with cross-sectional data there is only a between variation.

The pooled OLS is the Best Linear Unbiased Estimator (BLUE) if the following two assumptions hold:

1. The zero conditional mean assumption holds (the error term must be uncorrelated with the variable of interest).
2. There has not to be serial correlation ( $Corr(X) = 0 \forall t \neq s$ ).

## zero conditional mean assumption

If the zero conditional mean assumption holds it is possible to estimate a parameter equal on expectation to the true population parameter.

This can happen only if both components of the error term are not correlated with a variable.

□ A pooled OLS will give us an unbiased estimate with panel data if we can assume that the individual heterogeneity and the idiosyncratic shock are uncorrelated with the variable of interest.

If they are not correlated, if we use an OLS, we will get an unbiased parameter for the variable of our interest.

## serial correlation

Having no serial correlation (or autocorrelation) is the second condition required to make the pooled OLS the best estimation technique.

If error terms are uncorrelated with each other, then whatever we estimate as a standard error with an OLS technique will be the appropriate standard error that we can get.

This means that conditional on  $X$ , the error term in two different time periods should be uncorrelated with itself.

If these the two conditions hold, as well as other assumptions such as having full ranked data, not having multiple correlation between variables, etc., then, by using OLS, we will get the best estimates possible. Pooled OLS under these circumstances will always be the best estimation technique, even in the field of panel data.

However, in the case of panel data, we have two issues to discuss to know if pooled OLS is the best estimating choice.

## endogeneity



Are we in a circumstance where the zero conditional mean assumption holds? If we are able to assume that the error term is uncorrelated with the variable of interest, pooled OLS might be the best option and we should test to see if it is also the most efficient option

But as we saw for cross-sectional data, being able to assume the zero conditional mean assumption is very unlikely in observational studies: the treatment allocation (why people have a certain variable), is often correlated with time-variant and time-invariant unobserved characteristics that also influence the outcome.

The term "endogeneity" means that there is no zero conditional mean assumption □ if this is the case, pooled OLS estimates will be biased. So even if there is no serial correlation, whether there would be endogeneity, pooled OLS are not an estimator to use, as the estimates will be biased.

If there is endogeneity, the other estimation techniques we choose have to account for the sources of endogeneity.

As possible solutions, we will see fixed effects and first differences as well as difference in differences. Also instrumental variables can be considered a solution, since they can also be coupled with panel data.

## serial correlation (continuation)

The second issue we have to face in choosing the pooled OLS as our estimation technique with panel data, is the possible presence of serial correlation.

In the case of panel data, it is very unlikely that the error terms are uncorrelated with each other over time, so that  $Corr(X) = 0$

The reason for this is that each error term for an individual  $i$  at time  $t$ , is based on a common individual heterogeneity that is similar in all observations of an individual:

$$Corr(\mu_{it}, \mu_{it-1}) \neq 0 \{ \mu_{it} = \alpha_i + \varepsilon_{it} \quad \mu_{it-1} = \alpha_i + \varepsilon_{it-1}$$

The consequence is that it is likely that error terms within an individual will be serially correlated. As this leads to a different structure in the error terms, if we use an OLS, the standard error that we will get will be invalid □ whatever standard error we get it will be incorrect. This has an impact on the significance and the confidence intervals obtained. It is possible that estimates may still be unbiased, if the zero conditional mean assumption holds, but the standard errors will be incorrect.

So, when the estimates are unbiased (as the zero conditional mean assumption holds) but we have serial correlation, we must opt for solutions that address the structure of the error term; the one we will see first is random effects (which is preferred when working with panel data as it addresses the structure of the error terms).

Another solution is the clustered standard errors (which we will not analyze).

## Lecture 24 (3.6) – Fixed effects, Least Squared Dummy Variables and First Differences

As we saw in the previous lecture, OLS might be biased as an estimation technique (it often is the case when working with panel data). For this reason, we can also count on other estimation techniques, that take advantage of the panel nature of the data to try to address part of the correlation that might be leading to biases.

We are going to study three new estimation techniques:

- Fixed effects
- Least Squared Dummy Variables
- First Differences

These estimation techniques address the correlation that might exist between the unobserved heterogeneity and the variable of interest.

We will see the intuition behind these estimations, their implementation and the assumptions required for them to be unbiased. At the end of the lectures, we will try to understand how to choose between these different techniques.

### intuition for FE, LSDV, and FD

The main challenge we can have, as we previously saw, in pooled OLS, is that it might be biased if the zero conditional mean assumption does not hold, which occurs when the error term is correlated with our variable of interest.

In the case of panel data this can occur if either the unobserved heterogeneity is correlated with the variable of interest, or if the idiosyncratic shock is.

In practice we cannot know for sure which is the case. In the case of panel data, we are able to exploit the multiple observations that we have for each unit, to account

for the correlation that might exist between the unobserved heterogeneity and our variable of interest.

All three of the methods we analyze will eliminate all of the between variations that might be the source of this correlation and will only focus on within variations.

That is the reason why these methods are called within effects estimators (as they only exploit the within variations).

## fixed effects method: implementation

The main idea is that we are going to apply a transformation to the data, called **time demeaning**, to eliminate the unobserved heterogeneity  $\alpha_i$ .

In essence, this method takes into account the correlation between the unobserved heterogeneity and the variable of interest by eliminating the unobserved heterogeneity.

We know the original model:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit} + \alpha_i + \varepsilon_{it}$$

If we now calculate the average over time for each component of the equation (the between effect estimator), we will have the following equation:

$$\bar{Y}_{it} = \beta_0 + \sum_{c=1}^C \beta_c \bar{X}_{cit} + \alpha_i + \bar{\varepsilon}_{it}$$

The average of the unobserved heterogeneity will be the unobserved heterogeneity itself. Similarly, the average of a time-invariant characteristic will be that time invariant characteristic itself.

This is important, as by doing the transformation we are now able to eliminate the unobserved heterogeneity  $\alpha_i$ :

$$Y_{it} - \bar{Y}_{it} = \sum_{c=1}^C \beta_c (X_{cit} - \bar{X}_{cit}) + (u_{it} - \bar{u}_i) \leftrightarrow \ddot{Y}_{it} = \sum_{c=1}^C \beta_c \ddot{X}_{cit} + (\varepsilon_{it} - \bar{\varepsilon}_i)$$

This is done by subtracting the second equation from the first.

This subtraction of the averages is what is called a time demeaning transformation.

In this case the new formula will be such as that our outcome is now a transformation, where we subtracted each mean for each observation in a unit. The same will happen for each variable  $X$ , the unobserved heterogeneity, and the idiosyncratic shock. As a result of this transformation, we will have a new model,

given by the outcome of unit  $i$  at time  $t$  as a linear function of the transformed  $X$  variable for unit  $i$  at time  $t$ , plus the idiosyncratic shock.

When we did the subtraction (the time demeaning), we eliminated the unobserved heterogeneity. We are now only left with only time variant characteristics. Each of these characteristics' transformed values, represent a deviation from the unit mean:  $\ddot{Y}_{it}$  is the deviation with regards to unit mean.

By doing this transformation we eliminated the unobserved heterogeneity  $\square$  we now do not care if there was a correlation between the unobserved heterogeneity and the variable of interest as we have been able to eliminate that unobserved heterogeneity.

It is important to note that we do not have eliminated the idiosyncratic shock, which is still in the model.

## least square dummy variables method: implementation

In contrast to the fixed effects method, we now account for the correlation between the unobserved heterogeneity and the variable of interest by actually estimating a parameter for each unobserved heterogeneity.

Instead of eliminating that value, we now estimate a value for each unobserved heterogeneity. In essence, we have the same original model as before:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit} + \alpha_i + \varepsilon_{it}$$

Now in STATA we are going to include a dummy variable for each unit in our sample, equal to 1 if that observation belongs to the unit, and 0 if not.

We now include all of the dummies in our variable and implement an OLS strategy. By this approach, we do not eliminate the unobserved heterogeneity, but estimate what that parameter is for each unit dummy. It does not matter whether the unobserved heterogeneity was correlated with the variable of interest, as including that heterogeneity as a variable in our model, we are now able to account for that correlation.

This is often referred to as a fixed effects dummy  $\square$  we are adding a unit fixed effects to the model:

$$Y_{it} = \delta_i + \sum_{c=1}^C \beta_c X_{cit} + \varepsilon_{it}$$

$\delta_i$  □ it is the dummy variable for each unit in the sample.

## first differences method: implementation

This approach has a similar perspective as the fixed effects approach: we account for the correlation between the unobserved heterogeneity and the variable of interest by differentiating away the unobserved heterogeneity.

This exploits the fact that we have multiple observations per unit and that the unobserved heterogeneity does not change across time.

We have an example with two time periods (notice that we can also apply this procedure with infinite time periods):

$$Y_{i1} = \beta_0 + \sum_{c=1}^C \beta_c X_{ci1} + \alpha_i + \varepsilon_{i1}$$

$$Y_{i2} = (\beta_0 + \delta_0) + \sum_{c=1}^C \beta_c X_{ci2} + \alpha_i + \varepsilon_{i2}$$

$\delta_0$  □ it indicates how much the intercept changes over time (this will be relevant later in the lectures).

If we now subtract the time period 1 from the time period 2, we are able to see that we are eliminating all of the time invariant variables as well as the unobserved heterogeneity. Because the unobserved heterogeneity does not change over time, if we subtract one period from the next, we will eliminate it.

We are left with a first difference model where we now estimate what is the relationship between a change in the outcome  $Y$  as a linear function of a change in the variable  $X$  and a change in the idiosyncratic shock:

$$(Y_{i2} - Y_{i1}) = \delta_0 + \sum_{c=1}^C \beta_c (X_{ci2} - X_{ci1}) + (\alpha_i - \alpha_i) + (\varepsilon_{i2} - \varepsilon_{i1})$$

$$\Delta Y_i = \delta_0 + \sum_{c=1}^C \beta_c \Delta X_{ci} + \Delta \varepsilon_i$$

By doing these first differences we have been able to eliminate the unobserved heterogeneity. Whatever correlation might have existed between the unobserved heterogeneity and the variable of interest, it is now gone, and irrelevant to us.

With this method we added  $\delta_0$  (as described above). This concept is not included in the two previous methods.

In STATA we will see that this method gives us the option to include or not  $\delta_0$ .

## FE, LSDV, and FD: assumptions

All of these models use within variations to estimate the parameter of interest. With the various transformation we got rid of everything regarding the between variations. Any between variations across units will be eliminated  $\square$  this is the source of the unobserved heterogeneity.

These methods require a **strict exogeneity assumption**  $\square$  this means that FE, LSDV and FD will be unbiased if and only if  $\varepsilon_{it}$  is uncorrelated with our variable of interest  $X_{cit}$  in any time period.

We do not need to know whether the unobserved heterogeneity  $\alpha_i$  is correlated with the variable of interest because we can now account for this correlation  $\square$  this addresses all observed and unobserved time-invariant sources of bias.

Even if we are not able to collect data on important time-invariant sources of bias, all these methods will account for it.

## FE, LSDV, and FD: challenges

Because FE, LSDV, and FD only produce estimates based on within variation, it would be impossible for us to estimate the effect of time-invariant characteristics  $\square$  it is impossible to estimate parameters for time-invariant characteristics with no within variations. So, if, for example, the aim of our research is exploring the effect of some time-invariant characteristics such as gender, education in older people, etc., we will not be able to use these estimation techniques: STATA will drop these variables from the model.

In addition to this, because these models only use within variations, they will also be less efficient, as they get rid of a lot of sources of variations and focus only on within variations. This is the reason why, if the zero conditional mean assumption holds (therefore making us able to use the pooled OLS), using a FE, LSDV, and FD technique will be less efficient than a pooled OLS.

Also now with FE, LSDV, and FD, measurement errors problems are usually much more important. The reason for this is that, once you get rid of all the between variation, a greater part of the signal will be the error term  $\epsilon$  whatever it is left in our variation that we can actually estimate, the measurement term, might become a greater share than what we had in the model that used both within and between variations.

The biggest challenge, which is also the assumption, of these techniques, is that FE, LSDV, and FD will be unbiased if and only if the idiosyncratic shock  $\epsilon_{it}$  is uncorrelated with the variable of interest  $X_{cit}$  at any time point.

Whenever we are using these methods, there is a series of questions that we have to ask us to see if we can assume strict exogeneity or not:

- Where is the within variation of our variable of interest coming from?
- Is it a random variation or is it endogenous?
- Is it driven by other time-variant unobserved variables that could lead to bias?
- Can we assume that the idiosyncratic shock is uncorrelated with our variable of interest?

If we can assume that the last question is true, these are probably the methods that we should be using. If not, we have to implement other techniques to address this sort of biases.

One of these techniques is using an instrumental variable with panel data. Another one is using difference in differences (will be discussed later). Then we have some experimental or quasi-experimental methods to account for the source of endogeneity.

How can we choose between these estimation techniques?

All these methods provide the same result and the same parameter. However, the final choice should be made relying on the efficiency of each one in a given situation.

## comparing FE and LSDV

Both methods provide similar results, and both will be unbiased under the same assumptions. The only difference between these two, is that LSDV will provide us an estimate for the unobserved heterogeneity; at the same time, this will require a major computational effort, as for each dummy we have to estimate a parameter: if we have a sample with 1000 individuals, we will have to estimate a parameter for

each of those units. The more and more units we have, the more STATA will take to estimate those parameters. In essence, it is often preferred to choose fixed effect as it computationally more efficient, unless the aim of our research is to estimate the unobserved heterogeneity (in this case we have to choose for LSDV; in any other case FE will be a better choice).

## comparing FE and FD

We have a similar situation as the previous one. Both of them provide similar results, both of them will be unbiased under the same assumptions. The choice we will make will rely on the structure of the idiosyncratic error. If we have a data set with only two time periods, then fixed effects and first differences will be identical: they will give the exact same estimates and the same standard errors. If the data set has more than two time periods, we have to evaluate the structure of the idiosyncratic error to see which of these two techniques will be more efficient. If the idiosyncratic error is serially uncorrelated, the FE will be more efficient. If it is not, then FD will be the best choice for dealing with that structure.

