

STAT S251F  
2020 Spring Take-home Assignment  
(Final Examination UG)

Name: Jiawei Wang

Student Number: S12395870

Question 1 (12 marks)

(a) (5(i) 2(ii))

(i) Because the sample size ( $n=100 > 30$ ), and we need to get the standard deviation of the population we use Z test:

$$\left( \bar{X} - Z_{\frac{\alpha}{2}} \frac{6}{\sqrt{n}}, \bar{X} + Z_{\frac{\alpha}{2}} \frac{6}{\sqrt{n}} \right) = (172.3, 174.1) \quad \textcircled{1}$$

from the Z-table:                      we can get from the question that:

$$Z_{0.025} = 1.96 \quad \textcircled{2}$$

$$n = 100 \quad \textcircled{3}$$

combine  $\textcircled{1}$   $\textcircled{2}$   $\textcircled{3}$ , we can get

$$\begin{cases} \bar{X} = \frac{172.3 + 174.1}{2} = 173.2 \\ 6 = \frac{9}{1.96} \approx 4.5918 \end{cases}$$

(ii) from the Z-table, when  $CI = 99\%$ ,  $Z_{0.005} = 2.575$

the 99% confidence interval for the mean height is:  $\left( \bar{X} - Z_{0.005} \frac{6}{\sqrt{n}}, \bar{X} + Z_{0.005} \frac{6}{\sqrt{n}} \right)$ .

$$CI = (172.0176, 174.3824) = \left( 173.2 - 2.575 \cdot \frac{4.5918}{\sqrt{100}}, 173.2 + 2.575 \cdot \frac{4.5918}{\sqrt{100}} \right)$$

(b) (2)

from the question, we can know that:                      from Z-table:

$$2 \times Z_{0.05} \cdot \frac{6}{\sqrt{n}} \leq 15 \quad \textcircled{1}$$

$$Z_{0.05} = 1.645 \quad \textcircled{3}$$

$$6 = 12 \quad \textcircled{2}$$

combine  $\textcircled{1}$   $\textcircled{2}$   $\textcircled{3}$ :  $n > 6.9274 \therefore$  the least number of tests is 7.

1

### Question 1 (c) (3)

the sample mean  $\bar{x} = \frac{\sum x_i}{n} = \frac{973.5}{110} = 8.85$

the sample standard deviation  $S = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{\sum x_i^2 + n(\bar{x})^2 - 2\bar{x} \sum x_i}{n-1}} = \sqrt{1.5165} \approx 1.2315$

Since we don't know the population standard deviation ( $\sigma$ ), and the sample size  $n > 30$ .

We use t-test

from the t-table  $t(0.005, 109) = 2.6217$  ①

$\therefore$  the 99% confidence interval for the population mean  $\mu$  is  $\left( -t_{0.005, 109} \cdot \frac{S}{\sqrt{n}} + \bar{x}, t_{0.005, 109} \cdot \frac{S}{\sqrt{n}} + \bar{x} \right)$

$$CI_{99\%} = (8.85 - 0.3078, 8.85 + 0.3078) = (8.5422, 9.1578)$$

### Question 2 (4 marks)

(a) It is a random sampling method.

Because every alumni who attend this gathering has the same probability to be choosed from the box and get the gifts.

(b) I think it is not a random sampling method.

Because this questionnaire is only for the person who take the MTR, not for the all Hong Kong citizens.

In other words, all the Hong Kong citizens do not have the

same probability to poll about the HK government's performance during COVID-19.

### Question 3 (5 marks)

(a) (4)

Brandy Band	Rank (1 <sup>st</sup> )	Rank (2 <sup>nd</sup> )	$ d  =  x - y $	$d^2$
1	2	3	1	1
2	3.5	2	1.5	2.25
3	1	4.5	3.5	12.25
4	5	4.5	0.5	0.25
5	3.5	1	2.5	6.25

$$r_s = 1 - \frac{6 \sum d^2}{n(n^2-1)} = 1 - \frac{6 \cdot 22}{5(25-1)} = -0.1$$

(b) (1)

the  $r_s$  is  $-0.1$ , which means that by normal standards, the association between the two variables would not be considered statistically significant.

