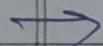
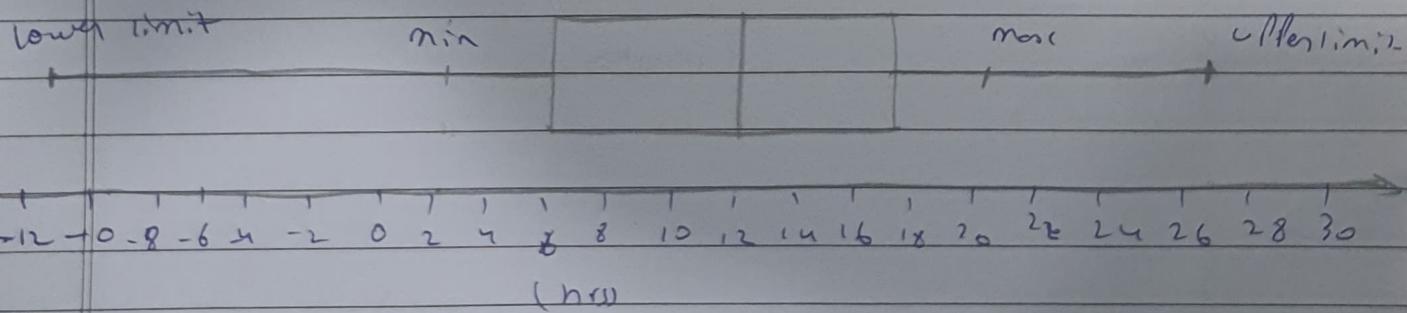


SM assignment~~PART Part A~~PART - ABox - Plot

- Q.1] A teacher collects data on the number of hours that her students spend studying per night. The data is as follows →

$$2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20.$$


- Minimum = 2
- Maximum = 20
- Median = $\frac{11 + 12}{2} = 11.5 \approx 12$ ($\left(\frac{N/2^{\text{th}} + (N/2+1)^{\text{th}}}{2} \right)$)
- $Q_1 = \frac{N}{4} = \frac{18^{\text{th}}}{4} = 6$
- $Q_3 = \frac{3N}{4} = \frac{54^{\text{th}}}{4} \approx 17$
- Inter Quartile Range = $17 - 6 = 11$
- Lower whisker = $Q_1 - 1.5 \times IQR = 6 - (1.5 \times 11) = -10.5$
- Higher whisker = $Q_3 + 1.5 \times IQR = 11 + (1.5 \times 11) = 27.5$
- Since max and min are written lower and upper whisker, there are no outliers.



Q.2] The following dataset representing three test scores (out of 100) : 76, 78, 84, 85, 88, 88, 89, 90, 90, 84, 88.

→ 76, 78, 84, 85, 88, 88, 89, 90, 90, 84, 88

$$\rightarrow \text{minimum} = 76$$

$$\text{maximum} = 95$$

$$\rightarrow \text{median} = 88$$

$$\rightarrow Q_1 = \frac{N}{4} = 84$$

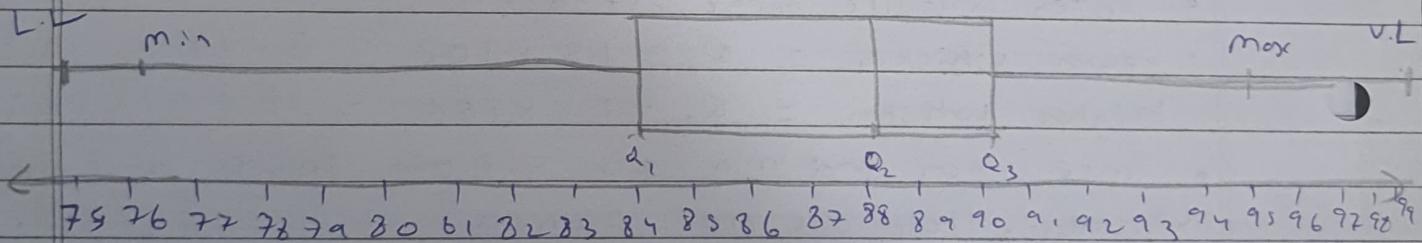
$$Q_3 = \frac{3N}{4} = 90$$

$$\rightarrow \text{IQR} = Q_3 - Q_1 = 90 - 84 = 6$$

$$\rightarrow \text{lower whisker} = Q_1 - (1.5 \times \text{IQR}) = 84 - (1.5 \times 6) = 75$$

