

## 12/13 Regression

①

2011/2011

1. D. The point estimate for the average value of delay in innovation for company of size 60m is given by

$$y' = b_0 + b_1(60)$$

A 95% CI here has the usual interpretation

If we were to get the delay in innovation time for several companies of size 60m, then ~~the~~ we are 95% confident that the true mean/average of delay would lie within the interval

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 value is much greater than the variability of the average.

$$\text{Single observation} = \text{mean} + \text{change variation}$$

This is why P.I. are wider than 95% C.I.

A 60m euro company took: (single instance)  
1.20 weeks PI = (4.95, 24.02)

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

2

20 can be accepted as a good point estimate for <sup>True</sup> prediction

ii 35 weeks

- If we were to test  $H_0: \mu = 35$  vs  $H_1: \mu \neq 35$  then from PI we can reject  $H_0$

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

E.

$y_i \Rightarrow$  observation  $y$

$\bar{y} \Rightarrow$  sample mean

$\bar{y} \Rightarrow$  mean of  $y_i$  for all  $y$

$$SSTO = SSR + SSE$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\bar{y} - \bar{y})^2 + \sum_{i=1}^n (y_i - \bar{y})^2$$

$$R-sq = \frac{SSR}{SSTO}$$

It is the proportion of total variability about the mean  $\bar{y}$ , as explained by the regression line

2011 - 2012

Q3

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

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

For  $\beta_1$   $p = 0.008$  which is less than 0.05, reject  $H_0$ .  
 $\Rightarrow$  There is a relationship. Sweetness index and peltin are related as per the data.

C For peltin content = 250, Sweetness Index =  $6.25 - 0.00231(250)$

For reliability we look at  $r-sq$  value = 22.9%.  $R^2$  is really small  
 $\Rightarrow$  ~~proportion~~ Proportion of total variability about the mean as explained by the regression line = 0.23 implying a bad fit  
 Unrealistic prediction

Q3 a) 1 2012 Q3 Re

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

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

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

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

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

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

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

$$= SSE + SSR$$

b) Height = x      Weight = y

$\sum x_i = 765$	$\sum y_i = 285$	$\sum x_i y_i = 43956$
$\bar{x} = 153$	$\bar{y} = 57$	$SSC = 3646$
$\sum x_i^2 = 117903$	$\sum y_i^2 = 16425$	$n = 5$

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

$$\frac{\sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + n \bar{x} \bar{y}}{\sum x_i^2 - \bar{x} \sum x_i - \bar{x} \sum x_i + n \bar{x}^2}$$

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

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

$$b_1 = \frac{43956 - \frac{765 \cdot 285}{5}}{117903 - \frac{(765)^2}{5}} = \frac{351}{858} = \frac{9}{22} = 0.4090$$



- $b_1$  is an estimate of  $\beta_1$
- It is the slope of the regression line
- For each unit increase in height, there is a  $\frac{9}{22}$  increase in weight.

Intercept =  $b_0$

Calculated by  $\bar{y} - b_1 \bar{x}$   
 $57 - \frac{9}{22}(153) = -\frac{123}{22} = -5.5909$

- This is the value of weight when height is 0.
- However we cannot have a negative weight so we disregard this value.

ii.  $E[y] = -\frac{123}{22} + \frac{9}{22}x_i$

$x_i = 150 \Rightarrow = -\frac{123}{22} + \frac{9}{22}(150)$   
 $= \frac{1227}{22} \approx 55.77 \text{ (kg)}$

Prediction interval =  $\hat{y}' \pm t_{\text{critical}} (S.e(x'))$   
 $= \hat{y}' \pm 6.04753 \sqrt{\frac{\sum(y_i - \hat{y}_i)^2}{n-2} \left( 1 + \frac{1}{n} + \frac{(x' - \bar{x})^2}{\sum x_i^2 - n(\bar{x})^2} \right)}$

$= \frac{1227}{22} \pm 2.353 \sqrt{\frac{(36.44)^2}{3} \left( 1 + \frac{1}{5} + \frac{(150 - 153)^2}{117903 - 5(153)^2} \right)}$   
 $= \frac{1227}{22} \pm (2.353) \left( \frac{34837714}{117903 - 5(153)^2} \right) \left( 1 + \frac{1}{5} + \frac{9}{858} \right)$   
 $= \frac{1227}{22} \pm 49372.98197$   
 $(-49372.98197, 105741.5)$

$\frac{1227}{22} \pm 8.215513326$

$(47.56, 63.99) = \text{prediction interval for}$   
height = 150 and weight =  $\frac{1227}{22}$

3  
[1] 2012 Rg Q3

Prediction interval tell us where we can expect to see the next data point sampled. Assume that the data really are randomly selected from Gaussian Normal dist.

Collect a sample of data and calculate a prediction interval

then sample a new value from the population

If you do this many times, we expect that the next value will lie within that prediction interval in 95% of the cases.

Key point is that the prediction interval tells you about the distribution of values, not the uncertainty in determining the population mean.

Prediction interval must account for both the uncertainty in knowing the value of the population mean, plus data scatter.

So a prediction interval is always wider than a confidence interval

- A prediction interval defines a range of values within which a response is likely to fall given a specific value of a predictor.

- The uncertainty represented by a prediction interval includes not only the uncertainty (variance) associated with the population mean and the new observation, but the uncertainty associated with the regression parameter as well.

