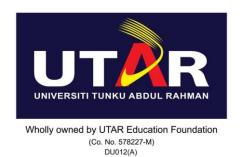
UNIVERSITI TUNKU ABDUL RAHMAN ACADEMIC YEAR 2019/2020 JANUARY 2020 TRIMESTER FINAL ASSESSMENT

ANSWER SCRIPT

Candidate is required to fill in ALL the information below:

Name: (as stated in Student Identity Card)	Ngu Yi Hui		
Faculty /Institute/ Centre:	FSc	Programme:	Statistical Computing And Operations Research
Index No. (in numbers):	A00082DBSCF	Index No. (in words):	A Zero Zero Eight Two DBSCF
Course Code :	UDPS2223	Course Description:	Applied Regression Analysis
Submission Date :	5 th MAY 2020	Submission Time :	9:30am - 12pm

	FOR EXAMINER'S USE ONLY		
QUESTION NUMBER	MARKS		
	Internal	External	
Q1			
Q2			
Q3			
Q4			
TOTAL MARKS			



DECLARATION STATEMENT

I, Ngu Yi Hui (Name), Student ID No. 18ADB01438, hereby solemnly and fully declare and confirm that during my programme of study at Universiti Tunku Abdul Rahman, I shall abide and comply with all the rules, regulations and lawful instructions of Universiti Tunku Abdul Rahman and endeavour at all times to uphold the good name of the University.

I hereby declare that my submission for this Final Assessment is based on my original work, not plagiarised from any source(s) except for citations and quotations which have been duly acknowledged. I am fully aware that students who are suspected of violating this pledge are liable to be referred to the Examination Disciplinary Committee of the University.

Programme:	Statistical Computing And Operations Research
(Digital) Signature:	HUI
Student's I.C / Passport No.:	991110-14-6378
Index No:	A00082DBSCF
Date of Submission:	5 th MAY 2020

Course Code: UDPS2223

Index Number (in figure): A00082DBSCF

Page: 3

QI.

(a) Define
$$x = amount$$
 of protein intake per day (in grams)

 $y = diastolic$ blood pressure.

 $\Sigma \times = 762$ $\Sigma y = 66.7$ $\overline{y} = 7.411$
 $\Sigma \times^2 = 64.868$ $\Sigma \times y = 5758.2$ $\overline{\times} = 84.667$
 $S_{XX} = \overline{\Sigma} \times^2 - \frac{(\Sigma \times)^2}{n} = 64.868 - \frac{76.2^2}{9} = 352$
 $S_{XY} = \Sigma \times y - \frac{(\Sigma \times)(\Sigma y)}{n} = 5.758.2 - \frac{(762)(66.7)}{9} = 110.933$
 $\hat{\beta}_1 = \frac{S_{XY}}{S_{XX}} = \frac{110.933}{352} = 0.315$
 $\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \, \overline{x} = 7.411 - 0.315 \left(84.667 \right) = -19.259$
 $\therefore \text{ Linear Regression : } \hat{y} = -19.259 + 0.315 \times$

Obtained
$$ZXY = 255274$$
, $ZX^2 = 125068$

$$\hat{\beta} = \frac{255274}{125068} = 2.041 = a,$$

For example, we fit a regression line which the dependent variable is height of children and the independent variable is age of children. When the age of children is equal to zero, his height should not be a negative value and should be zero value. Thus, we will fit a regression line pass through origin which the model equation y= a,x is adequate.

Index Number (in figure): A00082DBSCF Course Code: UDPS2223 Page: 4

Q2.

(ii)
$$R^2 = \frac{SSR}{TSS} = \frac{4283062.960}{479.9789.5} = 0.8923$$

$$R_a^2 = \left[-\frac{n-1}{n-k-1} \left(1 - R^2 \right) \right] = \left[-\frac{31}{29} \left(1 - 0.8923 \right) \right] = 0.8849$$

R2 = 0.8923, shows that 69.23% of the variation of the auction price (4) is interpreted by number of bidders (x.) and age of clocks (xz) in the model.

