**SOLUTION**                                    **// 60**
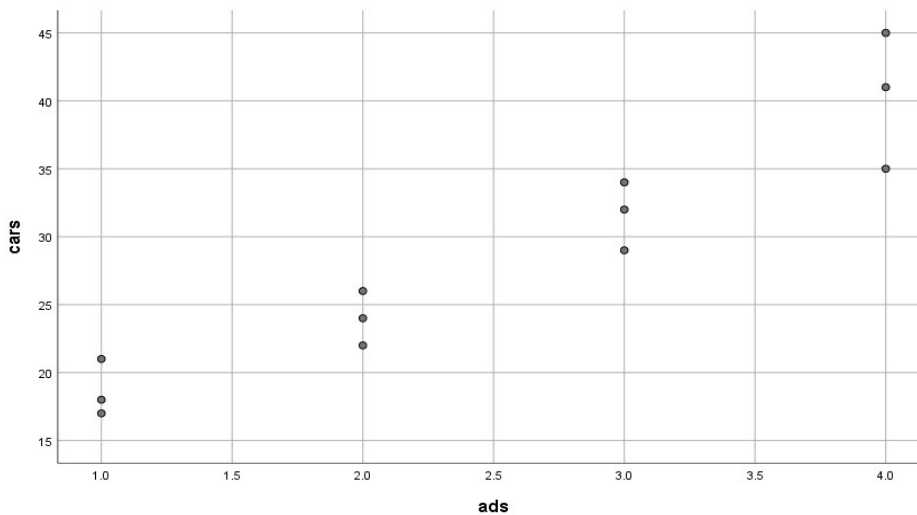
**1.  [42 marks]**

**[1]**  a)   Identify independent (*x*)  and dependent (*y*) variables.

> **x = # of ads per day  [1/2]**
> **y = # of cars sold  [1/2]**

**[2]**  b)      **[1]**



**Scatter plot indicates [1/2] <u>approximately straight line</u> (or linear relationship) with [1/2] <u>positive slope.</u>**

**[3]**  c)

**Model:**  $y = \beta_0 + \beta_1 x + \varepsilon$  **[1/2],  n = 12**

**Assumptions: (i)  *x*'s are observed without error  [1/2]**
> **(ii) *y*'s (or $\varepsilon$'s) are <u>independently</u> [1/2] <u>distributed</u> with <u>mean</u>** $E(y) = \beta_0 + \beta_1 x$
> **(or $E(\varepsilon) = 0$ ) [1/2]**
> **(iii) variance of *y*'s (or $\varepsilon$'<u>s) is constant [1/2]</u>,  $\sigma^2$  for all *x*'s**
> **(iv)  *y* ~ $N\left(E(y), \sigma^2\right)$ [1/2] for any value of *x*  (or  $\varepsilon \sim N\left(0, \sigma^2\right)$ for any value of *x*)**

NOTE: Assumptions (ii) – (iv) can be summarized also as $y \overset{i.i.d.}{\sim} N\left(E(y), \sigma^2\right)$ **(or** $\varepsilon \overset{i.i.d.}{\sim} N(0, \sigma^2)$**)**

