

8.30

15 Friday

January '99

9.00 S. Estimate the variance of  $u_i$  i.e.,  $\sigma_u^2$

9.00 Show that Unbiased estimator of  $\sigma_u^2$  is unknown  
 10.00 and constant i.e., the variance of disturbance  
 10.30 term is unknown.

11.00 Let the equation be  $y_i = \alpha + \beta x_i + u_i$

Mean value of  $y_i = \bar{y} = \alpha + \beta \bar{x} + \bar{u}$

11.30 Estimate value of  $y_i = \hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$   
 We know that, Error( $e_i$ ) = Actual value - Estimated value.

12.00

$$= y_i - \hat{y}_i$$

12.30

$$\text{or } e_i = \hat{y}_i - \hat{\beta} x_i$$

1.00 Now,  $y_i - \bar{y} = \beta(x_i - \bar{x}) + (u_i - \bar{u})$

1.30  $\hat{y}_i = \beta x_i + (u_i - \bar{u})$

2.00  $\therefore e_i = \beta x_i + (u_i - \bar{u}) - \hat{\beta} x_i \quad [\because \hat{y}_i = \hat{\beta} x_i]$

2.30  $= -(\hat{\beta} - \beta)x_i + (u_i - \bar{u})$

3.00  $E[e_i^2] = E[-(\hat{\beta} - \beta)x_i + (u_i - \bar{u})]^2$

3.30  $= (\hat{\beta} - \beta)^2 E[x_i^2] + E[(u_i - \bar{u})^2] - 2(\hat{\beta} - \beta) E[x_i(u_i - \bar{u})]$

4.00  $E[\sum e_i^2] = E[(\hat{\beta} - \beta) \sum x_i^2 + \sum (u_i - \bar{u})^2 - 2(\hat{\beta} - \beta) \sum x_i(u_i - \bar{u})]$

4.30

→ ①

5.00 Now,

i)  $E[(\hat{\beta} - \beta) \sum x_i^2] = \frac{\sigma_u^2}{\sum x_i^2} \cdot \sum x_i^2 = \sigma_u^2$

5.30

ii)  $E[(u_i - \bar{u})^2] = E[\sum u_i^2 + n\bar{u}^2 - 2\bar{u} \sum u_i]$

Evening

$= E[\sum u_i^2 + \frac{(\sum u_i)^2}{n} - 2 \frac{(\sum u_i)^2}{n}]$

$= E[\sum u_i^2 - \frac{(\sum u_i)^2}{n}]$

$= n\sigma_u^2 - \frac{n\sigma_u^2}{n}$

$\therefore E(u_i^2) = \sigma_u^2$   
 by assumption

$= \cancel{n\sigma_u^2} = \sigma_u^2(n-1)$

8.30

$$(iii) 2E[(\hat{\beta} - \beta) \varepsilon x_i (\bar{u} - \bar{u})]$$

January '99

$$= 2E \left[ \frac{\sum u_i x_i}{\sum x_i^2} (\bar{u} x_i - \bar{u} \sum x_i) \right]$$

10.00

$$\therefore \hat{\beta} = \beta + \sum u_i x_i \\ \therefore \hat{\beta} - \beta = \frac{\sum u_i x_i}{\sum x_i^2}$$

10.30

$$= 2E \left[ \frac{(\sum u_i x_i)^2}{\sum x_i^2} - \frac{\bar{u} \sum x_i \sum u_i x_i}{\sum x_i^2} \right]$$

11.00

$$= 2E \left[ \frac{(\sum u_i x_i)^2}{\sum x_i^2} \right]$$

11.30

$$\therefore \sum \varepsilon x_i = 0$$

12.00

$$= 2 \cdot \frac{\bar{u} \sum x_i^2}{\sum x_i^2}$$

12.30

$$= 2\bar{u}^2$$

1.00

$$\therefore E(\sum u_i x_i)^2 \\ = \bar{u}^2 x_1^2 + \bar{u}^2 x_2^2 + \dots + \bar{u}^2 x_n^2 \\ = \bar{u}^2 \sum x_i^2$$

1.30

Hence, subtracting the value of (i), (ii) and (iii)

we get,

$$E[\sum \varepsilon x_i^2] = \bar{u}^2 + n\bar{u}^2 - \bar{u}^2 - 2\bar{u}^2 \\ = (n-2)\bar{u}^2$$

2.00

$$\therefore E \left[ \frac{\sum \varepsilon x_i^2}{n-2} \right] = \bar{u}^2$$

$$E(\hat{\sigma}_u^2) = \bar{u}^2$$

$$\text{where, } \hat{\sigma}_u^2 = \frac{\sum \varepsilon x_i^2}{n-2}$$

$\therefore \frac{\sum \varepsilon x_i^2}{n-2}$  is an unbiased estimator of  $\bar{u}^2$

$n$  is the number of observations.