In practice, people will test both approaches and see whether the results are too sensitive to the choice of the method. If the results are very sensitive to the choice of the method, we will have to go into more details and evaluate what type of structure the idiosyncratic shock has (we will not analyze this matter in this course).

## Lecture 25 (3.7) – Random Effects and Correlated Random Effects

If pooled OLS is biased, we can consider going into the within effects estimation techniques.

If pooled OLS is unbiased, is there serial correlation in the error terms, or no?

If there is, we should opt for a solution like random effects.

## intuition for random effects

When using panel data, pooled OLS may be inefficient if zero correlation occurs between error terms. If there is a correlation between error terms, at any given time



point in our data □ using pooled OLS is inefficient □ they will give us incorrect standard errors and we will have a potential false rejection of the null hypothesis

To address this problem, we can use the structure of panel data to account for the serial correlation between the error terms that may exist. This exploits the fact that in panel data we have multiple time periods for the same unit.

We can think about Random Effects as being somewhere between an OLS and a fixed effects, or first difference estimation techniques.

Random effects will, as pooled OLS, use both between and within variation.

The advantage with regards to pooled OLS will be that random effects accounts for the serial correlation that might exist in the error term.

The idea on which random effects is based on, is that we use a quasi-demeaning transformation to eliminate the serial correlation.

We saw the time demeaning transformation, which occurs when we subtract the average. Now we are doing a similar transformation, called **quasi-demeaning**.

Looking at the error terms structure in panel data, we can also analyze the serial correlation as follows:

$$\text{Corr}(u_{it}, u_{is}) = \frac{\sigma_{\alpha}^2}{\sigma_{\alpha}^2 + \sigma_{\varepsilon}^2}, t \neq s$$

The serial correlation is the share of the total variation given by the individual heterogeneity.

This formula represents how correlated are the error terms in panel data.

By using this, we can define a GLS transformation that can eliminate the correlation in the error terms:

$$\theta = 1 - \left[ \frac{\sigma_{\varepsilon}^2}{T\sigma_{\alpha}^2 + \sigma_{\varepsilon}^2} \right]^{\frac{1}{2}} = 1 - \left[ \frac{1}{\left( T \left( \frac{\sigma_{\alpha}^2}{\sigma_{\varepsilon}^2} \right) + 1 \right)} \right]^{\frac{1}{2}}$$

With  $0 < \theta < 1$

It is important to understand the implication of this formula, not the formula per se. We refer at the GLS transformation as the  $\theta$ . It is the value we use for the quasi-demeaning transformation □ once the theta is calculated, we can transform

the model into a quasi-demeaning data that will subtract from each observation in our sample the quasi-demeaned value:

$$Y_{it} - \theta \bar{Y}_i = \beta_0(1 - \theta) + \sum_{c=1}^C \beta_c (X_{cit} - \theta \bar{X}_{ci}) + (\alpha_i - \theta \alpha_i) + (\mu_{it} - \theta \bar{\mu}_i)$$

In this equation, the original model is transformed by subtracting the average of each component, so the average of the outcome for the unit, the average of the variables for the unit, the average of the unobserved heterogeneity and the average of the idiosyncratic shock, and we multiply these averages for the quasi-demeaning factor.

If theta assumes a value equal to 1, we will have the demeaning function we saw previously for within effects estimation techniques.

For any type of model we do not know the real value of theta, but with the data it is possible to have a consistent estimate of it.

Once we calculate theta, we can have two scenarios, that will let us know whether the random effects is closer to a fixed effects or to an OLS.

In the first scenario, if we have a greater variation between units rather than within units, it means that we need to account to a greater extent for what is driven by that unit heterogeneity. So, if units are more different between each other than within each other, we need to account for a lot of that unobserved heterogeneity.

This process should remind us of fixed effects, where we try to eliminate as much as between variation as we can.

In this case we will see that theta tend into the value of 1:

$$\sigma_{\alpha}^2 \gg \sigma_{\varepsilon}^2 \Rightarrow T \left( \frac{\sigma_{\alpha}^2}{\sigma_{\varepsilon}^2} \right) + 1 \rightarrow \infty \Rightarrow \theta = 1 - \left[ \frac{1}{T \left( \frac{\sigma_{\alpha}^2}{\sigma_{\varepsilon}^2} \right) + 1} \right]^{\frac{1}{2}} \rightarrow 1$$

□ The more variation we have between units the more that the random effects theta will go closer to one.

In that circumstance, the random effects estimator will be closer to a fixed effects estimation technique than to a pooled OLS. The reason for this is that because unobserved unit heterogeneity is relatively important, we will try to account for it as much as possible and eliminate as much as between variation as we can.

A similar finding can happen if the time periods go to infinity  $\square$  the more and more time periods that we have, the more the within variation will become important, and the more we will need to account for it.

$$Y_{it} - \theta \bar{Y}_i = \beta_0(1 - \theta) + \sum_{c=1}^C \beta_c (X_{cit} - \theta \bar{X}_{ci}) + (\alpha_i - \theta \alpha_i) + (\mu_{it} - \theta \bar{\mu}_i)$$

In an extreme case scenario theta could assume value 1: it is the exact same time demeaning transformation. Anyway, in practice theta will never assume value 1.

In the second scenario the variation between units is less important. In this case, the idiosyncratic shock is far more important, as we do not care as much about the unobserved heterogeneity, as it is a very small share of the total variation.

In this case the theta tends towards 0:

$$\sigma_\alpha^2 \ll \sigma_\varepsilon^2 \Rightarrow T \left( \frac{\sigma_\alpha^2}{\sigma_\varepsilon^2} \right) + 1 \rightarrow \infty \Rightarrow \theta = 1 - \left[ \frac{1}{T \left( \frac{\sigma_\alpha^2}{\sigma_\varepsilon^2} \right) + 1} \right]^{\frac{1}{2}} \rightarrow 0$$

We are therefore closer to a pooled OLS rather than to a fixed effects model.

In this case, the idiosyncratic shock is much more important, so we can leave as much as the between variation as we want, as it is less important.

Similarly with the previous scenario, we could think of an extreme case, where the theta assumes the value of 0. In this case we left with just a regular pooled OLS:

$$Y_{it} - \theta \bar{Y}_i = \beta_0(1 - \theta) + \sum_{c=1}^C \beta_c (X_{cit} - \theta \bar{X}_{ci}) + (\alpha_i - \theta \alpha_i) + (\mu_{it} - \theta \bar{\mu}_i)$$

But, as before, in practice theta will never be equal to 0.

## RE: assumptions

One important characteristic of a random effects model is that, even if it is a panel data estimation technique, this model will use both the between and within variation to estimate the parameters of interest. It will distribute the share of how important each of these are and, in the case where theta tends towards 1, we will eliminate more and more of the between variation, but, in any case, the random effects model will always use both sources of variation.

One advantage of this approach is that therefore we can estimate the parameter for both time-variant and time-invariant characteristics □ we are not in the same situation as we were in fixed effects models where we can only estimate parameters for time variant characteristics.

Random effects will be unbiased if the error term is uncorrelated with the variable of interest at any time point; it is the same assumption as for pooled OLS □ if pooled OLS is biased, random effects will be as well, and vice versa.

Both pooled OLS and random effects do not address the potential issues of the potential endogeneity that might exist.

The only thing that random effects does in contrast to pooled OLS is that it will correct for the serial correlation that exists in the error term, it will not address the issue of bias.

## correlated random effects (CRE)/mundlak

One extension of the random effects model that has been widely developed in recent years, is the correlated random effects (also referred to as the Mundlak estimate).

The correlated random effects is an extension from random effects, where rather than assuming that the unobserved heterogeneity is uncorrelated with the variable of interest, we will model that relationship.

If we want to use a random effects estimation technique, we need to assume that the whole error term is uncorrelated with a variable of interest. This implies, in panel data, that also the unobserved heterogeneity is uncorrelated with the variable of interest. Correlated random effects works on panel data by, instead of assuming that unobserved heterogeneity is uncorrelated with a variable of interest, we are going to model part of that relationship.

We can think that the correlated random effects is somewhere between a random effects model and a fixed effects model.

We will model that relationship and account for part of the correlation that exists between the unobserved heterogeneity and the variable of interest.

The main idea behind the correlated random effects is that we can model the unobserved heterogeneity as a function of the averages of time-variant variables and include these in the model. Since the average of time-variant variables might represent partly that unobserved heterogeneity, by including these modelled

unobserved heterogeneity we will be able to account for the modelled correlation between the unobserved heterogeneity and the variable of interest.

## CRE: implementation

For each time-variant variable, we are calculating the average over the units. So, in this case, for each time variance characteristic, we are calculating the average and we are going to define the unobserved heterogeneity as a linear function of these variables and a residual unobserved heterogeneity, as represented by the following equations:

$$\bar{X}_{ci} = T^{-1} \sum_{t=1}^T X_{cit}, \forall C$$

$$\alpha_i = \alpha + \sum_{c=1}^C \gamma_c \bar{X}_{ci} + r_i$$

Once we have been able to model this and calculate the averages for each time-variant characteristic, we are now going to be able to include it in the random effects model and have the following equation:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit} + \alpha + \sum_{c=1}^C \gamma_c \bar{X}_{ci} + r_i + \varepsilon_{it}$$

The outcome for unit  $i$  at time  $t$  is a linear function of our original parameters and the variables we are estimating, plus the modelled part of the unobserved heterogeneity that is constructed based on the averages of the time-variant characteristics. At the end, instead of having the full heterogeneity as we had previously, we have a residual unobserved heterogeneity which is  $r_i$  and, as always, the idiosyncratic shock.

The hope behind this method is that by including this modelled unobserved heterogeneity, we can account for a great share of the correlation that would have existed in the original unobserved heterogeneity and the variable of interest.

## CRE: assumptions

Correlated random effects is in this similar to random effects as we are using both the between and within variations to estimate the parameter of interest. But because we have been able to model part of the structure of the unobserved heterogeneity, we are going to account for a part of the correlation that exists between the unobserved heterogeneity and the variable of interest. In this setting, the correlated random effects will be unbiased if we are able to assume that the residual part of the unobserved heterogeneity and the idiosyncratic shock are uncorrelated with the variable of interest at any point.

When doing this modeling we hope that what is left of that unmodelled unobserved heterogeneity is a much smaller part and account for less of the correlation that would have existed between the unobserved heterogeneity and the variable of interest.

Therefore  $\square$  correlated random effects is an extension of random effects that starts to consider for part of that correlation that could exist, but there is still a concern that whatever is not modelled could still lead to biases.

If we can assume that whatever was left of the error term is uncorrelated with the variable of interest, then we know that the estimates we would get from a fixed effects model and a correlated random effects model will be the same for time-variant variables, so in that circumstance we have that the correlated random effects will be as good as a fixed effects model. However, the second advantage and why we might in this case prefer the correlated random effects, is that using correlated random effects will allow us to include time-invariant characteristics. This is an advantage from fixed effects, with which we would not have been able to do this.

As long as we can assume that the way we have modelled the unobserved heterogeneity leads to the fact that the rest is uncorrelated with the variable of interest, the estimates we get for a fixed effects model and the estimates for correlated random effects model would be the same for time-variant variables ( $\hat{\beta}_{C-FE} = \hat{\beta}_{C-CRE}$ ).  $\square$  using a correlated random effects might actually be better than a fixed effects model because it will actually be more efficient by using the between and within variation.

But in reality, it is likely very difficult to assume that whatever is left in the unmodelled part of the unobserved heterogeneity is still uncorrelated with the variable of interest.

A second advantage of the correlated random effects, that we will analyze more in detail in the next lecture, is that correlated random effects is one of the techniques that we can use to test whether we should use a random effects model or a fixed effects model.

In the end, a correlated random effects can give us some information about how likely is it that we can assume that the unobserved heterogeneity is uncorrelated with the variable of interest (we will see this more in detail in the next lecture).

## Lecture 26 (3.8) – Estimator choice

How can we choose the most appropriate estimation technique?

We have already seen how under some circumstances we could choose an estimation technique over another.

We are going to analyze two methods to choose between a random effects and a within effects estimation technique, which are both correlated random effects and the Hausman Test.

These methods are useful to test how likely is that the unobserved heterogeneity is uncorrelated with the variable of interest. In both methods, if we have evidence that it is unlikely to hold  $\square$  we will opt for the within effects estimation techniques (FE, LSDV and FD). In the case that we are likely to assume heterogeneity  $\square$  we will opt for a random effects model as it is more efficient. None of these methods tell us anything about the possible correlation that might exist between the idiosyncratic shock and the variable of interest.

These methods only tell us whether we can assume if the unobserved heterogeneity is uncorrelated with the variable of interest.

Looking at the models we have already seen, we know that FE, LSDV, FD and RE require, to be unbiased, the idiosyncratic shock to be uncorrelated with the variable of interest.

If this is not the case, all of these methods would be biased, and we need to rely on other techniques.

The difference between the within estimation techniques and random effects is that random effects additionally requires that the unobserved heterogeneity has to be uncorrelated with the variable of interest.

Therefore, the choice on the estimation technique hinges upon whether we can assume this assumption.

The two tests can verify how likely it is that this assumption holds. In practice, we can never verify those assumptions because we can never observe the unobserved heterogeneity, but both these tests give us some ideas of how likely this assumption holds.

## estimator choice: CRE

As we previously saw, the correlated random effects is an approach by which we model part of the unobserved heterogeneity as a linear function of the average of time-variant characteristics:

$$Y_{it} = \beta_0 + \sum_{c=1}^C \beta_c X_{cit} + \alpha + \sum_{c=1}^C \gamma_c \bar{X}_{ci} + r_i + \varepsilon_{it}$$

We then have a model where we have our outcome for units  $i$  at time  $t$  as a linear function of our parameters  $\beta$ , the variables of interest in the original model, and a linear function of the parameter's gammas and the averages of time-variant characteristics and the residual part of the unobserved heterogeneity and the idiosyncratic shock.

The idea behind the correlated random effects model as a choice mechanism, is that we are going to verify how likely it is that those parameter gammas are all equal and equal to 0:

$$H_0: \gamma_1 = \gamma_2 = \dots = \gamma_c = 0$$

If we can reject the hypothesis that all the gammas are equal and equal to zero  $\square$  at least one component of the time variance averages influences the outcome. After this, if we know that there is one component of the time-variant averages that influences the outcome and therefore omitting it would lead to bias, is therefore likely that there are many other unobserved components that also matter and that we cannot capture.

