

## Exercise – Understanding the omitted variable bias using simulation

### 1. Aim of the exercise

In regression analysis, we often encounter situations where we cannot observe all the independent variables that are potentially correlated with the main variable of interest and hence should be controlled for. When this happens, the zero conditional mean assumption is violated. Econometrics classes demonstrate the resulting bias in the coefficient estimate of the main variable of interest theoretically, which makes the understanding of this bias somewhat abstract. This exercise illustrates the omitted variable bias through simulation. It shows how the sampling distribution of the OLS estimator, in particular its mean, is affected when a variable relevant to the main variable of interest is omitted from the regression.

### 2. Theory

Consider the linear regression model

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \varepsilon_i.$$

where for simplicity we do not consider a constant. Suppose that

$$\mathbb{E}[\varepsilon_i | x_{i1}] = 0,$$

and

$$\mathbb{E}[\varepsilon_i | x_{i2}] = 0.$$

Suppose that we do not observe  $x_{i2}$  so that it enters the error. The model becomes

$$y_i = x_{i1}\beta_1 + \varepsilon_i^*$$

where

$$\varepsilon_i^* = x_{i2}\beta_2 + \varepsilon_i.$$

Then,

$$\begin{aligned}\mathbb{E}[\varepsilon_i^* | x_{i1}] &= \mathbb{E}[x_{i2}\beta_2 | x_{i1}] + \mathbb{E}[\varepsilon_i | x_{i1}] \\ &= \beta_2 \mathbb{E}[x_{i2} | x_{i1}] + 0 \\ &\neq 0\end{aligned}$$

if  $\beta_2 \neq 0$  and  $\mathbb{E}[x_{i2} | x_{i1}] \neq 0$ .  $\beta_2 \neq 0$  means that  $x_{i2}$  should enter the model.  $\mathbb{E}[x_{i2} | x_{i1}] \neq 0$  means that  $x_{i1}$  and  $x_{i2}$  are correlated. The zero conditional mean assumption in this case is violated for  $\varepsilon_i^*$ .

What is the implication of

$$\mathbb{E}[\varepsilon_i^* | x_{i1}] \neq 0$$

for the OLS estimator  $\hat{\beta}_1$ ? The OLS estimator  $\hat{\beta}_1$  when  $x_{i2}$  is omitted from the regression is given by

$$\begin{aligned}\hat{\beta}_1 &= (\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \mathbf{y} \\ &= (\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 (\mathbf{x}_1 \beta_1 + \mathbf{x}_2 \beta_2 + \boldsymbol{\varepsilon}) \\ &= \beta_1 + (\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \mathbf{x}_2 \beta_2 + (\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \boldsymbol{\varepsilon}.\end{aligned}$$

$\hat{\beta}_1$  represents the impact of  $x_1$  on  $y$  where  $y$  is in fact driven not only by  $x_1$  but also  $x_2$  according to the true DGP. This means that we explain  $y$  only with  $x_1$  whereas we should explain it also with  $x_2$ .

Taking the expectation conditional on  $\mathbf{X}$ , we have

$$\mathbb{E} \left[ \hat{\beta}_1 \mid \mathbf{X} \right] = \beta_1 + (\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \mathbf{x}_2 \beta_2$$

since  $\mathbb{E}[\varepsilon \mid \mathbf{X}] = \mathbf{0}$  in the true model. In two cases  $\hat{\beta}_1$  is an unbiased estimator. First, if

$$(\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \mathbf{x}_2 = 0,$$

meaning that there is no correlation between  $\mathbf{x}_1$  and  $\mathbf{x}_2$  in the sample. Realize that the stated expression is the OLS estimate of the coefficient of  $\mathbf{x}_1$  from the regression of  $\mathbf{x}_2$  on  $\mathbf{x}_1$ . Second, if

$$\beta_2 = 0,$$

meaning that  $\mathbf{x}_2$  does not enter the true model. Otherwise  $\hat{\beta}_1$  is subject to the omitted variable bias.

### 3. Application

#### 3.1. Clear the memory

Clear the memory from possible calculations from an earlier session.

```
11 % 3.1. Clear the memory
12 clear;
```

#### 3.2. Set the number of simulations

Set the number of simulations to be carried out.

```
14 % 3.2. Set the number of simulations
15 N_sim = 1000;
```

#### 3.3. Set the sample size

Assume a linear regression model with `N_obs` observations available for the variables of this model.

```
17 % 3.3. Set the sample size
18 N_obs = 1000;
```

#### 3.4. Set true values for the coefficients of the model

Assume that the linear regression model contains two independent variables and, for simplicity, no constant. Assume also that this is the true DGP and that we know the true values of the coefficients.

```

20 % 3.4. Set true values for the coefficients
21 B_true = [0.5 0.75]';

```

### 3.5. Create the constant term

Create the constant term.

```

23 % 3.5. Create the constant term
24 x_0 = ones(N_obs,1);

```

### 3.6. Create a vector of covariances between two independent variables

A problem related to the systematic part of the DGP is omitting a relevant independent variable that is part of the true DGP. One potential reason is that the researcher does not know that a particular variable belongs in the specification or cannot collect data on that variable. If the omitted independent variable is uncorrelated with all of the independent variables that are included in the regression model, leaving it out will not bias the estimated coefficients. The omitted variable could, however, still explain some part of the variation in the dependent variable. If the omitted independent variable is correlated with one or more of the independent variables, this will bias the coefficient estimates of those variables.

