

# Regression Analysis

Regression analysis is a statistical method used to examine the relationship between one dependent variable and one more independent variables.

The goal of regression analysis is to understand how the dependent variable changes when one or more independent variables are altered, and to create a model that can predict the value of the dependent variable based on the values of the independent variables.

1. Define the research question: Identify the dependent variable (the variable you want to predict or explain) and the independent variable(s) (the variables that you think influence the dependent variable).
2. Collect and prepare data: Gather data for the dependent and independent variables. The data should be organized in a tabular format, with each row representative an observation and each column representing a variable. It's essential to clean and pre-process the data to handle missing value, outliers, and other potential issue that may affect the analysis.

3. Visualize the data: Before fitting a linear regression model, it's helpful to create scatter plots to visualize the relationship between the dependent variable and each independent variable. This can help you identify trends, outliers and any potential issue with the data.

4. Check assumption:- Linear regression model has some underlying assumption, including linearity, independence of errors, homoscedasticity (constant variance of errors), and normality of errors. You can use diagnostic plots and statistical test to check whether these assumption hold your data.

5. ~~Fit~~ Fit the Linear Regression model.

6. Interpret the model: Analyse the estimated regression coefficient their standard errors, t-value, and p-value to determine the statistical significance of the relationship between the dependent and independent variables. The R-squared and adjusted R-squared value can provide insight into the goodness-of-fit of the model and the proportion of variation in the dependent variable explained by the independent variables.

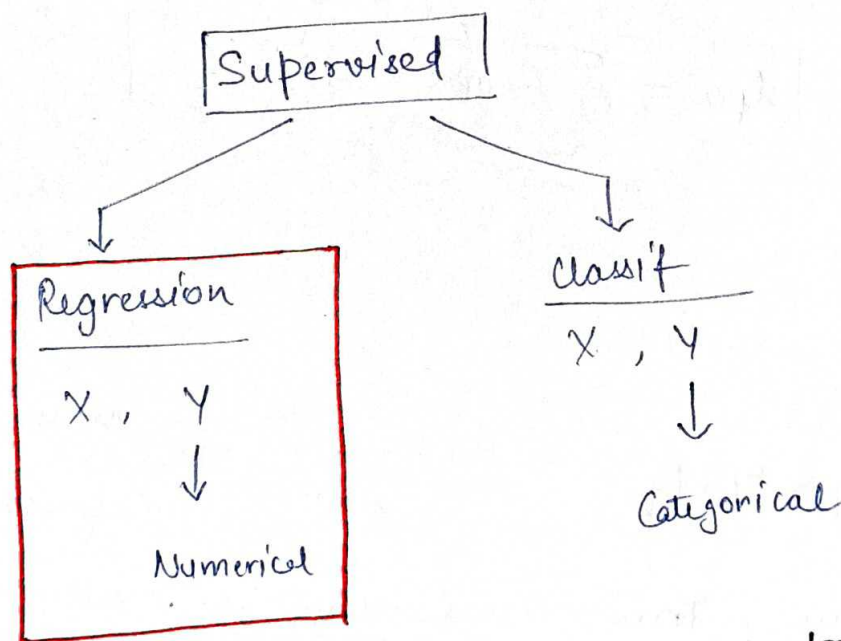
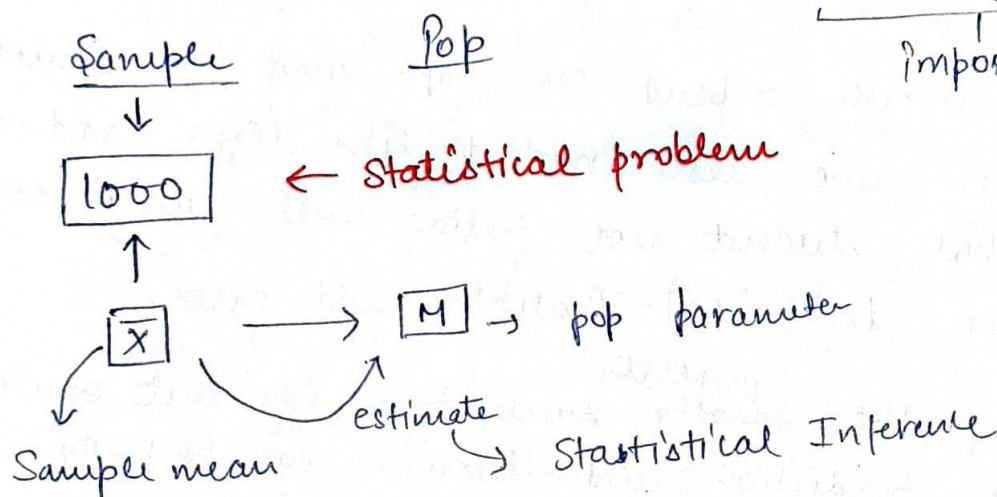


