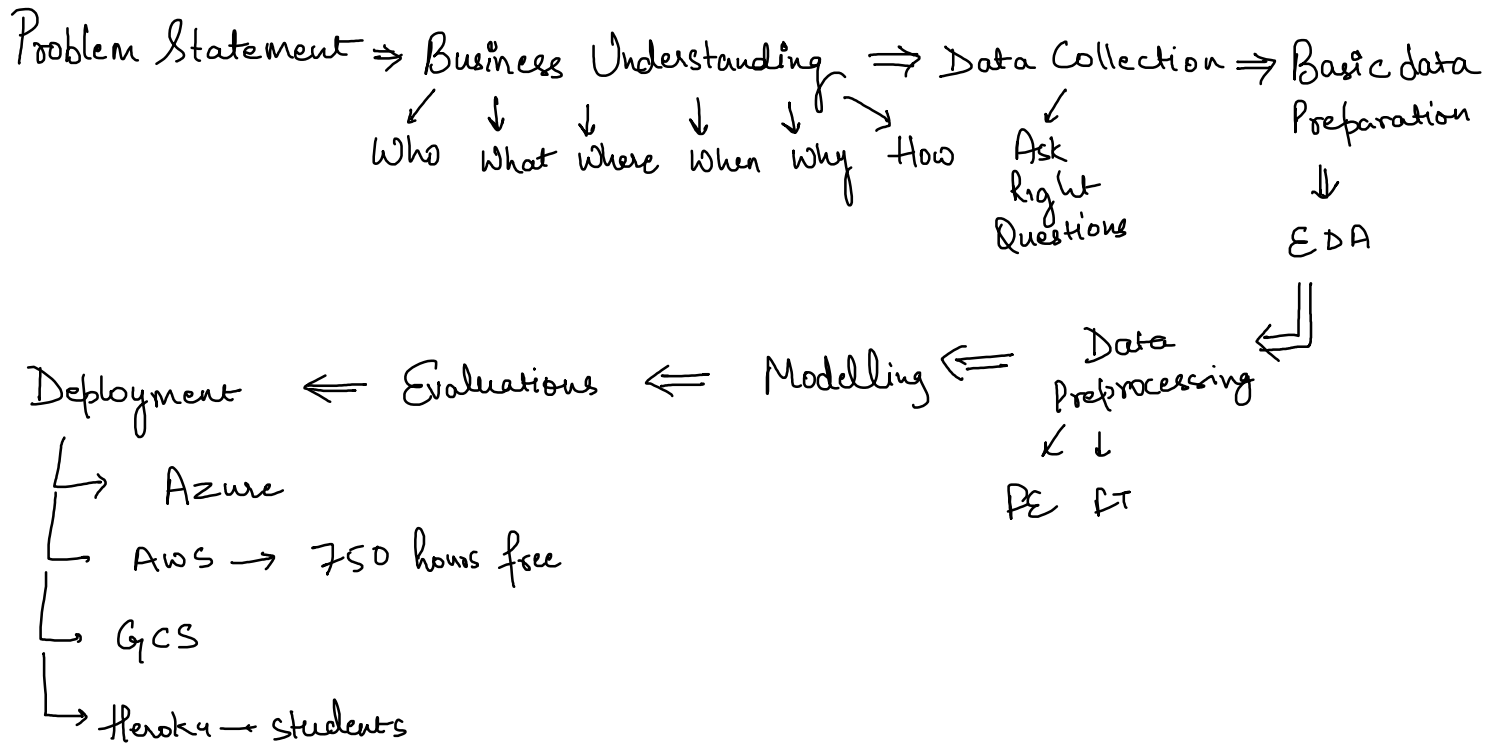


Lifecycle of Data Science Project (3 months - 1 year)

vague

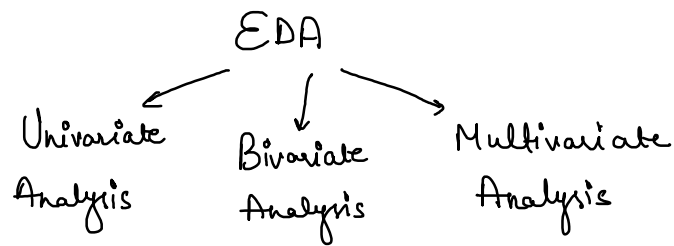


Starting Project

- \rightarrow Import libraries & load the data
- \rightarrow Basic Info about data \Rightarrow `df.info()`
 - metadata
 - Rows
 - columns
 - datatypes
 - nulls & non-nulls
- \rightarrow Basic description of the data \Rightarrow `df.describe()`
- \rightarrow Basic data study \Rightarrow `unique`, `nunique` & `value-counts()`

