

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]

- NoSQL 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
(Periodic Data)

↓
Continuous 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 :-

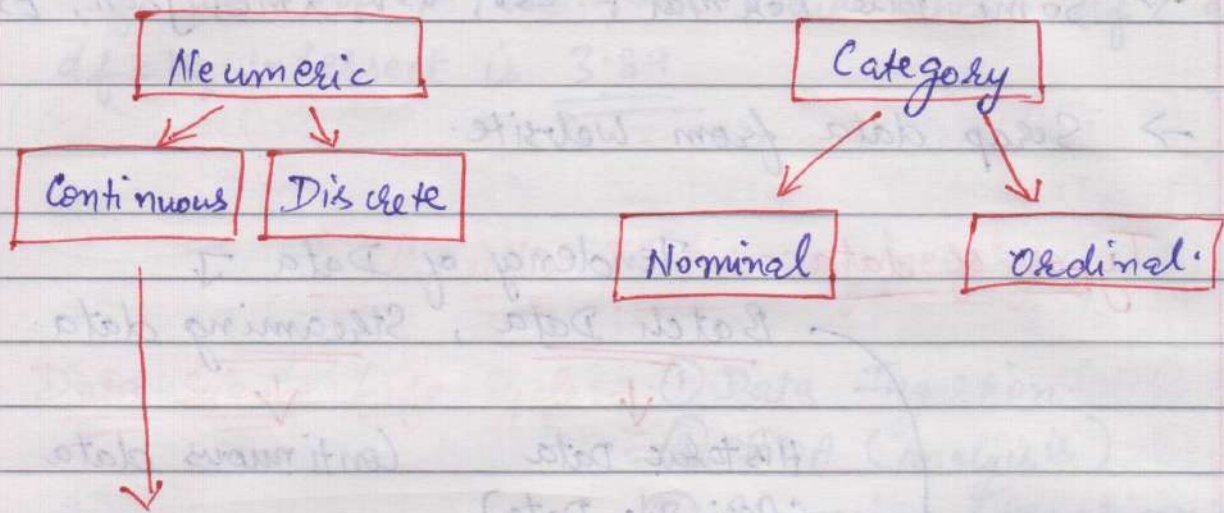
| ① | ② | Feature ③ |
|--------|--------|-----------|
| Weight | Height | BMI |
| 70 | 170 | 22 |
| 80 | 180 | 24 |
| 90 | 190 | 26 |
| 100 | 200 | 30 |
| 60 | 160 | 21 |

Continuous continuous Continuous.

Structure data can be divided into two parts.

→ Numeric data.

→ Category data.



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 8th class.

Category → male Black
Female white.

Nominal :-> order does not matter.
male } order does not matter.
Female }

Ordinal :-> order matters.
Example:- Degree:-

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

Dataset Student Performance:-

| Name | Age | Height | Sex | Weight | Education. |
|--------|-----|--------|--------|--------|------------|
| Sunny | 25 | 170 | Male | 70 | UG |
| Akshit | 30 | 180 | Male | 80 | PG |
| Priyam | 35 | 160 | Male | 60 | UG |
| Priya | 20 | 150 | Female | 55 | PHD |
| Aaditi | 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 }

Bivariate:- Two columns { If we want to check height with respect to Age then it is bivariate }

Multivariate:- More than two columns.

- If we want to check height and Age with respect to Sex then it is multivariate.

Independent / Dependent

Suppose we have

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

So weight → Dependent

Age, Height, Sex → Independent.

Core ML Pipeline

- ① Data Ingestion
- ② EDA → Analysis
- 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.
- Renaming 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 D type.

Statistics based Analysis:- (Interpretation)

- ① Variance of the column.
- ② Co Variance of the column.
- ③ Standard Deviation.
- ④ Correlation of the data set in 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.
- ② Scatter Plot.
- ③ Pie chart.
- ④ Histogram.
- ⑤ KDE.
- 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.

Scatter Plot:- 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
- Scaling 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.