

LECTURE 6

Introduction to Econometrics

Omitted variables, Multicollinearity & Heteroskedasticity

December 1st, 2023

ON PREVIOUS LECTURES

- ▶ We discussed the specification of a regression equation
- ▶ **Specification** consists of choosing:
 1. correct independent variables
 2. correct functional form
 3. correct form of the stochastic error term

ON TODAY'S LECTURE

- ▶ We will talk about omitted variables
- ▶ We will finish the discussion of the choice of independent variables by talking about **multicollinearity**
- ▶ We will start the discussion of the correct form of the error term by talking about **heteroskedasticity**
- ▶ For both of these issues, we will learn
 - what is the nature of the problem
 - what are its consequences
 - how it is diagnosed
 - what are the remedies available

OMITTED VARIABLES

- We omit a variable when we

forget to include it

do not have data for it

- This misspecification results in

not having the coefficient for this variable

biasing estimated coefficients of other variables in the equation \rightarrow **omitted variable bias**

OMITTED VARIABLES

- Where does the omitted variable bias come from?
- True model:

$$y_i = \beta x_i + \gamma z_i + u_i$$

- Model as it looks when we omit variable z :

$$y_i = \beta x_i + \tilde{u}_i$$

implying

$$\tilde{u}_i = \gamma z_i + u_i$$

- We need that $Cov(\tilde{u}_i, x_i) = 0$, but:

$$Cov(\tilde{u}_i, x_i) = Cov(\gamma z_i + u_i, x_i) = \gamma Cov(z_i, x_i) \neq 0$$

- The classical assumption is violated \Rightarrow biased (and inconsistent) estimate!!!

OMITTED VARIABLES

- For the model with omitted variable:

$$E(\widehat{\beta}^{\text{omitted model}}) = \beta + \text{bias}$$

$$\text{bias} = \gamma * \alpha$$

Coefficients β and γ are from the true model

$$y_i = \beta x_i + \gamma z_i + u_i$$

Coefficient α is from a regression of z on x , i.e.

$$z_i = \alpha x_i + e_i$$

- **The bias is zero if $\gamma = 0$ or $\alpha = 0$ (not likely to happen)**

OMITTED VARIABLES

- o **Intuitive explanation:**

if we leave out an important variable from the regression ($\gamma \neq 0$), coefficients of other variables are biased unless the omitted variable is uncorrelated with all included dependent variables ($\alpha \neq 0$)

the included variables pick up some of the effect of the omitted variable (if they are correlated), and the coefficients of included variables thus change causing the bias

- o Example: what would happen if you estimated a production function with capital only and omitted labor?

OMITTED VARIABLES

- Example: estimating the price of chicken meat in the US

$$\hat{Y}_t = 31.5 - \frac{0.73}{(0.08)} PC_t + \frac{0.11}{(0.05)} PB_t + \frac{0.23}{(0.02)} YD_t$$

$$R^2 = 0.986 \quad , \quad n = 44$$

Y_t ... per capita chicken consumption
 PC_t ... price of chicken
 PB_t ... price of beef
 YD_t ... per capita disposable income

OMITTED VARIABLES

- When we omit price of beef:

$$\hat{Y}_t = 32.9 - \frac{0.70}{(0.08)} PC_t + \frac{0.27}{(0.01)} YD_t$$

$$R^2 = 0.895 \quad , \quad n = 44$$

- Compare to the true model:

$$\hat{Y}_t = 31.5 - \frac{0.73}{(0.08)} PC_t + \frac{0.11}{(0.05)} PB_t + \frac{0.23}{(0.02)} YD_t$$

$$R^2 = 0.986 \quad , \quad n = 44$$

- We observe positive bias in the coefficient of PC (was it expected?)