→ Basic Data Preparation → data quality check / data assessment

⇒ EDA

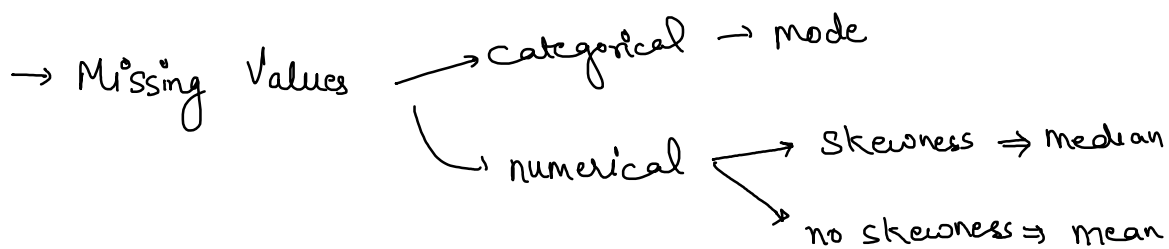
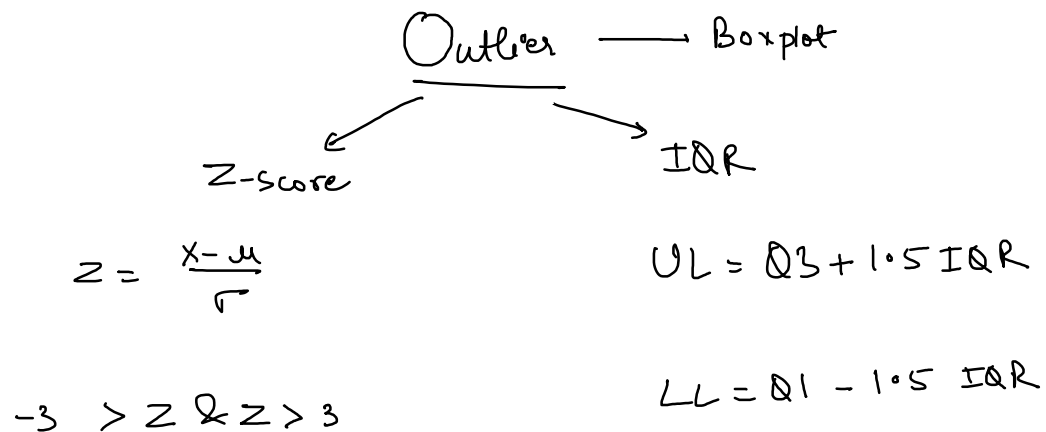


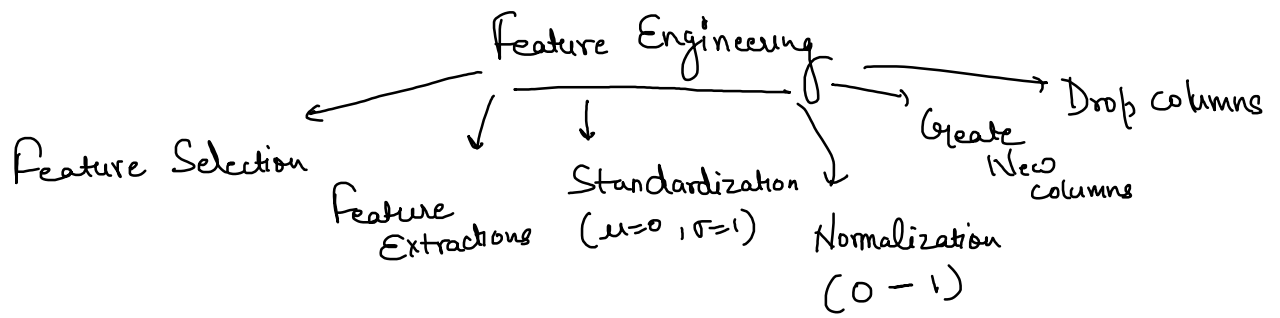
Univariate ⇒ histogram, countplot, boxplot, kdeplot

Bivariate ⇒ line, lmpplot, scatter, bar, pie

Multivariate ⇒ "hue", heatmap, pairplot

UNIVARIATE → BIVARIATE → MULTIVARIATE (Project EDA flow)





Encoding \Rightarrow OHE (Nominal), Ordinal Encoding, Label Encoding
 (One hot encoding) (g/p variables) (O/P variable)
 pd.get-dummies

