

U4-(BI)

-90-21287

Data preparation and related to initial level of data quality
 ↳ cleansing & transforming raw data prior to processing of analysis.

steps of Data Preprocessing -

- 1) Data Cleaning (filling missing data, smooth noise, resolve inconsistencies)
- 2) Data Integration (put together due to from diff set resolve conflicts)
- 3) Data Transformation (normalize, aggregate, generalize)
- 4) Data Reduction (reduce representation dimensions)
- 5) Data Discretization (reduce numerical values of continuous attrib. into intervals)

Data cleaning -

make invalid data valid by correcting it by:

- 1) filling missing values.
- 2) unified data format
- 3) convert to network data.
- 4) remove noise, identify & correctly remove outliers.
- 5) correct inconsistent data.

Incomplete data

→ lack of data

→ presence of missing values

→ techniques to adapt:

- 1) Elimination
- 2) Inspection
- 3) Identification
- 4) Substitution

~~Data affected by Noise~~ (Handle missing values) →

- ~~ignore the record~~
- fill missing values manually
- use global constant
- use attribute avg, mean, max, mode.

* Data affected by noise

- noise is random error or variance in measured variables values.
- data collection, human error, data transmission error.
- random disturbance within value.

- 1) Numeric Noise \rightarrow plot a boxplot, scatter plot or dispersion
- 2) categorical data \rightarrow clustering. Outliers away from all cluster

Data smoothing techniques

- 1) binning \rightarrow bin median, bin max, bin min.
- 2) Regression - factor values
- 3) outlier analysis.

* Data Transformation

- transforming or consolidation of data into appropriate form.
- Techniques -
- standardization (Normalization)

\hookrightarrow entire data set values have particular property is transform it to have a particular distribution or centred around a particular value.

A) Decimal Scaling -

$$\text{transformation } x_{ij} = \frac{x_{ij}}{10^n} \quad n = \text{smallest values that range from } [-1, 1]$$

B) Min Max Scaling -

values are transformed to fit a certain scale [range] $[0, 1]$

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} (x_{\text{newmax}} - x_{\text{newmin}}) + x_{\text{newmin}}$$

$$\text{eg } \rightarrow x = 5 \quad x_{\text{max}} = 10 \quad x_{\text{min}} = 2.$$

$$x' = \frac{5-2}{10-2} (1-0) + 0 = \frac{3}{8} (1) + 0 = 0.375$$

C) Z-score transformation \rightarrow based on mean & standard deviation

$$x' = \frac{x - \bar{x}}{\sigma} \quad \text{when } \sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

- ② Feature Extraction
 ↗ Data Reduction → Sampling → random sampling.
 ↗ Feature Selection → redⁿ in no. of attributes by selection & projection.

Filter	Wrapper	embedded
• generic method no ML algo	• evaluates specific ML algo	• done by observing each item
• fast	• slow (high computation)	• medium
• less prone to overfitting.	• high chance of overfitting	• used to reduce overfitting
• correlation based selection.	eg - forward selection, backward elimination, stepwise selection, etc.	• penalizing coefficient
• chisquare, t-test, ANOVA, etc.		• eg → LASSO, elastic net, ridge regression, etc.
• more features, more complex the model, diff to train & interpret.		

* Principal Component Analysis

- unsupervised ML algo
- used for dimensionality reduction
- converts observation of correlated function into set of linearly features
- new transformed features → principal component
- popular tool → EDA & predictive modeling
- technique to draw strong patterns from reducing variance
- Appln → Image preprocessing, movie recommendation system.
- contains imp variables & deletes / drops least important variables

Steps -

$$1) \text{ Standardization} \rightarrow z = \frac{x - \bar{x}}{\sigma}$$

2) Covariance matrix computation -

- calculate strength of joint variability of 2 more variable
 how much they change in redata to each other.

