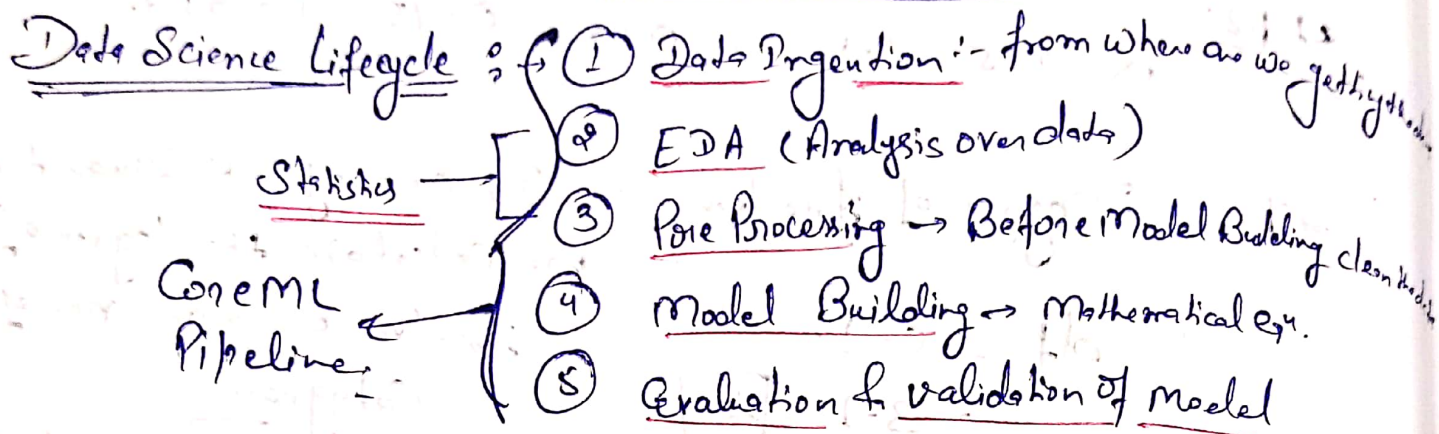
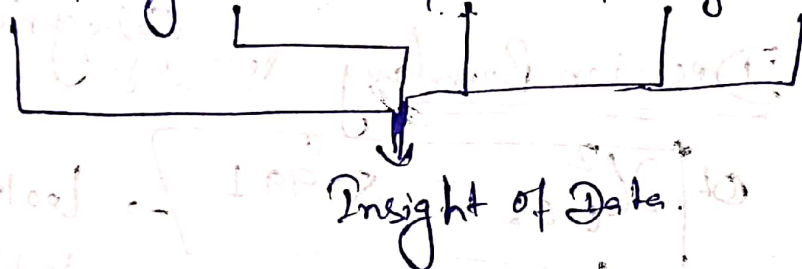


## Core M.L Pipeline



Statistics :- Collect, Organise, interpretation, analysis of data.

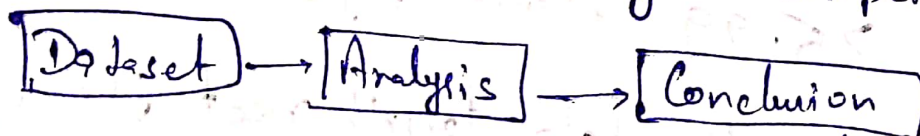


\* In every field like, Scientific, healthcare, social problems etc we use Statistics to get the insight about data.

\* In a company → Sales of product is going down. ↓ Sales

Reason:

- ① Product is not good
- ② Not Paying attention to Customer
- ③ Leadership is not good with the Project
- ④ Marketing Strategy is not good
- ⑤ Not looking to Competitors



why sale is down?

# EDA & Feature Engineering

- ① Product Manager
  - ② Business Analytics
  - ③ Data Scientist
- They will look into the Dataset and give the conclusion why sale is down and will give benefit to company.

Conclusion:- Any domain requires EDA & Feature Engg.

Data Science Lifecycle :-

① Data Ingestion :- Retrieve data from big data tools, remote location (SQL, NoSQL), some file format, CSV, TSV, XML, JSON, excel, scrape from a particular website.

