

## Core ML Pipeline

- 1 Data Colled
- 2 EDA (Analysis)
- 3 Preprocessing or FE
- 4 model building
- 5 evaluation matrix or validation

## EDA

- ① Profile
- ② Stats. based analysis
- ③ graph based analysis (Python)

## Preprocessing

- ① missing values
- ② Outlier handling
- ③ Scale
- ④ transformation
- ⑤ encoding
- ⑥ handle imbalanced
- ⑦ feature selection
- ⑧ dimension reduction (PCA, LDA, tSNE)
- ⑨ duplicate value / duplicate col
- ⑩ split | merge | drop | Add → Data

Preprocessing

model

Way of performing feature eng

## ① missing value handle

- ① Random
- ② forward filling / backward fill
- ③ statistical approach
  - mean, median, mode
- ④ end of the distribution
- ⑤ drop that row
- ⑥ knn-imputer
- ⑦ can we take that ml algo which missing value
- ⑧ own ml model you can predict.

## ② Outlier

detect

2-score

IQR

box plot

Scatter Plot

Violin Plot

handling

drop

median

replace

trimming

## ③ transformation

- box-Cox
- Power transformation
- log
- square
- Cube
- Yeo Johnson

## ④ Scaling

Standardization  
min-max  
Unit Scaling

## ⑤ encoding

- One hot
- label encoding
- (Binary encoding)
- target guided encoding
- Hash encoding

## ⑥ imbalanced

- ① Colled more.
- ② under sampling
- ③ over sampling
- ④ Cluster based over sampling

Dynamic

Data → EDA

Preprocessing

model → 75%

- ① missing →
- ② outliers →
- ③ Scale →
- ④ encoding →

77% → 80%

Structure = Image - text

- ① Profile
- ② stats
- ③ graph.

10 dataset

One dataset

read data

Dataframe

EDA

Observation

- ① univariate
- ② bivariate
- ③ multivariate

interias

= 90% → EDA Preprocessing

Single folder

github

↓

② Preprocessing

1. missing (10 - 62)
2. outlier (Step 5)
3. Scaling / transform
4. encoding
5. Feature selection
6. Imbalance data

Automate EDA

Pandas  
Sweetviz  
autoviz (Pip)