

→ Exploratory Data Analysis Explore the data
meanh analyze and Investigate the data.
Summarized their main characteristics.
Employing visualization method.

1. Univariate Non graphical.

- Univariate meanh Single feature or column.
- Description or Summary of this feature.

- Mean
- Median
- Mode
- Count
- Quantile
- Percentile
- Min
- Max
- Standard deviation
- Skewness
- Outliers detection
- Variance.

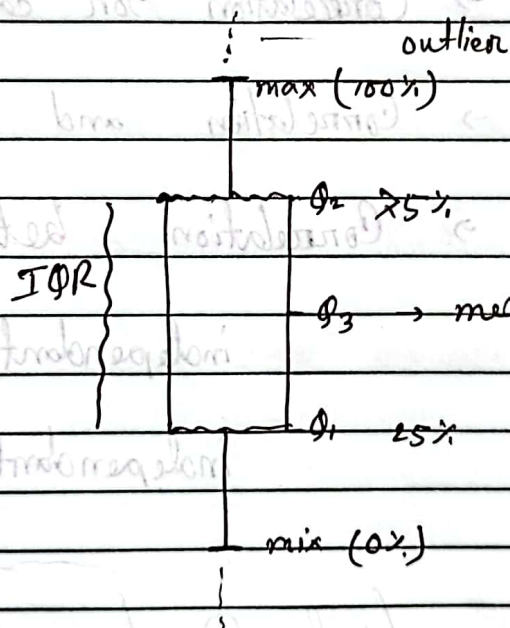
2. Univariate Graphical \rightarrow Non Graphical can't provide full picture of data.

\rightarrow Summarize or Describe data in a Graphical way.

\rightarrow Histogram.

\rightarrow Stem and leaf plot

\rightarrow Boxplot



five number
Summary.

\rightarrow Bar - plot.

\rightarrow Pie-chart.

3. Bivariate means two feature.

Multivariate means multiple feature.

→ Multivariate non-graphical.

→ Describe the data non graphical way.

→ Relationship between two or more variables. (tabular format).

→ Correlation for categorical variable

→ Correlation and Covariance.

→ Correlation between

independent vs independent

independent vs dependent.

→ Multivariate graphical: Display the Relationship between two or more variables.

→ Describe the information in a graphical way.

→ Scatter plot.

→ Pair plot

→ Heatmap (Correlation)

→ Facetgrid

→ Clustermap

→ Side by side box-plot.

→ Contingency table is one of the techniques for exploring two or more variables in a tabular or matrix format.

[Multivariate non graphical]

Preprocessing

- Data cleaning.
- Data integration.
- Data Reduction.
- Data Transformation.

Data Quality -

- Accuracy.
- Completeness.
- Timeliness.
- Consistency.
- Reliability.
- Interpretability.

Accuracy
Completeness
Consistency

Data cleaning,

- Remove outlier, handle missing value.
- smoothout - noisy data
- Redundancy
- Remove - inconsistency

Data Reduction. Reduce the data size or dimension.

- Dimensionality Reduction
- Numerosity Reduction.
- Compressed.

Data Transformation .. Discretization,

Normalization

Standardization

→ Raw data values for attributes are replaced by Range or higher Conceptual levels.

Age → youths, adult, Senior.

Handle missing value,

→ Mean, Median, Mode,

→ Constant,

→ Most frequent.

→ Remove or skip or ignore

→ Manually.

→ Build a model for ^{Predict} missing value. (etc)

Noisy : Noise is a Random error or variance in a measured variable.

→ Binning

→ Regression analysis

→ Outlier Analysis

Data cleaning Process,

① → Discrepancy detection,

- human error data entry
- poorly design data " forms
- inconsistency data representation
- System Error.

→ Data Reduction

→ Data Reduction strategy