

3

At $(\text{DOF}, \alpha(p))$

From χ^2 distribution table, value at $(1, 0.001)$ is 10.828

Since χ^2 value which is computed is 508, which is above this value, the hypothesis can be rejected.

i.e. Gender & preferred reading { are not independent
 \Rightarrow These two attributes are strongly co-related

If $\chi^2 >$ value at $(\text{DOF}, \alpha(p))$,
 strongly co-related

For Numeric, it is correlation coefficient

$$\text{i.e. } \gamma_{A,B} = \frac{\sum_{i=1}^N (a_i b_i) - N \bar{A} \bar{B}}{N \sigma_A \sigma_B} \quad \begin{array}{l} N \rightarrow \text{no. of tuples} \\ A, B \rightarrow \text{Attributes (Numeric)} \\ a_i \rightarrow \text{values of A} \\ b_i \rightarrow \text{values of B} \end{array}$$

$$-1 \leq \gamma_{A,B} \leq 1$$

If $\gamma_{A,B} > 0$, A & B are positively correlated
 (i.e. if A ↑, B also increases)

If $\gamma_{A,B} = 0$, A & B are independent
 (No correlation between them)

If $\gamma_{A,B} < 0$, A & B are negatively correlated
 i.e. if A ↑, B ↓

Data Transformation

(4)

- The Data are transformed from original state to an appropriate new set of consistent state so that the old value can be identified by the one of the new values.

Methods for Transformation

- Smoothing - Remove noise from data using techniques as Binning, Regression, clustering
- Aggregation - Summarization, Aggregation opn. applied
- Generalization - Hierarchy concept applied
- Normalization - Attributes, ^{values} are scaled to fall within a small specified range as -1.0 to 1.0 or 0 to 1.0
- Attribute / feature construction - new attributes are constructed added to existing attributes to help mining process.

Smoothing, Aggregation - same as in Data Cleaning
↓
↓, generalization process.

form of
Data Cleaning

Serve as forms of Data Reduction.

Normalization -

useful for classification algorithms involving neural networks

✓ Normalizing outputs that help speed up the learning phase

* useful for Distance measurements

Methods -

Min-Max normalization

Z-Score Normalization
(zero-mean Normalization)

Decimal scaling

$$v' = \left(\frac{v - \text{min}_A}{\text{max}_A - \text{min}_A} \right) (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

$$v' = \frac{v - \text{avg}_A}{\sigma_A}$$

$$v' = \frac{v}{g}$$

where g is the
smallest integer
such that $\text{Max}(v') < g$

$$\text{avg}_A = \frac{1}{N} \sum_{i=1}^N m_i \rightarrow \text{mean}$$

$$\sigma_A^2 = \frac{1}{N} \sum_{i=1}^N (m_i - \text{avg})^2 \rightarrow \begin{array}{l} \text{variance} \\ \text{standard} \\ \text{dev} \end{array}$$

Q) Using Min-Max, Normalize the income from \$12,000 to \$98,000 into the range [0.0, 1.0].

Ans: Q) Suppose minimum, maximum values of the attribute "Income" are \$12000 & \$98,000 resp. which is mapped to the range [0.0, 1.0]. Find the transformation of income \$73,600.

Ans:-

$$\text{min}_A = 12000, \quad \text{Max}_A = 98000 \\ \text{new_min}_A = 0, \quad \text{new_max}_A = 1$$

value which is to be mapped, $v = 73600$

Mapped value on the range [0.0, 1.0] is v'

$$\Rightarrow v' = \frac{v - \text{min}_A}{\text{Max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A \\ = \frac{73600 - 12000}{98000 - 12000} (1 - 0) + 0 \\ = 0.716$$

Q) Normalize the following data by min-max
by z-score

200 300 400 600 1000 [within [0, 1] or
min-max]