Here we examine the omitted variable problem across a range of correlations between an included independent variable and the omitted independent variable. We define correlation levels using 10 different covariance values, ranging from 0 to 0.99. Line 27 creates a vector array containing these values. We assume that each variable has a variance of 1. Line 28 defines the column dimension of this vector. We will use this variable when iterating over different correlation scenarios in our simulation.

```

26 % 3.6. Create a vector of covariances between two independent variables
27 sigma_x_1_x_2 = 0:0.11:0.99;
28 N_sig = size(sigma_x_1_x_2,2);

```

### 3.7. Preallocate a matrix to store the simulated OLS coefficient estimates

Preallocate a matrix that will store the coefficient estimates from repeated sampling, simulated at different correlation levels between the included and omitted independent variables. The matrix is  $N_{\text{sim}} \times N_{\text{sig}}$  because we will simulate  $N_{\text{sim}}$  samples, at  $N_{\text{sig}}$  different levels of correlation.

```

30 % 3.7. Preallocate a matrix for storing OLS estimates from all samples
31 B_hat_1_sim = NaN(N_sim,N_sig);

```

### 3.8. Define an input argument for the multivariate normal random number generator

In the next section, we will draw random numbers from the multivariate normal distribution to create two independent variables that are correlated with each other. To do this, we will make use of the built-in MATLAB `mvnrnd` function. The function accepts three input arguments. The first input argument is the mean vector of the distribution. The second input argument is

the covariance matrix of the distribution. The third input argument specifies the number of observations to be drawn for each random variable of the distribution. The third input argument is defined above. Here we define the first input argument. The second input argument will be defined within the simulation in the next section because in each iteration of the simulation the covariance matrix will be updated in accordance with the different correlation levels between two two variables.

```
34 % 3.8. Define mean vector for the multivariate random number generator
35 mu = [0 0];
```

3.9. Create the sampling distribution of the OLS estimator at different correlation levels between the included and omitted independent variables

Here we consider a nested loop structure, with one outer loop and one inner loop. The inner for loop simulates `N_sim` coefficient estimates from repeated sampling. The outer for loop repeats this simulation for 10 different correlation levels between the included and the omitted variables. Consider the inner for loop. Line 37 defines the index of the for loop. Line 39 defines the covariance matrix of the included and omitted variables at given covariance values. Variances of each independent variable are set to 1. In line 41 we supply the `mvnrnd` function with the defined covariance matrix, and with two other input arguments defined in the previous section. The function generates random values for the included and omitted variables from the multivariate normal distribution so that the variables are correlated.

In lines 42 and 43 we define the included (`x_1`) and omitted (`x_2`) variables.

In line 45 we generate random values for the error term. Note that we rule out heteroskedasticity by setting the standard deviation of the error term to 1. In line 46 we generate values for the dependent variable using both `x_1` and `x_2` following the true DGP. In line 47 we use the function `exercisefunction` to estimate the coefficient of `x_1` using this dependent variable but `x_1` as the only explanatory variable, incorrectly ignoring `x_2`. In line 48 we collect the simulated coefficient estimates at different levels of correlation between the included and the omitted variable in the matrix array `B_hat_1_sim(i,j)`.

```
36 % 3.9. Create the sampling distribution of the biased OLS estimator
37 for j = 1:N_sig
38     for i = 1:N_sim
39         Sigma = reshape([1 sigma_x_1_x_2(:,j) sigma_x_1_x_2(:,j) 1] ...
40                        ,2,2);
41         x_1_x_2_mvn = mvnrnd(mu,Sigma,N_obs);
42         x_1 = x_1_x_2_mvn(:,1);
43         x_2 = x_1_x_2_mvn(:,2);
44         X = [x_1 x_2];
45         u = random('Normal',0,1,[N_obs 1]);
46         y = X*B_true+u;
47         LSS = exercisefunctionlss(y,x_1);
48         B_hat_1_sim(i,j) = LSS.B_hat(1,1);
49     end
50 end
```

