

**The International College of Economics and Finance**  
**Econometrics. Mid-year exam. 2018 October 25. Part 2. (120 minutes).**

**SECTION A. Answer ALL questions 1-3 from this section.**

**1. [15 marks]** You are asked to study the factors that determine employee monthly earnings  $E_i$  (in thousands of rubles) of a big corporation (60 observations in total), such as education  $ED_i$ , professional experience  $EXP_i$  (both in years), age  $AGE_i$  (also in years), and  $AGE_i^2$  (in years squared):

$$E_i = \beta_1 + \beta_2 ED_i + \beta_3 EXP_i + \beta_4 AGE_i + \beta_5 AGE_i^2 + u_i; \quad i = 1, 2, \dots, N. \quad (1)$$

Let estimated equation be  $\hat{E}_i = -632.01 + 1.94ED_i + 2.17EXP_i + 42.36AGE_i - 0.74AGE_i^2$   $R^2 = 0.18$  (1\*)  
 (800.86) (0.75) (0.86) (55.95) (0.97)

Another estimated equations is  $\hat{E}_i = -27.54 + 2.10ED_i + 2.39EXP_i$   $R^2 = 0.17$  (2\*)  
 (14.69) (0.69) (0.76)

**(a) [3 marks]** Outline briefly how to test whether education matters in determination of the monthly earnings using equation (1\*). How would you test whether experience influence positively on the earnings? Is it possible to test whether the earnings of all employees are constant independently of their education level, professional experience and age?

**(b) [6 marks]** Explain the meaning of coefficients  $\hat{\beta}_2$  and  $\hat{\beta}_3$  in (1\*). The coefficient  $\hat{\beta}_3$  is a little bigger than  $\hat{\beta}_2$ . Can you conclude from here that spending a year acquiring the professional experience is more useful for future earnings in the corporation than a year of study? Why or why not? How can you test the hypothesis  $\beta_3 = \beta_2$  against  $\beta_3 \neq \beta_2$  using F-test? How can you test the hypothesis  $\beta_3 = \beta_2$  against  $\beta_3 > \beta_2$  using appropriate test?

**(c) [6 marks]** Explain the reason for inclusion variable of  $AGE_i^2$  into the model? (You may use a graphic illustration in addition to the text). What is marginal effect of age in the model for the employee of 45 years old? Test whether age influence earnings and comment obtained results. Test whether this influence is linear (relationship between earnings and age is a linear function)?

**2. [15 marks]** A researcher is trying to build the best model for the total expenditure on private education  $PRIV_i$  in USA (in billions of dollars) for 1993-2017 as a function of disposable personal income  $DPI_i$  (also in billions of dollars), relative price index for private education  $PRELPRIV_i$  and  $TIME_i$  (equal to 1 in 1993, 2 – in 1994 and so on).

**(a) [5 marks]** First a linear and log-linear models were estimated:

$$PRIV_i = 29.72 + 0.018 \cdot DPI_i - 0.38 \cdot PRELPRIV_i + 0.32 \cdot TIME_i \quad R^2 = 0.97 \quad (1)$$

(11.81) (0.0086) (0.14) (0.25) RSS = 5.54

$$\log(PRIV_i) = 6.92 + 0.0021 \cdot DPI_i - 0.071 \cdot PRELPRIV_i + 0.067 \cdot TIME_i \quad R^2 = 0.94 \quad (2)$$

(1.66) (0.0012) (0.020) (0.036) RSS = 0.11

Give interpretation to the models and their coefficients (no justification for your interpretation is required). Discuss the significance of coefficient of  $DPI_i$  in equation (2) using 2-tail and 1-tail tests when needed. What assumptions should be taken to use of 1-tail test? Is it possible to choose between models (1) and (2) on the base of available information? Why or why not?

