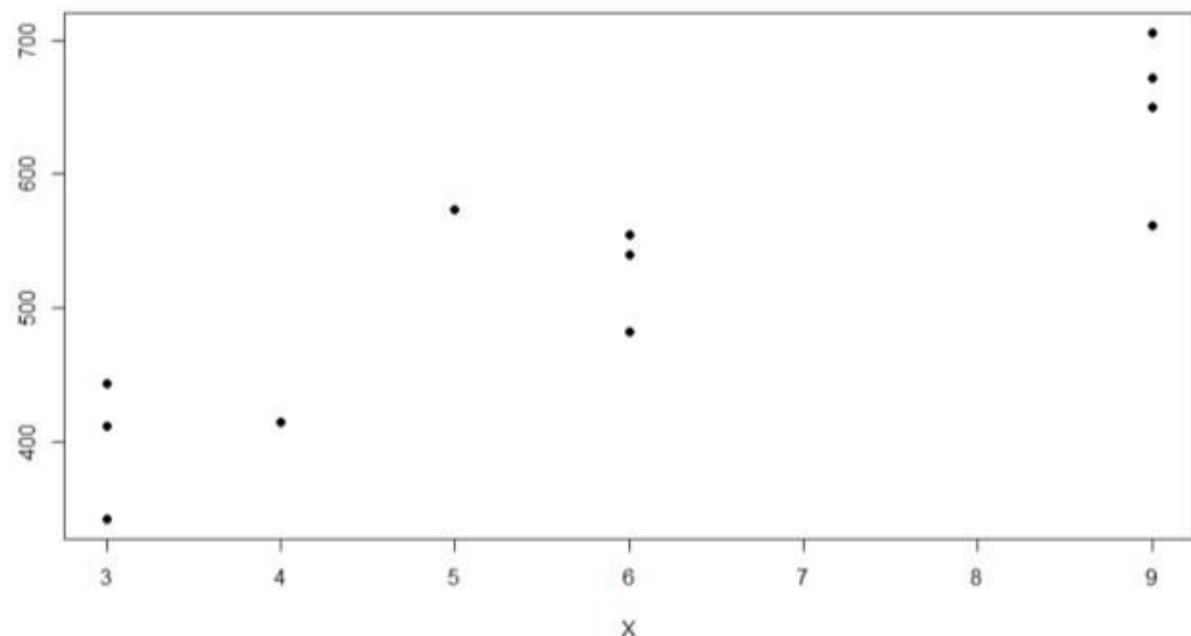


Ans 1)

a)

Scatter Plot for x and y



Theory Assignment 3

Ans 1-

- b) X and Y have a positive linear relationship.

c)

X	Y	X^2	Y^2	XY	$X - \bar{X}$	$(X - \bar{X})^2$
6	540	36	291600	3240	0	0
4	415	16	172225	1660	-2	4
6	555	36	308025	3330	0	0
9	650	81	422500	5850	3	9
3	412	9	169744	1236	-3	9
9	562	81	315844	5058	3	9
6	482	36	232324	2892	0	0
3	443	9	196249	1329	-3	9
9	706	81	498436	6354	3	9
5	574	25	329476	2870	-1	1
3	342	9	116964	1026	-3	9
9	672	81	451584	6048	3	9
$\Sigma X =$	$\Sigma Y =$	$\Sigma X^2 =$	$\Sigma Y^2 =$	$\Sigma XY =$	$\Sigma (X - \bar{X}) \neq \Sigma (X - \bar{X})^2 =$	
72	6353	500	3504971	40893	= 0	68

$$\begin{aligned}
 \text{Correlation, } r &= \frac{n \sum xy - (\sum x)(\sum y)}{\sqrt[n]{(\sum x^2) - (\sum x)^2} \sqrt[n]{(\sum y^2) - (\sum y)^2}} \\
 &= \frac{(12 \times 40893) - (72 \times 6353)}{\sqrt[(12 \times 500) - (72 \times 72)] \sqrt[(12 \times 3504971) - (6353)^2]} \\
 &= \frac{490716 - 457416}{\sqrt[816 \times 1699043]} = 0.89432
 \end{aligned}$$

(suggesting strong positive linear relation)

d)

