

Exercise – Understanding the normality assumption of the regression error

1. Aim of the exercise

In finite samples, the errors of the linear regression model are assumed to follow a normal distribution. This exercise examines the implications for the sampling distribution of the OLS estimator when the regression errors do not follow a normal distribution.

2. Application

2.1. Clear the memory

Clear the memory from possible calculations from an earlier session.

```
1 % 2.1. Clear the memory
2 clear;
```

2.2. Set the number of simulations

Set the number of simulations to be carried out.

```
1 % 2.2. Set the number of simulations
2 N_sim = 1000;
```

2.3. Set the sample size

Assume that we have a linear regression model that contains a constant term and an independent variable. Assume also that we have `N_obs` observations for the variables of this model.

```
1 % 2.3. Set the sample size
2 N_obs = 350;
```

2.4. Set true values for the coefficients

Assume that we know the true values of the coefficients of the variables of the linear regression model we consider, and that these values are as indicated at the end of the section.

```
1 % 2.4. Set true values for the coefficients
2 B_true = [0.2 3.5]';
```

2.5. Define the number of coefficients to be simulated

Define the number of coefficients to be simulated.

```
1 % 2.5. Define the number of coefficients to be simulated
2 N_par = 1;
```

2.6. Create the systematic component of the regression equation

Create the constant term. Draw a set of random numbers from the uniform distribution, and require the numbers to be in the range $[-1, 1]$. Consider this vector as the independent variable of the regression model. Create the systematic component of the regression equation, and call it X .

```
1 % 2.6. Create the systematic component of the regression equation
2 x_0 = ones(N_obs,1);
3 x_1 = random('Uniform',-1,1,[N_obs 1]);
4 X = [x_0 x_1];
```

2.7. Preallocate matrices for storing statistics to be simulated

Preallocate matrices that will store simulated coefficient estimates generated under different distributional assumptions for the error, using different estimators. Each matrix is $N_{\text{sim}} \times N_{\text{par}}$ because we simulate N_{sim} times the coefficient of the independent variable x_1 , and we have N_{par} coefficients to simulate.

```
1 B_hat_1_sim_OLS_normal = NaN(N_sim,N_par);
2 B_hat_1_sim_IRLS_normal = NaN(N_sim,N_par);
3 B_hat_1_sim_OLS_t = NaN(N_sim,N_par);
4 B_hat_1_sim_IRLS_t = NaN(N_sim,N_par);
```

2.8. Degrees of freedom of the t distribution

The t distribution depends the degrees of freedom parameter. The parameter controls the kurtosis of the distribution. Create this parameter, and assume an experimental value of 2 for it. We will use this parameter to generate random draws from the t distribution later in our simulation.

```
11 % 2.8. Define the degrees of freedom for t distribution
12 t_df = 2;
```

2.9. Create sampling distributions for the OLS and IRLS estimators based on errors with different distributions

We will examine the implications of violating the normality assumption. There are many ways in which this assumption can be violated, but here we focus on only one. In particular, we study the consequences for the sampling distribution of the OLS estimator when the true error distribution has heavier tails than the normal distribution. Distributions with heavy tails are important to social scientists because they are more likely to generate observations that are outliers compared with what one would expect from the normal distribution. In this example, we compare the performance of the OLS estimator with that of an alternative estimator that is robust to outliers. One distribution that has heavier tails than the standard normal distribution is the t distribution.

How does assuming that the error term follows a t distribution effect estimation? To answer this question we will compare the performance of the OLS estimator with that of an alternative estimator. This alternative estimator is the Iteratively Reweighted Least Squares estimator (IRLS), which produces robust estimates when outlying observations are present: <https://>

en.wikipedia.org/wiki/Iteratively_reweighted_least_squares. By minimizing the sum of absolute residuals rather than squared residuals like OLS, IRLS is not disproportionately influenced by outliers. Indeed, while outliers can have a substantial effect on OLS estimates, their impact on those of IRLS is smaller. As a result, IRLS is a more efficient estimator than the OLS estimator under this circumstance. IRLS allows the analyst to handle heavy-tailed data without giving special treatment to outliers or just deleting them.

