

Data Mining Notes

Contributor: Deepanshu

[KMV CS(H)]

Computer Science Notes

Download **FREE** Computer Science Notes, Programs, Projects, Books for any university student of BCA, MCA, B.Sc, M.Sc, B.Tech CSE, M.Tech at
<https://www.tutorialsduniya.com>

Please Share these Notes with your Friends as well

facebook



1. Introduction to Data Mining, Pang-Ning Tan, Michael Steinbach, Vipin Kumar, Pearson

Data are raw facts which are not processed and mining when we try to get useful info. So in data mining we process our data and try to get maximum useful information from it. Data Mining is a technology that blends traditional data analysis methods with sophisticated algorithms for processing large volumes of data. It's the process of discovering useful information all useful patterns in a large data repositories that might otherwise remain unknown.

* APPLICATIONS :-

- 1) Customer Profiling :- Customer service records are collected from Call Centres, Web logs from e-Commerce Websites etc.
- 2) Targeted Marketing :- To find a suitable group of customers for their product and does marketing selectively.
- 3) Store layout
- 4) Fraud Detection :- To track an unusual transaction in an account.
- 5) Medicines :- It's used to address biological challenges such as protein structure prediction, modelling of biomedical pathways etc.

* INFORMATION RETRIEVAL

→ is looking up for individual records using a DBMS or finding particular Webpages using a query.

Data mining techniques are used to enhance information retrieval systems.

* Types of Data :- (according to structure)

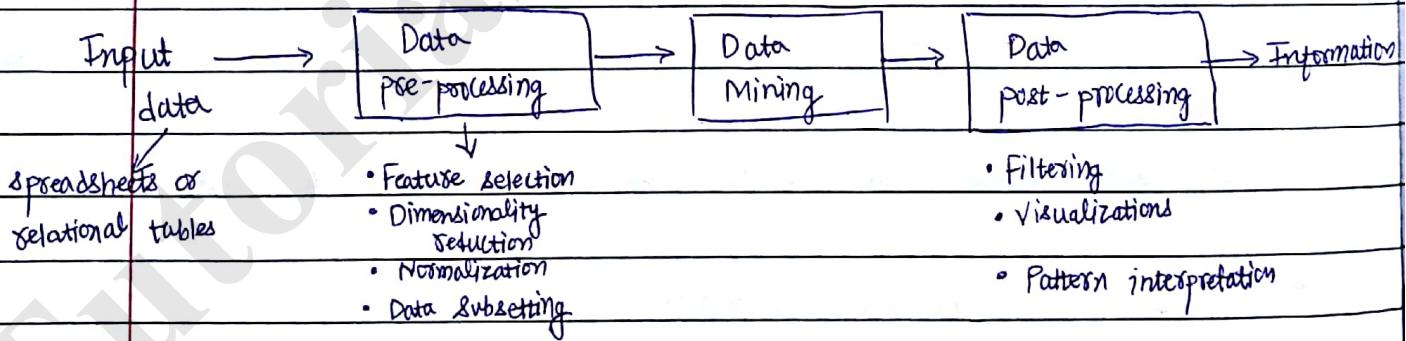
- 1) Structured :- Data that has a defined datatype or format. e.g - Word or excel files.
- 2) Semi-structured :- Textual data with a noticeable pattern. e.g - XML files.
- 3) Quasi-structured :- Textual data which requires effort, tools and time for formatting.
e.g - Web server logs.
- 4) Unstructured :- This data has no structure and is stored as different types of file.
e.g - pdf.

* The mixed information must be :-

- 1) New:- The extracted information should give us new patterns or relationships among data entity. e.g - if we have a dataset regarding the people who smoke and if we get the result smoking causes Cancer it's not something new. As it's not new so there is no purpose investing so much time, money and efforts.
- 2) Useful:- The mixed info from the data should be useful and relevant to us. for e.g - in the smoking dataset if we concludes about the longevity of the person who smoked and who doesn't, it would prove to be more useful.
- 3) Correct:- The mixed needs to be evaluated for its correctness before we use it for any other purpose. for e.g - in the dataset of stock prices of all companies we predict Company $X \rightarrow 20\%$ but after sometime it actually got decreased by 10% so it's always necessary to validate information before reaching to any conclusion.

* Data Mining and Knowledge Discovery

→ Knowledge discovery is the process of converting raw data into useful information which consists of transformation steps from data preprocessing to post processing of data mining results.



* Data Preprocessing

Feature selection:- is also known as attribute or variable selection is the process of selecting a subset of relevant features for use in model construction.

Dimensionality reduction:- includes feature selection and feature extraction.

Normalization:- to reduce data redundancy and improve data integrity.

Data Subsetting:- process of retrieving parts of large files.

- Post Processing :-

Filtering patterns :- It's a design pattern that enables the developers to filter a set of objects using different criteria.
Visualizations -

Pattern interpretation :- evaluation and interpretation of pattern to decide what qualifies as knowledge.

Post processing ensures that only valid and useful results are incorporated into decision support systems.

Data Warehousing - Data from multiple sources and in multiple formats is stored together at one single location is called data warehousing:

Data Mining - Intelligent operations such as clustering, classification, association, regression are applied to extract patterns. It's the most important step in knowledge discovery.

R

is an open source statistical environment modelled by Robert Gentleman and Ross A. Ihaka of statistics department of University of Auckland in 1995. It's currently maintained by the R core development team consisting of international team of volunteer developers.

To see what packages are currently available → `installed.packages()`.

You can install many more packages using :-

Step1 - using the command `install.packages (" ")`

Step2 - loading using `library (<packagename>)`

Once you have some extra packages installed they're not automatically ready for use and we must load them to make library of core routines available for use. The command to see which packages are currently loaded is `search()`.

→ If we want to remove or unload a package we command `detach(package: packagename)`

* R used as calculator

* Storing results of calculation

`ans1 <- 3+9-2+3`

`ans1`

[1] 13

Using the combined Command for making data
data1 = c(1, 2, 3, 4).

data1

[1] 1 2 3 4

Entering text data :-

data2 = c("a", "b", "c", "d")

data2

[1] "a" "b" "c" "d"

data2 = c(1, 2, 3, "a")

data2

[1] "1" "2" "3" "a"

* Using scan Command for making the data

When using the c Command typing all those command to separate the values can be a bit tedious so we can use another command scan() to do a similar job.

data5 = scan()

1: 2 3 4 5 ↵

2: 6 7 8 ↵

3: ↵

Read 7 items

> data5

[1] 2 3 4 5 6 7 8

> day2 = scan(what = 'character')

1: "a" "b" ↵

2: "c" ↵

3: ↵

Read 3 items

> day2

[1] "a" "b" "c"

data6 = scan(sep = ',')

1: 1, 2, 3, 4 ↵

2: ↵

Read 4 items

* Challenges faced by Traditional data analysis technique that motivated the development of data mining :-

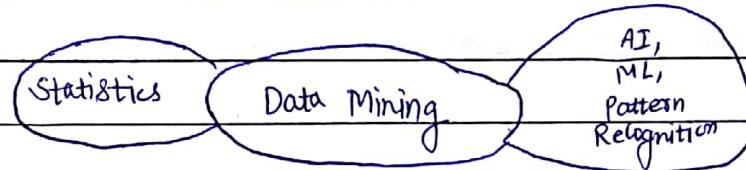
- 1) scalability:- As datasets with large sizes are becoming common, data mining algorithms are used which employs special strategies to handle exponential search problems also out of the 4 algorithms become necessary when processing a dataset that can't fit into main memory. e.g - climate data, hospital data. etc.
- 2) High dimensionality:- Now-a-days datasets have 100's and 1000's of attributes instead of a few handful attributes. for e.g - in bioinformatics the gene expression data involves 1000's of features, climate data contains measurements of temperature at various locations for an extended period of time so traditional data analysis technique don't work well for such high dimensional data.
- 3) Heterogeneous and Complex data:- Traditional methods deal with the datasets containing attributes of same types but recent years have seen emergence of more complex data objects. for e.g - DNA data with sequential and 3D structure, climate data which consists of time series measurements.
- 4) Data ownership and distribution:- Sometimes the data needed for analysis is geographically distributed among the resources belonging to multiple entities which requires the development of distributed data mining techniques.

* What is Supervised and Unsupervised learning ?

Unsupervised:- These are no predefined labels we just have the observations and the system clusters them based on their similarity this type of learning is called unsupervised. for e.g - there are mixed observations consisting of different animals and birds. so based on their similarities the system would be able to cluster them into 2 and only after clustering we get to know that one of the clusters is of the animals and second is of birds.

Supervised:- In this we have predefined labels and we have to classify our record in one of the predefined labels.

* Origins of Data Mining :-

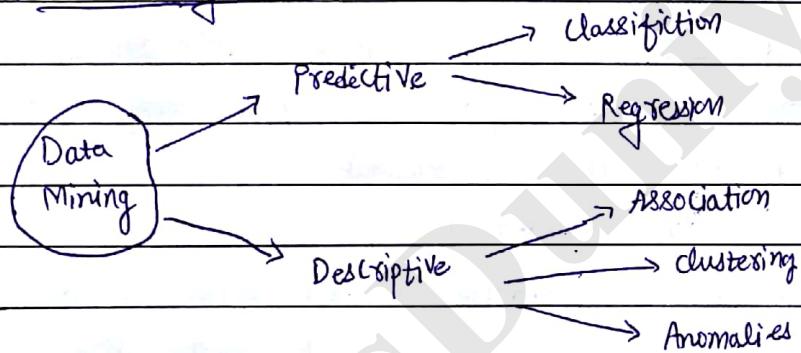


Database Sys, Parallel Comp, Distributed Computing

Data mining is made up of :-

- 1) Sampling, estimation and hypothesis testing from statistics.
- 2) Search algorithms, modelling techniques and learning theories from AI.
- 3) Database system, parallel computing and distributed computing are needed to provide support for efficient storage, indexing and query processing.

* Data Mining Tasks



R Programming

ui repository

Data folder → iris.csv

setwd("C:/Users/Desktop/Data") → set working directory

getwd() → get working director

read.csv("iris.csv")

read.csv("iris.csv", sep=',', header= TRUE / FALSE)

data1 = scan(file="abc.txt")

data1

dir("c1--")

data2 = scan(file.choose())

text file can be read

data2 = scan(file.choose(filename), sep=',', what='character')

or
'char'

data3 = read.table(file.choose(), sep = ',', header = FALSE)
↓
in this the header is FALSE by default.
for tab sep = '/t'

Most Class 1	Class 2	data	class
12	8	12	class1
15	9	15	class1
17	7	17	class1
11	9	11	class1
15	NA	15	class1
		8	class2
		9	class2
		7	class2
		9	class2

Preprocessing

* Data Mining Tasks Continued

* Predictive Tasks

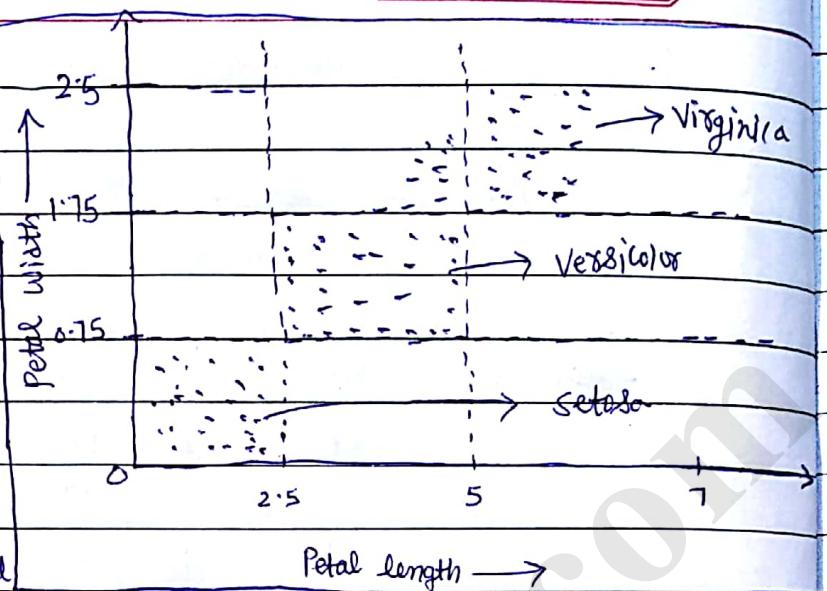
To predict the value of a particular attribute based on the value of other attributes. The attribute to be predicted is known as target or dependent variable while the attribute used for making the prediction are known as explanatory or independent variables.

Predictive modelling refers to the task of building a model for target variable as a function of independent variables.

1) Classification:- It's used for discrete target variables. e.g - iris data set is having 150 data items and 4 attributes Sepal length, Sepal width and petal length, petal width. The 3 species are Setosa, Versicolor and Virginica. Petal width has values as low, medium, high same as petal length. Based on petal length and width following rules can be derived:-

Petal length	Petal width	
Low [0, 2.5)	Low [0, 0.75)	→ Setosa
Med [2.5, 5)	Med [0.75, 1.75)	→ Versicolor
High [5, ∞)	High [1.75, ∞)	→ virginica

These rules help in classifying most of the flowers. Setosa species are well separated from versicolor and virginica but the latter two somewhat overlap with each other w.r.t. these attributes.



The goal is to create a model

that minimizes the error between true and predicted values of target variable.

- 2) Regression:- is used for continuous target variable. for e.g - forecasting the future price of a stock.

* DESCRIPTIVE TASKS

Here the objective is to derive patterns that summarize the underlined relationships in data.

- 1) Association Analysis:- It's used to discover patterns that describe strongly associated features in an efficient manner. The discovered patterns are represented in the form of implication rules and feature subsets.

↳ Application :-

① Identifying Webpages that are accessed together e.g - for buying an electronic item the lowest price of the product is displayed at all possible e-commerce sites, when we search for a research article on a particular topic all the related articles are shown one after other in search results.

② Finding groups of genes that have related functionality. It's used by scientists to improve the features of mutated plants.

③ Market basket analysis:- grocery stores apply association analysis to find the items that are frequently bought together by customers.

e.g - {Milk, butter} → {bread} whenever milk and butter is bought bread is also bought.

This type of rule can be used to identify cross selling opportunities among

selected items.

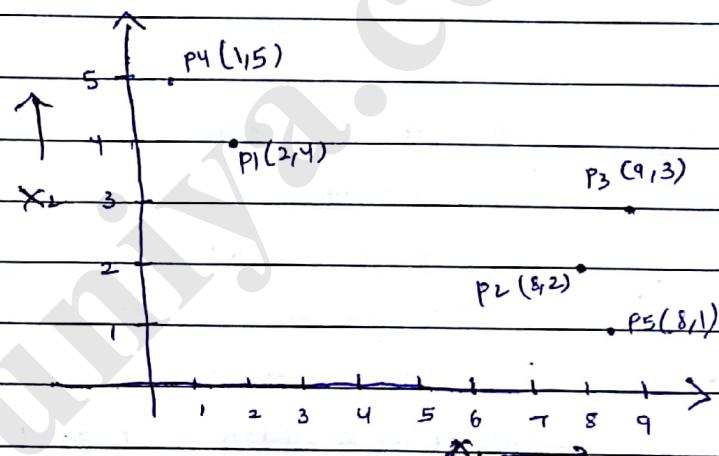
2) e1

2) Cluster Analysis:- To find the groups of closely related information, observations that belong to the same cluster are more similar to each other than the observations that belong to other clusters.

e.g.- it's used in compression of data, to find groups of related customers, collection of news articles can be grouped on the basis of repetitive topics.

Daily expenditure on food (x_1) and clothing (x_2) of 5 persons are shown:-

x_1	x_2	Person
2	4	P1
8	2	P2
9	3	P3
1	5	P4
8.5	1	P5



P1 and P4 in one cluster

P2, P3 and P5 in another.