$$\text{Cov}(x_1, x_2) = \frac{\sum (x_{1j} - \bar{x}_1)(x_{2j} - \bar{x}_2)}{n-1}$$

$$>1 = \uparrow\uparrow \\ <1 = \uparrow\downarrow$$

$0 = \text{no reln}$

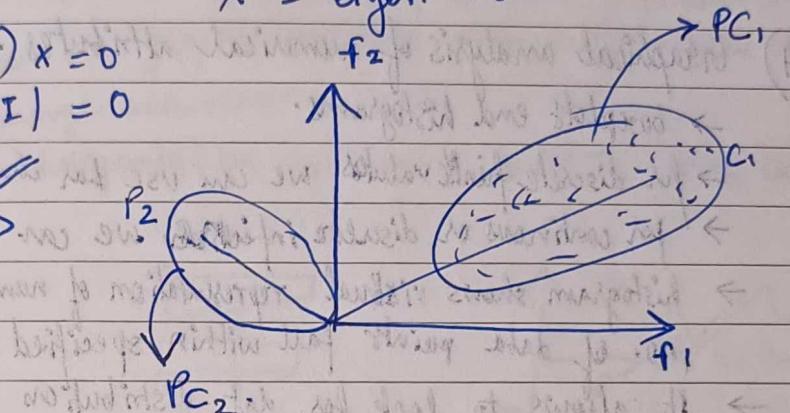
3) Calculate Eigen values & Eigen vectors.

$$A\mathbf{x} = \lambda \mathbf{x} \quad \lambda = \text{eigen vector}$$

\mathbf{x} = eigen vector.

$$(A - \lambda I)\mathbf{x} = 0 \\ |A - \lambda I| = 0$$

$$\lambda \approx \\ \mathbf{x} \Rightarrow$$



3) Data Discretization - reduction in number of values through discretization & aggregation.

- decreases number of distinct values.

- continuous values are converted to categorical attributes

e.g. \rightarrow weekly expenditure. $\rightarrow 101, 103.5, 45, 35$.
instead we put it in range low (0-20), high (100-400) etc.

techniques \rightarrow

- 1) Subjective Subdivision (based on expertise & judgement of experts)
- 2) subdivision into classes (automated classes based on equal size or width)
- 3) Hierarchical discretization (applied to categorical attributes)

Methods -

- 1) Binnings \rightarrow freq. equal width (list).
- 2) Histogram Analysis.
- 3) Decision Tree Analysis
- 4) Correlation Analysis

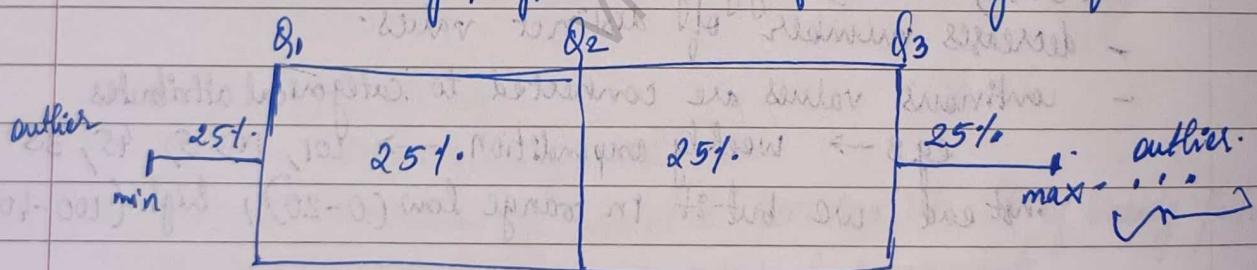
* Data Exploration -

- 3 phase →
- 1) Univariate → study each attribute's property are investigated
 - 2) Bivariate → pair of attributes are considered
 - 3) Multivariate → rel' b/w subject of attributes

A) Graphical analysis of numerical attributes -

→ boxplots and histograms.

- for discrete finite values we can use bar chart.
- for continuous or discrete intervals we can use histogram by grouping values into bins.
- histogram shows visual representation of numerical data by showing no. of data points fall within specified range of values (bins).
- it allows to look for data distribution.
- Box plot → summarizes data statistics, demonstrate like locality, typical skewness of data through quartile



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

$$\text{Lower limit} = Q_1 - 1.5 \times \text{IQR}$$

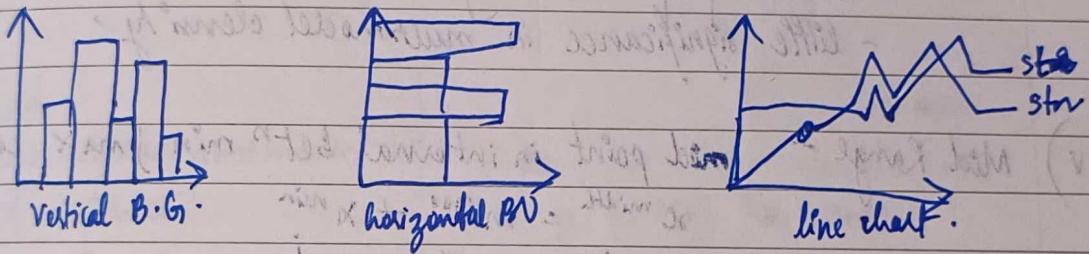
$$\text{Upper limit} = Q_3 + 1.5 \times \text{IQR}$$