To illustrate this, we simulate a model with errors drawn from standard normal distribution, and one with errors drawn from the t distribution. The t distribution has one parameter: degrees of freedom. Above, we have set it to 2. After creating a DGP with standard normal errors and a DGP with t errors, we estimate models on both types of data using the OLS and IRLS estimators. For this, we use the built-in `robustfit` function of MATLAB that offers to produce standard OLS estimates as well as IRLS estimates if the option for the estimation method is specified accordingly in the function syntax of the `robustfit` function.

Consider the for loop presented at the end of the section. We specify a DGP with errors that have a standard normal distribution, and another with errors that have a t distribution. We then estimate the regression with standard normal errors, using the OLS estimator. Next, we estimate the same regression using the IRLS estimator. Do you expect these estimators to produce similar coefficient estimates? In the remaining lines, we estimate the regression with errors that have a t distribution, using the OLS and IRLS estimators. In this case, do you expect these estimators to produce similar coefficient estimates?

```

1 % 2.9. Create sampling distributions for the OLS and IRLS estimators
2 for i = 1:N_sim
3     u_normal = random('Normal',0,1,[N_obs 1]);
4     y_normal = X*B_true+u_normal;
5     u_t = random('t',t_df,[N_obs 1]);
6     y_t = X*B_true+u_t;
7     OLS = robustfit(x_1,y_normal,'ols');
8     B_hat_1_sim_OLS_normal(i,1) = OLS(2,1);
9     IRLS = robustfit(x_1,y_normal,'bisquare');
10    B_hat_1_sim_IRLS_normal(i,1) = IRLS(2,1);
11    OLS = robustfit(x_1,y_t,'ols');
12    B_hat_1_sim_OLS_t(i,1) = OLS(2,1);
13    IRLS = robustfit(x_1,y_t,'bisquare');
14    B_hat_1_sim_IRLS_t(i,1) = IRLS(2,1);
15 end

```

3. Plot example distributions of errors with different distributional assumptions

In the previous section, in the for loop, we have generated two types of errors. One with a standard normal distribution, and another one with a t distribution. Plotted here are kernel smoothed histograms of these errors. The distributions result from the last draw in the for loop. Looking at these two distributions, what do you conclude?

```

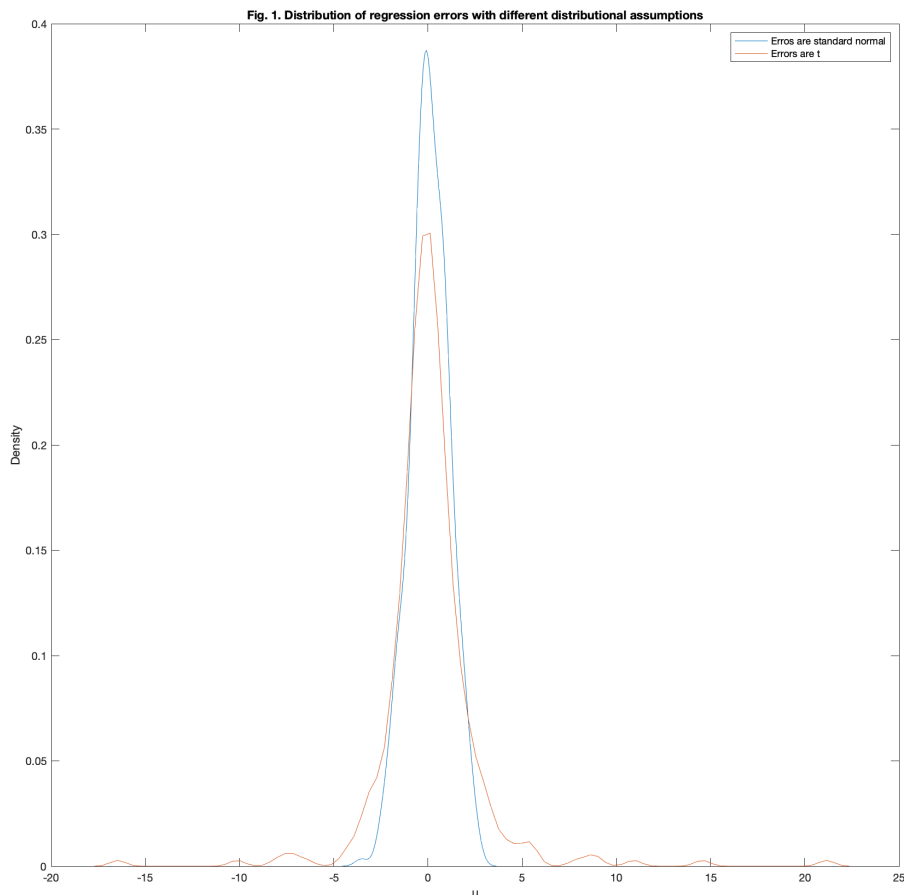
1 %% 3. Plot errors with different distributional assumptions
2 ksdensity(u_normal)
3 hold on
4 ksdensity(u_t)
5 title(['Fig. 1. Distribution of regression errors with different ' ...