$$\text{upper whisker} = Q_3 + (1.5 \times \text{IQR}) = 90 + (1.5 \times 6) = 99$$

(since max and min are within the whiskers, there are no outliers)



(marks)

3

Q.3) In a study of how students use mobile phones, the phone usage of random sample of 11 student was examined for a particular week.

The total length of calls, in mins for 11 students were

12, 23, 35, 36, 51, 53, 54, 55, 60, 72, 10

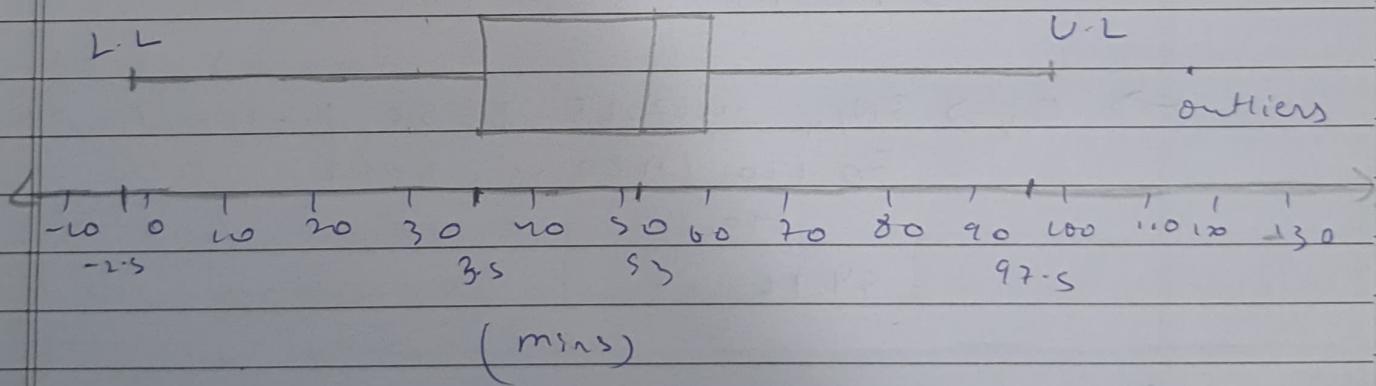
$$\rightarrow \text{median} = 53 \quad \rightarrow \text{minimum} = 12$$

$$\rightarrow Q_1 = 35 \quad (\text{N/4}) \quad \rightarrow \text{max} = 10$$

$$\rightarrow Q_3 = 60 \quad (3N/4) \quad \rightarrow \text{IQR} = 60 - 35 = 25$$

$$\begin{aligned} \rightarrow \text{Lower whisker} &= Q_1 - (1.5 \times \text{IQR}) \\ &= 35 - (1.5 \times 25) \\ &= -25 \end{aligned}$$

$$\begin{aligned} \rightarrow \text{Upper whisker} &= Q_3 + (1.5 \times \text{IQR}) \\ &= 60 + (1.5 \times 25) \\ &= 97.5 \end{aligned}$$



Qn) There are 543 members in a village, with the following age distribution:

<u>Age (yr)</u>	<u>f</u>	<u>Cf</u>
10-20	3	3
20-30	61	64
30-40	132	196
40-50	153	349
50-60	140	489
60-70	51	540
70-80	3	543

$$\rightarrow N = 543$$

$$\rightarrow Q_2 = 1 + \frac{h}{6} \left(\frac{N}{2} - C \right)$$

$$= 10 + \frac{10}{153} (271.5 - 196)$$

$$= 44.93$$

$$\rightarrow Q_3 \Rightarrow 3N/4 = 407.25, C = 50, f = 140, l = 349$$

$$= 50 + \frac{10}{140} (407.25 - 349)$$

$$= 54.16$$

$$\rightarrow Q_1 \Rightarrow N/4 = 135.75, C = 30, f = 132, l = 64$$

$$= 30 + \frac{10}{132} (135.75 - 64)$$

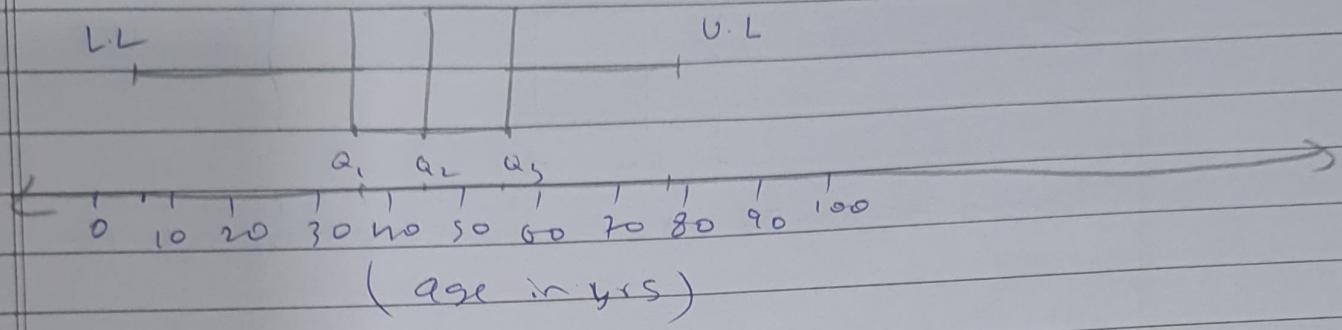
$$= 35.43$$

$$\rightarrow IQR = Q_3 - Q_1 = 18.72$$

$$\rightarrow \text{Lower limit} = Q_1 - (1.5 \times IQR) = 7.35$$

$$\text{Upper limit} = Q_3 + (1.5 \times IQR) = 82.24$$

\rightarrow No outliers



Q5) A professor analyses the scores of 20 courses \rightarrow A & B in maths. Create respective box plots and analyse if there are any outliers.

$A = 85, 90, 92, 95, 97, 98, 100, 102, 105, 110, 112, 120,$
 $125, 130, 135, 140$
 $B = 75, 78, 80, 82, 85, 88, 92, 95, 96, 98, 100, 105,$
 $108, 112, 115, 120.$

\rightarrow for class A

$$\text{median} : \left(\frac{N}{2} + \frac{N+1}{2} \right) = \frac{102 + 103}{2} = 102.5$$

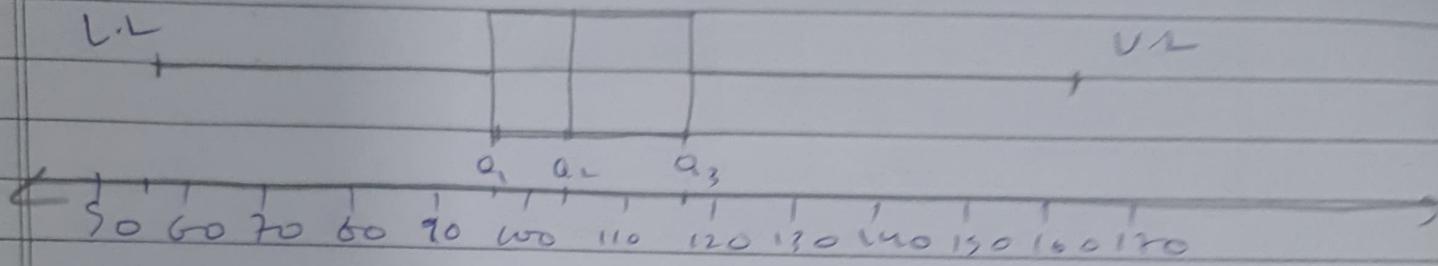
$$Q_1 = \cancel{\frac{Q - (1.5 \times IQR)}{2}} \quad Q_3 + Q_1 = 96$$