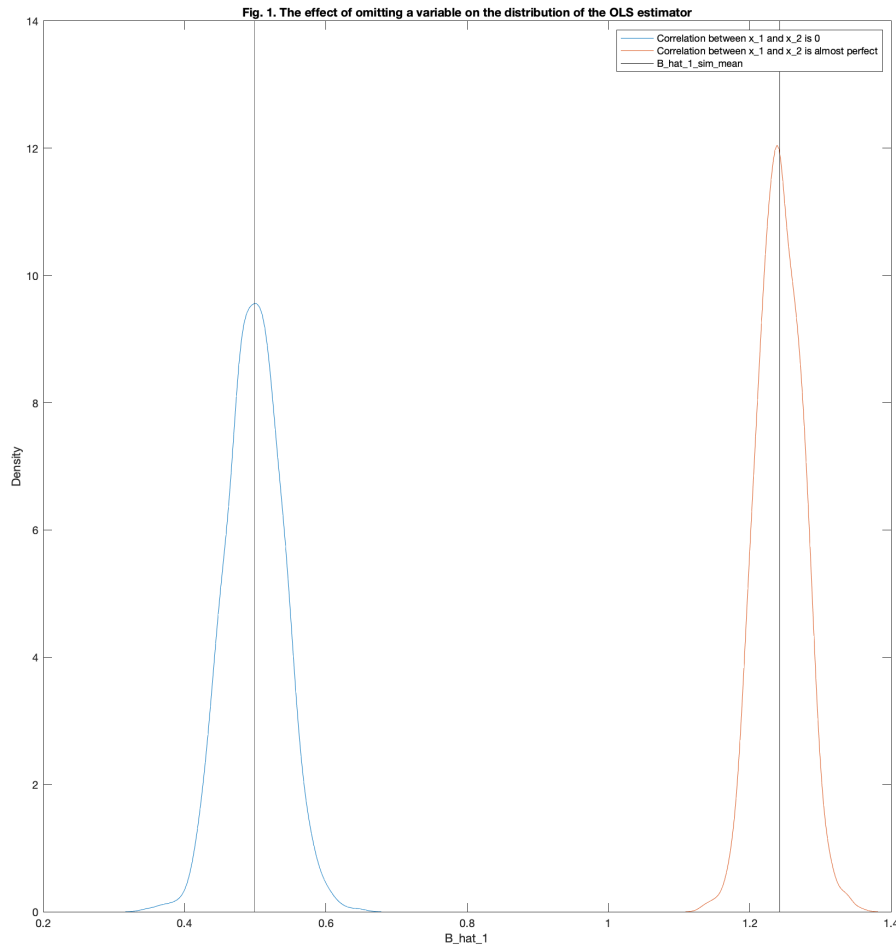
4. Plot the sampling distribution of the biased OLS estimator

The plot produced here shows the density estimate of the 1000 OLS coefficient estimates of  $x_1$ , at a correlation of 0 and 0.99 between  $x_1$  and  $x_2$ . The distribution of estimates at correlation 0 is centered right at 0.5, indicating no bias. This demonstrates that omitting a variable that is not correlated with an included variable does not affect parameter estimates in the case of OLS. In contrast, the distribution of estimates when the correlation is 0.99 shows a substantial amount of bias. In fact, recall from above that the coefficient on the omitted variable was set to 0.75. The mean of the distribution with correlation 0.99 is 1.242, which is 0.5 (true coefficient of  $x_1$ ) plus 0.742. Hence, at near-perfect correlation with the omitted variable, almost all of the true effect of that omitted variable is incorrectly attributed to  $x_1$  through the biased estimate of the coefficient of  $x_1$ .

```

52 %% 4. Plot the sampling distribution of the biased OLS estimator
53 ksdensity(B_hat_1_sim(:,1))
54 hold on
55 ksdensity(B_hat_1_sim(:,10))
56 hold on
57 line([mean(B_hat_1_sim(:,1)) mean(B_hat_1_sim(:,1))],ylim, ...
58      'Color','black')
59 hold on
60 line([mean(B_hat_1_sim(:,10)) mean(B_hat_1_sim(:,10))],ylim, ...
61      'Color','black')
62 title(['The Effect of Omitting a Variable on the Distribution ' ...
63      'of the OLS estimator'])
64 legend('Correlation between x\1 and x\2 is 0',['Correlation ' ...
65      'between x\1 and x\2 is almost perfect'], ...
66      'B\_hat\1\_sim\_mean')
67 ylabel('Density')
68 xlabel('B\_hat\1')

```



## 5. Other simulation experiments

One can explore how changing aspects of the simulation conducted above changes the nature of the results. For example, increase the sample size for each draw in the simulation. What do you conclude? This kind of exploration helps to see the real nature of the DGP and how the chosen statistical estimator performs.

## 6. Final notes

This file is prepared and copyrighted by Tunga Kantarcı. Parts of the simulation exercise are based on Carsey, T. M., and Harden, J. J., 2014. Monte Carlo simulation and resampling methods for social science. SAGE Publications. This file and the accompanying MATLAB files are available on GitHub and can be accessed via this [link](#).