

1) Box - transform

Box - Cox test

Zarembki transformation

$$\textcircled{1} \quad \lambda = \frac{y^\lambda - 1}{\lambda} \rightarrow \begin{cases} y-1 & \lambda = 1 \\ \rightarrow \log y & \lambda = 0 \end{cases}$$

$$y^* = \frac{y^{\lambda_1} - 1}{\lambda_1} \quad x^* = \frac{x^{\lambda_2} - 1}{\lambda_2}$$

lin

vs

$$y^* = \beta_0 + \beta_1 x^* + u$$

log-linear

1) estimate $\hat{\beta}_0$ $\hat{\beta}_1$ using OLS

2) $\hat{\lambda}_1$ and $\hat{\lambda}_2$ est. using NLS

t-test:

$$\lambda_1 = 1, \lambda_2 = 1$$

linear

$$\lambda_1 = 1, \lambda_2 = 0$$

lin-log mod

...

linear vs log-linear :

Zarembki Test.

1) Transformation using geom. means

$$\bar{y} = \sqrt[n]{\prod_{i=1}^n y_i} = \left(e^{\log \prod_{i=1}^n y_i} \right)^{1/n} = e^{\frac{1}{n} \cdot \sum \log y_i} = e^{\overline{\log y_i}}$$

$$y_i^* = \frac{y_i}{\bar{y}}$$

H_0 : no dif.

H_a : is a sig. dif

$$\begin{array}{lcl} 2) \text{ Est. } & y_i^* \mid X_i & \rightarrow \text{RSS}_1 \\ & \log y_i^* \mid X_i & \rightarrow \text{RSS}_2 \end{array}$$

$$\frac{n}{2} \cdot \left| \log \frac{\text{RSS}_1}{\text{RSS}_2} \right| \sim \chi_1^2$$

Problem 1. (UoL Exam). The rise in prices for public transport leads to lower corporate earnings, as people tend to choose cheaper alternatives. The student tries to find the best form of dependence of the volume of transportation T_i of some 50 transportation companies (in millions of dollars) from the prices of transportation P_i (in cents per one kilometer of transportation). She runs regressions (1-4) (linear, logarithmic and semi-logarithmic functions), she also runs two auxiliary regressions (5-6) performing Zarembka transformation (variable TZ_i is defined as $TZ_i = T_i / \sqrt{T_1 \cdot T_2 \cdot \dots \cdot T_n}$):

	(1)	(2)	(3)	(4)	(5)	(6)
Dependent variable	T_i	T_i	$\log(T_i)$	$\log(T_i)$	TZ_i	TZ_i
Independent variable\Constant	8.74	12.26	2.175	2.635	1.171	1.641
P_i	-0.339	-	-0.0045	-	-0.0045	-
$\log(P_i)$	-	-1.362	-	-0.179	-	-0.179
R^2	0.638	0.738	0.665	0.755	0.638	0.738
RSS	4.481	3.247	0.068	0.051	0.080	0.058

- (a) Explain the differences in the values of a slope coefficient in regression (1) and (4) giving interpretation to both regressions.
- (b) Explain the differences in the values of a slope coefficient in regression (2) and (3) giving interpretation to both regressions.
- (c) Explain using some math why your interpretation of regression (4) is correct using different methods. Do the same for regressions 2-3.
- (d) Which pairs of regression are comparable directly without Zarembka transformation). Which regressions becomes comparable after Zarembka transformation? Compare some regressions performing appropriate tests.

$$Y = p_0 + p_1 \ln x$$

$$\frac{dy_i}{dx_i} = \beta_1$$

$$T \uparrow p \cdot 100\%$$

$$y^* \ln x$$

$$\ln y \ln x$$

c) Algebraic approach

$$\log Y = a + b \cdot \log X$$

$$\log T = 2.6 - 0.2 \cdot \log P$$

$$X \uparrow 1\% \quad X \left(1 + \frac{1}{100}\right)$$

$$\log X + \log \left(1 + \frac{1}{100}\right)$$

$$\ln(1+x) \approx x$$

$$\log X + \frac{1}{100}$$

$$\log(Y + \Delta Y) = a + b \left(\log X + \frac{1}{100} \right)$$

$$\log(Y + \Delta Y) - \log Y = \frac{b}{100}$$

Problem 1. (UoL Exam). The rise in prices for public transport leads to lower corporate earnings, as people tend to choose cheaper alternatives. The student tries to find the best form of dependence of the volume of transportation T_i of some 50 transportation companies (in millions of dollars) from the prices of transportation P_i (in cents per one kilometer of transportation). She runs regressions (1-4) (linear, logarithmic and semi-logarithmic functions), she also runs two auxiliary regressions (5-6) performing Zarembka transformation (variable TZ_i is defined as $TZ_i = T_i / \sqrt{T_1 \cdot T_2 \cdot \dots \cdot T_n}$):