```

```

6     'distributional assumptions'])
7 legend('Errors are standard normal','Errors are t')
8 ylabel('Density')
9 xlabel('u')
10 hold off

```



4. Plot the sampling distributions of the OLS and IRLS estimators when errors are normal

Plot the sampling distributions of the OLS and IRLS estimators when the errors of the regression are assumed to be standard normal. The two distributions are close to each other. Is this a surprising result?

```

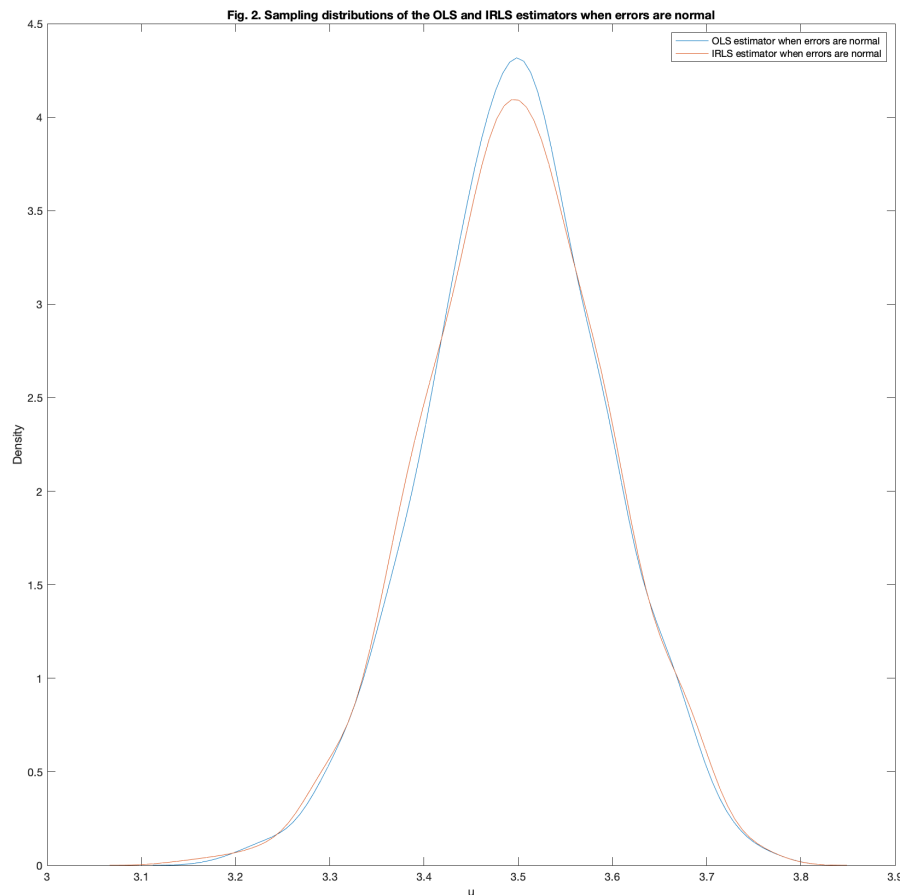
11 %% 4. Plot the sampling distributions of estimators when errors are normal
12 ksdensity(B_hat_1_sim_OLS_normal(:,1))
13 hold on
14 ksdensity(B_hat_1_sim_IRLS_normal(:,1))
15 title(['Fig. 2. Sampling distributions of the OLS and IRLS ' ...
16       'estimators when errors are normal'])
17 legend('OLS estimator when errors are normal','IRLS estimator ' ...