Types of Data :-

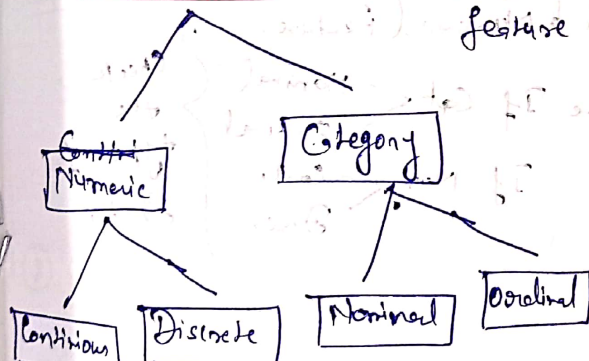
- Batch Data
  - Historical Data or (Periodic Data)
  - Mini Batch data (little more frequent)
- Streaming Data
  - Continuous / Live Data

② Data :-

- ① Structured Data :- Table → M.L (Machine Learning)
- ② Unstructured Data :- Videos, Images, Voice, Sound, text → DL (Deep Learning)
- ③ Semi-structured Data :- XML, JSON

③ We can perform EDA & FE in all types of data.

① Structure Data :-



Feature 1	Feature 2	Feature 3
Weight	Height	BMI
70	170	22
80	180	24
90	190	26
100	200	30
60	160	21

Continuous    Continuous    Cont

Continuous Data :- Continuous in Nature. Eg:- Height could be 160, 160.5, 160.88 means b/w 2 values suppose b/w 160 & 161 there could be infinite no. of values i.e. called Continuous data.

Discrete data :- Whole no. Eg:- no. of students in class i.e. 10, 20, 50, 100 etc. We can't say 10.5 students in class.



Categorical Data : → Different. different Category. Suppose

Gender < Male  
Female

Colour < Black  
white

Nominal : → Order doesn't matter. eg:- male  
female  
Here order doesn't matter whether we write  
Male first or female first.

Ordinal : → Order matters. eg:- Degree:- 10th  
12th  
Grad.  
PG  
PHD  
Order matters.

## Practical Implementation of Dataset

Feature

Student Performance

Multivariate  
Univariate  
Bivariate

Name	Age	Height	Sex	Weight	Education
Sunny	25	170	Male	70	UG
Amit	30	180	Male	80	PG
Pooja	35	160	Male	60	UG
Priya	20	150	Female	55	PHD
Aditi	27	145	Female	58	PG

while select  
give more  
weightage to  
Categorical  
UG - 0  
PG - 1  
PHD - 2  
Categorical

Categorical  
↓  
Nominal

Numerical  
↓  
Cont.

Numerical  
↓  
Cont.

Cat.  
↓  
Nominal

Numerical  
↓  
Cont.

Cat.  
↓  
ordinal

## EDA

First Level : → Identify Categorical & Numerical Feature

2nd Level : → Further Segregate If Cat < Nominal  
Ordinal  
If Num < Cont  
Disc

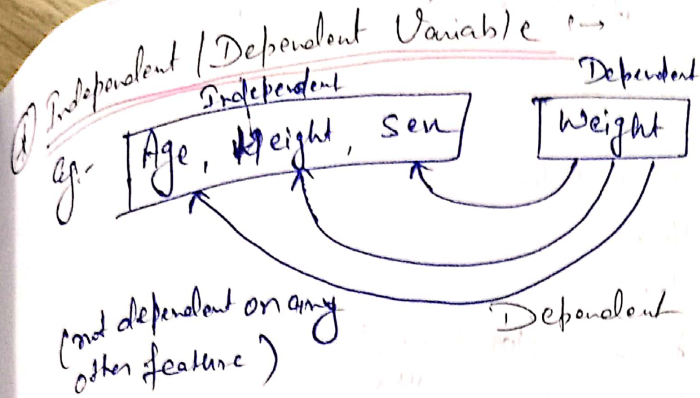
Check the type of Data

① Univariate → Single column  
Bivariate → Two column  
Multivariate → More than 2 columns

Height → Univariate Analysis

Height with Age → Bivariate Analysis

Height, Age & Sex → Multivariate Analysis



Weight is dependent in all these 3 feature. So weight is a Dependent variable

- One ML Pipeline
- ① Data Indagation
  - ② Data PEDA → Analysis of data. Based on feature
  - ③ Preprocessing → (Feature Engg.) → cleaning or wrangling or temping of data.
  - ④ Model building
  - ⑤ Evaluation or validation of Model