	(1)	(2)	(3)	(4)	(5)	(6)
Dependent variable	T_i	T_i	$\log(T_i)$	$\log(T_i)$	TZ_i	TZ_i
Independent variable/Constant	8.74	12.26	2.175	2.635	1.171	1.641
P_i	-0.339	-	-0.0045		-0.0045	
$\log(P_i)$	-	-1.362	-	-0.179	-	-0.179
R^2	0.638	0.738	0.665	0.755	0.638	0.738
RSS	4.481	3.247	0.068	0.051	0.080	0.058

- (a) Explain the differences in the values of a slope coefficient in regression (1) and (4) giving interpretation to both regressions.
- (b) Explain the differences in the values of a slope coefficient in regression (2) and (3) giving interpretation to both regressions.
- (c) Explain using some math why your interpretation of regression (4) is correct using different methods. Do the same for regressions 2-3.
- (d) Which pairs of regression are comparable directly without Zarembka transformation). Which regressions becomes comparable after Zarembka transformation? Compare some regressions performing appropriate tests.

Model (1) & (3) log-linear vs linear mod

$$\chi^2 = \frac{50}{2} \cdot \left| \log \frac{0.08}{0.068} \right| = 4.063$$

$$\chi^2_{1, 0.95} = 3.84$$

$\Rightarrow H_0$ is H_1 - dif \Rightarrow model (3)
log-linear

Model (2) & (4) lin-log vs log-log

$$\chi^2 = \frac{50}{2} \cdot \left| \log \frac{0.068}{0.051} \right| = 3.2$$

$\Rightarrow H_0$ is not rej \Rightarrow no sig. dif
between \Rightarrow choose both

An employee of a real estate agency in a Russian city with a developed subway network is interested in estimating the influence of the distance from the city center $CENTER_i$ (in kilometers) on the price of a two-room apartment in millions of rubles. Based on the data of 21 apartments sold during a period under consideration she runs a regression.

(1)

- (a) ☐ Is the regression coefficient significant (take into account that the realtor did not know exactly the sign of its coefficient before the regression calculation)?

→ ☐ Are the results of the estimation compatible with the hypothesis that true regression coefficient is positive?

□ Are the results of the estimation compatible with the hypothesis that true regression coefficient is 0.1?

□ How the conclusion on significance of the slope would change if the manager could use the assumption that the influence of the $CENTER_i$ on the apartment price is not positive?

□ Is intercept of the equation significant? Summarize all information on the test results and discuss economic meaning of the equation (1).

The realtor, not satisfied with the obtained result, decided to take into account the additional factor – the distance to the nearest subway station $|METRO|$ (also in kilometers).

(2)

During the discussion at the workshop, the realtor received advice from a colleague to use Ramsey's test for this equation. Since the realtor was not experienced enough in econometrics, a colleague helped her calculate appropriate equation (using in the right side of (3) estimated values \hat{PRICE}_i from equation (2):

(3)

Then the colleague helped her to estimate a new equation

(4

and did Ramsey's test again (using in the right side of (5) estimated values $-\log PRICE_i^{**}$ from equation (4):

(5)

- (b)** ☐ Help the realtor to understand the logic of her colleague in estimating these equations.

- Explain what the Ramsey test is, what is the null hypothesis and what statistics it uses; use them to perform the necessary calculations.

She estimated non-linear regression (4) using logarithm of dependent variable

(4)

and evaluates Ramsey test again

(5)

- What conclusions can be drawn from the results in this part of the study?

Reset Test (Ramsey test)

↳ functional form test

$$\text{Ho: } y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + \epsilon_i \quad (1)$$

H₀: functional (A) is wrong
(omitted var. or quadz
should be added)

$$\rightarrow y_i \mid x_1, \dots, x_k \Rightarrow \hat{y}_i$$

$$2) \quad y_i = \beta_1 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i} + \gamma_1 \hat{y}_i^2 + \dots + \gamma_p \hat{y}_i^{p+1} + \varepsilon_i$$

$$p=1 \quad \text{or} \quad p=2$$

Ramsey test : F -stat. $\gamma_1 = \dots = \gamma_p = 0$

$$F \sim F_{p, n-k-p}$$

The realtor, not satisfied with the obtained result, decided to take into account the additional factor – the distance to the nearest subway station METRO_{*i*} (also in kilometers).

$$\hat{PRICE}_i = 13.71 + 0.22 \cdot CENTER_i + 0.58 \cdot METRO_i \quad R^2 = 0.37 \quad (2)$$

