

# Lecture 9

## Multiple Regression IV

# Outline

- Dummy Variable Trap
- Level- Level model
- Log-Level model
- Level-Log model
- Log- Log model
- Interaction effects

# The notion of control

- Meaning of control #1:
  - When you control for variables you hold that variable constant or fixed
  - Similar notion to a lab experiment—keep other factors held fixed (*ceteris paribus*)
- Meaning of control #2:
  - e.g.  $Health = \alpha + \beta Insured + \gamma_1 Age + \gamma_2 Educ + \gamma_3 Income + \varepsilon$
  - Within group comparison: compare mean health status of the insured with the uninsured for those with same age, education and income. It's more proper than across group comparison
- Meaning of control #3:
  - “Partialling Out”
  - $\beta$  tells us the causal effect of whether insured on health after the confounding effect of age, education and income has been partialled (or netted) out

# Good, bad and the useless

- Good controls are variables that you include in the regression in order to kill dirty identifying variation
- Useless controls are variables that you include in the regression that aren't correlated with treatment variable D.
  - Useless controls can still explain Y
  - Often it's worthwhile to include these controls to reduce your standard errors, but fundamentally, these controls don't do that much
- Bad controls are variables that are caused by D. Controlling for these variables is called "over controlling"

## Example: Immigration

- Suppose immigration positively impacts crime rate
- If I regress crime on immigration, I may find a positive effect
- $Crime = \alpha + \beta \text{ Number of Immigrants} + \varepsilon, \hat{\beta} > 0$
- Is this relationship biased?
  - Possibly, because some other factors (e.g. big city indicator) is left in the error term. Once controlling for big city, the effect may change
- Note, when we include a control, we are more concerned whether it affects treatment **variable D**, not whether it affects Y.

## Example: good, bad and useless controls

- Variable “big city” affects both immigration and crime rate, therefore it is a **good** control.
- Now you include another variable number of rainy days. You think number of rainy days will affect crime. But it has nothing to do with immigration. Therefore it is a **useless** control. (Although the adjusted R-square will increase, it doesn't eliminate any “dirty” variation)
- Now consider the variable-local unemployment rate. This variable is likely to be caused by immigrants and therefore may be a **bad** control. Because it will take away the “clean” variation too. You will over-control by including it
  - Note: when you try to think of a bad control, think about the mechanisms why immigrants affect crime rate, controlling for these mechanisms are usually over-controlling

# Criteria for good, bad and useless controls

- Good controls affect Y and D (variable of interest)
- Useless controls affect Y but not D
- Bad controls affects Y, and are caused by D (Bad controls are often mechanisms through which D affect Y)

## Example: code categorical variables into dummies

- Sometimes, you want to turn a categorical variable into dummy variables to model the relationship more flexibly
- e.g. Education: 0 if HS dropout, 1 if HS grad, and 2 if college grad.

Recode it into three dummy variables:

- less\_HS:1 if education=0, 0 otherwise
- HS\_grad:1 if education=1, 0 otherwise
- college\_or\_above:1 if education=2, 0 otherwise



## Dummy variable trap

- If we run  $\text{Health} = \beta_0 + \beta_1 \text{Less\_HS} + \beta_2 \text{HS\_grad} + \beta_3 \text{College\_or\_above} + \varepsilon$ 
  - The regression will not work
  - This is called perfect multi-collinearity

Constant	less_HS	HS_grad	college	health status
1	1	0	0	3
1	0	1	0	4
1	0	1	0	4
1	0	0	1	3
1	0	0	1	2
1	1	0	0	1
1	0	1	0	4

**Less\_HS+HS\_grad+college\_or\_above  
=constant  
Cannot estimate the regression!**

## How to get out of dummy variable trap(1)?

- There are two ways of getting out of the dummy variable trap
- Way #1: omit one category of the dummy variables
- e.g. if we run a regression:
- $\text{Health} = \beta_0 + \beta_1 \text{HS\_grad} + \beta_2 \text{College\_or\_above} + \varepsilon$
- We omit HS dropouts dummy and treat it as the baseline. The categories we include are compared to the category we exclude
- How do we interpret  $\beta_1$ ?
  - The average health status of HS grads relative to HS dropouts.
  - **It is the difference in average health between HS grad and HS dropouts**
- Exercise: Interpret  $\beta_0, \beta_2$

## How to get out of dummy variable trap(2)?

- Way#2: omit the constant term
- $\text{Health} = \beta_1 \text{LessHS} + \beta_2 \text{HS\_grad} + \beta_3 \text{College\_or\_above} + \varepsilon$ 
  - How does this regression differ from the previous?
- How do we interpret  $\beta_1$ ?
  - The average health status of HS dropouts.
- Exercise: Interpret  $\beta_2$  and  $\beta_3$

## Level-Level

- A “Level-level” regression specification.
- $y = \beta_0 + \beta_1 x + u$
- This is called a “level-level” specification because raw values (levels) of  $y$  are being regressed on raw values of  $x$ .
- How do we interpret  $\beta_1$ ?
- We interpret it as the increase in  $y$  when 1 unit of  $x$  increases.