### Question 4 (18 marks)

(a) (1)

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{495}{15} = 33$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{2066}{14}} \approx 12.1479$$

(b) (1)

the population mean is defined as the average difference between the recorded amount and the audited amount.

(c) (1) < Question 4 >

$$H_0: \mu = 25$$

$$H_1: \mu > 25$$

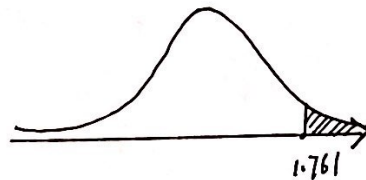
(d) (2)

since  $\sigma^2$  is not known and  $n = 15 < 30$ , we may resort to T test.

(e) (1)

$$\alpha = 0.05$$

from the t-table:  $t_{0.05; 14} = 1.761$



from the graph above, the rejection region is  $T > 1.761$ .

(f) (2)

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{33 - 25}{12.1479/\sqrt{15}} \approx 2.5506$$

(g) (3): from (f), we can know that on 5% level of significance.  $t = 2.5506 > 1.753$ .

$\therefore$  reject  $H_0$  at the 5% level.

$\therefore$  we conclude that, at the 5% level, the average difference is greater than \$25.

(h) (3)

since we do not know  $\sigma$ , and both the recorded and audited amounts of inventory are both normally distributed, we may resort to T test.

from t table:

$$t_{(0.025, 14)} = 2.145$$

the 95% confidence interval for the mean time  $\mu$  is  $(\bar{x} - t_{0.025; 14} \cdot \frac{s}{\sqrt{n}}, \bar{x} + t_{0.025; 14} \cdot \frac{s}{\sqrt{n}})$

$$CI_{95\%} = \left( 33 - 2.145 \cdot \frac{12.1479}{\sqrt{15}}, 33 + 2.145 \cdot \frac{12.1479}{\sqrt{15}} \right) = (26.2720, 39.7280)$$



Question 4.

(i) (2)

simplify the statement:

$$e = z_{\frac{\alpha}{2}} \cdot \frac{6}{\sqrt{n}} = z_{0.005} \cdot \frac{5}{\sqrt{n}} = 2.575 \cdot \frac{12.1479}{\sqrt{15}} \approx 8.0767$$

Because  $e < 10$ , we can use this formula to calculate the sample size, with 99% confidence within an error  $e = \pm \$10$ .

(j) (2)

there exists relationship between confidence interval and the level of significance.

like in (h), the range CI of  $n$  is always greater than \$25 with 95% CI.

And in (g), we conclude that, at the 5% level, the  $n$  is greater than \$25.

Question 5 (12 marks)

$$(a) (1) \quad \bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{546}{30} = 18.2$$

$$S = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2} = \sqrt{\frac{3180.8}{29}} \approx 10.4730$$

(b) (2) the sample distribution of sample means is the distribution which formed by this sample of Kong Lung Airlines's days between the receipt of the complaints and the resolution of the complaints.

(C) (5) (Question 5)

(i) Since  $\sigma$  is unknown, but the sample size is  $30 \geq 30$ .

We use T-test to get CI.  
from t-table

94%:

96%

98%

$$t_{0.03;29} = 1.957$$

$$t_{0.02;29} = 2.150$$

$$t_{0.01;29} = 2.462$$

$$CI_{94\%} = \left( \bar{x} - t_{0.03;29} \frac{s}{\sqrt{n}}, \bar{x} + t_{0.03;29} \frac{s}{\sqrt{n}} \right) \quad CI_{96\%} = \left( 18.2 - 2.150 \frac{10.4730}{\sqrt{30}}, 18.2 + 2.150 \frac{10.4730}{\sqrt{30}} \right)$$

$$= \left( 18.2 - 1.957 \frac{10.4730}{\sqrt{30}}, 18.2 + 1.957 \frac{10.4730}{\sqrt{30}} \right) \quad CI_{98\%} = \left( 18.2 - 2.462 \frac{10.4730}{\sqrt{30}}, 18.2 + 2.462 \frac{10.4730}{\sqrt{30}} \right)$$

