

ScPoEconometrics Advanced

Recap 2

Bluebery Planterose SciencesPo Paris 2023-01-31

Recap 2

- Last time, we refreshed our basic OLS knowledge
- Today we continue and look at more than one explanatory variable, and associated problems

- But, why more than one variable?
- Like, **how many** other variables?
- And, above all: which ones ? 🤔



Recap 2

- Last time, we refreshed our basic OLS knowledge
- Today we continue and look at more than one explanatory variable, and associated problems

- But, why more than one variable?
- Like, **how many** other variables?
- And, above all: which ones ? 🤔

We will remember what we meant by **a model**.



Back to the STAR Experiment

- Remember what we learned about the STAR Experiment
- What is the causal impact of class size on test scores?

 $ext{score}_i = eta_0 + eta_1 ext{classize}_i + u_i \ ?$

• We use a **model** to order our thoughts about how a causal impact is determined.



•

Back to the STAR Experiment



Multiple Variables

Let's augment our model with more variables:

$$y=eta_0+eta_1x_1+eta_2x_2+eta_3x_3+u$$









Omitted-variable bias (OVB) arises when we omit a variable that

- 1. affects our outcome variable y
- 2. correlates with an explanatory variable x_j

As it's name suggests, this situation leads to bias in our estimate of β_j .

Omitted-variable bias (OVB) arises when we omit a variable that

- 1. affects our outcome variable y
- 2. correlates with an explanatory variable x_j

As it's name suggests, this situation leads to bias in our estimate of β_j .

Note: OVB Is not exclusive to multiple linear regression, but it does require multiple variables affect *y*.

Example

Let's imagine a simple model for the amount individual i gets paid

$$\operatorname{Pay}_i = \beta_0 + \beta_1 \operatorname{School}_i + \beta_2 \operatorname{Male}_i + u_i$$

where

- School_i gives *i*'s years of schooling
- $Male_i$ denotes an indicator variable for whether individual i is male.

thus

- β_1 : the returns to an additional year of schooling (*ceteris paribus*)
- β_2 : the premium for being male (*ceteris paribus*) If $\beta_2 > 0$, then there is discrimination against women—receiving less pay based upon gender.

Example, continued

From our population model

$$\operatorname{Pay}_i = \beta_0 + \beta_1 \operatorname{School}_i + \beta_2 \operatorname{Male}_i + u_i$$

If a study focuses on the relationship between pay and schooling, *i.e.*,

$$egin{aligned} ext{Pay}_i &= eta_0 + eta_1 ext{School}_i + (eta_2 ext{Male}_i + u_i) \ & ext{Pay}_i &= eta_0 + eta_1 ext{School}_i + arepsilon_i \end{aligned}$$

where $arepsilon_i=eta_2\mathrm{Male}_i+u_i.$

We used our exogeneity assumption to derive OLS' unbiasedness. But even if E[u|X] = 0, it is not true that $E[\varepsilon|X] = 0$ so long as $\beta_2 \neq 0$.

Specifically, $oldsymbol{E}[arepsilon| ext{Male}=1]=eta_2+oldsymbol{E}[u| ext{Male}=1]
eq 0.$

Example, continued

From our population model

$$\operatorname{Pay}_i = \beta_0 + \beta_1 \operatorname{School}_i + \beta_2 \operatorname{Male}_i + u_i$$

If a study focuses on the relationship between pay and schooling, *i.e.*,

$$egin{aligned} ext{Pay}_i &= eta_0 + eta_1 ext{School}_i + (eta_2 ext{Male}_i + u_i) \ & ext{Pay}_i &= eta_0 + eta_1 ext{School}_i + arepsilon_i \end{aligned}$$

where $arepsilon_i=eta_2\mathrm{Male}_i+u_i.$

We used our exogeneity assumption to derive OLS' unbiasedness. But even if E[u|X] = 0, it is not true that $E[\varepsilon|X] = 0$ so long as $\beta_2 \neq 0$.

Specifically, $m{E}[arepsilon|\mathrm{Male}=1]=eta_2+m{E}[u|\mathrm{Male}=1]
eq 0$. Now OLS is biased.