$Q_1 \Rightarrow$ Median of 0 to Q_2 (median) 25%

$Q_3 \Rightarrow$ median of Q_2 to n (75% quartile)
outliers by opt. of lower & upper limit.

b) Graphical analysis of categorical attributes

- most natural representation vertical bar chart.
- vertical bar chart value indicates categories & respective frequencies.
- sometimes horizontal bar chart is also used.
- pie chart is used to show empirical density function values in +.
- sometimes line graphs are used when one wants to see how values change over time.
- Adv of horizontal bar chart → easier to display long labels.



c) Measuring Central Tendency - (describe main loc" statistics).

- i) Mean - avg of data points $M = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$
- strongly affected by extreme values (outliers) hence not robust
- deviation - sum of diff b/w each values of mean = 0 $\sum_{i=1}^m (x_i - M) = 0$

- it minimizes sum of squared deviations from constant ref. values.

$$\sum_{i=1}^m (x_i - \mu)^2 = \min_{\mu} \sum_{i=1}^m (x_i - \mu)^2$$

- It minimizes sum of weighted mean $\rightarrow M = \frac{w_1 x_1 + w_2 x_2 + \dots + x_n w_n}{w_1 + w_2 + \dots + w_n} = \frac{\sum w_i x_i}{\sum w_i}$

$x \rightarrow$ shipping cost of product A

$w \rightarrow$ no. of A products shipped.

ii) Median - central value of all observations.

$$x_{\text{med}} = \frac{x_{m+1}}{2} \quad \text{if } m \text{ n odd} \quad x_{\text{med}} = \frac{x_{m/2} + x_{m+2/2}}{2}$$

- if affected by no. of elements in sample due to two is more robust mean as not affected by outliers.

- suitable measure for arithmetic distributions.

- if data is not concentrated at central part median if no outliers

iii) Mode - value corresponding to peak of density curve is most freq occurring values.

- little significance in multimodal density.

iv) Mid Range - mid point in interval between min & max column

$$x_{\text{mid}} = \frac{x_{\text{max}} + x_{\text{min}}}{2}$$

$$x_{\text{min}} = \min x_i$$

$$x_{\text{max}} = \max x_i$$

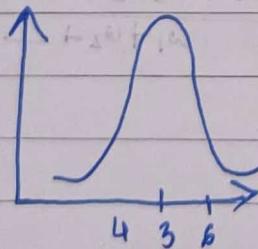
- not robust as effected by outliers.

v) Geometric Mean - m^{th} root of product of m observations of attribute.

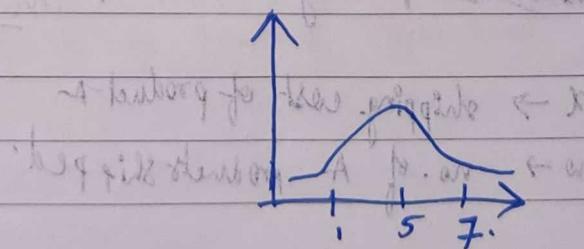
$$g_{\text{geom}} = \sqrt[m]{x_1 x_2 \cdots x_m}$$

D) Measures of dispersion for numeric attributes-

dispersion of data represents based of variability expressed by observation wrt to central values (location said in central)



less dispersion
(4, 5).



more dispersion.

i) Range = diff b/w 2 extreme obs. \therefore range = $x_{\max} - x_{\min}$

ii) Variance = measure of variability of data \therefore simple variance = $(\sigma^2) = \frac{\sum (x - \bar{x})^2}{n-1}$
 pop variance $\sigma^2 = \frac{\sum (x - \mu)^2}{N}$

iii) standard Deviation (σ) (s) - positive sqrt of variance

$$\sigma = s = \sqrt{\frac{(x - \bar{x})^2}{n-1}}$$

std

variance

- measures of dispersion from mean
- how far nos are from average
- sort of variation
- avg in diff from mean
- spread betw nos in dataset
- avg degree to which each part diff from mean
- same unit of data
- squared unit as data or percents
- $\sqrt{\frac{(x - \bar{x})^2}{n}}$

iv) Mean Absolute Deviation (Deviation or spread). $S_i = \sum_{i=1}^m |x_i - \bar{x}|$.