Simple Linear Regression is given by

$$y = a_0 + a_1 x$$

Using table from the previous part for below calculations

$$\text{Here, } a_0 = \frac{(\sum y)(\sum x^2) - (\sum x)(\sum xy)}{n(\sum x^2) - (\sum x)^2}$$

$$= \frac{(6353 \times 500) - (72 \times 40893)}{(12 \times 500) - (72 \times 72)}$$

$$= \frac{232204}{816} = 284.5637$$

$$a_1 = \frac{n(\sum xy) - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2}$$

$$= \frac{(12 \times 40893) - (72 \times 6353)}{(12 \times 500) - (72 \times 72)}$$

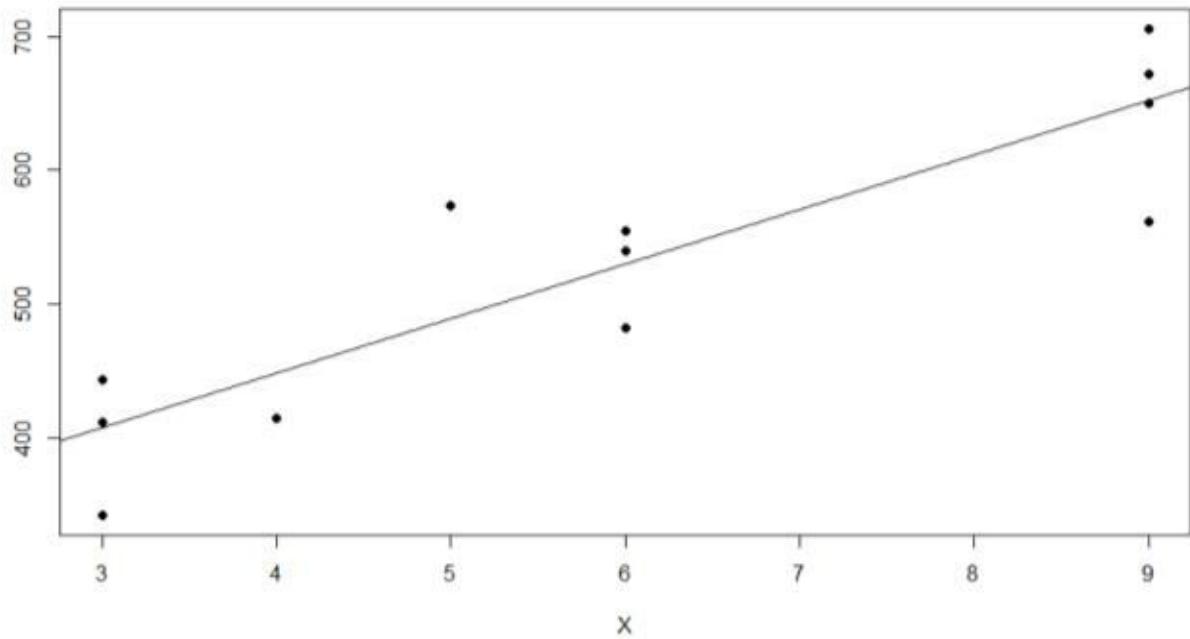
$$= \frac{33300}{816} = 40.8088$$

Hence, the equation :

$$y = 284.5637 + 40.8088x$$

e)

Scatter Plot for x and y



f)

Fitted values and residuals for each observation

x_i	y_i	$\hat{y}_i = a_0 + a_1 x_i$	$E = y_i - \hat{y}$	$(y_i - \hat{y})^2$
6	540	529.4165	10.5835	112.01047
4	415	447.7990	-32.7990	1075.7744
6	555	529.4167	25.5833	654.5052
9	650	651.8431	-1.8431	3.3970
3	412	406.9902	5.0098	25.0980
9	562	651.8431	-89.8431	8071.7826
6	482	529.4167	-47.4167	2248.3434
3	443	406.9902	36.0098	1296.7057
9	706	651.8431	54.1568	2932.9589
5	574	488.6078	85.3921	7291.8107
3	342	406.9902	-64.9903	4223.7390
9	672	651.8431	20.1569	406.3006

$$\sum (y_i - \hat{y}) = 28342.4259$$

$$\begin{aligned} \sum_{i=1}^{12} E &= 0 (10.5835 - 32.7990 + 25.5833 - 1.8431 + 5.0098 \\ &\quad - 89.8431 - 47.4167 + 36.0098 + 54.1568 + 85.3921 \\ &\quad - 64.9903 + 20.1569) \\ &= 0 \end{aligned}$$

g)