If we can find something that matters, it is likely that there are other things that matter that we cannot observe.

So, there will be some parts of the component of our  $i$  that it is likely to be correlated with  $x$ , and in that case, it reinforces the fact that we need to use the fixed effects model.



If we reject the hypothesis and have some gammas are significantly different than zero  $\square$  it is likely that other time-variant or other time-invariant characteristic influences the outcome  $\square$  we need to account for it.

The only way to account for all time-invariant characteristics is by using the within effects estimator that will get rid of all between variation.

If we fail to reject the hypothesis, so if we find no evidence that the parameter gammas are significantly different than zero  $\square$  at least all of the time-invariant component we have added, that are the averages of time-variant variables, do not matter. We can then extrapolate that it is likely that other time invariant components do not matter as well. This is based on the essence that usually the variables you observe are often the most important for the model, and therefore the average of these variables should be more important than anything we do not observe  $\square$  under this circumstance we are able to assume that the whole unobserved heterogeneity is uncorrelated with the variable of interest and therefore, as long as we can assume that the idiosyncratic shock is also uncorrelated, we will know that random effects will be unbiased, and will be a much more efficient choice than FE, LSDV and FD.

## estimator choice: Hausman test

We are going to run both a fixed effects model and a random effects model, comparing how the coefficients of time-variant characteristics look like in both models.

If we find that the time-variant coefficients are very different in the fixed effects model than in the random effects, we are able to say that the time-invariant characteristics that we cannot control for in the random effects model but can control for in the fixed effects model matter. If they matter, it is likely that random effects is biased. This is a way of testing between both by negating the other. If the coefficients are very different  $\square$  whatever we cannot observe in the random effects model that is time-invariant matters  $\square$  fixed effects is a more appropriate estimation technique that will account for all time-invariant characteristics.

But if we run a fixed effects model and a random effects model and we compare the coefficients of the time-variant variables and find no significant difference  $\square$  the time-invariant characteristics, that are not controlled for in the random effects model, do not necessarily matter or influence the other coefficients. We are able to assume in this case that the random effects is likely to be unbiased  $\square$  random effects

will be a more appropriate and efficient technique rather than a fixed effects estimation technique.

In the STATA application we can see how to develop these two tests with the data and how to compare and obtain the coefficients.

In this course we must be aware that, unlike in other courses when we can decide between models by verifying the goodness of fit, in the case of panel data, the choice of an estimator is based purely on bias and efficiency. For any type of output we will get, we will always get measures of goodness of fit, and they will be distinguished in:

- Between  $R^2$
- Within  $R^2$
- Overall  $R^2$

These measures should never be used to choose estimators. Each estimation technique, pooled OLS, RE, or the within effects variations, will optimize the different sources of variation □ in each case we will get a different  $R^2$  value for the three types of variations.

- The goodness of fit measures are never to be used when choosing between estimators and panel data.
- The choice should only be based purely on the concepts of bias and efficiency, nothing else and nothing more.

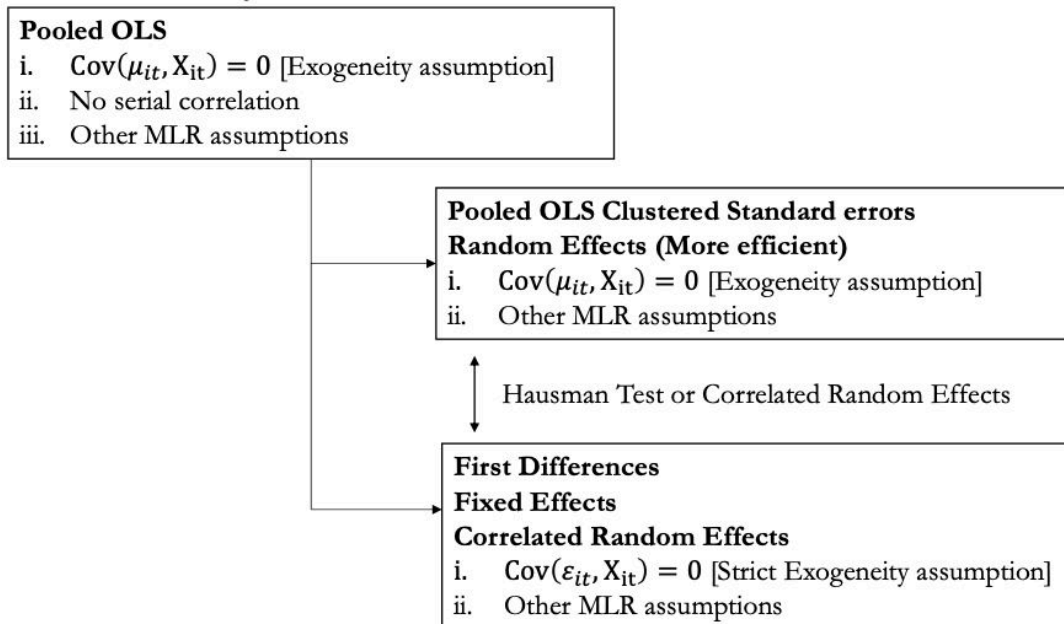
So, when choosing which is the best and most appropriate estimation technique: When we work with panel data, the first estimation technique that comes to mind is a pooled OLS, the best option if we can assume exogeneity (the complete error term is uncorrelated with the variable of interest) and there is not serial correlation. If these two conditions hold □ pooled OLS are the best linear unbiased estimate, and also most efficient estimator.

If we cannot assume full exogeneity, the most appropriate choice is using the within effects estimation techniques. In this case, FE, FD and CRE will be more appropriate than a pooled OLS, as they account for the correlation that might exists between the unobserved heterogeneity and our variable of interest. In this sense, the choice between pooled OLS, FD, or within effects estimation techniques is not based on efficiency, but on the unbiasedness: FD, FE, LSDV and CRE will be less biased than pooled OLS.

However, if we are able to assume exogeneity, but we have serial correlation, the most efficient choice is to use random effects model. Random effects models and pooled OLS are unbiased under similar situations, but random effects will use the quasi-demeaning transformation to account for this structure in the error term and it will be more efficient.

Figure n. 8<sup>30</sup> (Carlos Riumalló Herl, 2022)

## Summary



The test we have just seen are useful to see how likely it is for the full exogeneity assumption to hold. By using the Hausman Test or correlated random effects, we are able to choose between a random effects model and within effects estimation techniques.

Again, in practice it really difficult to use pooled OLS as it is extremely unlikely to have no serial correlation when working with panel data.

## Lecture 27 (3.9) – Attrition

The data structure in panel data can be balanced or unbalanced. When data is unbalanced (for some reason we are not able to follow people up for the whole

amount of time), it is important to understand why we are missing some information, as this could affect the validity of our analysis.

Attrition □ the loss of follow up of people in a panel data survey.

## attrition and item non-response

Panel data are not always complete, and in this case, they are unbalanced. There could be two reasons for this:

1. **Attrition** □ unit dropped from the sample (the individuals that were present at the start of the research project might not stay until the end, or maybe they might start later than the beginning of the research). There could be many reasons for dropping from the sample:
  - o We could not find the individuals, as they could have moved.
  - o Death (relevant for aging surveys).
  - o People might refuse to remain in the sample, for whatever reason.

Attrition not only is an issue for samples that follow people up; if we are following up a hospital, if it would go over bankruptcy, it would not be possible for us to follow it up.

2. **Non-response** □ individuals might not want to respond to certain questions. This could have occurred more prominently in cross-sectional data, when we had missing answers. But. We could also have this issue with panel data.

Attrition and item non-responding are challenging in panel data as we need to understand the source and the reason of that attrition or of the item non response.

Is the incompleteness missing at random or is it a selected sample?

In the case of attrition and in the case where we are following people up over time, are the people that are refusing to answer the surveys of a particular nature or are they just a sample of the population?

If healthier persons, for example, start dying more frequently in a longitudinal aging, then the sample that we are being left off with is a healthier sample, and the following analysis could be biased because of that. Another issues that can happen in panel data that affects both item non-response and attrition, is that individuals might modify their behavior and reporting across waves.

Because we are collecting data over time and we are serving the same individuals or firms over time, there might be some learning effects or gaming as to which question answer. One example of this from the aging survey is that there are cognitive

questions that are asked of individuals to measure their cognitive capacity (for example individuals are asked to count back from a hundred to one) on seven units intervals. The first time people are faced with the challenge they find it very difficult, but the second time they follow the survey, they might have an easier time because of their experience □ this is the learning effect.

This is important to keep in mind for panel data analysis, as for certain questions there might be a learning effect that biases the measures of variable towards upward or downward.

Many of the longitudinal surveys are quite long, some time they take some hours. If people understand how to answer the survey in order to make it faster, they might actually game the responses, they might give you the responses that make it faster to answer or quicker to end.

Another side of learning effect which refers more to panel conditioning, is that individuals might change their behavior because they belong to a panel. For example, if we are doing a health longitudinal survey, we might be giving information through the questions that makes people change their behavior. For example, if we an individual a lot about healthy food, healthy lifestyle or sports, they might actually modify their behavior after the survey because we have been talking about that a lot.

All of this also combines into affecting the representative of the sample. Whenever we are working with panel data and an issue of attrition or item non-response is that the sample we are working with might change over time. We have to reflect about the representativeness of the sample.

In some cases, the sample can become very different to the population because we might lose some particular people.

This is a concern if and only if at the beginning we started with a representative sample. That might not be the case, and in that scenario, this might not be an issue.

## testing for attrition

Whenever we are working with panel data analysis we need to understand if the attrition that is occurring is random or not. We have three methods to test if attrition is occurring in the sample and if it might not be random. These methods are based on the different construction of an indicator variable. They are the following:

1. We construct one variable that will assume value 1 if that unit is available in all waves, and 0 if otherwise.

2. The second method will be based on whether the unit is available in the next week or not.
3. This method is a count of the number of waves that an individual is present.

□ Not all of these methods can be used with fixed effects.

## Method 1: All waves

With the first method (the all waves method) we are going to create, for each units in our sample a binary indicator whether that unit was present in all the waves or not. In this case, an indicator variable for unit  $i$  at time  $t$  will be equal to 1 if that unit  $i$  was in all waves or zero otherwise.

$$Ind_{it} = \{1 \text{ if unit } i \text{ is available all waves } 0 \text{ otherwise}$$

As we are creating a value that is the same whether the unit was in all waves or not, the value for the indicator variable for unit  $i$  at time 1 will be the same as the indicator variable for unit  $i$  at time 2 and so forth until time  $T$ :  $Ind_{i1} = Ind_{i2} = \dots = Ind_{iT}$

□ This indicator variable has no within variation □ this is a method that cannot be used with the within estimation techniques

Once we created the binary indicator, we have to run the model that we intend of running including the estimation technique that we think we will run and we are going to include the indicator variable as an additional covariate.

In the end we will have to interpret the correlation of that indicator variable with the outcome. In many cases attrition will happen anyways in our study, but, if we are able to find out that the indicator variable is significantly correlated with our outcome, this implies that attrition is not at random □ there are certain groups of the population that are falling out, or dropping from our sample, and that might lead to bias.

In any type of discussion concerning attrition, it is not enough to just say that missing is not at random or that attrition is not at random. We need to discuss and explain why the population that is dropping out could influence the outcomes and the results that we find afterwards. If we are looking at an example where we are evaluating the effect of health on income and we find that our indicator variable is

correlated with income (lower income people fall out of the survey), it is logical to think that if lower people are falling out of the survey and they are also those that are unhealthier, then the actual relationship we find with our data is going to be on downward biased estimates. This is because the people left in the sample are wealthier and healthier.

## Method 2: Unit is available next wave

The second method relies on another idea. It follows the same first step:

$$Ind_{it} = \{1 \text{ if unit } i \text{ is available in wave } t + 1 \text{ 0 otherwise}$$

This indicator variable will therefore change over time  $\square$  the indicator variable for unit  $i$  at time 1 will not necessarily be equal to the indicator of unit  $i$  at time 2. If, for example, a person drops out at wave three, the indicator at time 1 will be equal to 1, but the indicator at time 2 will be equal to 0.

Having a potential within variation is the only method of the three that can be used with fixed effects or within effects estimation techniques.

Steps 2 and 3 are the same in the method 1.

## Method 3: Number of waves

We create a new indicator variable which is going to be a number of continuous waves where a person was present:

$$Ind_{it} = \sum_{t=1}^T I(W_{it} = 1)$$

In this case, the indicator variable for the unit  $i$  at time  $t$  will be equal to the number of waves that we can observe that person. In the case of a person that was in our survey for three time periods, this indicator variable would take value 3, if it is a person that was in the survey for five time periods, would take value 5. That value will remain the same within a unit  $\square$  the indicator variable for unit  $i$  at time 1 is the same as the indicator variable for unit  $i$  at time 2 and so on.

Therefore, there is no between variation  $\square$  it is a method that cannot be used with within effects estimators.

Steps 2 and 3 are the same as previously.

# Lecture 28 (3.10) – Introduction to Difference-in-Differences

One of the challenges behind a FE, LSDV and FD models is that they are only unbiased if and only if the idiosyncratic error is not correlated with the variable of interest at any time point.

Where is the within variation coming from? Is that change in that variable random or is it endogenous? Are there time varying unobserved variables that influence our treatment?

All these questions are useful to assume or not if the idiosyncratic shock is uncorrelated.

In the DiD we will use a quasi-experimental situation to compare two groups, a treatment one and a control one.

By comparing the groups over time, we will get an unbiased estimate if we can assume that all unobserved time-varying changes are common to both groups. So, if for the exception of our actual treatment, all other changes that could occur over time are common to both groups, it is irrelevant whether we can observe them or not, because a DiD estimation technique will address them and give us unbiased parameters anyways for the treatment.

The DiD complements the methods we saw in panel data as it accounts for unobserved time-varying changes that might lead to bias.

However, a DiD is not exclusive to panel data, and we can do it with repeated cross sections as well. The main idea behind DiD estimator is that we can use these two comparison groups to differentiate out any sort of bias that is constant and common for those both groups.

Even if the sources of bias change over time, as long as they are the same for the treatment and control group, we can differentiate out and get the causal effect of a treatment.

