

→ Regression Analysis - I

↳ It is statistical method used to examine the relationship b/w one dependent var. and one or more ind. variable.

↳ Its goal is to understand how dep. variable changes when one or more dep. variable altered

Flow of:-

- Define the research quest.
- Collect & prepare data
- Visualize the data
- Check assumptions
- Fit linear regression model
- Interpret the model (Regression Analysis)
- Validate the model
- Report results

→ Statistical inference

↳ by using sample we try to conclude about population

C9PA	99	Package 110
-	-	-
-	-	-
-	-	-

$$f(Lp9) \rightarrow Lp9 = \underbrace{f(C9PA, L2)}_{\text{calculate it}} + \underbrace{\varepsilon}_{\text{irreducible error}}$$

parametric

↳ we assume it is following linear or other regression

non-parametric

↳ we don't assume anything rather start from scratch

↳ $Lp9 = \beta_0 + \beta_1 C9PA + \beta_2 L2$

↳ while calculating we get values $b_0, b_1, b_2 \neq \beta_0, \beta_1, \beta_2$

Q. $f'(C9PA, L2)$

So error $\Rightarrow f() - f'()$

↳ ~~fixed~~ reducible error

$f()$ → True result of x & y for population

$f'()$ → estimated " " " " given sample

∴ Prediction — Suppose we have C9PA & L2 & LPA
now we want to predict other people's LPA

$$lpa = \beta_0 + \beta_1 cgra + \beta_2 l2$$
 → any person come we put its cgra & l2 and get lpa

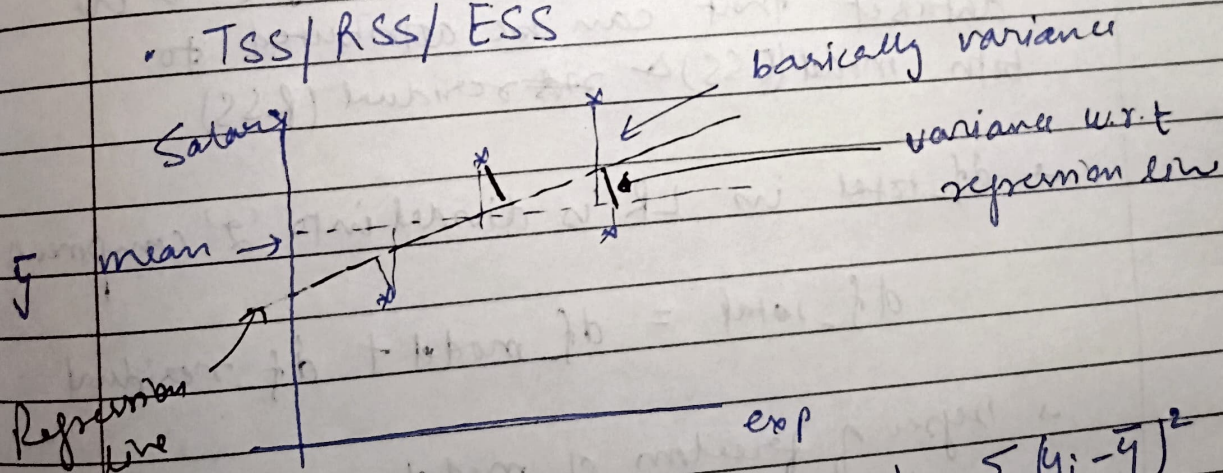
- Inference → In this we do relationship study, like select b/w lpa & cgra now it changes etc

- Linear regression is good in inference
 if inference good then predict work & vice versa

• F-Statistics & Prob(F-Statistics) → They tell us whether there is any select b/w X (input columns) & Y (output)

For this we have to do hypothesis test called F-test for overall significance

• TSS / RSS / ESS



$$TSS \pm \text{Total Sum of squares} = \sum (y_i - \bar{y})^2$$
 ↳ It tells overall variance in the data

RSS (Residual Sum of Squares)

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

↳ It tells variance present in data w.r.t regression line (means even after knowing the exp. we can't tell why variance occur)

• ESS (Explained Sum of Squares)

$$= TSS - RSS$$

$$TSS = ESS + RSS$$

• Degree of freedom -
 ↳ # rows - 1 = n - 1

$$df_{total} = n - 1$$

↳ It represent overall variability (TSS) in dataset that can be attributed to both model (ESS) & ~~residual~~ residual (RSS).

→ df_{total} in LR is divided into '2' components -

$$df_{total} = df_{model} + df_{residual}$$

→ Degree of freedom of model -
 ↳ Equal to no. of independent variables (k)

$$\underline{x_1 \ x_2} \mid y \quad \rightarrow k=2$$

→ Degree of freedom of residual -
 ↳ $n -$ no. of estimated parameters including intercept

$$= n - (k + 1)$$

⇒ F-test for overall Significance is a statistical test used to find whether there is relation by columns.

exp | Salary

⇒ whether salary is affected by experience or not