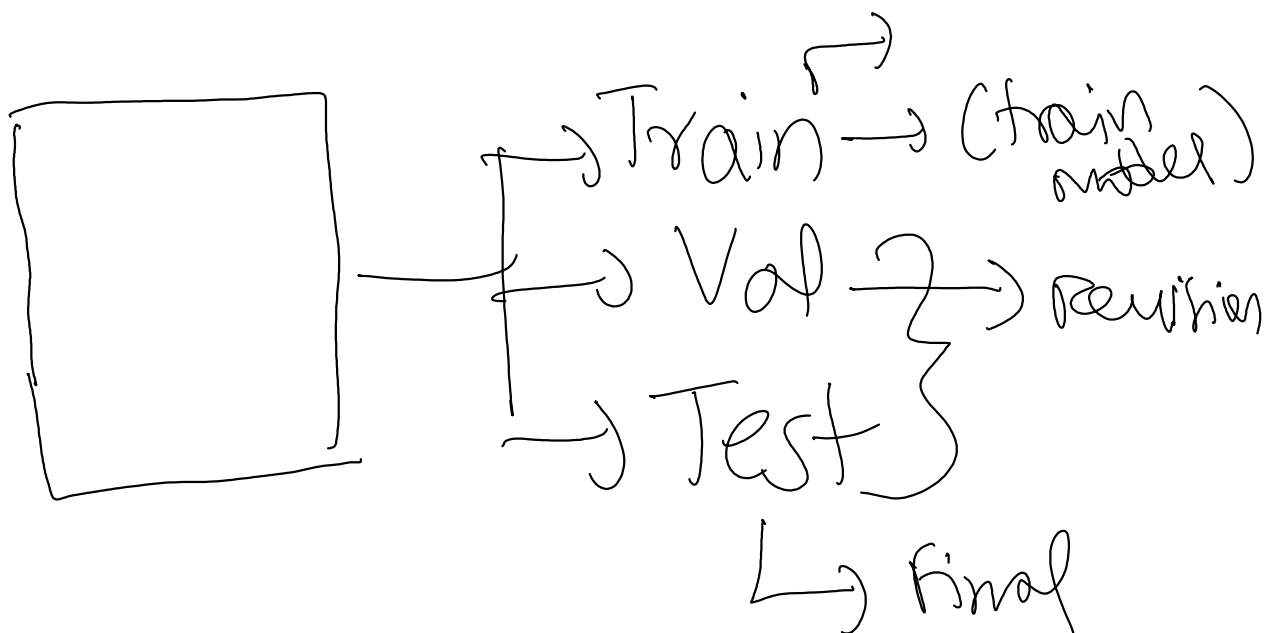# Session 3 - EDA

## EDA - Exploratory Data Analysis:

- It's the practice of exploring your dataset, by utilizing visualization tools and statistical measures to understand and extract underlying patterns and information within the data.
- This gives a clearer picture of the data, helps us make data-informed decisions and solve crucial business problems with much fewer assumptions and more facts.
- This step is the backbone of any Data Science project and takes up a major chunk of the project timeline.

## Additional Note:

- The term EDA was coined by the late mathematician, John Tukey
- First introduced in his book "*Exploratory Data Analysis*" (1977)
- Mr. Tukey also introduced the "Box Plot"

Val $\longrightarrow$ Tune
      $\hookrightarrow$ Best parameters

Retrain
      $\hookrightarrow$ Best param
            $\hookrightarrow$ Train + Val Set