Ans: a) $\text{min}_A = 200, \quad \text{new_min}_A = 0$
 $\text{Max}_A = 1000, \quad \text{new_max}_A = 1$

a) $v = 200$

$$\Rightarrow v' = \frac{200 - 200}{1000 - 200} (1 - 0) + 0 = 0$$

b) $v = 300$

$$\Rightarrow v' = \frac{300 - 200}{1000 - 200} (1 - 0) + 0 = 0.125$$

$$c) V = 400$$

$$V' = \frac{400 - 200}{800} = 0.25$$

$$d) V = 600$$

$$V' = \frac{600 - 200}{800} = 0.5$$

$$e) V = 1000$$

$$V' = \frac{1000 - 200}{800} = 1$$

Using Z-Score :- $V' = \frac{V - \bar{x}_A}{\sigma_A}$

$$\bar{x}_A = \frac{200 + 300 + 400 + 600 + 1000}{5} = 500$$

$$\sigma_A^2 = \frac{1}{5} \left[(200 - 500)^2 + (300 - 500)^2 + (-100)^2 + (100)^2 \right] \\ = \frac{40,000}{5} = 8000$$

$$\Rightarrow \sigma_A = 283$$

$$a) V = 200$$

$$\Rightarrow V' = \frac{V - \bar{x}_A}{\sigma_A} = \frac{200 - 500}{283} = -1.06$$

$$b) V = 300$$

$$\Rightarrow V' = \frac{300 - 500}{283} = -0.706$$

$$c) V = 400$$

$$\Rightarrow V' = \frac{400 - 500}{283} = -0.35$$

$$d) V = 600$$

$$\Rightarrow V' = \frac{600 - 500}{283} = 0.35$$

$$e) V = 1000$$

$$\Rightarrow V' = \frac{1000 - 500}{283} = 1.76$$

Data Reduction -

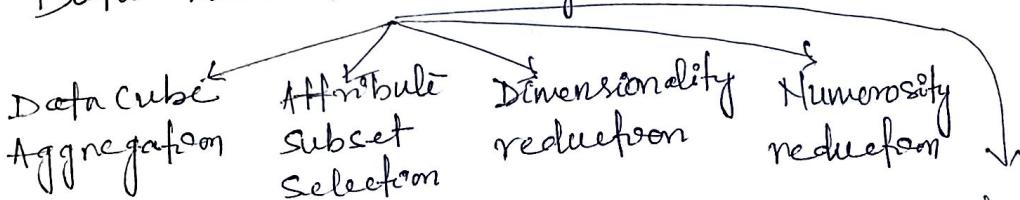
Why data reduction?

- 1) Data set may contain too terabytes of data
- 2) Since, huge data set,
 - Data Analysis } may be complex/tedious
 - Data mining }
 - Since complex,
 - it may take a very long time to analyze & mining.
- 3) Since complex,
 - it may take a very long time to analyze & mining.

Data Reduction -

- applied to obtain a reduced representation of data set that is much smaller in volume & maintain the Integrity of the original data

Data Reduction Strategies -



↓
Discretization

Data Cube Aggregation — Here aggregation operations are applied to the data to form a data cube.

- Data cube is defined by dimensions & facts
- It stores multidimensional analyses of aggregated information.

↓
Data can be aggregated so that it can be viewed on the abstract level without loss of information & also it takes less volume.

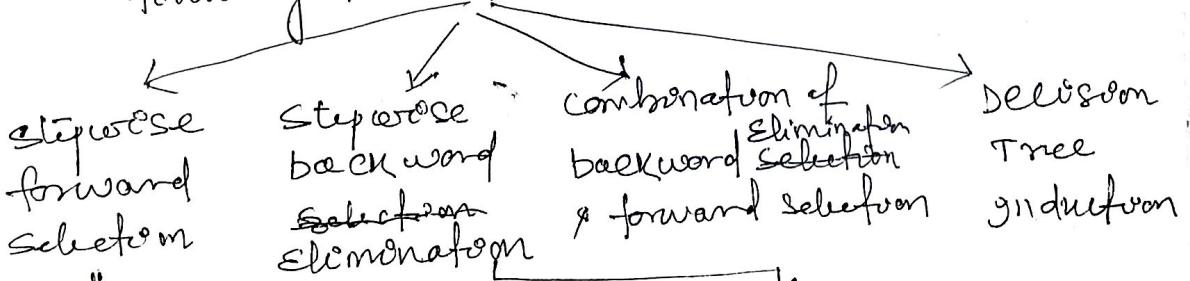
- Data cubes provide faster access to precomputed summarized data

- The cube created at the lowest level of abstraction is called Base cuboid.
- At the highest level of abstraction, called Apex cuboid.

