

SOLUTION

- Q.1. [1] a) the response variable, y , is lung damage (1)
 [1] b) the explanatory variable, x , is yrs of smoking (1)
 [6] c) SLR model: $y = \beta_0 + \beta_1 x + \epsilon$ (1), $n=10$

- where (i) x_i 's are observed without error (1)
 (ii) ϵ_i 's are independently (1) distributed (with mean $E(\epsilon) = \beta_0 + \beta_1 x$) (1)
 (iii) variance of ϵ 's is constant (1) for all x 's
 (iv) $y \sim N(E(y), \sigma^2)$ for any value of x or, equivalently

- (i) x 's observed without error
 (ii) ϵ_i 's are indep. distributed with mean $E(\epsilon) = 0$
 (iii) variance of ϵ 's is constant for all x 's
 (iv) $\epsilon \sim N(0, \sigma^2)$ for any value of x

d) [5] $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} =$
 $= \frac{18055 - \frac{(319)(530)}{10}}{11053 - \frac{(319)^2}{10}} = \frac{1148}{876.9} = 1.309157 =$
 $\underline{\underline{1.309}}$ (1)

$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{\sum y_i}{n} - \hat{\beta}_1 \frac{\sum x_i}{n} =$
 $= \frac{530}{10} - (1.309157) \left(\frac{319}{10} \right) = 11.23788 = \underline{\underline{11.238}}$ (1)

∴ the least-squares fitted regr. line is given by

$$\hat{y} = 11.23\beta + 1.309x \quad (1)$$

e) $s^2 = \frac{SSE}{n-2}$, where $SSE = \sum y_i^2 - \frac{(\sum xy)^2}{\sum x_i^2} =$
 $[6] \quad = \left[\sum y_i^2 - \frac{(\sum y_i)^2}{n} \right] - \frac{\left[\sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} \right]^2}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} =$
 $= \left[30600 - \frac{(530)^2}{10} \right] - \frac{(1148)^2}{876.9} = 2510 - 1502.913$
 $= 1007.087 \quad (1)$

∴ $s^2 = \frac{SSE}{n-2} = \frac{1007.087}{8} = 125.8859 \quad (1)$

∴ $s = \sqrt{s^2} = 11.21989$

f) $H_0: \beta_1 = 0 \quad (1)$; $\alpha = 0.10 \Rightarrow \alpha/2 = 0.05$
 $H_a: \beta_1 \neq 0 \quad (1)$

- test-stat: $t = \frac{\hat{\beta}_1}{s/\sqrt{S_{xx}}} = \frac{1.309157}{11.21989/\sqrt{876.9}} = 3.45524 \quad (1)$

R.R. - we reject H_0 if $t < -t_{\alpha/2; n-2} = -t_{0.05; 8} = -1.860$
 or $t > t_{\alpha/2; n-2} = t_{0.05; 8} = 1.860 \quad (1)$

- since $t = 3.45524 > 1.860$, we reject H_0 and
 conclude that at 10% level of signif., there is
 an evidence that a linear relationship
 between yrs of smoking and lung damage
 exists. (12)

g) 90% C.I. for β_1 :
[4]

$1-\alpha = 0.90$
 $\alpha = 0.10$
 $\alpha/2 = 0.05$

$$\rho_1 \in \left(\hat{\beta}_1 \pm t_{\alpha/2; n-2} \times \frac{S/\sqrt{S_{xx}}}{\sqrt{S_{xx}}} \right) =$$

$$= \left(1.309 \pm \underbrace{t_{0.05; 8}}_{1.860} \times \frac{11.21999}{\sqrt{976.9}} \right) =$$

$$= (1.309 \pm 0.704736) = (0.604264, 2.013736)$$

ie. we are 90% confident that in repeated sampling the true value of the population slope β_1 will lie in the interval (0.604, 2.014).

h) TSS = $\sum y^2 = 2510$ (given)
SSE = $\sum (y - \hat{y})^2 = 1007.087$ (calculated in part c)
SSR = TSS - SSE = 1502.913

$MSR = \frac{SSR}{1} = 1502.913$
 $MSE = \frac{SSE}{n-2} = \frac{1007.087}{8} = 125.8859$
 $F = \frac{MSR}{MSE} = 11.93869$

ANOVA table:

Source of Variation	d.f.	SS	MS	F
Regression	1	1502.913	1502.913	11.93869
Error	8	1007.087	125.8859	
Total	9	2510		

$$H_0: \beta_1 = 0 \quad \text{①} \quad \alpha = 0.10$$

$$H_a: \beta_1 \neq 0 \quad \text{①}$$

- test-stat. $F = \frac{MSR}{MSE} = 11.93269$ ①

- R.R. - we reject H_0 if $F > F_{\alpha}(1, n-2) = F_{0.10}(1, 9) = 3.46$ ①

- since $F = 11.93269 > 3.46$, we reject H_0 and conclude that at 10% level of signif. there is an evidence that a linear relationship between yrs of smoking and lung damage exists. ①

i) $r = \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}} = \frac{1142}{\sqrt{(276.9)(2570)}} = 0.773802 \approx 0.774$ ①

ie. the lung damage and yrs of smoking are positively correlated (related) with the strength of the relationship approx. 77.4%. ①②

$$r^2 = \frac{SSR}{TSS} = 0.59277 \approx 0.599$$
 ①

ie. approx. 60% of the total variation in the data is explained by the regression line (and 40% is due to error). ①②

ie. Model is not good as 60% is quite low. ①

i) 95% c.I. for $E(y)$ where $x_p = 40$ ①

$$\Rightarrow \hat{y} = 11.232 + 1.309(40) = 63.592$$
 ①

$$\text{and } 1 - \alpha = 0.95$$

$$\alpha = 0.05$$

$$\alpha/2 = 0.025$$

$$\Rightarrow E(y) \in \left(\bar{y} \pm t_{\alpha/2; n-2} \times s \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}}} \right) =$$

$$= \left(63.598 \pm \underbrace{t_{0,025; 8}}_{2.306} \times (11.21989) \sqrt{\frac{1}{10} + \frac{(40 - 31.9)^2}{876.9}} \right) =$$

$$= (63.598 \pm 10.81792) = (52.78008, 74.41592) =$$

$$= \underline{(52.78, 74.42)} \quad (1)$$

io. We are 95% confident that an average lung damage of people who smoked for 40 years falls in interval (52.78, 74.42). (1)

k) 95% P.I. for y when $x_p = 40$

[5] $\Rightarrow \bar{y} = \underline{63.598} \quad (1)$

$$y \in \left(\bar{y} \pm t_{\alpha/2; n-2} \times s \sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{S_{xx}}} \right) =$$

$$= \left(63.598 \pm \underbrace{t_{0,025; 8}}_{2.306} \times (11.21989) \sqrt{1 + \frac{1}{10} + \frac{(40 - 31.9)^2}{876.9}} \right) =$$

$$= (63.598 \pm 28.04359) = (35.55441, 91.64159) =$$

$$= \underline{(35.55, 91.64)} \quad (1)$$

ie. We are 95% confident that the lung damage for an individual who smoked for 40 years falls in interval (35.55, 91.64). (1)

Q.2. a) $SSE = S_{yy} - \frac{(S_{xy})^2}{S_{xx}} =$
 $[5] = \left[\sum y_i^2 - \frac{(\sum y_i)^2}{n} \right] - \frac{\left[\sum x_i y_i - \frac{(\sum x_i)(\sum y_i)}{n} \right]^2}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} =$
 $n=12$
 $= \left[2854 - \frac{(178)^2}{12} \right] - \frac{\left[373 - \frac{(27)(178)}{12} \right]^2}{\left(64.5 - \frac{(27)^2}{12} \right)}$
 $= 12.000004 \approx 12$

and $SSE = SSPE + SSLF$, where $SSPE = 7.3333333$
 So $SSLF = SSE - SSPE =$
 $= 4.6666707$

[6] b) $y = \beta_0 + \beta_1 x + \epsilon$, $\alpha = 0.05$

H_0 : model is appropriate
 H_a : model is not appropriate

-test-stat: $F = \frac{MSLF}{MSPE} = \frac{SSLF / [n-2 - \sum (n_i-1)]}{SSPE / \sum (n_i-1)}$
 $= \frac{4.6666707 / (10-8)}{7.3333333 / 8} = \frac{2.3333354}{0.9166666} = 2.5454564$

-D.R. - we reject H_0 if $F > F_{\alpha, [n-2 - \sum (n_i-1), \sum (n_i-1)]} = F_{0.05, [2, 8]} = 4.46$

-since $F = 2.5454 < 4.46$, we do not reject H_0 and conclude that at 5% level of signif. there is not enough evidence to say that the lin. model is not appropriate.

Q.3:

a) [3] violation of the assumption of constant variance since the residual values are increasing with x.

b) [3] violation of the assumption of errors being normally distributed since the histogram of residuals is not bell-shaped nor is it symmetric (It is in fact negatively skewed).

c) [3] violation of linearity (or independence), since curve-linear pattern.

d) [3] no violations, since residuals are randomly scattered around their mean (i.e. no pattern).