1

**[4]** d)

$$[1/2]\ \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\displaystyle\sum_{i=1}^{n} x_i y_i - \frac{\left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)}{n}}{\displaystyle\sum_{i=1}^{n} x_i^2 - \frac{\left(\sum_{i=1}^{n} x_i\right)^2}{n}} = \frac{969 - \frac{(30)(344)}{12} \quad [1/2]}{90 - \frac{(30)^2}{12} \quad [1/2]} =$$
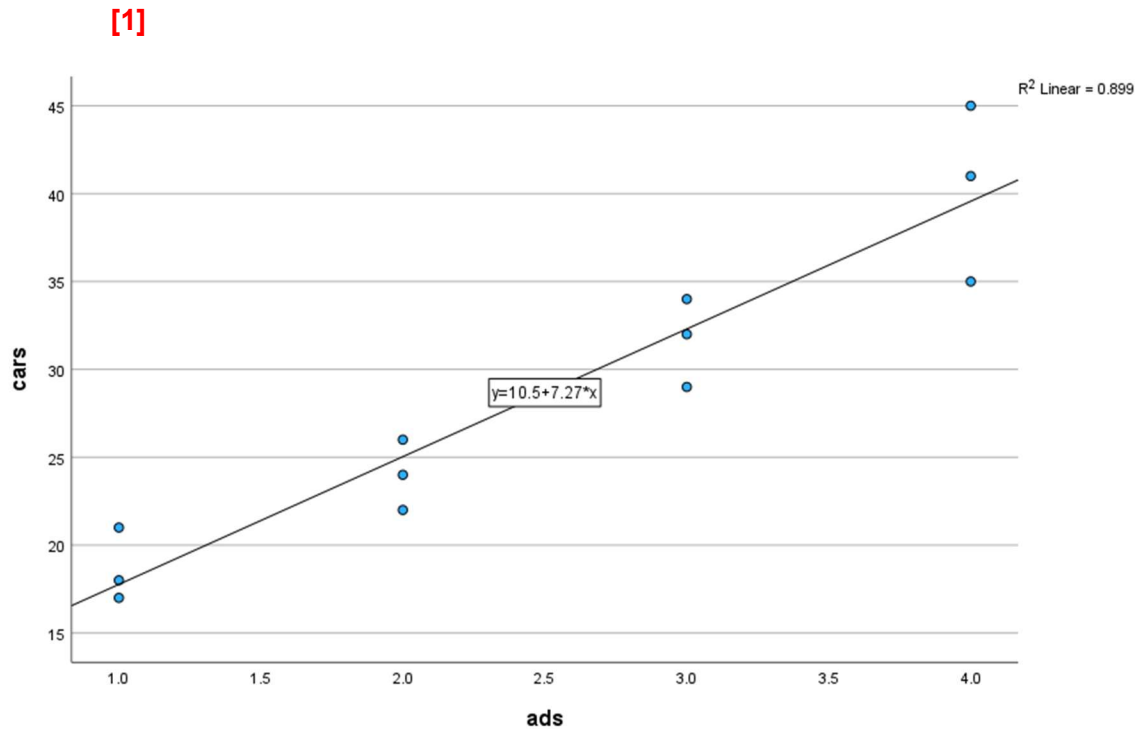
$$= \frac{109}{15} = \underline{\underline{7.266666667 \doteq 7.2667}} \quad [1/2]$$

$$[1/2]\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1\bar{x} = \frac{\displaystyle\sum_{i=1}^{n} y_i}{n} - \hat{\beta}_1\left(\frac{\displaystyle\sum_{i=1}^{n} x_i}{n}\right) = \frac{344}{12} - (7.266666667)\left(\frac{30}{12}\right) =$$

$$= 28.66666667 - 18.16666667 = \underline{10.5} \quad [1/2]$$

$\therefore$ **the least squares fitted regression line is given by:** $\hat{y} = \underline{10.5 + 7.2667}\,x$ **[1]**

**[1]** e)

**[1]**

**[4]** f)

$$[1] \quad s^2 = \frac{SSE}{n-2} = \frac{S_{yy} - \dfrac{S_{xy}^2}{S_{xx}}}{n-2} = \frac{\left[\displaystyle\sum_{i=1}^{n} y_i^2 - \dfrac{\left(\displaystyle\sum_{i=1}^{n} y_i\right)^2}{n}\right] - \dfrac{\left[\displaystyle\sum_{i=1}^{n} x_i y_i - \dfrac{\left(\displaystyle\sum_{i=1}^{n} x_i\right)\left(\displaystyle\sum_{i=1}^{n} y_i\right)}{n}\right]^2}{\left[\displaystyle\sum_{i=1}^{n} x_i^2 - \dfrac{\left(\displaystyle\sum_{i=1}^{n} x_i\right)^2}{n}\right]}}{n-2} =$$

$$= \frac{\left[10\,742 - \dfrac{(344)^2}{12}\right] - \dfrac{(109)^2}{15}}{10} = \frac{\overset{[1/2]}{880.6666667} - \overset{[1/2]}{792.0666667}}{10\,[1/2]} = \frac{88.6}{10} = \underline{\textbf{8.86}} \quad [1/2]$$

$$\therefore \ [1/2] \ s = \sqrt{s^2} = \underline{\textbf{2.976575213}} \cong \underline{\textbf{2.9766}} \ [1/2]$$