(0.97) (0.09) (0.25) $RSS = 79.29$

During the discussion at the workshop, the realtor received advice from a colleague to use Ramsey's test for this equation. Since the realtor was not experienced enough in econometrics, a colleague helped her calculate appropriate equation (using in the right side of (3) estimated values \hat{PRICE}_i^* from equation (2):

$$\hat{PRICE}_i = 0.023 + 0.13 \cdot CENTER_i + 0.35 \cdot METRO_i + 0.07 \cdot (\hat{PRICE}_i^*)^2 \quad R^2 = 0.51 \quad (3)$$

(6.04) (0.18) (0.47) (0.033) $RSS = 60.64$

Then the colleague helped her to estimate a new equation

$$\log \hat{PRICE}_i = 2.62 + 0.019 \cdot CENTER_i + 0.059 \cdot METRO_i \quad R^2 = 0.32 \quad (4)$$

(0.10) (0.0095) (0.026) $RSS = 0.8448$

and did Ramsey's test again (using in the right side of (5) estimated values $\log \hat{PRICE}_i^{**}$ from equation (4):

$$\log \hat{PRICE}_i = 0.62 + 0.030 \cdot CENTER_i + 0.084 \cdot METRO_i + 0.012 \cdot (\log \hat{PRICE}_i^{**})^2 \quad R^2 = 0.39 \quad (5)$$

(1.53) (0.039) (0.11) (0.0088) $RSS = 0.7672$

(b) Help the realtor to understand the logic of her colleague in estimating these equations.

□ Explain what the Ramsey test is, what is the null hypothesis and what statistics it uses; use them to perform the necessary calculations.

She estimated non-linear regression (4) using logarithm of dependent variable

$$\log \hat{PRICE}_i = 2.62 + 0.019 \cdot CENTER_i + 0.059 \cdot METRO_i \quad R^2 = 0.32 \quad (4)$$

(0.10) (0.0095) (0.026) $RSS = 0.8448$

and evaluates Ramsey test again

$$\log \hat{PRICE}_i = 0.62 + 0.030 \cdot CENTER_i + 0.084 \cdot METRO_i + 0.012 \cdot (\log \hat{PRICE}_i^{**})^2 \quad R^2 = 0.39 \quad (5)$$

(1.53) (0.039) (0.11) (0.0088) $RSS = 0.7672$

□ What conclusions can be drawn from the results in this part of the study?

$$t\text{-test} \quad t = \frac{0.07}{0.03} =$$

$$= 2.12$$

$$t^{97.5} = 2.1$$

\Rightarrow some non-lin. is present

Question 1. (17 marks)

Two students developed different models for their course papers: student (a) was convinced that the coefficient of variable Z_i is equal to 1, while student (b) believed that this coefficient is the opposite of the coefficient for the variable X_i

$$Y_i = \beta_1 + \beta_2 X_i + Z_i + u_i \quad (a)$$

$$Y_i = \beta_1 + \beta_2 X_i - \beta_2 Z_i + v_i \quad (b)$$

They think that in fact both models will give the same estimates of β_2 , but both faced unexpected difficulties when trying to evaluate these models (what difficulties?). So they decided instead to calculate sample variances and covariances on the base of the same data (n observations): $\text{Var}(Y) = 4$, $\text{Var}(X) = 3$, $\text{Var}(Z) = 5$, $\text{Cov}(Y, X) = 6$, $\text{Cov}(Y, Z) = 1$, $\text{Cov}(X, Z) = 2$.

- (a) ☐ Help the students to find the least squares estimates of β_2 for their models, indicating all necessary steps.
☐ Are these estimates really the same as students think?

1) (a) $\beta_3 = 1$ $\beta_2 = -\beta_2$

$$Y_i - Z_i = \beta_1 + \beta_2 X_i + u_i \Rightarrow \hat{\beta}_2 = \frac{\text{Cov}(X_i, Y_i - Z_i)}{\text{Var}(X_i)}$$

(b) $Y_i = \beta_1 + \beta_2 (X_i - Z_i) + v_i$ RSS_R

(a) $\hat{\beta}_2 = \frac{\text{Cov}(X, Y - Z)}{\text{Var}(X)} = \frac{\text{Cov}(X, Y) - \text{Cov}(X, Z)}{\text{Var}(X)} =$

$$(6 - 2) / 3 = 4/3$$

(b) $\hat{\beta}_2 = \frac{\text{Cov}(X - Z, Y)}{\text{Var}(X - Z)} = \frac{6 - 1}{3 + 5 - 2 \cdot 2} = \frac{5}{4}$

(b) The scientific advisor told the students that both their models are restricted versions of the more general model

UR: $Y_i = \beta_1 + \beta_2 X_i + \beta_3 Z_i + w_i$ (c) $R: Y_i = \beta_1 + \beta_2 X_i + Z_i + u_i$ (a)
 $R: Y_i = \beta_1 + \beta_2 X_i - \beta_2 Z_i + v_i$ (b)

☐ What are restrictions in each case?