Example, continued

Let's try to see this result graphically.

The population model:

 $\mathrm{Pay}_i = 20 + 0.5 imes \mathrm{School}_i + 10 imes \mathrm{Male}_i + u_i$

Our regression model that suffers from omitted-variable bias:

$$\mathbf{Pay}_i = \hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_1 \times \mathbf{School}_i + e_i$$

Finally, imagine that women, on average, receive more schooling than men.

Example, continued: $\operatorname{Pay}_i = 20 + 0.5 \times \operatorname{School}_i + 10 \times \operatorname{Male}_i + u_i$

The relationship between pay and schooling.



Example, continued: $Pay_i = 20 + 0.5 \times School_i + 10 \times Male_i + u_i$

Biased regression estimate: $\widehat{\mathrm{Pay}}_i = 31.3 + -0.9 imes \mathrm{School}_i$



Example, continued: $Pay_i = 20 + 0.5 \times School_i + 10 \times Male_i + u_i$

Recalling the omitted variable: Gender **female** and **male**



Example, continued: $Pay_i = 20 + 0.5 \times School_i + 10 \times Male_i + u_i$

Recalling the omitted variable: Gender **female** and **male**



Example, continued: $Pay_i = 20 + 0.5 \times School_i + 10 \times Male_i + u_i$

Unbiased regression estimate: $\widehat{\mathrm{Pay}}_i = 20.9 + 0.4 \times \mathrm{School}_i + 9.1 \times \mathrm{Male}_i$



Solutions

- 1. Don't omit variables 🤤
- 2. Instrumental variables and two-stage least squares (coming soon): If we could find something that **only** affects x_1 but *not* the omitted variable, we can make progress!
- 3. Use multiple observations for the same unit *i*: panel data.

Warning: There are situations in which neither solution is possible.

Solutions

- 1. Don't omit variables 🤤
- 2. Instrumental variables and two-stage least squares (coming soon): If we could find something that **only** affects x_1 but *not* the omitted variable, we can make progress!
- 3. Use multiple observations for the same unit *i*: panel data.

Warning: There are situations in which neither solution is possible.

- 1. Proceed with caution (sometimes you can sign the bias).
- 2. The key is to have a mental map of *should* belong to the model.

13 / 68

Continuous variables

Consider the relationship

$$\operatorname{Pay}_i = eta_0 + eta_1\operatorname{School}_i + u_i$$

where

- Pay_i is a continuous variable measuring an individual's pay
- School_{*i*} is a continuous variable that measures years of education

Interpretations

- β_0 : the *y*-intercept, *i.e.*, Pay when School = 0
- β_1 : the expected increase in Pay for a one-unit increase in School

Continuous variables

Consider the model

$$y=eta_0+eta_1\,x+u$$

Differentiate the model:

$$rac{dy}{dx}=eta_1$$

Task 1: Interpretation (4 minutes)

1. Load the wage1 dataset from the wooldridge package. you may have to install this first.

- 2. Run <a>skimr::skim on the dataset to get an overview. what is the fraciton of nonwhite in the data?
- 3. Regressing wage on education and tenure, what is the interpretation of the tenure coefficient? You may need to consult ?wage1 here.

Categorical variables

Consider the relationship

$$\operatorname{Pay}_i = eta_0 + eta_1 \operatorname{Female}_i + u_i$$

where

- Pay_i is a continuous variable measuring an individual's pay
- Female $_i$ is a binary/indicator variable taking 1 when i is female

Interpretations

- β_0 : the expected Pay for males (*i.e.*, when Female = 0)
- β_1 : the expected difference in Pay between females and males
- $\beta_0 + \beta_1$: the expected Pay for females

Categorical variables

Derivations

$$oldsymbol{E}[ext{Pay}| ext{Male}] = oldsymbol{E}[eta_0+eta_1 imes 0+u_i] \ = oldsymbol{E}[eta_0+0+u_i] \ = eta_0$$