2 is the unknown parameter in this model.

Memo

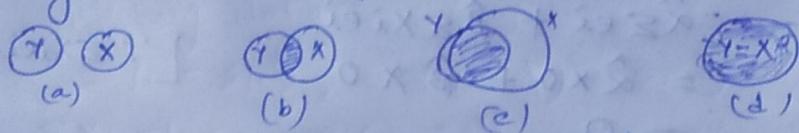
$$\therefore S^2 = \frac{1}{n-2} \sum (y_i - (\bar{y} - \hat{\beta} \bar{x}))^2 = \frac{1}{n-2} \sum (y_i - \hat{\beta} x_i)^2$$

$$\text{we know, } \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}, \sum x_i y_i = \hat{\beta} \sum x_i^2 = \frac{1}{n-2} \sum (y_i^2 + \bar{y}^2 - 2\bar{y} \sum x_i y_i)$$

$$\therefore S^2 = \frac{1}{n-2} \sum (y_i^2 - \hat{\beta} \sum x_i y_i)$$

## 'Goodness of Fit' / Coefficient of determination ( $R^2$ )

The coefficient of determination is a summary measure that tells how well the sample regression line fits the data.



In diagrams,  $\textcircled{Y}$  represents the variation of dependent variable and  $\textcircled{X}$  is the variation of explanatory or independent variable. The overlap of the two circles indicates the extent to which the variation in  $Y$  is explained by the variation in  $X$ . The greater the ~~overlap~~  $\rightarrow$  overlap, the greater the variation in  $Y$  is explained by  $X$ .

a) when there is no overlap,  $r^2$  or  $R^2 = 0$

In (b) and (c) as overlap increases  $\rightarrow$  explained variation of  $Y$  also increases.

(d) when the overlap is complete  $r^2$  or  $R^2 = 1$ .

means 100% of the variation of  $Y$  is explained by  $X$ .

The formula for,  $R^2 = \frac{\text{Explained variation}}{\text{Total variation}}$ .

Compute the value of  $R^2$

$$\text{If } Y_i = \alpha + \beta X_i + e_i \quad (22T)$$

$$\text{and } \hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i \quad (21) \text{ sol}$$

$$e_i = Y_i - \hat{Y}_i \quad 228 + 223 = 22T$$

$$\text{or, } Y_i = \hat{Y}_i + e_i \quad \text{sol at } 22T$$

$$\text{or, } Y_i - \bar{Y} = \hat{Y}_i - \bar{Y} + e_i \quad (\text{taking both sides by } -\bar{Y})$$

$$\text{or, } (Y_i - \bar{Y})^2 = \{(\hat{Y}_i - \bar{Y}) + e_i\}^2 \quad (\text{both sides take perfect square})$$

$$\text{or, } \hat{Y}_i^2 = \hat{Y}_i^2 + e_i^2 + 2\hat{Y}_i e_i \quad (ii)$$

$$\text{or, } \sum \hat{Y}_i^2 = \sum \hat{Y}_i^2 + \sum e_i^2 + 2 \sum \hat{Y}_i e_i \quad (\text{taking both sides by } \Sigma)$$

$$\begin{aligned}
 \text{(a) Now } \sum \hat{y}_i e_i &= \sum (\hat{y}_i - \bar{y}) e_i \\
 &= \sum y_i e_i - \bar{y} \sum e_i \\
 &= \sum (\hat{\alpha} + \hat{\beta} x_i) e_i - 0 \quad [\because \sum e_i = 0] \\
 &= \hat{\alpha} \sum e_i + \hat{\beta} \sum x_i e_i \\
 &= \hat{\alpha} \times 0 + \hat{\beta} \times 0 \quad [\because \sum x_i e_i = 0]
 \end{aligned}$$

$$\begin{aligned}
 \text{and } \sum \hat{y}_i^2 &= \sum (\hat{y}_i - \bar{y})^2 \\
 &= \sum (\hat{\alpha} + \hat{\beta} x_i - \bar{y} - \bar{\beta} \bar{x})^2 \\
 &= \hat{\beta}^2 \sum (x_i - \bar{x})^2 \\
 &= \hat{\beta}^2 \sum x_i^2 \Rightarrow \text{Explained variation} \\
 \text{and } \sum e_i^2 &= \sum (y_i - \hat{y}_i)^2 \Rightarrow \text{Unexplained variation} \\
 \text{and } \sum y_i^2 &= \sum (y_i - \bar{y})^2 \Rightarrow \text{Total variation}
 \end{aligned}$$

