

holding other  
 ↑ factor unchanged!

3.2 How could the multiple regression analysis enable *ceteris paribus* analysis?

- Consider a multiple regression function of *wage*

$$wage = \beta_0 + \beta_1 educ + \beta_2 inc + u \quad (4.4)$$

to know impact  
 only 1 variable

Here,

- $\beta_0$  is the intercept.
- $\beta_1$  measures the change in *wage* with respect to *educ*, holding other factors fixed.
- $\beta_2$  measures the change in *wage* with respect to *inc*, holding other factors fixed.
- What if the function of *wage* is, instead written as

$$wage = \beta_0 + \beta_1 educ + \beta_2 inc + \beta_3 educ^2 + u$$

Then,

$\beta_0$  is the intercept.

The change in *wage* with respect to *educ* (holding other factors fixed) is measured by:

$$\frac{\partial wage}{\partial education} = \beta_1 + 2\beta_3 educ$$

The change in *wage* with respect to *inc* (holding other factors fixed) is measured by:

$$\frac{\partial wage}{\partial income} = \beta_2$$

TBC :-)

## A story to Ceteris Paribus analysis

20/02/2020

- A student in the EE489 class wants to find the impact of parking space on Condo's price.

So, he proposes  $\text{price} = \beta_0 + \beta_1 \cdot \text{parking} + u$

⇒ what do you think should be the sign of  $\beta_1$ ? (+) or (-) Why?

(+) because better quality condominiums usually allocate more space for parking

- But then, the student said, this may not be the case because condos that are closer to BTS & MRT usually have lower % parking / room but their prices are higher!

⇒ This student said he expects the sign of  $\beta_1$  to be (-) !!!

⇒ Ajarn said, definitely. If you estimate a simple regression (1) having % parking as the only X variable, your analysis will not be able to analyze the impact of % parking, holding

distance in MRT, BTS **other factors** constant

- so, to make  $\beta_1$  measure only the impact of % parking, you should include distance to the nearest BTS, MRT in the regression.

We can expect  $\beta_1$  to be (+) in this equation

$$\text{Price}_i = \beta_0 + \beta_1 \cdot \% \text{ parking} + \beta_2 \text{ dist. BTS} + \beta_3 \text{ dis. MRT}$$

By including BTS, MRT into the equation, we can conduct partial effects analysis of % parking, holding BTS, MRT constant.

4 Expected Value of the OLS Estimators in multiple regression

- Under assumptions MLR 1 to 4 (see Wooldridge),  $\hat{\beta}_{OLS}$  are unbiased.
- 2 issues should be considered regarding the biasedness of  $\hat{\beta}_{OLS}$

ex) distance from condo to Japan to explain condo price

4.1 Issue #1: Including Irrelevant Variable (Overspecifying the Model) Too many variable than necessary

- Suppose we specify the model

if include irrelevant  
var equation biased?

assume that  
 $X_2$  is irrelevant

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u \tag{4.5}$$

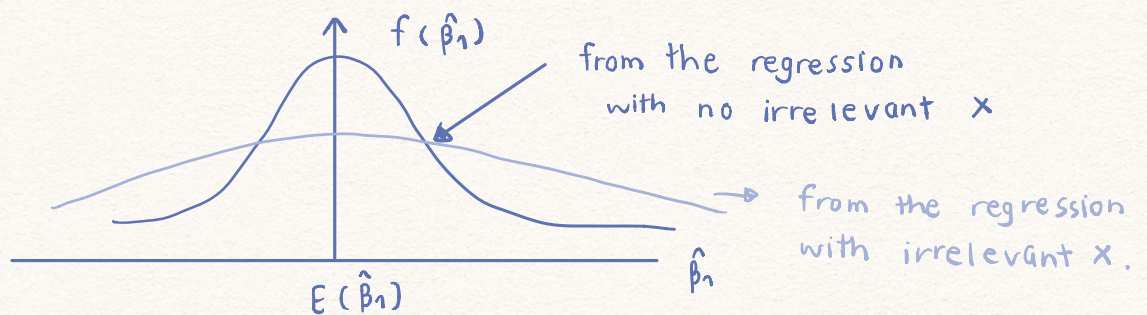
and this model satisfies the multiple regression assumptions 1 to 4

- If the true value of  $\beta_2$  is "0"  $\Rightarrow$  no need to put  $X_2$  in the multiple regression (Because  $X_2$  does not explain  $Y$ )
- If we estimate a model with  $X_2$  anyway,

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$$

$\Rightarrow$  The estimated OLS parameters  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$  will still be unbiased because we do not violate assumption MLR 1-4

- However, including irrelevant variable ( $X_2$  in this case) would make the variance of  $\hat{\beta}_0, \hat{\beta}_1$  become unnecessarily LARGE. so,  $\hat{\beta}_0, \hat{\beta}_1$  not be the most efficient.



