# DAY 2 - EDA & FE

25th/9/22

Recall

## ML Pipeline

1) Data Collection
2) EDA (analysis FE)
3) Feature Engineering
4) Model building
5) Model evaluation (validation)

## EDA

1) Profile the data
2) Stat based Analysis
3) Graphical Analysis

## Pre processing (or) FE

1) Missing values - handle
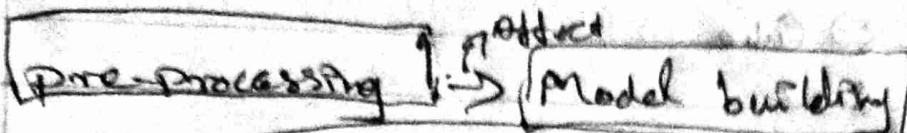2) outlier
3) Scaling data
4) transforming the data

5.) encoding
6.) Imbalance data handling
7.) Feature Selection
8.) Dimension Reduction [PCA, LDA, t-SNE]
9.) duplicate handling
10.) split / merge / drop / add.

Affect
Pre-processing | → Model building
                way of performing FE

→ Missing value - Handle

In 2 we have different types of technique
to handle the missing values
   1.) to add the Random values
   2.) forward filling / backward filling
   2) Stat (mean, model, mode)
   4.) end of the Distribution value
   5) Drop the Particular rows
   6) Knn - Imputer
   7.) ML Algorithm Which Support
       missing value handling
                    ⋮
2.) outlier handle        etc.

   1.) Detect the outlier
   2.) handle the outlier


Detect                    handle
1) z-score                1.) drop
2.) IQR                   2.) help of median
3.) Box, Scatter          3.) replace avg value for outlier
   violin plot            4.) trimming

3) Transformation

    1) box-cox 2)
    2) power Transformation
    3) log
    4) square
    5) cube
    6)
      etc

4) Scaling

    1) Standardization
    2) Min-mame Scalar
    3) unit Scaling
        etc

5) encoding

    1) one hot
    2) label encoding
    3) binary encoding
    4) Target guider encoding
    5) Hash encoding

6) Embalance Data Technique

    ★ 1) Collect more data
      2) under Sampling
      3) over Sampling
      4) Cluster based over Sampling

Data → EDA ⟶ Preprocessing → model test

1) missing ↩

2) outlier →

3) scale →

4) encoding ⟩

If we want 50% -%

then we have to

choose the preprocess

Technique