Therefore we can write that

$$\begin{aligned}
 \sum y_i^2 &= \sum \hat{y}_i^2 + \sum e_i^2 \\
 \text{or, } \sum y_i^2 &= \hat{\beta}^2 \sum x_i^2 + \sum e_i^2 \\
 \text{Total variation} &= \text{Explained variation} + \text{Unexplained variation}
 \end{aligned}$$

$$\begin{aligned}
 \text{or, Total sum of Squares} &= \text{Explained sum square} + \text{Unexplained sum square} \\
 (\text{TSS}) &\qquad\qquad\qquad (\text{ESS}) \qquad\qquad\qquad (\text{Sum squares})
 \end{aligned}$$

$$\begin{aligned}
 \text{or, } \text{Var}(Y_i) &= \text{Var}(\hat{Y}_i) + \text{Var}(e_i) \\
 \therefore \text{TSS} &= \text{ESS} + \text{RSS}
 \end{aligned}$$

- i) TSS represents the Total sum of squared deviations from  $\bar{Y}$ .
- ii)  $\text{ESS} = \hat{\beta}^2 \sum x_i^2$  represents the estimated effect of  $X$  on the variations in  $Y$ .
- iii)  $\text{RSS} = \sum e_i^2$  represents the variations in  $Y$  which unexplained by the estimated relationship.

This decomposition of total variations in  $y$  leads to a measure of the "goodness of fit"/( $R^2$ ) coefficient of determination.

$$\text{where, } R^2 = \frac{\text{Explained variation}}{\text{Total variation}} = \frac{\text{Var}(\hat{y}_i)}{\text{Var}(y_i)}$$

$$= \frac{\hat{\beta}^2 \sum x_i^2}{\sum y_i^2}$$

Prove that  $R^2$  lies between +1 and -1

$$\text{Since, } \text{Var}(y_i) = \text{Var}(\hat{y}_i) + \text{Var}(e_i)$$

$$\text{and } 0 \leq \text{Var}(\hat{y}_i) \leq \text{Var}(y_i)$$

$$\text{or, } 0 \leq \frac{\text{Var}(\hat{y}_i)}{\text{Var}(y_i)} \leq 1$$

$$\text{or, } 0 \leq R^2 \leq 1 .$$

$R^2 = 0$  when  $\text{Var}(\hat{y}_i) = 0$  i.e.,  $\sum e_i^2 = \sum y_i^2$   
 and  $R^2 = 1$  when  $\text{Var}(\hat{y}_i) = \text{Var}(y_i)$  i.e.,  $\sum e_i^2 = 0$

100

106

108

Econometrics.H.W.

- i. A sample of 20 observations corresponding to regression model  $Y_i = \alpha + \beta X_i + u_i$ , gives the following data:

$$\sum Y_i = 210.9, \quad \sum (Y_i - \bar{Y})^2 = 86.9$$

$$\sum X_i = 186.2, \quad \sum (X_i - \bar{X})^2 = 215.4, \quad \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

- i) Obtain the estimated value of  $\hat{\alpha}$  &  $\hat{\beta}$ .
- ii) Estimate  $\text{Var}(\hat{\beta})$  &  $\text{Var}(\hat{\alpha})$ .
- iii) Find the value of  $R^2$ .

$$\therefore R^2 = 0.935.$$

This suggests that 93.5 per cent of the sample observations of  $Y$  can be attributed to the variations of the fitted value of  $Y$  i.e.,  $\hat{Y}_i$  or we say that our regression line fits the given data well.

Thus  $R^2$  measures the proportion of variations in the dependent variable that is explained by the independent variables.

**Example : 2.5.** A sample of 20 observations corresponding to the regression model  $Y_i = \alpha + \beta X_i + u_i$  where  $u_i$  is normally distributed with mean zero and unknown variance  $\sigma_u^2$ , gives the following data :

$$\sum_{i=1}^n Y_i = 21.9, \quad \sum_{i=1}^n (Y_i - \bar{Y})^2 = 86.9, \quad \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = 106.4$$

$$\sum_{i=1}^n X_i = 186.2, \quad \sum_{i=1}^n (X_i - \bar{X})^2 = 215.4, \quad n = 20$$

Obtain the usual regression results.

