

Thursday, 4 December 2025 6:00 PM

* Re-learning EDAEDA

- Univariate
- Bi-variate
- but very few know the purpose behind them.
 - why do we perform univariate analysis
 - bi-variate analysis
- what exactly should we analyze under univariate & bi-variate
- what mistakes we should avoid.

3. major points.

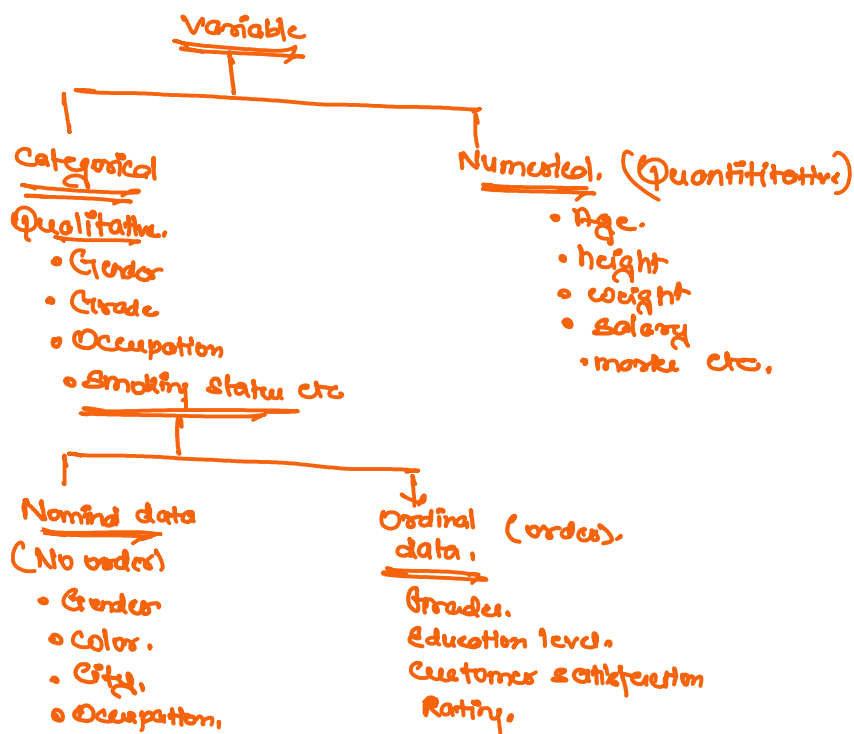
- What comes under univariate analysis.
- What comes under bi-variate analysis.
- Why each part is important from the modelling perspective.

Uni-variate

Uni → one
 Variate → variable.

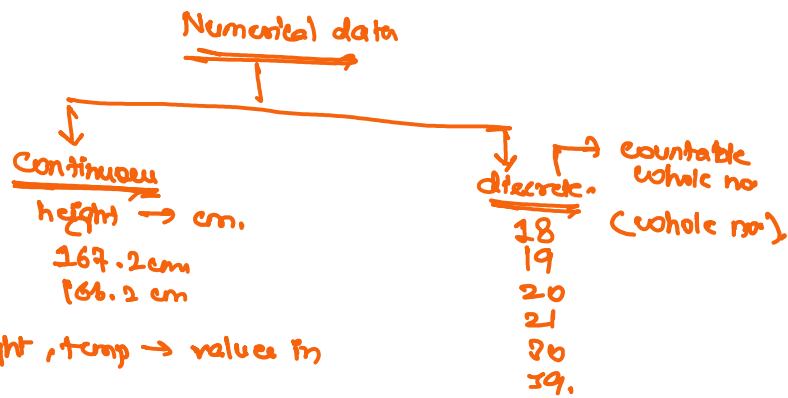
Analyze one - single variable

- Type of Analysis.
 - Visual
 - Non-visual analysis (Statistics)

Bi-variate

This distinction is very important because it will affect

- How we visualize.
- What statistical measure we calculate
- How we create these data (columns) before modelling.



* Uni-variate analysis

Categorical column

- Non-visual analysis
- Visual analysis

* Non-visual analysis

✓ Count of value in the column

count() or df.shape()

✓ Most frequent value → mode

✓ Number of unique values → nunique()

Before doing any kind of distribution analysis

How many unique categories are present.

• nunique → no. of distinct categories.

Technical term for no. of unique categorical values.

Cardinality.

High Cardinality → many unique values.

