

Day 2 of EDA & F.E

Revision of Day-1:-

Core ML Pipeline

- Data ingestion
- EDA (Analysis)
- Feature engineering or Pre-Processing
- Model Building
- Evaluation matrix or validation

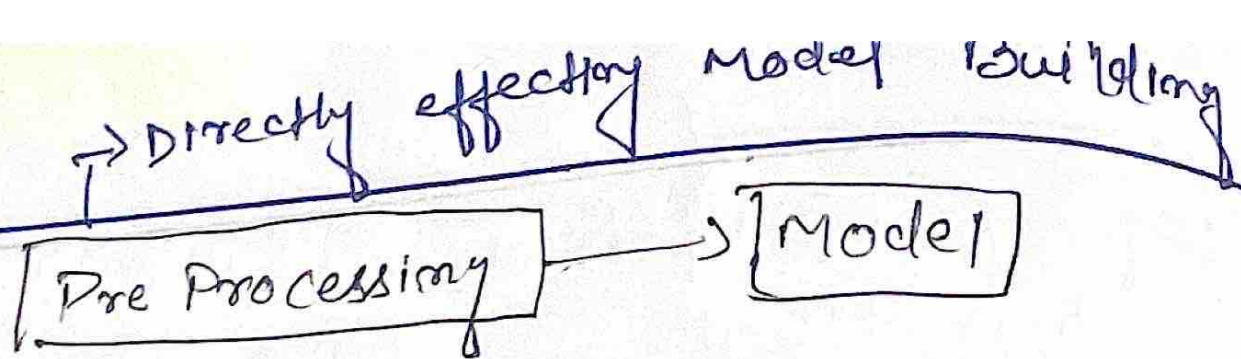
EDA

1. Profile of the DATA
2. Statistical Based Analysis
3. Graph Based Analysis (important)

PreProcessing

1. Missing value handle
2. Outliers handle
3. Scaling of Data
4. Transformation (Log, Cox-Box, Square, Cube)
5. Encoding
6. Handle Imbalanced Data
7. Feature Selection
8. Dimension Reduction
9. Duplicate value / Duplicate column
10. Split / Merge / Drop / Add (column)





## ways of Performing Feature engineering

### ① Missing value Handle

- Fill with Random no.
- Forward / Backward filling
- Statistical Approach (mean, median, mode)
- With the help of end of Distribution, fill the missing value
- Drop the Row
- Impute with KNN (KNN-IMPUTER)
- ML Algorithm for missing value
- Build own ML Model to Predict missing value.

### ② Outlier Handling

Detect the outlier

- Z-Score
- IQR
- Box
- Scatter-Plot
- Violin-Plot

Handle outlier

- Drop
- Fill with median
- Replace / Imputing



### ③ Transformation

- ① Box Cox Transformation
- ② Power Transformation
- ③ log
- ④ Square
- ⑤ cube

### 4. Scaling of Data

- 1) Standardization
- 2) Min Max Scaling
- 3) Unit Scaling

### ⑤ Encoding Various Method

- ① one-hot encoding
- ② Label - encoding
- ③ Binary encoding
- ④ Target Guided Encoding
- ⑤ Hash Encoding

### 6. Imbalanced Dataset (Treatment various method):-

- ① collect more Data
- ② under Sampling
- ③ over Sampling
- ④ Cluster Based over Sampling



Q. How to Find the Best Model Accuracy  
various Method?

⇒ To increase the Accuracy, we need to change the Pre-Processing technique & use different Method or step from the above.

⇒ we need to use each & every Pre-Processing step & find Best Accuracy.

Q. How do we transform the Data?

⇒ import numpy as np  
np.log(df)  
sns.distplot(df)

Q. Does outlier effect Mean?

⇒ Yes, outlier effect the Mean.