

**Instructions**

Handwritten:  
6304641472

- (1) Please read the instruction carefully. Also take this habit with you into the exam room.
- (2) Please read each question carefully and answer the questions straightforwardly. Always provide economic reasons at least a paragraph for your analysis, or a graph when necessary, even when the question does not indicate so.
- (3) Handing and submitting assignments are only available via BE Moodle.

**Answering the questions and preparing answer sheets**

- (1) Answers are to be handwritten, in either digital or analog form, in a blank canvas or any clean paper. Make sure that your handwriting is clearly visible and readable.
- (2) There is no need to rewrite the question. Just indicate the question number clearly for each of the answer, such as 1.a).
- (3) Default decimal point is 4.
- (4) Choose precise wordings, especially when you want to interpret the meaning of a test, confidence interval, or coefficients.
- (5) When done, for the digital case, collage all the pages into a single PDF file. For those who write on sheets of paper, take photo of all pages then convert all of them into a single PDF file as well.
- (6) Name your PDF file as StudentID\_YourNickname, such as 640123456\_Bo.

**Submitting your answers**

- (1) Make sure your file does not exceed 10MB. This is the maximum file size for BE Moodle upload.
- (2) Login to BE Moodle, head into the course, then the assignment topic.
- (3) Choose your file to submit. Done. There will be timestamp for your upload date and time, so please make sure to not submit later than that.

**For all questions, answer up to 4 decimal places**

**Question 1. (15 points)** Given this information

$$\begin{aligned}
 n &= 18 & \sum_{i=1}^n X_i &= 388.00 & \sum_{i=1}^n Y_i &= 50.90 \\
 \sum_{i=1}^n (X_i)^2 &= 9,620.00 & \sum_{i=1}^n X_i Y_i &= 1,254.90 \\
 \sum_{i=1}^n (X_i - \bar{X})^2 &= 211.00 & \sum_{i=1}^n (Y_i - \bar{Y})^2 &= 2.5844 \\
 \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) &= 20.58 & \sum_{i=1}^n \hat{u}_i^2 &= 0.5781
 \end{aligned}$$

Use the above sample information to answer all the following questions. Show explicitly all formulas and calculations.

- From regression model:  $Y_i = \beta_1 + \beta_2 X_i + u_i$ ,  $u_i \sim NIID(0, \sigma^2)$ , **find the estimators** of  $\beta_1$  and  $\beta_2$  with OLS method. Interpret the intercept and slope coefficients.
- Compute the value of  $R^2$  and explain its meaning.
- If  $X_i = 30$ , estimate the value of  $\hat{Y}_i$  and explain its meaning.
- Calculate the estimators of  $\text{var}(u_i)$ ,  $\text{var}(\hat{\beta}_1)$  and  $\text{var}(\hat{\beta}_2)$ .
- What are the 90-percent confident intervals for  $\beta_2$ ? Interpret the meaning.
- Test the hypothesis whether the slope coefficients are different from zero at 0.05 level of significance.

**Question 2.** Using the 2015 Health and Welfare Survey from the National Statistical Office, a simple linear regression is modeled as follows,

$$outp_i = \beta_1 + \beta_2 age_i + u_i$$

where  $outp_i$  is how many times person  $i$  has visited hospital in 2015, from 0 to 7 times  
 $age_i$  is how old is person  $i$ , from 0 to 97 years.

We assume that both  $outp_i$  and  $age_i$  are continuous, the estimation results in the following table. Answer the following questions and show your work.