**Solution :** On the basis of the given information we have to fit a linear relation between  $Y$  (dependent variable) and  $X$  (explanatory variable).

(i) Estimation of  $\hat{\alpha}$  and  $\hat{\beta}$  :

$$\text{We know that, } \hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\therefore \hat{\beta} = \frac{106.4}{215.4} = 0.494 \quad \text{Where } \bar{Y} = \frac{\sum_{i=1}^n Y_i}{n} = \frac{21.9}{20} = 1.095$$

$$\text{and } \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

$$= 1.095 - 0.494 \times 9.31 \quad \text{and } \bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{186.2}{20} = 9.31 \\ = 1.095 - 4.60 = -3.505$$

Thus we have,  $\hat{\alpha} = -3.505$  and  $\hat{\beta} = 0.494$  and our estimated regression line is :

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i \Rightarrow \hat{Y}_i = -3.505 + 0.494 X_i$$

(ii) Estimation of variances :

Since we know that,  $\text{var}(\hat{\alpha}) = \sigma_u^2 \left( \frac{\sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2} \right)$  and  $\text{var}(\hat{\beta}) = \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n x_i^2}$ .

Here we see that  $\sigma_u^2$  is not known and hence we replace it by its

unbiased estimator  $\hat{\sigma}_u^2 = \sum_{i=1}^n e_i^2 / n - 2$ .

Thus we have,  $\text{var}(\hat{\alpha}) = \hat{\sigma}_u^2 \left( \frac{\sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2} \right)$  and  $\text{var}(\hat{\beta}) = \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n x_i^2}$

Again we know that,  $\sum_{i=1}^n e_i^2 = \sum_{i=1}^n y_i^2 - \hat{\beta} \sum_{i=1}^n x_i^2$

$$\therefore \sum_{i=1}^n e_i^2 = 86.9 - (0.494)^2 \times 215.4$$

$$= 86.9 - 52.56 = 34.34$$

$$\text{Now } \hat{\sigma}_u^2 = \sum_{i=1}^n e_i^2 / n - 2$$

$$= \frac{34.34}{20-2} = \frac{34.34}{18} = 1.908$$

$$\text{Where } \sum_{i=1}^n y_i^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\text{and } \sum_{i=1}^n x_i^2 = \sum_{i=1}^n (X_i - \bar{X})^2$$

$$\text{Now } \text{var}(\hat{\alpha}) = \hat{\sigma}_u^2 \left( \frac{\sum_{i=1}^n X_i^2}{n \sum_{i=1}^n x_i^2} \right) = \frac{1.908 \times 1948.922}{20 \times 215.4} = 0.8631$$

$$\left[ \therefore \sum_{i=1}^n (X_i - \bar{X})^2 = 215.4 \text{ or, } \sum_{i=1}^n X_i^2 - n \bar{X}^2 = 215.4 \right]$$

thus we conclude that OI

100

106

108

111

### The Simple Linear Regression Model

$$\begin{aligned}
 \text{or, } \sum_{i=1}^n X_i^2 &= 215 \cdot 4 + n \bar{X}^2 = 215 \cdot 4 + 20 \times (9 \cdot 3)^2 \\
 &= 215 \cdot 4 + 1733 \cdot 522 \\
 &= 1948 \cdot 922 ] \\
 \therefore \text{var}(\hat{\alpha}) &= 0 \cdot 8631
 \end{aligned}$$

$$\text{Similarly, } \text{var}(\hat{\beta}) = \frac{\hat{\sigma}_u^2}{\sum_{i=1}^n x_i^2} = \frac{1 \cdot 908}{215 \cdot 4} = 0 \cdot 0089$$

$$\text{Now, } SE(\hat{\alpha}) = \sqrt{\text{var}(\hat{\alpha})} = \sqrt{0 \cdot 8631} = 0 \cdot 929$$

$$SE(\hat{\beta}) = \sqrt{\text{var}(\hat{\beta})} = \sqrt{0 \cdot 0089} = 0 \cdot 094$$

#### (iii) Construction of confidence intervals :

Now we like to set up a confidence interval for  
 $\alpha = 0 \cdot 95$  (i.e., 5% level of significance) and (b)  $P = 0$   
 of significance)

In other words, we like to find the value of ' $t$ '  
 $0 \cdot 025$  and (b)  $0 \cdot 995$ .