v) Quartile Deviation (QDR) $Q_3 - Q_1$

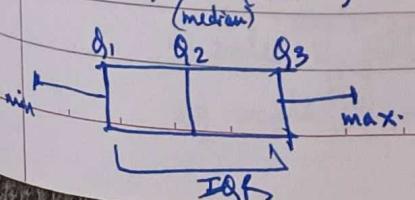
vi) Coefficient of variation.

$$CV = 100 \cdot \frac{\sigma}{\bar{x}} \quad \text{ratio of same st. deviation to sample mean expressed at \%}$$

o) Identification of Outliers.

1. using Z-index $Z \text{ score} = \frac{x - \bar{x}}{\sigma}$

2. box plot (box & whiskers plot).



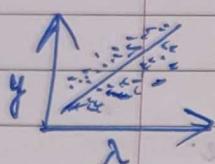
	lower bound	upper bound
inner outliers	$* 1.5 \quad Q_1 - (1.5 \times IQR)$	$Q_3 + (1.5 \times IQR)$
outer outliers	$* 3 \quad Q_1 - (3 \times IQR)$	$Q_3 + (3 \times IQR)$

2] Bivariate Analysis -

- helpful in testing simple hypotheses of association.
- 3 distinguish cases that can occur.
 - i) both attribute are numeric
 - ii) one attribute is numeric, one is categorized
 - iii) both attributes are categorical

B) Graphical Analysis -

i) Scatter Plot

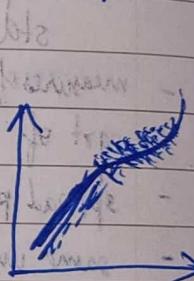


- intuitive representation of relationship b/w 2 numeric attributes
- pair of values as treated as pts & plotted on cartesian graph.

ii) Loess curve (local regression)

based on scatter plot, used for both numeric & categorical

states from scatterplot, adds a curve to express functional



trend where can be attend obtained using local regression testing R^2 .
curve is regulated by 2 parameters.

1) degree $\lambda > 0$ of polynomial that represents are usually 1 or 2 is assigned to λ .

2) usually regularization constant of $\alpha \approx 70$ regulates size of neighbourhood if $\alpha \uparrow$, neighbourhood \uparrow , more obs. used to calculate regression.

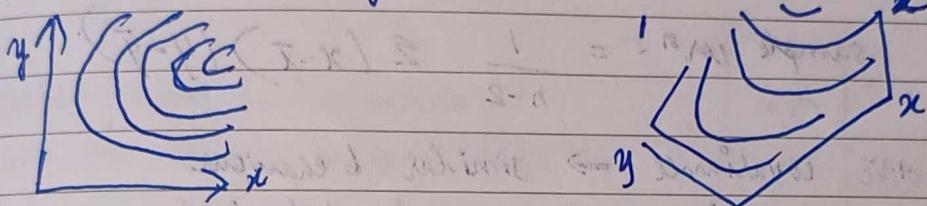
α in range of $1/4$ to 1 is chosen normally

* note that α does not affect R^2

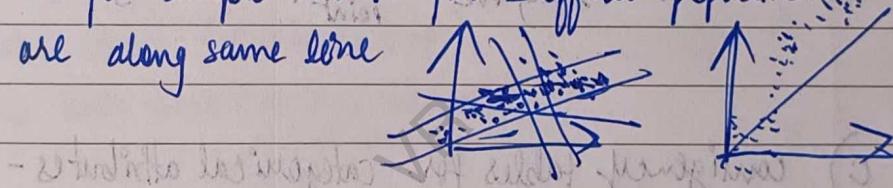
(very sensitive to noise) highly sensitive



- iii) level curves (identical to control lines as graph / maps) -
- further development of scatter plot against numeric attribute.
 - highlight value by 3rd numerical attribute a_2 as attribute $a_1 \& a_2$ placed
 - connecting to each other by points in plot share values of 3rd or excess of other attribute obtains curved lines representing geometric locus of points for attribute a_2 assumes given value -