7. Validating the model:- If you have a sufficient large dataset, you can split it into a training and testing set. Fit the linear regression model to the training set, and then use the model to predict the dependent variable in the testing set. Calculate the mean squared error, root mean squared error, or another performance metric to assess the predictive accuracy of the model.

8. Report results: Summarize the findings of the linear regression analysis in a clear and concise manner, including the estimated coefficients, their interpretation, and any limitations or assumptions that may impact the results.

# Why ML problems are a Statistical Inference Problems? [with Example]

India  $\rightarrow$  avg salary  
impossible



100 data

cgpa	iq	lpa
-	-	-
-	-	-
-	-	-

machine learning problem

input

output

$$lpa = f(cgpa, iq)$$

\* Study 100 data of Students and predict for all the student and understand the mathematical relation

$$lpa = f(cgpa, iq) + \epsilon$$

error  $\rightarrow$  irreducible error

$\rightarrow$  why add error?

$lpa$  is not directly depend on  $cgpa$  and  $iq$ . Sometimes other factors are also involved like ( $cgpa$  and  $iq$  is good but student not feeling well during interview so he got low  $lpa$ ). That's why add error.

\* If we predict ~~exactly~~ <sup>perfectly</sup> ~~same~~ ~~the~~  $lpa$  but error add in prediction and become not perfectly.

$$lpa = f(cgpa, iq) + \epsilon$$

True relationship

Parameter

$$lpa = \beta_0 + \beta_1 cgpa + \beta_2 iq$$

non-parameter  
 $\downarrow$   
decision Tree

\* Let assume, True relationship is really ~~True~~ Linear and relationship is  $lpa = \beta_0 + \beta_1 cgpa + \beta_2 iq$

apply ~~relate~~ Linear Regression and get coefficient is  $b_0, b_1, b_2$ .

which is not exactly equal to population coefficient  
 $b_0, b_1, b_2 \neq \beta_0, \beta_1, \beta_2$



because if dataset change, then coefficient is also change. parameter doesn't match with current dataset's parameter.

Try to close the perfect output and we call

$$f'(cgpa, iq) \rightarrow b_0 \quad b_1 \quad b_2$$

↳  $f$  dash

$$\text{True coefficient} \rightarrow \beta_0 \quad \beta_1 \quad \beta_2$$

beez we have sample not whole dataset

$$\boxed{f() - f'()} \rightarrow \text{reducible error}$$

$$lpa = f'(cgpa, iq) + \text{reducible} + \epsilon_{\text{irreducible}}$$

estimate of  
x and y  
based on given  
data

$$f'() \approx f()$$

↑  
True of x, y for pop

Regression - Analysis of iPyub

# Inference Vs Prediction

[why regression Analysis is required]

gpa | iq | lpa      ↙ linear  
—      —      —  
—      —      —  
—      —      —  
—      —      —  
—      —      —

$$lpa = \beta_0 + \beta_1 \text{gpa} + \beta_2 \text{iq}$$

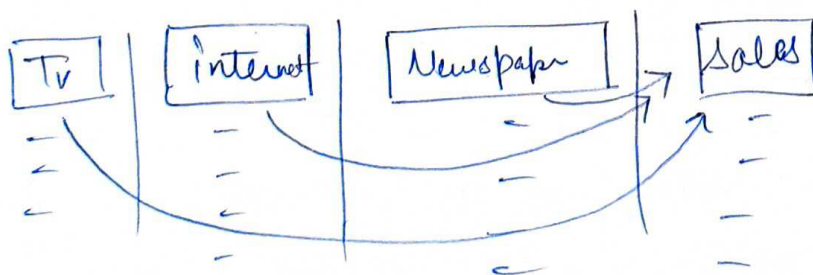
Inference (relationship study)

$lpa \rightarrow \text{gpa, iq}$   
y                      x  
relationship?