3) Anomaly detection:- To identify the observations called anomalies or outliers whose characteristics are significantly different from the rest of the data. Good anomaly detection algorithms must have a high detection rate and low false alarm rate (to call an observation an anomaly just below it's little different from rest of the observations).

Applications:-

a) detection of fraud:- a credit company builds a profile of every user by recording the legitimate transactions made by them along with some information like credit limit, name, income, transaction history. When a new transaction arrives it's compared against the profile of the user if the characteristics of the transaction are very different from the payment history (transactions) the transaction is flagged off as potentially fraudulent.

b) Network intrusion:- identifying strange patterns in network traffic that could signal hack.

c) System health monitoring

Prior to data processing we must ensure that either the outliers are removed or some adjustments are made before analysis of data.

R programming

`ls()` → lists all the name, variable or items that are available.

* Searching the patterns

`ls(pattern = 'b')` → it will display all variable names containing b

`ls(pattern = '^b')` → " " " " " starting with b

`ls(pattern = 'b$')` → " " " " " ending with b.

`ls(pattern = '^ [ba]')` → " " " " " named starting either with b or a.

`ls(pattern = 'b . c')` → ----- of form bac, bac1 etc. with 1 letter between b and c.

* Vectors as told

* DataFrames

`dim(variablename)` lists no. of columns and rows in the data file.

`colnames(Variable) = c (---)`

or
`names(data) = c ("SepalL", "SepalW", "PetalL", "PetalW")`] → names the column

`colnames(data)` or `names(data)` to display name of column

* Useful Summary of Dataset

`summary(data)`, `str(data)`, `glimpse(data)`

↳ for using it install dplyr package.

`head(data)` → first 6 rows

`tail(data)` → last 6 rows

`head(data, 15)` → first 15 rows

> `hist(data[, "petal-len"])`

> `plot(data[, "petal-len"], data[, "petal-width"])`

To select first two attributes of data

`data2 = data[c(1,2)]`

and to select only 50 rows

`data2 = head(data2, 50)`

DATA

* Types of Data:- (Qualitative / Quantitative)

These are different types of data sets for e.g - the attributes used to describe the data objects can be either qualitative or quantitative.

Colors, grades, \leftarrow
& kill & marks etc. \rightarrow age, marks, price, salary etc.

Types of data also determines the tools and techniques that can be used to analyze the data.

* Quality of Data

Data quality issues:-

- 1) Presence of noise or outliers
- 2) missing, inconsistent or duplicate data
- 3) data that's biased

While most DM techniques can tolerate some level of imperfection, in data improving data quality improves the quality of resulting analysis.

* Processing steps to make the data more suitable for DM

Raw data must be processed:-

- 1) to improve data quality.
- 2) so that it fits better in a specified DM technique or tool.

e.g - length attribute may be transformed into an attribute with discrete values. e.g - short, medium or long.

No. of attributes in a dataset is reduced so that techniques applied can be more effective.

* Analyzing the data in terms of its relationships

→ Data analysis can be done by finding the relationship among the data objects and then performing the remaining analysis using these relationships. for e.g - we can compute the similarity or distance b/w pair of objects and then perform analysis (clustering, classification or anomaly detection) based on these similarities or distances.

* TYPES OF DATA

A dataset is a collection of data objects (records, event, case, sample, observation, entity).

Attributes are basic characteristics of an object (also known as field, feature, dimension).

for e.g - in dataset of students each row corresponds to a student and each column describes some aspects of a student such as student id, name, address, gpa & score etc.

Generally symbolically attributes have small number of possible values while numeric attributes have potentially unlimited no. of values. e.g. - employee data set has attribute employeeage and employeeid. The age attribute is very much related to employee while the id attribute is only used to ensure that they are distinct employees so knowing the type of attribute is important in understanding whether the values of the attributes are consistent while underlined properties of the attributes.

Lab

```
iris_1 = data
```

```
iris_1$class <- NULL
```

removes
creates column class from ~~data~~. iris_1.

→ fix(iris_1) shows the changes made in iris_1

Check for missing values:-

complete.cases(iris_1) shows which row is complete as TRUE or which is not as FALSE.
length(which(complete.cases(iris_1))) shows no. of complete rows without NA.
which(complete.cases(iris_1)) row no. of complete rows.
which(!complete.cases(iris_1)) row no. of incomplete rows (NA rows).

Store rows with no NA:-

```
com <- which(complete.cases(iris_1))
```

```
data3 <- iris_1[com, ]
```

or

```
data4 <- na.omit(iris_1)
```

```
> data4
```

```
data5 = cut(data1$sepal_len, 3)
```

divides into 3 levels

```
q1 = quantile(data1$petal_len, (0:3)/3)
```

divides into 3 quantiles

* Converting one variable type to another

```
vec <- c(1, 2, 3, 4, 5, 6)
```

class(vec) → "numeric" [tells data type of variable]

```
as.character(vec) → converts vec into character
```

vec1 <- c("f", 2)

class(vec1) \Leftrightarrow "character"

vec2 <- c(TRUE, 2)

Vec2 [1] 1 2

class(vec2) "numeric"

vec2 <- c(TRUE, "f")

Vec2 - "TRUE" "f"

class(vec2) "character"

Function to remove duplicate rows from data.

distinct(iris) // removes duplicate rows

Read about \rightarrow dplyr, ggplot2, tidyR, tm, stringR.

OUTLIER DETECTION (boxplots)

boxplot(data) \rightarrow used to detect outliers and removing them
 \hookrightarrow part of prepossessing data

boxplot(data, horizontal = T) \rightarrow for horizontal boxplot.

A boxplot is a graphical representation of statistical data based on minimum, first quartile, median, third quartile and maximum.

The term boxplot comes from the fact that the graph looks like a rectangle with lines extending from the top and bottom. because of the extending lines this type of graph is sometimes called box-and-whisker plot.

The relative vertical spacing between the labels reflect the values of the variable in proportion. A boxplot can be placed on a coordinate plane resembling the Cartesian system so that the 5 values arranged vertically one above the other run parallel to y axis.

Program

SP

$\tau_1 = \text{which}(!(\text{ok} = \text{Age} \geq 22 \& \text{Age} \leq 150))$

τ_1

$\text{length}(\tau_1)$

TutorialsDuniya.com

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

Please Share these Notes with your Friends as well

facebook

WhatsApp 

twitter 

Telegram 

Properties of numeric attributes

- (1) Distinctness
- (2) Order
- (3) Addition / Subtraction
- (4) Multiplication / Division

on the basis of these properties we have 4 types of attributes

- (1) Nominal
- (2) Ordinal
- (3) Interval
- (4) Ratio

(1) → The value of nominal attribute provide enough information to distinguish one object from another. E.g.: - EmpId, Pincode. It requires one to one mapping of values in case of transformation. If all employee id are rearranged it will not make any difference.

(i) Mode, entropy, χ^2 test, etc.

↳ Attribute with more information gain increases predictability which reduces entropy. χ^2 tests (chi-square-test) → is used to access the goodness of fit between a set of observed values and those expected theoretically.

(2) Ordinal attribute :- We can order the objects on the basis of information provided by the attributes. for e.g. - length attribute will have the values small, medium and large, grade attribute, hardness of minerals good, better, best. While transforming their must be order preserving change of values.

$$\text{new_value} = f(\text{old_value})$$

operations that it supports → median, percentile, rank, co-relation etc.

(3) Interval type of attributes

There exists a unit of measurement, the differences between the values are meaningful.

$$T(^{\circ}\text{F}) = \frac{9}{5} T(^{\circ}\text{C}) + 32$$

- For transforming a new value the following conversion formula would be used
new value = $a * \text{old-value} + b$
where a, b are constants.
- c and k are two different units of measurement which differ in their zero value and unit size.
- operations that these attributes support:- mean, standard deviation etc.

(A) Ratios

this type of attributes deal with numeric data. Here the transformation function used will be

$$\text{new-value} = a + \text{old-value}$$

where a is constant.

e.g - count, name, length, salary etc.

Here the attributes don't differ in their zero value. for eg - the statistical operation when performed on length attribute will give the same result whether length is in meter or foot.

operations:- geometric mean, percent variation etc.

- Nominal and ordinal attributes are collectively referred to as categorical or qualitative attributes. These attributes though represented by numbers must be treated like symbols.
 - They lack most of the properties of numbers.
- Interval and ratio attributes are collectively referred to as quantitative or numeric attributes. They can be integer values or continuous.

Important properties of these attributes

- Each attribute type possesses all the properties and operations of the attribute types about it. i.e. any operation that is valid for nominal, ordinal and interval attributes is also valid for ratio attribute.

Describing attributes on the basis of number of Values

We have two types of attributes:-

- Discrete:- It has a finite set of values. For e.g - categorical attributes such as zip code, id no, count attributes. They generally have 2 values e.g - TRUE or FALSE.

YES or NO, 0 or 1 etc. They are represented either by using boolean or integer variables.

- 2) Continuous:- Its values are real numbers and are represented by floating point variables. E.g - temperature, height, weight etc.

* Asymmetric Attributes :-

e.g - Consider a dataset with student records where each attribute records whether or not a particular student took a certain subject at the university. 1 value will signify that student has taken that particular subject and 0 specify that student hasn't taken subject. As the students will take small no. of all available subjects most of the values in dataset would be 0 so it's more imp to focus only on non-zero values. these types of attr. are called asym. binary attributes and are particularly imp. for association analysis.

* Types of DataSet

General characteristics of datasets :-

- (1) Dimensionality :- is the no. of attributes that the objects in a dataset possess. Data with lesser no. of dimensions tends to be qualitatively better than the moderate or high dimensional data. The difficulties associated with analyzing high dimensional data is referred to as "Curse of dimensionality".
- (2) Sparcity :- In some datasets in asymmetric features most attributes of the object have 0 values → practically it helps in saving computation time and storage bcz only non 0 values need to be stored and manipulated.
- (3) Resolution :- The properties of data are different at different resolutions. e.g - surface of earth seems very uneven at a resolution of a few metres but is relatively smooth at the resolution of a few km's if the resolution is too fine a pattern may not be visible or may be buried in noise. if resolution is too coarse the pattern may disappear.

* Types of Datasets (Generally Asked)

- (1) Record Data :- Here dataset is collection of records each of which consists of a fixed set of data themes. It has further subtypes:-
 - a) Transaction / Market Basket Data → columns are items and rows are transactions
 - b) Data Matrix
 - c) Sparse Data Matrix → e.g -

- a) Transaction data - is a special type of record data where each record involves a set of items purchased by a customer. This type is also called market basket data. Attributes may be binary indicating whether or not an item was purchased or can be discrete or continuous such as the no. of items purchased or the amount spent on those items.
- b) Data Matrix:- In it the data objects is a collection of data having the same fixed set of numeric attributes a set of such objects can be interpreted as an $M \times N$ matrix where there are M rows and N columns. The data matrix is the standard data format for most statistical data.
- c) Sparse data matrix:- Special case of data matrix in which the attributes are of some type and are asymmetric i.e. only non-zero values are important. for e.g. document data where document can be represented as a vector and each word as the attribute. the value of each component is no. of items the corresponding word occurs in the document. It's also called document term matrix.

② GRAPH BASED DATA

It is convenient and powerful representation of data. It has two specific types:

- a) the graph captures the relationships among the data objects
- b) the data objects themselves are represented as graphs.

→ The relationships among objects convey important info. Where the data objects are mapped to nodes on the graph while the relationships among the objects are captured by the links b/w the objects and link properties. for e.g.- Webpages on WWW contain both text and links to other pages

→ Data with objects as graphs:- If objects have structure i.e. the objects contain sub objects that have relationships then such objects are frequently represented as graphs for e.g. the structure of chemical compounds can be represented by graphs where the nodes are atoms and the link b/w the nodes

Axle chemical bond.



A graph representation makes it possible to determine which substructures occur frequently in a set of compounds and to ascertain whether the presence of any of these substructures is associated with the presence or absence of certain chemical properties.

Substructure Mining \Rightarrow is a branch of data mining that analyses such data.

③ ORDERED DATA

For some types of data the attributes have relationships that involve order in time or space

- a) Sequential data:- It's also referred as temporal data where each record has a time associated with it. for e.g- in a retail transaction data the time of the transaction is also stored. This time information makes it possible to find the patterns regarding which products sale peaks during a particular season or festival.
- b) Sequence data:- consists of a dataset that is sequence of individual entities such as sequence of words or letters there are no timestamps involved. e.g- the genetic info. of plants and animals can be represented in the form of genes.
- c) Time series data:- it's a special type of sequential data in which each record is a time series i.e. a series of measurement taken over time for e.g- (1) financial dataset that contains the time series of the daily prices of various stocks. (2) time series of the average monthly temperature of a particular city during the years 2006 - 2010 . In temporal data there's temporal autocorrelation i.e. if two measurements are close in time then the values of those measurements are often very similar.
- d) Spatial Data:- It's a data collected at different geographical locations for e.g- Weather data In spatial data there's spatial autocorrelation i.e. objects that are physically close tend to be similar e.g- 2 points on the earth that are geographically close to each other have similar values for temperature and rainfall.

* HANDLING NON-RECORD DATA

Record Oriented Techniques can be applied to non-record data by extracting features from the data objects and using these features to create a record corresponding to each object. For e.g. given a set of common substructures each compound can be represented as a record with binary attributes that indicate whether a compound contains a specific substructure.

* DATA QUALITY

Data mining focuses on :- 1) the detection & correction of data quality problems also called data cleaning

2) the use of algorithms that can tolerate poor data quality

Measurement and data collection issues

There may be problems due to human error, limitations of measuring devices or flaws in the data collection process, values or even entire data objects may be missing (erroneous). There may be spurious or duplicate objects, all the data may be present but there may be inconsistencies. For e.g. a person with height 2m has weight of only 2 kg.

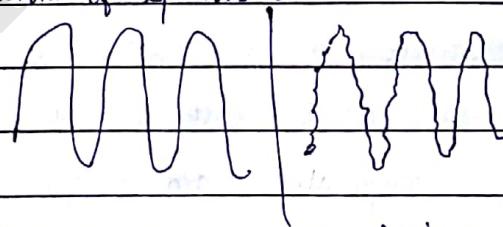
1) Measurement Error :- It occurs when the value recorded in the measurement process differs from the true value to some extent.

2) Data collection Error :- refers to errors such as omitting the data objects or attribute values or inappropriately including a data object.

3)

Noise and Artifacts

Noise is the random component of a measurement error involving distortion of a value or addition of spurious objects.



The elimination of noise is difficult in DM so our main focus is on making robust algorithms that produce acceptable results even when the noise is present. Data errors which are repeated for every record such as a ^{strange} mark in the same place on a set of photos such distortions of the data are often referred to

(a) artifacts.

→ Precision, Bias and accuracy

In statistics & experimental science the quality of measurement process is measured by precision and bias.

Precision:- The closeness of repeated measurements to one another. Precision is often measured by the standard deviation of a set of values.

Bias:- A systematic variation of measurements from the quantity being measured. It is measured by taking difference b/w the mean of a set of values & known value of the quantity being measured.

Accuracy:- The closeness of measurement to the true value of the quantity being measured.