**[4.5]** g)

$H_0 : \beta_1 = 0$ **[1/2]** $\qquad \alpha = 0.05 \Rightarrow \alpha/2 = 0.025$

$H_a : \beta_1 \neq 0$ **[1/2]**

**test-statistics:** **[1/2]** $\ t = \dfrac{\hat{\beta}_1}{s/\sqrt{S_{xx}}} = \dfrac{7.266666667}{2.976575213/\sqrt{15}} = \underline{\textbf{9.45505}} \cong \underline{\textbf{9.455}}$ **[1/2]**

**R.R:** we reject $H_0$ if $t < -t_{\alpha/2;n-2} = -t_{0.025;10} = \textbf{- 2.228}$

$\qquad$ or $t > t_{\alpha/2;n-2} = t_{0.025;10} = \textbf{2.228}$ $\qquad$ **[1]**

**[1/2]**
Since *t* = $\underline{\textbf{9.455}}$ > 2.228, <u>we reject</u> $H_0$ **[1/2]** and conclude that at 5% level of significance there is an evidence to say that the No. of ads per day and the No. of cars sold are linearly related. **[1/2]**

**[2.5]** h) $\quad 1 - \alpha = 0.95 \Rightarrow \alpha = 0.05 \Rightarrow \alpha/2 = 0.025$
**[1/2]**

**[1/2]**
$$\beta_1 \in \left(\hat{\beta}_1 \pm t_{\alpha/2;n-2} \ \frac{s}{\sqrt{S_{xx}}}\right) = \left(7.2667 \pm t_{0.025;10} \ \frac{2.976575213}{\sqrt{15}}\right) = \left(7.2667 \pm 2.228 \left(0.768548415\right)\right) =$$

$$= \left(7.2667 \pm 1.712325869\right) = \left(5.554374131 \ , \ 8.979025869\right) \cong \underline{\left(5.5544 \ , \ 8.979\right)} \ [1]$$

**(1/2 mark for each correct interval value)**

3

i.e. We are 95% confident that in repeated sampling the true value of the population slope would lie in the interval (5.5544 , 8.979). **[1/2]**

**[12]** i)

**[1/2]** $TSS = S_{yy} = \sum_{i=1}^{n} y_i^2 - \dfrac{\left(\sum_{i=1}^{n} y_i\right)^2}{n} = \underline{880.6666667}$ **[1/2]** **(as calculated in part (f))**

**[1/2]** $SSR = \dfrac{S_{xy}^2}{S_{xx}} = \underline{792.0666667}$ **[1/2]** **(as calculated in part (f))**

**[1/2]** $SSE = TSS - SSR = \underline{88.6}$ **[1/2]** **(calculated in part (f))**

**[1/2]** $MSR = \dfrac{SSR}{1} = \underline{792.0666667}$ **[1/2]**

**[1/2]** $MSE = \dfrac{SSE}{n-2} = \dfrac{88.6}{10} = \underline{\underline{8.86}}$ **[1/2]** **(= $s^2$)**  **(as calculated in part (f))**

**[1/2]** $F = \dfrac{MSR}{MSE} = \underline{\underline{89.39804}}$ **[1/2]**

| Source | d.f. | SS | MS | F |
|---|---|---|---|---|
| Regression | 1 | 792.0666667 | 792.0666667 | 89.398 |
| Error | 10 | 88.6 | 8.86 | |
| Total | 11 | 880.6666667 | | |
| | **[1/2]** | **[1/2]** | **[1/2]** | **[1/2]** |

$H_0 : \beta_1 = 0$  $\alpha = 0.05$
$H_a : \beta_1 \neq 0$  **[1]**

**one mark for each column, if values**
**are entered correctly**

**test-statistics:** $F = \dfrac{MSR}{MSE} = \underline{\underline{89.398}}$  **[1/2]**

**R.R:**  we reject $H_0$ if $F > F_{\alpha(1, n-2)} = F_{0.05(1,10)} = $ **4.96**  **[1]**

**[1/2]**
Since **$F$ = 89.398** > 4.96, <u>we reject</u> $H_0$ **[1/2]** and conclude that at 5% level of significance there is an evidence to say that a linear relationship between the No. of ads per day and the No. of cars sold exists. **[1/2]**