Categorical variables

Derivations

$$oldsymbol{E}[ext{Pay}| ext{Male}] = oldsymbol{E}[eta_0+eta_1 imes 0+u_i] \ = oldsymbol{E}[eta_0+0+u_i] \ = eta_0$$

$$egin{aligned} oldsymbol{E}[ext{Pay}| ext{Female}] &= oldsymbol{E}[eta_0+eta_1 imes1+u_i] \ &= oldsymbol{E}[eta_0+eta_1+u_i] \ &= eta_0+eta_1 \end{aligned}$$

Categorical variables

Derivations

$$oldsymbol{E}[ext{Pay}| ext{Male}] = oldsymbol{E}[eta_0+eta_1 imes 0+u_i] \ = oldsymbol{E}[eta_0+0+u_i] \ = eta_0$$

$$oldsymbol{E}[ext{Pay}| ext{Female}] = oldsymbol{E}[eta_0+eta_1 imes1+u_i] \ = oldsymbol{E}[eta_0+eta_1+u_i] \ = eta_0+eta_1$$

Note: If there are no other variables to condition on, then $\hat{\beta}_1$ equals the difference in group means, *e.g.*, $\overline{x}_{\text{Female}} - \overline{x}_{\text{Male}}$.

Note₂: The *holding all other variables constant* interpretation also applies for categorical variables in multiple

Categorical variables

 $y_i = eta_0 + eta_1 x_i + u_i$ for binary variable $x_i = \{0, 1\}$





Categorical variables

 $y_i = \beta_0 + \beta_1 x_i + u_i \text{ for binary variable } x_i = \{0, 1\}$ $\hat{\beta}_0 + \hat{\beta}_1 = \text{mathrmGroup}$ $\hat{\beta}_0 = \text{mathrmGroup}_0$

Task 2: Categorical Variables (3 Minutes)

- Continue with the wage1 dataset.
- Now regress wage on female. What is E[wage|male]?
- Add married to the regression. Now what is E[wage|female,not married]?

Interactions

Interactions allow the effect of one variable to change based upon the level of another variable.

Examples

- 1. Does the effect of schooling on pay change by gender?
- 2. Does the effect of gender on pay change by race?
- 3. Does the effect of schooling on pay change by experience?
Interactions

Previously, we considered a model that allowed women and men to have different wages, but the model assumed the effect of school on pay was the same for everyone:

 $\operatorname{Pay}_i = eta_0 + eta_1 \operatorname{School}_i + eta_2 \operatorname{Female}_i + u_i$

but we can also allow the effect of school to vary by gender:

 $\operatorname{Pay}_i = eta_0 + eta_1\operatorname{School}_i + eta_2\operatorname{Female}_i + eta_3\operatorname{School}_i imes\operatorname{Female}_i + u_i$

Interactions

The model where schooling has the same effect for everyone (**F** and **M**):



Schooling

Interactions

The model where schooling's effect can differ by gender (**F** and **M**):



Schooling

Interactions

Interpreting coefficients can be a little tricky with interactions, but the key[†] is to carefully work through the math.

$$\operatorname{Pay}_i = eta_0 + eta_1\operatorname{School}_i + eta_2\operatorname{Female}_i + eta_3\operatorname{School}_i imes\operatorname{Female}_i + u_i$$

Expected returns for an additional year of schooling for women:

$$oldsymbol{E}[ext{Pay}_i| ext{Female} \wedge ext{School} = \ell + 1] - oldsymbol{E}[ext{Pay}_i| ext{Female} \wedge ext{School} = \ell] = \ oldsymbol{E}[eta_0 + eta_1(\ell+1) + eta_2 + eta_3(\ell+1) + u_i] - oldsymbol{E}[eta_0 + eta_1\ell + eta_2 + eta_3\ell + u_i] = \ eta_1 + eta_3$$

+ As is often the case with econometrics.