Without understanding the accuracy of data & results an analyst can commit serious data analysis blunders.

$$\text{For Standard Deviation } \sigma = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}$$

e.g.- Consider a sample of IQ scores 96, 104, 126, 134, 140

$$\begin{array}{r} 126 + -30 -22 +0 +8 +14 \\ \hline 5 \end{array}$$

$$\begin{array}{r} 122 \\ -96 \\ \hline 6 \end{array}$$

$$126 + -16 = 120$$

$$\begin{array}{r} (26)^2 + (18)^2 + 4^2 + 12^2 + 18^2 \\ \hline 4 \end{array}$$

$$\begin{array}{r} (34)^2 + (16)^2 + 6^2 + 14^2 + 20^2 \\ \hline 4 \end{array} = 19.13$$

8

Find SD of first 10 natural no.

$$\begin{array}{r} (-4.5)^2 + (3.5)^2 + (2.5)^2 + (1.5)^2 + (0.5)^2 + (-0.5)^2 + (1.5)^2 + (2.5)^2 + (3.5)^2 \\ + (4.5)^2 \\ \hline 9 \end{array}$$

$$= 3.027$$

Suppose that we have a standard laboratory weight with mass of 1g and want to assess the precision and bias of the laboratory scale we weigh the mass 5

times and values 1.015, 0.990, 1.013, 1.001, 0.986

$$\bar{x} = 1.001$$

1.000

0.001 → bias

1.001

- 0.990
0.011

1.013

- 1.001

.12

- 0.013095 [Precision?]

→ * OUTLIERS

are the data objects that have the characteristics which are different

from most of the ^{other} data objects in the dataset. The value of an

attribute are unusual w.r.t. the typical values of that attribute.

1.001

- 0.986

.015

0.000/7/5

Outliers are legitimate data objects or values that may sometimes be of interest

for e.g. - in fraud and network intrusion detection the goal is to find unusual scores or events from among a large no. of normal ones.

→ * Missing VALUES

It is not unusual for an object to have one or more missing attribute values

1) Reasons:-

- 1) The info. was not collected. for e.g. - some people decline to give their age or weight.
- 2) Some attributes are not applicable for all the objects e.g. - in library there's a student record which has the attributes student name, id, mobile no., course, booksissued, returndate and fine. If a student hasn't issued any book booksissued, fine and returndate columns are invalid to him.

There are several strategies for dealing with missing values:-

- 1) Eliminate the data objects or attributes. A simple and effective strategy is to eliminate the objects with missing values. If partially specified data object contains some info. and if many objects have missing values then reliable analysis can be impossible. If we eliminate attributes that have missing values then the eliminated attributes may be the ones that are critical to analysis.

- 2) Estimate missing values - sometimes missing data can be reliably estimated for e.g. Consider a dataset that has many similar data points. the attribute values of the points closest to the point with a missing value are often used to estimate

the missing value. If the attribute is continuous average attribute value of the nearest neighbours is used if it's categorical then most commonly occurring attribute value can be taken.

- 3) Ignore the missing values during analysis - Many DM approaches can be modified to ignore the missing values. Suppose the objects are being clustered & the similarity b/w pairs of data objects needs to be calculated so here the approx similarity can be calculated by using only those attributes that don't have those missing values. This degree of inaccuracy may not matter much provided the total no. of attributes are large and the no. of missing values isn't very high.

(6) → Inconsistent Values

Data can contain inconsistent values. for e.g. Consider an address field where the pincode area specify is not present in the city. So maybe a digit was missed when the info was scanned so it's important to detect and if possible correct such problems. Some inconsistencies are easy to detect. for e.g. a person's height and age shouldn't be -ve. Sometimes there are discrepancies in a dataset but that doesn't mean the data shouldn't be used only the analyst should consider the potential impact of such discrepancies on the data mining analysis.

(7) → Duplicate Data

A dataset may include data objects that are duplicates. for e.g. many people receive duplicate mails bcz they appear multiple times in a DB with slightly different names. To detect and eliminate such duplicates inconsistent values must be resolved also care needs to be taken to avoid accidentally combining the data objects that are similar but not duplicates. Such as 2 distinct people with identical names. The term D-Duplication is often used to refer to the process of dealing with these issues.

Data Quality

→ * ISSUES RELATED TO APPLICATION OF COLLECTED DATA

(1) Timeliness:-

(2) Relevance

(3) Knowledge about data.

(1) Timeliness:- Some data starts to age as soon as it has been collected. for e.g - the data of purchasing behaviour of customers or web browsing patterns represents reality for only a limited time. If the data is out of date then so are the models and patterns based on it. for e.g - stock market trend for the next day.

(2) Relevance:- The available data must contain the info. necessary for application. for e.g - the task of building a model that predicts accident rate for drivers if the info about age and gender of the driver is absent or omitted then it's likely that the model will have limited accuracy. for e.g - a survey was held to see how many people would prefer to have fast food in lunch instead of home cooked food that survey didn't include the individuals from age group 15-25. So the survey won't serve the purpose as we have excluded an important part of our sample.

Another problem is "Sampling bias" which occurs when a sample doesn't contain different types of objects in proportion to their actual occurrence in the population. Sampling bias will result in an erroneous analysis because the results of data analysis will reflect only the data that's present.

(3) Knowledge about the data:- The datasets are accompanied by documentation that describes different aspects of data like precision, types of attributes etc. The quality of the documentation can help or hinder the subsequent analysis. for e.g - if the documentation fails to tell us that missing values for a particular field are indicated with 9999 then our analysis of the data may be faulty.

* DATA PREPROCESSING

1) Aggregation:- means combining of two or more objects into a single object. for e.g - consider a dataset consisting of transactions recording the daily sales of products in various store locations. One way to aggregate transactions for this dataset is to replace all the transactions of a single store with a single store wide transaction. this reduces the 100's and 1000's of transactions to a single daily transaction and the no. of data objects is reduced to the

No. of stores. & How an aggregate transaction is created i.e. how the values of each attribute are combined?

Ans → Quantitative attributes such as price are aggregated by taking sum or average, qualitative attributes such as items can be summarized as the set of all the items that were sold at that location. We can also reduce the possible values for dates from 365 days to 12 months. This type of aggregation is commonly used in OLAP (online Analytical Processing) systems. OLAP tools enable users to analyze multidimensional data from multiple perspectives.

Advantages of Aggregation :-

- a) The smaller datasets resulting from data reduction require less memory and processing time which permits the use of more expensive data mining algorithms.
- b) It provides a higher level view of data.
- c) The behaviour of group of objects or attributes is now more stable than that of individual objects or attributes. for e.g - the average yearly precipitation has less variability than the average monthly precipitation for any given particular locⁿ over ^{given} a ~~particular~~ period of time.

Disadvantages of aggregation is the potential loss of interesting details for e.g - Which day of the week had the highest sale for a particular product

2) Sampling

Sampling approach is used for selecting a subset of the data objects to be analyzed. A sample must be representative i.e. having the same property as original set of data. for e.g - mean of the sample data must be the same as the mean for ^{the} entire dataset. We must choose a sampling scheme that guaranteed a high probability of getting a representative sample.

Different Sampling approaches :-

- a) Simple Random Sampling
 - b) Stratified Sampling
 - c) Progressive Sampling
- ```
graph TD; A[Simple Random Sampling] --> B[With replacement]; A --> C[Without replacement]
```

## ① Simple Random Sampling

There is an equal probability of selecting any particular item.

It's of two types

- a) Sampling without replacement → as each item is selected it's removed from the set of all the objects that constitute the population.
- b) " With replacement → The objects are not removed from the population as they are selected for the sampling. In this method the same object can be picked more than once. In this the probability of selecting any object remains constant during the sampling process.

Problem in Simple Random Sampling :-

When the population consists of diff. types of objects with widely different no's of objects simple random sampling can fail to adequately represent those type of objects that are less frequent. This can cause problems when the analysis requires

proper representation of all object types. So in stratified sampling, equal no. of objects are drawn from each group even though the groups are of different sizes. A slight variation can be done in this method such that the no. of objects drawn from each group is proportional to the size of that group.

## ② Stratified Sampling

### a) How to choose the correct sample size?

Ans It's necessary to choose correct sample size larger ss increases the probability that a sample will be representative but also eliminate the advantage of sampling with smaller sample sizes patterns may be missed or erroneous patterns <sup>may</sup> not be detected.

To determine proper ss it requires a methodological approach. When we are

as follows:

Stratified Sampling Suppose we're given a set of data that consists of almost equal sized groups then we need to find atleast 1 representative point for each of the groups assuming the objects in each group are highly similar to each other.

### ③ Progressive Sampling

The proper SS can be difficult to determine so adaptive or progressive sampling schemes are used these approaches start with a small sample and then increase the SS until a sample of sufficient size has been obtained. This technique eliminates the need to determine the correct SS initially but there must be a way to evaluate the sample to judge if it's large enough.

In progressive sampling how will we come to know that we've got correct SS :-

$$2, 3, \textcircled{5}, 7, 9, \textcircled{4}, \textcircled{6}, 8, \textcircled{7}, 5, 4, 6 \\ \frac{66}{12} = 5.5$$

Mean = 5.5

5, 4, 6, 7 → Mean = 5.5 So it's a representative sample

If we're trying to find a pattern in which we've total 3000 records so first we take 500 records as our sample and get a very faded pattern so we take some more samples say 300 so new SS = 800 now we get somewhat clear pattern so again we need to check whether it's correct SS or we're not missing out on imp. part of our pattern so we further increase our SS by 200 now we're getting the same pattern as before but is more clear so we'll come to know that the SS = 800 was adequate to form a representative sample.

### 3) DIMENSIONALITY REDUCTION

Datasets can have a large no. of features for e.g. - in document dataset

The advantages :-

- Many DM algorithms work better if the dimensionality is lower below dimensionality "reduc" can eliminate irrelevant features and reduce noise.
- A reduction of dimensionality can lead to a more understandable model as it may allow the data to be more easily visualised
- The amount of time and memory required by the DM algorithm is reduced with reduction in dimensionality.

DR can be done in 2 ways:-

- Creating new attributes that are combination of old attributes. e.g. - the student marks in which we can combine all the marks column to 1 column (GPA).
- Selecting new attributes that are a subset of old is known as features

# **TutorialsDuniya.com**

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

**Please Share these Notes with your Friends as well**

**facebook**

**WhatsApp** 

**twitter** 

**Telegram** 

Subset selection or feature selection.

Curse of Dimensionality :- (VVImp)

As dimensionality increases the data becomes increasingly sparse in the space that it occupies as a result many clustering and classification algorithms have reduced classification accuracy and poor quality clusters.

#### 4) FEATURED SUBSET SELECTION

This is generally useful if redundant and irrelevant features are present.

Redundant features duplicate much or all of the info. Contained in one or more attribute for e.g - age and date of birth attributes.

Irrelevant features contain almost no useful information for the DM tasks. for e.g - student id's are irrelevant to the task of predicting student's marks.

The ideal approach to featured selection is to try all possible subsets of features as input to DM algorithms and then take the subset that produces the best result. Since the no. of subsets involving  $n$  attributes is  $2^n$  such an approach is impractical in most situations. There are 3 approaches to featured selection:-

- a) Embedded Approach :- When the DM algorithm itself decides which attributes to use and which to ignore. e.g - Decision Tree classifiers
- b) Filter Approach :- Features are selected before the DM algorithm is run using an approach that's independent of DM tasks. e.g - Selection of the set of attributes whose pairwise correlation is as low as possible.
- c) Wrapper Approach :- These methods use the target DM algorithm to find the best subset of attributes without enumerating all possible subsets. We may identify and run algorithm only on some of the best possible subsets.

#### 5) FEATURE CREATION

It's possible to create a new set of attributes from the original attributes that captures the important information in a dataset much more efficient effectively. The no. of new attributes can be smaller than that of original attributes. One of the methods for creating new attribute is FEATURE EXTRACTION. The creation

If new set of features from the original raw data is known as feature EXTRACTION. It's application is found in image processing where each photograph is to be classified according to whether or not it contains a human face.

### DISCRETIZATION

#### b) DIGITIZATION AND BINARIZATION

| Categ. Value | Integer value | $x_1$ | $x_2$ | $x_3$ |              | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|--------------|---------------|-------|-------|-------|--------------|-------|-------|-------|-------|-------|
| awful        | 0             | 0     | 0     | 0     |              | 1     | 0     | 0     | 0     | 0     |
| poor         | 1             | 0     | 0     | 1     | Binarization | 0     | 1     | 0     | 0     | 0     |
| OK           | 2             | 0     | 1     | 0     |              | 0     | 0     | 1     | 0     | 0     |
| good         | 3             | 0     | 1     | 1     |              | 0     | 0     | 0     | 1     | 0     |
| great        | 4             | 1     | 0     | 0     |              | 0     | 0     | 0     | 0     | 1     |

$(0, 110) \xrightarrow{10^3 \text{ split}} (0-10) (10-20) \dots (100-110)$

A      B      G

Discretization

Some DM algorithms especially classification algorithms requires the data to be in the form of categorical attributes. algorithms that find association patterns require that the data be in the form of binary attributes so it's often necessary to transform a continuous attribute into a categorical attribute called DISCRETIZATION and transforming both continuous and discrete attributes into one or more binary attributes called BINARIZATION.

BINARIZATION → A simple technique to binarize a categorical attribute is if there are m categorical values then uniquely assign each original value to an integer in the interval  $[0, m-1]$  if the attribute is ordinal then order must be maintained by the assignment.

In 1st table the transformation can cause complications such as creating unintended relationships among the transformed attributes. Here the attributes  $x_2$  and  $x_3$  are irrelevantly correlated. for association problems it's therefore necessary to introduce one binary attribute for each categorical value.

DISCRETIZATION OF CONTINUOUS ATTR. — The transformation of cont. attrs. to categorical attribute involves 2 subtasks :-

- a) Deciding how many categories to have
- b) how to map the values of cont. attr. to these categories. In the first step the values are divided into n intervals

by specifying  $n-1$  split points. In the 2nd step all the values in one interval are mapped to the same categorical value.

### 7) Variable Transformation

A variable transformation refers to a transformation that is applied to all the values of a variable.

In variable transformation a simple mathematical function is applied to each value individually.

If  $x$  is a variable then examples of transformation include :-

$$x^k, \log(x), e^x, \sqrt{x}, \frac{1}{x} \text{ etc.}$$

Variable transformations especially  $(\sqrt{x}, \log(x), \frac{1}{x})$  are often used to transform

the data that doesn't have the Gaussian or Normal distribution into the data that does.

- Speciality in Normal distribution :- [How normal distribution is different from others]
- In normal distribution the data tends to be around a central value with no bias left or right. This distribution has bell shaped density curve described by the mean and standard deviation. The density curve is symmetrical about the mean.

e.g. of Variable transformation -

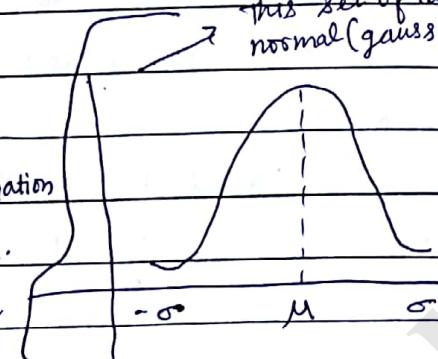
If the no. of data bytes that are transferred in a session ranges from 1 to  $10^{10}$  billion then it may be advantageous to compress by using  $\log_{10}$  transformation

$$\begin{array}{cc} I & II \\ 10^8 & 10^9 \\ \downarrow & \downarrow \\ 10^0 & 10^3 \end{array} \quad |8-9| = 1 \quad \text{So } I \text{ & } II \text{ are more similar in terms of no. of bytes transferred per session below they represent}$$

transf. of larger files.

Variable transformn should be applied with caution since they can change the ratio of data - for e.g. the transformation  $\frac{1}{x}$  reduces the magnitude of values

this set of data is in normal (gaussian) form



that are 1 or larger but increases the magnitude of values b/w 0 - 1.

- Normalization or Standardization → Their goal is to make an entire set of values have a particular property.
- e.g. Consider comparing people on basis of two variables age and income. for any 2 people the difference in income will likely be much higher than the difference in age. If the diff. in the range of values of age and income are not taken into account the comparison b/w the people will be dominated by the differences in income so transformation is often necessary to avoid having a variable with large values dominate the result of calculation.
- The mean and SD are strongly affected by outliers so mean generally replaced by median and SD by absolute SD.