**[5]** j)

**[1/2]** $r = \dfrac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \dfrac{109}{\sqrt{(15)(880.6666667)}} = \underline{\underline{0.94836}} \cong \underline{\underline{0.95}}$ **[1/2]**

i.e. the No. of ads per day and the No. of cars sold are <u>positively</u> **[1/2]**correlated (related) with the <u>strength of their relationship approx. 95%.</u> **[1/2]**

**[1/2]** $r^2 = \dfrac{SSR}{TSS} = \underline{\mathbf{0.89939}} \cong \underline{\mathbf{0.90}}$ **[1/2]**

**i.e. approximately 90% of the total variation in the data is explained by the regression line (and approx. 10% is due to error). [1]**

**The model is a very good fit (or it is a very good model). [1]**

**[3]** k)

### Model Summary[b]    [1]

| Model | R | R Square | Adjusted R Square | Std. Error of the Estimate |
|---|---|---|---|---|
| 1 | .948[a] | .899 | .889 | 2.977 |

← S

a. Predictors: (Constant), ads

b. Dependent Variable: cars

**[1]**

p-value < 0.05
⇒ reject Ho

### ANOVA[a]

| Model | | Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|---|
| 1 | Regression | 792.067 | 1 | 792.067 | 89.398 | .000[b] |
| | Residual | 88.600 | 10 | 8.860 | | |
| | Total | 880.667 | 11 | | | |

a. Dependent Variable: cars

b. Predictors: (Constant), ads

**[1]**

### Coefficients[a]

| Model | | Unstandardized Coefficients | | Standardized Coefficients | | | 95.0% Confidence Interval for B | |
|---|---|---|---|---|---|---|---|---|
| | | B | Std. Error | Beta | t | Sig. | Lower Bound | Upper Bound |
| 1 | (Constant) | 10.500 | 2.105 | | 4.989 | .001 | 5.810 | 15.190 |
| | ads | 7.267 | .769 | .948 | 9.455 | .000 | 5.554 | 8.979 |

a. Dependent Variable: cars

$\hat{\beta}_0$ (arrow to B 10.500)

$\hat{\beta}_1$

t test statistics

$\beta_1 \in (5.554, 8.979)$

p-value < 0.05 ⇒ reject Ho

5

**2. [9 marks]**

**[6]** a) **95% C.I. for *E(y)* when $x_p$ = 0:**

$\hat{y}$ = 10.5 + 7.2667 (0) = **10.5 [1/2]** and $\quad 1-\alpha = 0.95 \Rightarrow \alpha = 0.05 \Rightarrow \alpha/2 = 0.025$

**[1/2]** $\therefore E(y) \in \left( \hat{y} \pm t_{\alpha/2;n-2} s \sqrt{\dfrac{1}{n} + \dfrac{(x_p - \bar{x})^2}{S_{xx}}} \right) = \left( 10.5 \pm t_{0.025;10} (2.976575213) \sqrt{\dfrac{1}{12} + \dfrac{(0-2.5)^2}{15}} \right) =$

**[1/2]** $= (10.5 \pm 2.228(2.104756518)) = (10.5 \pm 4.689397522) = (5.810602478 \ , \ 15.18939752) \cong$

$\cong (5.8106 \ , \ 15.1894)$  **[1] (1/2 mark for each correct interval value)**

**i.e. We are 95% confident that in repeated sampling the <u>average value</u> of the No. of cars sold when the 0 ads were run, will fall in the interval (5.8106 , 15.1894). [1/2]**

**and**

**95% P.I. for *y* when $x_p$ = 0:**

$\hat{y}$ = 10.5 + 7.2667 (0) = **10.5** and $\quad 1-\alpha = 0.95 \Rightarrow \alpha = 0.05 \Rightarrow \alpha/2 = 0.025$

**[1/2]** $\therefore y \in \left( \hat{y} \pm t_{\alpha/2;n-2} s \sqrt{1 + \dfrac{1}{n} + \dfrac{(x_p - \bar{x})^2}{S_{xx}}} \right) = \left( 10.5 \pm t_{0.025;10} (2.976575213) \sqrt{1 + \dfrac{1}{12} + \dfrac{(0-2.5)^2}{15}} \right) =$

**[1/2]** $= (10.5 \pm 2.228(3.645545226)) = (10.5 \pm 8.122274764) = (2.377725236 \ , \ 18.62227476) \cong$

$\cong (2.377 \ , \ 18.622)$  **[1] (1/2 mark for each correct interval value)**

**i.e. We are 95% confident that in repeated sampling the No. of cars sold when the 0 ads were run, will lie in the interval (2.377 , 18.622). [1/2]**