# Components

## **What is done in EDA stage?**

## Look at and analyze numbers and plots
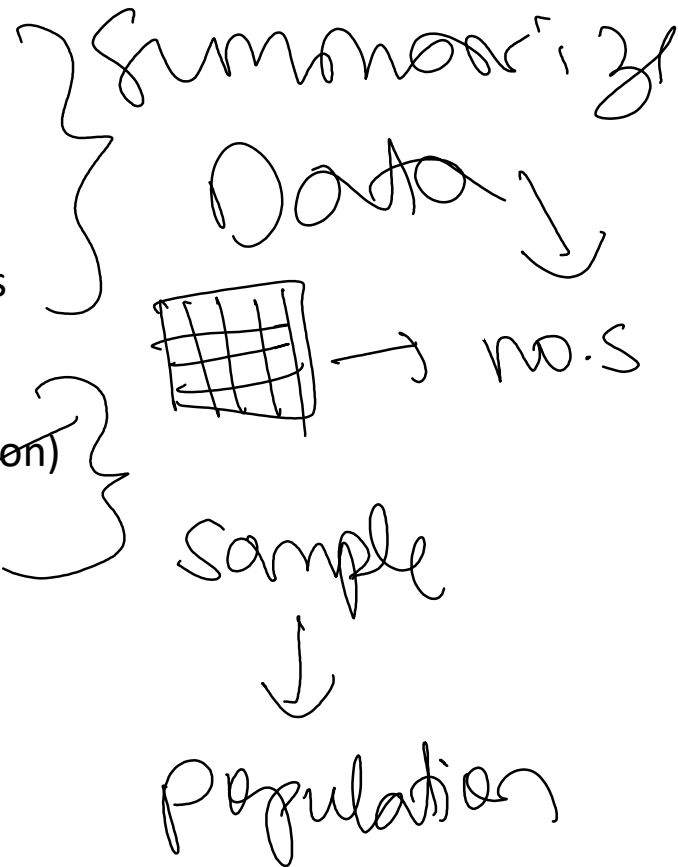
- **Descriptive Statistics**
    - Central Tendency
    - Dispersion / Spread
    - Distribution of Variables
    - Symmetry and Shape of variables

- **Inferential Statistics**:
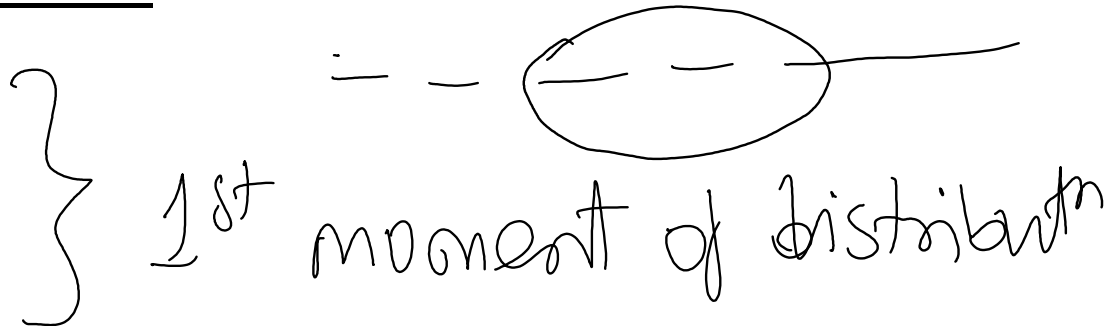    - Strength of Association (correlation)
    - Hypothesis Testing

- **Plots / Graphs**:
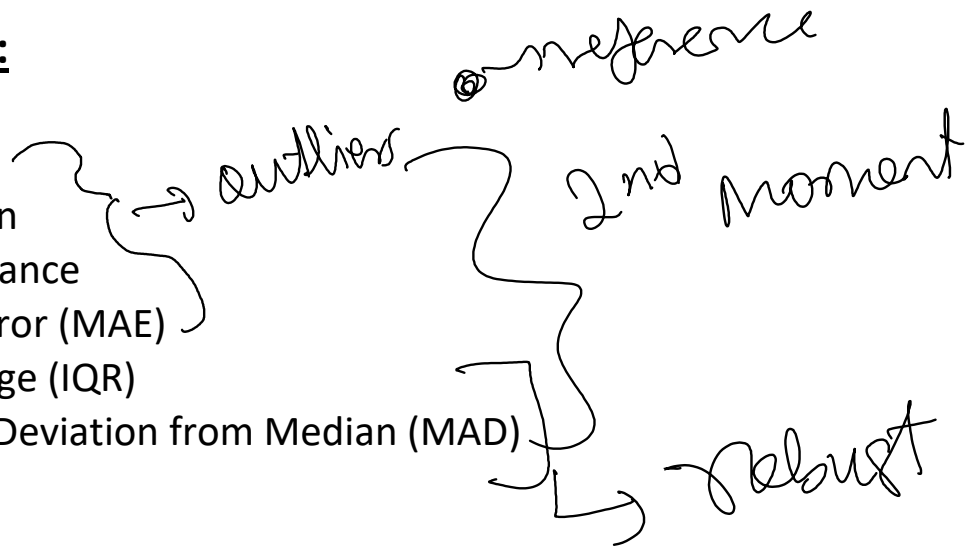    - Univariate
    - Bivariate
    - Multivariate