### \* Dissimilarities b/w Data Objects

→ Distances :- The Euclidean distance b/w 2 points  $x \geq y$  in n-dimensional space is

$$d(x, y) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

Where  $n$  = no. of dimensions and  $x_k$  and  $y_k$  are the kth attributes of  $x \geq y$ .

The Euclidean distance is generalised by Minkowski distance matrix

$d(x, y) = \left( \sum_{k=1}^n (x_k - y_k)^\gamma \right)^{1/\gamma}$  Where  $\gamma$  is a parameter. Euclidean special case of Minkowski when  $\gamma = 2$

When  $\gamma = 1$  it's called Hamming (or Manhattan) Distance which is the no. of bits that are different b/w two objects that have only binary attributes

When  $\gamma = 2$  it's Euclidean distance also called  $L_2$  norm

When  $\gamma = \infty$  it's supremum distance  $d(x, y) = \lim_{\gamma \rightarrow \infty} \left( \sum_{k=1}^n |x_k - y_k|^\gamma \right)^{1/\gamma}$

$\hookrightarrow L_\infty$  or  $L_\infty$  distance

$$d(x, y) = \max_{k=1}^n |x_k - y_k|$$

The Euclidean, Manhattan and supremum distances are defined for all values of  $n$  and specify different ways of combining the differences in each dimension into an overall distance.

If  $d(x,y)$  is distance b/w 2 points  $x \neq y$  then following properties hold:-

- (1) Positivity a)  $d(x,y) \geq 0 \quad \forall x, y$   
 b)  $d(x,y) = 0 \quad \text{if } x = y$

- (2) Symmetry

$$d(x,y) = d(y,x) \quad \forall x, y$$

- (3) Triangle Inequality

$$d(x,z) \leq d(x,y) + d(y,z) \quad \forall x, y, z$$

Measures that satisfy all the three properties are known as METRIC.

| Point          | x Coordinate | y Coordinate | $d(P_1, P_2) =  0-2  +  2-0  = 4$ |
|----------------|--------------|--------------|-----------------------------------|
| P <sub>1</sub> | 0            | 2            | $d(P_1, P_2)$                     |
| P <sub>2</sub> | 2            | 0            |                                   |
| P <sub>3</sub> | 3            | 1            |                                   |
| P <sub>4</sub> | 5            | 1            |                                   |

L<sub>1</sub> norm matrix

|                | P <sub>1</sub> | P <sub>2</sub> | P <sub>3</sub> | P <sub>4</sub> |
|----------------|----------------|----------------|----------------|----------------|
| P <sub>1</sub> | 0              | 4              | 4              | 6              |
| P <sub>2</sub> | 4              | 0              | 2              | 4              |
| P <sub>3</sub> | 4              | 2              | 0              | 2              |
| P <sub>4</sub> | 6              | 4              | 2              | 0              |

$$\begin{array}{r} 5 \\ 5 \\ 26 \\ 25 \\ \hline 1 \end{array}$$

(L<sub>2</sub> norm) matrix

|                | P <sub>1</sub>          | P <sub>2</sub>         | P <sub>3</sub>         | P <sub>4</sub> |
|----------------|-------------------------|------------------------|------------------------|----------------|
| P <sub>1</sub> | 0                       | $\sqrt{8} = 2\sqrt{2}$ | $\sqrt{10}$            | $\sqrt{26}$    |
| P <sub>2</sub> | $\sqrt{8} = 2\sqrt{2}$  | 0                      | $\sqrt{2} = 1\sqrt{1}$ | $3\sqrt{1}$    |
| P <sub>3</sub> | $\sqrt{10} = 3\sqrt{1}$ | $\sqrt{2} = 1\sqrt{1}$ | 0                      | 2              |
| P <sub>4</sub> | $\sqrt{26}$             | $3\sqrt{1}$            | 2                      | 0              |

(L<sub>∞</sub> norm)

|                | P <sub>1</sub> | P <sub>2</sub> | P <sub>3</sub> | P <sub>4</sub> |
|----------------|----------------|----------------|----------------|----------------|
| P <sub>1</sub> | 0              | 2              | 3              | 5              |
| P <sub>2</sub> | 2              | 0              | 1              | 3              |
| P <sub>3</sub> | 3              | 1              | 0              | 2              |

B Calculate L, distance matrix, euclidean, supremum distance matrix for the following points

|    | x | y |
|----|---|---|
| P1 | 2 | 1 |
| P2 | 3 | 4 |
| P3 | 0 | 1 |
| P4 | 2 | 4 |

Hamming

$\gamma = 1$

|    | P1 | P2 | P3 | P4 |
|----|----|----|----|----|
| P1 | 0  | 4  | 2  | 3  |
| P2 | 4  | 0  | 6  | 1  |
| P3 | 2  | 6  | 0  | 5  |
| P4 | 3  | 1  | 5  | 0  |

$\frac{1 \cdot 4}{4 \ 2 \ 3}$

Euclidean

$\gamma = 2$

|    | P1                 | P2                 | P3                  | P4                  |
|----|--------------------|--------------------|---------------------|---------------------|
| P1 | $\sqrt{10} = 3.16$ | $\sqrt{10} = 3.16$ | 2                   | 3                   |
| P2 | $\sqrt{10} = 3.16$ | 0                  | $\sqrt{2} = 4.28$   | 1                   |
| P3 | 2                  | $3\sqrt{2} = 4.28$ | 0                   | $\sqrt{13} = 3.605$ |
| P4 | 3                  | 1                  | $\sqrt{13} = 3.605$ | 0                   |

Supremum

|    | P1 | P2 | P3 | P4 |
|----|----|----|----|----|
| P1 | 0  | 3  | 2  | 3  |
| P2 | 3  | 0  | 3  | 1  |
| P3 | 2  | 3  | 0  | 3  |
| P4 | 3  | 1  | 3  | 0  |

## NON MATRIX DISSIMILARITIES

Set differences

$$A = \{1, 2, 3, 4\}$$

$$B = \{2, 3, 4\}$$

$$A - B = \{1\}$$

$$B - A = \{2\} = \emptyset$$

$$d(A - B) \neq d(B - A)$$

$$d(A, B) \neq d(B, A)$$

so  $A \approx B$  are not similar.

Test → 18th Feb

$$d(A, B) = d(A-B) + d(B-A)$$

$$d(B, A) = d(A, B) \quad \text{now symmetry properties will hold}$$

→ If mark question to check whether set  $A \sqsupseteq B$  are similar if they are not then make them similar.

Time

$$d(t_1, t_2) = \begin{cases} t_2 - t_1, & \text{if } t_1 \leq t_2 \\ 24 + (t_2 - t_1), & \text{if } t_1 > t_2 \end{cases}$$

$$\begin{aligned} d(1\text{PM}, 2\text{PM}) &= 1 & \rightarrow \text{non-symmetric} & d(t_1, t_2) = |t_1 - t_2| \\ d(2\text{PM}, 1\text{PM}) &= 23 & \rightarrow \text{now it becomes symmetric} \end{aligned}$$

### \* Similarity Between Data Objects

If  $\delta(x, y)$  is similarity b/w the points  $x \approx y$  then the properties of symmetry are as follows:-

- (1)  $\delta(x, y) = 1$  if  $x = y$
- (2)  $\delta(x, y) = \delta(y, x), \forall x \approx y$

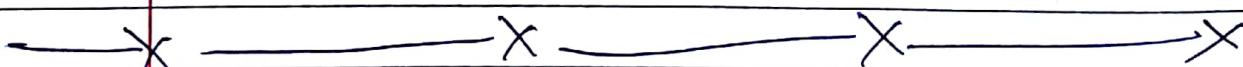
e.g. Consider an experiment in which people are asked to classify a small set of characters as they flash on screen. Here each character is classified as itself or as another character. For e.g. "O" appeared 160 times and was classified as O 160 times and as O 40 times. Suppose O appeared 200 times and was classified as <sup>small</sup>O 170 times but as O 30 times. Find out a similarity measure that is symmetric.

$$\delta(O, O) = (160, 40) = 160 + 40 = 200$$

$$\delta(O, O) = (170, 30) = 170 + 30 = 200$$

New  $\rightarrow$  similarity measure

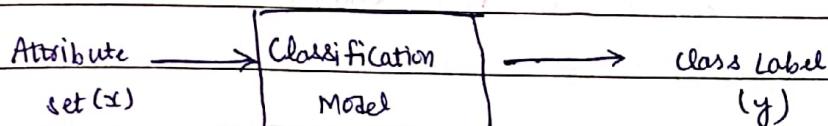
$$\delta'(x, y) = \delta'(y, x) = \delta(x, y) + \delta(y, x)$$



It's the end of test syllabus.

## CLASSIFICATION

Classification is the task of assigning objects to one of the several predefined categories e.g. detecting spam e-mail messages based upon the message header and content. The i/p data for a classification task is a collection of records. Each record is characterized by a tuple  $(x, y)$  where  $x$  is the attribute set and  $y$  is a special attribute called as CLASS LABEL.



Consider a sample dataset Vertebrates data set

| Name   | Bodytemp | Gives Birth | Has Legs | Aerial Creature | Class Label |
|--------|----------|-------------|----------|-----------------|-------------|
| Human  | Warm     | Yes         | Yes      | No              | Mammal      |
| Frog   | Cold     | No          | Yes      | No              | Amphibian   |
| Bat    | Warm     | Yes         | Yes      | Yes             | Mammal      |
| Pigeon | Warm     | No          | Yes      | Yes             | Bird        |
| :      | :        | :           | :        | :               | :           |
| :      | :        | :           | :        | :               | :           |
| :      | :        | :           | :        | :               | :           |

A sample dataset used for classifying Vertebrates into one of the following category:- (Mammal, bird, fish, reptile, amphibian) The attr. set includes the properties such as body temperature, ability to fly, method of reproduction etc. The class label must be a discrete attribute as it's a key characteristic that distinguishes classification from Regression. In regression class label is continuous.

→ Classification is task of learning a target function  $f$  that maps each attribute set  $x$  to one of the predefined class label  $y$ . The target func<sup>n</sup> is also known as classification model. A classif<sup>n</sup> model is useful for :-

- 1) descriptive modelling :- A CM can serve as an explanatory tool to distinguish b/w objects of different classes.
- 2) predictive modelling :- This model can be used to predict class label of unknown

Yellosd.

predictive  
modelling

A CM can be treated as a blackbox that automatically assigns a class label when presented with an attribute set of an unknown record.

→ Classification techniques are most suited for predicting or describing the data with binary or nominal categories and are less effective for ordinal categories. Below they don't consider the implicit order among the categories.

→ A classification technique is a systematic approach to building classification models from an input dataset. Some examples are decision tree classifiers, rule based classifiers, neural networks, support vector machines, naive Bayes classifiers. Each technique employs a learning algorithm to identify a model that best fits the relationship b/w the attribute set & class label of the input data and correctly predicts the class label of records.

→ For solving classification problems first a training set consisting of records whose class labels are known must be provided. The training set is used to build a CM which is subsequently applied to the test set, which consists of records with unknown labels. Evaluation of the performance of a classifier is based on the counts of test records predicted correctly or incorrectly by the model. The counts are tabulated in a table known as CONFUSION MATRIX. Each entry  $f_{ij}$  in the table denotes the no. of records from class i predicted to be class j.

Confusion Matrix

|        |               | Predicted     |               | TP → True Positive  |
|--------|---------------|---------------|---------------|---------------------|
|        |               | Class = 1     | Class = 0     |                     |
| Actual | Class = 1 (t) | $f_{11}$ (Tp) | $f_{10}$ (Fn) | FN → False Negative |
|        | Class = 0 (f) | $f_{01}$ (Fp) | $f_{00}$ (Tn) | TN → True Negative  |

To determine how well a classification model performs we may summarise the information with a single no. that would make it more convenient to compare the performance of different models. This can be done using performance metric.

$$\text{Accuracy} = \frac{f_{11} + f_{00}}{f_{00} + f_{01} + f_{10} + f_{11}} = \frac{\text{no. of correct predictions}}{\text{Total no. of co predictions}}$$

(E) Error rate =  $\frac{\text{no. of wrong predictions}}{\text{Total no. of predictions}}$  =  $\frac{f_{00} + f_{10}}{f_{00} + f_{11} + f_{01} + f_{10}}$

|   |   |
|---|---|
| 6 | 4 |
| 2 | 8 |

$$\text{Accuracy} = \frac{14}{20} = \frac{7}{10} = 0.7$$

$$\text{Error rate} = \frac{3}{10} = 0.3$$

9.09

Q Suppose there are 2 predicted classes Yes or NO. Classifier made total of 165 predictions out of which the classifier predicted Yes  $\rightarrow$  110 times and No  $\rightarrow$  55 times. In reality 105 patients in the sample had the disease and 60 doesn't. The classifier correctly predicted the no. of patients who had the disease as 100. Calculate accuracy and error rate.

|            |     | Predicted |     | Actual Yes | Actual No | accuracy = $\frac{150}{165} = \frac{10}{11} = 0.909$ |
|------------|-----|-----------|-----|------------|-----------|------------------------------------------------------|
|            |     | Yes       | No  |            |           |                                                      |
| Actual Yes | 100 | 50        | 105 |            |           |                                                      |
| Actual No  | 50  | 60        | 110 |            |           |                                                      |
|            | 110 | 55        | 165 |            |           | error = $\frac{15}{165} = \frac{1}{11} = 0.0909$     |

Q Evaluating Spam classifier Consider a classical problem of predicting spam and non-spam email by using binary classification model. Classifier made 100 predictions are modelled. Classified 95 emails  $\rightarrow$  85 as non-spam and 10 as spam. 5 emails were actually spam but predicted as non-spam. Calculate error rate and accuracy.

|                 |    | Predicted |          | Actual spam | Actual non-spam | accuracy = $\frac{95}{100} = 0.95$ |
|-----------------|----|-----------|----------|-------------|-----------------|------------------------------------|
|                 |    | spam      | non-spam |             |                 |                                    |
| Actual spam     | 10 | 85        | 95       |             |                 |                                    |
| Actual non-spam | 0  | 85        | 100      |             |                 | error = 0.05                       |

95

Q A binary classification model predicted it to be a cat 55 times & dog 45 times. Actual both 50. The model predicted the cat incorrectly 10 times. Find error and accuracy rate.

$$\text{Accuracy} = 0.75$$

$$\text{Error rate} = 0.25$$

|        |     | Predicted |     |    |
|--------|-----|-----------|-----|----|
|        |     | Cat       | Dog |    |
| Actual | Cat | 15        | 40  | 50 |
|        | Dog | 40        | 10  | 50 |
|        |     | 55        | 45  |    |

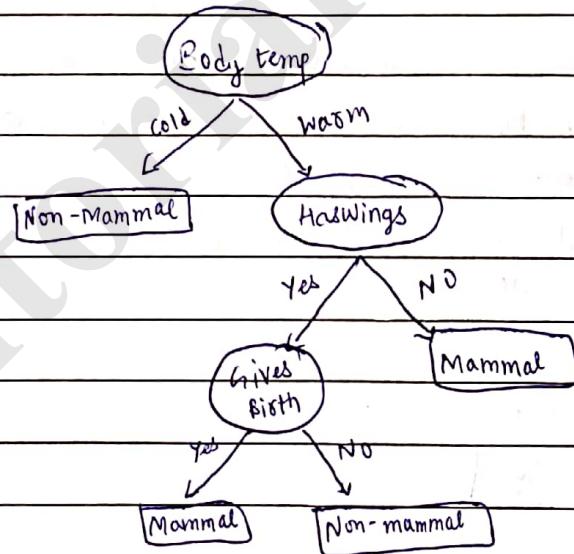
$$\text{ER} = 0.25$$

$$0.25$$

Q A binary classification model predicted fraud 7 times and not fraud 93 times in reality (6, 94). The model correctly predicted fraud as 4 times.

