# EDA AND Feature engineering:-

Data Science Life Cycle:- ① Data Ingestion.
        ② EDA (Analysis)
        ③ Processing (Preprocessing)
        ④ Model Building.
        ⑤ Evaluate & Validate.

EDA:- Exploratory Data Analysis

Statistics:- - Collect the Data.
     - Organise the Data.
     - Interpretation.
     - Analysis of Data.

Kaggle is used to get different types of datasets.
  - Titanic dataset.
  - Diabetes dataset.

① Data Ingestion:-
      Get data from:-
  → Big data tools - Data can be at HDFS [Hadoop
       distributed File System]
      - No SQL Database.
      - Kafka [Streaming Data]
      - Spark Streaming

  → Remote location → SQL, NoSQL

→ Some File Format :- csv, tsv, xml, json, Excell

→ Scrap data from Website.

Types of data :- Tendency of Data ⅂
                 Batch Data, Streaming data.
                         ↓              ↓
                  Historic Data    Continuous data
                  (Periodic Data)

Mini Batch
Data [ little more freq ]

Data Wheather it is Batch Data or Streaming data
can be divided into two parts.

① Structured data. → Table.
② Unstructured data. → video, images, voice, text
③ Semi Structured data → xml, json.

Structure Data :-

| ① weight | ② Height | Feature ③ BmI |
|---|---|---|
| 70 | 170 | 22 |
| 80 | 180 | 24 |
| 90 | 190 | 26 |
| 100 | 200 | 30 |
| 60 | 160 | 21 |
| Continuous | continuous | Continuous. |

Structure data can be divided into two parts.

→ Neumeric data.
→ Category data.

```
┌──────────────┐                    ┌──────────────┐
│   Neumeric   │                    │   Category   │
└──────────────┘                    └──────────────┘
     ↙      ↘                            ↙        ↘
┌────────────┐ ┌──────────┐    ┌──────────────┐ ┌──────────────┐
│ Continuous │ │ Discrete │    │   Nominal    │ │   Ordinal.   │
└────────────┘ └──────────┘    └──────────────┘ └──────────────┘
     │
     ↓
```

Continuous [means Neumerical data that can have decimal value]

Eg:- Height [160, 160.5, 160.55]

Discrete data means no decimal value. [whole No.]
    10, 100, 200 students in st class.

Category → male     Black
           Female   white.

Nominal:→ order does not matter.
          male    } order does not matter.
          Female  }

Ordinal :→ sOrder matters.
       Example:- Degree:-

First we do 10th → 12th → Gradution → Post
           Graduation → Phd.

Dataset Student Performance :-

| Name | Age | Height | Sex | Weight | Education. |
|------|-----|--------|-----|--------|------------|
| Sunny | 25 | 170 | Male | 70 | UG |
| Arijit | 30 | 1.80 | Male | 80 | PG |
| Priyam | 35 | 160 | Male | 60 | UG |
| Priya | 20 | 150 | Female | 55 | PHD |
| Aditi | 27 | 145 | Female | 58 | PG |

Categorical   Numerical   Num   Cat   Num   Cat.

Nominal   Continuous   Continuous   Nominal   Continuous   ordinal.

Univariate :- Single column - { If we want to check height then it is univariate }

Bivariante :- Two columns { If we want to check height with respect to Age then it is bivariante

Multivariate :- More the two columns.
- If we want to check height and Age with respect to Sex then it is multivariate.

Independant /Dependant

Suppose we have
     Age, height, Sex of a person and we can define weight by knowing Age, height, Sex.

So weight → Dependant
Age, Height, Sex → Independent.

(1) Data Ingestion

(2) EDA → Analysis.

**Core ML Pipeline**

(3) Preprocessing → Feature Engineering

(4) Model Building

(5) Evaluation or Validation of model.

EDA → Preprocessing ⟹ model.

Will impact

**EDA :-** Based on the given feature, we are going to perform the analysis of the data.

**Preprocessing / Feature Engineering :-**
- Cleaning of the Data.
- Rengling of the Data.
- Preparing of the Data.

Is Preprocessing and Feature Engineering is same?
- Yes.

| Name | AGE | Education | Salary | Experience. |
|------|-----|-----------|--------|-------------|
| Sunny | 25 | UG | 25K | 2 |
| Deepak | 30 | PG | 30K | 3 |
| Rushi | 40 | UG | 40K | 5 |
| Priyam | 50 | PHD | 50K | 10 |
| Shalini | 20 | UG | 35K | 1 |

**EDA (Analysis)** →
(1) Profile of the data
(2) Statistical Analysis.
(3) Graph based analysis.

## Profile of the data:-

① No. of Rows.
② No. of Columns.
③ Missing values.
④ How many Categorical column
⑤ How many Numerical column.
⑥ Is there Duplicate value.
⑦ Dtype.

## Statistics based Analysis:- (Interpretation)

① Variance of the column.
② Covariance of the column.
③ Standard Deviation.
④ Correlation of the data b/w two column.
⑤ Perform chi square test
⑥ Perform t-test.
⑦ Perform Z-Test.
⑧ Perform Anova Test.
⑨ Mean / median / mode.

## Graph based analysis

① Box plot.              ③ Pie chart.        ⑤ KDE
② Scatter plot           ④ Histogram.        ⑥ Count bar.
⑦ Heat map.

**Box Plot :-** With the help of box plot we can find the outlier, distribution.

**Count Bar :-** Check how many rows and column is there

**Heat Map :-** We can check the correlation.

**Histogram :-** We can check the distribution.

<u>Scatter Plot</u> :- We can check the outlier of the data,
We can check data is linear or not.


By EDA we can [Preprocessing]
- Handle the missing value.
- Handle the outlier.
- Scalling of data.
- transformation (log, Box cox, square, Cube)
- encoding
- We can handle imbalance data.
- Feature Selection
- We can do Dimension reduction [PCA, tSnE)


Automated tool in python For EDA.
- Pandas Profiling.
- Mito
- Knime.


Books For Feature Engineering :-
① Feature Engineering and selection : A Practical Approach for Predictive models.
② Python Feature Engineering Cookbook.
③ Feature engineering for machine learning.