# EDA (Exploratory Data Analysis) & Feature Engineering

EDA - EDA is used by data scientists to analyze and investigate data sets and summarize their main characteristic, often employing data visualization methods.

Data Science Life cycle -
1) Data ingetion ———— Project
2) EDA (analysis)
3) Processing (Pre)
4) Model
5) Evaluate and validate

(big data Tools), remote location (sql, nosql), Some file format CSV, tsv, Xml, json, website
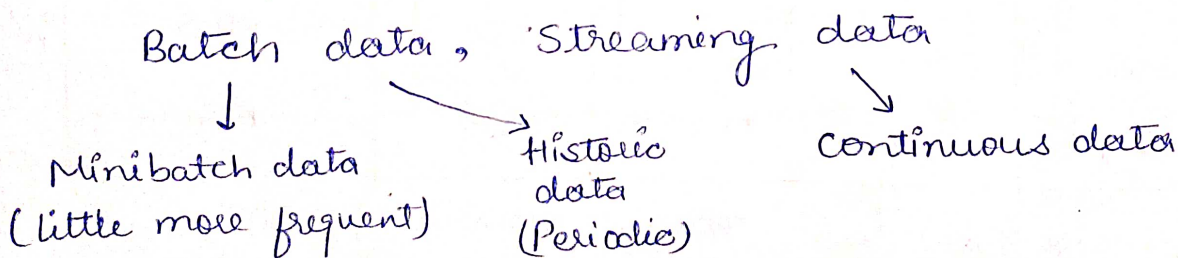
→ HDFS, NOSQL, Kafka, spark

Statistics -

Collect, Organise, Interpretation, Analysis

Insight

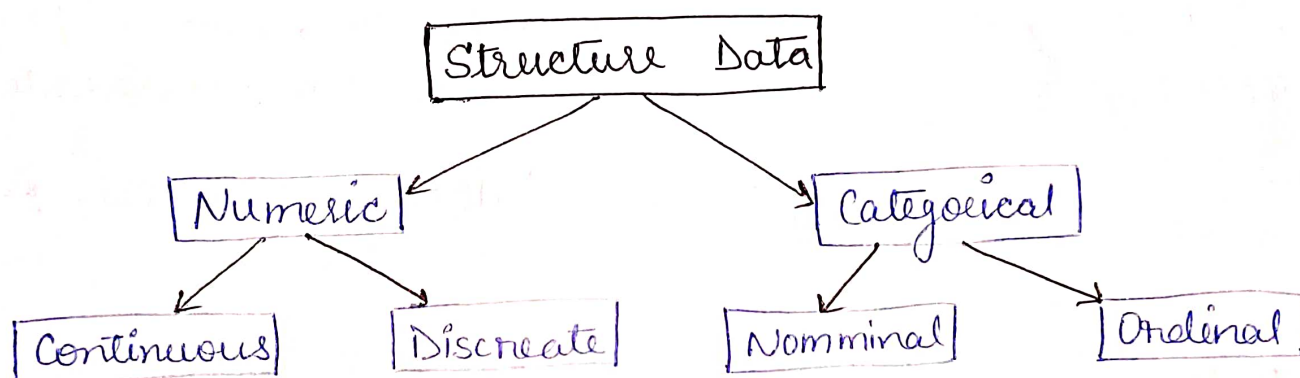(Scientific, healthcare, Social problem)

Data Types -

Batch data, Streaming data

Minibatch data
(little more frequent)

Historic data
(Periodic)

Continuous data

1.) Structure data        — Table
2.) Unstructure data      — Video, Images, text
3.) Semi structure data   — XML, Json

EDA + FE →

Structure Data -

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

Structure Data
↙ ↘
Numeric ← → Categorical
↙ ↘ ↙ ↘
Continuous Discreate Nomminal Ordinal

STUDENT PERFORMENCE —

Multivariate
Bivariate
Univariate

| FEATURE Name | Age | Height | Sex | Weight | Education |
|--------------|-----|--------|-----|--------|-----------|
| Sunny | 25 | 170 | Male | 70 | UG |
| Arijit | 30 | 180 | Male | 80 | PG |
| Priyam | 35 | 160 | Male | 60 | UG |
| Priya | 20 | 150 | Female | 55 | Phd |
| Aditi | 27 | 145 | Female | ·58 | PG |
| Categorical ↓ Nominal | Numerical ↓ Continuous | Numerical ↓ Continuous | Categorical ↓ Nominal | Numerical ↓ Continuous | Categorical ↓ Ordinal |

EDA - TYPE OF DATA

Univariate — Single coloumn
Bivariate — Two Coloumn
Multivariate — More than two coloumn

Independent / Dependent —

Age, height, Sex — [Dependent [Weight]]

Age, height — Independent

Data - Analysis

Core ML Pipeline
1) Data ingestion
2) EDA
3) Preprocessing → FE
4) Model building
5) Evaluation & validation

1) Missing value
2) Oulliers        } Feature/
3) Scaling            Coloumn

First EDA is required on FE or PP? — EDA → PP/FE

| 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 |

Steps -
EDA (Analysis)
→ Profile of the data
→ Statistical analysis
→ Graph based analysis

# Profile of the data

| Profile of the data | Graph based analysis | Stats based Interpretation |
|---|---|---|
| Row | Box plot | Varience |
| Coloumn | Scatter plot | Covarience |
| No of missing value. | Pie plot | Standard deviation |
| Categorical | Histogram | correlation |
| Numerical | KDE [Kernal density] | Chi-squre |
| duplicate | Bar Chart | T-test |
| Data type | Heatmap | Z-test |
| RAM | | Anova test |
| Data Size | | mean/median/mode |

* Based on EDA can we do processing of the data?

    Yes

## STEPS FOR FEATURE ENGINEERING

1.) Missing value handle

2.) Outliers handle

3.) Scaling of data

4.) Transformation (log, Boxcox, Squre, Cube)

5.) encoding

6) Imbalance data

7.) Feature selection

8.) Dimention reduction (PCA, ESnE).

9.) Duplicate value/coloumn

10.) Split/merge/drop/add

# WAY OF PERFORMING FEATURE ENGINEERING —

**1.) Missing Value handle**

1) Random no filling
2) Forward filling / backward filling
3) Statistical approach — mean, median, mode
4) end of the distribution
5) drop the row
6) Knn — imputer
7) Can we take that ML algorithm which missing value.
8) Can create own ML model and predict missing value.

**2.) Outlier**

| detect | handling |
|---|---|
| Z- Score | drop |
| IQR | median |
| box - plot | replace / trimming |
| Scatter plot | |
| Violin plot | |

**3.) Transformation / Sc**

box - cox
power transformation
log
square
Cube
Yeo Johnson

**4.) Scaling**

Standrization
Min / Max
Unit scaling

**5.) Encoding**

One hot
label encoding
binary encoding
target guided encoding
hash encoding

**6.) Imbalanced**

collect more data
undersampling
oversampling
Cluster - based oversampling