

Statistical Data Analysis

2020 - 21

Take-home Assignment 2 (Mid-term Test).

Student Name: Angold, Jiawei WangStudent ID: S12395870 Mark: _____

Question 1 (24 marks)

Part (a) (15 marks)

$$(i) \text{ Standard deviation: } S = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} = \sqrt{\frac{(10.7-10.7)^2 + \dots + (10.6-10.7)^2}{5}} = \sqrt{\frac{0.025}{5}} = \sqrt{0.005} = 0.7071.$$

$$\text{Sample mean: } \bar{x} = \frac{\sum_{i=1}^N x_i}{N} = 10.7 \text{ (seconds).}$$

$$(ii) H_0: \mu = \mu_0 = 10.6 \text{ (seconds), (is equal).}$$

$$H_1: \mu > \mu_0 = 10.6 \text{ (seconds), (is greater).}$$

(iii) t test, because the SD of the population is known, but SD is unknown.
and the sample size $n < 30$, but follows normal distribution

$$(iv) \text{ From the t-table: } df = 5-1 = 4$$

$$\text{Take } \alpha = 5\% \quad t_{0.05} = 2.132 \quad \therefore \text{ rejection region: } t < -2.132.$$

$$(v) \text{ test statistic} \quad t = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{10.7 - 10.6}{0.7/\sqrt{5}} = \frac{0.1}{0.0791 \cdot \sqrt{5}} = 0.5628.$$

$$S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{(10.7-10.7)^2 + \dots + (10.6-10.7)^2}{5-1}} = \sqrt{0.00625} = 0.0791$$

(vi) $\because t > -2.132 \therefore$ we conclude that, at 5% level, the athlete's running time
is not greater than his past mean value.

Part (b) (9 marks)

(i) Because the sample size is 15 < 30.

∴ Using Unpaired t-test.

(ii) $H_0: \mu_{\text{Home}} = \mu_{\text{minor}}$ (μ_{Home} represents the average score of living in Home Environment).

$$H_1: \mu_{\text{Home}} < \mu_{\text{minor}}$$

$$(iii). \bar{X}_{\text{Home}} = \frac{1075}{15} = 71.6667 \quad S_{\text{Home}} = \sqrt{\frac{273.33}{15}} = 4.269$$

$$\bar{X}_{\text{minor}} = \frac{1145}{15} = 76.3333 \quad S_{\text{minor}} = \sqrt{\frac{549.33}{15}} = 4.826$$

From the question, we can know that the population variances are not equal.

$$\therefore \hat{\sigma}_{\text{HMIN}} = \sqrt{\frac{S_{\text{Home}}^2 + S_{\text{minor}}^2}{n_{\text{Home}} n_{\text{minor}}} = \sqrt{\frac{622.66}{15 \times 15}} = 1.6635}$$

Under H_0 , the test statistic is

$$t = \frac{(\bar{X}_{\text{Home}} - \bar{X}_{\text{minor}}) - (\mu_{\text{Home}} - \mu_{\text{minor}})}{\hat{\sigma}_{\text{HMIN}}} = \frac{71.6667 - 76.3333}{1.6635} = -2.8053.$$

the critical value is: $-t(28, 0.02) = -2.457$.

(iv). ∵ $t = -2.8053$

$$\therefore \text{d.f.} = n_{\text{Home}} + n_{\text{minor}} - 2 = 28.$$

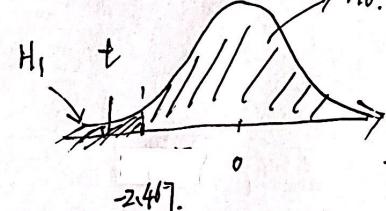
$$\therefore t(0.01, 28) = 2.763.$$

$$\therefore t < -t(0.01, 28) = -2.763.$$

∴ reject H_0 at the 1% level.

∴ At 1% level, there is insufficient evidence to support the claim that the μ_{Home} is equal to μ_{minor} .

which means, at 1% level, the living in the minority environment can lead to higher mean scores of social attitudes.



Question 2 (33 marks).

(a) $H_0: \mu_{\text{July}} = \mu_{\text{Sep}}$ (μ_{July} means the average price (\$) of rice in July).

$H_1: \mu_{\text{July}} \neq \mu_{\text{Sep}}$.

(b). t-test

Because $\textcircled{1}$ the sample size is 27 or 25 < 30. $\textcircled{2}$ We do not know the standard deviations of the population (both).