Interactions

Interpreting coefficients can be a little tricky with interactions, but the key[†] is to carefully work through the math.

$$\operatorname{Pay}_i = eta_0 + eta_1\operatorname{School}_i + eta_2\operatorname{Female}_i + eta_3\operatorname{School}_i imes\operatorname{Female}_i + u_i$$

Expected returns for an additional year of schooling for women:

$$oldsymbol{E}[ext{Pay}_i| ext{Female} \wedge ext{School} = \ell + 1] - oldsymbol{E}[ext{Pay}_i| ext{Female} \wedge ext{School} = \ell] = oldsymbol{E}[eta_0 + eta_1(\ell+1) + eta_2 + eta_3(\ell+1) + u_i] - oldsymbol{E}[eta_0 + eta_1\ell + eta_2 + eta_3\ell + u_i] = eta_1 + eta_3$$

Similarly, β_1 gives the expected return to an additional year of schooling for men. Thus, β_3 gives the **difference in the returns to schooling** for women and men.

+ As is often the case with econometrics.

Task 3: Interactions (4 minutes)

- Same dataset!
- Regress wage on experience, female indicator and their interaction. What is the interpretation of all the coefficients here? Can you distinguish them from zero?
- What is the expected wage for a male with 5 years of experience?

Log-linear specification

In economics, you will frequently see logged outcome variables with linear (non-logged) explanatory variables, *e.g.*,

```
\log(\operatorname{price}_i) = eta_0 + eta_1 \operatorname{bdrms}_i + u_i
```

This specification changes our interpretation of the slope coefficients.

Log-linear specification

Interpretation

- A one-unit increase in our explanatory variable increases the outcome variable by *approximately* $\beta_1 \times 100$ percent.
- *Example:* An additional bedroom increases sales prices of a house by *approximately* 16 percent (for $\beta_1=0.16$).



Log-linear specification

Consider the log-linear model

$$\log(y)=eta_0+eta_1\,x+u$$

and differentiate

$$rac{dy}{y}=eta_1 dx$$

So a marginal change in x (i.e., dx) leads to a $\beta_1 dx$ percentage change in y.

Log-linear specification

What about that **approximation** part?

An additional bedroom increases sales prices of a house by *approximately* 16 percent (for $\beta_1 = 0.16$).

- $\%\Delta ypprox 0.16 imes 100=16\%.$
- Good approximation as long as Δy is not too big.
- We approximate

$$\log\!\left(rac{\Delta y}{y_0}+1
ight)pproxrac{\Delta y}{y_0}$$

Log-linear specification

What about that **approximation** part?

An additional bedroom increases sales prices of a house by *approximately* 16 percent (for $\beta_1 = 0.16$).

- $\%\Delta ypprox 0.16 imes 100=16\%.$
- Good approximation as long as Δy is not too big.
- We approximate

$$\log\!\left(rac{\Delta y}{y_0}+1
ight)pproxrac{\Delta y}{y_0}$$



Log-linear specification

What about that **approximation** part?

An additional bedroom increases sales prices of a house by *approximately* 16 percent (for $\beta_1 = 0.16$).

- $\%\Delta ypprox 0.16 imes 100=16\%.$
- Good approximation as long as Δy is not too big.
- We approximate

$$\log\!\left(rac{\Delta y}{y_0}+1
ight)pproxrac{\Delta y}{y_0}$$

• The **exact** formula is

$$\%\Delta y = 100 imes (\exp(\Delta x eta) - 1)$$

• In our case:

$$\%\Delta y = 100 imes (\exp(0.16) - 1) = 17.3$$

Task 4

- same Dataset!
- Now regress *log wage* on education and tenure. How does the interpretation of the coefficient on education change?

Log-log specification

Similarly, econometricians frequently employ log-log models, in which the outcome variable is logged *and* at least one explanatory variable is logged

```
\log(	ext{price}_i) = eta_0 + eta_1 \, \log(	ext{sqrft}_i) + u_i
```

Interpretation:

- A one-percent increase in x will lead to a β_1 percent change in y.
- Often interpreted as an elasticity.