DiD is often referred to as a quasi-natural experimental design, as it often exploits the rolling out or the implementation of a particular policy.

We have a population of interest that is then split into two:

- A treatment group that will receive a treatment later on.
- A control group that will not receive the treatment.



The definition of treatment groups is based on policy design. For example, we could have income tax credits that are only given to single mothers, poverty benefits that are only given to people below a certain income, etc.

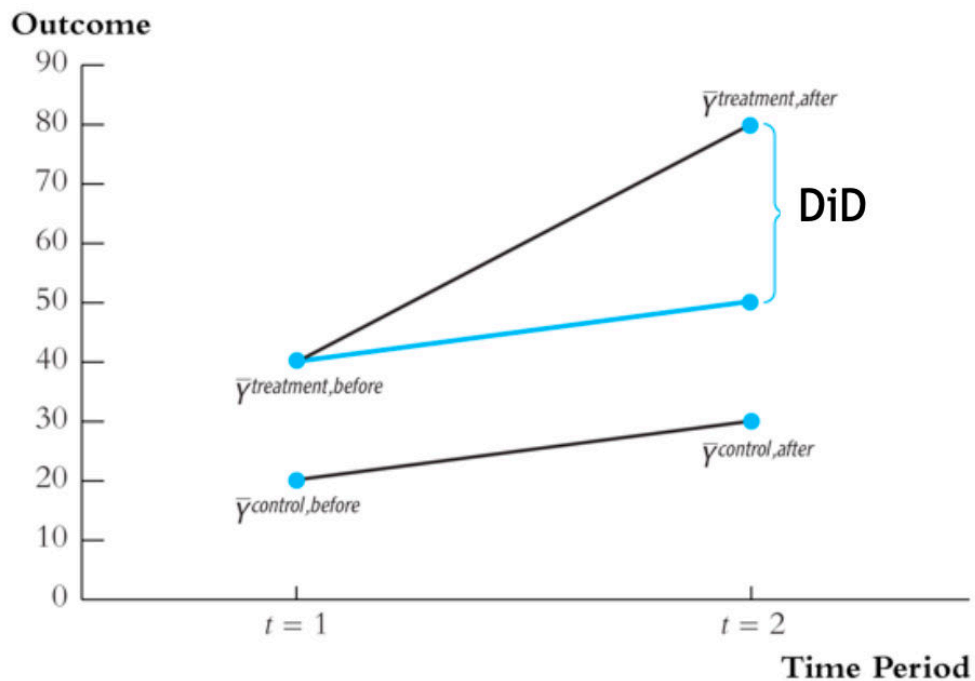
After the distinction of the two groups, we have a pre-intervention period where we are going to have the data of our outcome. In this phase, neither the treatment group nor the control group are being given the treatment.

Then there is the intervention phase, and it is when the benefits are starting to be given out.

Then there is the post intervention phase, where we are going to collect data on our outcomes or other variables. Now the treatment group will actually receive the treatment and the control group will not.

So we are going to follow the treatment and the control group, before the actual treatment start to take place and after it has taken place.

Figure n. 9<sup>31</sup> (Carlos Riumalló Herl, 2022)



We can visualize this. We have two time periods and two intervention or treatments that starts to take place between the two time periods. On the y axis we have our outcome of interest.

What is the effect of the treatment on the outcome?

We have the treatment and the control group. At the bottom we can see the trends to the control group before and after the intervention. At time period 1 we have the outcome of the control group before the intervention, and at time period 2 the outcome of the control group after the intervention.

In the case of the treatment group (we can visualize it with the line above), with the black line we can visualize what happens to the treatment group. We have the outcome of the treatment before the intervention actually takes place and we have the outcome of the treatment after the intervention has taken place.

If we were only seeing this line without comparing with the control group, one could conclude, incorrectly, that the increase in the outcome was only due to the treatment. We know that if there is some reason why people are being treated, the line could be biased, so people in the treatment group could have seen their outcome increased in any case.

DiD incorporates the control group that gives us what the trend of the treatment group would have been in absence of the treatment. The trend of the line of the control group is the same that the line of the treatment group would have if they did not take the treatment. Using this information, we can calculate what the outcome should have been for the treatment group in the second time period, in a circumstance where the intervention would have not taken place. Then any difference in the actual value we observe that is with the intervention and that counterfactual value that is given by the information in the control group, reflects the DiD estimate.

If we can assume that the control group is a valid counterfactual group, then this simple comparison of the outcomes before and after the intervention and between the treatment and control group will give us the causal effect of that intervention on our outcome.

This is shown by the following equation:

$$DiD\ Estimator = (Y_{T2} - Y_{T1}) - (Y_{C2} - Y_{C1})$$

In this equation we compare the outcomes for the treatment group before and after the intervention against the difference in outcomes of the control group before and after the intervention. This is the reason why this model is called Difference-in-Differences: we estimate the difference between the differences we observe in the two treatment groups.

## DiD: assumptions

The assumption for the DiD to work, so that the control group is a valid counterfactual trend for the treatment group, is that the correlation between the idiosyncratic error and any time-varying variable in our treatment of interest is constant. If the correlation is constant, by comparing before and after and between groups we can eliminate the source of bias.

This assumption implies:

- If there is constant bias, the trend of the control and the treatment group should be similar in absence of the treatment. This is called the **parallel trends assumption/constant bias assumption**. It can be seen in late terms as saying that there are no other factors that will change differentially over time for the treatment group relative to the control group that would impact the outcome. This is an assumption, and we cannot verify whether the trend of the treatment group in absence of the treatment would be the same as the control group because we cannot observe that. There are ways to test how likely this assumption holds. In these cases, if we have data more than one time period of data for the pre-intervention phase, we can observe whether the trends before the intervention were similar before the treatment or the control. If we can say that over time, before the intervention, the trends were similar, it is likely that the trends would have remained similar going forward. Therefore, any deviation from the trend can only be attributed to the intervention that has now taken place.
- The units (the treatment and the control group) are stable. This means that there is not a relevant interaction between the treatment and control groups □ there is no **spillover** between the treatment and the control groups. If receiving the treatment influences the behavior and responses of the control group, then there will be a violation to the constant bias, as the control group will change in a differential way to the treatment group.

This condition cannot be tested in our implications, and we only need to discuss why treatment and control groups are relatively independent.

## DiD: advantages (and disadvantage)

The DiD estimation has a set of advantages:

- If coupled with panel data, we can account for all time-invariant unobserved factors that differ between the treatment and the control group. As we saw

with the within effects estimators, using a DiD with individual effects will also account for all time-invariant factors, whether they are observed or unobserved.

- The advantage of the DiD with regards to the within effects estimators, is that DiD can also account for all time-varying unobserved factors that are the same for the treatment and the control group, even if do not observe them. In this situation, if we have changes that we can observed but that happened simultaneously for both the treatment and control group, even if they are correlated with the outcome, it will not influence our estimates because we will be able to eliminate that with the DiD technique.

The only drawback for the DiD estimation technique is that cannot account for time-varying unobserved factors that are different between the treatment and the control group. If there are changes that are different for treatment and control groups and that are unobserved and influence the outcome, then it is impossible for us to distinguish how much of the change in the outcome is due to our actual intervention or to these other changes.

This implies arguing why there are no other concurrent events with our intervention that would be different for the treatment and the control group.

# applied microeconomics - Module 4 – Introduction to empirical methods: Binary data

## Lecture 29 (4.1) – Binary data: Introduction

For this module we will focus on the use of public transportation<sup>32</sup> (Garcia-Gomez, 2022).

### **Examples of binary dependent variables**

**Binary dependent variables** are common in many applications of Behavioural Economics, Health Economics, and so on. We can think about models to know

whether individual participate in the labour market (if they are employed or not, whether they work full time or part time or whether an individual is retired or not). There are also some examples in the field of health and healthcare: if we are interested in estimate the determinants of an existing chronic condition, or whether an individual has any chronic condition, whether they had been visited by doctors or not, whether they have been to the hospital.

In all these examples we can start noting that the dependent variable has always two possible options.

More examples could be:

Does an individual smokes or not? What are the effects of changes in taxes on the probability that someone smokes? And if it is provided some incentives for someone to exercise or not? Drink or not?

Why do some people become entrepreneurs?

What is the effect of their parental background in becoming an entrepreneur?

If someone uses public transportation or not, or if he uses a bike or not.

To answer these questions, the dependant variable could also either be yes or no, so 1 or 0.

If we want to understand if a firm is innovative and whether new policies provide incentives for firms to innovate, we can be interested in the mode of transportation.

If we reflect on people choices and on how they make decisions under uncertainty, providing our subjects with a risky and a safe option, and we want to see whether the way in which we show the information has an effect on the probability that that person choose the risky or the safe option.

In all these examples we can look at the dependent variable (the one we want to explain) and define it as a **dummy variable**. So, it will take value 1 if, for example, the individual participates in the labour market, if another individual has a chronic disease, and so on, and 0 in the other case.

We focus on the use of public transportation, using the Dutch Mobility Survey (the 2010 survey). The sample is composed by 3000 individuals. The dependant variable is ptuse, and it assumes value 1 if the individual uses public transportation, 0 otherwise. As explanatory variables we have a dummy variable (urban, it assumes value 1 if the individual lives in an urban area and 0 otherwise), and a continuous variable (for age; individuals' age can also be divided into four categories: age 0 to 19, age 20 to 39, age 40 to 59 and 60 plus).

Figure no. 1<sup>33</sup> (Garcia-Gomez, 2022)

Variable	Obs	Mean	Std. Dev.	Min	Max
ptuse	3000	.942	.2337824	0	1
age	3000	38.69833	21.98864	0	93
urban	3000	.593	.4913568	0	1
age20_39	3000	.22	.4143154	0	1
age40_59	3000	.3326667	.4712468	0	1
age60_plus	3000	.1923333	.3941992	0	1

In the figure above we can find various information about the dependent and explanatory variables.

The minimum and the maximum are useful to do some cleaning because we expect them to have a minimum of 0 and a maximum of 1.

The mean of the binary variables is the proportion of 1's in the data set. So, for example, if we look at the probability of using public transport, 94% of the individuals in the sample use public transportation.

The average age of the respondents is 38.7 years.

How can we estimate the models with binary outcomes? Later we will see that we can use the linear regression model, the ordinary least squares, and specific non-linear models that accounts for the specificities of this kind of variables.

## Lecture 30 (4.2) – Linear regression model for binary data (LPM)

In a linear regression we assume that the dependent variable  $y$  is a linear function of the explanatory variables  $x_1$  and  $x_2$ :

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

If the zero conditional mean assumption holds, so if  $E(\epsilon | x_1, x_2) = 0$  (which means that the expected value of the error conditional on  $x_1$  and  $x_2$  is equal to 0), then using the expected value of  $y$  conditional on  $x_1$  and  $x_2$ , the linear function can be translated as:

$$E(y | x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

What does this imply when the dependent variable assume value 0 or 1?

The expected value of  $y$  is equal one times the probability that  $y$  is equal to one plus zero times the probability that  $y$  equals to zero:

$$E(y) = 1 \times \Pr(y = 1) + 0 \times \Pr(y = 0)$$

Clearly,  $0 \times \Pr(y = 0) = 0$ , so this means that:

$$E(y) = \Pr(y = 1)$$

The expected value of  $y$  is the probability that  $y = 1$ .

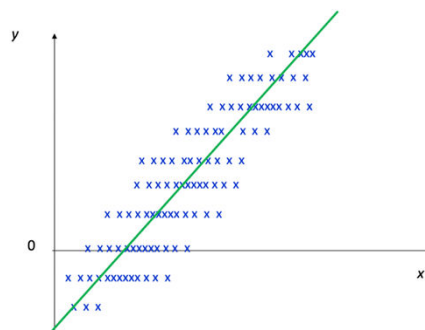
The **linear probability model** means that:

$$\Pr(y = 1 | x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

So, the probability that  $y$  equals 1 is a linear function of  $x_1$  and  $x_2$ .

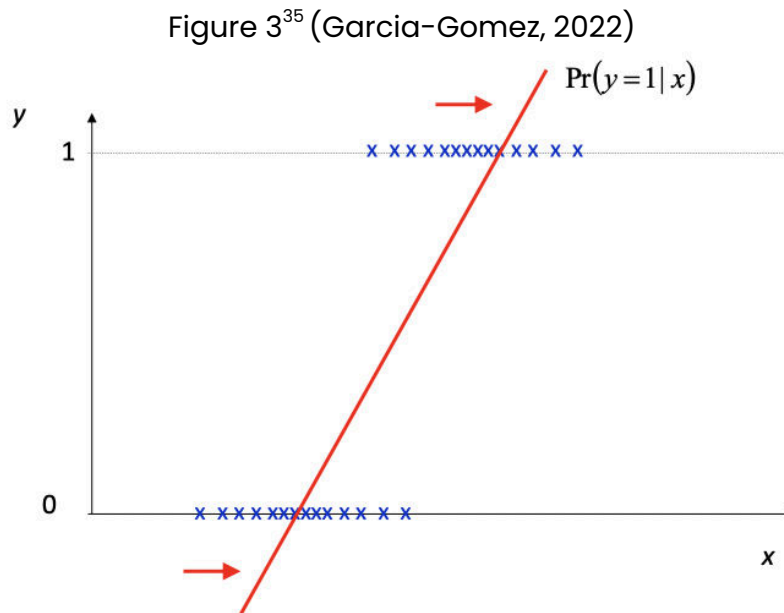
When we have only a continuous dependent variable (and this is not the case) we only have one  $y$  and one  $x$  variables and  $x$  could be any observed value any individual in our population. We could then fit the following line:

Figure 2<sup>34</sup> (Garcia-Gomez, 2022)



And if the dependant variable is not continuous but binary instead?  
 There are not values between 0 and 1 (as people either use public transportation or not).

The OLS regression will fit a line as before:



This already evidences one of the drawbacks of using a linear probability model in the case of binary outcomes: from the fitted line we either get values larger than 1 or smaller than 0, but it is impossible to have either a larger probability than one or a larger probability than zero.

