

## 2012 Regression Exam Part

B Sample size = 20

Mean of Sample = 2.1

$$SE(\text{mean}) = \frac{1.552}{\sqrt{20}} = 0.347$$

$$95\% \text{ CI} = 2.1 \pm 0.347(1.96) \\ = (1.374, 2.826)$$

95% confident that the true mean value ( $\mu$ ) is in (1.374, 2.826)

$$D \text{ CI } \hat{y} \pm t_{\text{critical}} \sqrt{MSE \left[ \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right]}$$

$$PI \hat{y}_i \pm t_{\text{critical}} \sqrt{MSE \left[ 1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right]}$$

- CI is model for estimating the mean value of  $y$ ,  $E(y)$  for a specific value of  $x$ .

- PI predicts a particular  $y$  value for a given  $x$ , predicting outcomes of single experiment given  $x$  value.

CI is narrower why?

- The error in estimating the mean value of  $y$ ,  $E(y)$ , for a given  $x$ , is the distance between the least squares line and the true line of mean  $E(y) = \beta_0 + \beta_1 x_i$ .

- Error shown in fit  $[ \hat{y} - E(y) ]$

- In contrast the error  $(y_i - \hat{y}_i)$  in predicting some future value of  $y$  is the sum of 2 errors: the error of estimating the mean of  $y$ ,  $E(y)$  plus the random error that is a component of the value of  $y$  to be predicted.

- Consequently the error of predicting a particular value of  $y$  will always be larger than the error of estimating the mean value of  $y$  for particular  $x$ .

- The further  $x_p$  lies from  $\bar{x}$ , the larger will be the errors of estimation and prediction.

$$2A \quad \begin{aligned} \sum x_i &= 2150 & \sum y_i &= 1430 & n &= 5 \\ \bar{x} &= 430 & \bar{y} &= 286 \\ \sum x_i y_i &= 618500 \\ \sum x_i^2 &= 931100 & \sum y_i^2 &= 411900 \end{aligned}$$

$$b_1 = \frac{\sum x_i y_i - \frac{\sum x_i \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}} = \frac{618500 - \frac{2150(1430)}{5}}{931100 - \frac{(2150)^2}{5}} = \frac{6}{11}$$

$$b_0 = \bar{y} - \bar{x} b_1 = 286 - \frac{6}{11} 430 = 514.5$$

$$E[y_i] = 514.5 + \frac{6}{11} x_i$$

$$B \quad MSE = \frac{SSE}{N-2} = \frac{\sum (y_i - \hat{y}_i)^2}{N-2} = \frac{S_{yy} - b_1 S_{xy}}{n-2}$$

$$\begin{aligned} &= \frac{\sum y_i^2 - n(\bar{y})^2 - (b_1 (\sum x_i y_i - n\bar{x}\bar{y}))}{n-2} \\ &= \frac{411900 - 5(286)^2 - (\frac{6}{11} (618500 - 2150 \cdot 1430))}{3} \\ &= \frac{10720}{3} = 3573.33 \end{aligned}$$

$$C \quad r^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}} = \frac{(260)^2}{680(290)} = 0.6724$$

-  $R^2$  measures the proportion of total variance about the mean  $\bar{y}$ , explained by the regression

-  $R^2$  explains 67% of the total variance in the data about the average  $\bar{y}$ , goes from 0 to 1

D.  $0.6724 = R^2$  is a strong figure, roughly 70% of error is due to variance about the mean  $\bar{y}$ .  
- With each increase in number of men, the variance decreases at a rate proportional to  $1/n$ .

## 2012 Regression Exam Prep

## 2 E. Assumptions:

- $X_i$  is the  $i^{th}$  value of the predictor variable which is a known constant for all  $i$ .
  - The observations  $y_i$  or  $\epsilon_i$  are independent.
  - At any given  $X_i$ ,  $y_i$  or  $\epsilon_i$  is normally distributed.
  - The observations or  $\epsilon_i$  have constant variance.
  - Mean of  $y_i$  can be joined by a straight line:
 
$$E(y_i) = \beta_0 + \beta_1 x_i$$
- $\beta_1$  - slope of regression  
 $\beta_0$  - intercept

## Q3A. Test if there is a relationship between 2 variables

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

$$T_{stat} = -2.55$$

$P_{value} = 0.008$  tell him  $\alpha$  of 0.05, reject  $H_0$ , a relationship exists between the 2 variables

B. Slope =  $-0.00231$ , the mean distribution rate of Sweetbushes reduced by  $0.00231$  for each unit increase in price

Intercept = 6.25, this is the value of the Sweetbush rate for a price value of 0.

$$C. 6.25 - (0.00231)(210) = 5.67 \text{ Sweetbush max}$$

D.  $R^2$  value of 22.4 suggests a bad fit for the model. Proportion of total variability due to regression is 0.23 implying a bad fit  $\rightarrow$  unreliable prediction



## Regression 2012 Exam Paper

Variable	N	mean	std dev	SE Mean	95% CI
Benzar	20	2.1	1.552	0.307	(1.374, 2.826)

- Sample size  $N=20$

- The mean of the sample is 2.1

$$SE_{mean} = \frac{1.552}{\sqrt{20}} = 0.307$$

95% of time the true mean lies in (1.374, 2.826)

This is created by  $mean \pm t_{(0.95, 19)} SE(\text{mean})$

### D Confidence versus prediction interval

$$CI: \bar{y} \pm t_{(n-2, 1-\frac{\alpha}{2})} S \cdot \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

$$PI: \bar{y} \pm t_{(n-2, 1-\frac{\alpha}{2})} S \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

- CI use model for estimating the mean value of  $y$ ,  $E(y)$  for a specific value of  $x$ .