Log-log specification

Consider the log-log model

$$\log(y) = eta_0 + eta_1\,\log(x) + u$$

and differentiate

$$rac{dy}{y}=eta_1rac{dx}{x}$$

which says that for a one-percent increase in x, we will see a β_1 percent increase in y. As an elasticity:

$$rac{dy}{dx}rac{x}{y}=eta_1$$

Task 5

- Load the hprice1 dataset from the wooldridge package.
- Regress log price on log sqrft. What is the interpretation on log(sqrft)?
- What is the E[price|sqrft = 115] (Caution! not log price!)

Log-log specification

- a 1% increase in square footage of the house leads to a 0.873% increase in sales price.
- Notice the absence of *units* here (it's all in **percent** terms of both variables involved).

Log-linear with a binary variable

Note: If you have a log-linear model with a binary indicator variable, the interpretation for the coefficient on that variable changes.

Consider again

$$\log(y_i)=eta_0+eta_1x_1+u_i$$

for binary variable x_1 .

The *approximate* interpretation of β_1 is as before:

When x_1 changes from 0 to 1, y will change by $100 imes eta_1$ percent.

```
#>
#> Call:
#> lm(formula = log(price) ~ log(lotsize) + log(sqrft) + bdrms +
      colonial, data = hprice1)
#>
#>
#> Residuals:
       Min
                 10 Median
                                  30
#>
                                          Max
#> -0.69479 -0.09750 -0.01619 0.09151 0.70228
#>
#> Coefficients:
#>
               Estimate Std. Error t value Pr(>|t|)
#> (Intercept) -1.34959
                          0.65104 -2.073
                                           0.0413 *
#> log(lotsize) 0.16782
                         0.03818
                                  4.395 3.25e-05 ***
#> log(sqrft)
                0.70719
                        0.09280 7.620 3.69e-11 ***
                        0.02872 0.934
#> bdrms
         0.02683
                                           0.3530
#> colonial 0.05380
                                  1.202
                          0.04477
                                           0.2330
#> ---
#> Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 0.1841 on 83 degrees of freedom
#> Multiple R-squared: 0.6491, Adjusted R-squared: 0.6322
#> F-statistic: 38.38 on 4 and 83 DF, p-value: < 2.2e-16
```

Approximate

• When *colonial* changes from 0 to 1 (i.e. house *becomes* colonial), y will change by $100 \times \beta_1 = 5.37$ percent.

Exact

• When colonial changes from 1 to 0, y will change by $100 imes \left(e^{eta_1}-1
ight)=5.52$ percent.

Is there more?

Up to this point, we know OLS has some nice properties, and we know how to estimate an intercept and slope coefficient via OLS.

Our current workflow:

- Get data (points with *x* and *y* values)
- Regress y on x
- Plot the OLS line (i.e., $\hat{y}=\hat{eta}_0+\hat{eta}_1$)
- Done?

But how do we actually **learn** something from this exercise?

Linkup with Intro Course

This is related to *Intro Course material*:

- 1. Sampling
- 2. Hypothesis Testing
- **3.** Regression Inference

There is more

But how do we actually **learn** something from this exercise?

- Based upon our value of $\hat{\beta}_1$, can we rule out previously hypothesized values?
- How confident should we be in the precision of our estimates?
- How well does our model explain the variation we observe in *y*?

We need to be able to deal with uncertainty. Enter: Inference.

Learning from our errors

As our previous simulation pointed out, our problem with **uncertainty** is that we don't know whether our sample estimate is *close* or *far* from the unknown population parameter.[†]

However, all is not lost. We can use the errors $(e_i = y_i - \hat{y}_i)$ to get a sense of how well our model explains the observed variation in y.

When our model appears to be doing a "nice" job, we might be a little more confident in using it to learn about the relationship between y and x.

Now we just need to formalize what a "nice job" actually means.

+: Except when we run the simulation ourselves—which is why we like simulations.