Summarize
Data

no.s

Sample

Population

# Descriptive Statistics

## 1. **<u>Measure of Location:</u>**

- Mean
- Mode
- Median
- Percentiles
- Quartiles

} 1st moment of distribution

## 2. **<u>Measure of Spread:</u>**

- Variance
- Standard Deviation
- Coefficient of Variance
- Mean Absolute Error (MAE)
- Inter Quartile Range (IQR)
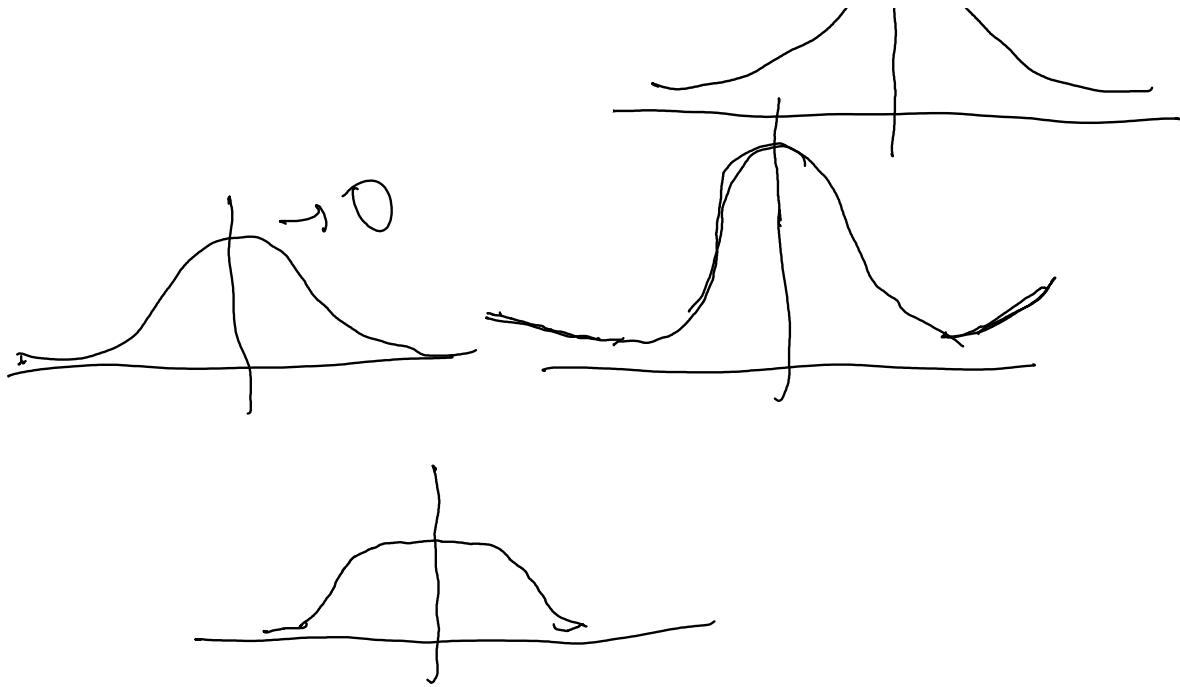- Median Absolute Deviation from Median (MAD)