Variance of y as explained by x is given by

$$\begin{aligned} \text{Coefficient of Determination, } r^2 &= (\text{Correlation coeff.})^2 \\ &= 0.7998 \end{aligned}$$

h)

$$\sigma_{est} = \sqrt{\frac{\sum (y_i - \hat{y})^2}{N-2}}$$
$$= \sqrt{\frac{28342.4259}{10}}$$

$$= \frac{-48.5990}{53.2376}$$

Using the $(y_i - \hat{y})^2$ values from the previous table

i)

Hypothesis Testing to test for significance of r'

$$H_0: \rho = 0 \text{ (no correlation)}$$

$$H_a: \rho \neq 0 \text{ (correlation exists)}$$

Test statistic, $t = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$ (with $n-2$ degrees of freedom)

$$= \frac{0.89432}{\sqrt{\frac{1-(0.89432)^2}{10}}}$$

$$= \frac{0.89432}{\sqrt{0.0200}} = 6.3238$$

From t-table, at $\alpha = 0.01$

$$t_{critical} = t_{0.01} = 3.169$$

Since, ~~t~~ $t > t_{crit}$, we reject the null hypothesis and say that correlation exists b/w the variables x and y .

j)

Prediction interval

$$\hat{y}^* \pm t_{\alpha/2} s_{\text{pred}} \quad \text{where } \hat{y}^* \text{ is the predicted value of } y$$

& s_{pred} is the standard deviation to the predicted value of y

$$\text{Here, } \hat{y} = a_0 + a_1 x \Rightarrow 284.5637 + 40.8088x$$

$t_{\alpha/2}$ with 10 degrees of freedom is ~~3.169~~

$$s_{\text{pred}} = s^2 + s_{y^*}^2 \quad \text{where } s^2 \text{ is the standard error of estimate}$$

and $s_{y^*}^2$ is variance because of using \hat{y}^*

Here, $s^2 = \text{standard error of estimate}$

$$\text{and } s_{y^*}^2 = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

$$\text{Here, } s = 53.2376$$

$$s_{y^*}^2 = 53.2376 \Rightarrow \sqrt{\frac{1}{12} + \frac{(x-6)^2}{68}}$$

$$s_{\text{pred}} = (53.2376)^2 + 53.2376 \sqrt{\frac{1}{12} + \frac{(x-6)^2}{68}}$$

~~3.169~~

Prediction intervals,

$$\left[(284.5637 + 40.8088x) \pm 3.169 \left((53.2376)^2 + 53.2376 \sqrt{\frac{1}{12} + \frac{(x-6)^2}{68}} \right) \right]$$

Here, substituting different values of the independent variable (x) corresponding will give the different ~~value~~ prediction intervals for the corresponding y values.

Ans 2-a) Assumptions of Simple Linear Regression:

- i) Linear Relationship \rightarrow The outcome Y has a roughly linear relationship with the independent variable X .
- ii) Homoscedasticity \rightarrow For each value of X , the distribution of residuals has the same variance.
- iii) Independent errors \rightarrow Residuals should be uncorrelated (independent of each other).
- iv) Normally distributed residuals \rightarrow Residuals should be normally distributed.



Assumptions of Multiple Regression (Multivariable Regression) are same as that of Simple linear Regression with some additions.

- i) No multicollinearity \rightarrow This occurs when independent variables are too highly correlated to each other. Hence, the assumption states that the independent variables shouldn't be correlated to each other.

Assumptions of Logistic Regression:

- i) Logistic Regression requires that the dependent variable should be binary.
- ii) It requires that the observations should be independent of each other.
- iii) It requires that there should not be any multicollinearity among the independent variables.
- iv) It assumes a linear relationship b/w independent

variables and log odds.

- v) It requires a large sample size.
- b) Violation of assumptions:-
 - i) Linearity violation for simple linear regression → If we fit ~~a~~ a linear model to data which is are nonlinearly related, the predictions are likely to be seriously in error, especially when we extrapolate beyond the range of sample data.
 - ii) Multicollinearity violation for multiple linear regression → If we have independent variables that are highly correlated to each other, then it is hard to distinguish b/w the effects of the two or more related independent variables to the dependent variable. 1 unit increase in ^{one of the} ~~in~~ independent variable causes an increase in the other independent variable as well as the dependent variable, thus making it hard to understand which independent variable actually caused the change in value of the dependent variable.
 - c) Independence of outcome violation for logistic regression → If the outcomes are duplicated,

the error values will be similarly correlated and errors will follow a particular pattern OR trend.