Source	SS	df	MS	Number of obs	=	27,886
Model	77.5444409	1	77.5444409	F(1, 27884)	=	186.96
Residual	11565.0627	27,884	.414756231	Prob > F	=	0.0000
				R-squared	=	0.0067
				Adj R-squared	=	0.0066
Total	11642.6072	27,885	.417522223	Root MSE	=	.64402

outp	Coefficient	Std. err.	t	P> t	[95% conf. interval]
age	.0031338	.0002292			.0026846 .003583
_cons	.4279898	.0140339			.4004828 .4554969

- Test if both parameters are significantly different from zero or not. Use  $\alpha = 0.05$ .
- Interpret the meaning of  $\hat{\beta}_2$ . Does the sign of  $\hat{\beta}_2$  make economic sense? Explain.
- If  $outp_i$  is turned into natural logarithmic scale (ln), how would you reinterpret the relationship between  $\hat{\beta}_2$  and  $\widehat{outp}_i$ , assumed that the given coefficient given in the table above can be used to interpret this new functional form.
- If  $age_i$  variable is divided by 10, how does it affect both the coefficients, standard errors, and confidence intervals? Answer the changes of both the constant and slope (if there is).
- Find the confidence interval of mean prediction at the age of 50 years old, given that  $var(\hat{Y}_0) = 0.00002$  and  $\alpha = 0.01$ .

**Question 3.** Discuss in a short paragraph why the confidence interval for both the mean prediction and individual prediction get larger as the  $X_0$  is further away from  $\bar{X}$ .

-----

Question 1. (15 points) Given this information

$$\begin{aligned}
 n &= 18 & \sum_{i=1}^n X_i &= 388.00 & \sum_{i=1}^n Y_i &= 50.90 \\
 \sum_{i=1}^n (X_i)^2 &= 9,620.00 & \sum_{i=1}^n X_i Y_i &= 1,254.90 \\
 \sum_{i=1}^n (X_i - \bar{X})^2 &= 211.00 & \sum_{i=1}^n (Y_i - \bar{Y})^2 &= 2.5844 \\
 \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) &= 20.58 & \sum_{i=1}^n \hat{u}_i^2 &= 0.5781
 \end{aligned}$$

Use the above sample information to answer all the following questions. Show explicitly all formulas and calculations.

a) From regression model:  $Y_i = \beta_1 + \beta_2 X_i + u_i$ ,  $u_i \sim NID(0, \sigma^2)$ , find the estimators of  $\beta_1$  and  $\beta_2$  with OLS method. Interpret the intercept and slope coefficients.

$$\hat{\beta}_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{20.58}{211} = 0.0975$$

$$\bar{y} = \hat{\beta}_1 + \hat{\beta}_2 \bar{x} \rightarrow \hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x} \rightarrow \hat{\beta}_1 = 2.8278 - (0.0975)(21.5556) = 0.7261$$

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n} = \frac{50.90}{18} = 2.8278$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{388}{18} = 21.5556$$

(expect)

Thus,  $[\hat{\beta}_1 = 0.7261]$  when  $x = 0$ ,  $y = 0.7261$   
 $[\hat{\beta}_2 = 0.0975]$   $x_i \uparrow 1$  unit,  $y \uparrow$  on average by 0.0978 unit.

b) Compute the value of  $R^2$  and explain its meaning.

c) If  $X_i = 30$ , estimate the value of  $\hat{Y}_i$  and explain its meaning.

$$b) R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^n \hat{u}_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{0.5781}{2.5844} = 0.7763$$

$R^2 \rightarrow 0.7763$ :  $x$  variable can explain 77.63% of variation in  $y$

$$\begin{aligned}
 c) \hat{y}_i &= \hat{\beta}_1 + \hat{\beta}_2 x_i \\
 &= 0.7261 + 0.0975(30) \\
 \hat{y}_i &= 3.6511 \text{ on average when } x_i = 30
 \end{aligned}$$

d) Calculate the estimators of  $\text{var}(u_i)$ ,  $\text{var}(\hat{\beta}_1)$  and  $\text{var}(\hat{\beta}_2)$ .

e) What are the 90-percent confident intervals for  $\beta_2$ ? Interpret the meaning.

f) Test the hypothesis whether the slope coefficients are different from zero at 0.05 level of significance.