**(b) [5 marks]** The researcher also evaluates variable  $PRIVZ_i$  dividing values of  $PRIV_i$  by their geometric mean for the whole period  $PRIVZ_i = PRIV_i / \sqrt[n]{PRIV_1 \cdot \dots \cdot PRIV_n}$  and runs new regression

$$PRIVZ_i = 2.97 + 0.0018 \cdot DPI_i - 0.038 \cdot PRELPRIV_i + 0.032 \cdot TIME_i \quad R^2 = 0.97 \quad (1*)$$

(1.18) (0.0008) (0.014) (0.025) RSS = 0.055

Use Box-Cox test to choose between linear and logarithmic specifications of the model. Specify the formula for test statistic, its distribution and corresponding critical values. What is your conclusion? Which model can be selected as the best for further economic and econometric analysis? Explain your choice.

**(c) [5 marks]** The colleague of the researcher mentioned that equation (1\*) cannot be compared with log-linear regression (2) using Box-Cox test as Zarembka transformation was applied only to linear regression, while equation (2) uses untransformed dependent variable, and recommended to apply Zarembka transformation also to (2) and use regression (2\*) instead of (2):

$$\log(\hat{PRIV}_i) = 4.62 + 0.0021 \cdot DPI_i - 0.071 \cdot PRELPRIV_i + 0.067 \cdot TIME_i \quad R^2 = 0.94 \quad (2^*)$$

(1.66) (0.0012) (0.020) (0.036)  $RSS = 0.11$

Comment on her advice.

She also recommended to run double logarithmic regression (3)

$$\log(\hat{PRIV}_i) = -2.38 + 2.30 \cdot \log(DPI_i) - 2.24 \cdot \log(PRELPRIV_i) - 0.02 \cdot TIME_i \quad R^2 = 0.98 \quad (3)$$

(5.19) (0.22) (1.13) (0.02)  $RSS = 0.028$

instead of (2) and compare it with regression (1\*) using Box-Cox test. Comment on this advice, do recommended test and make your conclusion. Give interpretation to the coefficients of regression (3).

**3. [15 marks]** September 2018, elections for the mayor of Moscow were held. The winner was Sergei Sobyenin, his main competitor was Vadim Kumin (KPRF party). A student studying the factors influencing the victory of candidates got exit poll data from one of the polling stations. Each polled voter indicated the candidate for whom he voted (in this polling station all votes were distributed between Sobyenin and Kumin), and also reported how many campaign programs in favor of each of these two candidates (from 0 to 3) he had seen on television. His working hypothesis was that the decision of voters was affected by the number of views of campaign TV programs in favor of each of the candidates. As a result, the student has collected 16 observations:  $(S=0, K=0, VS_1)$ ,  $(S=1, K=0, VS_2)$ , ...,  $(S=3, K=3, VS_{16})$ , (where  $S$  is the number of views for Sobyenin,  $K$  is the number of views for Kumin,  $VS_i$  is the percentage of votes cast in each group for Sobyenin. For example in the group of voters with 3 TV views in favor of Sobyenin and with 1 TV view in favor of Kumin 96% of voters voted for Sobyenin. He also calculated the total number of TV views for each group  $T_i = S_i + K_i$  and calculated sample covariance between  $S$  and  $K$ :  $Cov(S, K) = 0$ .

Then he runs several regressions

$$\hat{VS}_i = 67.35 + 6.85S_i \quad R^2 = 0.26 \quad (1)$$

(5.74) (3.06)

$$\hat{VS}_i = 79.05 + 6.85S_i - 7.80K_i \quad R^2 = 0.60 \quad (2)$$

(5.60) (2.34) (2.34)

$$\hat{VS}_i = 79.05 + 14.65S_i - 7.80T_i \quad R^2 = 0.60 \quad (3)$$

(5.6) (3.30) (2.34)

**(a) [5 marks]** Give interpretation to the coefficients of the variable  $S_i$  in equations (1-3). Test the significance of the coefficients of these regressions and the significance of the equations as a whole. (*Do not compare coefficients or other characteristics of equations 1-3, you will be asked to do this later in b and c*)