Find accuracy and error rate.

|       |           | Fraud                              | non-fraud |    |
|-------|-----------|------------------------------------|-----------|----|
| Fraud | Fraud     | 4                                  | 12        | 16 |
|       | non-fraud | 3                                  | 91        | 94 |
|       |           | 7                                  | 93        |    |
|       |           | accuracy = $\frac{95}{100} = 0.95$ |           |    |
|       |           | error rate = 0.05                  |           |    |



→ Decision tree has 3 types of nodes :-

- 1) Root node which has no incoming edges and 0 or more outgoing edges.
- 2) Internal nodes each of which has exactly 1 incoming edge and 2 or more outgoing edges
- 3) Leaf or terminal nodes each of which has exactly 1 incoming edge and no outgoing edges.

Classifying a test record is straight forward once a decision tree has been constructed. starting from the root node. we apply the test condition to the record and follow the appropriate branch based on outcome of the test. This will lead us either to another internal node for which a new test cond<sup>n</sup> is applied or to a leaf node. The class label associated with the leaf node is then assigned to the leaf node.

## Q HOW TO BUILD A DECISION TREE ?

Ans Hunt Algorithm. There are exponentially many decision trees that can be constructed from a given set of attributes finding an optimal tree is computationally infeasible because of the exponential size of the search space. Efficient Algos have been developed that employ a greedy strategy to grow a decision tree by making a series of locally optimum decisions about which attribute to use for partitioning the data. Hunt Algorithm is the basis of many existing decision tree algorithms including ID3, C4.5 and CART

### \* HUNT'S ALGORITHM

→ Here a decision tree is grown in recursive fashion by partitioning the training records into successively purer subsets.

$D_t \rightarrow$  set of training records associated with a node  $t$

$y_1, - , y_3$  will be the class labels.

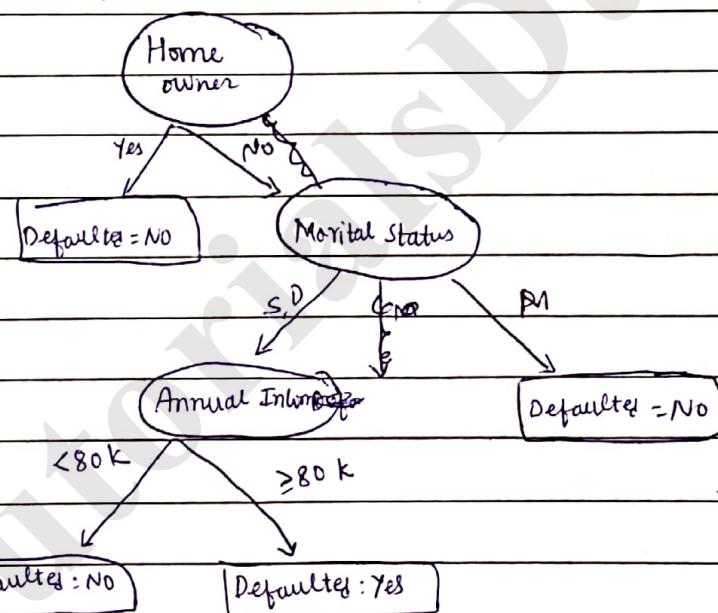
Step1 → If all records in  $D_t$  belong to the same class  $y_t$  then  $t$  is a leaf node labelled as  $y_t$ .

Step2 → If  $D_t$  contains records that belong to more than 1 class an attribute test cond<sup>n</sup> is selected to partition the records into smaller subsets. A child node is

Created for each outcome of the test cond" and records in Dt are distributed based on the outcomes. The algorithm is then recursively applied to each child node.

The problem of predicting whether a loan applicant will repay his or her loan obligation or becomes defaulter

| Tid | Home owner | Marital Status | Annual Income | Defaulted Borrower |
|-----|------------|----------------|---------------|--------------------|
| 1   | Yes        | S              | 125 K         | No                 |
| 2   | No         | M              | 100 K         | No                 |
| 3   | Yes        | S              | 70 K          | No                 |
| 4   | No         | M              | 120 K         | No                 |
| 5   | No         | D              | 95 K          | Yes                |
| 6   | No         | M              | 60 K          | No                 |
| 7   | Yes        | D              | 220 K         | No                 |
| 8   | No         | S              | 85 K          | Yes                |
| 9   | No         | M              | 75 K          | No                 |
| 10  | No         | S              | 90 K          | Yes                |



Hunt Algorithm will work if every combination of attr. values is present in training data and each combination has a unique class label but in most practical situations additional conditions are needed to handle following situations

- If none of the training records have the combination of attribute values associated with some node then the node is declared a leaf node with the same

Class label as the majority class of training records associated with its parent node.

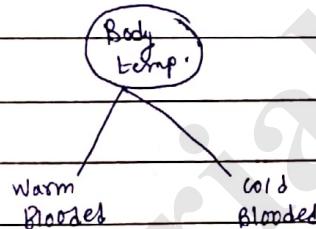
2) If all the records associated with  $D_t$  have identical attribute values but different class labels then it's not possible to split records any further and the node is declared a leaf node with the same class label as the majority class.

### Design issues of decision tree Algo

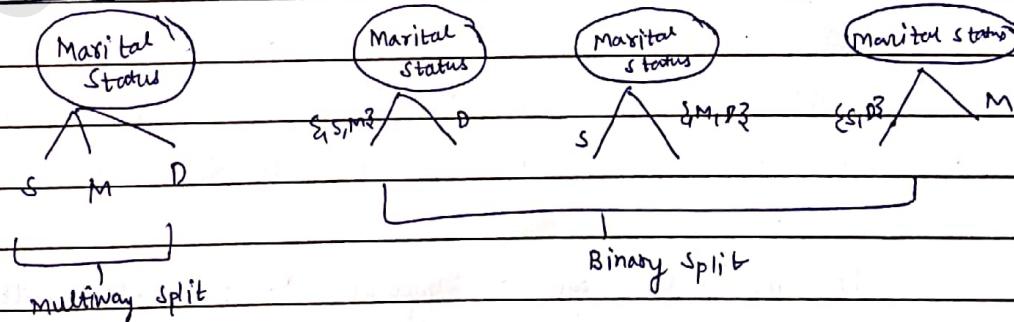
- 1) How should the training records be split each recursive step must select an attr. test condition to divide the records into smaller subsets. It must also define an objective measure for evaluated the goodness of each test condition.
- 2) How should the splitting procedure stop? A possible strategy is to continue expanding a node until either all the records  $\in$  same class or all the records have identical attribute values or other criteria may be imposed to terminate the tree growing procedure earlier.

### \* Methods for Expressing Attribute Test Conditions

#### 1) BINARY ATTRIBUTES



#### 2) NOMINAL ATTRIBUTES



Since a nominal attribute can have many values its test condition can be expressed in 2 ways first is a multiway split in which the no. of outcomes depends on no. of distinct values for the corresponding attribute.  $\oplus$   
second is some decision tree algo such as CART produce only binary splits

# **TutorialsDuniya.com**

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

**Please Share these Notes with your Friends as well**

**facebook**

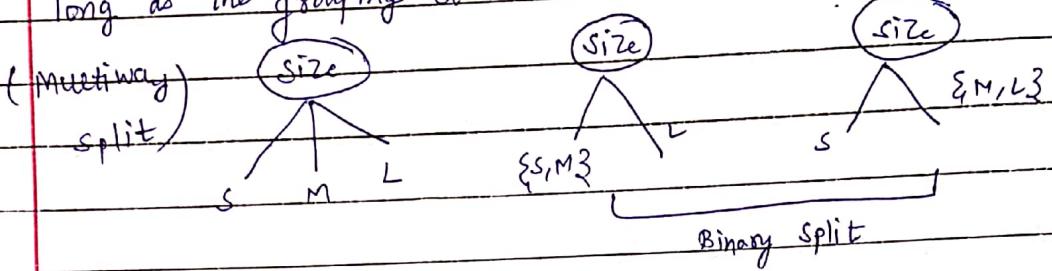
**WhatsApp** 

**twitter** 

**Telegram** 

### 3) ORDINAL ATTR.

They can also produce binary or multiway splits they can be grouped as long as the grouping doesn't violate the order property of attribute values.



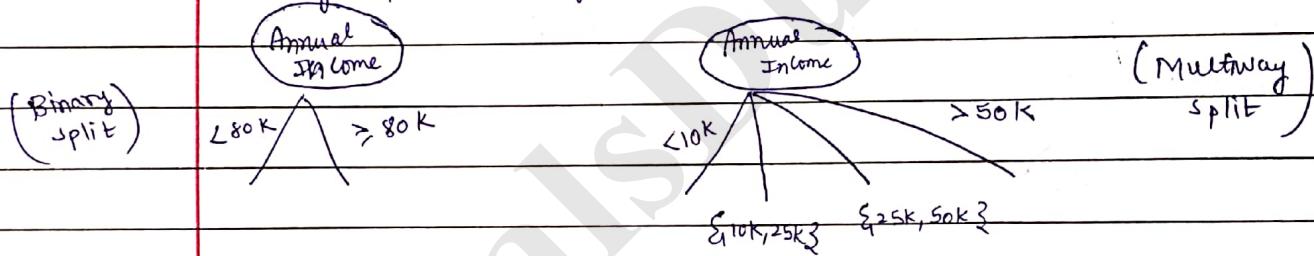
### 4) CONTINUOUS ATTRIBUTES

for contin. attr. the test condition can be expressed as a comparison test

$A < \vee$  or  $A \geq \vee$  with binary outcomes or a range query with outcomes

of the form  $\forall i \leq A < \forall i+1, i=1\dots k$

for the binary case the decision tree algo. must consider all the possible split positions  $\forall i$  and selects the one that produces the best partition. for multiway split the algo. must consider all possible ranges of continuous values.



### \* MEASURES FOR SELECTING THE BEST SPLIT

Let  $P(i|t)$  denote the fraction of records belonging to class  $i$  at a given node  $t$ . In a 2 class problem the class distribution at any node can be written as  $(p_0, p_1) \in [0, 1]^2$ ,  $p_0 = 1 - p_1$ . The measures developed for selecting the best split are based on degree of impurity of the child nodes.

The smaller the degree of impurity the more skewed the class distribution is. A node with class distribution  $(0, 1)$  has 0 impurity whereas a node with uniform C.D.  $(0.5, 0.5)$  has the highest impurity.

Different impurity measures :-

$C \rightarrow$  no. of classes

$$(1) \text{ Entropy } (t) = - \sum_{i=0}^{C-1} p(i|t) \log_2 p(i|t)$$

$[0 \log_2 0 = 0]$

$$(2) \text{ Gini } (t) = 1 - \sum_{i=0}^{C-1} [p(i|t)^2]$$

$$(3) \text{ classification error} = 1 - \max_i [p(i|t)]$$

| Node N <sub>1</sub> | Count |
|---------------------|-------|
| class = 0           | 0     |
| class = 1           | 6     |

Find Entropy

$$\text{Entropy} = - \left[ \frac{0}{6} \log_2 \frac{0}{6} + \frac{1}{6} \log_2 \frac{1}{6} \right] = 0$$

(pure set)

$$\text{Gini} = 1 - \left[ \frac{0^2 + 1^2}{6} \right] = 0 \quad (\text{pure subset})$$

$$p(i|t)_{\max} = 1$$

$$\text{Classification error} = 1 - 1 = 0 \quad (\text{pure subset})$$

| Node N <sub>1</sub> | Count |
|---------------------|-------|
| class = 0           | 1     |
| class = 1           | 5     |

$$\begin{aligned} \text{Entropy} &= - \left[ \frac{1}{6} \log_2 \frac{1}{6} + \frac{5}{6} \log_2 \frac{5}{6} \right] \\ &= - \left[ \frac{1}{6} [-\log_2 6] + \frac{5}{6} [\log_2 5 - \log_2 6] \right] \end{aligned}$$

$$\text{Gini} = 1 - \left[ \frac{\frac{1}{36}}{36} + \frac{25}{36} \right] = \frac{10}{36} = \frac{5}{18} = 0.278$$

$$\text{Classification error} = 1 - \frac{5}{6} = \frac{1}{6} = 0.167$$

| N <sub>3</sub> | Count |
|----------------|-------|
| C = 0          | 3     |
| C = 1          | 3     |

$$\begin{aligned} \text{Gini Entropy} &= - \left[ \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right] \\ &= - \left[ \frac{1}{2} - \frac{1}{2} \right] = 0 \end{aligned}$$

$$\text{Gini} = 1 - \left[ \frac{1}{4} + \frac{1}{4} \right] = \frac{1}{2} = 0.5$$

$$\text{Classification error} = 1 - \frac{1}{2} = \frac{1}{2} = 0.5$$

So Node N<sub>1</sub> is the best based on these calculations as it

Contributor: Deepanshu has the lowest impurity value followed by N<sub>2</sub> and N<sub>3</sub>.

$I(\text{parent})^{0.6}$

min entropy or diff from parent

$$\begin{array}{ccc} \textcircled{1} & \textcircled{2} & \checkmark \\ I(v_1) & I(v_2) & 0.2 \\ 0.3 & & \end{array}$$

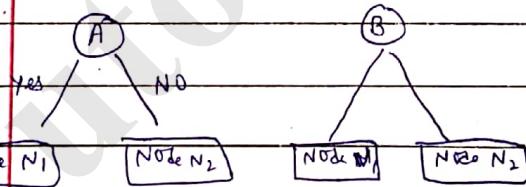
To determine how well a test condition performs we need to compare the degree of impurity of parent node before splitting with the degree of impurity of the child nodes after splitting. The larger the difference the better the test condition. This difference called  $\Delta$  (gain).

$$\Delta(\text{gain}) = I(\text{parent}) - \sum_{j=1}^K N(v_j) I(v_j)$$

$N \rightarrow$  total no. of records at parent node ,  $K \rightarrow$  no. of attributes values  
 $N(v_j) \rightarrow$  no. of records associated with child node  $v_j$ . Decision tree algo choose a test condition that maximizes the gain  $\Delta$ . When entropy is used as impurity measure in the above formula the  $\Delta$  is known as  $\Delta_{\text{info}}$  or information gain.

#### \* Splitting of binary attributes

|       | Parent |
|-------|--------|
| $C_0$ | 6      |
| $C_1$ | 6      |



|       | $N_1$ | $N_2$ |  | $N_1$ | $N_2$ |
|-------|-------|-------|--|-------|-------|
| $C_0$ | 4     | 2     |  | 1     | 5     |
| $C_1$ | 3     | 3     |  | 4     | 2     |

If attribute A is chosen

$$1) \text{ Gini index for } N_1 = 1 - \left[ \frac{16}{49} + \frac{9}{49} \right] = \frac{24}{49} = 0.489$$

$$\text{Gini index for } N_2 = 1 - \left[ \frac{4}{25} + \frac{9}{25} \right] = \frac{12}{25} = 0.48$$

$$\begin{aligned} \text{Weighted average of Gini index for attr. A} &= \frac{7(0.489) + 5(0.48)}{12} \\ &= 0.486 \end{aligned}$$

$$\begin{aligned} \text{attr. A} \rightarrow \Delta(\text{gain}) &= 0.5 - 0.486 = I_{\text{parent}} - 0.486 \\ &= 0.014 \end{aligned}$$

If attribute B is chosen

$$\text{Gini index for } N_1 = 1 - \left[ \frac{1}{25} + \frac{16}{25} \right] = \frac{8}{25} = 0.32$$

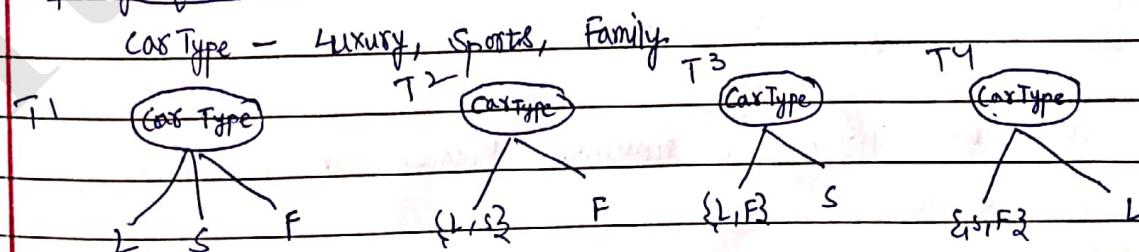
$$\text{Gini index for } N_2 = 1 - \left[ \frac{25}{49} + \frac{4}{49} \right] = \frac{20}{49} = 0.408$$

$$\begin{aligned} \text{Weighted average of Gini index for attr. B} &= \frac{5(0.32) + 7(0.408)}{12} \\ &= 0.376 \end{aligned}$$

$$\begin{aligned} \text{attr. B} \rightarrow \Delta(\text{gain}) &= 0.5 - 0.376 \\ &= 0.124 \end{aligned}$$

So we will choose attribute B as next partitioning attribute.