$$0.06017109$$

$$d) \text{var}(u_i) = \sigma^2 = \frac{\sum \hat{u}_i^2}{n-k} = \frac{0.5781}{16} = 0.0361$$

$$\text{var}(\hat{\beta}_1) = \frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2} \cdot \sigma^2 = \frac{9620}{18(211)} \cdot 0.0361 = 0.0914$$

$$\text{var}(\hat{\beta}_2) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{0.0361}{211} = 0.0002$$

$$e) \text{se}(\hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_2)} = \sqrt{0.0002} = 0.0141$$

$$\alpha = 0.1 \rightarrow t_{16, 0.05} = 1.746 : \hat{\beta}_2 = 0.0975$$

hence,  $P(\hat{\beta}_2 - t_{16, 0.05} \text{se}(\hat{\beta}_2) \leq \beta_2 \leq \hat{\beta}_2 + t_{16, 0.05} \text{se}(\hat{\beta}_2)) = 1 - \alpha$

$$P(0.0975 - (1.746 \cdot 0.0141) \leq \beta_2 \leq 0.0975 + (1.746 \cdot 0.0141)) = 0.9$$

$$\hookrightarrow P(0.0729 \leq \beta_2 \leq 0.1221) = 0.9$$

Thus,  $\beta_2$  will be  $0.0729 \leq \beta_2 \leq 0.1221$

f) Test the hypothesis whether the slope coefficients are different from zero at 0.05 level of significance.

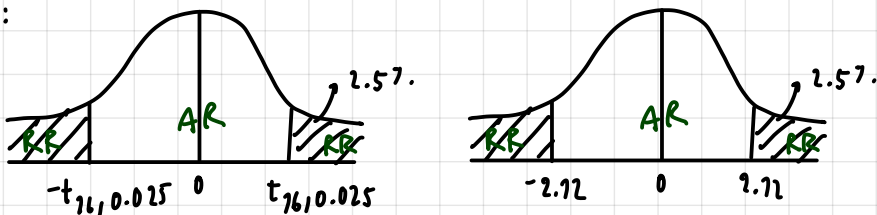
Step 1:  $H_0 : \beta_2 = 0$

$H_1 : \beta_2 \neq 0$

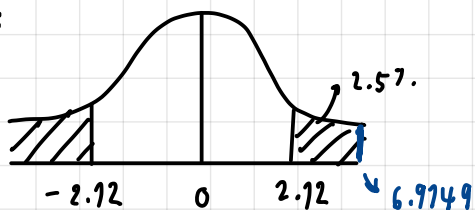
Step 2:  $\alpha = 0.05$

Step 3:  $t = \frac{0.0975 - 0}{0.0149} = 6.9149$

Step 4:



Step 5:



since  $t_{(9)} > t_{76, 0.025}$ , reject the null hypothesis and can conclude that with 95% confident interval, the  $\beta_2 \neq 0$

Question 2. Using the 2015 Health and Welfare Survey from the National Statistical Office, a simple linear regression is modeled as follows,

$$outp_i = \beta_1 + \beta_2 age_i + u_i$$

where  $outp_i$  is how many times person  $i$  has visited hospital in 2015, from 0 to 7 times  
 $age_i$  is how old is person  $i$ , from 0 to 97 years.

We assume that both  $outp_i$  and  $age_i$  are continuous, the estimation results in the following table. Answer the following questions and show your work.

Source	SS	df	MS	Number of obs	=	27,886
Model	77.5444409	1	77.5444409	F(1, 27884)	=	186.96
Residual	11565.0627	27,884	.414756231	Prob > F	=	0.0000
				R-squared	=	0.0067
				Adj R-squared	=	0.0066
				Root MSE	=	.64402

outp	Coefficient	Std. err.	t	P> t	[95% conf. interval]
age	.0031338	.0002292		.0026846	.003583
_cons	.4279898	.0140339	Omitted	.4004828	.4554969

- Test if both parameters are significantly different from zero or not. Use  $\alpha = 0.05$ .
- Interpret the meaning of  $\hat{\beta}_2$ . Does the sign of  $\hat{\beta}_2$  make economic sense? Explain.