**(b) [6 marks]** Why the intercepts in equations (2) and (3) are identical while in (1) it is different? Explain why the coefficient of the variable  $T_i$  in equations (3) is identical to the coefficient of the variable  $K_i$  in equation (2)? Explain why the coefficient of the variable  $S_i$  in equations (3) is bigger than in equation (2)? Explain why the coefficients of the variable  $S_i$  in equations (1) and (2) are identical? (*You may use the formulas for estimators of regression coefficients or any other method for your explanations*).

**(c) [4 marks]** Explain why determination coefficients  $R^2$  in equations (2) and (3) are identical, while in equation (1)  $R^2$  is about three times less. Explain why in equation (2) the standard errors of both slope coefficients are identical. Suggest probable reason why the standard error of the coefficient for the variable  $S_i$  in equation (2) is less than corresponding standard error in equation (1).

**Bonus question:** what drawbacks can you indicate in the student's study?

**SECTION B.** Answer **ONE** question from this section (4 OR 5).

**4. [30 marks]** At a lecture on econometrics the professor said to the students that OLS estimator for the coefficient  $\alpha$  in a simple regression model without constant  $Y_t = \alpha X_t + u_t$ ,  $t = 1, 2, \dots, T$  is given by

$$\hat{\alpha}_{OLS} = \frac{\sum_{t=1}^T X_t Y_t}{\sum_{t=1}^T X_t^2}. \text{ But this estimator is not the only one, there is another estimator of } \alpha, \hat{\alpha}_1 = \frac{\bar{Y}}{\bar{X}} = \frac{\frac{1}{T} \sum Y_t}{\frac{1}{T} \sum X_t}.$$

**(a) [8 marks]** Show that  $\hat{\alpha}_1 = \frac{\bar{Y}}{\bar{X}}$  is also an unbiased estimator of  $\alpha$  in  $Y_t = \alpha X_t + u_t$ . What assumptions are needed for this result being valid? Is  $\bar{Y} / \bar{X}$  also a good estimator for  $\beta_2$  in  $Y_t = \beta_1 + \beta_2 X_t + u_t$ ? Explain.

**(b) [7 marks]** Derive population variance of the estimator  $\hat{\alpha}_1 = \frac{\bar{Y}}{\bar{X}}$ . What assumptions are used to derive your result?

**(c) [8 marks]** Giving a lecture an absent-minded professor by mistake wrote on a board the formula  $\hat{\alpha} = \overline{Y_t / X_t}$  instead of  $\hat{\alpha} = \bar{Y} / \bar{X}$  ( $\overline{Y_t / X_t} = \frac{1}{T} \sum_{t=1}^T (Y_t / X_t)$ ). Does it change the conclusion that the estimator is unbiased? Show that population variance of  $\hat{\alpha}_2 = \overline{Y_t / X_t}$  is given by the expression  $\text{var}(\hat{\alpha}_2) = \frac{\sigma^2}{T^2} \sum_{t=1}^T (1 / X_t^2)$ . State the assumptions needed for your work.

**(d) [7 marks]** Which of the two estimators  $\hat{\alpha}_1 = \bar{Y} / \bar{X}$  and  $\hat{\alpha}_2 = \overline{Y_t / X_t}$  is more efficient? Explain what you understand by efficiency.

*Hint: use inequality for the four means (quadratic, arithmetic, geometric, harmonic):*

$$\sqrt{\frac{1}{T} \sum_{t=1}^T x_t^2} \geq \frac{\sum_{t=1}^T x_t}{T} \geq \sqrt[T]{x_1 \cdot x_2 \cdot \dots \cdot x_T} \geq \frac{T}{\left(\sum_{t=1}^T \frac{1}{x_t}\right)}$$

**Question 5. [30 marks]** A researcher has data on income  $Y$ , capital  $K$  and labor  $L$  indices for 1991-2017 related to the economy of some developing country (time series specific problems are out of consideration in this question). She estimated relationship between these variables using different production functions.

