# EDA & FE

Exploratory Data Analysis, Feature Engineering

Data Science Life cycle:-   1. Data ingestion [Project] or big data tools
2. EDA [Analysis].
3. Processing [Pre Processing] ← (remote location,
                                 Some file format CSV, TSV,
                                 xml, excel website
4. Model
5. Evaluate (Validate)    (math)

HDFS, Nosaw, Kafka,
Spark Streaming.

Statistics:- Collect, Organize, interpretation, analysis
                 |         |           |              |
                        insight

scientific, Ascientific, social problem, healthcare

Sales of Product :- Sales is going down
                      Pint

                    Neathace
                     y
                    Ads

Product, Paying to customer, Leadership, marketing, completed.

dataset → analysis → conclusion

Project Manager
Business analy.  } domain expert
date scientist

Types of data :- Batch data. Streaming data } data

min batch data (little more frea)
histric data (Periodic)
continuous data ↓ live

1. Structured data → tables (R×C) ⇒ ML
2. Unstructured data → Video's, images, voice, sound, text, ⇒ DL
3. Semi structure data → xml, Json

1. Structured data :-

based on some charactie-ristics

Numeric

Category

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

↓ continuous
↓ continuous
↓ continuous
Segregate the feature

Discrete
↓ whole number

continuous
↓ decimal, whole numbers eg 10.5, 7.3

nominal

ordinal

Continuous :- Continuous in nature

eg : Height   160   160.5,   160.55

[160  161]

nominal :- Oder doesn't matter

eg   male,  female

Ordinal :- eg : Qualitication education in order

10th
12th
degree
PG
Phd

order matters here

Practical implementation :-    Student Performance

Feature

| Name | Age | Height | Sex | Weight | education |
|------|-----|--------|-----|--------|-----------|
| Sunny | 25 | 170 | Male | 70 | UG |
| Anujit | 30 | 180 | Male | 80 | PG |
| Prijam | 35 | 160 | Male | 60 | UG |
| Pujja | 20 | 150 | Female | 55 | Ph.d |
| Aditi | 27 | 145 | Female | 58 | PG |

↓                    ↓            ↓          ↓              ↓              ↓
                 Numerical    numeric   Categorical   numerical    Categorical
Categorical         ↓           ↓           ↓            ↓             ↓
    ↓           continuous   continuous  Nominal     continuy     Ordinal → [here order
nominal                                                                          matter]
(you other                                                                        ↓
  matter)        EDA  ⟶  types of data ?                            Ex    UG — 0
                                                                          PG — 1
                                                                          Phd — 2

Univariate  →  Single Column

bivariate  →  two Column

Multivariate  ⟶  more than 2 column

Eg  if  we want check only Height → Univarite
    if   "    "    "    Height, with age → bivariate
    if   "    "    "    Height, Weight & age → multivariate

EDA ⟹ Analysis of data

FE / preprocessing ⟹ cleaning the data

| NAME | AGE | Education | Salary | experience |
|------|-----|-----------|--------|------------|
| 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 |

EDA → Analysis of data

. Create the Profile of data

. Statical analysis

. Graph based analysis

. Create the Profile of data

→ Rows
→ columns
→ missing
→ Categorical
→ numerical
. duplicate
. D type
. Row

Statis band (interpretation)

- Variance
- Covariance
- standard deviation
- correlation
- chi square test
- t- test
- z- test
- Anova test
- mean / median / mode

} → uni, bi, multivariate

## Graph based analysis :-

- Box plot → Plotting → Outlier, distribution,
- Scatter plot → Outlier, linear
- Pie chart
- Histogram → distribution
- Kde plot → kernal Density Estimater
- Heat map → (corr)
- count bar B,C

} → dashboarding
data analysis

Observation / conclusion : Plott. → Univariate, bivariane, multivariate

Based on EDA, can we do a processing of the data ?

Yes :-

→ Can we handle missing value handlee

Pre Processing of data.
- Outlier handle
- scally of data
- Transformation of data (log, Box. cox, Square, cube)
- encoding
- Imbalance data
- feature selection

→ feature Engineering

- feature
- Dimension reduction (PCA, tsne)

Missing null value → missing value handle

EDA                          P.D

Outlier          → handle

Cat, (man, woman) → encoding + conversion of categorial into numerical

Skewed range     → scale (within a certain range)

count of feature → handle imbalance

feature selection

Dimention reduction (PCA, tsne)

```
106.                ex  1000 → subset
1. 10                      1
2. 40                     (200)
3. 20
4. 10
5. 20
```

$x_1, x_2$     X

EDA   1. Profile
          . States
          . graph

Preprocessing
. Missing / null value
. Outlier
. Scale value
. Transform
. encode
. Imbalance
. drop / duplication
. Feature selection
. dimentionality reductions.

Automated tool in Python

Pandas Profiling , autoviz , auto viz
miro              , sweet ypys , sweev vit
Knime

atleast use 3 automated tools, max 5. with respect to 1 dataset.

EDA + FE => 10 in depth analysis => your git hub repo.

Sunny.Savita@ineuron.ai
Krish.naik@ ineuron.ai

## 4. Automated EDA Libraries by Krish naik Video :-

- Dtale
- Pandas Profiling
- Sweet viz
- Auto viz
- LUX → Paid only
- DataPrep
- Pandas-visual-analysis

25/09/22 → Core ML Pipelin

- Data collection
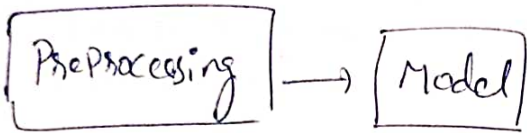- EDA & Analysis
- Preprocessing /FE
- Model building
- Evaluation matrix /validation.

EDA.

1. Profile
2. Stats based analysis
3. Graph based .

Preprocessing:

- Missing Value
- Outlier handle
- Scale
- Transformation

- encoding
- Handle imbalance
- Feature selection
- Dimention reduction (PCA, LDA, TSNE) → Principle component analysis.

- duplicate value /duplicate cols
- Split / Merge /drop /add

$$\boxed{\text{PreProcessing}} \longrightarrow \boxed{\text{Model}}$$

ways of Performing FE

1. Missing value handle :-

- Random
- forward filling / backward filling
- Statical approach (mean, median, mode)
- End of the distribution
- drop that row
- KNN - inputer
- Can we take that ML algorithm which missing value
- Create your own ML model, you can Predict the missing value

2. Outlier handle :-

- detect the Outlier

  handling

- Z - Scale
- IQR
- Box Plot
- Scatter plot
- Violin Plot

  - Drop
  - fill with median
  - Replace / with any value
    trimming

3. Transformation / scaling of data

- box- cox transformation
- Power - transformation
- log
- Square
- Cube
- Yeo Johnson

scaling

- Standardization
- Min-max scales
- Unit Scaling

Encoding:

- One hot
- Label encoding
- Binary
- Target guided encoding
- Hash encoding

Data → EDA

Dynamic
Preprocessing → Model → 75%,
/),
80%,
77%

- Missing
- Outlies
- Scale
- encoding

- Structured data

- image data / text data / unstructed data

Imbalanced

- Collect more data
- Under Sampling
- Over sampling
- Cluster based over sampling