Looking at the example on the use of public transportation:

Figure 4<sup>36</sup> (Garcia-Gomez, 2022)

Variable	Obs	Mean	Std. Dev.	Min	Max
ptuse	3000	.942	.2337824	0	1
age	3000	38.69833	21.98864	0	93
urban	3000	.593	.4913568	0	1

94% of the individuals in the sample use public transport. The average age is 38.7 years and 59% of the respondents live in an urban area.



We estimate a linear regression model, we get the ordinary least square estimators, and we proceed in the same way we have seen in the other modules.

Figure 5<sup>37</sup> (Garcia-Gomez, 2022)

```
regress ptuse urban age, robust
```

```
Linear regression                Number of obs    =    3,000
                                F(2, 2997)       =    35.34
                                Prob > F               =    0.0000
                                R-squared              =    0.0165
                                Root MSE           =    .23192
```

	ptuse	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
→	urban	-.0116806	.008454	-1.38	0.167	-.0282568 .0048955
→	age	.0013389	.0001646	8.13	0.000	.0010161 .0016616
→	_cons	.8971152	.0108157	82.95	0.000	.8759081 .9183222

Now it is important to be cautious about the coefficients' interpretation: we have to think in the units of our variables. The explanatory variable urban, is a dummy, so it will always be interpreted compared to the reference category. Age is a continuous variable, in years. But what about ptuse?

In this case the dependent variable is a probability, and the units of probabilities are **percentage points**. Looking at the estimated probability of using public transportation, as we used the coefficients before:

$$\hat{Pr}(urban, age) = 0.8971152 - 0.0116806 \cdot urban + 0.0013389 \cdot age$$

So, we could get predictions in the same way as the other modules with OLS. In this case we will have the prediction of the use of public transportation.

If we interpret the coefficient of age, we can see it is in percentage points. So, we have to multiply it by 100, and then we can see that as an individual grows older, for every additional year of age, the probability of using the public transport increases by 0.1 percentage points, ceteris paribus (if the zero conditional mean assumption holds, else we could say "keeping urban fixed").

As we did before, we can also see that age is statistically significantly different from 0 at a 1% significance level, looking at either the t-statistic or the p-value.

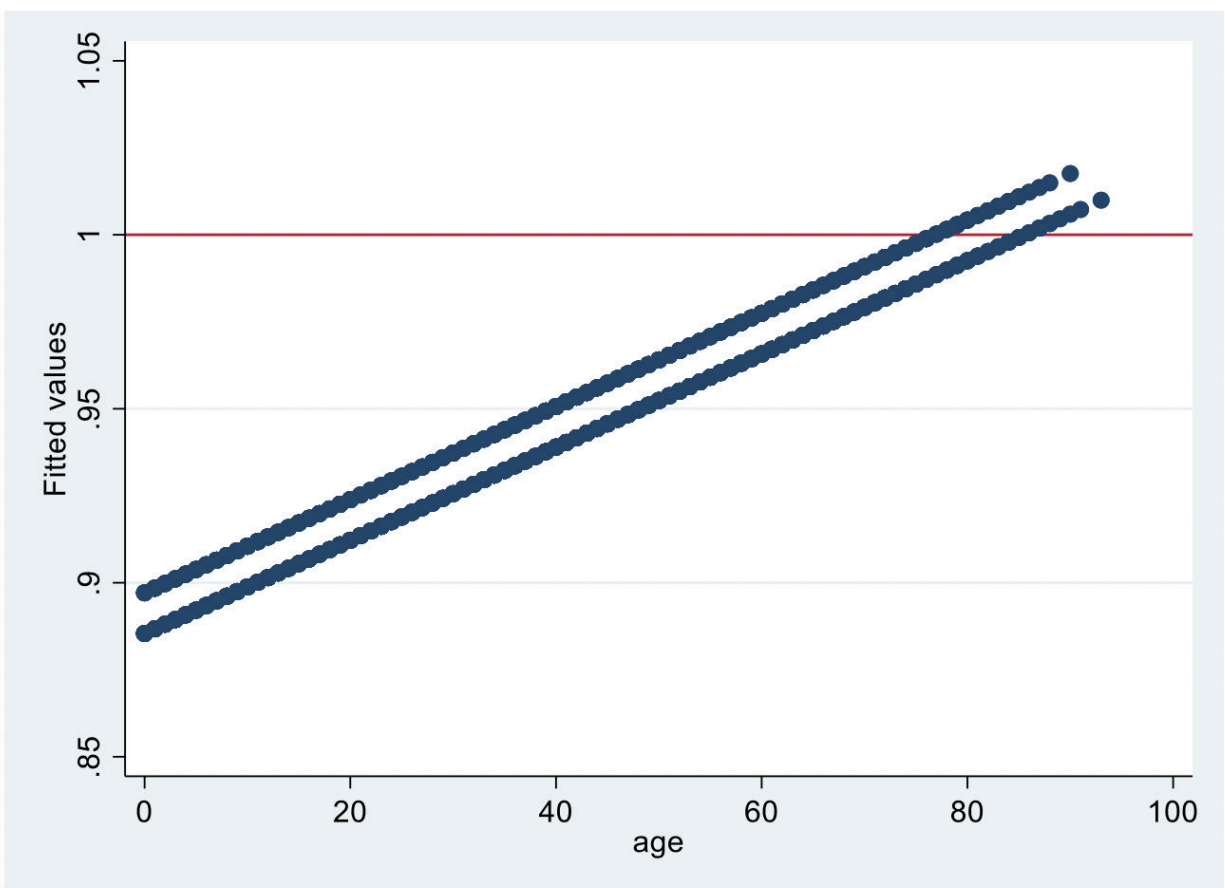
Urban is a dummy variable, so to interpret its coefficient we need to compare it. The probability of using the public transport is 1.16 percentage points lower for an

individual that lives in an urban area compared to an individual non-living in an urban area, *ceteris paribus*, or keeping age fixed.

What is the constant telling us? It tells us the expected value of public transport use when all the  $x$ 's are equal to zero (so when the individual does not live in an urban area and the individual has an age of 0). But is this meaningful? Are there people that do not live in an urban area? Are there individuals with zero years of age? In our examples, we know that there are people living in rural areas and that there are individuals with 0 years of age. So, 89.7% of those that live in a rural area with zero years of age use public transport (which is the same as saying that the probability of using the public transport is 89.7 percentage points).

If we use this model, we get the predicted probability and then we can plot the predicted probability against age. We have then two fitted lines:

Figure 6<sup>38</sup> (Garcia-Gomez, 2022)



This is because we have the dummy variable urban that can assume value 1 or 0. As we have seen before, we get predictions higher than 1: these values would refer to very old individuals, but we are predicting a probability of using public transportation larger than one, which is impossible.

And now we do not see probabilities below 0. Is there something wrong with the model? No. It is not true that we always get prediction larger than one or below zero, we just may get them.

In general, when the sample mean is very high, in this case it is close to one, we are going to get predictions above one, but it is very unlikely that we will get predictions below zero.

On the other hand, when the sample mean is very close to zero, we will get predictions below zero, but hardly we will get predictions above 1.

When the sample mean is around 0.5, it is quite common not to get predictions that are larger than one or below zero.

This is one of the main disadvantages of using the linear probability model. One of the advantages, is that it is really easy to interpret: we know the magnitude of effect just by looking at the coefficients.

## Lecture 31 (4.3) – Nonlinear Models for Binary Data

We are going to analyse models that will not have predictions larger than one or below zero.

We will focus on:

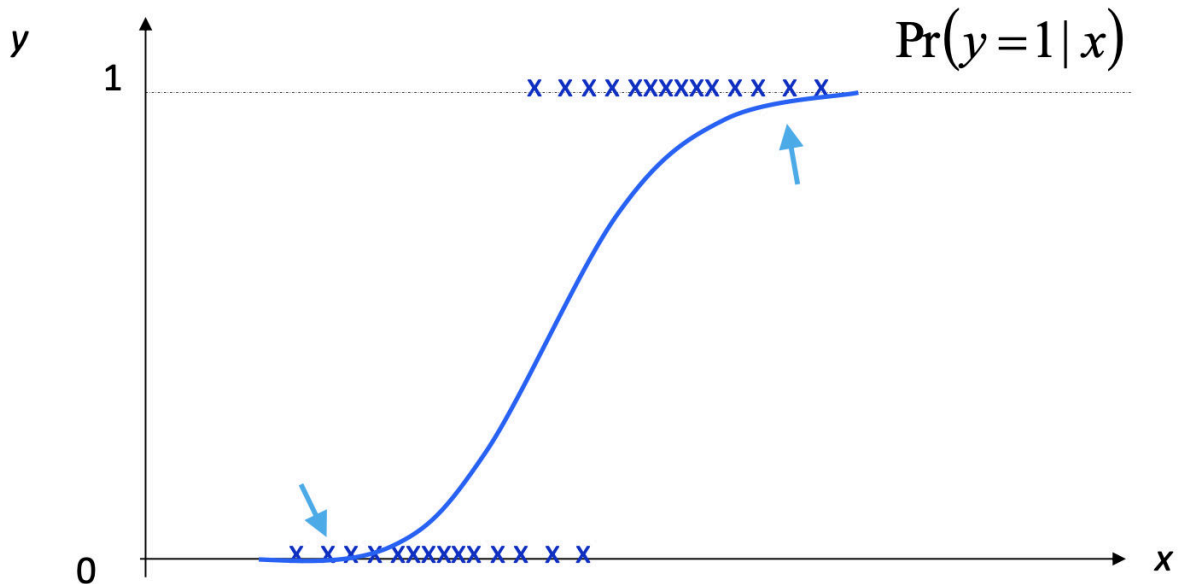
- The logit model.
- The probit model.

The main problem from the last lectures is that when we run a regression with a binary dependent variable, we cannot have predictions outside the interval between 0 and 1.

Our objective is to avoid these larger predictions, and at the same time we want that the effect of  $x$  on the probability that it is 1 or 0, so on its tails, decreases as we get closer to these two limits.

This is what is done by and **S-curve**.

Figure 7<sup>39</sup> (Garcia-Gomez, 2022)



Starting from the observed values, if we plot the distribution of the predicted probability that  $y$  is equal to 1 conditional on  $x$ , we see that this curve is always between the interval  $[0, 1]$ . As it increases, so as we are closer to one, an additional unit of  $x$  has a smaller effect on the probability compared to an increase in one addition unit of  $x$  in the middle of the curve.

This particular S-curve in the graph above is for a case where  $x$  has a positive effect on the outcome. If we would have a negative effect, so that when  $x$  increases, the probability decreases, we would have an inverse S-curve. In this case, with low values of  $x$  we are closer to  $y$  and as  $x$  increases our S-curve tends to zero.

For the S-shape curve we have **non-linear binary models**: we model this probability as a function of  $x_1$  and  $x_2$ , instead of having a linear function:

$$\Pr Pr (x_1, x_2) = F(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

There are various options for the  $F$  functions. There are only two conditions we need to impose:

1.  $F(\cdot)$  is defined in such a way that is impossible to get predictions outside  $[0, 1]$ , so it is  $0 < F(\cdot) < 1$ .
2.  $F(\cdot)$  has always an S-shape.

The most common models that satisfy these conditions are the probit and the logit.

Probit and logit models

The **probit model** is defined as:

$$\Pr Pr(x_1, x_2) = \Phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

With  $\Phi(\cdot)$  which is the **cumulative distribution function** of the standard normal distribution we know from other statistics courses.

The **logit model** is defined as:

$$\Pr Pr(x_1, x_2) = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}$$

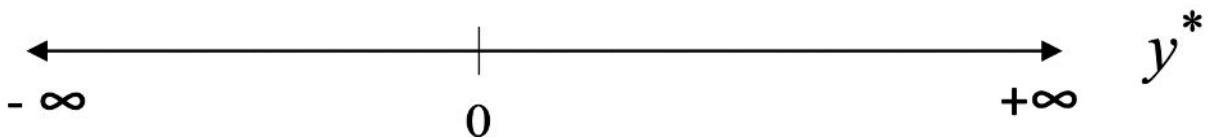
When we talk about these models, we talk about latent variable models.

### Latent and observed variables

We call a **latent variable**  $y^*$ . For example, when we think about the probability of someone using or not the public transport, this could be some sort of latent propensity of someone to use the public transport. Maybe some people, because of their preferences, are more likely to use the public transport, while others less. We can think as this **unobserved latent propensity** as a sort of continuous variable that goes between minus infinite and plus infinite.

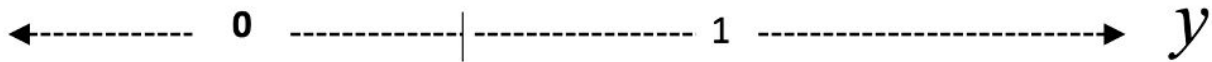
The same applies if we are thinking about someone that has a chronic condition or not, as we could have some sort of underlying health. Then after a given threshold, the person has a chronic condition because we cannot really observe this latent health, as could be a continuous that goes from minus infinite to plus infinite. So, these kinds of variables are latent variables, and we do not observe them.

Figure 8<sup>40</sup> (Garcia-Gomez, 2022)



We observe this other variable, that only takes value zero and one (when the latent has crossed a given threshold that we set at zero for convenience and it take value 0 otherwise).

Figure 9<sup>41</sup> (Garcia-Gomez, 2022)



We have an underlying propensity to use public transport when this underlying propensity is large enough (so, when it is above 0) we use the public transport, while when it is below zero we do not use it.

So,  $y^*$  is a latent variable and  $y$  is what we observe. The relationship between these two variables is  $y=1$  when  $y^*>0$ , and 0 otherwise.

With this relationship we can express this so then we see that this is a continuous variable. So, then we can express the latent as linear function of  $x_1$  and  $x_2$  and the parameters:

$$Pr Pr (x_1, x_2) = Pr Pr (x_1, x_2)$$

The observed variable is equal to 1 when the probability of the latent is larger than 0.

If we assume that the errors follow a normal distribution, then we have a probit. This means that this probability that the latent is larger than zero is equal to the  $\Phi$  function of  $x_1$  and  $x_2$ :

$$Pr Pr (x_1, x_2) = Pr Pr (x_1, x_2) = \Phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$$

This is the reason why we use this function, and there is a similar relationship when we use a logit model.