4.2 Issue #2: Excluding Relevant Variable (Underspecifying the Model → omitted variable bias. This is a serious problem!)

omitted  $\beta_2, \beta_3$   
 omitted  $\beta_1$   
 biased

- Suppose we the **TRUE** model is actually

so, we omitted the factor that explains  $y$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u \rightarrow \text{TRUE MODEL}$$

where none of the  $\beta$  is zero and this model satisfies the multiple regression assumptions 1 to 4.

- But we omit variable  $X_2$  and estimate the following equation using OLS

$$Y = \beta_0 + \beta_1 X_1 + v \rightarrow (v = \beta_2 X_2 + u)$$

and we estimate the wrong model :

$$\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 x_1$$

let's check if  $\tilde{\beta}_1$  will be biased by checking if we have violated assumption MLR 1-4

$$\Rightarrow E(v) = E(\beta_2 X_2 + u) = \beta_2 E(X_2) + E(u)$$

$$= \beta_2 E(X_2) \neq 0 \Rightarrow E(v) \text{ in the wrong model } \neq 0$$

MRL 1- linear parameter

MRL-2 random sampling

MRL-3 multi

MRL-4 zero conditional

$\Rightarrow$  Since  $X_2$  is in the error term of the wrong model ( $v$ ),  $\text{Cov}(X_1, X_2)$  has to be "0" for  $\tilde{\beta}_1$  to be unbiased (or for  $\text{Cov}(X_1, v) = 0$ )

☺ some proof omitted see textbook ☺  $E[\tilde{\beta}_1] = \beta_1$  (unbiased)

result from doing OLS of the wrong model

we get  $E[\tilde{\beta}_1] = \beta_1 + \beta_2 E\left(\frac{\sum_i x_{1i} x_{2i}}{\sum_i x_{1i}^2}\right)$ ;  $x_{1i} = x_{1i} - \bar{x}$   
 $x_{2i} = x_{2i} - \bar{x}$

$$= \beta_1 + \beta_2 \frac{\text{Cov}(X_1, X_2)}{\text{Var}(X_1)} \rightarrow \text{Biased}$$

so,  $E(\tilde{\beta}_1) \neq \beta_1 \Rightarrow$  it will be biased by

the size of  $\beta_2 \frac{\text{Cov}(X_1, X_2)}{\text{Var}(X_1)}$

### Direction of Bias

If  $\beta_2 > 0$ , then  $\text{Cov}(X_1, X_2) > 0 \Rightarrow \tilde{\beta}_1$  will be upward biased

$\text{Cov}(X_1, X_2) < 0 \Rightarrow \tilde{\beta}_1$  will be downward biased

$\beta_2 < 0$ , then if  $\text{Cov}(X_1, X_2) > 0 \Rightarrow \hat{\beta}_1$  will be downward biased

$\text{Cov}(X_1, X_2) < 0 \Rightarrow$  upward biased

### Example

Suppose the real population regression is

$$\text{price} = \beta_0 + \beta_1 \% \text{ parking} + \beta_2 \text{ dist\_BTS\_MRT} + u$$

$\% \text{ parking}$   
↑  
to # rooms

$\text{dist\_BTS\_MRT}$   
↑  
distance to the nearest  
BTS or MRT station

But we estimate

$$\text{price} = \beta_0 + \beta_1 \% \text{ parking}_i + v_i$$

$$\text{Then, } E(\tilde{\beta}_1) = \beta_1 + \beta_2 \frac{\text{Cov}(\% \text{ parking}, \text{distance\_BTS\_MRT})}{\text{var}(\% \text{ parking})}$$

(+)

(-)

સિંગલ નેગેટિવ બબવા

var (% parking)

(+)

★  $\tilde{\beta}_1$  is going to be smaller

than expected if we omitted

the dis\_MRT-BTS variable!

(-)

we have a downward  
biased here!

— more ex in Gujarati —

## 5 Variance of the OLS Estimators

- The  $\hat{\beta}_{OLS}$  would be the most efficient among the linear unbiased estimators if assumption 5 is satisfied

- Multiple Linear Regression (MLR) assumption 5: Homoskedasticity

The error term  $u$  has the same variance given any values of the explanatory variables.

$$\text{Var}(u|X_1, X_2, \dots, X_k) = \sigma^2 \leftarrow \text{constant regardless of the value of } x$$

- Example:

$$\text{salary} = \beta_0 + \beta_1 \text{profits} + \beta_2 \text{sales} + u$$

homoskedasticity require that  $\text{var}(u/\text{profit}, \text{sales}) = \sigma^2 \Rightarrow \text{MLR 5}$

- If the MLR assumption 5 is true, then  $\rightarrow$  with some prove (in the textbook)

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_i (x_{ij} - \bar{x}_j)^2 (1 - R_j^2)} = \frac{\sigma^2}{\text{SST}_j (1 - R_j^2)}$$

where  $j \in 1, 2, \dots, k^{\text{th}}$  explanatory variable)

- $R_j^2$  is the  $R^2$  from regression  $x_j$  on all other  $X_s$  (except  $x_j$  itself)

i.e.  $R_j^2$  is the  $R^2$  from

$$x_j = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k$$

not including  $x_j$

$$\text{or } x_2 = \alpha_0 + \alpha_1 x_1 + \alpha_3 x_3 + \dots + \alpha_k x_k$$