$\therefore$

$$94\% \text{ CI} = (14.4580, 21.9420) \quad 96\% \text{ CI} = (14.0890, 22.3110)$$

$$98\% \text{ CI} = (13.4924, 22.9076)$$

(ii) from (ci), we can know that: The wider the confidence interval, the higher the confidence.

cd) (4) since we still do not know  $\sigma$ , because <sup>from (ci)</sup> when confidence equal to 94% the  $d$  is equal to  $\pm 3.7419$ . so the  $n$  must bigger than 30. we can use t-test, let the sample size be  $n_s$

$$\text{we can get that: } t_{0.05;29} \cdot \frac{s}{\sqrt{n_s}} \leq 3$$

$$n_s \geq 35.1792$$

At least 36 samples should be given.

6

# Question 6 (20 marks)

(a) (2) from the table:  $\bar{X} = \frac{\sum_{i=1}^{10} X_i}{10} = 50$   $\bar{Y} = \frac{\sum_{i=1}^{10} Y_i}{10} = 110$

$SS_X = 3400$   $SP = 6800$

$b = \frac{SP}{SS_X} = 2$   
 $a = -b\bar{X} + \bar{Y} = 10$

$\therefore$  the  $\beta_1$  is equal to 2,  $\beta_0$  is equal to 10  
 regression model:  $\hat{y} = 2x + 10$

(b) (6)  $MSB = \frac{SSB}{k-1} = \frac{18000}{1} = 18000$   $F\text{-stat} = \frac{MSB}{MSE}$   
 $MSE = \frac{SSE}{N-k} = \frac{17060.0205}{18} = 947.7789$   $= \frac{18000}{947.7789} = 18.9918$

ANOVA

Sources	DF	SS	MS	F-stat	P-value
Between Groups	1	18000	18000	18.9918	0.0004
Within Groups	18	17060.0205	947.7789		
Total	19	35060.0205			

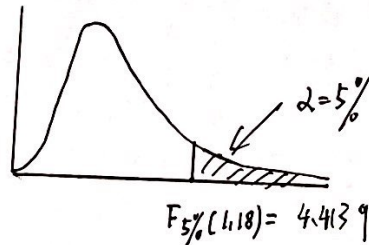
(c) (4)  $H_0: \beta_1 = 0 \longleftrightarrow H_1: \beta_1 \neq 0$

At 5% level,  $F_{5\%}(1, 18) = 4.4139$

$\therefore F_0 = 18.9918 > 4.4139$

$\therefore$  reject  $H_0$  at the 5% level

We conclude that there exists a linear relationship between  $X$  and  $Y$ .



(d) (2)

$R^2 = \frac{SSB}{SST} = \frac{18000}{35060.0205} \approx 0.5134$

Around 51.34% of the total variation in  $Y$  can be explained by  $X$ .



Question 6. (cont'd)

(e) (3)

(i) from Q6(a), we can get the fitted regression model:

$$\hat{y} = 2x + 10$$

$$e_1 = y_1 - \hat{y}_1 = 73 - (2 \times 30 + 10) = 3$$

$$e_6 = y_6 - \hat{y}_6 = 108 - (2 \times 50 + 10) = -2$$

$$e_2 = y_2 - \hat{y}_2 = 50 - (2 \times 20 + 10) = 0$$

$$e_7 = y_7 - \hat{y}_7 = 135 - (2 \times 60 + 10) = 5$$

$$e_3 = y_3 - \hat{y}_3 = 128 - (2 \times 60 + 10) = -2$$

$$e_8 = y_8 - \hat{y}_8 = 69 - (2 \times 30 + 10) = -1$$

$$e_4 = y_4 - \hat{y}_4 = 170 - (2 \times 80 + 10) = 0$$

$$e_9 = y_9 - \hat{y}_9 = 148 - (2 \times 70 + 10) = -2$$

$$e_5 = y_5 - \hat{y}_5 = 87 - (2 \times 40 + 10) = -3$$

$$e_{10} = y_{10} - \hat{y}_{10} = 132 - (2 \times 60 + 10) = 2$$

(ii) from (i):