→ outliers

@ reference

2nd moment

→ robust

## 3. **<u>Measure of Symmetry:</u>**

- Skewness

3rd moment

right left skewed

→ right left

## 4. **<u>Measure of Shape:</u>**

- Kurtosis

4th moment

→ symmetric

# Inferential Statistics

$x_1$  $x_2$

## 1. **<u>Strength of Association:</u>**

- Pearson's Correlation
- Spearman's Rank Correlation
- Cramer's V

2 Categorical variables

## 2. **<u>Hypothesis Testing:</u>**

- ***<u>Test for Normality:</u>***

  - Shapiro-Wilk Test
  - Anderson-Darling Test

  } Normality → numeric

- ***<u>Test for Association:</u>***

  - <u>Numeric Variables:</u>

    - Pearson's Test
    - Spearman's Test

  - <u>Categorical Variables:</u>

    - Chi-Square Test

  - <u>Numeric - Categorical Variable:</u>

    - One-way ANOVA Test
    - Kruskal-Wallis Test

- ***<u>Steps Involved:</u>***

ANOVA

Cat 1 _ m1

Cat 2 _ m2

Cat 3 _ m3

○ **_Steps Involved:_**

- State the hypotheses:

  □ *Null Hypothesis* $H_0$

  □ *Alternate Hypothesis* $H_A$

  num. var

- Determine significance level (alpha ~ 5%) $\alpha$

- Determine which test to perform — Test

- Collect necessary data (sample) — Data

- Obtain Critical values ⟶ (based on ②)
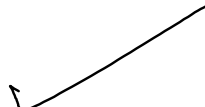
- Compute Test Statistic (and p-value) ↳ (sample)

- Compare:

  □ *Significance level vs p-value; or*

  □ *Critical values vs Test Statistic*

- State conclusions

frequency Table

↓

Chi-Square Test

$\downarrow$

$X^2$ Test-Statistic

Frequency Table $\longrightarrow$ Chi-Squar Test $\longrightarrow$ $X^2_2$

Chi-square Statistic

$\downarrow$

Cramer's V

# Plots

## 1. <u>Univariate:</u>

- ***Numeric:***

  - Histrogram
  - KDE Plot $\rightarrow$ (Smoothed Histogram)
  - Rug Plot
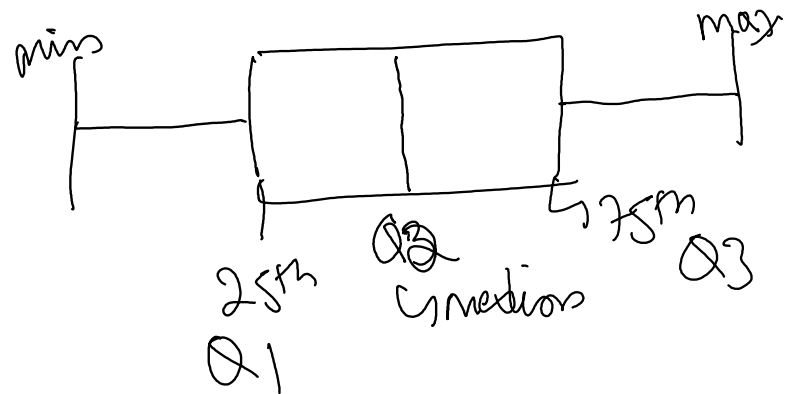  - Box Plot
  - Violin Plot (box plot + kde plot)
  - Q-Q Plot