(iii)
$$H_0: \beta_2 = 0$$
 $H_i: \beta_2 \neq 0$

$$+ \frac{\hat{\beta}_2 - 0}{se(\hat{\beta}_2)} = \frac{12.741}{0.905} = 14.078$$

+0.025,29 = 2.045

Since 1+*1 > to.025,29, reject Ho.

- :. At a = 0.05, there is significant evidence to conclude that the age of clock (x.) is related to auction price (y) in the model, when number of bidders (x,) is already in the model.
- (b) The coefficient estimators should be standardized, because the sample size is small (n=17) and the unit of x, seems larger than x,. Besides that, since the degree of freedom (df) of sum square error (SSE) is 15, means that the number of variable in the model is only one (17-15-1) Thus, we will choose a better variable between the two variables which explained the more and is significant to the dependent variable.

Index Number (in figure): A00082DBSCF Course Code: UDPS2223 Page: 5

Q3

(a) Procedure:

- @ Initially, there is no regressor variable in the model except intercept.
- @ We calculate the t-statistics and p-values for all the independent variables. Then select the variable that has the largest t-statistics, which is x3 in this case. We compared its t-value with to, and found that It = 7.23 > tn = 1.677. Therefore, we add x3 into our model.
- @ Next, we calculate the t-statistics and p-values for the rest independent variables again. XI has the largest t-statistics among the others, and its t-value is larger than tin, given that x3 is already in the model. H=7.25 > fin=1.678. Thus, we add x, into our model.
- @ Furthermore, we calculate the t-statistics and p-values for the rest independent variables again. X2 has the largest t-statistics among the others, and its t-value is higher than tim, given that x3 and x, already in the model. It = 6.58 > tin = 1.679. Hence, we add x2 into the model.
- B Then, we calculate the t-statistics and p-values for the rest independent variables again. We found that both t-values of x4 and x5 are smaller than tin = 1.679. So, x4 and xs are not allowed to be added into the model. We stop the procedure and obtain the final model.
- :. Final model: $\hat{y} = 2.331 + 3.288 \times 1 + 3.251 \times 1 + 1.038 \times 3$

Course Code: UDPS2223 Index Number (in figure): A00082DBSCF Page: 6

Q3.

(b) (i) Model A:
$$\hat{y} = a + bx$$
,
Model B: $\hat{y} = a + bx$, $+ cx$, $+ dx$,

Firstly, since Model B includes more predictors than Model A, it will definitely explain more variation of y compared with Model A.

Furthermore,
$$R^2 = \frac{SSR}{TSS} = 1 - \frac{SSE}{TSS}$$

Since SSE for Model B minus more contents than Model A, the SSEModel® smaller than SSE model A. And cause the RA & RB.

(ii) If predictors x2 and x3 are having a negative relationship with y and the predictor x, is having a positive relationship with y, the b in the model B can be increased.

Index Number (in figure): A00082DBSCF

Course Code: UDPS2223

Page: 7

Q4.

(a)
$$\Sigma h_{ii} = 4$$

 $h_{44} = 4 - (0.231 + 0.235 + 0.711 + 0.173 + 0.4 + 0.882 + 0.213 + 0.499 + 0.428)$
= 0.228

Determine leverage points:

$$cut-off$$
 pant = $\frac{2p}{n} = \frac{2(4)}{10} = 0.8$
 $h_{77} > 0.8$
 $Point 7$ is a leverage point.

Determine influence points:

Point 1: P(F(4,6) < 0.012) = 0.00037

Point 2: P(F(4,6) < 0.012) = 0.00037

Paint 3: P(F(4,6) < 0.255) = 0.10347 (>0.1)

Point 4: P(F(4,6) < 0.011) = 0.00031

Point 5: P(F(4,6) < 0.008) = 0.00017

Paint 6 : P(F(4,6) & 0.033) = 0.0027

Point 7: P(F(4,6) & 1.9) = 0.77017 (>0.5)

Point 8: P(F(4,6) & 0.01) = 0.00026

Point 9: P(F(4.6) < 0.06)= 0.00843

Point 10: P(F(4.6) < 0.039) = 0.00372

Point 3 is an influence point with small influence.

Point 7 is an influence point with major influence.

(b) No, because the plot shows not much problem on linearity, normality and constant variance assumptions.