# Lecture 32 (4.4) – Estimation of Nonlinear Models for Binary data

How can we estimate nonlinear models for binary data? We will focus on how maximum likelihood estimation works (it is the method used to interpret the coefficient in probit and logit models), and how we can (or not) interpret the coefficients in the logit and probit models.

## Maximum likelihood estimation

First of all, we must define the density function of  $y$ . The density is the probability that  $y$  takes certain values (in this case 0 or 1).

The **density function** is the following:

$$f(x_1, x_2) = \left[ F(\beta_0 + \beta_1 x_1 + \beta_2 x_2) \right]^y \times \left[ 1 - F(\beta_0 + \beta_1 x_1 + \beta_2 x_2) \right]^{1-y}, \quad y = 0, 1$$

Once we have defined the density function, we can define the **log-likelihood function** for observation  $i$ :

$$\log \log l_i(\beta_0, \beta_1, \beta_2) = \log \log [f(x_{1i}, x_{2i})]$$

For every individual in our sample, we can define the log-likelihood, which is the log of the density function.

We do this for every observation in the sample; then we can do the log-likelihood function across all observations:

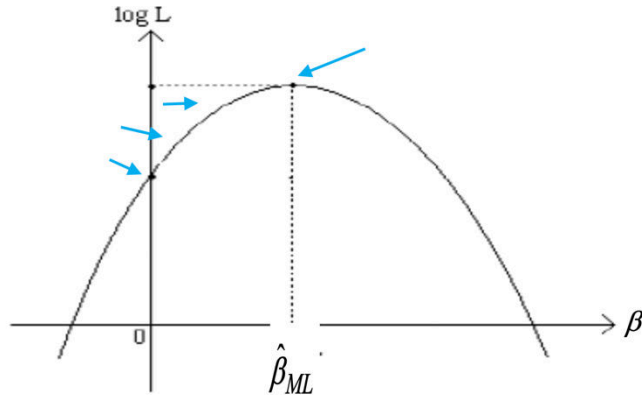
$$\log \log L(\beta_0, \beta_1, \beta_2) = \sum_{i=1}^n \log \log l_i(\beta_0, \beta_1, \beta_2)$$

We obtain the estimated parameters  $\beta_0, \beta_1, \beta_2$  by maximizing the  $\log \log L(\beta_0, \beta_1, \beta_2)$ . So, the maximum likelihood estimator maximizes log-likelihood function to obtain the estimate coefficients  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$ .

## Maximum likelihood estimation graphically

With OLS we have a formula that minimize the sum of the squares of the residuals. Then we apply this formula to get the estimated coefficients. In maximum likelihood estimation we do not have a formula to do this, but we have an iterative procedure.

Figure 10<sup>42</sup> (Garcia-Gomez, 2022)



This could be a representation of the log-likelihood function. We start from the intersection between the curve and the y axis. Then with a numerical procedure we can find the value of  $\beta$  that maximizes the log-likelihood function. We start at 0 and then, in incremental steps that Stata computes (we see Step 0, Step 1, Step 2...), we reach the maximum. Once we are there, at  $\hat{\beta}_{ML}$  we know that this is the estimated coefficient that we could get from maximum likelihood. This is our estimated coefficient for the variable of interest.

We don't have a formula, but we do it in steps.

### Example PT use Probit

For this example, we use a probit model. In Stata we get:

Figure 11<sup>43</sup> (Garcia-Gomez, 2022)

```
probit ptuse urban age
```

```
Probit regression
```

```
Number of obs = 3000
LR chi2(2) = 54.65
Prob > chi2 = 0.0000
Pseudo R2 = 0.0411
```

```
Log likelihood = -636.96328
```

	ptuse	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
→	urban	-.1178419	.0782359	-1.51	0.132	-.2711814 .0354976
→	age	.013089	.0018674	7.01	0.000	.009429 .016749
	_cons	1.201824	.0834919	14.39	0.000	1.038183 1.365465



Stata will always give us the log-likelihood, which will always be negative, because the log of density is always negative since the density is always between 0 and 1 and the log of something between 0 and 1 is negative. So, a higher log-likelihood value is closer to zero, while more negative values are worst likelihood functions. We can compare models through their log-likelihood, but the value per se is not telling us anything.

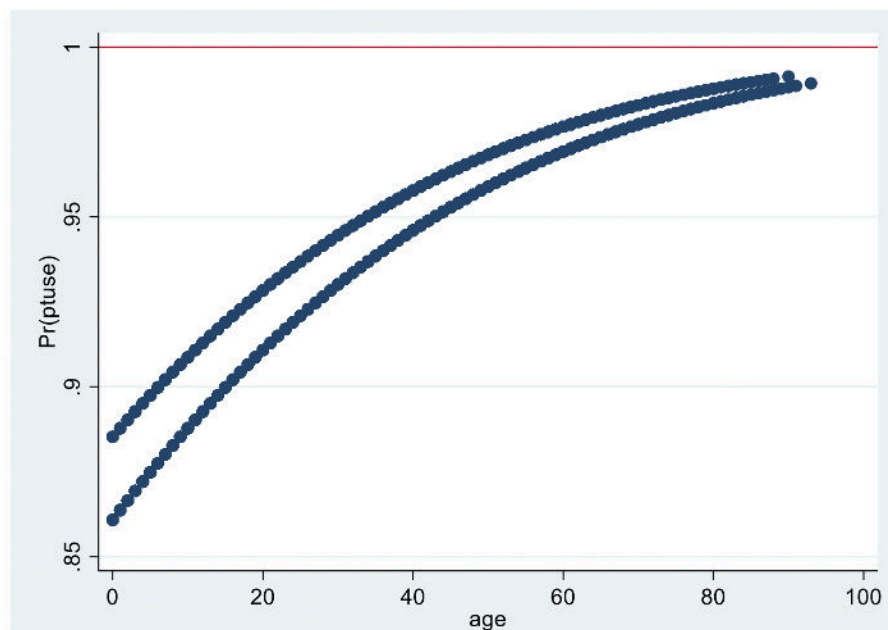
In this example we cannot interpret the coefficients as they are in units of the latent variable. We do not know the units of the latent variable, so we need to operate some transformations to interpret those coefficients in percentage point (the units of probability).

We can still interpret the signs of the coefficients: if a variable has a negative coefficient, then we know that this variable will decrease the probabilities of the outcome (in the example, people who live in an urban area are less likely to use public transport than people who live in a rural area, *ceteris paribus*, while for age the likelihood of using public transport increases with age, *ceteris paribus*).

We can also look at the significance and say whether these associations/effects are statistically significantly different from 0 or not, in order to use these as reference values.

### **Predicted probability**

Figure 12<sup>44</sup> (Garcia-Gomez, 2022)



We do not get predictions above 1 or below 0, and the effect gets smaller as an individual gets older (in this example). And this latter fact is exactly what we wanted (we wanted the effect to be smaller as we are approaching the 1 or the 0). The curves are not parallel, but we will discuss it later.

## Lecture 33 (4.5) – Interpretation with Graphs

How can we interpret the effects of a binary data model using graph?

We are interested in interpreting the size of the estimated effect using graph for both categorical and continuous variables.

We will continue to use the example in which we are interested in the effect of two explanatory variables:

- A dummy variable that assumes value 1 if the individual lives in an urban area and 0 otherwise.
- A continuous variable that is age in years.

We are interested in the effects of these two variables on a binary outcome, the use of public transportation, that takes value 1 if the individual uses public transportation and 0 if not.

We model this probability as a nonlinear function of  $x_1$  and  $x_2$ . The functional form of this function depends on whether we use a probit or a logit model:

$$\Pr Pr(x_1, x_2) = \{\Phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2) \text{ in probit model } \frac{e^{(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}{1 + e^{(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}} \text{ in logit model}\}$$

Once we estimated these models, so once we get the estimated coefficients  $\beta_0$ ,  $\beta_1$ ,  $\beta_2$ , we can obtain the predictions and then we can plot these predicted probabilities against  $x_2$ .

The graphs are useful to understand how we can get the estimated effects. Actually, we could get the exact numbers, obtaining what are referred to as the marginal effects of  $x_1$  and  $x_2$  on the probability that  $y$  is equal to 1, conditional on  $x_1$  and  $x_2$ .

For this example, we will use a logit model:

Figure 13<sup>45</sup> (Garcia-Gomez, 2022)

logit ptuse urban age

Logistic regression	Number of obs	=	3000
	LR chi2(2)	=	50.72
	Prob > chi2	=	0.0000
Log likelihood = -638.92556	Pseudo R2	=	0.0382

ptuse	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
urban	-.2239737	.1641711	-1.36	0.172	-.5457431	.0977957
age	.0257286	.0038261	6.72	0.000	.0182296	.0332275
_cons	2.068369	.1699202	12.17	0.000	1.735332	2.401407

We cannot interpret the magnitude of these coefficients, but only whether they are positively or negatively associated with the probability that someone uses public transport (if the coefficient is negative, the probability decreases, and vice versa).

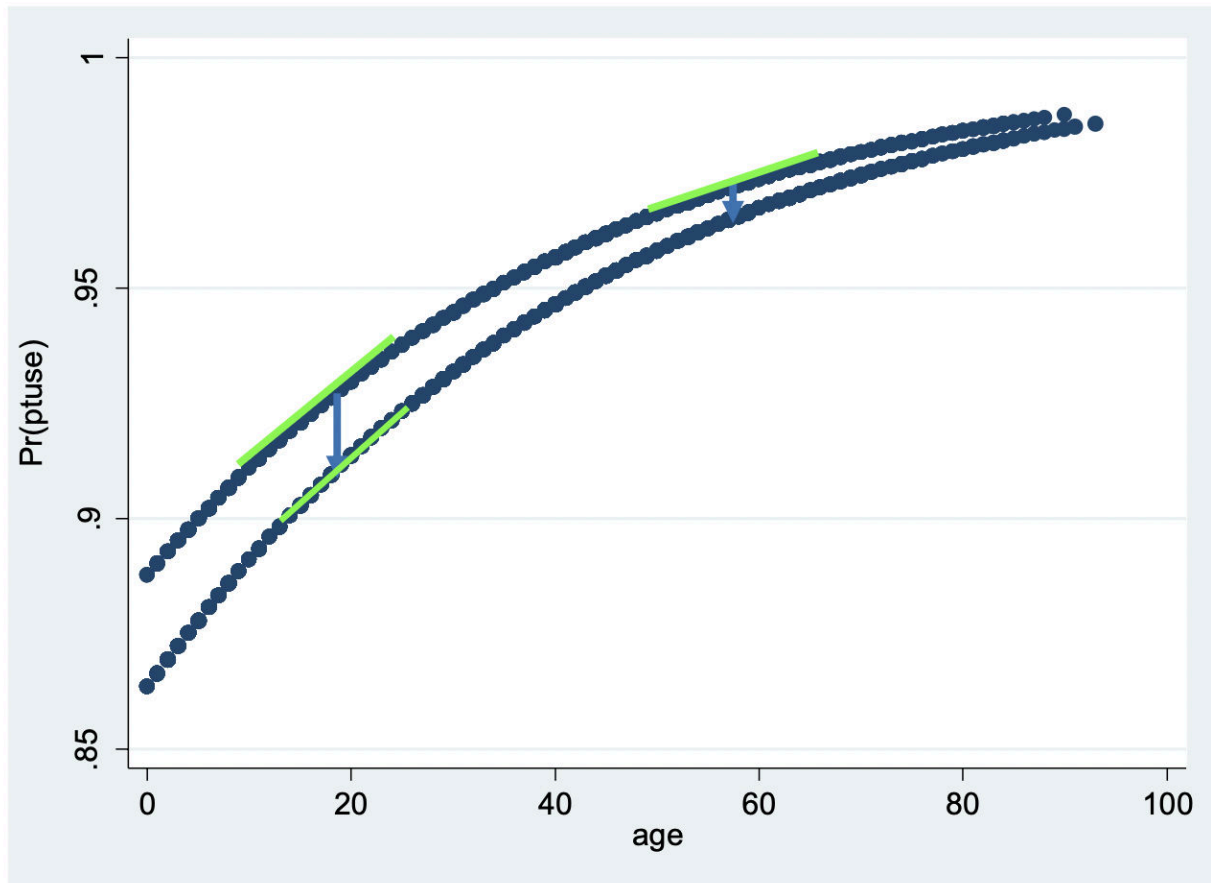
We can get the estimated coefficients and then we know the formula of the logit model. We can then compute the predicted probability for every individual in our sample, computing the following formula for each individual:

$$(age, urban) = \frac{\exp(\exp(2.068369 - 0.2239737 \cdot urban + 0.0257286 \cdot age))}{1 + \exp(\exp(2.068369 - 0.2239737 \cdot urban + 0.0257286 \cdot age))}$$

This way we can get the exact probability for every individual in the sample, conditional on his/her age and on whether they live in an urban or in a rural area by plugging in their values.

In the following graph we can see how the predicted probability looks like:

Figure 14<sup>46</sup> (Garcia-Gomez, 2022)



We have two lines, as people either live in an urban or a rural area. Each one of the two lines tells us how age translates into different probability of public transport use. The effect of age on the probability of public transport use is non-linear. Depending on the age the effect is going to be different, but it is non-linear depending on age, but we also have to pick one of the two curves, as the effect of age is different as an individual lives or not in an urban area.

The coefficient of urban is negative  $\square$  the probability of public transport use is smaller for those that live in an urban area. The line below is the one for the individuals that live in an urban area, while the line above is for the people that live in a rural area. So, the effect of age is different depending on where an individual live and according to that we have to pick a curve; after we picked a curve, the effect is going to be different based on the age of the individual.

The distance between the lines tells us the difference in the effect whether an individual lives in an urban or rural area. Contrarily on what we saw in a linear

regression model, the distance between the lines now is not constant, but it differs based on the age.

What is the effect of urban at age 20? The effect of living in a urban area and living in a rural area is the difference between the line, and we need to check the distance between the curves for when an individual is aged 20 (in figure 14 the blue arrow on the left).

We do the same reasoning for age 60. At that age the effect of urban is much smaller than at age 20, which means that the effect of living in an urban area at age 60 is not so different compared to living in a rural area.