$$\sum_{i=1}^{10} e_i = 3 + 0 - 2 + 0 - 3 - 2 + 5 - 1 - 2 + 2 = 0$$

f(3)

we need to find  $E = \sum_{i=1}^n (y_i - ax_i - b)^2$  minimum

then it is an Quadratic function with  $b$ . when  $b$  is equal to:  $-\frac{\sum_{i=1}^n (y_i - ax_i)}{2n} = \frac{\sum_{i=1}^n (y_i - ax_i)}{n}$

the  $E$  get minimum value

$$= \bar{y} - a\bar{x}$$

For any regression line, if it do not pass  $(\bar{x}, \bar{y})$ . then the error  $E$  would not be minimum, then we cannot get the accurate Linear regression equation.



# Question 7 (13 marks)

(a) (2)

$$\bar{X}_1 = \frac{\sum_{i=1}^9 X_i}{9} = \frac{2260}{9} \approx 251.1111$$

$$\bar{X}_2 = \frac{\sum_{i=1}^{10} X_i}{10} = \frac{2633}{10} = 263.3$$

$$S_1 = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} = \sqrt{523.61} \approx 22.8826$$

$$S_2 = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} = \sqrt{312.9} \approx 17.689$$

(b)  $H_0: \mu_1 = \mu_2 \Leftrightarrow H_1: \mu_1 < \mu_2$ .

(c) Since we do not know the value of  $\sigma$ , but samples are independent and the  $\sigma$  are equal but unknown, we use  $t$ -test (unpaired).

(d) the sample is normally distributed.

(e)

$$d.f = n_1 + n_2 - 2 = 9 + 10 - 2 = 17$$

$\therefore$  critical value

$$t(0.025; 17) = -2.110$$

Because of the variance of  $X_1, X_2$  were equal:

$$(f) S_p = \sqrt{\frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}} = \sqrt{\frac{8 \times 523.6 + 9 \times 312.9}{17}} \approx 20.2991$$

$$b_{12} = S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = 20.2991 \cdot \sqrt{\frac{1}{9} + \frac{1}{10}} \approx 9.3268$$

Under  $H_0$ , the test statistic is

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{b_{12}} = \frac{-12.1889}{9.3268} \approx -1.3069$$

(g)

$$t = -1.3069 > t(0.025; 17) = -2.110$$

$\therefore$  do not reject  $H_0$  at 5% level

which means the new drug do not have significantly effect than the placebo in reducing cholesterol levels.

Question 8 (16 marks)

Assuming there exists a linear relationship between the Exercise level and Body weight.

$$Y = \beta_1 X + \beta_0 + \varepsilon$$

the Test is in order to test the significance of linearity of regression line.

We may perform the usual F-test <i>

The Hypotheses are

$$H_0: \beta_1 = 0 \leftrightarrow H_1: \beta_1 \neq 0 \text{ <ii>}$$

ANOVA

Sources	DF	SS	MS	F-stat
Between Groups	1	193442.7	193442.7	240.7332
Within Groups	28	22499.5785	803.5564	
Total	29	215942.2785		

<v>

$$MSB = \frac{SSB}{K-1} = \frac{193442.7}{1} = 193442.7$$

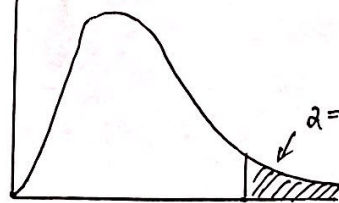
$$MSE = \frac{SSE}{N-K} = \frac{22499.5785}{28} = 803.5564$$

$$F\text{-stat} = \frac{MSB}{MSE} = 240.7332$$

At 5% level,  $F_{5\%}(1, 28) = 4.1960$

the rejection region:  $F\text{-stat} > 4.1960$

<iv>



$$F_{5\%}(1, 28) = 4.1960$$

<iii>

$$\therefore F\text{-stat} = 240.7332 > 4.1960$$

$\therefore$  reject  $H_0$  at the 5% level. <vi>

We conclude that there exist a relationship between the body weights and exercise level.

<vii>

END