sta302midterm-10101-crowdmark-assessm

#1 1 of 7



Colution.

## STA302/1001 Section L0101- Midterm Exam

June 6, 2017, 2:10pm- 3:40pm at EX200

															_	UVV	VVI V	
U of T e-mail:	·	@mail.utoronto.ca																
Surname (Last name):																		
Given name (First name):									=									
Student ID:																		
UTORID: (e.g. lihao8)																		

## Instructions:

- You have 90 minutes for 3 questions with multiple parts. Keep these papers closed on your desk until the start of the test is announced.
- Use a benchmark of  $\alpha = 5\%$  for all inference, unless otherwise indicated
- You may use a calculator. For numerical answer, please round it off to 4 decimal digits.
- Full mark: 50. Total pages (include the cover): 7.
- Write your answers in the given space only. You cannot use blank space for other questions nor can you write answers on the back. Your entire answer must fit in the designated space provided immediately after each question.

 $b_{1} = \frac{\sum_{i=1}^{n} (X_{i} - \overline{X}) (Y_{i} - \overline{Y})}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}} = \frac{\sum_{i=1}^{n} X_{i} Y_{i} - n \overline{X} \overline{Y}}{\sum_{i=1}^{n} X_{i}^{2} - n \overline{X}^{2}} \qquad b_{0} = \overline{Y} - b_{1} \overline{X}$   $Var \{b_{1}\} = \frac{\sigma^{2}}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}} \qquad Var \{b_{0}\} = \sigma^{2} \left(\frac{1}{n} + \frac{\overline{X}^{2}}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}}\right)$   $Cov \{b_{0}, b_{1}\} = -\frac{\sigma^{2} \overline{X}}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}} \qquad SSTO = \sum_{i=1}^{n} (Y_{i} - \overline{Y})^{2}$   $SSE = \sum_{i=1}^{n} (Y_{i} - \hat{Y}_{i})^{2} \qquad SSR = \sum_{i=1}^{n} (\hat{Y}_{i} - \overline{Y})^{2} = b_{1}^{2} \sum_{i=1}^{n} (X_{i} - \overline{X})^{2}$   $\sigma^{2} \{\hat{Y}_{h}\} = Var \{\hat{Y}_{h}\} = \sigma^{2} \left(\frac{1}{n} + \frac{(X_{h} - \overline{X})^{2}}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}}\right) \qquad r = \frac{\sum_{i=1}^{n} (X_{i} - \overline{X}) (Y_{i} - \overline{Y})}{\sqrt{\left[\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}\right]\left[\sum_{i=1}^{n} (Y_{i} - \overline{Y})^{2}\right]}}$   $\sigma^{2} \{pred\} = Var \{Y_{h} - \hat{Y}_{h}\} = \sigma^{2} \left(1 + \frac{1}{n} + \frac{(X_{h} - \overline{X})^{2}}{\sum_{i=1}^{n} (X_{i} - \overline{X})^{2}}\right)$ 

sta302midterm-10101-crowdmark-assess

#1 2 of 7



- Q1 (12 pts) Short answer questions. (SLR stands for Simple Linear Regression:  $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$ )
  - (a) (1 pt) Suppose I have **X=10**. Running **class(X)** on R console, what is the output?

## "numeric" (1)

(b) (2 pts) There is a generic function "tapply()" in R. You can get help information about it by running a command on R console, what is it?

(c) (2 pts) Suppose I have two vectors  $\mathbf{Y} = \mathbf{c}(3, 5, 1, 10, 12, 6)$ ;  $\mathbf{X} = \mathbf{c}(0,1,1,0,0,1)$  and I want to print out all the Y values where X takes 1. What R code achieves this?



(d) (2 pts) True or false and justify your answer: "in lecture we shown that  $\sum_{i=1}^{n} e_i = 0$  if a SLR is fitted to a set of n cases by method of Least squares. We then have  $\sum_{i=1}^{n} \epsilon_i = 0$  where  $\epsilon_i$  is the error term in SLR model."

False. (1)

Ei's are RVs with mean 0, we can't actually observe them. There is no reason that they should sum to 0.

sta302midterm-10101-crowdmark-assessr#1 3 of 7



(e) (2 pts) True or false and justify your answer:  $Var(b_0+b_1\bar{X}) = \sigma^2/n$  under the Gauss-Markov conditions.

True.  
Since 
$$Y_i = b_0 + b_1 X_i$$
 always goes through  $(\hat{X}, \hat{Y})_1$ , so  
We have  $\hat{Y} = b_0 + b_1 \hat{X}$   
=)  $Var(b_0 + b_1 \hat{X}) = Var(\hat{Y}) = Var(\hat{T}_1 \hat{T}_1 \hat{T}_1) = \frac{\sigma^2}{h}$ .

(f) (3 pts) Under the normal error model assumption, what is the distribution of b1 and why? (Please specify the distribution and the model parameters).

