

IT3011 - Take home Assignment III
Regression Analysis.

Q1) Take data in the following table give the heart rate at rest (Y) and body weight (X) in kilograms.

X	Y
90	62
86	45
67	40
89	55
81	64
75	53

X [independent] = body.weight

Y [dependent] = heart.rate

a) Graph these data. Does it appear that there is a linear relationship between body weight and heart rate at rest?

For R studio, first insert above data into the CSV files and write below code to import the CSV file.

```
HRBW <- read.csv("C:\\users\\MSI\\Desktop\\HR_Bv.csv", header=TRUE)
```

```
attach(HRBW)
```

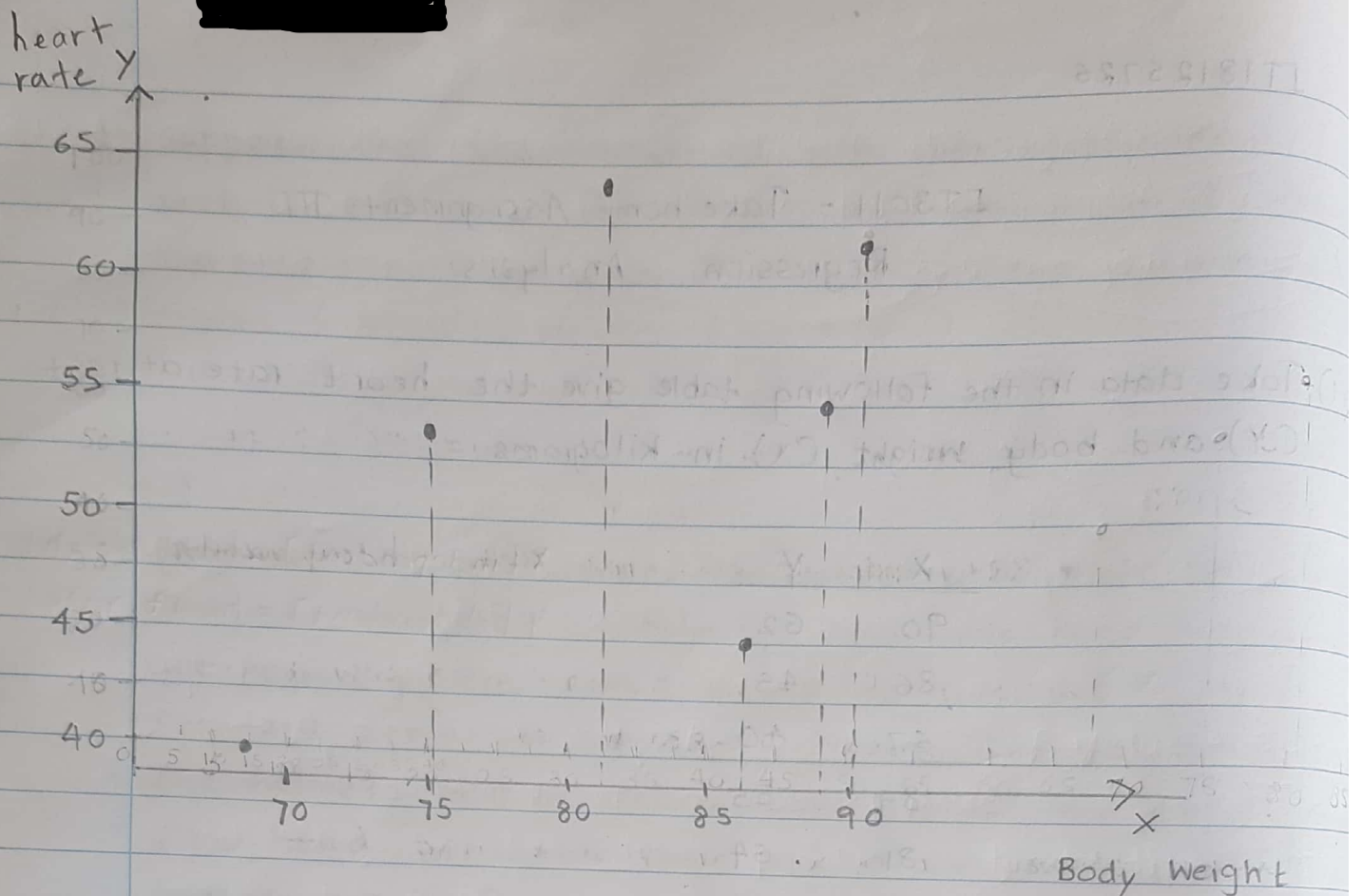
To view the scatter plot used below code

```
plot(body.weight, heart.rate, pch=21, col="blue",  
      bg="red")
```

below shows the scatter plot corresponding to the data given

T. N. P →

T. N. P →



according to the graph x and y have a linear relationship and to check how strong their relationship we can use below R code.

```
cor (body.weight, heart.rate)
= 0.568705 [Moderate positive linear Relationship]
```