**(a) [8 marks]** First she runs a linear model

$$\hat{Y}_t = 10.83 + 0.22 \cdot K_t + 0.72 \cdot L_t \quad R^2 = 0.95 \quad (1)$$

(14.0) (0.025) (0.12)  $RSS = 3096.4$

Noticing that the sum of the slope coefficients is close to unity she decided to test the restriction  $\beta_K + \beta_L = 1$ , and decided to run the restricted regression

$$Y = \alpha + \theta \cdot K + (1 - \theta) \cdot L + u, \quad (2)$$

How equation (2) could be estimated using OLS (you may use results of this estimation  $R^2 = 0.85$  and  $RSS = 3137.9$ )? How to test the restriction? Is it possible to run F-test, based on comparison of determination coefficients of restricted and unrestricted regressions? Why or why not? What is economic meaning of this restriction (if any)? Would you recommend to accept the restriction?

**(b) [7 marks]** Then the researcher runs auxiliary regression (1a),

$$\hat{Y}_t = -51.36 + 0.48 \cdot K_t + 1.30 \cdot L_t - 0.003 \cdot (Y_t^*)^2 \quad R^2 = 0.96 \quad (1a)$$

(26.6) (0.10) (0.24) (0.001)  $RSS = 2375$

where  $(Y_t^*)^2$  are squared estimated values of  $\hat{Y}_t$  from the equation (1):  $\hat{Y}_t^* = 10.83 + 0.22 \cdot K_t + 0.72 \cdot L_t$

Comment on the meaning of researcher's actions, run the appropriate test and draw the conclusion. What implicit economic assumptions are used in the specification of equation (1)? To what extent can they be considered relevant to economic theory and practice? How to improve the specification of the equation (1)?

**(c) [8 marks]** Next she runs two logarithmic regressions

$$\log(Y) = 0.48 + 0.37 \cdot \log(K) + 0.53 \cdot \log(L) \quad R^2 = 0.97 \quad (3)$$

(0.39) (0.05) (0.12)  $RSS = 0.0763$

and

$$\log(Y/L) = -0.0050 + 0.33 \cdot \log(K/L) \quad R^2 = 0.82 \quad (4),$$

(0.0183) (0.03)  $RSS = 0.0811$

Give the interpretation of both models and their parameters. Prove that equation (4) is a restricted version of the equation (3). What is the restriction? Is it significant? What equation would you choose and why?

**(d) [7 marks]** Now the researcher introduces time ( $T$  equal to 1 in the first year of the period, 2 – in the second and so on), and estimates Cobb-Douglas (CD) production function using two different methods:

OLS applied to linearized model:

$$\log(Y) = 2.53 - 0.067 \cdot \log(K) + 0.51 \cdot \log(L) + 0.03 \cdot T \quad R^2 = 0.97 \quad (5-OLS)$$

(1.17) (0.24) (0.11) (0.016)  $RSS = 0.0664$

NLS (Non-Linear Least Squares) applied to original model (all estimates except capital elasticity are significant) :

$$\hat{Y} = e^{3.8} \cdot K^{-0.046} \cdot L^{0.086} \cdot e^{0.008T} \quad R^2 = 0.97 \quad (5-NLS)$$

Comment on parameters of both equations, paying special attention to the coefficients of variables  $K$  and  $T$ . Why estimates of equations (5-OLS) and (5-NLS) are different?

**GENERAL INSTRUCTIONS:** Start answering each question 1-5 on the blank with corresponding number (ask for extra paper if necessary). Structure your answers in accordance with the structure of the questions. Testing hypotheses always state clearly null and alternative hypotheses provide critical value used for test, mentioning degrees of freedom and the significance level chosen for the test.