Why: 
$$b_1 = \frac{\sum(X_1 - \overline{X})Y_1}{S_{XX}} = \frac{\sum_{i=1}^{n}(X_i - \overline{X})^2}{S_{XX}}$$

Why:  $b_1 = \frac{\sum(X_1 - \overline{X})Y_1}{S_{XX}} = \frac{\sum_{i=1}^{n}(X_1 - \overline{X})^2}{S_{XX}}$ 
 $Y_i = \beta_0 t \beta_1 X_i + \epsilon_{vi}$ 
 $\Rightarrow tid_1 N(0, 0^2)$ 
 $\Rightarrow h_1 = \frac{1}{2} \sum_{i=1}^{n} k_i Y_i$  is also normal distributed

 $E(b_1) = \beta_1 \in b_1$  is BLUE.

 $Var(b_1) = \frac{0^2}{S_{XX}}$ 



- Q2 (15 pts) A simple linear regression model is fit on n observed data points. Assume Gauss-Markov conditions hold, coefficients are estimated by least squares method.
  - (2.a) (3 pts) In class, we show  $b_1 = \sum_i k_i Y_i, b_0 = \sum_i w_i Y_i$  where  $k_i = \frac{X_i \bar{X}}{S_{xx}}$ , and  $w_i = \frac{1}{n} \frac{(X_i \bar{X})\bar{X}}{S_{xx}}$ . Now find  $h_{ij}$  s.t.  $\hat{Y}_i = \sum_{j=1}^n h_{ij} Y_j$  and show that  $\sum_{j=1}^n h_{ij} = 1$

$$\begin{array}{l}
\hat{x}_{i} = b_{0} + b_{1} \hat{x}_{i}, \quad b_{i} = \sum_{j=1}^{n} k_{j} Y_{j}, \quad b_{0} = \sum_{j=1}^{n} W_{j} Y_{j} \\
= \sum_{j=1}^{n} \left( \frac{1}{h} - \frac{(x_{j} - \bar{x}_{j}) \bar{x}_{i}}{s_{xx}} \right) Y_{j} + X_{i} \sum_{j=1}^{n} \frac{x_{i} - \bar{x}_{i}}{s_{xx}} Y_{j} \\
= \sum_{j=1}^{n} \left( \frac{1}{h} - \frac{(x_{j} - \bar{x}_{i}) \bar{x}_{i}}{s_{xx}} + \frac{(x_{j} - \bar{x}_{i}) \bar{x}_{i}}{s_{xx}} \right) Y_{j} \\
= \sum_{j=1}^{n} \left( \frac{1}{h} + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} \right) Y_{j} = \sum_{j=1}^{n} h_{ij} Y_{j} \\
\text{where } h_{ij} = \frac{1}{h} + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} \\
= \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} (x_{j} - \bar{x}_{i}) \\
= \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} (x_{j} - \bar{x}_{i}) \\
= \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{j} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx}} = \sum_{j=1}^{n} \left( \frac{1}{h} \right) + \frac{(x_{i} - \bar{x}_{i}) (x_{i} - \bar{x}_{i})}{s_{xx$$

(2.b) (2 pts) Explain why the result in (a) implies that if  $h_{ii}$  is close to 1 then  $Y_i$  is close to  $\hat{Y}_i$ .

2.a implies that

$$\widehat{f}_{i} = \underset{\text{hi}}{\text{hi}} \underbrace{f_{i} + h_{i}} \underbrace{f_{i} + \cdots + h_{i}} \underbrace{f_{i} + \cdots + h_{i}} \underbrace{f_{i}} \underbrace{f_{$$

sta302midterm-10101-crowdmark-assess



Solno (2.c) (5 pts) Show  $cov(b_1, \bar{Y}) = 0$ .

Solno  $cov(b_1, \bar{Y}) = av(b_1, b_0 + b_1 \bar{X})$   $= av(b_1, b_0) + \bar{X} \quad Var(b_1)$   $= -\frac{\sigma^2 \bar{X}}{Svv} + \bar{X} \quad \frac{\sigma^2}{Svv}$ 

=0 (1)

in cover page )

 $b_1 = \sum_{i=1}^{n} \frac{|x_i - \overline{x}|}{|s_{xx}|} |x_i|$   $= \sum_{i=1}^{n} \frac{|x_i - \overline{x}|}{|s_{xx}|} |x_i|$   $= \sum_{i=1}^{n} \sum_{j=1}^{n} \frac{|x_i - \overline{x}|}{|s_{xx}|} |x_j|$   $= \sum_{i=1}^{n} \frac{|x_i - \overline{x}|}{|s_{xx}|} |x_i|$   $= \sum_{i=1}^{n$ 