OMITTED VARIABLES

- **Determining the direction of bias: $\text{bias} = \gamma * \alpha$**

Where γ is a correlation between the omitted variable and the dependent variable (the price of beef and chicken consumption)

γ is likely to be positive

Where α is a correlation between the omitted variable and the included independent variable (the price of beef and the price of chicken)

α is likely to be positive

- **Conclusion: Bias in the coefficient of the price of chicken is likely to be positive if we omit the price of beef from the equation.**

OMITTED VARIABLES

- In reality, we usually do not have the true model to compare with

Because we do not know what the true model is

Because we do not have data for some important variable

- We can often recognize the bias if we obtain some unexpected results
- We can prevent omitting variables by relying on the theory
- If we cannot prevent omitting variables, we can at least determine in what way this biases our estimates

IRRELEVANT VARIABLES

- A second type of specification error is including a variable that does not belong to the model
- This misspecification
does not cause bias
but it increases the variances of the estimated coefficients of
the included variables

IRRELEVANT VARIABLES

- True model:

$$y_i = \beta x_i + u_i \quad (1)$$

- **Model as it looks when we add irrelevantz:**

$$y_i = \beta x_i + \gamma z_i + \tilde{u}_i \quad (2)$$

- **We can represent the error term as $\tilde{u}_i = u_i - \gamma z_i$**
- **but since from the true model $\gamma = 0$, we have $\tilde{u}_i = u_i$ and there is no bias**

SUMMARY OF THE THEORY

- Bias - efficiency trade-off:

	Omitted variable	Irrelevant variable
Bias	Yes*	No
Variance	Decreases *	Increases*

* As long as we have correlation between x and z

PERFECT MULTICOLLINEARITY

- ▶ Some explanatory variable is a perfect linear function of one or more other explanatory variables
- ▶ Violation of one of the classical assumptions
- ▶ OLS estimate cannot be found

Intuitively: the estimator cannot distinguish which of the explanatory variables causes the change of the dependent variable if they move together

Technically: the matrix $\mathbf{X}'\mathbf{X}$ is singular (not invertible)

- ▶ Rare and easy to detect

EXAMPLES OF PERFECT MULTICOLLINEARITY

Dummy variable trap

- ▶ Inclusion of dummy variable for each category in the model with intercept
- ▶ Example: wage equation for sample of individuals who have high-school education or higher:

$$wage_i = \beta_1 + \beta_2 high_school_i + \beta_3 university_i + \beta_4 phd_i + e_i$$

- ▶ Automatically detected by most statistical softwares

IMPERFECT MULTICOLLINEARITY

- ▶ Two or more explanatory variables are highly correlated in the particular data set
- ▶ OLS estimate can be found, but it may be very imprecise

Intuitively: the estimator can hardly distinguish the effects of the explanatory variables if they are highly correlated

Technically: the matrix $\mathbf{X}'\mathbf{X}$ is nearly singular and this causes the variance of the estimator $Var(\hat{\beta}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$ to be very large

- ▶ Usually referred to simply as “multicollinearity”

CONSEQUENCES OF MULTICOLLINEARITY

1. Estimates remain unbiased and consistent (estimated coefficients are not affected)
2. Standard errors of coefficients increase

Confidence intervals are very large - estimates are less reliable

t -statistics are smaller - variables may become insignificant

DETECTION OF MULTICOLLINEARITY

- ▶ Some multicollinearity exists in every equation - the aim is to recognize when it causes a severe problem
- ▶ Multicollinearity can be signaled by the underlying theory, but it is very sample depending
- ▶ We judge the severity of multicollinearity based on the properties of our sample and on the results we obtain
- ▶ One simple method: examine correlation coefficients between explanatory variables

if some of them is too high, we may suspect that the coefficients of these variables can be affected by multicollinearity

REMEDIES FOR MULTICOLLINEARITY

- ▶ Drop a redundant variable
 - when the variable is not needed to represent the effect on the dependent variable
 - in case of severe multicollinearity, it makes no statistical difference which variable is dropped
 - theoretical underpinnings of the model should be the basis for such a decision
- ▶ Do nothing
 - when multicollinearity does not cause insignificant t -scores or unreliable estimated coefficients
 - deletion of collinear variable can cause specification bias
- ▶ Increase the size of the sample
 - the confidence intervals are narrower when we have more observations