```

```

18     'when errors are normal'])
19 ylabel('Density')
20 xlabel('u')
21 hold off

```



5. Plot the sampling distributions of the OLS and IRLS estimators when errors are t

Here we plot the sampling distributions of the OLS and IRLS estimates when the errors of the regression follow a t distribution. Given the samples size of `N_obs` observations, the OLS estimator appears to be less efficient.

Remember that the OLS estimator is BLUE (when the variance-covariance matrix has a scalar form). That is, apart from being an unbiased estimator, it is the most efficient estimator. But the plot suggests that it is not more efficient than the IRLS estimator when the errors are t distributed. So how can we still claim that the OLS estimator is the most efficient? We need to realize that the OLS estimator is a linear estimator while the IRLS estimator is a nonlinear estimator. So we cannot compare the two estimators based on their efficiency because they are not both linear or non-linear estimators.

```

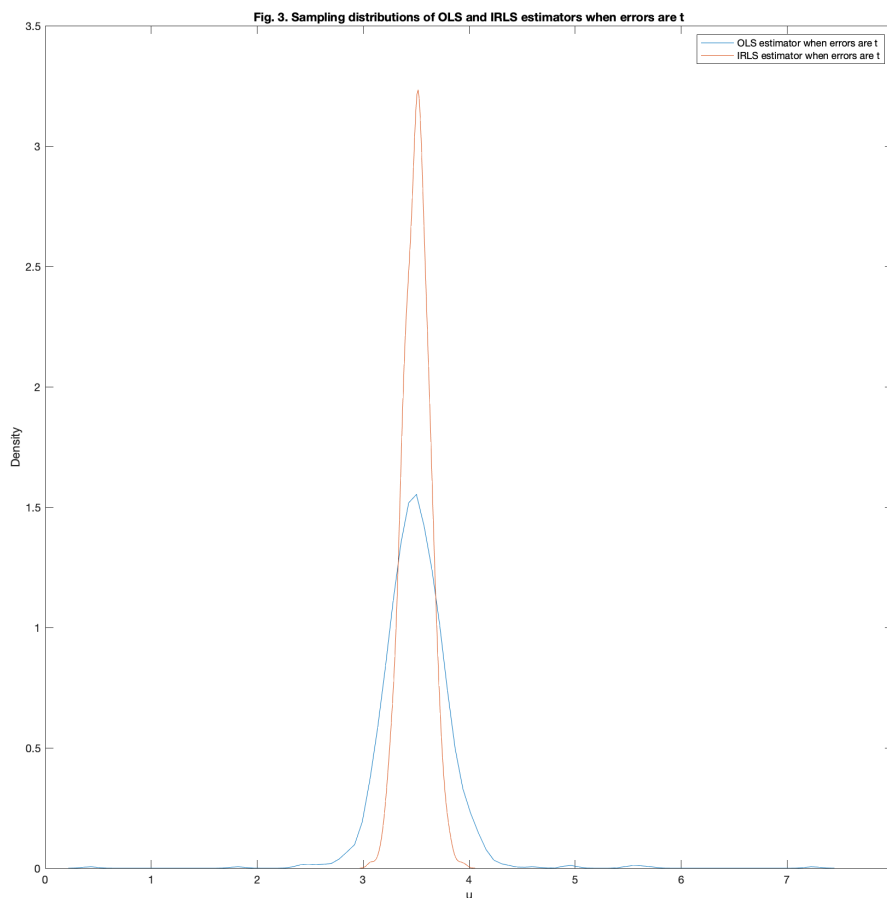
1 %% 5. Plot the sampling distributions of estimators when errors are t

