

## Regression Analysis - 2

Step 1 → for finding whether there is relationship b/w exp & salary we use hypothesis.

$$\text{Hypo} \quad Y = \beta_0 + \beta_1 X$$

$\frac{Y}{\text{Salary}}$        $\frac{\beta_0}{\text{Intercept}}$        $\frac{\beta_1}{\text{Exp}}$

$$H_0: \bar{X} = 0 \quad H_1: \bar{X} \neq 0$$

Step 2 find linear regression model of data  
estimating regression coefficients (intercept & slopes)

Step 3 calculate TSS, ESS, RSS

Regression

Step 4 Mean Square Error (MSE)

$$= \frac{\text{ESS}}{\text{df.model}} = \frac{\text{ESS}}{K}$$

↳ also called Avg explained variance per independent feature

Step 5 Mean Square Error (MSE)

$$= \frac{\text{RSS}}{\text{df.residual}} = \frac{\text{RSS}}{n-(k+1)}$$

↳ also called avg unexplained variance per dof

Page No.	
Date	

Step 5: calculate f statistic

$$= \frac{MSR}{MSE}$$

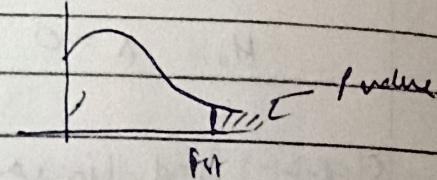
$$MSR = \frac{\sum ESS}{K} = \frac{TSS - RSS}{K}$$

$$MSE = \frac{RSS}{n-K-1} = \frac{\sum (y_i - \hat{y}_i)^2}{n-K-1}$$

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

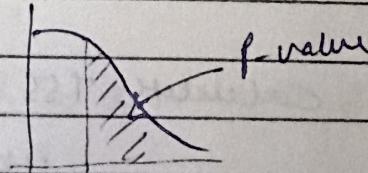
Step 6 As we got f-statistics now we will find p-value  
 $p\text{-value} < 0.05$  (reject null hypo)

$F \text{ stat} \Rightarrow \frac{\frac{ESS}{K}}{\frac{RSS}{n-K-1}}$   $\rightarrow$  suppose  $F \text{ stat} >> 1$   
means explaining SS is by



$p\text{-value}$  will be small, reject  $H_0$

$F \text{ stats} \ll 1$



$p\text{-value}$  will be large so can't reject  $H_0$

if we have more than one  $B_1, B_2$  then rejecting  $H_0$  means atleast one of  $B_1, B_2, B_3$  is non-zero

### Strength of

Now finding<sup>↑</sup> relationship bw X & Y

$$\text{↳ } R^2 = \frac{ESS}{TSS} = \frac{TSS - LSS}{TSS} = 1 - \frac{LSS}{TSS}$$

→ It → is used tells how much % or proportion of variance present in dependent variable can explained by independent var.

→ Suppose R<sub>2</sub> score of CGPA / Salary is 0.4

- If we add a relevant cov<sup>m</sup> → R<sub>2</sub> Score = 0.6

But if we add irrelevant cov<sup>m</sup> → R<sub>2</sub> Score = 0.4

But we want R<sub>2</sub> Score to decrease if add irrelevant cov<sup>m</sup>.

→ Adjusted R<sup>2</sup> Score

$$= 1 - \left[ \frac{(1-R^2) + (n-1)}{n-k-1} \right]$$

• T-Statistics -

$$\text{T-test} \Rightarrow y_{QR} = \beta_0 + m\beta_1$$

for slope → Null hypothesis :  $\beta_1 = 0$

Alternative :  $\beta_1 \neq 0$

for intercept :  $H_0 : \beta_0 = 0$

$$H_A : \beta_0 \neq 0$$

$$t\text{-Statistic} = \frac{\beta_1 - 0}{SE(\beta_1)}$$

because we assume  $\beta_1 = 0$

$$\text{for } \beta_1 \Rightarrow \frac{b_0 - 0}{SE(b_0)}$$

$$SE(b_0) = \sqrt{\frac{\sum (y_i - \bar{y}_i)^2}{(n-2) \sum (x_i - \bar{x})^2}}$$

$$SE(b_0) = \sqrt{\frac{\sum (y_i - \bar{y}_i)^2}{(n-2)} \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]}$$

• Calculating confidence interval -

$$CI_{b_0} = b_0 \pm t\text{-value} * SE(b_0)$$

$$CI_{b_1} = b_1 \pm t\text{-value} * SE(b_1)$$

• Polynomial Regression -

input	output	
x	y	(1, 2)
35	100	
35	100	
35	100	

Suppose we do degree 2

$$\text{then } x^0, x^1, x^2 | y \\ \underline{1} \quad \underline{35} \quad \underline{1225} \quad | \quad \underline{100} \quad (1, 4)$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

→ we have to choose degree in proper way.

If degree is too high overfitting  
less underfitting