$$Q_2 = \frac{120 + 125}{2} = 122.5$$

$$IQR = Q_3 - Q_1 = 122.5 - 96 = 26.5$$

$$\text{Lower limit} = Q_1 - (1.5 \times IQR) \\ = 56.25$$

$$\text{Upper limit} = Q_3 + (1.5 \times IQR) \\ = 162.25$$



(score)

There are no outliers in Class A

→ For Class B

$$\text{median} = \frac{a_5 + a_6}{2} = 95.5$$

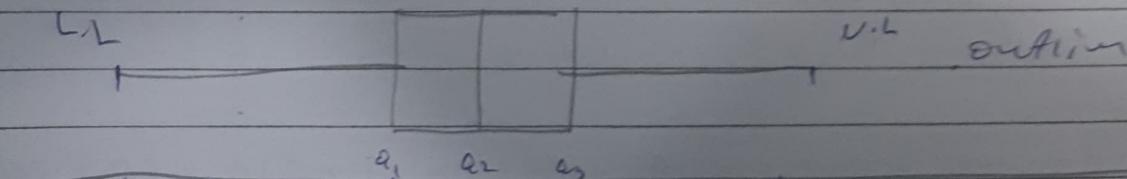
$$Q_1 = \frac{82 + 85}{2} = 83.5$$

$$Q_3 = \frac{105 + 108}{2} = 106.5$$

$$\text{IQR} = Q_3 - Q_1 = 106.5 - 83.5 = 23$$

$$\text{Lower limit} = Q_1 - (1.5 \times \text{IQR}) = 83.5 - (1.5 \times 23) = 49$$

$$\text{Upper limit} = Q_3 + (1.5 \times \text{IQR}) = 106.5 + (1.5 \times 23) = 141$$



(score)

∴ In Class B, there is 1 outlier = 150

Assignment - B

Q.1) Check whether a linear regression model is appropriate for the following data

$$x: 60 \ 70 \ 80 \ 85 \ 95 = 390$$

$$y: 70 \ 65 \ 70 \ 95 \ 85 = 385$$

$$\hat{y} = 65.41 \ 71.80 \ 70.70 \ 81.5 \ 87.94 = 384.91$$



$$\text{residual } (y - \hat{y}) : 4.59 \ - 6.84 \ - 3.28 \ - 13.5 \ - 29.4$$

$$(y - \hat{y}) : 21.0681 \ 16.7856 \ 68.5584 \ 182.25 \ 8.6436$$

$$\star \text{SS}_{\text{res}} \{ (y - \hat{y})^2 \} = 327.3052$$

$$x_y : 4000 \ 4550 \ 5600 \ 8075 \ 8075 = 30500$$

$$\text{SS}_{x_y} = 30500 - \frac{390 \times 385}{5}$$

$$= -470$$

$$\text{SS}_{x_x} = 31150 - 50420 \\ = 730$$

$$b_1 = \frac{\text{SS}_{xy}}{\text{SS}_{xx}} = 0.6438$$

$$b_0 = \frac{385 - (-0.6438) (390)}{5} \\ = 127.2164$$

$$y = 127.2164 - 0.6438x$$

$$S_e = \sqrt{\frac{327.3052}{5 - 2 - 1}} = \sqrt{\frac{327.3052}{2}}$$

$$= 12.7927$$

Q2) The equation of a regression line is, - the data is as follows; - solve for the residual and graph a residual plot. Do these data seem to violate any assumptions of regression?