Attribute Subset Selection :-

- It reduces the data size by removing irrelevant or redundant attributes or dimensions.
- It finds a minimum set of attributes that can be able to represent the original attributes approximately (as close as possible).
- So that it will be easier to understand a pattern.
- Heuristic methods are used that reduce the search space for attribute subset selection.
 (for n attributes, 2^n subsets can be formed.
 if we go on searching, this will be expensive
 & may take large search space)
 eg:- Greedy Approach (until getting the best choice at that time)

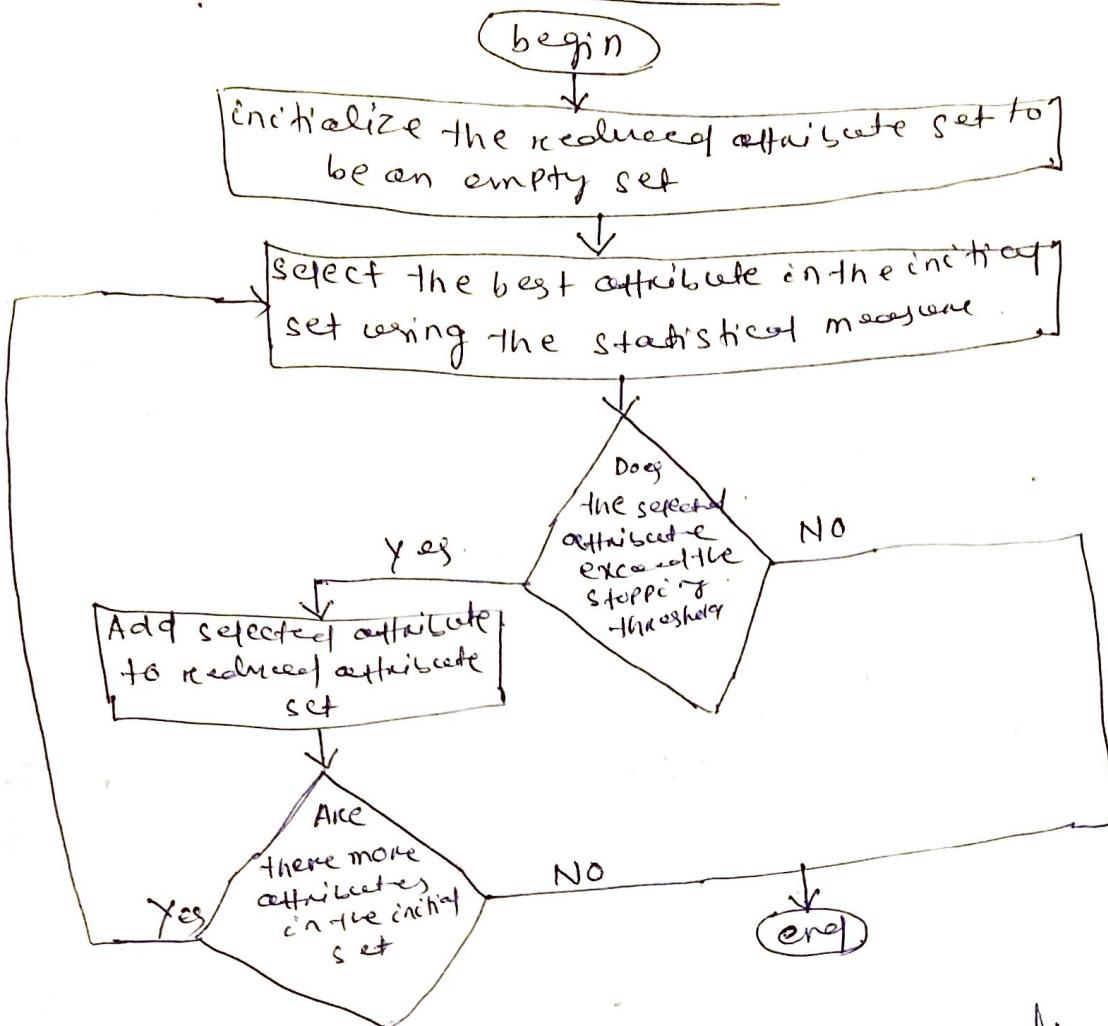
Various heuristic methods include the following techniques for Attribute subset Selection