- ***Categorical:***

  - Count Plot
  - Pie Chart

- ***Time-related:***

  - Line Plot
  - Aggregated Line Plot

$Q3 + 1.5 \times IQR$

$Q1 - 1.5 \times IQR$

## 2. <u>Bivariate:</u>

- ***Numeric - Numeric:***

  - Scatter Plot
  - Hexagonal Bin Plot
  - Contour Density Plot

- ***Numeric - Categorical:***

  - Bar Plot

- ○ Box Plot
- ○ KDE Plot
- ○ Violin Plot

- • *Categorical - Categorical:*

  - ○ Bar Plot
  - ○ Stacked Bar Plot
  - ○ Frequency Heatmap
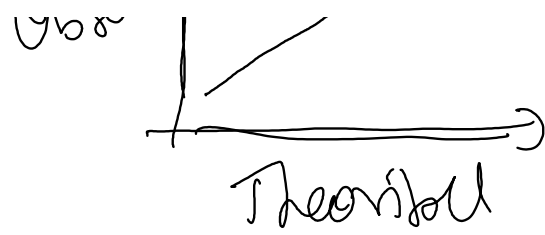
# 3. **Multivariate:**

- • *Pair Plots*

- • *Correlation Heatmap:*

  - ○ Pearson
  - ○ Spearman's Rank
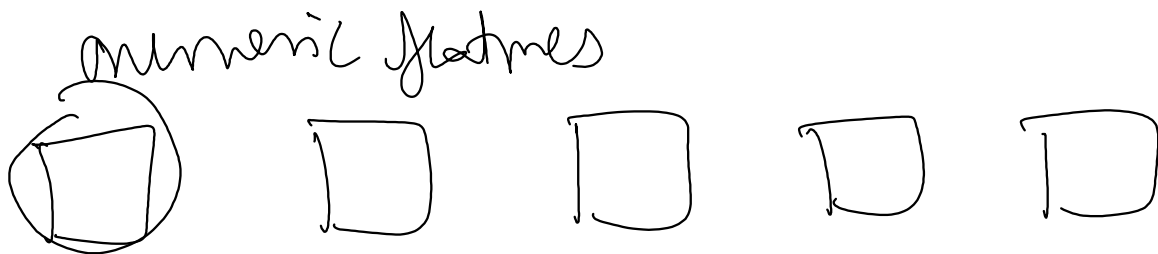  - ○ Cramer's V

- • *Facet Grid (Seaborn)*
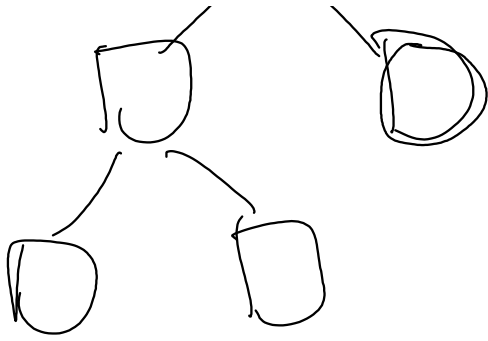
$v_{b}\sigma$

Theorical

# Steps

## **Sequence of Steps:**

- Import libraries

- Read the '*training*' data subset

    ○ Check data types
    ○ Fix data types (if applicable)

- Gather high-level summary of the data

    ○ .info() method
    ○ .describe() method on numeric and categorical features

- High-level analysis of missing values:

    ○ Bar plot
    ○ Count plot
    ○ Missingno

- High-level analysis of outliers:

    ○ Isolation Forest

- Pair plots

- Correlation Analysis (heatmaps)

    ○ Numeric (Pearson's / Spearman's)
    ○ Categorical (Cramer's V)

- Detailed Analysis of each Feature:

    ○ Summary

- ○ Univariate plots
- ○ Bivariate plots (w.r.t. the target variable)
- ○ Hypothesis Testing (normality, strength of association)
- ○ Multivariate plots
- ○ Inspect missing values and extreme values in-depth
  - ▪ Filter for necessary subsets
  - ▪ Inspect values of other features (plots, summary stats)
- ○ Note observations