The effect of age is the slope of the curve, so we need to pick one of the two lines and to fix an age. To understand the effect of age in one of the curves, for example the effect of age 20 for somebody who lives in an urban area, it is going to be the derivative at that point, which means that for the person aged 20 living in an urban area it is going to be the slope of the lower green line on the left in Figure 14. If we look at the effect of age 20 for an individual that lives in a rural area, now the slope will be slightly different (it will be the slope of the higher green line on the left).

The effect will be different based on the age distribution and on the values of the other variable.

In conclusion, for dummy variables the effect will always be the difference between the lines. The distance between the line will be different at different values of the other  $x$ 's.

For continuous variables the effect will be different at different points of the continuous variable but will also be influenced by the other variable as well, as we can see from this example.

## Lecture 34 (4.6) – Interpretation with Marginal Effects

We are interested in knowing how to interpret the **marginal effects**, and how to compute them for continuous and categorical variables. Then we will try to understand when we can and when we cannot compare logit and probit models.

In non-linear models we can interpret the sign, the significance but not the magnitude. As the coefficients measure the effect in units of the latent variable, we need a transformation to measure in percentage points (the units of probability).

To do this we need to compute marginal effects. How can we do this? We talk about marginal effects and discrete changes on the probability that  $y$  is equal to 1. This will be different whether we have a continuous or a discrete variable. We have seen that for the discrete (binary) variables, the difference was the distance between the lines, while for the continuous variables it was the derivative (the slope at that point).

If the marginal effect for a continuous variable is the slope at a given point, then it is a derivative of the probability with respect to  $x_2$ :

$$\frac{\partial \Pr(x_1, x_2)}{\partial x_2}$$

The effects (the discrete changes) of a categorical (dummy) variable  $x_1$  are:

$$\Pr(x_1 = 1, x_2) - \Pr(x_1 = 0, x_2)$$

For example, if we look at the effect of living in an urban area, it will be the difference in the probability that one individual lives in an urban area conditional on  $x_2$  minus the probability that the individual lives in a rural area conditional on  $x_2$ .

$x_1 = 1$  happens if the individual lives in an urban area,  $x_1 = 0$  when an individual lives in a rural area.

We could also focus on what we refer to as **risk ratio**, or the relative risk, which is the ratio between these two probabilities.

How can we compute this marginal effect?

$$\frac{\partial \Pr(x_1, x_2)}{\partial x_2} = \beta_2 \phi(\beta_0 + \beta_1 x_1 + \beta_2 x_2),$$

where  $\phi(\cdot)$  is the density of the standard normal distribution  $\beta_2 \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}{[1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)]^2}$

The first case is for a probit model, the second for a logit.

By looking at the equations we can understand why the marginal effect depends on the values of the other variables.

When we compute the marginal effect of  $x_2$  we compute this derivative. The coefficient is multiplying the function that depends on  $x_1$  and  $x_2$ . That's why it will always depend on the values of all the other variables.

It happens in both cases, whether we have a probit or a logit model.

We also know that we can interpret the sign by just looking at the coefficient. Here we can see that we have  $\beta_2$  multiplying a function which is always positive (it works both for the probit and logit models). The sign of the marginal effect is always going to be the same as the sign of  $\beta_2$  in this case.

When we have to compute the marginal effect, we need to choose values of  $x_1$  and  $x_2$ , but which ones do we have to use? In the past, commonly it has been used to calculate the marginal effects at the sample means of  $x_1$  and  $x_2$ . In the sample we find the average sample mean for  $x_1$  and for  $x_2$  and then we use these value in the formula above and we will get a marginal effect.

Now, with better statistical programs, we can do something else, as the sample means could not be representative of anyone in the sample (the mean of the urban variable, for example, is a number between 1 and 0: it would mean that someone lives partly in an urban area, which is nonsensical).

Now we can see that it is done for a reference individual, for example, someone that lives in an urban area and is aged 40. We use these values in the formula and find the marginal effect for someone aged 40 that lives in an urban area. This is a **conditional marginal effect**.

Another alternative would be to compute the formula for every individual in the sample. Once we have a value for each individual, we computer the average of these values and obtain the **average marginal effect**.

**Discrete change in  $P[x_1, x_2]$  due to change in dummy variable  $x_1$ .**

When we think about the change in the variable due to a change in the dummy variable, so we compute these discrete changes:

$$\Pr Pr(x_1 = 1, x_2) - \Pr Pr(x_1 = 0, x_2)$$

Or if we saw the relative ratio:

$$Relative\ Risk = \frac{P(x_1=1, x_2)}{PrPr(x_1=0, x_2)}$$

Then we already see that it will compute the probabilities at both  $x_1$  and  $x_2$ .

It will also compute the probabilities at both  $x_1$  equal to 1 and to 0, so just like for both urban and rural area. Then we have to fix only  $x_2$ , so then we can consider as before the sample mean of  $x_2$ . In our example it could be the average age or the reference individual. For example, we choose an individual aged 40 and use the following formula for a logit model. We then plug in  $age = 40$ .

$$\frac{\exp(\exp(2.1 - 0.22 + 0.03 \times 40))}{1 - \exp(\exp(2.1 - 0.22 + 0.03 \times 40))} - \frac{\exp(\exp(2.1 + 0.03 \times 40))}{1 - \exp(\exp(2.1 + 0.03 \times 40))}$$

In alternative we could just plug every single individual's value and then do the average across all the observation.

### Average marginal effects

This is what we do in Stata. We need to tell Stata that the dummy variables are dummy variables and the way in which we do this is by using the "i.".

Figure 15<sup>47</sup> (Garcia-Gomez, 2022)

```
logit ptuse i.urban age
margins, dydx(*)
```

```
Average marginal effects          Number of obs   =       3000
Model VCE      : OIM

Expression   : Pr(ptuse), predict()
dy/dx w.r.t. : l.urban age
```

	dy/dx	Delta-method Std. Err.	z	P> z	[95% Conf. Interval]	
urban						
Urban	-.0117787	.0084613	-1.39	0.164	-.0283625	.0048051
age	.001381	.0002193	6.30	0.000	.0009513	.0018108

Note: dy/dx for factor levels is the discrete change from the base level.

Stata will compute the marginal effect, which tells us that for factored levels it's a discrete change from the base level. The factor level refers to qualitative variables, or



categorical variables, and it is giving us the difference compared to the reference category. This is what Stata calls the base level and it is what is important to us: we want a discrete change, not a continuous one. If we do not do "i." in the command, Stata will think that this is a continuous variable, getting then the derivative (which makes no sense for dummy variables).

Once we have these coefficients, we can interpret the effects. How do we interpret that marginal effect? We can say that, on average, living in an urban area compared to living in a rural area decreases the probability of public transport use by 1.18 percentage points, *ceteris paribus*/keeping age fixed.

Similarly, thinking about the coefficient of the average effect of age, we can say that on average an additional year of age increases the probability of public transport use by 0.14 percentage points, *ceteris paribus*/keeping urban fixed.

### Conditional marginal effects

Sometimes we want to compute the conditional marginal effect in which we fix particular values of the covariates.

If we fix, for example, urban to 0, which means that an individual lives in a rural area, and we fix age to 80, we then get the conditional marginal effect again, it will always be the difference between the discrete variables. For age is going to be the derivative at that point.

Figure 16<sup>48</sup> (Garcia-Gomez, 2022)

```
margins, dydx(*) at(urban=0 age=80)
```

```
Conditional marginal effects          Number of obs   =       3000
Model VCE      : OIM
```

```
Expression      : Pr(ptuse), predict()
dy/dx w.r.t.    : l.urban age
```

```
at      : urban      =      0
        : age        =      80
```

	dy/dx	Delta-method Std. Err.	z	P> z	[95% Conf. Interval]	
urban						
Urban	-.0039077	.0029011	-1.35	0.178	-.0095938	.0017784
age	.0004022	.000054	7.45	0.000	.0002963	.000508

Note: dy/dx for factor levels is the discrete change from the base level.

We can say that living in an urban area decreases the probability of public transport use by 0.4 percentage point for an individual aged 80.

For age we would say that an additional year of age increases the probability of public transport use by 0.04 percentage points (as we have to multiply the coefficient by 100) for an individual that lives in an urban area and with age equal to 80 years old.

In the interpretation we always have to make clear at which values we are fixing it. When we talk about the binary variables, we have to notice that it is always the urban equal to 1 minus urban equal to 0, then it is for someone aged 80. So, it is always the distance between the two lines, but when we look at age we need to specify if we are looking for someone that lives in urban or rural areas.

### **Are the effects of logit and probit models comparable?**

We cannot compare the  $\beta$ 's directly. There is a formula in which we can roughly compare one to the other, although it is not really doing much. We can compare the set of statistics of the  $\beta$ 's, as we are dividing by the standard error: we are putting things in the same unit and more importantly we can compare the marginal effect and also the goodness of fit measure. We can use these measures to choose between a logit and a probit model.

## **Lecture 35 (4.7) – Marginal Effects with Interactions**

What is the effect of age? Is the effect of age different for people in urban or in rural areas? We could create a new variable and include an **interaction variable** in the models. We create a new variable called `urban_age` (which is `urban*age`) and then include it in the model. The coefficient will differ whether someone lives in a rural or in an urban area. When urban is 0 (the individual lives in a rural area), the coefficient of age will be 0.036. When urban equals to 1, the coefficient is  $0.036 - 0.015$ .

Figure 18<sup>49</sup> (Garcia-Gomez, 2022)

```
gen urban_age = urban*age
logit ptuse urban age urban_age
```

```
Logistic regression                Number of obs    =      3,000
                                   LR chi2(3)        =      54.13
                                   Prob > chi2         =      0.0000
Log likelihood = -637.22158        Pseudo R2       =      0.0407
```

ptuse	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
urban	.1870193	.2740625	0.68	0.495	-.3501334	.7241719
age	.0356552	.0068187	5.23	0.000	.0222908	.0490196
urban_age	-.0150395	.0082609	-1.82	0.069	-.0312307	.0011517
_cons	1.806555	.2137968	8.45	0.000	1.387521	2.225589

What is the problem? Now we want to estimate marginal effects. We can ask Stata to compute them. But Stata needs to understand that the interaction is really related to urban and age. If we are changing age, we also need that the urban\_age changes and that Stata understands that it is part of the same variable. Otherwise, Stata will compute the marginal effects for urban, for age and for urban age, as if this third variable was an independent unrelated variable.

To let Stata know that this is an interaction we compute: `logit ptuse i.urban ## c.age.`

Figure 19<sup>50</sup> (Garcia-Gomez, 2022)

```
Logistic regression                Number of obs    =      3,000
                                   LR chi2(3)        =      54.13
                                   Prob > chi2         =      0.0000
Log likelihood = -637.22158        Pseudo R2       =      0.0407
```

ptuse	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
urban	.1870193	.2740625	0.68	0.495	-.3501334	.7241719
Urban	.0356552	.0068187	5.23	0.000	.0222908	.0490196
age						
urban#c.age						
Urban	-.0150395	.0082609	-1.82	0.069	-.0312307	.0011517
_cons	1.806555	.2137968	8.45	0.000	1.387521	2.225589

In this way we include an interaction between urban (dummy variable, and we use "i.") and age (continuous variable, and we use "c."). If we look at the output that we

get, for the variable urban we get a coefficient for it, then we get a coefficient for age, and also one for the interaction.

These are the same coefficients we had before.

If we compare these outputs with the table before, nothing changes in the estimation of the coefficients. Now Stata recognizes those as being part of the same variable.

This information can be used when we compute the marginal effect.

What do we need to compute the marginal effects? We simply do as we did before.

Figure 20<sup>51</sup> (Garcia-Gomez, 2022)

```

Logistic regression                               Number of obs   =       3,000
                                                  LR chi2(3)      =       54.13
                                                  Prob > chi2     =       0.0000
Log likelihood = -637.22158                    Pseudo R2      =       0.0407
  
```

	ptuse	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
urban						
Urban		.1870193	.2740625	0.68	0.495	-.3501334 .7241719
age		.0356552	.0068187	5.23	0.000	.0222908 .0490196
urban#c.age						
Urban		-.0150395	.0082609	-1.82	0.069	-.0312307 .0011517
_cons		1.806555	.2137968	8.45	0.000	1.387521 2.225589

For example, we compute the average marginal effect and Stata gives us the average marginal effect of urban and age. While computing it, it accounts for the interaction.

One of the reasons to compute this interaction is because we are interested about the different marginal effects for people that live in urban and rural area, and the different marginal effect of age for people living in urban compared to rural areas. We think that this has to be different not only because the model is linear, but because there are some differences per se without an interaction, as we have two curves, and the slope of age will be different in each different curve. This allows for further differences and not only the ones that are imposed by the functional form of the probit or the logit model.

As we are interested in the different results between the different marginal effects of age for people that live in an urban and in a rural area, we could get those results through Stata:

Figure 21<sup>52</sup> (Garcia-Gomez, 2022)

```

margins, dydx(*) at(urban=1)
Average marginal effects           Number of obs   =       3,000
Model VCE      : OIM
Expression     : Pr(ptuse), predict()
dy/dx w.r.t.  : 1.urban age
at             : urban              =           1
-----
            |              Delta-method
            |              dy/dx   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
    urban |
Urban    |  -.0117479   .0084503   -1.39   0.164   -.0283101   .0048144
    age   |   .0011991   .0002829    4.24   0.000    .0006445   .0017537
-----
Note: dy/dx for factor levels is the discrete change from the base level.

```

```

margins, dydx(*) at(urban=0)
Average marginal effects           Number of obs   =       3,000
Model VCE      : OIM
Expression     : Pr(ptuse), predict()
dy/dx w.r.t.  : 1.urban age
at             : urban              =           0
-----
            |              Delta-method
            |              dy/dx   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
    urban |
Urban    |  -.0117479   .0084503   -1.39   0.164   -.0283101   .0048144
    age   |   .0016801   .0003594    4.68   0.000    .0009757   .0023844
-----
Note: dy/dx for factor levels is the discrete change from the base level.

```

We can compute the marginal effect for those who live in an urban area and then for those who live in a rural area. Once we have them, we can see if the average effect of age is larger or smaller for those that live in an urban area compared to those that live in a rural area.