Further we can perform Cor correlation test and get P value to Check Whether H_0 is rejected or not. for that we can use R code :-

```
cor.test (body.weight, heart.rate)
```

* From That we can get p-value as 0.2389. So at 5% Significance level we do not reject H_0 ($H_0 \rightarrow \rho_1 = 0$). \therefore we can Conclude the x and y have a relationship (Moderate positive linear)

- b) Compute $\hat{\beta}_0$ and $\hat{\beta}_1$ and write the regression equation for these data. plot the regression line on the graph obtained in part (a). Interpret the estimated regression coefficients.

To compute $\hat{\beta}_0$ and $\hat{\beta}_1$ we have to build the model using R studio

- `HRBWmodel <- lm(heart.rate ~ body.weight)`

To view the values corresponding to the $\hat{\beta}_0$ and $\hat{\beta}_1$:-

- `summary(HRBWmodel)`

using above R code we can view $\hat{\beta}_0$ and $\hat{\beta}_1$ -

$$\hat{\beta}_0 = 4.7990$$

$$\hat{\beta}_1 = 0.5947$$

Regression equation $\Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

$$\hat{y} = 4.7990 + 0.5947x //$$

To plot the regression line on graph used below R code -

`abline(4.7990, 0.5947)`

Interpret the estimated regression coefficients -

- A simple way to grasp regression coefficients is to picture them as linear slopes. The body weight coefficients in the regression equation is 0.5947. This coefficient represents the mean increase of heart rate in units, for every additional one unit in weight. If weight increases by 1 unit, the average heart rate increases by 0.5947 units.

Q1. d) Construct the ANOVA table and test the significance of the regression model at 5% level of significance [clearly state the hypothesis that you are going to test]

anova (HRBWmodel)

	Df	Sum of Sq	Mean Sq	F Value	Pr(>f)
body.weight	1	141.93	141.931	1.9122	0.2389
Residuals	4	296.90	74.225		

$\beta_1 = \text{body.weight}$

$H_0 \rightarrow \beta_1 = 0$ VS $H_1 \rightarrow \beta_1 \neq 0$ p value = 0.2389

Reject H_0 if $\alpha 0.05 > \text{pvalue}$

\therefore at 5% significance level p value is greater than α value

So we do not reject H_0 . $\beta_1 = 0$. therefore our first assumption of the linear relationship between x and y is false

Q1. C) Obtain the residuals and test the model assumption.

Residuals $\rightarrow e_1 = 3.679, e_2 = -10.942, e_3 = -4.643$

To test residuals \rightarrow plot $e_4 = -2.726, e_5 = 11.032, e_6 = 3.600$

To test residuals $\rightarrow \text{plot}(\text{HRBWmodel})$

1st graph \rightarrow Residuals vs fitted (dots are all over the graph and no pattern is visible. So we can say x & y are linearly related)

2nd graph \rightarrow Normal Q-Q (can't say exactly \therefore perform Shapiro test)

3rd graph \rightarrow Scale-Location (check the assumption of constant variance) No pattern is visible So we can say There is no constant variance \therefore assumption is not violated

4th graph \rightarrow Residuals vs Leverage (check the outliers) There are outliers in the obtained graph

\rightarrow

Shapiro. test (heart.rate)

$\Rightarrow p\text{-value} = 0.7116$

$H_0 = \mathcal{N}$ is normally distributed $\alpha = 0.05$

Reject H_0 if $P\text{value} < \alpha$ at 5% sig. level

\therefore At 5% significance level We do not reject H_0 . So we can conclude x and y are normally distributed, assumption is not violated.

e) What can you say about the fitted model?

1) `HRBWmodel <- lm(heart.rate ~ body.weight)` from Summary

`Summary (HRBWmodel), x.gcode`

\therefore x and y model is 0.06766 away from the fitted model.

2) using `plot(body.weight, heart.rate, pch=21, col="blue", bg="red")` and `abline(4.7990, 0.5947)` we can compare fitted model and our graph.