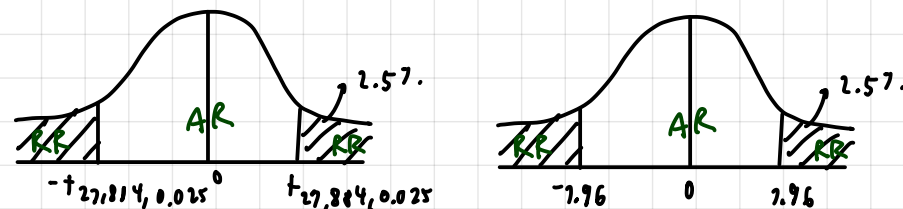
Step 1:  $H_0 : \beta_j = 0$  ( $\beta_j = \text{constant}$ )

$H_1 : \beta_j \neq 0$

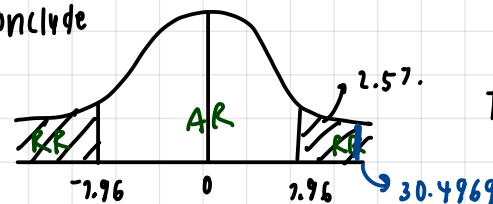
Step 2:  $\alpha = 0.05$

Step 3:  $t_{(9)} = \frac{0.4279898 - 0}{0.0140339} = 30.4969$

Step 4:



Step 5: conclude



$t_{(9)}$  falls in rejection region (RR)  
 Thus, we can reject null hypothesis  
 so, with 95% confident interval  
 $\beta_j$  (constant)  $\neq 0$

$$\beta_2 = \text{age}$$

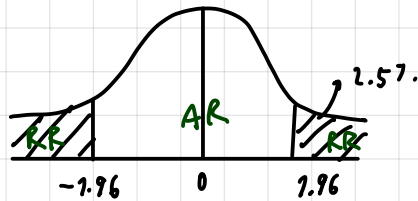
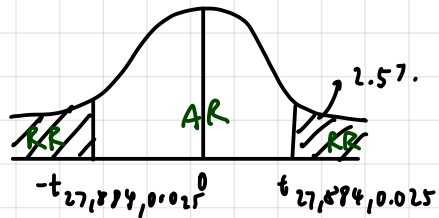
Step 1 =  $H_0: \beta_2 = 0$   
 $H_1: \beta_2 \neq 0$

Step 2 =  $\alpha = 0.05$

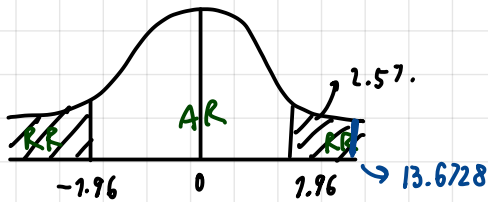
Step 3 =

$$t_{CG} = \frac{0.0031338 - 0}{0.0002292} = 13.6728$$

Step 4 =



Step 5



$t_{CG}$  falls in RR  $\rightarrow$  we can reject null hypothesis.  
 so, with 95% confidence interval,  $\beta_2(\text{age}) \neq 0$

- b) Interpret the meaning of  $\hat{\beta}_2$ . Does the sign of  $\hat{\beta}_2$  make economic sense? Explain.  
 c) If  $\text{out}_i$  is turned into natural logarithmic scale (ln), how would you reinterpret the relationship between  $\hat{\beta}_2$  and  $\widehat{\text{out}}_i$ , assumed that the given coefficient given in the table above can be used to interpret this new functional form.

b) when the age of people  $\uparrow$  1 yr, we can expect that the visit / yr will increase 0.0031 time on average.  
 This make economic sense! because we trend to use more medical service as we get older.