```

```

2 ksdensity(B_hat_1_sim_OLS_t(:,1))
3 hold on
4 ksdensity(B_hat_1_sim_IRLS_t(:,1))
5 title(['Fig. 3. Sampling distributions of OLS and IRLS estimators ' ...
6       'when errors are t'])
7 legend('OLS estimator when errors are t', ['IRLS estimator ' ...
8       'when errors are t'])
9 ylabel('Density')
10 xlabel('u')
11 hold off

```



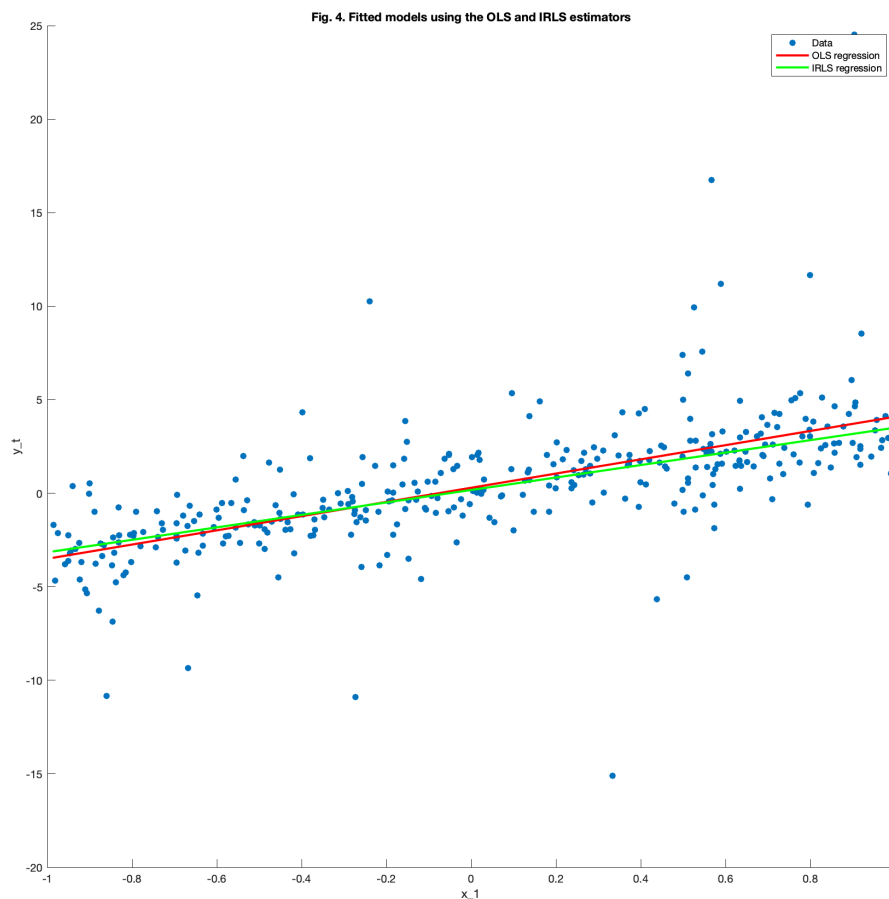
6. Plot the scatter plot and two regression lines fitted using the OLS and IRLS estimators

Here we overlay two types of plots. The first is a scatter plot of the dependent variable against the independent variable. The dependent variable is that of the DGP that imposes a t distribution on the errors. This scatter plot is overlaid by two regression lines fitted using the OLS and the IRLS estimators. The fitted line using the OLS estimator appears to be influenced by the outliers. On the other hand, the fitted line using the IRLS estimator is robust to the outliers.

Increase the sample size `N_obs` and execute the simulation again. The plot will show a change. What would explain this change?

Increase the degrees of freedom of the t distribution from 2 to, for example, 5. The plot will show a change. What would explain this change?