x	y	x^2	xy
57	41	3249	2677
11	38	121	418
12	32	144	384
19	24	361	456
25	22	625	550
$\sum x = 124$	$\sum y = 153$	$\sum x^2 = 4500$	$\sum xy = 4482$
(n)			

→ $n = 5$

$$\rightarrow S_{xy} = \sum xy - \frac{\sum x \cdot \sum y}{n} = 4482 - \left(\frac{124 \times 153}{5} \right) = 4482 - 4896 = -414$$

$$\rightarrow S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 4500 - \frac{124^2}{5}$$

$$\rightarrow b_1 = \frac{S_{xy}}{S_{xx}} = \frac{-414}{4500} = -0.0912$$

$$\rightarrow b_0 = \frac{\sum y}{n} - b_1 \frac{\sum x}{n} = \frac{153}{5} - \left(-0.0912 \times \frac{124}{5} \right) = 24.862$$

$$\rightarrow \hat{y} = b_0 + b_1 x$$

$$z = 24.8624 + 0.3120x$$



y	\hat{y}	$y - \hat{y}$
47	42.6464	4.3536
38	28.2944	-9.7036
32	28.6064	3.3936
24	30.7704	-6.7904
22	32.0624	-10.6824
		$8(y - \hat{y}) = 0$

→ from the graph we infer that the residuals are scattered randomly around the origin
 ∴ the model is an appropriate linear regression model

→ no, the data does not seem to be violating any assumptions of regression, the residuals are randomly dispersed as they should be and the sum of residual comes out to be zero

Assignment -3Case Study

Q.1]

(i) response variable : US ticket Sales

explanatory variables = x_0 = budget x_1 = stars x_2 = runtime

$$\text{Ans} \quad \hat{y} = -22.989025 + 1.1344238x_0 + 24.972358x_1 - 0.40329565x_2$$

(ii) x_0 : US ticket Sales is positively correlated to budget x_1 : US ticket Sales is positively correlated to stars x_2 : US ticket Sales is negatively correlated to runtime
(minutes)

$$\text{Ans} \quad x_0 = \$32 \text{ million}, x_1 = 3 \text{ stars}, x_2 = 100 \text{ minutes}$$

$$\therefore \hat{y} = -22.989025 + 1.1344238(32) + 24.972358(3) - 0.40329565(100)$$

$$\therefore \hat{y} = 72.8972 \text{ US\$}$$

(v)

F-test

T-test

H_0 : model is not significant
P-value ≥ 0.05

H_0 : x_0 is not significant
P-value > 0.05

H_1 : model is significant
P-value ≤ 0.05

H_1 : x_0 is significant
P-value < 0.05

(vi) f-test : P-value = < 0.0001 < 0.05
∴ model is significant

(vii) $\approx t$ -test : P-value of run time > 0.05
i.e. $0.113 > 0.05$

∴ we will drop runtime from the model and
rerun the model

(viii) $R^2 = 0.6737$

adj $R^2 = 0.4603$

∴ R^2 is more closer to 1 want to add R.

Q2)

(i) Response variable = netflix subscribers

Explanatory variables : x_0 = monthly Price

x_1 = no of movies

x_2 = Price of rival

(ii) $\hat{y} = 2905.9639 + 111.2130032x_0 + 10.08715151x_1 - 605.173348x_2$

(iii) x_0 : monthly price is positively correlated to netflix subscribers

x_1 = no of movies is positively correlated to netflix subscribers

x_2 = Price of rival is negatively correlated to netflix subscribers

(iv) $x_0 = 0, x_1 = 4250, x_2 = 7$

$\therefore \hat{y} = 2905.9639 + 111.2130032(0) + 10.08715151(4250)$

$\therefore -605.173348(7) = 42479.8424$

FOR EDUCATIONAL USE

3

(v) f-test

T-test

H_0 : model is not significant
 $P\text{-value} > 0.05$

H_1 : model is significant
 $P\text{-value} < 0.05$

H_0 : χ^2 is not significant
 $P\text{-value} > 0.05$

H_1 : χ^2 is significant
 $P\text{-value} < 0.05$

(vii) f-test : $P\text{-value} = 4.30416 \times 10^{-5} < 0.05$
 \therefore model is significant