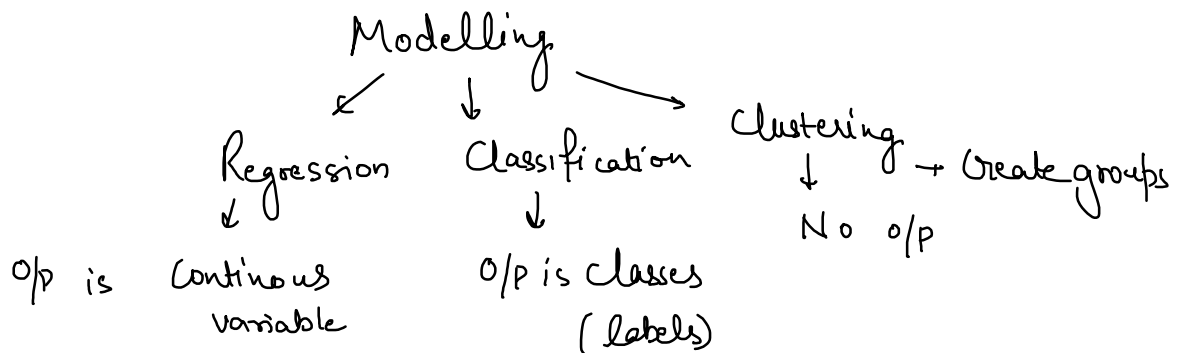
ordinal category

Grade	Encoding
A	1
B	2
C	3

One Hot Encoding

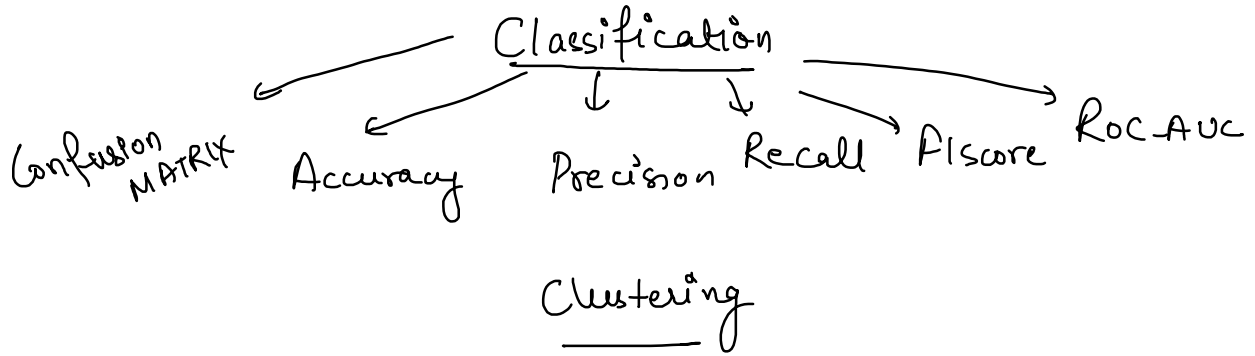
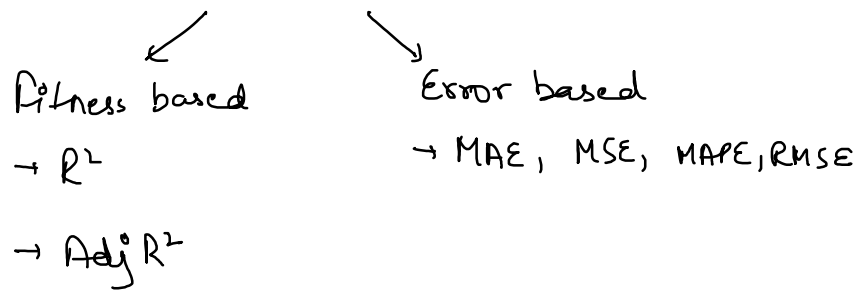
	A	B	C	D	E
A	0	0	0	0	0
B	0	1	0	0	0
C	0	0	1	0	0
D	0	0	0	1	0
E	0	0	0	0	1

drop first = True



\Rightarrow Evaluations

Regression



$$\text{Silhouette's score} = \frac{b - a}{\max(b, a)}$$

$[-1, 1]$

INDO - PAK Relations

Tanya

→ India → Hindu Majority
 Pak → Muslim Majority

→ Bigger land area of India.

Bharath

→ Health care system is better in India

→ Better Education in India

Monisha

→ Better foreign relations of India

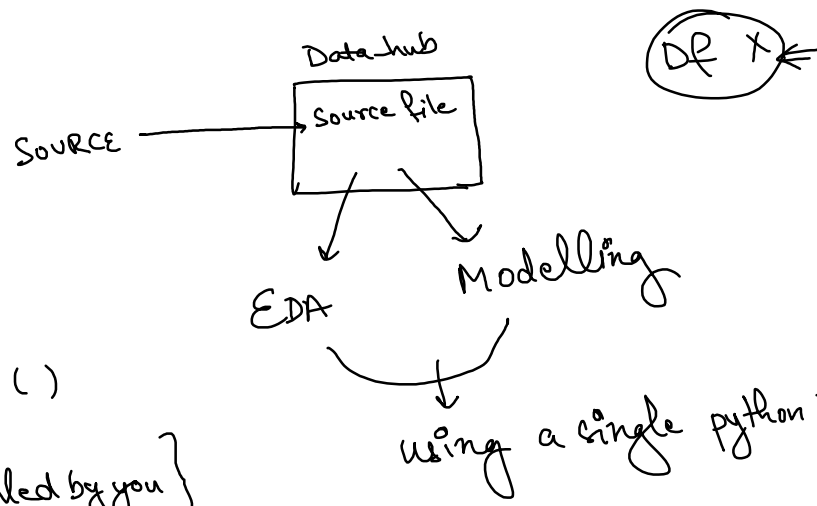
→ GDP is better (India)

PROJECT

AIR TICKET PRICE PREDICTION

Data hub
 0.1.1

DP X



```
def preprocess()  
{ to be filled by you }
```

```
return data_eda, data_model
```