Low cardinality → few unique values.

✓ Distribution of categories → value - count()

value - count() is most important non-visual analysis for categorical value.

Titanic

male 14





✓ Actual unique labels — unique ().

gender → [male and female]

gender.nunique — 2.

gender.unique — [male, female].

gender.value-count()

male
female

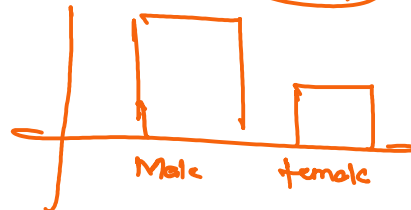
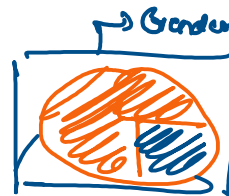
Visual analysis

- ✓ count plot } → frequency
- ✓ bar plot }
- ✓ pie chart → %

Univariate
— scatter plot
↓
2 column.

value-count → count plot / bar plot → frequency.

↳ pie chart → %



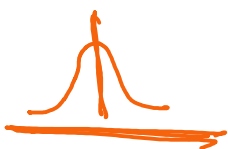
• Numerical

• Non-visual analysis

- ✓ Count → how many values are there in the column.
- ✓ Mean (Average) / Median → average value
- ✓ Min → minimum value
- ✓ Max → maximum value
- ✓ Standard deviation

$$\sigma^2 \text{ or } \text{Var} = \frac{\sum (x - \mu)^2}{N} \rightarrow \text{square of unit}$$

$$\text{std} = \sqrt{\sigma^2} = \text{Better measure} \rightarrow \text{same unit}$$



• A low standard deviation → values are close together around the mean

• A high standard deviation → values are widely spread around the mean.



• We can also check

- Quantile. / Percentile.
- Skewness
- Kurtosis.

95 percentile → means 95% of the total population who have attended the exam score lower than yours.
 $\frac{\%}{100} \rightarrow \frac{\text{your score} \times 100}{\text{total cr}}$

{ IQR →
MAD

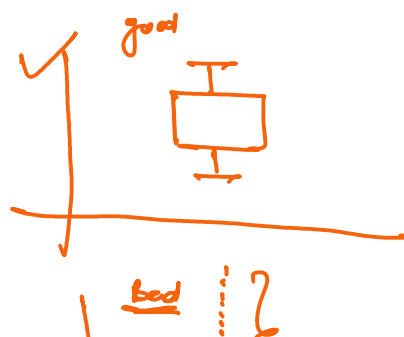
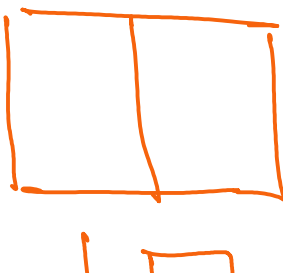
* Visualization for Numerical column

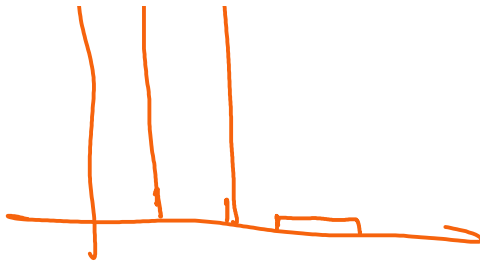
- ✓ Histogram. → frequency distribution
- ✓ KDE plot → smooth continuous probability distribution.
 ↳ good for understanding skewness
- ✓ Box plot
 ↳ Best for detecting.
 - Outliers
 - Median.
 - Quantile →
 - Spread (IQR).

WHY ?

- To understand the data. → help us to describe, summarize and understand the data.
- To learn about the distribution data.
 - Distribution tell us
 - Is the data is balanced or imbalanced?
 - Is it skewed?
 - Are there outliers.
 - Are the categories are evenly spread or highly dominated.
 - Is the data symmetric, bimodal, long tailed.
 - What is spread or range.

• Feature selection





• To detect Outliers

- Outliers need to be treated.
- Mixed our mean.

$$\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3.$$

$$\frac{1+2+3+4+100}{5} = \frac{110}{5} = 22.$$

- Affect ML model performance
- Create incorrect interpretation.

• To identify missing values

Reason to do Univariate Analysis

- Describe the data
- Learn data distribution
- Feature transformation
- Outlier detection
- Identify missing values
- Feature selection
- Business insight

