

Statistical and Econometric Models

Lectures: Lennart Oelschläger · Tutorials: Sebastian Büscher ✉

Problem Set – Week 4 (for tutorial on May 13, 2024)**Violations of MLR assumptions****Problem 4.1:**

We want to check the behaviour of the OLS estimator when the assumptions MLR.1 - MLR.6 are not exactly matched. For this, we will simulate data from the following underlying process:

$$y_i = 3 + 0.5 x_{1,i} + 7 x_{2,i} + 0.003 x_{2,i}^3 + 0.001 x_{3,i} + u_i, \quad u_i \sim \mathcal{N}(0, \sigma_i^2), \quad (1)$$

with $\sigma_i^2 = 0.002 x_{1,i}$.

We will, however, assume to not know the true data generating process and fit the model

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + v_i, \quad v_i \sim \mathcal{N}(0, \sigma^2). \quad (2)$$

- (a) Which assumptions of the classical linear model are violated by this model?
- (b) Draw 500 times from the data generating process (1). For this, draw the values of the regressors randomly, such that x_1 is drawn from a uniform distribution on the interval $[1, 3]$, x_2 is drawn from a normal distribution with mean zero and standard deviation three, and x_3 is drawn from a normal distribution with mean zero and standard deviation one.
- (c) Estimate model (2) on the simulated data to obtain the estimates of the parameters and the standard deviations of the estimates.
- (d) Repeat (b) and (c) 10000 times and save the estimates you obtain every time. Compare the mean on the 10000 estimates with the true parameters. Compare the standard deviation of the 10000 estimates of each parameter with the standard deviation of the estimators you calculated in (c).
- (e) Perform a RESET-Test with $p = 3$ on the model estimated in part (c). Does it detect the misspecification of the model?
- (f) Add now x_2^2 and x_2^3 to the model (2) and fit it to the data.
- (g) Repeat (d) for this new model and save the standard deviations of the 10000 obtained estimates.
- (h) Test the model obtained in (e) for heteroskedasticity. You can use a Breusch-Pagan test or a White test.
- (i) Based on the test you used to detect heteroskedasticity, calculate the Feasible Generalised Least Squares estimator. To do so, calculate the matrix L with $L^2 = \text{diag}(\hat{\sigma}_1^2, \dots, \hat{\sigma}_N^2)$ and estimate the rewritten model

$$L^{-1}y = L^{-1}X\beta + L^{-1}u.$$

- (j) Calculate the heteroskedasticity-consistent White estimator for the variance of the OLS estimator

$$\widehat{\text{Cov}}(\hat{\beta}) = (X'X)^{-1}X' \text{diag}(\hat{\varepsilon}_1, \dots, \hat{\varepsilon}_N)X(X'X)^{-1}.$$

- (k) Compare the estimates and their standard deviations of the FGLS estimation and the White estimator with those obtained in (f) and to the standard deviations of the 10000 estimates obtained in (g).

Instrumental Variables

Problem 4.2:

In the linear regression framework $y = X\beta + \varepsilon$, a central assumption states $E(\varepsilon | X) = 0$.

- (a) Show that this assumption implies $\text{Cov}(\varepsilon, X) = 0$ (this is referred to as “exogeneity” of X).
- (b) In the opposite case (“endogeneity”), X is correlated with ε . In this case, the OLS estimator $\hat{\beta}_{\text{OLS}}$ is biased. As a remedy, a technique called “instrumental variables estimation” was introduced in the lecture. What are “instrumental variables”?
- (c) The instrumental variables estimator for β is defined as

$$\hat{\beta}_{\text{IV}} = (Z'X)^{-1}Z'y.$$

Derive this equation by multiplying the equation $y = X\beta + \varepsilon$ by Z' on both sides, assuming that in the DGP it holds that $Z'\varepsilon = 0$ and $Z'X$ is invertible.

- (d) Consider the model

$$\log(y) = x'\beta + \varepsilon,$$

where y is the increase in sales of a company compared to last year and x are their expenditures for marketing. One could argue that this is a case of endogeneity. Why? Would the number of employees be a valid instrument? Is there a problem in using the age of the CEO as an instrument?