- • Feature Engineering

  - ○ Create new features
  - ○ Repeat above steps for newly created features
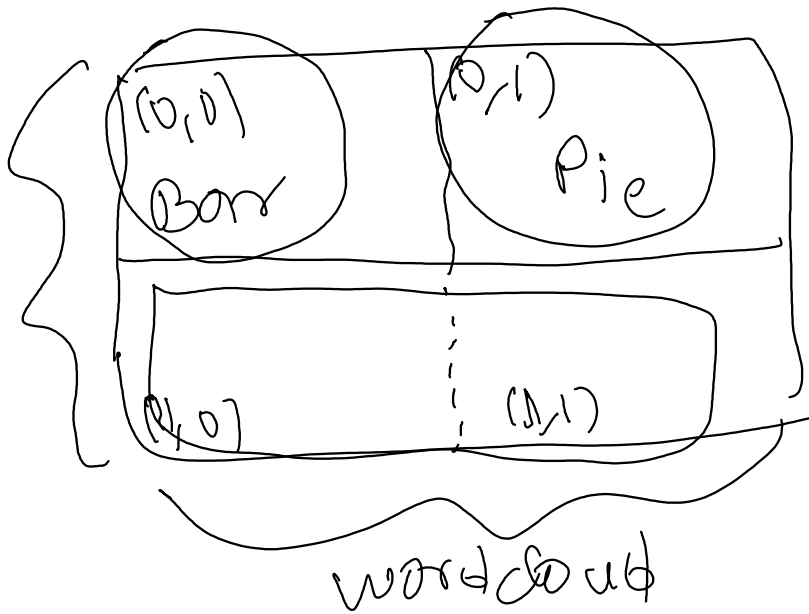
- • Repeat above steps iteratively

- • Note all observations

numeric features

1. Randomly pick one feature

2. Randomly pick one value

−1, 1

bar
step
pie

bar   step   pie

(0,0)
Bar

(0,1)
Pie

(1,0)          (1,1)

wordcloud

color

(r, g, b)

(0-1)

(0-255, 0-255, 0-255)

plt.figure()

axes → ax

plt. figure()

axes → ax

| False | False | F |
|-------|-------|---|
| False | False | F |
| F | F | P |

→

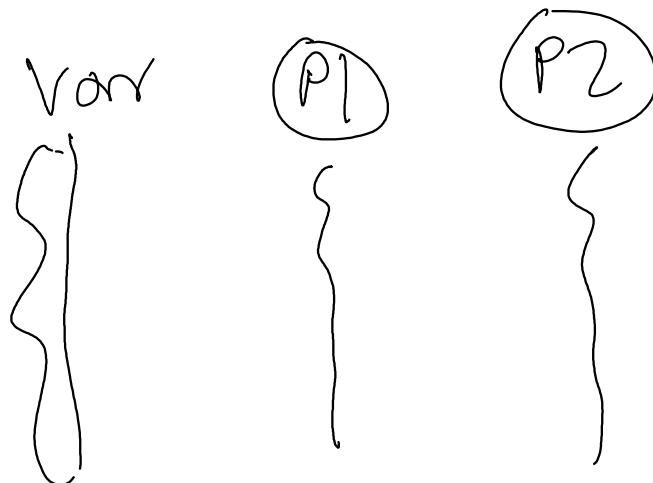| T | T | T |
|---|---|---|
| F | T | T |
| F | G | T |

"corr()"

A ——|—— Q1 —— Q2 —— Q3 —————— B

$$\text{IQR} = Q3 - Q1$$

$$B = Q3 + (1.5 \times \text{IQR})$$

$$A = Q1 - (1.5 \times \text{IQR})$$



P1  P2

Q1  Q2

|cat1  cat2  cat3  cat4

Var  (P1)  (P2)

{ { {

[ cond1, cond2, cond3, .....]
[ val1, val2, val7, .....)
defanlt = "abc"

[1,2, 10, 50, 80, 100, 120)
5% ⏜ ↓ mean 5%

1. Median
2. Var — median
3. |var — median|

# 4. median

var } impacted outliers
std

IQR } not imp. outliers
MAD



$$X$$

$$Power(X) = X^2$$

Power

price

blur

# Automated EDA

Tuesday, April 23, 2024     12:52 AM

- **Pandas Profiling (ydata-profiling)**

- **Sweetviz**

- **Autoviz**

- **D-Tale**