## Log-Level

- A “Log-level” regression specification
- $\log(y) = \beta_0 + \beta_1 x + u$
- What is  $\log(y)$ ? Log indicates logarithmic function.
- This is called a “log-level” specification because the natural log transformed values of  $y$  are being regressed on raw values of  $x$ .
- You might want to run this specification if you think that increases in  $x$  lead to a constant *percentage* increase in  $y$ . (e.g. wage on education)

## Log-Level

- A “Log-level” regression specification
- $\log(y) = \beta_0 + \beta_1 x + u$
- How do we interpret  $\beta_1$ ?
- One unit change in  $x$  leads to  $100 * \beta_1$  percent change in  $Y$

## Log-Level

- A “Log-level” regression specification
- $\log(y) = \beta_0 + \beta_1 x + u$
- What is  $\log(y)$ ? Log indicates logarithmic function.
- This is called a “log-level” specification because the natural log transformed values of  $y$  are being regressed on raw values of  $x$ .
- You might want to run this specification if you think that increases in  $x$  lead to a constant *percentage* increase in  $y$ . (e.g. wage on education)

## Level-Log

- A “Level-log” regression specification
- $y = \beta_0 + \beta_1 \log(x) + u$
- This is called a “level-log” specification because y is being regressed on natural log transformed values of x.
- You might want to run this specification if you think that *percentage* increases in x lead to a constant increase in y.



## Level-Log

- A “Level-log” regression specification
- $y = \beta_0 + \beta_1 \log(x) + u$
- How do we interpret  $\beta_1$ ?
- One **percent** change in  $x$  leads to  $\beta_1/100$  **unit** change in  $Y$

## Log-Log

- A “Log-log” regression specification
- $\log(y) = \beta_0 + \beta_1 \log(x) + u$
- This is called a “log-log” specification because natural log transformed values of  $y$  are being regressed on natural log transformed values of  $x$ .
- You might want to run this specification if you think that *percentage* increases in  $x$  lead to a constant *percentage* changes in  $y$ .

## Log-Log

- A “Log-log” regression specification
- $\log(y) = \beta_0 + \beta_1 \log(x) + u$
- How do we interpret  $\beta_1$ ?
- One percent change in x leads to  $\beta_1$  percent change in Y

Model	Equation	Interpretation
Level-Level Regression	$Y = \alpha + \beta X$	One unit change in $X$ leads to $\beta$ unit change in $Y$
Log-Linear Regression	$\log(Y) = \alpha + \beta X$	One unit change in $X$ leads to $100 * \beta$ percent change in $Y$
Linear-Log Regression	$Y = \alpha + \beta \log(X)$	One percent change in $X$ leads to $\beta/100$ unit change in $Y$
Log-Log Regression	$\log(Y) = \alpha + \beta \log(X)$	One percent change in $X$ leads to $\beta$ percent change in $Y$

## Examples

$$\log(\widehat{earnings}) = 2.805 + 0.0087Age$$

Interpret the coefficient of age.

## Examples

$$\log(\widehat{earnings}) = 2.805 + 0.0087Age$$

Interpret the coefficient of age.

Earnings are predicted to increase by 0.87 [ $0.0087 \times 100$ ]% for each additional year of age.

## Examples

$$\widehat{\text{test score}} = 557.8 + 36.42 \ln(\text{Income})$$

Interpret the coefficient of log(income).

## Examples

$$\widehat{\text{test score}} = 557.8 + 36.42\ln(\text{Income})$$

Interpret the coefficient of log(income).

A 1% increase in income is associated with an increase in test scores of  $36.42/100=0.36$  point.



## Examples

$$\widehat{\text{test score}} = 6.336 + 0.0554 \ln(\text{Income})$$

Interpret the coefficient of log(income).