EXAMPLE

- ▶ Estimating the demand for gasoline in the U.S.:

$$\widehat{PCON}_i = 389.6 - \underset{(13.2)}{36.5} TAX_i + \underset{(10.3)}{60.8} UHM_i - \underset{(0.043)}{0.061} REG_i$$
$$t = 5.92 \qquad - 2.77 \qquad - 1.43$$

$$R^2 = 0.924 \quad , \quad n = 50 \quad , \quad Corr(UHM, REG) = 0.978$$

- $PCON_i$... petroleum consumption in the i -th state
 TAX_i ... the gasoline tax rate in the i -th state
 UHM_i ... urban highway miles within the i -th state
 REG_i ... motor vehicle registrations in the i -th state

EXAMPLE

- ▶ We suspect a multicollinearity between urban highway miles and motor vehicle registration across states, because those states that have a lot of highways might also have a lot of motor vehicles.
- ▶ Therefore, we might run into multicollinearity problems. How do we detect multicollinearity?

Look at correlation coefficient. It is indeed huge (0.978).

Look at the coefficients of the two variables. Are they both individually significant? *UHM* is significant, but *REG* is not. This further suggests a presence of multicollinearity.

- ▶ Remedy: try dropping one of the correlated variables.

EXAMPLE

$$\widehat{PCON}_i = 551.7 - \frac{53.6}{(16.9)} TAX_i + \frac{0.186}{(0.012)} REG_i$$
$$t = -3.18 \qquad 15.88$$

$$R^2 = 0.866 \quad , \quad n = 50$$

$$\widehat{PCON}_i = 410.0 - \frac{39.6}{(13.1)} TAX_i + \frac{46.4}{(2.16)} UHM_i$$
$$t = -3.02 \qquad 21.40$$

$$R^2 = 0.921 \quad , \quad n = 50$$

HETEROSKEDASTICITY

- ▶ **Observations of the error term are drawn from a distribution that has no longer a constant variance**

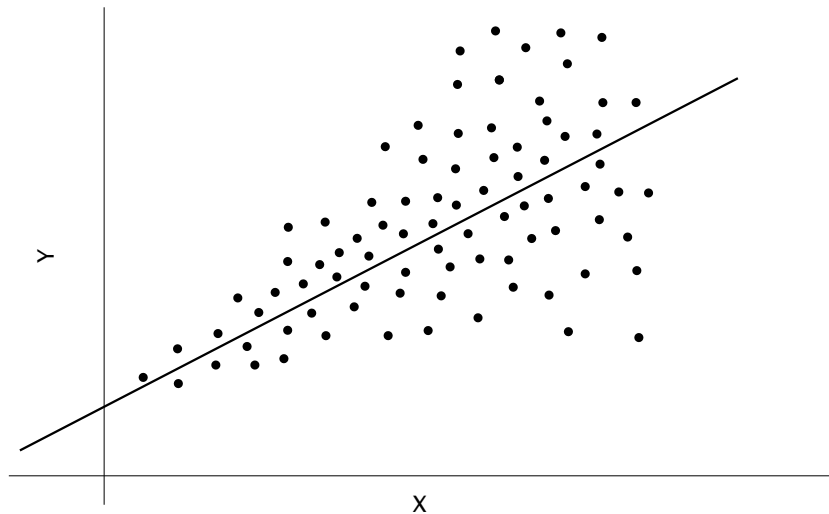
$$\text{Var}(\varepsilon_i) = \sigma_i^2, \quad i = 1, 2, \dots, n$$

Note: constant variance means: $\text{Var}(\varepsilon_i) = \sigma^2 (i = 1, 2, \dots, n)$