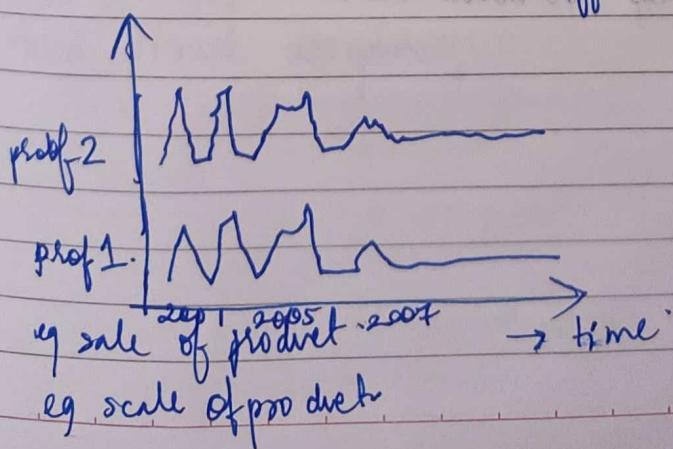
- iv) Q-Q plots (Quantile - Quantile) -
- compares distribution of some attributes for 2 diff characteristics of population or for sample extracted for 2 different populations.
 - if samples are along same line



- covariance
- extent to which 2 random variables change w.r.t to each other
 - means of correlation
 - indicates direction of linear relation
- correlation:
- how strongly 2 random variables related
 - scaled from covariance.
 - -1 to $+1$
 - directive of strength for linear relationship

v) Box Plot

vi) Time series - collection - of same attribute relative to different time instants



b) Measures of correlation for numeric attribute -

- covariance coefficient & covariance can be used to find variation b/w attributes.

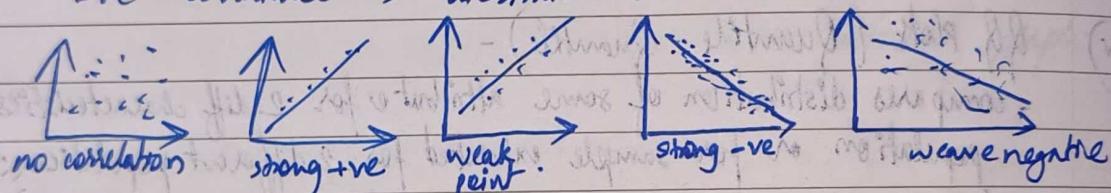
i) Covariance -

$$\text{cov}(x, y) = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) = E\{(xy) - \bar{x}\bar{y}\}$$

$$\text{sample cov}^n = \frac{1}{n-2} \sum (x_i - \bar{x})(y_i - \bar{y})$$

+ve covariance \rightarrow similar behaviour

-ve covariance \rightarrow dissimilar behaviour



c) Contingency tables for categorical attributes -

- tabular representation of categorical data, if represents frequencies for particular combination of values of attributes.
- defined as a metric + whose generic elements T_{ij} is indent freq of pair of values ($x_{ij} = v_s$) $\wedge \{x_{ijk} = 0\}$

		right hand	left hand
		male	female
male	40	11	20
female			5

- If summarizes into about data-

3) Multivariate analysis

- extends bivariate analysis in order to assess relationship among multiple attributes in dataset

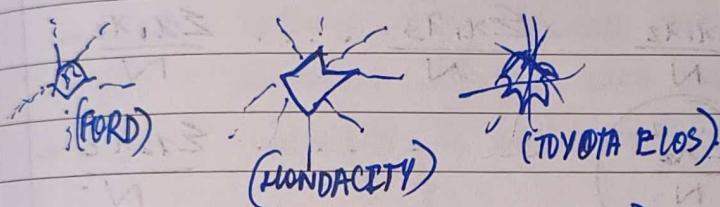
a) Graphical analysis -

i) matrix of scatter plot -

- roughly determines linear correlation b/w multiple variables.

ii) starplot (radar chart or web chart)

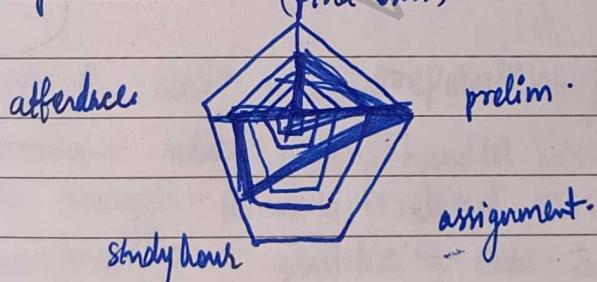
cars are analyzed for some attributes.