### 6 Estimator of the OLS Variance

- Since we don't know what  $\sigma^2$  is (population concept), we need to find an estimator of it.

$$\hat{\sigma}^2 = \frac{\sum_i \hat{u}_i^2}{n - \text{number of estimated parameter}}$$

$$= \frac{\sum_i \hat{u}_i^2}{n - (k+1)}$$

+1 intercept k slopes  
 $\beta_0 + \beta_1 + \dots + \beta_k$

Simple regress  $n-2$

so, variance of  $(\hat{\beta}_j) = \frac{\hat{\sigma}^2}{SST_j(1-R_j^2)}$

If simple regression,  $\text{Var}(\hat{\beta}_1) = \frac{\hat{\sigma}^2}{SST_x}$

Practice question on moodle in Cengage

- Thus, STATA's calculation of the std.err. of  $\hat{\beta}_j$  is

$$\widehat{\text{std.deviation}}.\hat{\beta}_j = \text{std.err.}\hat{\beta}_j = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^N (X_{ij} - \bar{X}_j)^2 (1 - R_j^2)}}$$

Comments:

\*  $A_j$ . will come back to this after the midterm exam

# Stata Lab 1 – Introduction

## 1 What is STATA?

- A statistical software package used mostly in economics, sociology, political science and epidemiology.
- Stata can be used to manage database, run regressions, generate graphics, do simulations, etc.
- The user should have their own dataset. The Stata data file is usually saved in the .dta format.
- Data of any other formats (like excel) can be imported and/or converted into .dta format.

### 1.1 STATA supports

- Stata's own website: <http://www.stata.com/support/faqs/>
- Stata program's help function: For example, suppose you would like to know more about the "regress" command, then... Open the stata program > in the "command" box > type "help regress" without the " "> press enter.
- Stata's official manual (can be found in the library and embeded in the program)
- Other Stata's user's manual: My favorite one is "An Introduction to Modern Econometrics Using Stata" by Christopher F. Baum.
- or... simply type your question(s) into a search engine.

### 1.2 Data files and Do-files

- The Stata's data file keeps all the data points. For example, each  $Y_i$  and the corresponding  $X_i, \forall i = 1, 2, 3, \dots, n$ .
- The do-file records all the commands that you use to analyze the data.
- Once the data is cleaned, it is best not to keep on re-saving the original data file. If several steps have to be done before analyzing the data (like running a regression), do it on the do-file.

## 2 Tutorial 1: Exploring the Data and Running a Simple Regression

- **Download Wooldridge datasets**

1. Download Wooldridge's data – go to thomsonedu's website (or go to your BE Moodle: EE325):

"[http://www.thomsonedu.com/aise/economics/wooldridge\\_2e\\_datasets/](http://www.thomsonedu.com/aise/economics/wooldridge_2e_datasets/)".

- **To open the STATA program**

1. Double click on the STATA icon.
2. Click on the "Do-file" icon on the top panel of the Stata program. "Save As" your Do-file (and name it "EE325") on your computer.

- **Open the data file using Do-file**

1. file -> open -> then, direct the program to the file "CEOSAL2.DTA".
2. On the command window, you will see a command to open this file. If you would like to open the file from your do-file in the future, you can use this command.

- **To explore and understand the data**

1. type: browse
2. type: describe
3. type: summarize
4. type: sum
5. type: codebook
6. type: describe salary
7. type: tabulate college
8. type: tab college
9. to use the "if" command to find conditional mean (average) type: sum if grad == 1
10. type: sum salary if age <= 40
11. type: correlate salary sales profits
12. type: correlate salary sales profits, covariance
13. type: plot salary profits
14. type: twoway scatter salary profits

- **To run a simple (OLS) regression (one explanatory variable)**

1. type: regress salary profits
  - **To create the fitted value ( $\hat{Y}_i$ ) and the residual ( $\hat{u}_i$ )**
    1. type: predict y\_hat, xb
    2. type: predict u\_hat, residual
  - **To see how well we do at finding a Sample Linear Function**
    1. type: twoway scatter salary profits || line y\_hat profits
    2. type: twoway scatter u\_hat profits
    3. To check if the OLS estimation makes  $X_i$  uncorrelated with  $\hat{u}_i$  (by the OLS calculation, they should not correlate), type: correlate sales u\_hat
  - **To execute mathematical operations**
    1. type: generate log\_salary = log(salary)
    2. type: gen log\_profit = log(profit)
    3. type: gen profit\_2 = 3+5\*profit
    4. type: regress salary profits profit\_2
    5. type: gen profit\_sq = profit^2
    6. type: regress salary profits profit\_sq
  - **To perform a multiple regression analysis**
    1. type: regress salary profits sales
    2. type: regress salary profits sales ceoten
  - **To exit your Stata**
    1. Save your do-file
    2. file -> exit -> don't save (never ever modify your master dataset!)
  - **To find out what all the above commands mean** – (type in the command box) help summarize, help predict, help twoway, etc etc.