\* Splitting of Nominal attributes



| Car Type           |     |   |
|--------------------|-----|---|
| F                  | S   | L |
| C <sub>0</sub>   1 | 8   | 1 |
| C <sub>1</sub>   3 | 0   | 7 |
|                    | + 8 | 8 |

enini index

$$F \rightarrow 1 - \left[ \frac{1}{96} + \frac{9}{16} \right] = \frac{6}{16} = 0.375$$

$$S \rightarrow 1 - [1+0] = 0$$

$$L \rightarrow 1 - \left[ \frac{1}{64} + \frac{49}{64} \right] = \frac{14}{64} = 0.21875$$

Weighted average =  $\frac{4}{20}(0.375) + \frac{8}{20}(0) + \frac{8}{20}(0.21875)$   
 of enini index T<sub>1</sub>  
~~256~~  
~~-130~~  
~~126~~

$$= 0.163$$

| Car Type           |      |   |
|--------------------|------|---|
| F                  | S, L | F |
| C <sub>0</sub>   5 | 9    | 1 |
| C <sub>1</sub>   7 | 3    |   |

$$\{S, L\} \rightarrow 1 - \left[ \frac{81}{256} + \frac{49}{256} \right] = 0.5080.492$$

$$\{F\} \rightarrow 1 - \left[ \frac{1}{16} + \frac{9}{16} \right] = 0.375$$

Weighted average =  $\frac{16}{20}(0.508) + \frac{4}{20}(0.375) = 0.468$   
 T<sub>2</sub>

| Car Type            |   |  |
|---------------------|---|--|
| S, L, F             | S |  |
| C <sub>0</sub>   2  | 8 |  |
| C <sub>1</sub>   10 | 0 |  |

$$\{L, F\} \rightarrow 1 - \left[ \frac{4}{144} + \frac{100}{144} \right] = 0.278$$

$$S \rightarrow 1 - [1+0] = 0$$

Weighted average T<sub>3</sub> =  $\frac{12}{20}(0.278) + \frac{8}{20}(0)$   
 $= 0.167$

So we will choose T<sub>1</sub> structure because in multivay split the attributes are separated.

\* Splitting of Continuous attributes

Annual Income

| class | No | No | No | Yes | Yes | Yes | No  | No  | No  | No |
|-------|----|----|----|-----|-----|-----|-----|-----|-----|----|
| 60    | 70 | 75 | 85 | 90  | 95  | 100 | 120 | 125 | 220 |    |
| 55    | 65 | 72 |    |     |     |     |     |     |     |    |
| 0     | 3  |    |    |     |     |     |     |     |     |    |
| 10    | 7  |    |    |     |     |     |     |     |     |    |

Annual Income  $\rightarrow$  high prob.  $\rightarrow$  high prob.

| Class | No   | No   | No     | Yes   | Yes   | Yes  | No   | No    | No    | No   |      |
|-------|------|------|--------|-------|-------|------|------|-------|-------|------|------|
| 60    | 70   | 75   | 85     | 90    | 95    | 100  | 120  | 125   | 220   |      |      |
| 55    | 65   | 72   | 80     | 87    | 92    | 97   | 110  | 122   | 172   | 230  |      |
| <= >  | <= > | <= > | <= >   | <= >  | <= >  | <= > | <= > | <= >  | <= >  | <= > |      |
| Yes   | 0    | 3    | 0      | 3     | 1     | 2    | 1    | 3     | 0     | 3    |      |
| No    | 0    | 7    | 1      | 6     | 2     | 5    | 3    | 4     | 3     | 5    |      |
| Uniq  |      |      |        |       |       |      |      |       |       |      |      |
|       | 0.42 | 0.4  | -0.375 | 0.342 | 0.416 | 0.4  | 0.3  | 0.313 | 0.375 | 0.4  | 0.42 |

65 0

$$\underline{E_{42}} \quad 1 - \left[ \frac{5}{9} \right] = \frac{4}{9} = 0.44$$

$$\underline{\underline{92}} \quad 1 - \left[ \frac{4}{25} + \frac{9}{25} \right] = \frac{12}{25} = 0.48$$

$$\underline{\underline{1 - \left[ \frac{1}{25} + \frac{16}{25} \right] = \frac{8}{25} = 0.32}}$$

$$\frac{1}{2}(0.48) + \frac{1}{2}(0.32)$$

$$\underline{\underline{72}} \quad 1 - \left[ \frac{9}{64} + \frac{25}{64} \right] = \frac{30}{64} = 0.469 \quad = 0.4$$

$$\underline{\underline{\frac{8}{10} (0.469) = 0.3}}$$

$$\underline{\underline{1 - \left[ \frac{9}{49} + \frac{16}{49} \right] = \frac{24}{49} = 0.489}}$$

$$\underline{\underline{1 - \left[ \frac{9}{64} + \frac{25}{64} \right] = \frac{30}{64} = 0.469}}$$

$$\underline{\underline{80}} \quad 0 \quad 1 - \left[ \frac{9}{49} + \frac{16}{49} \right] = \frac{24}{49} = 0.489$$

$$\frac{3}{9} + \frac{6}{9}$$

$$\underline{\underline{87}} \quad 1 - \left[ \frac{1}{16} + \frac{9}{16} \right] = 0.375$$

$$\frac{1}{3} + \frac{2}{3}$$

$$\underline{\underline{1 - \left[ \frac{1}{9} + \frac{4}{9} \right] = \frac{4}{9} = 0.44}}$$

$$\underline{\underline{\frac{2}{6}}} \quad 1 - \left[ \frac{1}{9} + \frac{4}{9} \right] = \frac{4}{9} = 0.44$$

$$0.4$$

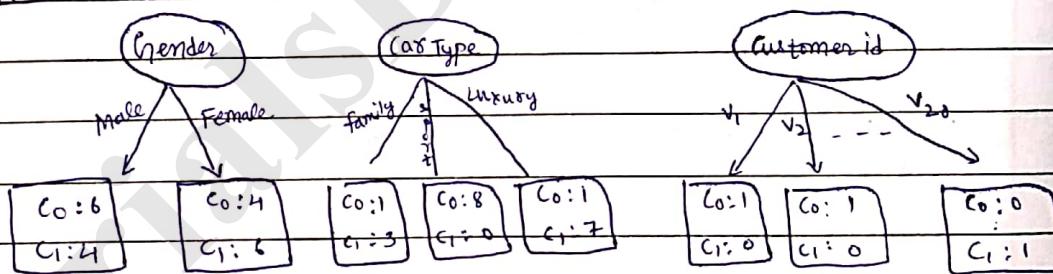
$$\underline{\underline{\frac{3}{10} (0.375) + \frac{6}{10} (0.44) = 0.376}}$$

$$\underline{\underline{1 - \left[ \frac{9}{400} + \frac{49}{100} \right] = }}$$

A brute force method for finding  $V$  is to consider every value of the attribute as a candidate split position for each candidate  $V$  the dataset is scanned once to count the no. of records with annual income less than or greater than  $V$ . We then compute the Gini index for each candidate and choose the one that gives the lowest value. This approach is computationally expensive because it requires  $O(n)$  time to compute the Gini index at each candidate split pos<sup>n</sup> since there are  $n$  candidates the overall complexity is  $O(n^2)$ . So to reduce the complexity the training records are sorted based on their annual income which requires  $O(n \log n)$  time Candidate & split positions are identified by taking the mid points b/w two adjacent sorted values.

This technique can be further optimized by considering only the candidate split pos<sup>n</sup> where located b/w 2 adjacent records with different class labels. This approach allows us to reduce the no. of candidate split positions from 11 to 2.

### \* GAIN RATIO



Here Customer id is not a predictive attribute because its value is unique for each record if the no. of records associated with each partition is too small it will not enable us to make any reliable prediction. So there are 2 ways of overcoming this problem:-

- 1) We can restrict the test conditions to binary splits only. This strategy is used by CART decision tree algorithm.
- 2) To modify the splitting criteria to take into account the no. of outcomes produced by the attribute test condition. In C4.5 decision tree algo. a splitting criterion known as gain ratio is used to determine the goodness of a split.

$$\text{Gain ratio} = \frac{\Delta \text{info}}{\text{Split info}}$$

$$\text{Split info} = -\sum_{i=1}^K P(v_i) \log_2 P(v_i)$$

$K \rightarrow$  total no. of splits

If an attribute produces a large no. of splits its split info will also be large which in turn reduces its gain ratio.

~~Each attribute has same number of splits~~

Ex: total no. of records are 21. There are 3 partitions & records each. calculate split info.

$$\text{Split info} = -3 \left( \frac{7}{21} \log_2 \left( \frac{7}{21} \right) \right)$$

$$= -3 \log_2 \left( \frac{7}{21} \right) = +\log_2 (3) = \frac{\log_{10} 3}{\log_{10} 2} = 1.585$$

|   | $a_1$ | $a_2$ | $a_3$ | Target class | Q: What are the information gains of $a_1, a_2$ and $a_3$ (compute for every possible split)? What is best split among $a_1, a_2$ and $a_3$ according to info. gain? |
|---|-------|-------|-------|--------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| x | 1     | T     | T     | +            |                                                                                                                                                                      |
| x | 2     | T     | F     | 6.0          | +                                                                                                                                                                    |
| x | 3     | T     | F     | 5.0          | -                                                                                                                                                                    |
| x | 4     | F     | F     | 4.0          | +                                                                                                                                                                    |
| x | 5     | F     | T     | 7.0          | -                                                                                                                                                                    |
| x | 6     | F     | T     | 3.0          | -                                                                                                                                                                    |
| x | 7     | F     | F     | 8.0          | -                                                                                                                                                                    |
| x | 8     | T     | F     | 7.0          | +                                                                                                                                                                    |
| x | 9     | F     | T     | 5.0          | -                                                                                                                                                                    |

$$\text{Parent entropy} = -\left[ \frac{4}{9} \log_2 \frac{4}{9} + \frac{5}{9} \log_2 \frac{5}{9} \right] = 0.991$$

$a_1$

|   | T | F | Entropy T = $-\left[ \frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4} \right] = 0.811$ |
|---|---|---|-------------------------------------------------------------------------------------------------------|
| + | 3 | 1 |                                                                                                       |
| - | 1 | 4 | Entropy F = $-\left[ \frac{1}{5} \log_2 \frac{1}{5} + \frac{4}{5} \log_2 \frac{4}{5} \right] = 0.721$ |

$$\text{Weighted} = \frac{4}{9} (0.811) + \frac{5}{9} (0.721) = 0.761$$

$$\frac{0.761}{230}$$

$$\Delta \text{info} = 0.230$$

$a_2$

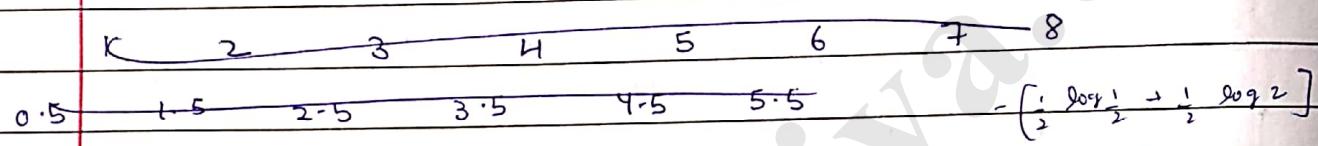
|   | T | F |
|---|---|---|
| + | 2 | 2 |
| - | 3 | 2 |

$$\text{Entropy } A = - \left[ \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right] = 0.970$$

$$\text{Entropy } F = - \left[ \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right] = 1$$

$$\text{Weighted entropy } a_2 = \frac{5}{9} (0.970) + \frac{4}{9} (1) = 0.983$$

$$\Delta \text{ info} = 0.008$$



|     | +  | -     | +     |       |       |       |       |       |    |
|-----|----|-------|-------|-------|-------|-------|-------|-------|----|
|     | 1  | 3     | 4     | 5     | 6     | 7     | 8     |       |    |
| 2   | <= | >     | <=    | >     | <=    | >     | <=    | >     | <= |
| +   | 1  | 3     | 1     | 3     | 2     | 2     | 2     | 3     | 1  |
| -   | 0  | 5     | 1     | 4     | 1     | 4     | 3     | 2     | 4  |
|     | 0  | 0.954 | 0.985 | 0.985 | 0.918 | 0.983 | 0.972 | 0.888 |    |
| → 0 |    | 0.848 | 0.988 |       |       |       |       |       |    |

$$\frac{3}{8} \log \frac{3}{8} + \frac{5}{8} \log \frac{5}{8}$$

$$\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{2}{3}$$

$$0.954$$

$$\frac{8}{9} (0.954)$$

$$0.991$$

$$\frac{-0.848}{0.143}$$

$$\text{Entropy} = 0.848$$

$$\Delta \text{ info} =$$

So  $a_1$  will be selected.

## \* Evaluating the Performance of a Classifier

These are 2 types of errors committed by a classification model

- 1) Training Error or Resubstitution error or Apparent Error - It's the no. of misclassification errors committed on training records.
- 2) Generalization Error :- expected error of the model on previously unseen records or the test set.

A good model must have low training error as well as low generalization error.

A model that fits the training data too well can have a poor generalization error than a model with higher training error such a situation is known as model overfitting.

After building the decision tree a tree-pruning step can be performed to reduce the size of decision tree. To prevent a phenomenon known as overfitting. Pruning improves the generalization capability of decision tree.

It's often useful to measure the performance of the model on the test set because it provides an unbiased estimate of its generalization error. However to do this the class label of the test records must be known. Some of the methods commonly used to evaluate performance of classifiers are :-

- 1) Holdout Method :- In this method the original data with labelled examples is partitioned into 2 distinct sets called the training and testing set. A classification model is then induced from the training set and performance is evaluated on the test set. The proportion of data for training and testing is either 50-50 or  $\frac{2}{3}$  rd for training and  $\frac{1}{3}$  rd for testing. The accuracy of the classifier is based on the accuracy of the induced model on the test data set.

Disadvantages :-

- 1) The induced model may not be as good as when all the labelled examples are used for training.
  - 2) The model may be highly dependent on composition of training and test set.
  - 3) The training and test sets are no longer independent of each other as they are subsets of the same data.
- 2) Random Subsampling :- The holdout method can be repeated several times to improve the estimation of classifier performance. This approach is

Known as random subsampling

$$\text{acc}_k = \frac{1}{K} \left( \sum_{i=1}^K \underline{\text{acc}}_i \right)$$

$\underline{\text{acc}}_i \rightarrow$  accuracy of model

### Disadvant

- 1) It doesn't utilize as much data as possible for training.
- 2) It has no control over the no. of ~~no~~ times each sample is used in training or testing set.

### 3) Cross Validation

In this approach each record is used the same no. of times for training and exactly once for testing.

