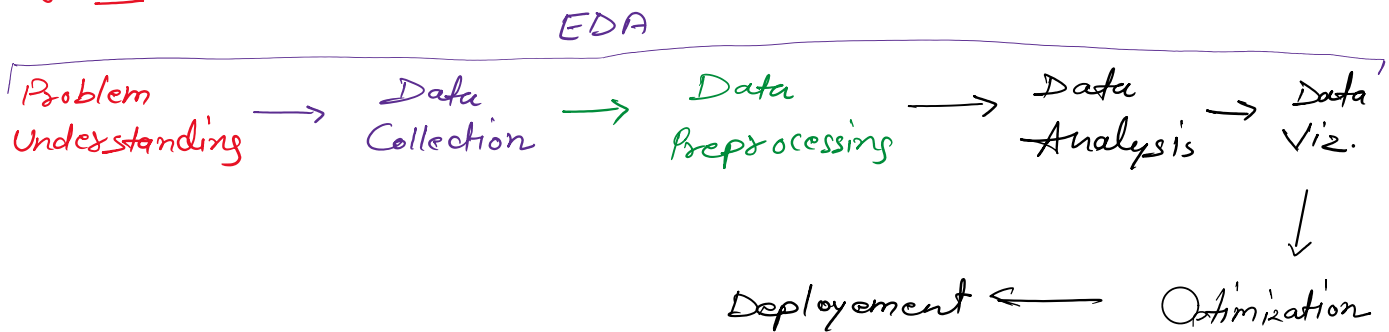


Lifecycle



EDA : Exploratory Data Analysis

↳ Understanding Data

① Univariate Analysis

↓ ↓
Single Column

Name	Age	Gender

② Bivariate

↓ ↓
Two Column

Cat + Cat → count

* Cat + Num → stats by category

* Num + Num → Relationship
[correlation]

③ Multivariate

↓ ↓
> 2 Column

Correlation
—
bar, col

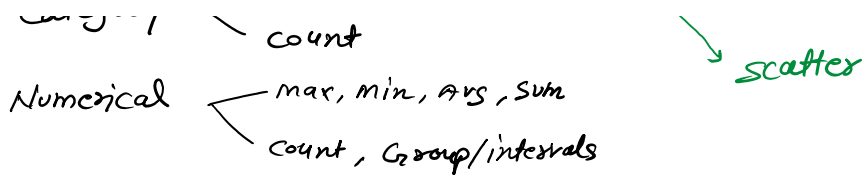
Analysis

Category — Unique category count

Numerical — max, min, avg, sum

column, bar, pie

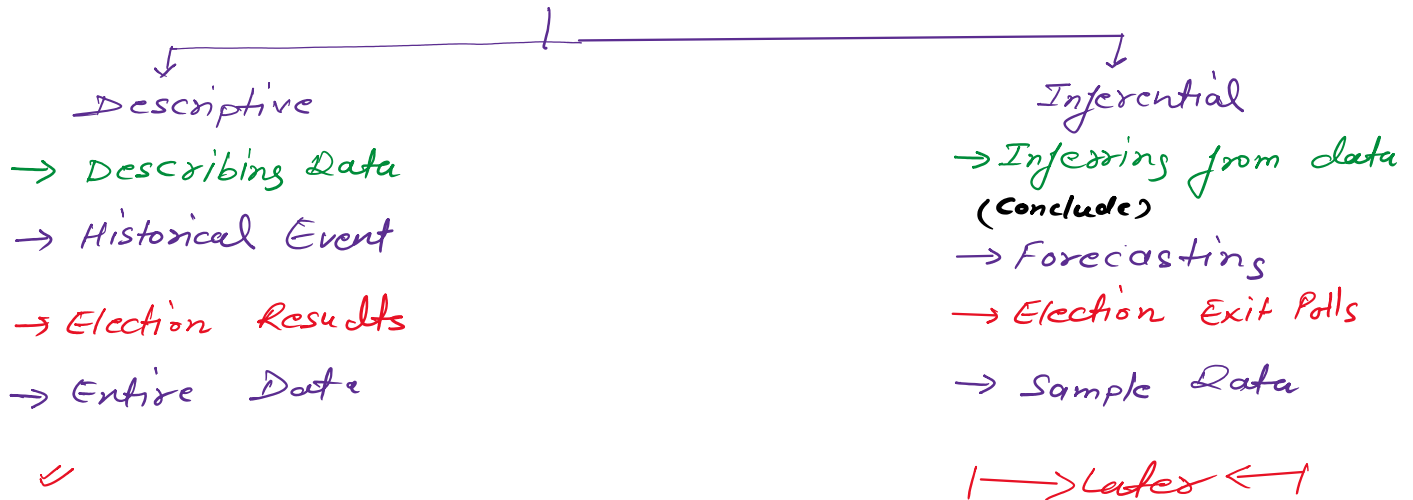
scatter



countplot(), histogram, boxplot

Statistics

Collect, Organize, Analyze, Interpret Data.

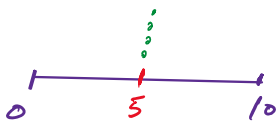


① Aggregate Measure

sum(), max(), min(), count(), average()

② Central Tendency Measure : exactly or approx into two parts.

mean	median	mode
$= \frac{\text{Sum}(\text{value})}{\text{Count}(\text{value})}$	$= \begin{array}{l} \text{Sort data} \\ \text{then} \\ \text{middle value} \end{array}$	$= \text{most repeated value}$



19, 20, 20, 21, 22, 20, 19, 35, 37

19, 19, 20, 20, 20, 21, 22, 35, 37

✗ Mean = 23.6

Median = 20 ✓ [Data contains outlier]

Mode = 20

⇒ To select right CTM, we need to look at data distribution.

③ Dispersion / Data distribution measure

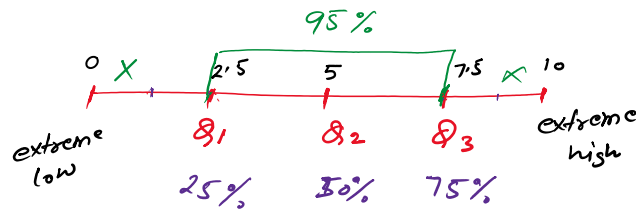
① Range

max, min



$$\text{Range} = \text{max} - \text{min}$$

② Interquartile Range [IQR]



$$\text{IQR} = Q_3 - Q_1$$

identify outliers by defining lower & upper

$$\text{Lower} = Q_1 - (1.5 * \text{IQR})$$

$$\text{Upper} = Q_3 + (1.5 * \text{IQR})$$

③ Variance / Standard Deviation

↳ distance of any value from its average.

$$\text{var} = (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2$$

$$= ()^2 \rightarrow \text{square unit}$$

$$\text{std} = \sqrt{\text{var}}$$

$$= () \rightarrow \text{same unit}$$

Higher Standard Deviation : Data more scattered