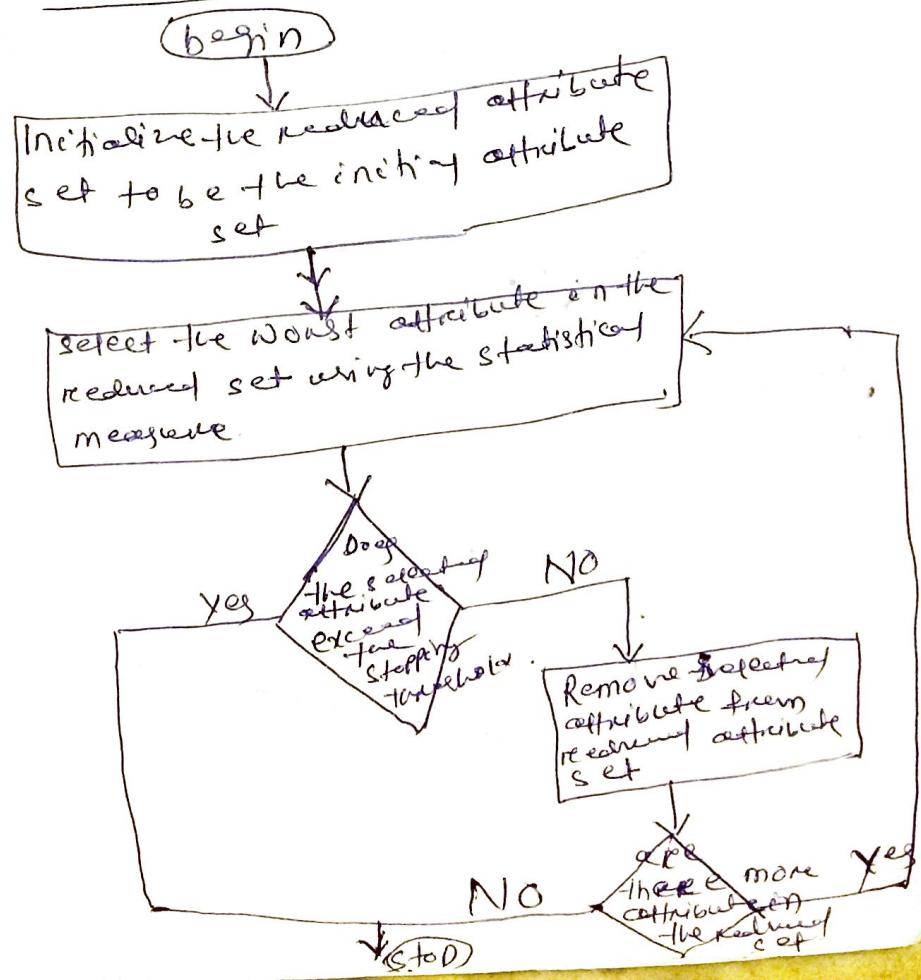
It starts with empty attribute sets. When best is determined it is added to attribute set.
 & so on.