- ▶ **Often occurs in data sets in which there is a wide disparity between the largest and smallest observed values**
Smaller values often connected to smaller variance and larger values to larger variance (e.g. consumption of households based on their income level)
- ▶ **One particular form of heteroskedasticity (variance of the error term is a function of some observable variable):**

$$\text{Var}(\varepsilon_i) = h(x_i), \quad i = 1, 2, \dots, n$$

HETEROSKEDASTICITY



CONSEQUENCES OF HETEROSKEDASTICITY

- ▶ Violation of one of the classical assumptions
1. Estimates remain unbiased and consistent (estimated coefficients are not affected)
 2. Estimated standard errors of the coefficients are biased
 - heteroskedastic error term causes the dependent variable to fluctuate in a way that the OLS estimation procedure attributes to the independent variable
 - heteroskedasticity biases t statistics, which leads to unreliable hypothesis testing
 - typically, we encounter underestimation of the standard errors, so the t scores are incorrectly too high

DETECTION OF HETEROSKEDASTICITY

- ▶ There is a battery of tests for heteroskedasticity

Sometimes, simple visual analysis of residuals is sufficient to detect heteroskedasticity

- ▶ We will derive a test for the model

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$$

- ▶ The test is based on analysis of residuals

$$e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 z_i)$$

- ▶ The null hypothesis for the test is no heteroskedasticity:

$$E(e^2) = \sigma^2$$

Therefore, we will analyse the relationship between e^2 and explanatory variables

BREUSCH PAGAN TEST FOR HETEROSKEDASTICITY

1. Estimate the equation, get the residuals e_i
2. Regress the squared residuals on all explanatory variables:

$$e_i^2 = \alpha_0 + \alpha_1 x_i + \alpha_2 z_i + v_i \quad (2)$$

3. Get the R^2 of this regression and the sample size n
4. Test the joint significance of (2): use F-test or alternative test statistic, $LM = nR^2 \sim \chi_k^2$, where k is the number of slope coefficients in (2)
5. If nR^2 is larger than the χ^2 critical value, then we have to reject H_0 of no heteroskedasticity

WHITE TEST FOR HETEROSKEDASTICITY

1. Estimate the equation, get the residuals e_i
2. Regress the squared residuals on all explanatory variables and on squares and cross-products of all explanatory variables:

$$e_i^2 = \alpha_0 + \alpha_1 x_i + \alpha_2 z_i + \alpha_3 x_i^2 + \alpha_4 z_i^2 + \alpha_5 x_i z_i + v_i \quad (2)$$

3. Get the R^2 of this regression and the sample size n
4. Test the joint significance of (2): test statistic $= nR^2 \sim \chi_k^2$, where k is the number of slope coefficients in (2)
5. If nR^2 is larger than the χ_k^2 critical value, then we have to reject H_0 of no heteroskedasticity

REMEDIES FOR HETEROSKEDASTICITY

1. Redefining the variables

in order to reduce the variance of observations with extreme values

e.g. by taking logarithms or by scaling some variables

2. Weighted Least Squares (WLS)

consider the model $y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \varepsilon_i$

suppose $Var(\varepsilon_i) = \sigma^2 z_i^2$

it can be proved that if we redefine the model as

$$\frac{y_i}{z_i} = \beta_0 \frac{1}{z_i} + \beta_1 \frac{x_i}{z_i} + \beta_2 + \frac{\varepsilon_i}{z_i} ,$$

it becomes homoskedastic

3. Heteroskedasticity-corrected robust standard errors

HETEROSKEDASTICITY-CORRECTED ROBUST ERRORS

- ▶ The logic behind:
 - Since heteroskedasticity causes problems with the standard errors of OLS but not with the coefficients, it makes sense to improve the estimation of the standard errors in a way that does not alter the estimate of the coefficients (White, 1980)
- ▶ Heteroskedasticity-corrected standard errors are typically larger than OLS s.e., thus producing lower t scores
- ▶ In panel and cross-sectional data with group-level variables, the method of **clustering** the standard errors is the desired answer to heteroskedasticity

SUMMARY

- ▶ Multicollinearity

 - does not lead to inconsistent estimates, but it makes them lose significance

 - if really necessary, can be remedied by dropping or transforming variables, or by getting more data

- ▶ Heteroskedasticity

 - does not lead to inconsistent estimates, but invalidates inference

 - can be simply remedied by the use of (clustered) robust standard errors

- ▶ Readings:

 - Studenmund Chapter 8 and 10

 - Wooldridge Chapter 8