|           |     |
|-----------|-----|
| 1 matters | 100 |
| as        | 106 |
|           | 108 |
|           | 111 |

Contents

7

### The Simple Linear Regression Model

59

In terms of our earlier example (Example 2.5) the estimated regression results can be written as :

$$Y_i = -3.505 + 0.494X_i \quad R^2 = 0.6048.$$

(0.929) (0.094)

Here  $\hat{\alpha} = -3.505$ ,  $\hat{\beta} = 0.494$ ,  $SE(\hat{\alpha}) = 0.929$ ,  $SE(\hat{\beta}) = 0.094$ ,

$$\text{and } R^2 = \frac{\hat{\beta}^2 \sum_{i=1}^n x_i^2}{\sum_{i=1}^n y_i^2} = \frac{(0.494)^2 \times 215.4}{86.9} = \frac{52.5653}{86.9} = 0.6048$$

This suggests that 60.48 per cent of the sample observations can be attributed to the variations of the fitted value of  $Y$  or we that our regression line fits the given data moderately (not very well).

Some econometricians report the  $t$ -ratios of the estimated coefficients in place of standard errors. This way of presentation makes the testing of hypothesis easier and direct.

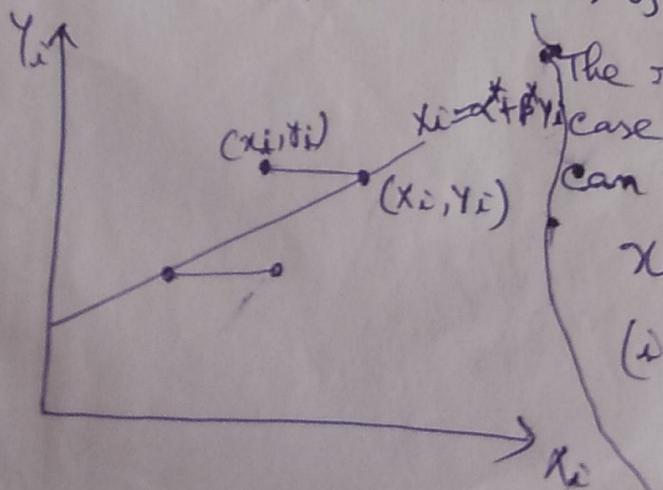
Thus the other form of presentation of results is as follows:

## Reverse Regression

In some cases where the usual regression provides biased estimates and in its place suggested an alternative procedure for obtaining ~~an~~ an unbiased estimate of the parameters on the basis of related econometric research (Vanhonacker and Day, 1987). This estimation procedure is ~~known as~~ referred to as "reverse regression".

In this method the role of endogenous and exogenous variables are reversed, i.e., exogenous variable are regressed on endogenous variables. The method of least squares estimates the parameters  $\alpha$  and  $\beta$  by minimizing the sum of squares of difference between observations and the line in the scatter diagram.

The reverse (or inverse) regression approach minimizes the sum of squares of horizontal distance between the observed data points and ~~the~~ line in the following scatter diagram to obtain the estimates of regression parameters.



The regression equation in case of reverse regression can be written as

$$x_i = \alpha^* + \beta^* y_i + \delta_i \quad (i=1, 2, \dots, n)$$

The regression equation in case of reverse

regression can be written as

$$x_i = \alpha^* + \beta^* y_i + \epsilon_i \quad (i=1, 2, \dots, n)$$

where  $\epsilon_i$ 's are the associated random error components and satisfy the assumptions as in the case of usual simple regression model.

The reverse regression estimates  $\hat{\alpha}^*$  of  $\alpha^*$  and  $\hat{\beta}^*$  of  $\beta^*$  for the model are obtained by interchanging the  $x$  and  $y$  in the direct regression estimators of  $\alpha$  and  $\beta$ , the estimates are obtained as

$$\hat{\beta}^* = \frac{\sum y_i x_i}{\sum y_i^2} = \frac{s_{xy}}{s_{yy}}$$

$$\text{and } \hat{\alpha}^* = \bar{x} - \hat{\beta}^* \bar{y}$$

The residual sum of squares in this case is

$$\begin{aligned} RSS^* &= s_{xx} - \frac{s_{xy}^2}{s_{yy}} \\ &= \sum x_i^2 - \frac{(\sum y_i x_i)^2}{\sum y_i^2} \end{aligned}$$

Coefficient of determination,  
and  $R^2 = \hat{\beta}^* \cdot \beta$

where  $\beta$  is the direct regression estimator of slope parameter,