

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 $\%$ 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 611. This is increase the sale price by one, the number value would need to increase by 2.

E. Assumptions:

- x_i is the i th value of the predictor variable, which is a known constant for all i .
- ϵ_i is a random error term with properties:
 - $E(\epsilon_i) = 0$.
 - $\text{Var}(\epsilon_i) = \sigma^2$.
 - ϵ_i and ϵ_j are uncorrelated so $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i, j \neq i$.
 - ϵ_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 β_0 and β_1 are unknown parameters such that:

- β_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 unreliable prediction.

d)

Regression 2011 Exam Paper

a. Assumptions:

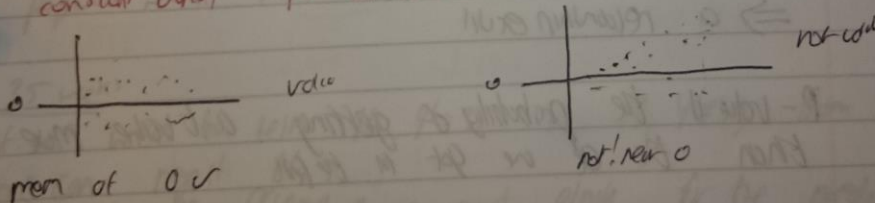
- X_i is the i th value of the predictor variable, which is a known constant for all i .
- The observations y_i (or ϵ_i) are independent.
- At any given X_i , y_i (or ϵ_i) is normally distributed.
- The observations y_i (or ϵ_i) have constant standard deviation.
- The mean of y_i can be joined by a straight line given as:

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

where β_0 and β_1 are unknown parameters

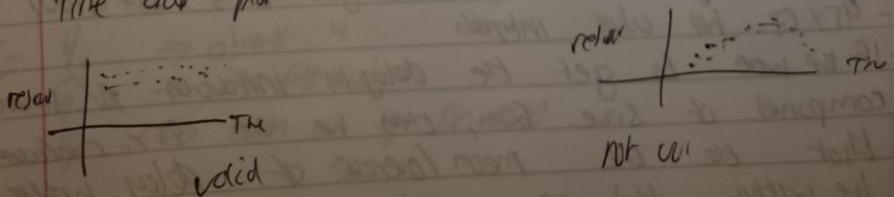
- β_1 is the slope of the regression line and indicates the change in the mean of y for one unit increase in x .
- β_0 is the intercept of the regression model. If it is sensible to think of a value of $X=0$ for a particular application, then β_0 gives the mean of the distribution of y at $X=0$, it is not always possible to have a physical explanation for each parameter.

B. 1. constant Varian for ϵ_i for all:



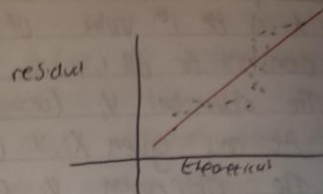
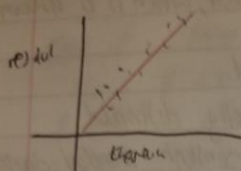
2. Independent of ϵ_i

Time order plot



2

3. Normality of ϵ_i
Q-Q plot



1. $H_0: \beta_0 = 0$ vs $H_1: \beta_0 \neq 0$

Compare t_{calc} with t_{crit}

$$t_{critical} = 2.306$$

$|t_{calc}| = 11.69 \Rightarrow t_{calc} > t_{critical}$, reject H_0 , thus we reject that the time taken to adopt innovation is 0 for company of size = 4

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

$$t_{critical} = 2.306$$

$$|t_{calc}| = 4.43 \Rightarrow t_{calc} > t_{critical}$$
 reject H_0

Reject that there is a zero change in time taken to adopt innovation for a unit increase in size
 \Rightarrow a relationship exists

P-value is the probability of getting a t-value more extreme than the one we got in the test.

- 1) The point estimate for the average value of delay in innovation for company of size 60 is given by $\hat{y} = b_0 + b_1(60)$
- 95% CI has usual interpretation
- If we were to get the delay in innovation for several companies of size 60, then we are 95% confident that the true mean (average of delay) would lie within this interval

Regression ? 2011 Exam Part

A point estimate for the delay time for a 'single' company of size 60m is also $y' = b_0 + b_1(60)$

But a 95% PI refers to the interval within which the true delay time of that single company of size=60m will lie with 95% confidence

They are different since the variability of delay time for a single firm is much greater than the variability of the average

Single observation = mean + change variance

This is why 95% PI are wider than 95% CI.

A 60m company took: Single observation

20 weeks $PI = (4.93, 24.02)$

If we were to test for $H_0: y' = 20$ vs $H_1: y' \neq 20$ then from PI we see that null hypothesis will not be rejected

$\rightarrow 20$ can be accepted as a good point estimate for the population

11th 35 weeks

If we were to test $H_0: y' = 35$ vs $H_1: y' \neq 35$, from PI

we can reject H_0

-35 cannot be accepted as a good estimate for the prediction of delay time

- $y_i \rightarrow$ observed y
- $\hat{y}_i \rightarrow$ fitted value
- $\bar{y} \rightarrow$ mean of y_i for all y

$$SSTO = SSR + SSE$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

$$R^2 = \frac{SSR}{SST}$$

R^2 is the proportion of total variability about \bar{y} explained by the regression line