

## 2013 Regression Rept

Q1A.  $b_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}$

$b_0 = \bar{y} - b_1 \bar{x}$

$E[y_i] = b_0 + b_1 x_i$

$b_1 = \frac{225864 - \frac{(1118)(3220)}{16}}{81452 - \frac{(1118)^2}{16}} = \frac{3486}{13377} = 0.2617 = 0.26$

$b_0 = \frac{3220}{16} - \frac{1118}{16} (0.26) = 182.9724 = 182.97$

$E[y_i] = 182.97 + 0.26x_i$

B.  $b_1$  - slope For each unit incre in the population density of the city, the mean distribution of robberies per 10,000 people incre by 0.26

$b_0$  - intercept: Sometimes a physical explanation is not possible, it is the mean distribution of robbery rate for population density of 0. In this context it does not make sense

C.  $MSE = \frac{SSE}{N-2} = \frac{S_{yy} - \frac{S_{xy}^2}{S_{xx}}}{N-2}$

$\sum y_i^2 - n\bar{y}^2 = \frac{(\sum x_i y_i - n\bar{x}\bar{y})^2}{\sum x_i^2 - n\bar{x}^2}$

SSR formula  $\sum (y_i - \hat{y})^2 = S_{yy} - \hat{\beta}_1 S_{xy}$

$\sum y_i^2 - n(\bar{y})^2 - (b_1 \sum x_i y_i - n\bar{x}\bar{y})$   
 $(649736 - 648025) - (0.26)(17422) = \frac{148441}{N-2} = 106.029$

2

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_1: \beta_1 \neq 0$$

$$t_{\text{calc}} = \frac{b_1 - \beta_1}{\text{se}(b_1)} = \frac{b_1}{\sqrt{\frac{\text{MSE}}{S_{xx}}}} = \frac{0.26}{\sqrt{\frac{106029}{3331.75}}} = \frac{0.26}{0.178}$$

$$t_{\text{calc}} = 1.457$$

$$t_{\text{critical}} (95\%, 14) = 2.141$$

$t_{\text{calc}} < t_{\text{critical}}$  Fail to reject  $H_0$ .  
There is not enough evidence to suggest higher robbery rate in more densely populated areas.

$$6. 95\% \text{ CI for } b_1 \quad b_1 \pm t_{\text{critical}} (\text{se}(b_1))$$

$$0.26 \pm 2.141 (0.178)$$

$$0.26 \pm 0.382$$

$$95\% \text{ CI for } b_0 \quad b_0 \pm t_{\text{critical}} (\text{se}(b_0))$$

$$\text{se}(b_0) = \sqrt{\text{MSE} \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} = \sqrt{106029 \left[ \frac{1}{16} + \frac{1249424}{3331.75} \right]}$$

$$182.97 \pm 2.141 (194.45)$$

$$182.97 \pm 417.354$$

3  
2013 Regression Exam Paper

Q2

A. Source	D.F.	SS	MS	F
Regression	1	4.9	4.9	4.9/1.1
Error	9	1.1	1.1/9	
Total (corrected)	10	6.0		

B.  $F(0.95, 1, 9) = F_{crit} = 5.3177$

$H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$

$F_{calc} = 5.3177$

$F_{calc} = 40.090$

$F_{calc} > F_{crit}$

reject  $H_0$ , relationship exists between advertising expenditure and sales revenue

C. No we cannot say whether it is strictly increasing/decreasing.

The F ratio/test is used as a ration of two mean square values, this gives no indication about positive/negative relationship only the relationship

The F figure can never be negative as it is a squared number

D.  $SSE = \sum (y_i - \hat{y}_i)^2$  Error Sum of Squares

This is the squared deviation of the predicted value from the observed by value. By the method of least squares we aim to minimize this value (variance) so as to get as close to the observed value as possible

E. - r=0 expresses a positive linear relationship

If correlation = 1 we have a perfect relationship.

The change in the 2 variables move opposite in same direction

No linear relationship corresponds to  $r=0$ . The two variables are uncorrelated. This does not imply independence

Warning!

- Does not capture nonlinear relationships
- A third variable may be factor in both
- Can be coincidental
- Correlation does not imply causation
- It shows direction but not slope  $\rightarrow$  job for regression



5.

## 2013 Regression Exam Paper

### Q5 A Assumptions:

1.  $X_i$  is the  $i$ th value of the predictor variable, which is known constant for all  $i$ .
2. The observations  $y_i$  or  $E_i$  are independent.
3. At any given  $X_i$ ,  $y_i$  or  $E_i$  are normally distributed.
4. The observations  $y_i$  or  $E_i$  have constant standard deviation.
5. The mean of  $y$  can be joined by a straight line.

$$E(y_i) = \beta_0 + \beta_1 x_i$$

$\beta_1$ : Slope of regression line, indicates the mean value of  $y$  distributed for one unit increase in  $x$ .

$\beta_0$ : Intercept - value of  $x=0$  for a particular application, or  $\beta_0$  gives the mean distribution of  $y$  at  $x=0$ . Not always possible to have an explanation.