3) `Summary (HRBWmodel)`

We can see the values of \rightarrow

Coefficients :-

Intercept = 4.7990 Body.weight = 0.5947

\therefore Model is fitted. So the accuracy level of the fitted model is $[R^2] = 0.3234 //$

Q2) for a particular brand of automobile tire, an experiment was carried out to examine the relationship between temperature (X , in $^{\circ}F$) and tread wear of a tire (Y). The following data were collected.

X [independent] = Temperature Y [dependent] = Wear

a) plot the data and interpret the graph.

```
autom <- read.csv("c:\\users\\MSI\\Desktop\\AutoMobile.csv",  
                  header = TRUE)
```

```
attach(autom)
```

```
autom
```

```
plot(Temperature, Wear, pch = 21, col = "blue", bg = "red")
```

```
cor(Temperature, Wear)
```

```
= -0.2907
```

This graph interpret weakly negative linear relationship

b) compute the Pearson Product Moment Correlation Coefficient and test its significance at 5% level. Interpret the results.

To get Correlation Coefficient used a R code below →

```
cor.test(Temperature, Wear)
```

From this it gives the p value as = 0.2575. So at 5% Significance level α value = 0.05, so $1 - \alpha < p$ value

We do not reject H_0

∴ We can conclude that the temperature and Wear does not have a linear relationship. So our first assumption of the linear relationship is false.

c) fit a simple linear regression model -

First we have to build the model in R Studio using below code.

```
AutoModel <- lm (Wear ~ Temperature)
summary (AutoModel)
```

After executing above code we can have a value for $\hat{\beta}_0$ and $\hat{\beta}_1$.

$$\hat{\beta}_0 = 2.5678$$

$$\hat{\beta}_1 = -0.00427$$

Equation for the simple linear regression model \rightarrow

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

$$\hat{y} = 2.5678 - 0.00427x //$$

d) perform the residual analysis and state your comments on the regression assumptions.

To test residuals \rightarrow plot (AutoModel)

1st graph \rightarrow Residuals vs fitted (dots are all over the graph and no pattern is visible so we can say X & Y are linearly related)

2nd graph \rightarrow Normal Q-Q (check data is normally distributed or not) can't exactly come to a conclusion so perform Shapiro test

3rd graph \rightarrow Scale Location (check the assumption of constant variance) No pattern is visible so we can say There is no constant variance \therefore assumption is not violated.

$T \rightarrow N \rightarrow P \rightarrow$

graph 04 - Residuals vs Leverage (Check whether there are any outliers)

There are outliers in the graphs which we obtained

• Shapiro. test (Wear)

p-value = 0.2734

$H_0 \Rightarrow \epsilon$ is normally distributed $\alpha = 0.05$

Reject H_0 if $P\text{value} < \alpha$ at 5% sig. level

at 5% significance level we do not reject H_0 . So we can conclude Temperature and Wear are normally distributed. assumption is not violated.

e) Construct the ANOVA table and test the significance of the regression model at 5% level of significance. [clearly state the hypothesis that you are going to test]

anova (Automodel)

Response : Wear		Sum of	Mean	F	Pr(>F)
	Df	Sq	Sq	value	Pr(>F)
Temperature	1	0.05495	0.05494		0.2575
Residuals	15	0.59495	0.03966		

$H_0 \rightarrow \text{Temperature}(\beta_1) = 0$ p value = 0.2575

Reject H_0 if $\alpha = 0.05 > P\text{value}$

Since pvalue (0.2575) $> \alpha$ (0.05) do not reject H_0 at 5% significant level. At 5% significant level we do not reject H_0 hence we have enough evidence to conclude that $\beta = 0$

f) Perform a Lack of fit test and test whether the simple linear regression model is adequate or not

AutoModel $\leftarrow \text{lm}(\text{Wear} \sim \text{Temperature})$ [This is the reduced model]

Summary(AutoModel)

[This is the full model]

anova(AutoModel)

fac.Temp $\leftarrow \text{as.factor}(\text{Temperature})$

Model.aov $\leftarrow \text{aov}(\text{Wear} \sim \text{fac.Temp})$ [This is the full model]

Summary(Model.aov)

anova(Model.aov)

anova(AutoModel, Model.aov)

Need execute the above commands to get the anova table of "AutoModel" model and "Model.aov" model, And get the below out put

Model 1: Wear ~ Temperature

Model 2: Wear ~ fact.Temp

Res.DF	Rss	DF	Sum of Sq	F	Pr(>F)
15	0.59495				
9	0.09778	6	0.49716	7.6265	0.003965

- Model 1 is the usual linear regression model; $SSE(R) = 0.59495$
- Model 2 is telling R to consider Temperature as a "factor" instead of a continuous variable.

