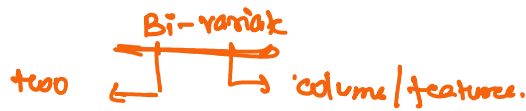


Friday, 5 December 2025 5:56 PM

• What is Bi-variate analysis?

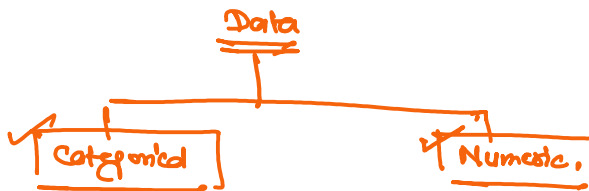


• So, here we analyze two variable together.

• Why two variable?

- How one variable affects another
- What kind of relationship or pattern exist b/w the two.
- Whether one variable depends on another.

3 possible combination in Bi-variate analysis.



• Possible scenarios.

• Numerical - Numerical. →

- Age v/e salary.
 - Height v/e weight.
 - Price v/e Distance.
- Both the columns are numerical
• Continuous/discrete.

• Numerical v/e Categorical

- Salary v/e Gender
 - Marks v/e department
 - Age v/e City.
- one column → numerical
other col → categorical.

• Categorical v/e Categorical

- Gender v/e Purchase
- City v/e Product Type
- Vehicle Type v/e Fuel Type

• Goal

• Univariate analysis

- Describe one variable → statistical summary.
- Identify the distribution, range, outliers, missing value etc.

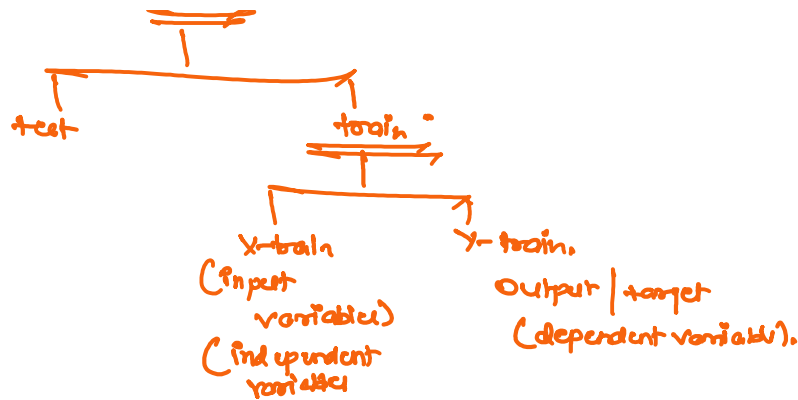
• Bi-variate analysis

- To identify the relationship / pattern b/w two variables.

• Why do we need it (Important)

In ML

Data



- What does the algorithm do?
 - It learns the relationship b/w X and Y.
 - It builds a model that expresses how X influence Y.

Model (data + algo) → trained algorithm, on the data.

- A model is nothing but a mathematical representation of relationship.
 - ↳ we must study the relationship manually using bi-variate analysis.

- Type of analysis under Bi-variate Variable

- ↳ Visual. → graphs or plots.
- ↳ Non-visual. → statistical and numerical measure

But which one to use & when to use depend upon.

- Whether the variables → numerical,
 - ↳ categorical
 - ↳ mixed.

- Why study relationship?

Suppose In a dataset



$$\text{price} = X_1 (\text{sqft}) + X_2 (\text{bedroom}) + X_3 (\text{locality}) + X_4 (\text{Age of})$$

- Numerical - Numerical

- Non-visual method → Correlation.

- Pearson Correlation Coefficient

- Measure linear relationship
- Works well data is normally distributed.
- Value range -1 to 1.

- Spearman Rank Correlation Coefficient

- Measure monotonic relationship. (non-linear)
- Works well data → not normally distributed

Both these methods help us to analyze.

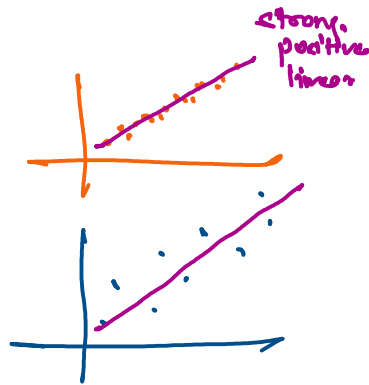
When $x \uparrow$, does $y \uparrow$
 $y \downarrow$
 y stay unrelated.

Visual approach \rightarrow Scatter plot.

Why?

Trends.

- Direction. (positive / -ve).
- Strength.
- Clusters.
- Outliers.
- Linear / non-linear patterns



Case-2 Categorical v/s Categorical.

Non-visual approach \rightarrow Crosstab.

It helps us to understand:

- Which category combination are common
- Which combination rarely happens
- How do categorical features interact with each other

Crosstab

Gender v/s Purchase \rightarrow this directly shows if male or female purchase more.

2. Visual method (with discrete zoom)

- Clustered bar chart
- Stacked bar chart
- Heatmap of Crosstab. \rightarrow b/w two variables
- Mosaic plot

Case-3 Numerical v/s Categorical.

Non-visual approach

- Groupby + descriptive statistics (Aggregate functions.)
- Pivot table.

Gender v/s Salary.

It tells us.

- average salary of male
- average salary of female

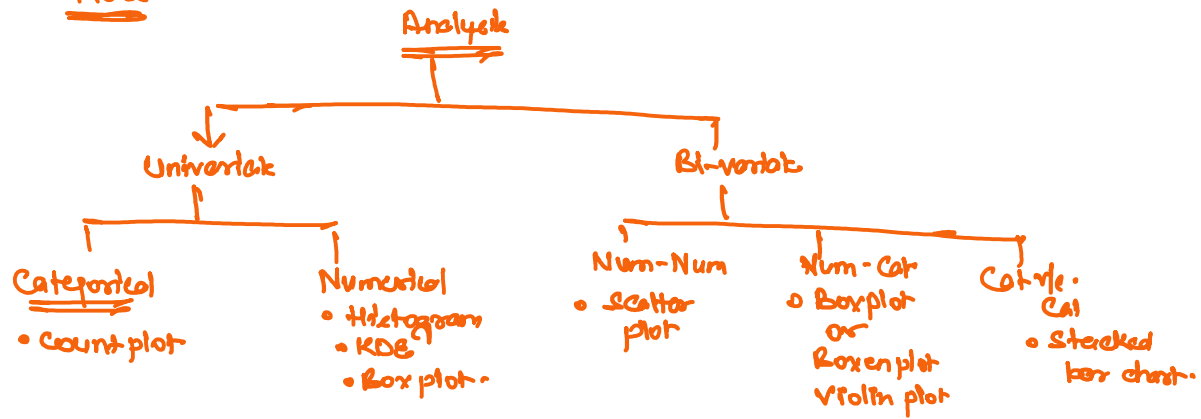
- min
- max
- median
- standard dev/s
- Count
- percentile
- value counts.

Visual approach

Box plot (Best choice)

• Boxen plot • (Alternative),

• Plots



Total plots we covered.

Univariate

- Count plot
- Histogram
- KDE
- Box plot

Bivariate

- scatterplot
- Box plot (Boxen plot) / violin.
- stacked bar chart
- Violin plot (optional).

(7-8) → extremely powerful plots