- a) 2 Fold Cross Validn → in this we partition our data into 2 equal sizes subsets. We choose one of the subset for training and other for testing. We then swap the role of subsets and the total error is obtained by summing up the errors for both the runs.
- b) K Fold Cross Validation → in this method the data is segmented into K equal size partitions during each run one of the partition is used for testing and others are used for training. This procedure is repeated K times so that each partition is used for testing at least once. The total error is obtained by summing up all the K errors.
- c) Leave one out approach :- In this  $K = N$  (no. of records), in this each test set contains only 1 record and it has the advantage of utilizing as much ~~as~~ data as possible for training. 1) It's computationally to repeat the procedure N times. 2) Since each test set contains only 1 record the variance of the performance tends to be high.

### 4) Bootstrap

In bootstrap approach the training records are sampled with replacement i.e. a record on an average a bootstrap sample of size  $n'$  contains about 63.2% of the records in the original data.

Probability that a record is chosen out of  $n'$  samples is  $\frac{1}{n}$

Prob that a record is not chosen out of  $n$  samples =  $1 - \frac{1}{n}$

Prob that it's not chosen ~~at~~  $n$  times =  $\left(1 - \frac{1}{n}\right)^n$

Prob that it's chosen at least once by a bootstrap sample =  $1 - \left(1 - \frac{1}{n}\right)^n$

also called .632 bootstrap

Records that are not included in the bootstrap sample become part of the test set

The model induced from the training set is then applied to the test set to obtain the <sup>estimate of</sup> accuracy of bootstrap sample ( $\hat{\epsilon}_i$ ). The sampling procedure is then repeated  $b$  times to generate  $b$  bootstrap samples.

$$\text{Accuracy} = \frac{1}{b} \sum_{i=1}^b (0.632 \times \hat{\epsilon}_i + 0.368 \times \text{acc}(x))$$

$\hat{\epsilon}_i \rightarrow$  accuracy of test set

$\text{acc}(x) \rightarrow$  accuracy obtained when we run model on training data.

## RULE BASED CLASSIFIER

\* Rule Based Classifier :  $\rightarrow$  It is a technique for classifying the records using a collection of if then rules. The rules for the model are represented in disjunctive normal form.

$R = (\gamma_1 \vee \gamma_2 \vee \dots \vee \gamma_k)$  Where  $R$  is the ruleset and  $\gamma_i$  are the classification rules each classification rule can be expressed as:

$\gamma_i : (\text{Condition}) \rightarrow y_i$

(GivesBirth = yes)  $\wedge$  (BodyTemp = WB)  $\longrightarrow$  Mammals  
Antecedent Consequent

Condition $_i = (A_1 \text{ op } V_1) \wedge (A_2 \text{ op } V_2) \wedge \dots \wedge (A_k \text{ op } V_k)$

A rule  $\gamma$  covers a record  $x$  if preCondition of  $\gamma$  matches the attributes of  $x$ .  
 $\gamma$  is said to be fired or triggered whenever it covers a given record

Accuracy

Quality of a classification rule can be evaluated using coverage and accuracy.

Coverage of a rule is defined as the fraction of records in  $D$  (dataset) that trigger the rule  $\gamma$ .

$$\text{Coverage} = \frac{|A|}{|D|}, \quad \text{Accuracy} = \frac{|A \cap Y|}{|A|}$$

$|A| \rightarrow$  Condition part same as rule

$|D| \rightarrow$  Total no. records

Accuracy  $\rightarrow$  defined as fraction of records triggered by rule  $\gamma$  whose class labels are equal to  $y$ .

| Name   | B.T | GivesBirth | Class Label |
|--------|-----|------------|-------------|
| Human  | WB  | Y          | Mammal      |
| Python | CB  | N          | Reptile     |
| Salmon | CB  |            |             |

$$|A| = 5, |Any| = 3$$

$$|D| = 16 \quad 0.3125$$

$$\text{Coverage} = \frac{5}{16} = 0.3125 \quad \text{Acc} = \frac{3}{5} = 0.6$$

## Q How a Rule based classifier Works ?

A rule based classifier classifies a test record based on the rule triggered by the record.

The two important properties of ruleset generated by a rule based classifier

- 1) Mutually Exclusive Rules - Each In this no 2 rules in R are triggered by same record. This property ensures that every record is covered by almost one rule in R.
- 2) Exhaustive Rule :- In this every record is covered by at least one rule in R. These properties ensure that every record is covered by exactly 1 rule.

If the ruleset isn't exhaustive then a default rule  $\bar{y}_d : () \rightarrow y_d$  must be added to cover the remaining cases. A default rule has an empty antecedent and  $y_d$  is the default class assigned to the training records not covered by existing rules.

→ If the ruleset isn't mutually exclusive then a record can be covered by several rules some of which may predict conflicting classes. To overcome this problem there are 2 ways:

1) ordered Rule → An ordered ruleset also known as decision list. When a test record is presented it's classified by the highest ranked rule that covers the record. This avoids the problem of having conflicting classes predicted by multiple  $\text{classif}^{in}$  rules.

2) unordered rules → It allows a test record to trigger multiple  $\text{classif}^{in}$  rules and consider the consequent of each rule as a vote for a particular class. The class label with the highest no. of votes is assigned to test record. Its advantages are:-

a) these are less susceptible to errors caused by ~~the~~ wrong rule being selected to classify a test record.

b) Model building is also & less expensive because the rules don't have to be kept in sorted order.

### Disadvantages

a) Classifying a test record can be an expensive task bcz of test record

# **TutorialsDuniya.com**

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

**Please Share these Notes with your Friends as well**

**facebook**

**WhatsApp** 

**twitter** 

**Telegram** 

must be compared against the precond<sup>n</sup> of every rule in ruleset

b)

#### \* Rule ordering schemes

- 1) Rule based ordering scheme :- This ensures that every test record is classified by the best rule covering it. the drawback of this scheme is that lower ranked rules are much harder to interpret.
- 2) Class based ordering scheme :- Here rules that belong to same class appear together in the ruleset. the relative ordering among the rules from the same class is not important. the problem is a high quality rule may be overlooked in favor of an inferior rule that predicts a higher rank class.

#### \* How to Build a rule based classifier

There are 2 methods for extracting classif<sup>n</sup> rules :-

- 1) Direct method - which extracts classif<sup>n</sup> rules directly from data
- 2) Indirect method which extract classif<sup>n</sup> rules from other classif<sup>n</sup> models such as decision trees

#### Direct Method for Rule Extraction

The sequential covering algorithm is often used to extract rules directly from data in a greedy fashion based on a certain evaluation measure the algorithm extracts the rules one class at a time and the class selected first depends on no. of factors such as class priority or the cost of misclassifying records from a given class.

During rule extraction all the training records for class y are considered to be +ve example while those that belong to other classes are considered to be -ve example. So a rule is desirable if it covers most of the +ve examples and none of the -ve. Once such a rule is found the training records covered by the rule are eliminated and a new rule is added to the ruleset. This procedure is repeated until the stopping criteria is met (i.e. all records are covered).

## NEAREST NEIGHBOUR CLASSIFIER

Some test records mayn't be classified bcz they don't match any training example so we find all the training examples that are relatively similar to the attributes of the test example these are known as nearest neighbours and can be used to determine the class label of the test example. A nearest neighbour classifier represent each e.g. as a data point in D dimensional space where D is the no. of attributes. The data point is classified based on the class label of its neighbours. Where the neighbours having more than 1 label the datapoint is assigned majority class of its nearest neighbour.

If K is chosen to be too small then the nearest neighbour classifier may be susceptible to overfitting bcz of noise in training data. If k is too large the nearest neighbour classifier may misclassify the test instance bcz it's list of nearest neighbours may include the points that are located far away from its neighbourhood.

The algorithm computes the distance b/w each test sample  $z = (x', y')$  and all the training examples  $(x_i, y_i) \in D$ . We need to determine the nearest neighbour list  $D_z$ . Such comput'n can be costly if the no. of training examples is large. Once the nearest neighbour list is obtained the test sample is classified based on the majority class of its nearest neighbour.

$$\text{Majority voting } y' = \arg \max \sum I(v=y_i)$$

$v \rightarrow$  class label  $y_i \rightarrow$  class label for nearest neighbour

For majority voting approach

$I \rightarrow$  Indicator function that will be 1 if the argument is true and 0 otherwise

In majority voting approach each neighbour has same impact on classification which makes algo sensitive to choice of k so we may reduce the impact of k by assigning weights to the nearest neighbour according to its distance

$$w_i = \frac{1}{d(x', x_i)^2}$$

As a result the training examples that are located far away will have a less impact on classification as compared to nearest ones

Majority voting  $y' = \arg \max_{(x_i, y_i) \in D} \sum_i w_i \times I(v=y_i)$

### Characteristics of nearest neighbour classifier

- 1) NNC is an instance based learning which uses specific training instances to make the predictions without having to maintain a model derived from data.
- 2) It's a lazy learner as it doesn't require model building but classifying a test example can be quite expensive because we need to compute the proximity value individually b/w test and training examples.  
Lazy learners often spend most of their computing resources for model building once a model has been constructed classifying a test example is extremely fast.
- 3) These classifiers make their predictions based on local information and are quite susceptible to noise.
- 4) These classifiers can produce wrong predictions if the scale of attributes isn't taken into account.

### Practical Lab Commands

Q5 \*

Implementing decision tree Algo.

library(rpart)

iris <- sample(150, 100)

iris\_train <- iris[8, ]

iris\_test <- iris[-8, ]

dtr <- rpart(species ~ ., iris\_train, method = "class")

rpart.plot(dtr)

lp <- predict(dtr, iris\_test, type = "class")

lp <- iris\_test[, 5]

lp <- as.factor(lp)

Confusion Matrix (as.factor(iris\_test[, 5], species))

not considering species  
→ other attributes

Class → classification

Q) Consider a binary classification problem with following set of attributes and attribute values:-

Air-Cond → Working, broken , Engine → Good or Bad , Mileage → High, med, low  
Rust → yes or no

A rule based classifier produces the following ruleset

Mileage = High → Value = Low

Mileage = Low → Value = High

Air-Cond = Working, Engine = good → Value = High

Air-Cond = Working, Engine = bad → Value = Low

Air-Cond = Broken → Value = Low

a) Are the rules mutually exclusive?

| Ans | Yes | No | Mileage | Air-Cond | Engine | Value    | Rule |
|-----|-----|----|---------|----------|--------|----------|------|
|     |     |    | High    | Working  | Good   | High/Low |      |

b) Are the rules exhaustive?

Ans Yes

c) Is ordering needed for this set of rules?

Ans Yes, as it's not mutually exclusive.

d) Do u need a default class for ruleset?

Ans No, becoz it's an exhaustive ruleset.

Q) For K nearest

|   |     |     |     |     |     |     |     |     |     |     |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| x | 0.5 | 3.0 | 4.5 | 4.6 | 4.9 | 5.2 | 5.3 | 5.5 | 7.0 | 9.5 |
| y | -   | -   | +   | +   | +   | -   | -   | +   | -   | -   |

consider 1D dataset as shown in the table

a) Classify the datapoint  $x=5.0$  according to 1, 3, 5 and 9 nearest neighbours

(using majority vote)

Ans  $1 \rightarrow 4.9 +$   
 $3 \rightarrow 4.9, 5.2, 5.3 -$

$5 \rightarrow +$

$9 \rightarrow -$

b) Repeat the previous analysis using the distance weighted Voting approach.

Ans

1 → +

$$3 \rightarrow \begin{array}{c} \frac{1}{(0.1)^2} (+) \\ \frac{1}{(0.2)^2} (-) \\ \frac{1}{(0.3)^2} (-) \end{array}$$

(100) 36.11

$$5 \rightarrow \begin{array}{ccc} \frac{1}{(0.4)^2} & + & \frac{1}{(0.5)^2} + \end{array} 11 (+)$$

$$9 \rightarrow \begin{array}{ccc} \frac{1}{(4.5)^2} (-) & \frac{1}{(3.0)^2} (-) & \frac{1}{(2)^2} (-) \end{array} (+)$$

### \* BAYESIAN CLASSIFIER

In many applications the class label of a test record can't be predicted with certainty even though its attribute set is identical to some of the training examples this situation may arise because of noisy data or presence of certain factors. Bayes classifier presents an approach for modelling probabilistic relationship b/w the attribute set and class variable.

$P(X=x)$  refers to probability that variable  $X$  will take on the value  $x$ .

Q How we derive the bayes formula?

$$\text{Ans } P(X, Y) = P(Y/X) \cdot P(X) \rightleftharpoons P(X/Y) \cdot P(Y)$$

$$P(Y/X) = \frac{P(X/Y) \cdot P(Y)}{P(X)}$$

Q Consider a football game b/w two teams T0 and T1. Suppose T0 wins 65% of times  $P(T_0) = 0.65$  and T1 wins remaining matches  $P(T_1) = 0.35$ . Among games won by T0 only 30% of them come from playing on T1's field on the other hand 75% of victories for T1 are obtained while playing at home. If T1 is to host next match. Which team will most likely emerge as winner?

Ans

$$P(T_0) = 0.30, \quad P(T_1) = 0.75$$

$$P(T_0) = 0.65 \quad P(T_1) = 0.35$$

$X \rightarrow$  team hosting match  $Y \rightarrow$  Winner of match

$$P(Y=0) = 0.65$$

$$P(Y=1) = 0.35$$

$$P\left(\frac{X=1}{Y=0}\right) = 0.3$$

$$P(X=1|Y=1) = 0.75$$

~~if  $X=1$~~

$$\begin{aligned} P(Y=1|X=1) &= \frac{P(X=1|Y=1) \cdot P(Y=1)}{P(X=1)} \\ &= \frac{P(X=1|Y=1) \cdot P(Y=1)}{P(X=1, Y=1) + P(X=1, Y=0)} \\ &= \frac{0.75 \cdot 0.35}{0.75 \cdot 0.35 + 0.3 \cdot 0.65} \\ &= \frac{0.75 \cdot 0.35}{0.75 \cdot 0.35 + 0.3 \cdot 0.65} \\ &\Rightarrow 0.2625 = 0.5738 \\ &0.4575 \end{aligned}$$

$$P(X=0|Y=1) = 1 - P(Y=1|X=1)$$

Team 1 will win the match as  $P(Y=1|X=1) > P(Y=0|X=0)$

\* Using the Bayes Theorem for Classification

$X \rightarrow$  attribute set  $Y \rightarrow$  class variable

If the class variable has non-deterministic relationship with attribute then we can treat  $X$  and  $Y$  as random variable and capture their relationship probability statistically using  $P(Y|X)$  which is known as posterior probability for  $Y$ .  
 $P(X), P(Y) \rightarrow$  prior

During the training phase we need to learn the posterior probabilities  $P(Y|X)$  for every combination of  $X \times Y$ . based on the information gathered from the training data so that the test record  $X'$  can be classified by finding the class  $Y'$  that maximizes the posterior probability  $P(Y'|X')$

$x = (\text{HomeOwner} = \text{No}, \text{MaritalStatus} = \text{Married}, \text{AnnualIncome} = \$120k)$   
 $P(\text{Yes}/x), P(\text{No}/x)$

Here we'll be using Bayes theorem :-

$$P(y/x) = \frac{P(x/y) \cdot P(y)}{P(x)} \quad P(x) = \alpha$$

$$P(y = \text{yes}) = 0.3 = \frac{3}{10}$$

$$P(y = \text{no}) = 0.7 = \frac{7}{10}$$

Naive ~~baise~~ based classifier estimates the class conditional probability by assuming that the attributes are conditionally independent given the class label.

$$P(x) = P(x|y=y) = \prod_{i=1}^d (P(x_i|y=y))$$

$$P(y/x) = \frac{P(y) \cdot \prod_{i=1}^d P(x_i|y=y)}{P(x)}$$

$$P(x = \text{Homeowner/No} | \text{Yes}) = \frac{3}{3} = 1$$

$$P(x = \text{Married, Homeowner NO} | \text{No}) = \frac{3}{7}$$

$$P(\text{MaritalStatus Married} | \text{No}) = \frac{4}{7}$$

$$P(AI = \$120k) = \frac{1}{7}$$

### \* Estimating Conditional probability for Continuous attribute

1) We can discretize each continuous attribute and then replace the continuous attribute value with its corresponding discrete interval.