Ans 3-

- a) In logistic regression, since the outcome variable is binary and only the predicted mean value (say m) is modeled, then there are only two possible values for error $m(1-m)$ or $(m-0)$.
- b) Variance in logistic regression is given by $P(1-P)$ where P is the probability of getting a binary 1 or True value. Hence, it varies with the probability of binary outcome 1 or true, being maximum when $P=0.5$ i.e. 0.25.

Ans 4-

- a) When x & y ~~are~~ have a non-linear relationship, we can use transformations (eg:- log transforms) to convert it into linear relation. Other transforms can also be applied depending on the distribution of data.
- b) To prove that if x_1 and x_2 are independent, then $\text{Cov}(x_1, x_2) = 0$

$$\text{Cov}(x_1, x_2) = E[(x_1 - E(x_1)) \cdot (x_2 - E(x_2))]$$

$$\begin{aligned}\text{Cov}(x_1, x_2) &= E[(x_1 - E(x_1)) \cdot (x_2 - E(x_2))] \\ &= E[x_1 \cdot x_2 - x_1 \cdot E(x_2) - x_2 \cdot E(x_1) + E(x_1) \cdot E(x_2)] \\ &= E(x_1 \cdot x_2) - E(x_1) \cdot E(x_2) - E(x_2) \cdot E(x_1) + E(E(x_1) \cdot E(x_2)) \\ &= E(x_1 \cdot x_2) - E(x_2) \cdot E(x_1) - E(x_2) \cdot E(x_1) + E(x_1) \cdot E(x_2) \\ &= E(x_1 \cdot x_2) - E(x_1) \cdot E(x_2)\end{aligned}$$

Now, if x_1 & x_2 are statistically independent
Then, $E(x_1 \cdot x_2) = E(x_1) \cdot E(x_2)$

Therefore,

$$\begin{aligned}\text{Cov}(x_1, x_2) &= E(x_1) \cdot E(x_2) - E(x_1) \cdot E(x_2) \\ &= 0\end{aligned}$$

Hence, proved.

parts a & b combined

Ans 5- a) & b) We have, $y_i = a_0 + a_1 x_i + \epsilon$
& $\hat{y}_i = \hat{a}_0 + \hat{a}_1 x_i$

$$\text{Now, } F(a_0, a_1) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$F(a_0, a_1) = \sum_{i=1}^n (y_i - (a_0 + a_1 x_i))^2$$

This is the SSE and should be minimum

Hence,

$$\cancel{\frac{\partial F}{\partial a_0}} = \sum_{i=1}^n 2(y_i - (a_0 + a_1 x_i)) = 0 \quad \text{--- (1)}$$

$$\frac{\partial F}{\partial a_1} = \sum_{i=1}^n 2(y_i - (a_0 + a_1 x_i))x = 0 \quad \text{--- (2)}$$

From eq (1),

$$\sum y_i - a_0 n - a_1 \sum x_i = 0$$

$$\Rightarrow \sum y_i = a_0 n + a_1 \sum x_i$$

Dividing by n

$$\frac{\sum y_i}{n} = a_0 + a_1 \frac{\sum x_i}{n} \quad \text{--- (3)}$$

$$\bar{y} = a_0 + a_1 \bar{x}$$

↳ This shows that (\bar{x}, \bar{y}) lie on the regression
(Part(b)) line.