(vii) $R^2 = 0.889212953$

adj $R^2 = 0.85597684$

R^2 is more closer to 1 than adj R^2

(viii)

\therefore t-test : $P\text{-value} : x_0 = 0.900029323 > 0.05$

$x_1 = 0.199721303 > 0.05$

\therefore we will drop monthly prices (x_{10})

Q.10)

(i) response variable : y

Explanatory variable : $x_0 = 1$,

$$x_1 = x_L$$

$$x_2 = x_S$$

$$x_3 = x_T$$

$$(ii) y = -55.9 + 0.0105x_1 - 0.107x_2 + 0.579x_3 - 0.8707,$$

$$(iii) x_1 = 9, x_2 = 11, x_3 = 1, x_4 = 80$$

$$\therefore y = -55.9 + 0.0105(9) - 0.107(11) + 0.579(1) - 0.8707 \\ = -117.0035$$

(iv) x_1 is positively correlated to y

x_2 is negatively correlated to y

x_3 is positively correlated to y

x_4 is negatively correlated to y

(v) F -test

T -test

if H_0 is true the model is not significant, H_1 : x_1 is not significant

P-value $> 10\%$

P-value $> 10\%$

H_0 : the model is significant, H_1 : x_1 is significant

P-value $< 10\%$

P-value $< 10\%$

$$(vi) R^2 = 80.2\%$$

$$\text{adj. } R^2 = 70.7\%$$

$R^2 \rightarrow$ more closer to 100% than adj. R^2

vii) f-test : P-value = 0.000 < 0.05

∴ model is significant

viii) t-test : P-value = $\gamma_1 = 0.619 > 0.05$
∴ we will drop γ_1 from the model

Q.7 Predictors : 2 (γ_1 and γ_2)

$$\hat{Y} = 203.3932 + 1.1151\gamma_1 - 2.2113\gamma_2$$

Strength:

$$\therefore R^2 = 0.663, \text{ adj } R^2 = 0.636$$

∴ R^2 is more closer to 1 than adj-R²

∴ model is strong

$$\therefore P-value = \gamma_1 = 0.4403 > 0.05$$

∴ γ_1 is not significant

$$\gamma_2 = 0.0006 < 0.05$$

∴ γ_2 is significant

$$\therefore \text{co-eff } \gamma_1 = 1.1151$$

$$\gamma_2 = -2.2113$$

∴ γ_2 has stronger negative correlation to the rest of the γ_i , with coefficients weaker positive correlations

Q.12)

(i) perform Variable 'y'

Explaining variable : x_1 = College x_2 = age - years x_3 = Income - k x_4 = Gender

$$(ii) \quad y = 1.945280922 + 8.683514224x_1 + 0.06811533 \\ + 0.810434532x_2 + 8.683514224x_3$$

(iii) all x_1, x_2, x_3 and x_4 are positively correlated to y
 so college, age (years), income (k) and gender are positively
 correlated to the amount of customer purchase.

$$(iv) \quad x_1 = 5, \quad x_2 = 20, \quad x_3 = 12$$

$$0^{\circ} \quad y = 1.945280922 + 8.683514224x_1 + 0.06811533 \\ + 0.810434532x_2 + 8.683514224x_3$$

$$= 51.62362797 + 8.683514224x_3$$

(v) F-test

FTEST TEST

 H_0 : model is not significant

F-value > 20.5

 H_1 : model is significant

F-value < 20.5

 H_0 : β_1 is not significant

P-value > 20.5

 H_1 : β_1 is significant

f-value < 20.5

(vi) $R^2 = 0.974610626$; and $R^2 = 0.954310076$
 R^2 is more closer to 1 from side, so

(vii) P-value : model = 0.00352721 (0.03)

∴ model is significant.