Learning from our errors

First off, we will estimate the variance of u_i (recall: $Var(u_i) = \sigma^2$) using our squared errors, *i.e.*,

$$s^2 = rac{\sum_i e_i^2}{n-k}$$

where k gives the number of slope terms and intercepts that we estimate (*e.g.*, β_0 and β_1 would give k = 2).

 s^2 is an unbiased estimator of $\sigma^2.$

Learning from our errors

We know that the variance of \hat{eta}_1 (for simple linear regression) is

$$\mathrm{Var}ig({\hateta}_1ig) = rac{s^2}{\sum_i ig(x_i - \overline{x}ig)^2}$$

which shows that the variance of our slope estimator

1. increases as our disturbances become noisier 2. decreases as the variance of x increases

Learning from our errors

More common: The **standard error** of $\hat{\beta}_1$

$$\hat{\operatorname{SE}}ig({\hateta}_1ig) = \sqrt{rac{s^2}{\sum_i ig(x_i - \overline{x}ig)^2}}$$

Recall: The standard error of an estimator is the standard deviation of the estimator's distribution.

Learning from our errors

Standard error output is standard in R's 1m:

Learning from our errors

We use the standard error of $\hat{\beta}_1$, along with $\hat{\beta}_1$ itself, to learn about the parameter β_1 .

After deriving the distribution of $\hat{\beta}_1$,⁺ we have two (related) options for formal statistical inference (learning) about our unknown parameter β_1 :

- **Confidence intervals:** Use the estimate and its standard error to create an interval that, when repeated, will generally^{††} contain the true parameter.
- **Hypothesis tests:** Determine whether there is statistically significant evidence to reject a hypothesized value or range of values.

+: *Hint:* it's normal with the mean and variance we've derived/discussed above)
 ++: *E.g.*, Similarly constructed 95% confidence intervals will contain the true parameter 95% of the time.

Confidence intervals

We construct (1-lpha)-level confidence intervals for eta_1

$${\hat eta}_1 \pm t_{lpha/2,{
m df}} \; {
m \hat{SE}} \Big({\hat eta}_1 \Big)$$

 $t_{lpha/2,{
m df}}$ denotes the lpha/2 quantile of a t dist. with n-k degrees of freedom.

Confidence intervals

We construct (1-lpha)-level confidence intervals for eta_1

$$\hat{{eta}}_1 \pm t_{lpha/2,{
m df}} \; \hat{
m SE} \Big(\hat{{eta}}_1 \Big)$$

For example, 100 obs., two coefficients (i.e., $\hat{\beta}_0$ and $\hat{\beta}_1 \implies k = 2$), and $\alpha = 0.05$ (for a 95% confidence interval) gives us $t_{0.025, 98} = -1.98$



Confidence intervals

We construct (1-lpha)-level confidence intervals for eta_1

$$\hat{eta}_1 \pm t_{lpha/2, ext{df}} \; \hat{ ext{SE}}\!\left(\hat{eta}_1
ight)$$

Example:

lm(y ~ x, data = pop_df) %>% tidy(conf.int = TRUE)

#>	#	A tibble: 2	× 7					
#>		term	estimate	std.error	statistic	p.value	conf.low	conf.high
#>		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
#>	1	(Intercept)	2.53	0.422	6.00	3.38e- 8	1.69	3.37
#>	2	X	0.567	0.0793	7.15	1.59e-10	0.410	0.724

Confidence intervals

We construct (1-lpha)-level confidence intervals for eta_1

$$\hat{eta}_1 \pm t_{lpha/2, ext{df}}\;\hat{ ext{SE}}igl(\hat{eta}_1igr)$$

Example:

lm(y ~ x, data = pop_df) %>% tidy(conf.int = TRUE) #> # A tibble: 2 × 7 estimate std.error statistic p.value conf.low conf.high term #> <chr> <dbl> <dbl> <dbl> <db1> <dbl> <dbl> #> #> 1 (Intercept) 2.53 0.422 6.00 3.38e- 8 3.37 1.69 #> 2 x 0.567 0.0793 7.15 1.59e-10 0.410 0.724

Our 95% confidence interval is thus $0.567 \pm 1.98 imes 0.0793 = [0.410, \, 0.724]$

Confidence intervals

So we have a confidence interval for β_1 , *i.e.*, [0.410, 0.724].

What does it mean?

Confidence intervals

So we have a confidence interval for β_1 , *i.e.*, [0.410, 0.724].

What does it mean?