**<u>Conclusion</u>:**

- **The P.I. is <u>wider</u> [1/2] than C.I. (as expected), since the variability in the error for predicting a single value of y is always greater than the variability of the error for the estimation of the mean/average value of y.**

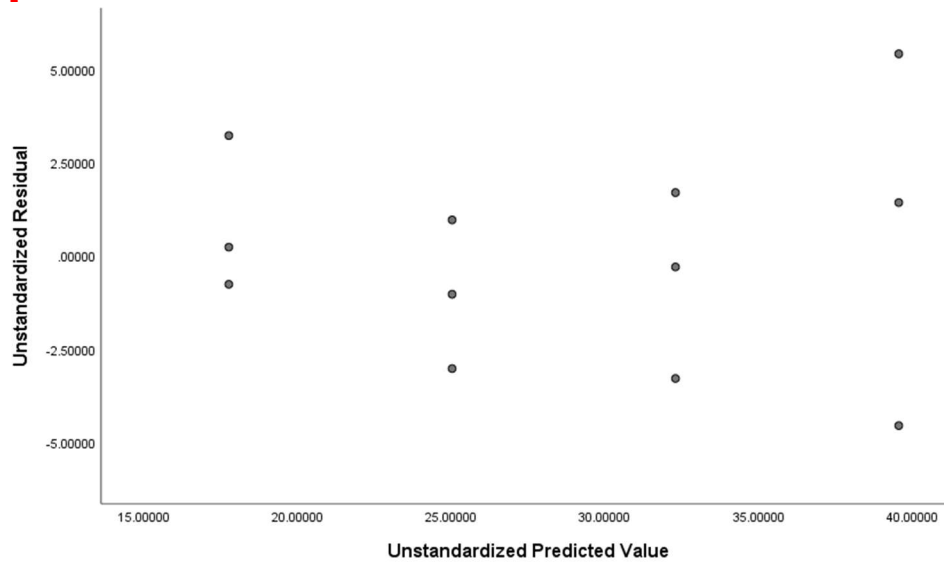**[3]** b) **95% C.I. for *E(y)* when $x_p$ = 0 : (5.81031 , 15.18969) [1]**

**95% P.I. for *y* when $x_p$ = 0 :    (2.37722 , 18.62278) [1]**

$\hat{y}$ = **10.50000** when $x_p$ = 0
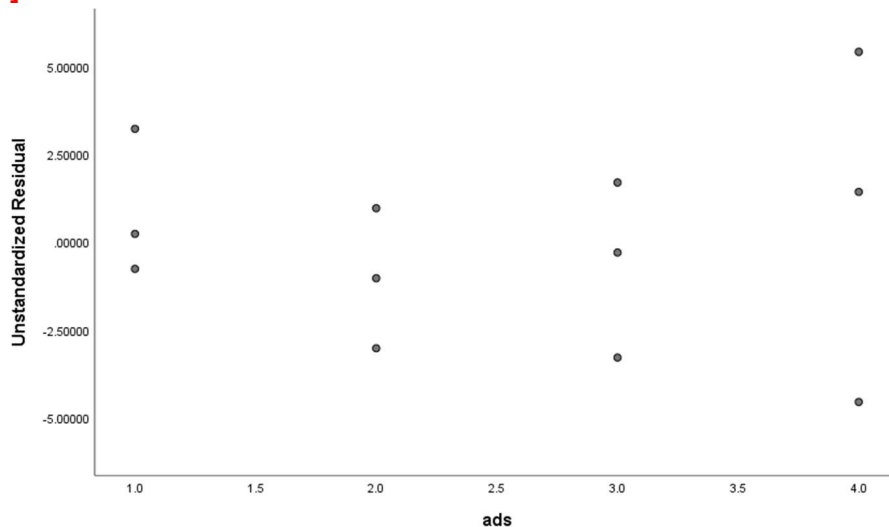**[1]**

## 3. [9 marks]

**[1]**



- **Residuals seem to be randomly scattered around zero (i.e. no pattern)** ⟹ **[1/2]** <u>no violations **of independence (and linearity)** [1/2]</u>