## Examples

$$\widehat{\text{test score}} = 6.336 + 0.0554 \ln(\text{Income})$$

Interpret the coefficient of  $\log(\text{income})$ .

A 1% increase in income is associated with an increase in test scores of 0.0554 percent.

## Percent vs. percentage point

- Most of us are comfortable with percentage increases and decreases.
- A hundred dollar skateboard goes on sale for 75 dollars and we can calculate easily enough that this is a 25 percent discount. The key feature of percentage change is that it provides a measure of change that is proportional to the original quantity (100 dollars in this case).

## Percent vs. percentage point

- Unfortunately, this simple setup can become confusing when the original quantity is itself expressed as a percentage.
- For example, I heard on the news that between 1995 and 2005, the percentage of Americans without health insurance rose from 60 percent to 69 percent.
- It is tempting to call this a 9 percent increase, but this understates the size of the increase. Sixty nine percent is actually a 15 percent increase over the original sixty percent. Try it. If we start with 60, and add to it 15 percent of 60, we get 69.
- To clarify this state of affairs, we say that the percentage of uninsured Americans rose by 15 percent. Alternatively, we may say that the percentage of uninsured Americans rose by 9 *percentage points*.

## Example

- Suppose you pick peaches and are paid 4 dollars per bushel. One day your boss announces that he is giving you a raise. You will now be paid 5 dollars per bushel.
- Question: what is the percentage increase in your wage?

## Example

- Suppose you have a student loan with an annual interest rate of 4 percent. One day your lender announces that the interest rate will soon increase to an annual rate of 5 percent.
- Question: what is the percentage increase in your interest rate?
- Question: what is the percentage point increase in your interest rate?

## Example

- Suppose you have a model:  $y = \beta_0 + \beta_1 x + u$
- $y$  is measured in %. For example  $y$  is cancer rate.
- Interpretation: when  $x$  increases by 1 unit,  $y$  increases by  $\beta_1$  *percentage point*.

## Interactions between independent variables

- Example: gender -> earnings
- Perhaps gender effects on earnings can be different depending on marital status
- Perhaps married women are more penalized in the labor market
- We assumed “constant” effect of being a female so far
- We want to allow for different effect of being a female depending on the marital status
- How can we do this?



## Interactions between independent variables

- $y = \beta_0 + \beta_1 \text{Female} + \beta_2 \text{Married} + \beta_3 \text{Female} * \text{Married} + u$
- $E(y | \text{Female} = 0, \text{married} = 0) = \beta_0$
- $E(y | \text{Female} = 1, \text{married} = 0) = \beta_0 + \beta_1$
- $E(y | \text{Female} = 0, \text{married} = 1) = \beta_0 + \beta_2$
- $E(y | \text{Female} = 1, \text{married} = 1) = \beta_0 + \beta_1 + \beta_2 + \beta_3$
- The effect of being a female for married people:
- The effect of being a female for non-married people:

## Interactions between independent variables

- $y = \beta_0 + \beta_1 \text{Female} + \beta_2 \text{Married} + \beta_3 \text{Female} * \text{Married} + u$
- $E(y | \text{Female} = 0, \text{married} = 0) = \beta_0$
- $E(y | \text{Female} = 1, \text{married} = 0) = \beta_0 + \beta_1$
- $E(y | \text{Female} = 0, \text{married} = 1) = \beta_0 + \beta_2$
- $E(y | \text{Female} = 1, \text{married} = 1) = \beta_0 + \beta_1 + \beta_2 + \beta_3$
  
- The effect of being a female for married people:  $\beta_1 + \beta_3$
- The effect of being a female for non-married people:  $\beta_1$
- The effect of being a female depends on marital status

## Interactions between independent variables

- Example: education -> earnings
- Perhaps the effects of education on earnings can be different for different gender
- We assumed “constant” effect of education so far
- We want to allow for different effect of being education depending on gender
- How can we do this?

## Interactions between independent variables

- $y = \beta_0 + \beta_1 Educ + \beta_2 Female + \beta_3 Educ * Female + u$
- $E(y | Female = 0) = \beta_0 + \beta_1 Educ$
- $E(y | Female = 1) = \beta_0 + \beta_2 + (\beta_1 + \beta_3) Educ$
- The effect of education can be different for males and females
- The intercept difference is  $\beta_2$
- Now we also have slope difference, which is  $\beta_3$

# Review for quizzes

- Know the dummy variable trap and how to get rid of it.
- Know how to interpret coefficients in different models (level- level; log-level; level-log, and log-log models)
- Know how to interpret interaction terms.