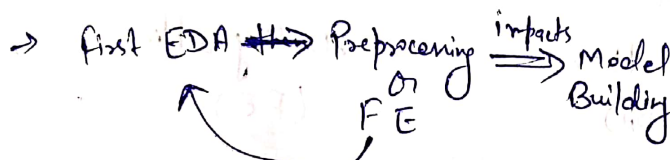
Data → Analysis

Feature / Column

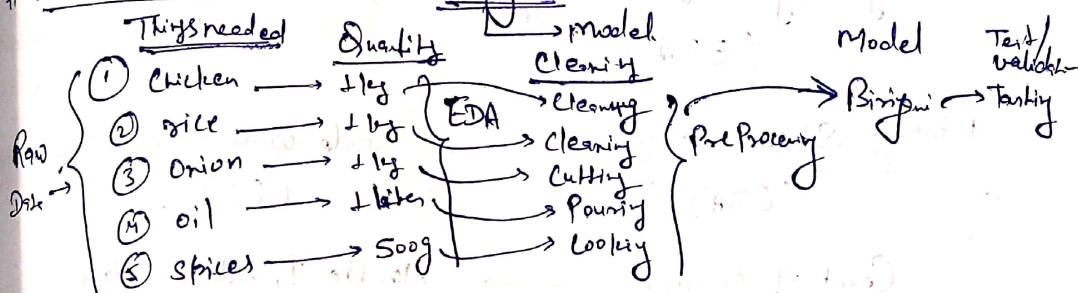
- ① Missing value
- ② Outlier
- ③ Scaling

Changes in data

Engg. on features  
↓  
Feature Engg.



# Real life example: → Cook Biryani



Are FE and EDA same? → YES

# EDA (Analysis) → Dataset: →

- Steps
- ① Profile of Data
  - ② Statistical Analysis
  - ③ Graph based Analysis

Name	Age	Education	Salary	Exp.
Sunny	25	UG	25k	2
Deepak	30	PG	30k	3
Rushi	40	UG	40k	5
Aman	50	PHD	50k	10
Shalini	20	UG	35k	1

① Profile of Data: → Person Profile

- i) no. of rows
- ii) no. of column
- iii) Missing values
- iv) Categorical variable
- v) Numerical values
- vi) Duplicates values
- vii) Datatype
- viii) Data size (RAM)

Univariate

Bivariate

Multivariate

etc.

etc.

② Statistical Analysis

i) Variation of Col.

ii) Covariance of Col.

iii) Standard deviation

iv) Correlation of Data (Col)

v) Chi square Test

vi) t-test

vii) z-test

viii) Anova test

Name

Age

Place

Designation

etc.

etc.

etc.

etc.

etc.

Statistical

Conclusion

Interpretation

of Data

Observation

etc.

etc.

etc.

etc.

③ Graph Based Analysis

i) Box Plot → outlier or dist. of data

ii) Scatter Plot → outlier, linear

iii) KDE

iv) Histogram → dist. of data

v) Pie chart

vi) Count Bar → no. of no. of Rows, Col.

vii) Heatmap → Corr. b/w Variable





Conclusion → Based on EDA we can do processing of the Data.

Steps for Preprocessing or Feature Engg. →

- Preprocessing of Data or Feature Engg.
- (i) Handle Missing Values
  - (ii) Handle outlier
  - (iii) Scaling of data
  - (iv) Transformation of Data (log, Box-Cox, Square, cube)
  - (v) Encoding of Categorical data
  - (vi) Imbalance Data
  - (vii) Feature Selection
  - (viii) Dimension Reduction (PCA, tSNE)
  - (ix) Drop Duplicates

EDA (only Analysis)  
Finding Missing Values (EDA) → Handling Missing Values (FE) (change the value of features)

Finding outlier (EDA) → Handling outlier (FE)

Finding Categorical variable (EDA)   
 Male  
 Female → Encoding (FE)

Finding Skewed Range (EDA) → Scale within a certain Range (FE)

Count of Feature (EDA) → Handle Imbalanced data (FE)

→ Feature Selection (Drop the least Correlated feature) (FE)

→ Dimension reduction (PCA, tSNE, LDA)

$\begin{matrix} X_1 & X_2 \\ \text{features} \end{matrix} \xrightarrow{\text{Combine}} \begin{matrix} X \\ \text{feature} \end{matrix}$

### Assignment 1

Perform EDA using Some Python Automated tools.

Ex:- Pandas Profiling  
. prime  
mito  
Sweetviz

} at least 3 auto or 5 Automated tool.

w.r.t only one Data

### Ass-2

EDA + FE  
Analysis

→ 10 Dataset → upload on kithub.