Informally: The confidence interval gives us a region (interval) in which we can place some trust (confidence) for containing the parameter.

Confidence intervals

So we have a confidence interval for β_1 , *i.e.*, [0.410, 0.724].

What does it mean?

Informally: The confidence interval gives us a region (interval) in which we can place some trust (confidence) for containing the parameter.

More formally: If repeatedly sample from our population and construct confidence intervals for each of these samples, $(1 - \alpha)$ percent of our intervals (*e.g.*, 95%) will contain the population parameter *somewhere in the interval*.
Confidence intervals

So we have a confidence interval for β_1 , *i.e.*, [0.410, 0.724].

What does it mean?

Informally: The confidence interval gives us a region (interval) in which we can place some trust (confidence) for containing the parameter.

More formally: If repeatedly sample from our population and construct confidence intervals for each of these samples, $(1 - \alpha)$ percent of our intervals (*e.g.*, 95%) will contain the population parameter *somewhere in the interval*.

Now back to our simulation...

Confidence intervals

We drew 10,000 samples (each of size n = 30) from our population and estimated our regression model for each of these simulations:

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

(repeated 10,000 times)

Now, let's estimate 95% confidence intervals for each of these intervals...

Confidence intervals



From our previous simulation: 97.8% of our 95% confidences intervals contain the true parameter value of β_1 .

That's a **probabilistic statement**:

- Could be more.
- Could be less.

Hypothesis testing

In many applications, we want to know more than a point estimate or a range of values. We want to know what our statistical evidence says about existing theories.

We want to test hypotheses posed by officials, politicians, economists, scientists, friends, weird neighbors, *etc.*

Examples

- Does increasing police presence **reduce crime**?
- Does building a giant wall **reduce crime**?
- Does shutting down a government **adversely affect the economy**?
- Does legal cannabis **reduce drunk driving** or **reduce opiod use**?
- Do air quality standards **increase health** and/or **reduce jobs**?

Hypothesis testing

Hypothesis testing relies upon very similar results and intuition.

While uncertainty certainly exists, we can still build *reliable* statistical tests (rejecting or failing to reject a posited hypothesis).

Hypothesis testing

Hypothesis testing relies upon very similar results and intuition.

While uncertainty certainly exists, we can still build *reliable* statistical tests (rejecting or failing to reject a posited hypothesis).

OLS *t* **test** Our (null) hypothesis states that eta_1 equals a value *c*, *i.e.*, $H_o:\ eta_1=c$

From OLS's properties, we can show that the test statistic

$$\hat{m{g}}_{\mathrm{stat}} = rac{\hat{m{eta}}_1 - c}{\hat{\mathrm{SE}}\left(\hat{m{eta}}_1
ight)}$$

follows the t distribution with n - k degrees of freedom.

Hypothesis testing

For an α -level, **two-sided** test, we reject the null hypothesis (and conclude with the alternative hypothesis) when

 $\left|t_{\mathrm{stat}}
ight|>\left|t_{1-lpha/2,\,df}
ight|$

meaning that our **test statistic is more extreme than the critical value**.

Alternatively, we can calculate the **p-value** that accompanies our test statistic, which effectively gives us the probability of seeing our test statistic *or a more extreme test statistic* if the null hypothesis were true.

Very small p-values (generally < 0.05) mean that it would be unlikely to see our results if the null hyopthesis were really true—we tend to reject the null for p-values below 0.05.

Hypothesis testing

R and statas default to testing hypotheses against the value zero.

Hypothesis testing

R and statas default to testing hypotheses against the value zero.

ln	n(y	/ ~ x, data =	= pop_df)	%>% tidy())	
#>	#	A tibble: 2	× 5			
#>		term	estimate	std.error	statistic	p.value
#>		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
#>	1	(Intercept)	2.53	0.422	6.00	3.38e- 8
#>	2	x	0.567	0.0793	7.15	1.59e-10

 $ext{H}_{ ext{o}}$: $eta_1=0$ vs. $ext{H}_{ ext{a}}$: $eta_1
eq 0$

Hypothesis testing

R and statas default to testing hypotheses against the value zero.