Viii) P-value : college = 0.022624984 < 0.05

$$\text{age (yr)} = 0.72386968830 \dots$$

$$\text{income (h)} = 0.000323453 < 0.05$$

$$\text{gender} = 0.036034438 > 0.05$$

∴ we will drop college (y₁), age (yr) y₁, and gender (x₂) from the model and run the model

Q. 3)

- (i) response variable = concentration of chitosan -
explanatory variables = temperature extraction, co-₂,
dpm, pH, sodium, O₂, bicarbonate,
(a, c, i, s, l, n, m, nO₂, NO₂, NO,
SO₄²⁻)

(ii)

$$y = -0.426 + 0.012x_1 - 0.02x_2 - 0.08x_3 - 0.015x_4 +$$

$$0.003x_5 + 0.013x_6 + 0x_7 + 0x_8 + 0x_9 + 0x_{10}$$

$$+ 0.003x_{12} + 0x_{13} + 0x_{14}$$

$$\therefore y = -0.426 + 0.012x_1 - 0.02x_2 - 0.08x_3 - 0.015x_4 +$$

$$0.003x_5 + 0.013x_6 - 0.003x_{12}$$

(iii) $x_1 = 7, x_2 = 11, x_3 = 1, x_4 = 20$

$$\therefore y = -0.426 + 0.012(7) - 0.02(11) - 0.08 - 0.015(20)$$

$$+ 0.003x_5 + 0.003x_{12}$$

$$= -1.756 + 0.003x_5 + 0.013x_6 - 0.003x_{12}$$

(*) (v) F-test

H_0 : Model is not significant
if P-value > 0.05

H_0 : Model is significant
if P-value < 0.05

T-test

H_0 : β_0 is not significant
if P-value > 0.05

H_0 : β_1 is significant
if P-value < 0.05

(vi) $R^2 = 0.82322601$; adj $R^2 = 0.748631046$
 R^2 is more closer to 1 than adj R^2

(vii) P-value: model = 0.05139974 > 0.05
∴ model is slightly insignificant

(viii) P-value: temperature = 0.25270.5

Gibbs free energy = 0.383 > 0.5

D₀ = 0.253 > 0.05

bicarbonate = 0.34 > 0.05

ammonium - N = 0.455 > 0.5

Nitrate - N = 0.322 > 0.05

Phosphate - P = 0.714 > 0.05

∴ we will drop temperature, extraction coeff., D₀, bicarbonate

Q.14) Predictors : β_0 (1 variable), x_1 variable 1, x_2 variable 2
 Observations: 1(4)

$$y = 657.053 + 957103x_1 - 0.4169x_2 - 3.42152x_3$$

Strength : model

- $R^2 = 0.710$; adj $R^2 = 0.670$

R^2 is more closer to 1 than adj R^2

∴ model is strong.

After Model : P-value = 0.0012 < 0.05

∴ model is strong

Significant : Predictors

- x_1 : P-value = 0.0087 < 0.05

- x_2 : P-value = 0.2122 > 0.05

- x_3 : P-value = 0.0349 < 0.05

∴ x_1 and x_3 are significant

Q.15) Strength : model

- model : P-value = 0.0214 < 0.05

∴ model is strong and significant

- $R^2 = 51.50\%$; adj $R^2 = 40.4\%$

R^2 is closer to 1 than adj R^2

∴ model is moderately strong.

10

Strength: Predictors

x_1 : P-value = 0.05 < 0.05

x_2 : P-value = 0.805 > 0.05

x_3 : P-value = 0.230 > 0.05

$\rightarrow x_1$ is the only significant Predictor

The model doesn't appear to be fit the data well.

Recommend: Prof x_1 and x_3 from the model one by one then model

Q.16) Predictors = 2 (x_1 and x_3)

$$y = 203.3232$$

Note Some data and related question
from a M from PDF and Q11 from these sheets

Q.17)

Note Some data and related question from
Q8 of AB PDF and Q15 from these sheets