## Data Preprocessing & EDA (Exploratory Data Analysis)
### (cleaning)

1) Data Preprocessing also called data cleaning.

Step 1
2) Constant Features

→ Whole column has only one value.

With one constant value only these are not needed, they are generally dropped.

Quasi-constant feature                                    → threshold defined.

→ A column which has 99% or so values are same and only one, two outliers.

i.e almost a constant feature. These are also generally dropped.

Step 2

3) Data filtering

i) → Drop columns which have very high issue of null..

df. isna(). mean() % 100 → (gives percentage of null in columns).

ii) Dropping duplicate values.

iii) dropna fn used to drop NoN values. This is used to drop rows.

iv) Dropping columns that have data, but it is not relevant to you.

v) If null values not dropped they are generally replaced by a string (or) number. Bcz null is not understood by data model later. [fillna used].

4) Date Time Conversion

→ HRC mainly handles time-series data
So, date-time conversion is very imp.

⇒ format = '%Y-%m-%dT%H:%M%Z' format

Syntax⇒ df 1["columnName"].dt.year
         └→ has data.    ↓date    ⇒ week, day, dayofweek etc
                                      also available.

5) Data Splitting → Missed bcz of some work. but basic.
                    is, → Train set ⟵ Data Set → Test set data
                         data            (80:20 or 70:30).

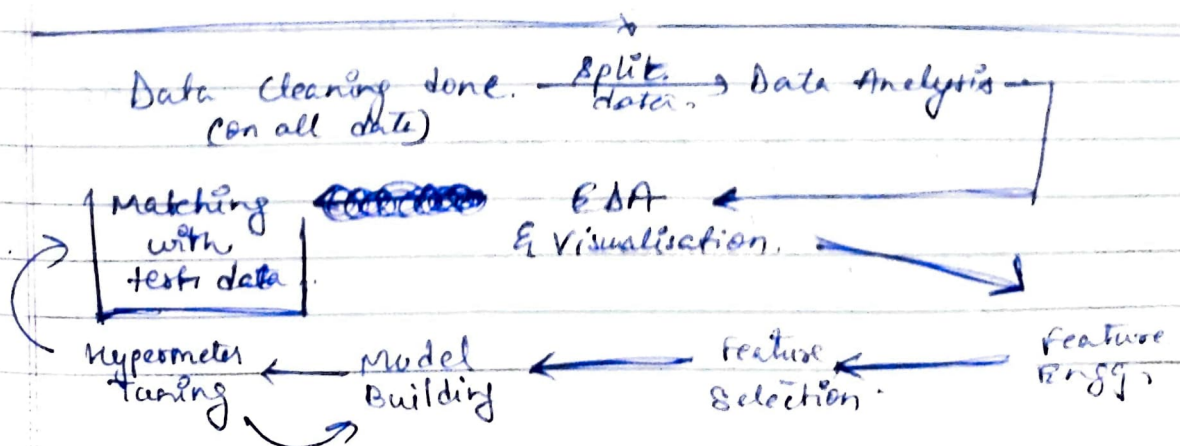Statistics

6) ⇒ Continuous Variables → A seq of possible values which are infinite
                                                         and uncountable.
   Categorical Variables →

7) ⇒ Distribution → Uniform, exponential, normal}→

                    Binomial → based on a histogram
                    └→ for a categorical variable

Data Cleaning done. → Split data → Data Analysis ⌐
(on all data)                                      |
┌Matching      ~~~~~~~~  EDA ⟵                     |
| with         & Visualisation.  ⟶                |
| test data |                        ⟶           ┘
Hyperparameter ⟵ Model ⟵ Feature ⟵ Feature
tuning          Building   Selection       Engg.

Data Analysis Maths overview
**DA** [we will only cover numerical analysis].
→

1.) **Univariant & Multivariant Analysis.**

→ Mean, Median or <u>mode</u>.
      ↳not generally used.    Univariant.
→ Variance, Standard deviation.    (with only one
             column).
→ Skewness, and Kurtosis.

<u>Studying multiple variables at once → multivariant.</u>

→ Correlation.
→ Covariance.
→ Principal Component Analysis (PCA).
        ↳ Not used in Fintech.

  ★ In fintech generally data is not modified.
    Just analysis is done.

2.) <u>Distribution & IQR (Inter Quartile Range).</u>

      → Population, Sample, Distributions.
      → Probability. concepts discussed
         here. ✓

 ★ <u>Skewness of distribution</u> → Left, Right discussed.
       ↳ measured by mean, median & mode.
         (Measures of central tendency).

Median → Positional Avg
Mean → Mathematical Avg
(AM, GM, HM)

3) <u>Outlier Detection</u> → Disrupt data & conclusions, generally taken out from data.

Detection generally done by boxplot.

IQR → Inter Quartile Range. (outside the IQR → Outlier).

———————————————×———————————————

**2.** <u>EDA (Expolatory Data Analysis)</u>.

↱ Private    ↱ client

① ⇒ Converting numbers to visual [Matplotlib, Seaborn &

② ≈ <u>Barplot</u>, <u>Scatterplot</u>, <u>Boxplot</u> } diff kinds of plots    Plotly ]
                                                        ↳ very high level.

③ ⇒ Import matplotlib as plt.
                    ↳ Like pandas is pd.
                       & numpy is np.

④ ⇒ Other plot → histogram, piechart, etc.
                    Distplot, Violin plot, Colon density plot.

———————————————×———————————————

<u>Final Project</u>

Clear date → when bill will be cleared. of each of
column in            50,000 invoices.
dataset.

We will predict clearing date.

PDP.

———————————————×———————————————