c)  $\ln(\widehat{\text{out}}_i) = \hat{\beta}_1 + \hat{\beta}_2(\text{age})$   
 $\ln(y_i) = \beta_1 + \beta_2 x_i$

$$\rightarrow \frac{d \ln(y_i)}{dx} = \hat{\beta}_2$$

$$\frac{dy/y}{dx} = \hat{\beta}_2$$

$$\frac{dy}{dx} \cdot \frac{1}{y} = \hat{\beta}_2 \rightarrow \frac{dy}{dx} = y \hat{\beta}_2 \rightarrow \text{slope of the line}$$

$$\frac{dy}{dx} \cdot \frac{x}{y} = x \hat{\beta}_2 \rightarrow \text{elasticity}$$

when age increase by 1 year, visiting hospital by people will increase by  $\hat{\beta}_2 \cdot 100 = 0.3133\%$ .

- d) If  $\text{age}_i$  variable is divided by 10, how does it affect both the coefficients, standard errors, and confidence intervals? Answer the changes of both the constant and slope (if there is).  
 e) Find the confidence interval of mean prediction at the age of 50 years old, given that  $\text{var}(\hat{y}_0) = 0.00002$  and  $\alpha = 0.01$ .

d) constant  $\hat{\beta}_1$ : no change  $\rightarrow y_i$  does not change, and thus if  $x = 0$ , the intercept still the same variable  $\hat{\beta}_2$ : if  $\hat{\beta}_2$ ,  $\text{se}_{\hat{\beta}_2}$ ,  $\text{CI}_{\hat{\beta}_2}$  change, they will be multiply by 10.

as the data is divided, the x axis 'scale' will then decrease [slope steeper]

new one:  $\hat{\beta}_2 = 0.031338$   
 $\text{se}_{\hat{\beta}_2} = 0.002292$   
 $\text{CI} = 0.0268465 \leq \beta_2 \leq 6.03583$

$$e) \text{ se } \hat{y}_0 = \sqrt{\text{Var}(\hat{y}_0)} = \sqrt{0.0002} = \underline{4.472135955 \times 10^{-3}}$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$\hat{y}_1 = 0.4279898 + 0.0031338(50) = 0.5846798$$

$$P(\hat{y}_{50} - t_{218, 0.005} \cdot \text{se } \hat{y}_0 \leq y_{50} \leq \hat{y}_{50} + t_{218, 0.005} \cdot \text{se } \hat{y}_0) = 1 - \alpha$$

$$P(0.5846798 - 2.576(4.4721 \times 10^{-3}) \leq y_{50} \leq 0.5846798 + 2.576(4.4721 \times 10^{-3}) = 0.99$$

$$\text{so } \Rightarrow \text{CI} = (0.5732 \leq y_{50} \leq 0.5962) = 0.99$$

**Question 3.** Discuss in a short paragraph why the confidence interval for both the mean prediction and individual prediction get larger as the  $X_0$  is further away from  $\bar{X}$ .

mean prediction :  $\text{Var}(\hat{y}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$

$$Pr \left[ \hat{y}_0 - \left( t_{\frac{\alpha}{2}} \cdot \text{se } \hat{y}_0 \right) \leq y \leq \hat{y}_0 + \left( t_{\frac{\alpha}{2}} \cdot \text{se } \hat{y}_0 \right) \right] = 1 - \alpha$$

on the other hand, individual prediction :  $\text{Var}(f_e) = \text{Var}(\hat{y}_0 - y_0) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$

$$Pr \left[ \hat{y}_0 - \left( t_{\frac{\alpha}{2}} \cdot \text{se } f_e \right) \leq y_0 \leq \hat{y}_0 + \left( t_{\frac{\alpha}{2}} \cdot \text{se } f_e \right) \right] = 1 - \alpha$$

as  $x_0$  is further away from  $\bar{x}$  [larger in  $x_0 - \bar{x}$ ], the  $\text{se}(\hat{y}_0)$  will be larger because the variance get larger.

→ when  $x_0$  move further away from  $\bar{x}$ , the lesser data points of  $x_i$ . Therefore the CI must get larger to tackle with the unknown.