↙ Individual  
 $lpa \rightarrow \text{gpa}$   
 $lpa \rightarrow \text{iq}$   
relationship?

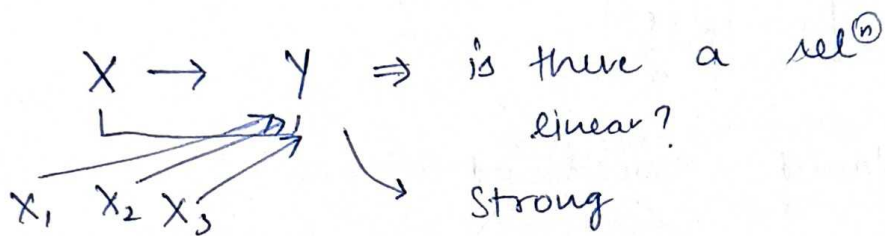
example:-

ads → sales



\* Not predicting just find relationship individually  
So, ~~that~~ all these relationship helps for  
future investment in ads.

# Stats model Linear Regression



Hypothesis test  $\rightarrow$  F-test for overall significance (ANOVA)  
this test help to find relationship bet<sup>n</sup> x and y.

## TSS, RSS and ESS

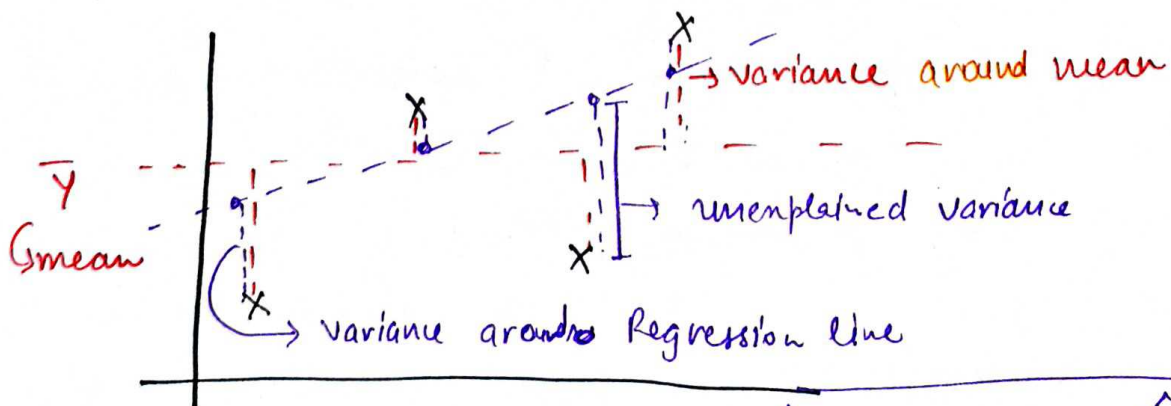
TSS  $\rightarrow$  Total sum of squared

RSS  $\rightarrow$  Residual sum of square

ESS  $\rightarrow$  Explained sum of square

example:

experience | salary  
X Y



$$TSS = \sum (Y_i - \bar{Y})^2$$

$$RSS = \sum (Y_i - \hat{Y}_i)^2$$

$\hat{Y}_i$  predict



$$TSS - RSS = \text{Explained Variance}$$

$$TSS \rightarrow \boxed{ESS} + \boxed{RSS}$$

↑                      ↑  
explained          unexplained

↓  
**Reducible**

↓  
**Irreducible**

↓  
**bad model built** or **already data have irreducible error**