- $SSE(F) = 0.09778 = SSE(PE)$ [F = full, PE = Pure error]

- The Lack of fit SSE is $SSE(LF) = SSE(R) - SSE(F)$

$$= 0.59495 - 0.09778$$

$$= 0.49725$$

$$F^* = \frac{0.49725 / 6}{0.09778 / 9} = 0.082875$$

$$0.010865$$

$$= 7.62770$$

Reject H_0 because the p-value = 0.003965, $0.01 = 0.00003965$

because p value < f value //

3. Consider the following data Set.

a) fit a multiple linear regression model.

- attach (Multi)
- Mmodel $\leftarrow \text{lm}(Y \sim X_1 + X_2 + X_3)$
- Summary (Mmodel)

using Summary (Mmodel) R code we can find values for as below

$$\hat{\beta}_0 = -4.5470$$

$$\beta_1 = 0.2589$$

$$\beta_2 = 0.5997$$

$$\beta_3 = 1.2035$$

Regression equation for above data

$$\begin{aligned}\hat{y} &= \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3 \\ &= -4.5473 + 0.2589 x_1 + 0.5995 x_2 + 1.2035 x_3\end{aligned}$$

b) According to the model you fitted, which predictor is the most influential in determining the response variable.

Coefficients :-

(Intercept)	Estimate
(Intercept)	-4.5470
<u>x_3</u>	<u>1.2035</u>
x_2	0.5997
x_1	0.2589

The standardized coefficients show that x_3 has the standardized coefficient with largest absolute value, followed by x_2 and x_1 . $\therefore x_3$ is the most influential in determining the response variable.

c. compute the adjusted R^2 value and interpret the result.

`Mmodel <- lm(Y ~ X1 + X2 + X3)`

`summary(Mmodel)`

Adjusted R-Squared : 0.9856 //

Multiple R^2 value = 0.9881 and Adjusted R^2 value = 0.9856

So Adjusted R^2 value decreases because of the predictor improves the model by less than expected.

d) Construct the ANOVA table and test the significance of the parameters.

`-anova(Mmodel)`

Analysis of Variance Table

Source of Variation	SS	df	mean Sum of Squares	F-value
Regression	9011.7	$p-1=3$	$9011.7/3 = 3003.9$	$F = .3003.9$
Error/Residual	108.1	$n-p=14$	$108.1/14 = 7.7214$	7.7214
Total	9199.8	$n-1=17$		$= .389.0356$

Regression SSR = 9011.7 > P_1 's p value

P_1 's p value = 0.9011.7 \therefore Reject H_0

Reject H_0 if $0.95 > P_2$'s p value

P_2 (total no. of variables) = 3 independent + 1 dependent

Reject H_0 if $0.95 > P_3$'s p value

P_3 n = 18 value = 0.68224 \therefore Reject H_0

Atlas

In the above ANOVA table we can find that f value is 389.0356 which is less than α . So $H_0 \neq 0$ we can reject H_0 .

$P-1 = 3 \rightarrow$ degree of freedom 1

$n-p = 14 \rightarrow df_2$

At 5% level of significance,

f -table value = 3.344

$$f_{3,14} = 3.344$$

Reject H_0 if $f_{cal} > f_{3,14}$

Since $389.0356 > 3.344 \therefore$ We can reject H_0 //

- 04) A student has fitted a multiple linear regression with 6 predictors using 100 observations. The ANOVA table constructed in the analysis is given below. Fill in the blanks of the table.

Source of Variation	Degrees of freedom	Sum of Squares	Mean Squares (MSR)	f
Regression	$P-1 = 5$	22364	4472.8	63.21130
Residual	$n-P = 94$	6651.4	70.7595	
Total	$n-1 = 99$	29015.4		

$p = \text{no. of variables} = 6$

$n = \text{no of observations} = 100$

1) $MSR = \frac{SSR}{P-1}$

$4472.8 \times 5 = SSR$

$SSR = 22364 //$

2) Total = $SSR + SSE$

$SSE = \text{Total} - SSR$

$= 29015.4 - 22364$

$= 6651.4 //$

3) $MSE = \frac{SSE}{n-P}$

$= \frac{6651.4}{100-6} = 70.759$

4) f-value = $\frac{MSR}{MSE}$

$= \frac{4472.8}{70.7595}$

$= 63.2113 //$