There are other case in which we need to pay attention about the relationship between the two variables. For example, we could have **polynomials** of continuous variables. We could estimate such a model in two ways:

1. As we did in the past, if we want, for example, to construct the variable age2 (age squared), which is age\*age, and then include the variable in our model.
2. We could tell Stata that we want to have age and age2 in our model, and that these variables are related.

The coefficients calculated in these two ways will be the same, but to compute the marginal effects we need to use the second option, otherwise Stata will think that

age and age2 are two independent variables, and it will not look at their joint coefficients when we want to compute the marginal effects.

These can be useful for other applications and alternatives which we will not analyse.

## Lecture 36 (4.8) – Odds ratio in logit

We are focusing on the **odds ratio** in the logit model. There is a common way to interpret the results in a logit model (it is not recommended to us, as it is often wrongly interpreted).

Odds ratios are commonly used in the interpretation of the logit model, especially in the past. Now we are just using the marginal effects.

What is an odds ratio?

First of all, we have to define the odds, which is a ratio between the probability of using public transport use in our example, when someone lives in an urban area, divided by the probability of not using public transport for someone who lives in an urban area, which is when y is equal to 0 conditional on  $x_1 = 1$ :

$$\frac{\text{PrPr}(x_1=1, x_2)}{\text{PrPr}(x_1=1, x_2)}$$

This expression is what we refer to as odds (these are the odds of using public transportation for urban). We can define the same expression for those that live in a rural area:

$$\frac{\text{PrPr}(x_1=0, x_2)}{\text{PrPr}(x_1=0, x_2)}$$

These are the odds of using public transportation for rural areas.

Once we have each of these two ratios, the odds ratio is the ratio of the odds:

$$\text{Odds Ratio} = \frac{\frac{\text{PrPr}(x_1=1, x_2)}{\text{PrPr}(x_1=1, x_2)}}{\frac{\text{PrPr}(x_1=0, x_2)}{\text{PrPr}(x_1=0, x_2)}}$$

This is hard to interpret. But we can see that it simplifies to the **exponential of  $\beta_1$** :

$$OR = \frac{\frac{\PrPr(x_1=1, x_2)}{\PrPr(x_1=1, x_2)}}{\frac{\PrPr(x_1=0, x_2)}{\PrPr(x_1=0, x_2)}} = \exp(\beta_1)$$

$\beta_1$  is the coefficient for  $x_1$  that is whether the dummy variable or whether the person lives in an urban or rural area. This is a simple transformation and now we can compute average marginal effects with the software we had, whether before researchers did not have Stata.

Once we have the estimate from the logit, we can apply the exponential and then we can get the odds ratio. This is appealing, but the problem is that if we say that the odds ratio increases or decreases and we give the amount, it is very hard to understand the meaning of the ratio.

An exception of this happens when the odds ratio can be interpreted as risk ratio. But this is an extreme event, as the probability that  $y$  is equal to 0 for  $x_1$  equal to zero must be close to one. It is an event that happens with a very low probability.

In formula, if  $\Pr(x_1 = 0, x_2) \rightarrow 1$  and  $\Pr(x_1 = 1, x_2) \rightarrow 1$ , then we have that  $OR = RR$ .

In papers we often find that people interpret the odds ratio as it was a risk ratio, but for events that do not have such a low incidence, therefore this interpretation is incorrect.

## Lecture 37 (4.9) – Hypothesis testing

We are focusing on hypothesis testing in nonlinear models for binary data, especially when we are interested in testing a single hypothesis and also multiple hypothesis.

We will start with **single hypothesis testing in logit and probit**. We talk about single hypothesis when we want to test the null hypothesis that the population parameter is equal to a given value versus the alternative, that is the different to that value:

$$H_0: \beta_k = \beta^* \text{ vs } H_1: \beta_k \neq \beta^*$$

Is similar to what we saw in the first module, as the test statistic works exactly as it worked with ordinary least squares:

$$z = \frac{\hat{\beta}_k - \beta^*}{\hat{\sigma}_{\hat{\beta}_k}} \sim N(0, 1)$$

We can use the same test statistic and then use the z-statistic comparing the value that we get with the critical values, or the p-values in the usual way.

Figure 22<sup>53</sup> (Garcia-Gomez, 2022)

```

Probit regression                               Number of obs   =       3000
                                                LR chi2(2)      =       54.65
                                                Prob > chi2     =       0.0000
Log likelihood = -636.96328                    Pseudo R2      =       0.0411

```

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
ptuse					
urban	-.1178419	.0782359	-1.51	0.132	-.2711814 .0354976
age	.013089	.0018674	7.01	0.000	.009429 .016749
_cons	1.201824	.0834919	14.39	0.000	1.038183 1.365465

This is something we already did in the previous lectures, so we can look at the coefficients and then if we are interested in seeing whether our given coefficient is statistically significantly different from 0 or not. We can just look at the z-statistic and then compare with the critical values (which are the same of when we were using OLS), or we can interpret the p-values in the usual way.

For example, we see that the differences in public transport use between those living in urban compared to those living in rural areas are not statistically significantly different from 0 at 10% significance level.

Or that urban is a statistically significant at a 10% significance level. These two things are comparable. While if we look at the effect of age on public transport use, remembering that we cannot interpret the size of the coefficient, we can look at the hypothesis testing on whether it is statistically significantly different from 0 or not. Then we see that in this case the p-value is 0.000 (so it is a very small p-value). We can then conclude that the effect of age on public transport use is statistically significantly different from 0 at a 1% significance level. We interpret it in the exact way we did with the OLS model.

For **multiple hypothesis testing in logit and probit** we are going to analyse the two most common tests done for multiple hypothesis testing. When we have multiple



hypothesis, we refer to two types of models, the restricted one ( $H_0$ ) and the unrestricted one ( $H_1$ ). Then all that null hypothesis is that the restricted model is true versus the unrestricted model is true.

The restricted model is a model in which we impose some restrictions such as, for example, that some of the coefficients are equal to zero, while in the unrestricted model they are not.

The **likelihood ratio test**, which is one of the two tests, is based on both the restricted and the unrestricted models, therefore we have to estimate them both. We start by estimating the restricted model, where we impose some restrictions on the parameters and then get a log-likelihood ratio of this model.

Then we also compute the unrestricted model and then we get the log-likelihood of that model.

We compare then those two log-likelihood ratios:

$$LR = 2(\log \log L_{UR} - \log \log L_R) \sim \chi^2(q)$$

Where  $q$  is the number of restrictions we have set in our model (e.g.,  $q=3$  means that we are testing whether three coefficients are equal to 0 in our null hypothesis).

If they are different enough, then we conclude that we reject the null hypothesis. If they are very similar, the unrestricted model is not really giving much compared to the restricted model, so we would not reject the null hypothesis that the restricted model is true.

We do this by computing the formula above: two times the difference of the log of the two models follows an **Akaike square distribution** with  $q$  degrees of freedom.

The other test is the **Wald test**, that it is only based on the unrestricted model:

$$W \sim \chi^2(q)$$

The exact formula is not analysed in this course. This is computed by Stata. The Wald test is also distributed following an Akaike square distribution with  $q$  degrees of freedom. In this case too,  $q$  is equal to the number of restrictions, and we will be able to interpret the results in the usual way.

## Lecture 38 (4.10) – Goodness of fit

We are focusing on the measures of goodness of fit, and in particular about the most commonly used measures of goodness of fit in logit and probit models.

Comparing goodness of fit is essential if our objective is to obtain the model that fits the data best. If we are interested in obtaining a causal effect, the same discussion, the same issues we had in previous topics also apply to this case.

The first measure we focus on is the **log-likelihood**, which we cannot interpret per se. We can compare different models and see which ones give us the largest log-likelihood ratio. The log-likelihood ratio will always be negative, so the higher likelihood ratio is going to be the one closer to 0. When just comparing log-likelihood, one of the problems is that the log-likelihood is always going to increase when we add additional variables, but sometimes the addition of a variable will increase it of just a little. As what we had in the linear case where we had the R-squared and the adjusted R-squared, where we were penalizing by the number of parameters, now we want to do the same. To do this we use the information criterion that penalizes the log-likelihood with a number of coefficients in the model. But how?

The **Akaike information criterion** computes this formula:

$$AIC = \frac{-2\log L + 2(k+1)}{N}$$

The “-2” in the formula means that the log-likelihood is negative. So now that Akaike is going to be positive and then we want the likelihood to be close to 0 as much as possible, so when we make it positive, the lower the Akaike the better our goodness of fit are.  $k$  is the number of coefficients we estimate in our model.

Our aim is to have the lowest Akaike possible, and we will pick the model with the lowest Akaike information criteria.

Another goodness of fit measure is the **percentage of explained variation**, or the **Efron’s  $R^2$** , which compares the variation that is explained by our data with the total variation in our explanatory variables.

$$1 - \frac{\sum_{i=1}^N (y_i - \hat{\pi}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

The  $\hat{\pi}$  is the predicted probability, so we have 1 minus the variation predicted by our model divided by total variation in the data.

Another measure is the **pseudo-R<sup>2</sup>** or **McFadden's R<sup>2</sup>**, which tells us the extent to which the model is an improvement over one with just a constant term.

$$1 - \frac{\log \log L_{UR}}{\log \log L_0}$$

We have this ratio in which we find the log-likelihood of the unrestricted model. We then compare this log-likelihood with the log-likelihood of a model that only includes a constant. If this model is not a real improvement over the one that only has a constant, the ratio will be equal to 1. Then one minus one will be 0. The larger our improvement (so the closer to zero), the higher our likelihood. Compare to the log-likelihood of the model with our constant that it is always going to be larger in absolute value, but both will be negative. Then the smaller this ratio, it means that our log-likelihood model of the unrestricted model is closer to zero, so then the difference will be closer to one. So, the closer the pseudo-R<sup>2</sup> is to 1, the better the fitting of our model.

The last measure is the **Count R<sup>2</sup> (percentage of correct predictions)**. It tells us the percentage of correct predictions. When Stata computes this, for every individual, for every observation in our sample, we have the observed  $y$  (so if the person used public transport, 1, or not, 0). Then for every single person we can obtain our predicted probability. We estimate the model and then we have seen how we can compute these predicted probabilities.

Figure 23<sup>54</sup> (Garcia-Gomez, 2022)

Observed $Y$	Model $\hat{P}[Y = 1]$	Predicted $Y$ (1, if $\hat{P}[Y = 1] > 0.5$ 0 otherwise)	Correct prediction?
0	0.23	0	yes
1	0.62	1	yes
0	0.05	0	yes
1	0.07	0	no
0	0.67	1	no
0	0.12	0	yes
0	0.02	0	yes
1	0.54	1	yes
1	0.42	0	no
0	0.17	0	yes
1	0.68	1	yes
...	...	...	...

The predicted  $y$  (and not the predicted probability), so whether a person is going to use public transport or not, is 1 if the predicted probability is larger than 0.5 and 0 otherwise. The first individual has a predicted probability of 0.23, which is smaller than 0.5, so that individual will not use the public transport. Is this a correct prediction? We can see the observed and note that it is correct. We have correct predictions until the fourth individual, for which we do not have a correct prediction. We can do this for every individual in the sample and then see the share of the predictions that were correct and the one of those that were not, and that gives us the Count  $R^2$ .

# Reference list

- Garcia-Gomez P. (2022) Lecture 1 – Introduction, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Lecture 2 – Estimation and Interpretation, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Lecture 4 – Assumptions for inference, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Lecture 5 – Inference – One parameter, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Lecture 7 – Interpretation categorical variables, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Module 2 Lecture 10 – Beyond statistical significance, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-linear-regression-models>
- Garcia-Gomez P. (2022) Module 2 Lecture 1 – Introduction, PowerPoint slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-endogeneity-and-instrumental-variables-estimation>
- Börsch-Supan A. (coordinator) (2004) The Survey of Health, Ageing and Retirement in Europe
- Wooldridge J.M. (2014) Introductory Econometrics – A Modern Approach
- Garcia-Gomez P. (2022) Module 2 Lecture 3 – Collider Bias, Powerpoint Slides, retrieved from: <https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-endogeneity-and-instrumental-variables-estimation>
- Cunningham S. (2020) <https://www.scunning.com/causalinferencenorap.pdf>

- Garcia-Gomez P. (2022) Module 2 Lecture 8 – Potential outcomes and DAGs, PowerPoint slides, retrieved from:  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-endogeneity-and-instrumental-variables-estimation>
- Riumalló Herl C. (2022), Module 3 Lecture 2 – What is panel data?, Powepoint slides, retrieved from  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-panel-data>
- Riumalló Herl C. (2022), Module 3 Lecture 3 – Notations, Powepoint slides, retrieved from  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-panel-data>
- Riumalló Herl C. (2022), Module 3 Lecture 8 – Estimator choice, Powepoint slides, retrieved from  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-panel-data>
- Riumalló Herl C. (2022), Module 3 Lecture 10 – Introduction to Difference-in-Differences, Powepoint slides, retrieved from  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-panel-data>
- Garcia-Gomez P. (2022) Module 4 Lecture 1 – Binary data: introduction, retrieved from:  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>
- Garcia-Gomez P. (2022) Module 4 Lecture 2 – Linear regression model for binary data (LPM), PowerPoint slides, retrieved from:  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>
- Garcia-Gomez P. (2022) Module 4 Lecture 3 – Nonlinear Models for binary data, PowerPoint slides, retrieved from:  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>
- Garcia-Gomez P. (2022) Module 4 Lecture 4 – Estimation of Nonlinear Models for binary data, PowerPoint slides, retrieved from:  
<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>
- Garcia-Gomez P. (2022) Module 4 Lecture 5 – Interpretation with Graphs, PowerPoint slides, retrieved from:

<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>

Garcia-Gomez P. (2022) Module 4 Lecture 6 – Interpretation with Marginal Effects, PowerPoint slides, retrieved from:

<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>

Garcia-Gomez P. (2022) Module 4 Lecture 7 – Marginal Effects with Interactions, PowerPoint slides, retrieved from:

<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>

Garcia-Gomez P. (2022) Module 4 Lecture 9 – Hypothesis testing, PowerPoint slides, retrieved from:

<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>

Garcia-Gomez P. (2022) Module 4 Lecture 10 – Goodness of fit, PowerPoint slides, retrieved from:

<https://canvas.eur.nl/courses/40137/pages/introduction-to-empirical-methods-binary-data>