☐ How to test the restriction for the model (a)? Indicate necessary steps.

☐ What model should be chosen if both restrictions are invalid? What model (a), (b) or (c) should be chosen if only one restriction is valid?

☐ Let both restrictions be valid. Which model out of (a) and (b) is preferable and why? Use numerical data above to choose between (a) and (b).

(a) $\beta_3 = 1$ (b) $\beta_3 = -\beta_2$

t-test $t = \frac{\hat{\beta}_3 - 1}{\text{se}(\hat{\beta}_3)}$

F-test

$$F = \frac{(RSS_R - RSS_{UR})/1}{RSS_{UR}/(n-3)}$$

$$(a) \quad \sigma^2_{\hat{\beta}_2} = \frac{\sigma_u^2}{n \cdot \text{Var}(X)}$$

$$(b) \quad \sigma^2_{\hat{\beta}_2} = \frac{\sigma_v^2}{n \cdot \text{Var}(X^*)} = \frac{\sigma_v^2}{n \cdot \text{Var}(X-2)}$$

$$\text{Var}(X) = 3$$

$$\text{Var}(X-2) = 3 + 5 - 2 \cdot 2 = 4$$

\Rightarrow given that $\sigma_u^2 = \sigma_v^2$ we prefer (b)

(2)

$$y_i = \beta \cdot x_i + u_i$$

$$TSS = ESS + ESS$$

Without
const

$$\Rightarrow R^2_{uc} = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2}$$

With
const

$$R^2 = 1 - \frac{\sum \hat{u}_i^2}{\sum (y_i - \bar{y})^2}$$

"TSS"

$$k \uparrow \Rightarrow R^2 \uparrow$$

$$- R^2_{adj} \downarrow \text{ if }$$

new var. is insignif.

$$- R^2_{adj} \uparrow \text{ if }$$

t-stat for new var.

is greater than 1

$$- R^2_{adj} \in (-\infty; 1)$$

Question 5. [25 marks] A lazy student found an interesting paper with econometric models describing how *COVID* - aggregate anti COVID expenditures, depend on *GDP*, aggregate gross national product, and *POP*, total population, for a sample of 70 countries in the second quarter of 2020. *COVID* and *GDP* are both measured in US\$ billion. *POP* is measured in million. *RSS* - Residual Sum of Squares. He decided to use this article in his course paper pretending that he got all equations himself using original data (what in fact was not true). He wrote the equations on a paper:

$$\rightarrow \log \frac{\widehat{COVID}}{POP} = -3.74 + 1.27 \log \frac{GDP}{POP} \quad R^2 = 0.90 \quad \underline{RSS = 15.45} \quad (1)$$

$$\log \widehat{COVID} = -3.60 + 1.27 \log GDP - 0.33 \log POP \quad R^2 = 0.95 \quad \underline{RSS = 13.90} \quad (2)$$

$$\rightarrow \log \frac{\widehat{COVID}}{POP} = -3.60 + 1.27 \log \frac{GDP}{POP} - 0.06 \log POP \quad R^2 = 0.91 \quad \underline{RSS = 13.90} \quad (3)$$

But when he wrote his coursework some details seemed a little strange to him and he began to doubt that he had correctly rewritten the equations on paper. He asked your advice.

(a) ☐ The student now believes that by mistake he repeated the same coefficient 1.27 in equations (1), (2) and (3), as well as he repeated the intercept - 3.60 in equations (2) and (3), but he does not remember the correct values. Are these coincidences really happened by mistake?

$$(1) \quad \log \frac{Cov}{Pop} = \beta_0 + \beta_1 \cdot \log \frac{GDP}{Pop}$$

$$\log Cov = \beta_0 + \beta_1 \log GDP + (1 - \beta_1) \log Pop$$

$$(2) \quad \log Cov = \beta_0 + \beta_1 \cdot \log GDP + \beta_2 \cdot \log Pop$$

$$(3) \quad \log Cov = \beta_0 + \beta_1 \cdot \log GDP + (1 - \beta_1 + \beta_2) \log Pop$$

☐ It seems strange to him that both coefficients of the variable $\log POP$ in equations (2) and (3) are negative but different in absolute value. Help the student to understand these.

$$(2) \quad \beta_2 = 1 - \beta_1 + \beta_2$$

$$\beta_2 = \beta_2 + \beta_1 - 1$$