From eqⁿ ②,

$$\sum_{j=1}^n - \left(a_0 n + a_1 \right)$$

$$\Rightarrow \sum_{j=1}^n y_i x_i - (a_0 \sum x_i + a_1 \sum x_i^2) = 0 \quad \left\{ \begin{array}{l} \text{Substituting } a_0 \\ \text{from eq}^n ③ \end{array} \right.$$
$$\sum_{j=1}^n y_i x_i = a_0 \sum x_i + a_1 \sum x_i^2$$
$$\sum_{j=1}^n y_i x_i = \left(\frac{\sum y_i}{n} - \frac{a_1 \sum x_i}{n} \right) \sum x_i + a_1 \sum x_i^2$$

$$\sum y_i \sqrt{\sum x_i^2} + a_1 (\sum x_i - \frac{\sum x_i}{n}).$$

$$\sum y_i x_i = \frac{\sum y_i \sum x_i}{n} - \frac{a_1 (\sum x_i)^2}{n} + a_1 \sum x_i^2$$

$$a_1 \left(\frac{(\sum x_i)^2}{n} - \sum x_i^2 \right) = \frac{\sum y_i \sum x_i}{n} - \sum y_i x_i$$

$$a_1 = \frac{\frac{\sum y_i \sum x_i}{n} - \sum y_i x_i}{\left(\frac{(\sum x_i)^2}{n} - \sum x_i^2 \right)}$$

Taking n common from above,

$$a_1 = \frac{\sum y_i \sum x_i - n \sum y_i x_i}{\left(\frac{(\sum x_i)^2}{n} - n \sum x_i^2 \right)} \quad \longrightarrow ④$$

Substituting a_1 to find a_0 .

$$a_0 = \frac{\sum y_i}{n} - \frac{a_1 \sum x_i}{n}$$

$$\begin{aligned}
 &= \frac{\sum y_i}{n} - \frac{\sum x_i}{n} \left(\frac{\sum y_i \sum x_i - n \sum y_i x_i}{(\sum x_i)^2 - n \sum x_i^2} \right) \\
 &= \frac{\sum y_i}{n} - \frac{\sum y_i (\sum x_i)^2 - n \sum x_i \sum y_i x_i}{n ((\sum x_i)^2 - n \sum x_i^2)} \\
 &= \frac{\sum y_i ((\sum x_i)^2 - n \sum x_i^2) - \sum y_i (\sum x_i)^2 + n \sum x_i \sum y_i x_i}{n ((\sum x_i)^2 - n \sum x_i^2)} \\
 &= \frac{-n \sum y_i \sum x_i^2 + n \sum x_i \sum y_i x_i}{n ((\sum x_i)^2 - n \sum x_i^2)} \\
 a_0 &\equiv \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}
 \end{aligned}$$

Ans 6- a) ① $H_0: \sigma = 15$
 $H_a: \sigma < 15$

Calculating using the t-value since sample size = 14

$$\begin{aligned}
 t &= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \\
 &= \frac{(9 - 15)}{\frac{15}{\sqrt{14}}}
 \end{aligned}$$

② Since the test is left tailed test, critical value for $\alpha = 0.05$ is obtained at the intersection of $(1-\alpha)$ and df. - i.e. 0.95 & 11 intersection

$$\chi^2_{\text{critical}} = 4.575$$

③ Obtaining the test value

$$t = (N-1) \left(\frac{s}{\sigma} \right)^2$$

$$= 11 \times \frac{9^2}{15}$$

$$= \frac{81}{15} = 5.46$$

④ Since the test value is less than χ^2_{critical} , the decision is to reject the null hypothesis. ie. there is enough evidence to support the claim that the standard deviation of number of aircrafts stolen every year in United States is less than 15.

b)

a-subpart) ① $H_0: \mu = 162.5$

$$H_a: \mu \neq 162.5$$

② Two-tailed test with $\alpha = 0.05$ meaning critical values are at $Z_{0.025}$ and $Z_{0.975}$ - ie. $Z_{\text{critical}} = -1.96, 1.96$

③ Finding the test statistic,

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{(165.2 - 162.5)}{6.9 / \sqrt{50}} = 2.77$$

④ Since the test value is within the critical region, we fail to reject the null hypothesis and hence there is no reason to believe that there is a change in average.

(4) Since the test value falls outside the critical region, we fail to accept the null hypothesis and believe that the average weight of the students has changed.

b-subpart)

① $H_0: \mu = 162.5$

$H_a: \mu \neq 162.5$

② Computing the value of test statistic

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = 2.77$$

③ Obtaining the area for $z=2.77$ from the z-table gives 0.0028.

Hence, p-value = 2×0.0028

p-value = 0.0056

④

~~Given, $\alpha = 0.05$~~

Given $\alpha = 0.05$, since $p\text{-value} < \alpha$, we fail to accept the null hypothesis and believe that there is a change in average weight of current batch of students.