- 2<sup>nd</sup> model PI, predicts a particular  $y$  value for a given  $x$ . Predicts outcome of single experiment given  $x$ -value

### CI is narrower why?

- The error in estimating the mean value of  $y$ ,  $E(y)$ , for a given  $x_p$ , is the distance between the least square line and the true line of mean,  $E(y) = \beta_0 + \beta_1 x_i$
- Error shown in figure  $[y - E(y)]$

- In contrast, the error  $(y_p - \hat{y})$  in predicting some future value of  $y$  is the sum of 2 errors: the error of estimating the mean of  $y$ ,  $E(y)$  plus the random error that is a component of the value of  $y$  to be predicted

2

-Consequently the error of predicting a particular value of  $y$  will always be larger than the error of estimating the mean value of  $y$  for particular  $x$

-The further  $x_p$  lies from  $\bar{x}$ , the larger will be the error of estimation and prediction

Q2 a  $b_1 = \frac{y - b_0}{x}$

$n = 5$

$\sum x_i = 2150$        $\bar{x} = 430$

$\sum y_i = 1430$        $\bar{y} = 286$

$\sum x_i y_i = 618500$

$\sum x_i^2 = 931100$        $\sum y_i^2 = 411900$

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

$$\frac{618500 - \frac{2150(1430)}{5}}{931100 - \frac{(2150)^2}{5}}$$

$$\frac{3600}{6600} = \frac{6}{11} = 0.5454$$

$$b_0 = 286 + \frac{6}{11}(430) = 51.4545$$

51.45

$$y_i = 51.45 + \frac{6}{11} x_i$$

b  $MSE = \frac{SSE}{n-2}$        $SSE = S_{yy} - \beta_1^2 S_{xx}$

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

$$\sum y_i^2 - n(\bar{y})^2 = b_1 [\sum x_i y_i - n\bar{x}\bar{y}]$$

$$411900 - 5(81796) = \frac{6}{11} [618500 - 5(430)(286)]$$

$$2920 = \frac{21600}{11}$$

$$= \frac{956.363}{3} = 318.787 = \frac{10520}{33} = 5^2$$

c  $r^2 = \frac{S_{xy} \times S_{xy}}{S_{xx} \times S_{yy}} = \frac{(3600)^2}{(6600)(2920)}$

$$\frac{12960000}{1927200} = \frac{540}{803} = 0.6726$$

### 3 Regression 202 Exam Page

2.1.  $R^2$  mean proportion of total variance about the mean  $\bar{y}$ , explained by the regression.

$r^2$  explains  $G\%$  of the total variation in the data about the average  $\bar{y}$ . Goes for  $0 \rightarrow 1$ .

If  $r^2$  was 20% I would be concerned, 80% of the error accounted for is due to chance, the other 80% is due to our model.

2.2  $67.24 = R^2$  is a decent figure, roughly 70% of error is due to variation about the mean  $\bar{y}$ .

With each increase in male value, the mean value for distribution of sale price rises by  $b_1$ . The  $b_1$  to increase the sale price by one, the number value would need to increase by  $1/b_1$ .

#### E. Assumptions:

- $x_i$  is the  $i$ th value of the predictor variable, which is a known constant for all  $i$ .
- $\epsilon_i$  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  $i, j \neq i$ .
  - $\epsilon_i$  are normally distributed  $N(0, \sigma^2)$ .

- The means of  $y_i$  can be joined by a straight line, given as:

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

where  $\beta_0$  and  $\beta_1$  are unknown parameters such that:

- $\beta_1$  is the slope of regression line and indicates the change in mean of  $y$  for unit increase in  $x$ .



-  $B_0$  is the intercept of the regression model.  $B_0$  gives the mean distribution of  $y$  at  $x=0$ .

Q3 Test if there is a relationship between the two variables or not.  
 $H_0: \beta_1 = 0$  vs.  $H_1: \beta_1 \neq 0$   
 $p\text{-value} = 0.018$  (less than 0.05, reject  $H_0$ , there is a relationship between Sweetroll index and price.  
Are related at price data.

b) Slope is  $-0.00231$ , for each unit increase in price, the mean dollar distribution of Sweetroll decreases by  $0.00231$  units.

$b_0 = 6.2521$ , this is the Sweetroll index value when price is set to zero.

c) 
$$= 6.25 - 0.00231(250) = 5.67235 \text{ Sweetroll}$$

$R^2$  value of 22.9 suggests a bad fit for the model.  
Proportion of total variability about the mean is explained about the regression line is 22.9, implying a bad fit  $\rightarrow$  unrealistic price.

d)