lr	n(y	y ~ x, data =	= pop_df)	%>% tidy())	
#>	#	A tibble: 2	× 5			
#>		term	estimate	std.error	statistic	p.value
#>		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
#>	1	(Intercept)	2.53	0.422	6.00	3.38e- 8
#>	2	Х	0.567	0.0793	7.15	1.59e-10

H_o: $eta_1=0$ vs. H_a: $eta_1
eq 0$

 $t_{
m stat}=7.15$ and $t_{0.975,\,28}=2.05$

Hypothesis testing

R and statas default to testing hypotheses against the value zero.

lr	n(y	y ~ x, data =	= pop_df)	%>% tidy())	
#>	#	A tibble: 2	× 5			
#>		term	estimate	std.error	statistic	p.value
#>		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
#>	1	(Intercept)	2.53	0.422	6.00	3.38e- 8
#>	2	x	0.567	0.0793	7.15	1.59e-10

 $ext{H}_{ ext{o}}:eta_{1}=0$ vs. $ext{H}_{ ext{a}}:eta_{1}
eq 0$

 $t_{
m stat} = 7.15$ and $t_{0.975,~28} = 2.05$ which implies p-value < 0.05

Hypothesis testing

R and statas default to testing hypotheses against the value zero.

lr	n(y	/ ~ x, data =	= pop_df)	%>% tidy())	
#>	#	A tibble: 2	× 5			_
#>		term	estimate	std.error	statistic	p.value
#>		<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
#>	1	(Intercept)	2.53	0.422	6.00	3.38e- 8
#>	2	X	0.567	0.0793	7.15	1.59e-10

```
	ext{H}_{	ext{o}}: eta_1=0 vs. 	ext{H}_{	ext{a}}: eta_1
eq 0
```

 $t_{
m stat} = 7.15$ and $t_{0.975,\,28} = 2.05$ which implies p-value < 0.05

Therefore, we **reject H**_o.

*F*tests

You will sometimes see F tests in econometrics.

We use F tests to test hypotheses that involve multiple parameters (e.g., $\beta_1 = \beta_2$ or $\beta_3 + \beta_4 = 1$),

rather than a single simple hypothesis (e.g., $\beta_1 = 0$, for which we would just use a t test).

*F*tests

Example

Economists love to say "Money is fungible."

Imagine that we might want to test whether money received as income actually has the same effect on consumption as money received from tax rebates/returns.

 $ext{Consumption}_i = eta_0 + eta_1 ext{Income}_i + eta_2 ext{Rebate}_i + u_i$

Ftests

Example, continued

We can write our null hypothesis as

$$H_o:\ eta_1=eta_2 \iff H_o:\ eta_1-eta_2=0$$

Imposing this null hypothesis gives us the **restricted model**

 $egin{aligned} ext{Consumption}_i &= eta_0 + eta_1 ext{Income}_i + eta_1 ext{Rebate}_i + u_i \ & ext{Consumption}_i &= eta_0 + eta_1 \left(ext{Income}_i + ext{Rebate}_i
ight) + u_i \end{aligned}$

Ftests

Example, continued

To this the null hypothesis $H_o: \ \beta_1 = \beta_2$ against $H_a: \ \beta_1 \neq \beta_2$, we use the F statistic

$$F_{q,\,n-k-1} = rac{\left(\mathrm{SSE}_r - \mathrm{SSE}_u
ight)/q}{\mathrm{SSE}_u/(n-k-1)}$$

which (as its name suggests) follows the F distribution with q numerator degrees of freedom and n - k - 1 denominator degrees of freedom.

Here, q is the number of restrictions we impose via H_o .

Ftests

Example, continued

The term SSE_r is the sum of squared errors (SSE) from our restricted model

 $ext{Consumption}_i = eta_0 + eta_1 \left(ext{Income}_i + ext{Rebate}_i
ight) + u_i$

and SSE_u is the sum of squared errors (SSE) from our **unrestricted model**

 $\text{Consumption}_i = eta_0 + eta_1 \text{Income}_i + eta_2 \text{Rebate}_i + u_i$





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