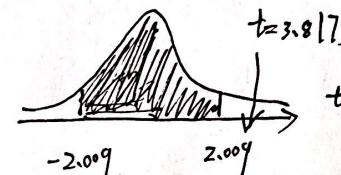
(c) (i) the critical value:

$$df = 27 + 25 - 2 = 50.$$

$$t_0 = t(0.05, 50) = 2.009.$$

rejection:

if t in $(-2.009, 2.009)$. H_0 is false



(ii). Because the pop variances are equal but unknown.

the standard error is

$$\hat{\sigma}_{JS} = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$= \sqrt{\frac{26 \cdot (0.24)^2 + 24 \cdot (0.27)^2}{50}} \cdot \sqrt{\frac{1}{27} + \frac{1}{25}} \approx 0.2548 \sqrt{0.07703} \approx 0.07073$$

Under H_0 , the test statistic is.

$$t = \frac{(\bar{X}_1 - \bar{X}_2)(\mu_1 - \mu_2)}{\hat{\sigma}_{JS}} = \frac{0.27}{0.07073} \approx 3.8173$$

(iii), $\because 3.8173 > 2.009$.

\therefore reject H_0 at 5% level.

The data support the contention that the mean prices in July and September are not the same at 5% level.

(d) (i). the critical value

$$\text{df} = n_1 + n_2 - 2 = 50.$$

$$t_0 = t(0.05, 50) = 2.009.$$

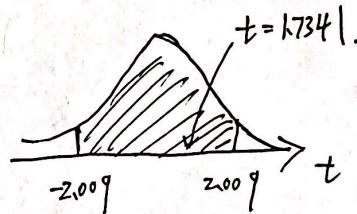
∴ Decision rule: Reject if $T_0 < -2.009$ or $T_0 > 2.009$.

(ii) Because the σ is unknown and unequal

$$\therefore \hat{\sigma}_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = \sqrt{\frac{(0.24)^2}{27} + \frac{(0.27)^2}{25}} \approx 0.1557$$

Under H_0 , the test statistic is.

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (M_1 - M_2)}{\hat{\sigma}_{\bar{x}_1 - \bar{x}_2}} = \frac{0.27}{0.1557} \approx 1.7341.$$



(iii). ∵ $1.7341 < 2.009$.

∴ do not reject H_0 at 5% level.

The data do not support the contention that the means in July and September are not the same.

Question 3 (43 marks).

(a). independent variable (X): Test Score

dependent variable (y): No. of Units sold.

(b). It can be seen from the scatter diagram:

Because except for a few adjacent two points, the rest of y increases with the increase of x .

We can think that x and y are roughly positively.

$$(c) r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \frac{292.6}{\sqrt{7.02 \cdot 1315.69}} = 0.9628.$$

① the r is positive, which means x and y are roughly correlated

② the r is bigger than 0.75 and very close to 1.0. which means x and y have strong positive correlation.

(d) (i). sample mean $\bar{X} = 3.4$. $\bar{y} = 137.1$.
 $\therefore a = \bar{y} - b\bar{X} = 137.1 - 3.4 \cdot 41.68 = -46.15$.

(ii). Let the fitted regression equation be like:

$$y = \beta_1 X + \beta_0$$

which $\left\{ \begin{array}{l} \beta_0 = b - a \cdot \bar{X} = 137.1 - 3.4 \cdot 41.68 = -46.15 \\ \beta_1 = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2} = 41.6809 / 41.68 \end{array} \right.$

\therefore the equation should be: $\hat{y} = 41.6809(X - 4.615)$.

(e) (i) When $X = 3.2$:

$$\hat{y} = 41.6809(3.2) - 4.615 = 128.763811999997$$

When $X = 6.0$:

$$\hat{y} = 41.6809(6) - 4.615 = 245.470359999997$$

(ii) from Q_{3C}, we can know that the $t \approx 0.9128$.

which means the coefficient of determination $R^2 = r^2 = 0.92698384$.

which means 92.698384% of the total variation in y can be explained by the fitted regression equation.

So it is not reliable, because x and y are not completely positively correlated.
 there must exist deviation.

ANOVA Table

(f).

SV	DF	SS	MS	F-ratio
Between Groups	$k-1=1$	89378.45	89378.45	122.2138
Within Groups.	$N-k=18$	13163.9133	731.3285	
Total	$N-1=19$	102542.3633		

$$SS_B = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2 \approx 89378.45$$

$$SS_W = \sum_{i=1}^k (n_i - 1) s_i^2 \approx 13163.9133$$

$$SS_T = SS_B + SS_W = 13163.9133 + 89378.45 = 102542.3633$$

(g). test for the linearity of the regression :

$$\text{suppose } \hat{y} = \hat{\beta}_1 \bar{x} + \hat{\beta}_0$$

$$(i) H_0: \beta_1 = 0 \Leftrightarrow H_1: \beta_1 \neq 0.$$

$$(ii). F_{5\%}(1, 18) = 4.4139.$$

(iii) if $F \leq 4.4139$. then, accept H_0 .

if $F > 4.4139$ then, reject H_0 .

(iv). From the ANOVA - table. we can know that

$$F = \frac{MS_B}{MS_W} = 122.2138$$

$$(v) \because 122.2138 > 4.4139$$

\therefore reject H_0 at 5% level.

(vi). We conclude that there exists a linear relationship between x and y .

(vii). From Q3 c, we can know that the r is 0.9628.

$$\therefore R^2 = r^2 = 0.92698384.$$

which means 92.698384% of the total variation in y can be explained by the fitted regression equation.