=00

 $\omega_{V}(Y_{i}-\hat{Y}_{i},\hat{Y}_{j}) = \omega_{V}(Y_{i},\hat{Y}_{j}) - \omega_{V}(\hat{Y}_{i},\hat{Y}_{j})$   $= \omega_{V}(Y_{i}, Z_{m_{2}}^{n}, h_{jm}Y_{m}) - \omega_{V}(b_{0}tb_{i}X_{i}, b_{0}tb_{i}X_{j})$   $= \omega_{V}(Y_{i}, h_{ji}Y_{i}) - v_{ar}(b_{0}) - x_{j} \omega_{V}(b_{0}, b_{i}) - x_{i} \omega_{V}(h_{i}, b_{0})$   $= h_{ij} v_{ar}(Y_{i}) - \sigma^{2}(\frac{1}{h} + \frac{\overline{X}^{2}}{3xx}) + (X_{i}tX_{j}) \frac{\sigma^{2}x}{5xx} - x_{i}X_{j} \frac{\sigma^{2}}{3xx}$   $= \sigma^{2}h_{ij} - \sigma^{2}(\frac{1}{h} + \frac{1}{3xx}(\overline{X}^{2} - \overline{X}(X_{i}tX_{j}) + X_{i}X_{j})]$   $= \sigma^{2}h_{ij} - \sigma^{2}(\frac{1}{h} + \frac{1}{3xx}(X_{i} - \overline{X})(X_{j} - \overline{X})]$   $= \sigma^{2}h_{ij} - \sigma^{2}h_{ij} - \sigma^{2}h_{ij}$   $= \sigma^{2}h_{ij} - \sigma^{2}h_{ij}$ 



Q3 (23 pts) Analysis of Handspan Data

À simplé linear regression model is fitted to the data where y = handspan(cm), X = Height (inch), for n = 167 students.

Call:

lm(formula = HandSpan ~ Height, data = HH)

Coefficients:

Estimate Std. Error t value Pr(>|t|)

(Intercept) -3.00161

1.69394 [A] 0.0782 .

Height

0.35057

<2e-16 \*\*\*

Residual standard error: [D] on [E] degrees of freedom Multiple R-squared: [F], Adjusted R-squared: 0.5442

F-statistic: 199.2 on 1 and [G] DF, p-value: < 2.2e-16

> anova(mod)

Analysis of Variance Table

Response: HandSpan

Df Sum Sq Mean Sq F value

Height

1 [H] [I]

< 2.2e-16 \*\*\*

Residuals 165 [J]1.69

3.a) (10 pts) Find the 10 missing values (A through H). Give mark for correct value only.

199.2

 $A = \frac{-3.0016}{1.6939} = -1.772$   $R = \sqrt{199.2} = 14.1138$ 

B= 0.3506/C=0.07/8 D= 57.69=1.3

E = -2 = 165 F = H/(H+J) = 0.5470

G = 165 H = 1. = 336.648

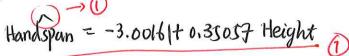
 $I = 1.69 \times 199.2 = 336.648$   $J = 165 \times 1.69 = 278.85$ 

3.b) (1 pt) What is the total sum of squares, SST= $\sum_{i=1}^{n} (Y_i - \bar{Y})^2$ ?

SST = H + J = 615.498

3.c) (2 pts) What is the fitted regression line? (define all terms in your

answer)





3.d) (3 pts) In the summary output with F value =199.2, what are the null and alternative hypotheses? And what do you conclude (in plain language instead of rejecting or failing to reject  $H_0$ ?

Ho: B=0 Ha: B=0. (1)

The observed F-value is 199,2 with pralue < 0.001. (1)

We reject the. We have very strong evidence that

there is a positive correlation between hardspan and height.

X = height

3.e) (2+2 pts) Find a 95% prediction interval of handspan when height is 20 inches. Also find the 95% confidence interval for the mean response at the same given height. Here are some quantiles from t-distributions which may be useful.

 $t_{0.95,1} = 1.6542, \; t_{0.95,165} = 1.6541, \; t_{0.95,166} = 1.6541; \; t_{0.95,167} = 1.6540$  $t_{0.975,1} = 1.9745, \ t_{0.975,165} = 1.9744, \ t_{0.975,166} = 1.9744; \ t_{0.975,167} = 1.9743$ 

X=20 Ŷ=-3,0016H 0,35057×20 = 4.0098 (1)

Denote In=Po+B,20+En, 95% PI for Yn is (9n ± to.975, 165. Spred) = (4.0098 ± 1.9744. J.69(Ht67+ 120-68,98)

= (40098±1.9744√3.1222) = (40098±1.9794×1.770)

52(6,)=0.02482=0.0006

= MSE = 1.69

959. CI for E(Y) at X=>0: (9±1.9744.5(9)) 0 = (4.0098 ± 1.9744 J 1.69 (+67 + 120-68.6719)2)

=) Sxx = 169 = 2742.76  $=(4.0098 \pm 1.9744 \times 1.968) = (1.6449, 6.3747) \leftarrow C2$ 3.f) (2 pts) If we assume height is a random variable, what is the es-

timate of the correlation between handspan and height? Are they positively or negatively correlated and why?

core(Y, X) = r = sign(b,) \( R^2 = + \sqrt{0.5470} = 0.7396

3.g) (1 pts) Answer only True or False for this statement: "For the Ftest in ANOVA, it can also be used for testing  $H_0: \beta_1 = 0$ ,  $vs H_a:$  $\beta_1 < 0$ "

False. (T)