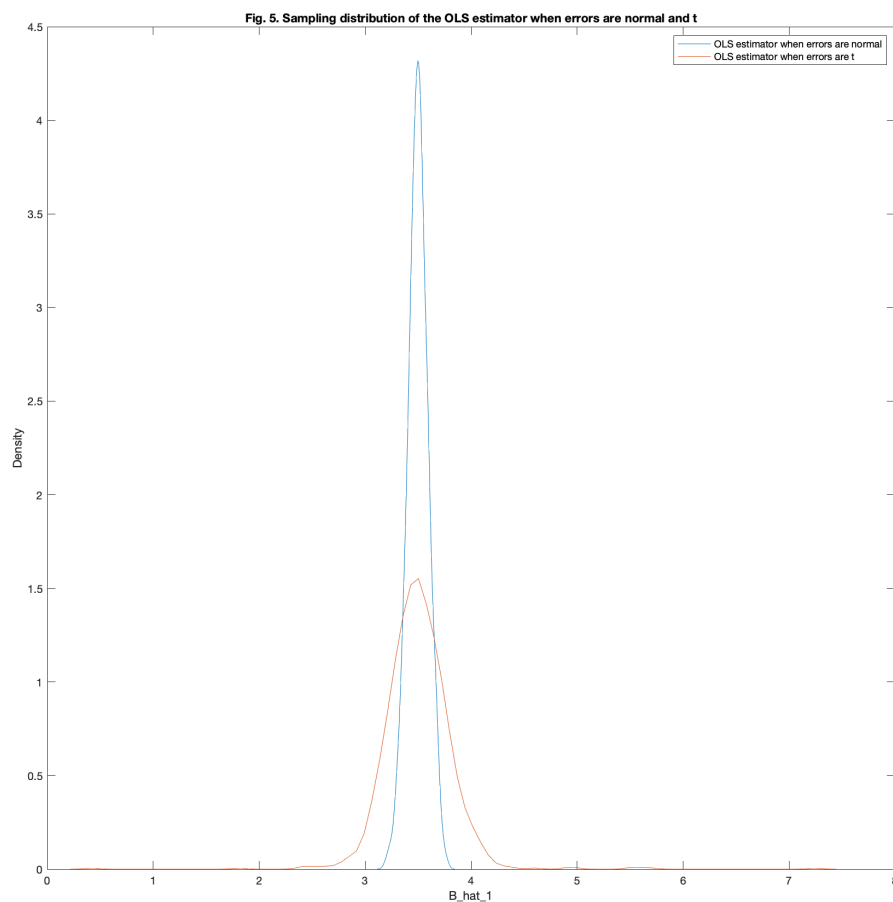
```
1 %% 6. Plot the fitted models using the OLS and IRLS estimators
2 scatter(x_1,y_t,'filled');
3 grid on;
4 hold on
5 y_t_hat_OLS = OLS(1)+OLS(2)*x_1;
6 y_t_hat_IRLS = IRLS(1)+IRLS(2)*x_1;
7 plot(x_1,y_t_hat_OLS,'red','LineWidth',2);
8 plot(x_1,y_t_hat_IRLS,'green','LineWidth',2)
9 title('Fig. 4. Fitted models using the OLS and IRLS estimators')
10 legend('Data','OLS regression','IRLS regression')
11 grid off;
12 ylabel('y\t')
13 xlabel('x\t_1')
14 hold off
```



7. Other simulation experiments

Plot the sampling distribution of the OLS estimator when regression errors are normal and when they are t. What do you conclude?

```
1 %% 7. Other simulation experiments
2 ksdensity(B_hat_1_sim_OLS_normal(:,1))
3 hold on
4 ksdensity(B_hat_1_sim_OLS_t(:,1))
5 title(['Fig. 5. Sampling distribution of the OLS estimator ' ...
6       'when errors are normal and t'])
7 legend('OLS estimator when errors are normal', ['OLS estimator ' ...
8       'when errors are t'])
9 ylabel('Density')
10 xlabel('B\_hat\_1')
11 hold off
```



8. Final notes

This file is prepared and copyrighted by Tunga Kantarcı. Parts of this simulation exercise is 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](#).