<u>NOTE:</u> if students saw and indicated a pattern, then there is a **[1/2]** <u>violation</u> of the <u>assumption of the independence of the errors</u> (and/or linearity) **[1/2]**.
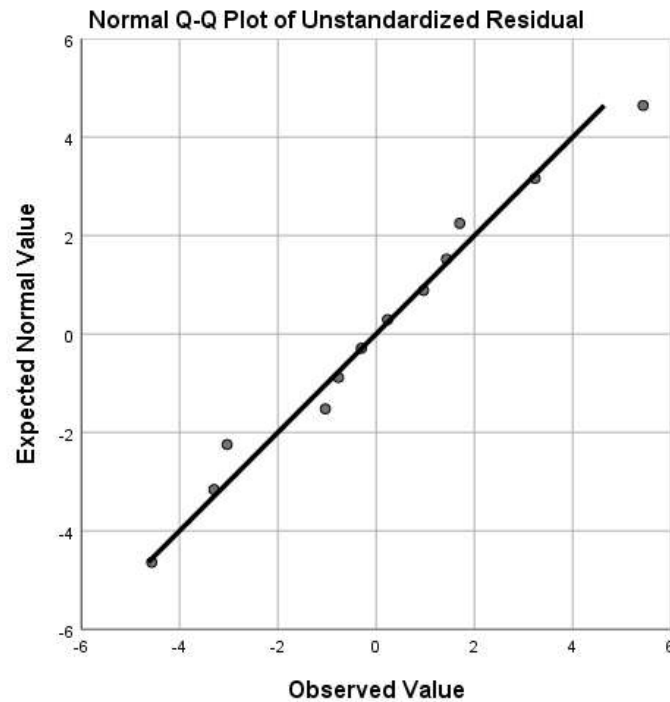
**[1]**



- **Residuals seem to be randomly scattered around zero (i.e. no pattern)** ⟹ **[1/2]** <u>**no violations of constant variance**</u> **[1/2]**
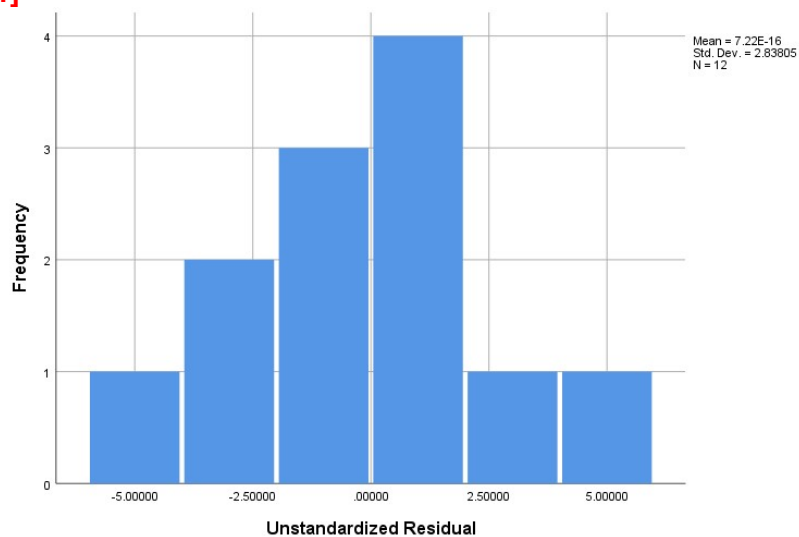
7

if students saw and indicated a pattern, then there is a **[1/2]** <u>violation</u> of the <u>assumption of the constant variance</u> **[1/2]**.

**[1]**



Normal Q-Q Plot of Unstandardized Residual

-   **Q-Q Plot of residuals shows approximately the straight line** $\Rightarrow$ **[1/2]** <u>no violations of normality of the errors</u> **[1/2]**

**[1]**



Mean = 7.22E-16
Std. Dev. = 2.83805
N = 12

**Histogram of the errors looks <u>approx. bell-shaped</u> [1/2]. It is <u>not really symmetric</u> [1/2] (but it is most likely due to small sample size, as *n* = 12). Therefore, since Q-Q plot did not show any violations, <u>errors are normally distributed</u>. [1/2]**

**All the plots suggest that the model assumptions are (reasonably) satisfied [1/2]**