- Because of the uncertainty associated with the population mean and the new observation or in dependent

if the observed used to fit the model the uncertainty estimate must be combined with sum of squares to yield  $MS_{\text{error}}$

- iii. Confidence interval tells how well we have done the mean. Assume data is randomly sampled from normal dist. If we do this many times, and calculate a confidence interval of the mean from each sample, we expect 95% of these intervals to include the true value of population mean. CI tells you about the likely location of the true population parameter.

A CI is an interval associated with a parameter and is a frequentist concept. The parameter is assumed to be non-random but unknown, and the CI interval is computed from the data. Because data is random, the interval is random. A 95% CI will contain the true parameter with prob 0.95. With a large number of repeated samples, 95% of the intervals would contain the true parameter.

- iv. As explored in above answers, the CI will be a much narrower interval

Reading 1 in formula.

$$\text{For CI se part} = \sqrt{MSE \left[ \frac{1}{n} \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}$$

$$\text{For PI se part} = \sqrt{MSE \left[ 1 + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}$$

For PI formula the extra or more  $\sqrt{\quad}$  bigger and hence a larger interval as other figures remain the same



2012 PAPER 1 Q3 Regress DAVID WEITBRECHT

Q3 A. Show that  $SST = SSE + SSR$

$$\begin{aligned} \sum (y_i - \bar{y})^2 &= \sum (\hat{y}_i - \bar{y} + y_i - \hat{y}_i)^2 \\ &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 + 2(\hat{y}_i - \bar{y})(y_i - \hat{y}_i) \\ &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 + 2\sum (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) \\ &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 + 2\sum (\hat{y}_i y_i - \bar{y} y_i - \hat{y}_i^2 + \bar{y} \hat{y}_i) \\ &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 + 2\sum \hat{y}_i y_i - 2\bar{y} \sum y_i - 2\sum \hat{y}_i^2 + 2\bar{y} \sum \hat{y}_i \\ &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 + 2\sum \hat{y}_i y_i - 2\bar{y} \sum y_i - 2\sum \hat{y}_i^2 + 2\bar{y} \sum \hat{y}_i \\ &\Rightarrow \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 = SST \end{aligned}$$

bi. Calculate and interpret slope and intercept.

Height = X	Weight = Y	$\sum x_i y_i = 43956$
$\sum x_i = 765$	$\sum y_i = 67$	$SSE = 36.41$
$\bar{x} = 153$	$\bar{y} = 13.4$	$n = 5$
$\sum x_i^2 = 117903$	$\sum y_i^2 = 909$	

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i y_i - \bar{x} \bar{y} - \bar{x} y_i + \bar{x} \bar{y})}{\sum (x_i^2 - 2\bar{x} x_i + \bar{x}^2)}$$

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

$$\sum x_i y_i - \frac{\sum y_i \sum x_i}{n} - \frac{\sum x_i \sum y_i}{n} + \frac{n \sum x_i \sum y_i}{n^2}$$

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

$$\text{Sub in values} \Rightarrow \frac{43956 - \frac{765(285)}{5}}{117903 - \frac{(765)^2}{5}} = \frac{351}{858} = \frac{9}{22} = 0.4090$$

2 DW

- $b_1$  is an estimate of population parameter  $\beta_1$ .
- It is the slope of the estimated regression line
- For each unit increase in height, there is a  $\frac{9}{22}$  increase in weight.

Y intercept =  $b_0$

Calculated by  $\bar{y} - b_1 \bar{x}$

$$57 - \frac{9}{22}(153) = \frac{-123}{22} = -5.5909$$

- This is the mean weight when height = 0.
- This is theoretical as we cannot have a negative weight.

ii. Prediction Interval for height = 150cm 95% confidence

$$\text{Prediction Interval} = \hat{y} \pm t_{\text{critical}} \cdot \text{se}(x')$$

$$= \hat{y} \pm t_{0.025, 13} \sqrt{\frac{2(4-9)^2}{n-2} \sqrt{1 + \frac{1}{5} + \frac{(x' - \bar{x})^2}{(2x^2 + 18)^2}}}$$

$$E[y] = -\frac{123}{22} + \frac{9}{22}x_i$$

$$x_i = 150 \Rightarrow = -\frac{123}{22} + \frac{9}{22}(150) = \frac{1227}{22} \approx 55.7727$$

$$\Rightarrow \frac{1227}{22} \pm 2.353 \sqrt{\frac{36.41}{5} \sqrt{1 + \frac{1}{5} + \frac{(150-153)^2}{(117.93 - 5(153)^2)}}$$

$$\frac{1227}{22} \pm 2.353(3.4837719)(1.100222)$$

$$\frac{1227}{22} \pm 8.215513326$$

$$95 \text{ PI} = (47.56, 63.99) \quad \text{height} = 150 \text{ cm and weight} = \frac{1227}{22} \text{ kg}$$



3

## MANG SCI PAPER 1 Q3 2012 Din

Prediction Interval

Prediction interval tells us where we can expect to see the next data point sampled. Tells us about distribution of values, not the uncertainty in determining the population mean.

Predicting a particular  $y$  for a given  $x$ . Predicting outcome of a single experiment given  $x$ -value.

iii. Confidence Interval

Estimate the mean value of  $y$  for a given  $x$ .

Confidence Interval (CI) tells us how well we have determined the mean. We expect 95% of time that these intervals will include the true value of the population value.

iv. Difference

It will always be larger, in the formula, PI has an extra 1 under the square root with all the figures the same meaning it will always be wider.

The error in estimating the mean value of  $y$ ,  $E[y]$ , for a given  $x$  say  $x_p$ , is the distance between the least square line and the true line of means  $E[y] = \beta_0 + \beta_1 x$ .

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

DW

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

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