It starts with the full set of attributes. The worst one is deleted & so on

Stepwise Forward Selection



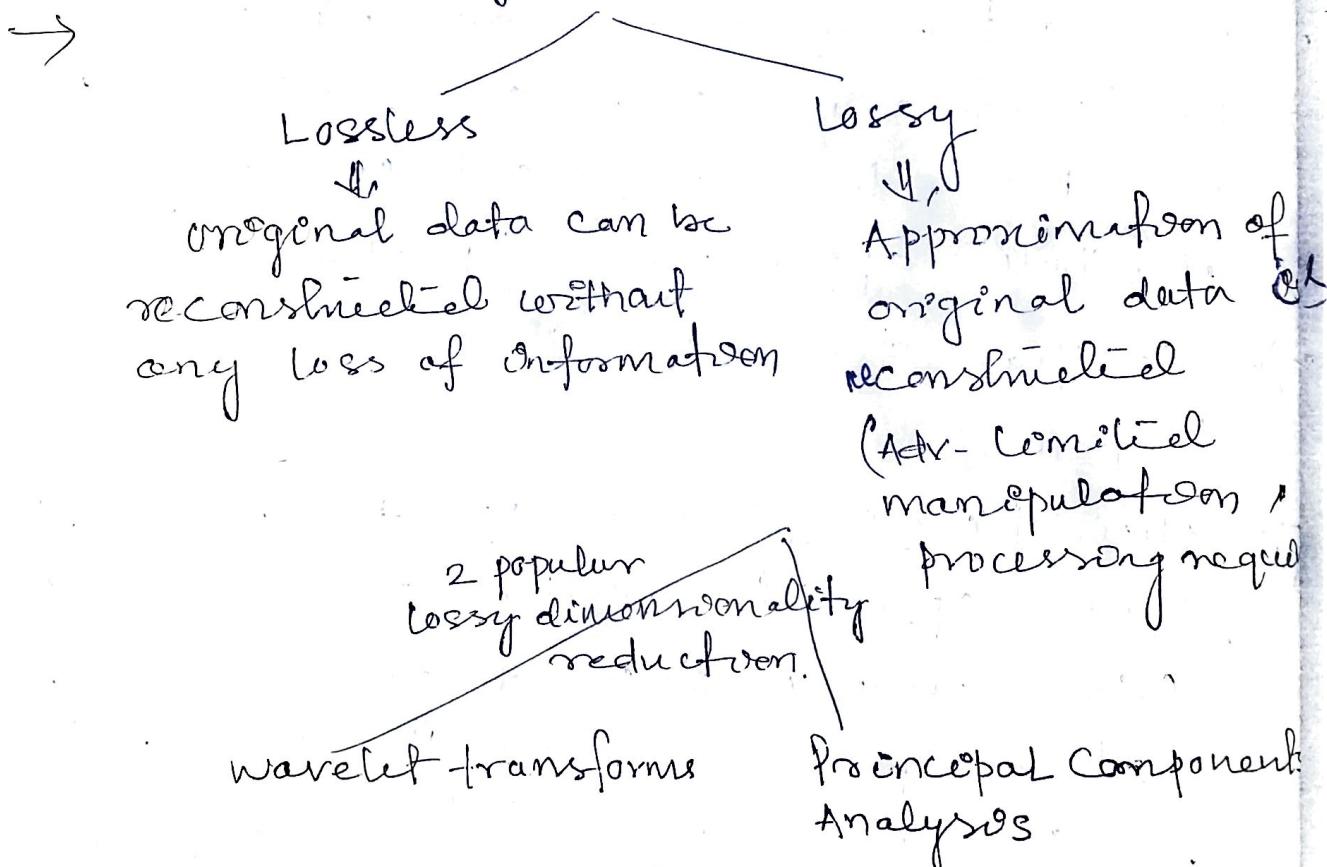
Stepwise backward elimination



- In Decision Tree induction, a tree is constructed from the given data.
- Set of attributes not appearing in tree are all irrelevant.
- Set of attributes appearing form the reduced subset of attributes.

Dimensionality Reduction :-

- Data are reduced/compressed by applying data encoding or transformation.



Numerosity Reduction :-

→ of reduces data volume by choosing smaller forms of data representation.

Techniques used

Parametric Method

- 1) it is assumed that data fits some model
- 2) The model estimates the data
- 3) only the data parameters are stored instead of actual data.

e.g.: Regression, log linear model

can be applied
to high dimensionality
data

Non-parametric method

- 1) no model is assumed.
- 2) it stores reduced representation of data in the form of either histograms or clustering or Sampling.

upto 10

Data Discretization :-

- If is a method of replacing numerous values of a continuous attribute by a small no. of interval labels.
- Thus it reduces & simplifies the original data.

Types / Category

Supervised

vs.
Un-Supervised
Discretization

Top-down (splitting)

vs.

Bottom-up (merging)
discretization

If the discretization process uses class information, it is called Supervised discretization otherwise un-supervised.

i.e if processes proceeds from top to down.

↓
of the process starts by first finding one or few interval to split the entire attribute range then repeats this recursively on the resulting interval it is called splitting

Topical Methods -

1) Partition -
Top-Down Split, unsupervised

2) Histogram Analysis -
Top Down split, unsupervised

3) Clustering Analysis -
Either top down or bottom up merge split
unsupervised.