iii) spirled web charts (Radar charts)

- display dots across unique elements.

eg -



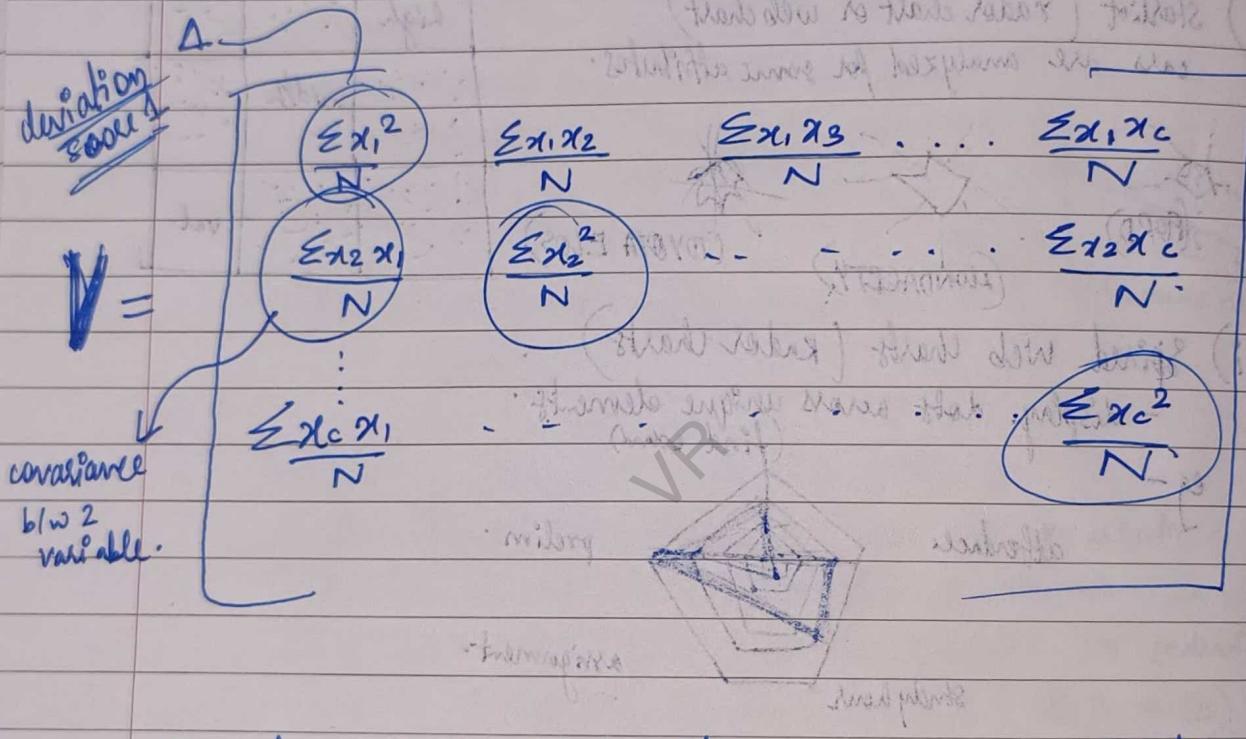
benefits -

- 1) make comparisons of strength & deficiencies with others.
- 2) clearly display important categories
- 3) define full performance in each category.
- 4) valid visual decipher

b) Measure of Correlation for Numerical Attributes -

Variance
covariance
matrix

- correlations & covariance matrix are calculated among pairs of attributes
- let V_{ij} be $N \times N$ matrix whose elements represent by covariance values & correlation values respect.
- covariance matrix V contains on its diagonal sample variance of each single value it called variance - covariance of matrix.



Univariate

- only summarize single variable at time
- don't deal with causality reln
- doesn't contain any dependent variable
- main purpose is describe
- eg \rightarrow height

Bivariate

- only summarize for two variables
- deal with cause & effect analysis is done
- does contain only one dependent variable
- main purpose is to explain
- eg \rightarrow temperature & ice sales in summer vacation

Multivariate

- summarizes more than 2 variables
- doesn't deal with causality & relationships & analysis is done
- similar to bivariate but it contains more than 2 variables
- main purpose is to study reln among them
- eg \rightarrow website user growth reln b/w variables