B  $H_0: \beta_1 = 0$

vs  $H_1: \beta_1 \neq 0$

$$T_{calc} = \frac{b_1 - \beta_1}{\text{se}(b_1)} = \frac{0.7 - 0}{\frac{\text{MSE}}{\sqrt{57.5}}} = \frac{0.7}{\frac{0.61}{\sqrt{10}}} = 2.834$$

$T_{calc} > 2.084$  reject  $H_0$

$T_{calc} > T_{critical}$  reject  $H_0$

C  $10 = b_0$

$0.5 = b_1$

Then if a 0.5 increase in the mean value for the distribution for  $y$  for each one unit increase of  $x$ .

$$\hat{y} = 10 + 0.5(10) = 10 + 5 = 15 \text{ correct } \checkmark$$

D The mean do seem to be zero  $\checkmark$

Constant variance is a given

Seem to be funnel shaped  $\leq$  possibly a mixture

The data seem to be funnel shaped indicating a non constant variance and non independent

## Regression Exam Paper 2013

$$\begin{aligned} & \sum (x_i - \bar{x})^2 \\ & \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) \\ & \sum x_i^2 - 2\sum x_i\bar{x} + \sum \bar{x}^2 \\ & \sum x_i^2 - 2\sum x_i\bar{x} + n\bar{x}^2 \\ & \quad - 2n\bar{x}\bar{x} + n\bar{x}^2 \\ & \sum x_i^2 - n\bar{x}^2 \end{aligned}$$

1. A.  $E[y_i] = b_0 + b_1 x_i$

$$b_0 = \bar{y} + b_1 \bar{x}$$

$$b_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}$$

$$b_1 = \frac{225864 - \frac{1118(3270)}{16}}{81452 - \frac{(1118)^2}{16}} = \frac{-827608715}{333175}$$

$$= 0.26157 = 0.26 = b_1$$

$$b_0 = \bar{y} + 0.26 \bar{x} = \frac{3270}{16} + 0.26 \frac{(1118)}{16} = 219.5275 = 219.53$$

$$y_i = 219.53 + 0.26 x_i$$

b.  $b_0$  can be interpreted as the robbery rate when the population is zero, however this does not make sense in this case as we cannot rob if there is no population.

$b_1$  is the increase in the mean distribution of robbery rate per (100,000) per unit increase in the population.

C MSE is  $\frac{SSE}{n-2}$   $\frac{\sum y_i - \bar{y}}{n-2}$

$\sum y_i = \sum \bar{y}$

$\sum y_i^2 - 2 \sum y_i \bar{y} + \sum \bar{y}$

$\sum y_i^2 - 2 \sum y_i^2 + \sum y_i^2$

$= \frac{\sum y_i - \sum y_i^2}{n-2}$

$= 3220 -$

$SSE = \sum y_i^2 - n\bar{y}^2 - \frac{(\sum x_i y_i - n\bar{x}\bar{y})^2}{\sum x_i^2 - n\bar{x}^2}$

$649736 - 16(40501585) - \frac{759512.25}{3331.75}$

$1711 - 227.9619$

$= \frac{1483.03843}{14}$

$MSE = 105.931288 = 105.93$

One tail

D  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$

$T_{calculated} = \frac{0.26}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{0.26}{\sqrt{\frac{105.43}{3331.75}}} = 0.46018$

$1.457$

$T_{critical} (0.05, 14) = 2.1447$

$|t| \leq t_{critical}$

no evidence against  $H_0$

$0.26 \pm 1.761(0.5650)$

$0.26 \pm 1.212$

0 lie on interval

### 3 Regression Exam Paper 2013

E.  $B_1: 0.26 \pm 1.212$

$B_0: 21953 \pm 2.145 \left( \sqrt{MSE \left[ \frac{1}{n} + \frac{4882.11625}{331.75} \right]} \right)$

$21953 \pm 2.145 (12.7222)$

$21953 \pm 27.289$

Q2 A

Source	DF	SS	MS	F
Regression	1	4.9	4.9	40.0049
Error	9	1.1	0.1222	
Total	10	6.0		

B  $F_{0.05, 1, 9} = 5.3177$   
 $F_{calc} > F_{critical}$  we reject  $H_0$ , there is a relationship between variables

C We cannot say, we can use of ANOVA table to calculate  $R^2$ , then the square root of this figure is the r-value, the correlation coefficient, which indicate either a strictly positive or negative relationship, because of the square root, we do not know if it is positive or negative

D SSE is the error sum of squares. This is actually the error term  $e_i$ , our model aim to minimize this value, it is the squared deviation of the predicted value from the actual value. By method of least squares this value should be at a minimum



6. For correlation  $r > 0$  we have positive linear correlation. If correlation  $= 1$ , we have perfect positive linear correlation. The change in the 2 variables move approx in same direction.

No linear relationship corresponds to  $r = 0$ . The two variables are uncorrelated. This does not imply that variables are independent.

- It does not capture nonlinear relationship
- It measures direction but not slope  $\rightarrow$  job is regression

Q 5.4 Assumptions:

$\epsilon_i$  is a random error term with properties:

- $E(\epsilon_i) = 0$
- $Var(\epsilon_i) = \sigma^2$
- $\epsilon_i$  and  $\epsilon_j$  are uncorrelated so  $Cov(\epsilon_i, \epsilon_j) = 0$  for all  $i, j \neq i$ .
- $\epsilon_i$ 's are normally distributed

B.  $N = 23$   $b_1 = 0.7$   $MSE = 0.61$   $S_{xx} = 10$   
 $H_0: B_1 = 0$  vs  $H_1: B_1 \neq 0$   
 $t_{calc} = \frac{0.7}{\sqrt{\frac{0.61}{10}}} = 2.834217$

$t_{critical} (0.95, 21) = 2.080$

$|t| > t_{critical}$

$\Rightarrow$  reject  $H_0$  there is a relationship

C. Conclusion: The mean value of the  $y$  distribution increase by 0.7 units each increase in  $x$  of unit by  $x$ .

But (or) when  $x$  is 10  $y = 15$ .

5.

Regression Exam Paper 2013

Q52. mean is zero ✓.  
constant variance ✓ in question with 3 values above  
the  $\pm 3\sigma$  range.

The data seems to be funnel shaped indicating a  
non-constant variance and non independent