$P(x_i|y=y)$  is estimated by the fraction of training records belonging to class  $y$  that falls within interval  $x_i$ .

If the no. of intervals is too large there would be very few training records in each interval and hence will not provide a reliable estimate.

for  $p(x_i | y)$ . If no. of intervals is too small then some intervals may aggregate records from different classes.

- 2) A gaussian distribution is usually chosen to represent the class Conditional probability for continuous attributes.

$$p(x = x_i | y = y_i) = \frac{1}{\sqrt{2\pi} \sigma_{ij}} e^{-\frac{(x_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

$\sigma \rightarrow$  Standard deviation

$\mu \rightarrow$  mean

$\frac{110}{7}$

$$\bar{x} = \frac{125 + 100 + 70 + 120 + 60 + 220 + 75}{7} = \frac{770}{7} = 110$$

$$\sigma = \sqrt{\frac{225 + 100 + 1600 + 100 + 2500 + 12100 + 1225}{6}}$$

$$\sigma = 54.54 \quad \sigma^2 = 2975$$

$$p(x = x_i | y = y_i) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi} \cdot 54.54} e^{-\frac{(120 - 110)^2}{2 \cdot 2975}} = 0.0072$$

$$p(x | NO) = \frac{4}{7} \times \frac{4}{7} \times 0.0072 = 0.0024$$

~~P(x | NO)~~

$$p(y = NO | x) = \frac{p(x | y = NO) \cdot p(y = NO)}{p(x)}$$

$$= \frac{0.0024 \times \frac{7}{10} \times 1}{\alpha}$$

$$p(y = NO | x) = 0$$

so  $x$  classlabel = NO

### \* CLASS IMBALANCE PROBLEM

|                      |   | Predicted |    |
|----------------------|---|-----------|----|
|                      |   | +         | -  |
| (True positive rate) | + | TP        | FN |
|                      | - | FP        | TN |

$$\text{Accuracy} = \frac{TP + TN}{\text{Total}}$$

$$\text{Error} = 1 - \text{Accuracy}$$

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

(γ)

$$TNR = \frac{TN}{TN + FP}$$

$$FNR = \frac{FN}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN} = TPR$$

$$F_1 = \frac{2}{\frac{1}{\gamma} + \frac{1}{P}} \rightarrow \text{Harmonic mean} = \frac{2\gamma P}{\gamma + P}$$

$$\text{Weighted Accuracy} = \frac{w_1 TP + w_4 TN}{w_1 TP + w_2 FP + w_3 FN + w_4 TN}$$

$$F_\beta = \frac{(\beta^2 + 1) TP}{(\beta^2 + 1) TP + \beta^2 FP + FN} \rightarrow \text{used to examine tradeoff b/w precision and recall}$$

|                | w <sub>1</sub> | w <sub>2</sub> | w <sub>3</sub> | w <sub>4</sub> |
|----------------|----------------|----------------|----------------|----------------|
| Recall         | 1              | 0              | 1              | 0              |
| Precision      | 1              | 0.1            | 0              | 0              |
| F <sub>β</sub> | $\beta^2 + 1$  | $\beta^2$      | 1              | 0              |
| Accuracy       | 1              | 1              | 1              | 1              |

ASSOCIATION ANALYSIS

In market basket transactions each row in the table corresponds to a transaction which contains a unique identifier transaction id. and a set of items bought by a given customer. Retailers analyse the data to learn about the purchase behaviour of the customers which can be used in marketing promotions to increase their profits and so on. It's also useful for discovering interesting relationships hidden in large datasets which can be represented in the form of association rules. 2 key issues must be taken care of when applying association analysis to market basket data.

- (1) Discovering patterns from a large transaction dataset can be computationally expensive.
- (2) some of the patterns may be spurious bcz they happen simply by chance.

\* BINARY REPRESENTATION

market basket data can be represented in binary form where each item treated as binary variable whose value is 1 if the item is present in transaction and 0 if not. thus an item is an asymmetric attribute

Let  $I = \{i_1, i_2, i_3, \dots, i_n\}$  be set of all the items in a market basket data and  $T = \{t_1, t_2, \dots, t_m\}$  be the set of all transactions.

Each transaction  $t_i$  contains a subset of items from itemset  $I$ .

If an itemset contains  $K$  items it's called  $K$ -itemset.

→ Transaction width:- A TW is defined as no. of items present in a transaction. A transaction  $t_j$  is said to contain an itemset  $X$  if  $X \subseteq t_j$

→ Support count:- refers to the no. of transactions that contain a particular itemset

$$\sigma(X) = |\{t_i | X \subseteq t_i, t_i \in T\}|$$

→ Association Rule:- It's of the form  $X \rightarrow Y$  where  $X$  and  $Y$  are disjoint item sets ( $X \cap Y = \emptyset$ ). The strength of an association rule can be measured in :-

1) support :- It determines how often a rule is applicable to a given dataset

$$S(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N} \quad C(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

Confidence determines how frequently the items we buy appear in the transaction that contains X.

$$\Sigma \text{Milk, Bread} \rightarrow \Sigma \text{Eggs}$$

$$X \longrightarrow Y$$

| TID | Milk | Bread | Eggs |
|-----|------|-------|------|
| 1   | 1    | 0     | 0    |
| 2   | 0    | 1     | 1    |
| 3   | 1    | 1     | 1    |
| 4   | 1    | 1     | 1    |
| 5   | 1    | 1     | 0    |

$$\text{Support} = \frac{3}{5} = 0.6$$

$$\text{Confidence} = \frac{3}{3} = 1.0$$

A rule that has very low support may occur by chance also it may be uninteresting from a business perspective because it may not be profitable to promote the items that customers seldom buy together. Support is used to eliminate uninteresting rules.

Confidence measures the reliability of the inference made by rule. The higher the confidence more are the chances of Y being present in transactions that contain X.

### \* ASSOCIATION RULE DISCOVERY

Given a set of transactions T find all the rules having the support  $\geq$  Minsup and Confidence  $\geq$  Minconf where Minsup and Minconf are corresponding support and confidence thresholds.

Brute-force approach for finding association rule is to compute the support and confidence for every possible rule. This approach is expensive because there are exponentially many rules that can be extracted from a dataset.

$$R = 3^d - 2^{d+1} + 1$$

Total no. of rules (R)      d (total no. of items)

$$27 - 16 + 1 \\ \Rightarrow 12$$

More than 80% of the rules are discarded after applying Minsup as 20% and Min Conf as 50%. To avoid reduce the no. of computations it would be useful if we can prune the rules before having to compute their support and confidence values. So a common strategy adopted to decompose the problem into 2 major subtasks.

- 1) Frequent itemset generation whose objective is to find all the frequent itemset that satisfy the Minsup threshold.
- 2) Rule generation whose objective is to extract all the high confidence strong rules from the frequent itemsets found in the previous step.  
The 1st step is computationally more expensive than 2nd.

### \* Frequent Itemset Generation

A dataset that contains  $K$  items can generate  $2^K - 1$ . A brute force approach for finding the frequent itemsets is to determine the support count for every candidate item set in the lattice set. Such an approach can be very expensive bcz it requires  $O(NMW)$  comparisons  $\rightarrow$  max no. of transactions,  $N \rightarrow$  total no of itemsets possible ( $2^K - 1$ )  $W \rightarrow$  transactions width.

To reduce the computational complexity 1) first way out is to reduce the no of candidate itemsets it's an effective way to eliminate some of the candidate itemsets without counting their support values.

2) Reduce no. of comparisons instead of matching each candidate itemset against each transaction we can reduce no. of comparisons by using some advanced DS either to store candidate itemsets or to compress the dataset.

### \* Apriori Principle

↳ If an itemset is frequent then all of its subsets must also be frequent.

If an itemset such as (a,b) is infrequent then all of its supersets must also be infrequent. It's based on a key property that the support for an itemset never exceeds the support of its subset. This property known as antimonotone property of the support measure.

e.g - If itemset (a,b) is infrequent then the entire subgraph containing

Supercells of  $(a,b)$  in lattice structure can be formed immediately.

Monotonicity Property:- Let  $I$  be a set of items and  $J = 2^I$  (powerset of  $I$ )

then a measure  $\mu$  is monotone if

$$\forall x, y \in J : (x \subseteq y) \rightarrow f(x) \leq f(y)$$

Antimonotone Property :-  $f(y) \leq f(x)$

Frequent itemset generation using Apriori algorithm.

| Candidate 1 items |       | Candidate 2 items |       | Candidate 3 items |                     |   |
|-------------------|-------|-------------------|-------|-------------------|---------------------|---|
| Item              | Count | Itemset           | Count | Itemset           | Count               |   |
| Beer              | 3     | Beers, Bread      | 2     | x                 | Beer, Diap, Bread   | 3 |
| Bread             | 4     | Beer, Diap        | 3     | x                 | {Bread, Milk, Diap} | 3 |
| Cola              | 2     | x                 |       |                   | {Beers, Diap, Milk} |   |
| Diapers           | 4     | Beer, Milk        | 2     | x                 | {Beer, Bread, Milk} |   |
| Milk              | 4     | Bread, Diap       | 3     | x                 |                     |   |
| Eggs              | 1     | x                 |       |                   |                     |   |
|                   |       | Bread, Milk       | 3     |                   |                     |   |
|                   |       | Diap, Milk        | 3     |                   |                     |   |
|                   |       |                   |       |                   | 6 x 5 x 4           |   |

~~Berry Bread~~  
~~Berry Dip~~  
Berry MILK  
~~Berry Diet~~  
~~Berry MILK~~  
~~Diet MILK~~

$$\text{Using brute force} = \frac{6}{c_1} + \frac{6}{c_2} + \frac{6}{c_3}$$

$$= 6 + 15 + 20 = \boxed{41}$$

$$\text{Using apriori} \sim {}^6C_1 + {}^4C_2 + {}^4C_3 = 6 + 6 + 4 = 16$$

Let  $C_k$  denote the set of candidate- $k$  itemsets and  $f_k$  denote set of frequent  $k$ -itemsets. The algorithm initially makes a single pass over the dataset to determine the support count of each item. Upon completion of this step, frequent one items will be known. The algo. will iteratively generate new candidate- $k$  itemsets using frequent  $k-1$  itemsets found in the previous iteration. The algo. terminates when there are no new frequent itemsets generated.

The frequent itemset generation of Apriori algo. has 2 imp. characteristics:

- 1) Candidate Generation:- At each iteration new candidate itemsets are generated from the frequent itemsets found in previous iteration

2) Candidate Pruning :- In this support for each candidate is counted and tested against the MinSup threshold. The total no. of iterations needed by the algo. is  $k_{max} + 1$

### Requirements for an effective Candidate Generation Procedure

- 1) A candidate itemset is unnecessary if atleast 1 of its subsets is infrequent (according to antimonotone property)
- 2) To ensure completeness the set of candidate itemsets must cover all the frequent itemsets.

$$\forall k : F_k \subseteq C_k$$

- 3) Generation of duplicate candidates leads to wasted computations and thus should be avoided for efficiency reasons.

### Brute Force method Complexity

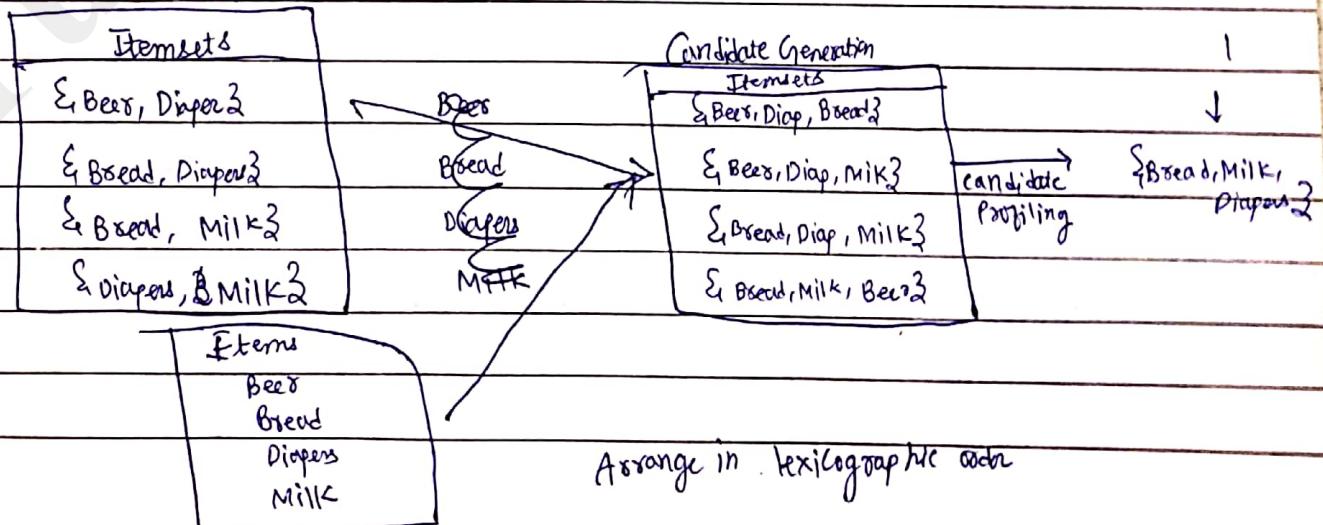
- 1) The no. of candidate itemsets generated at level  $k$  is  $\binom{d}{k}$  where  $d \rightarrow$  total no. of items
- 2) Candidate pruning requires  $O(k)$  time
- 3) Overall Complexity =  $O\left(\sum_{k=1}^d k \cdot \binom{d}{k}\right) = O(d \cdot 2^{d-1})$

### $F_{k-1} \times F_1$ method

It's an alternative method for candidate generation in which we combine frequent  $k-1$  itemsets with other frequent items.

$$\text{Complexity} = O\left(\sum_{k=1}^d |F_{k-1}| \cdot |F_1|\right)$$

This procedure is complete becoz every frequent  $k$  itemset is composed of frequent  $k-1$  and 1 itemset.



# **TutorialsDuniya.com**

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

**Please Share these Notes with your Friends as well**

**facebook**

**WhatsApp** 

**twitter** 

**Telegram** 

### Problem with this approach

(b) This approach doesn't prevent the same candidate itemset from being generated more than once. So the sol' to this is the items in each frequent itemset are kept sorted in their lexicographic order. Each frequent K-1 itemset is then combined with frequent items that are lexicographically larger than the items in the itemset.

### $F_{K-1} * F_{K-1}$ method

candidate generation  
This procedure merges a pair of frequent K-1 itemsets only if their first K-2 items are identical.

$$A = \{a_1, a_2, \dots, a_{K-1}\}$$

$$B = \{b_1, b_2, \dots, b_{K-1}\}$$

$A = \{a_i\}$  A and B be pair of frequent K-1 itemsets then  $A \geq B$  will

be merged if they satisfy following conditions

$$a_i = b_i \quad (i=1, 2, \dots, K-2), \quad a_{K-1} \neq b_{K-1}$$

This method illustrates both the completeness of candidate generation procedure

and advantages of using lexicographic ordering to prevent duplicate candidates

Also an additional candidate pruning step is required to ensure that remaining K-2 subsets of the candidate are frequent.

#### Itemsets

$\{\text{Bread, Diaper}\}$

$\{\text{Bread, Diaper}\}$

$\{\text{Bread, Milk}\}$

!

!

$\{\text{Bread, Diaper, Milk}\}$

Pruning



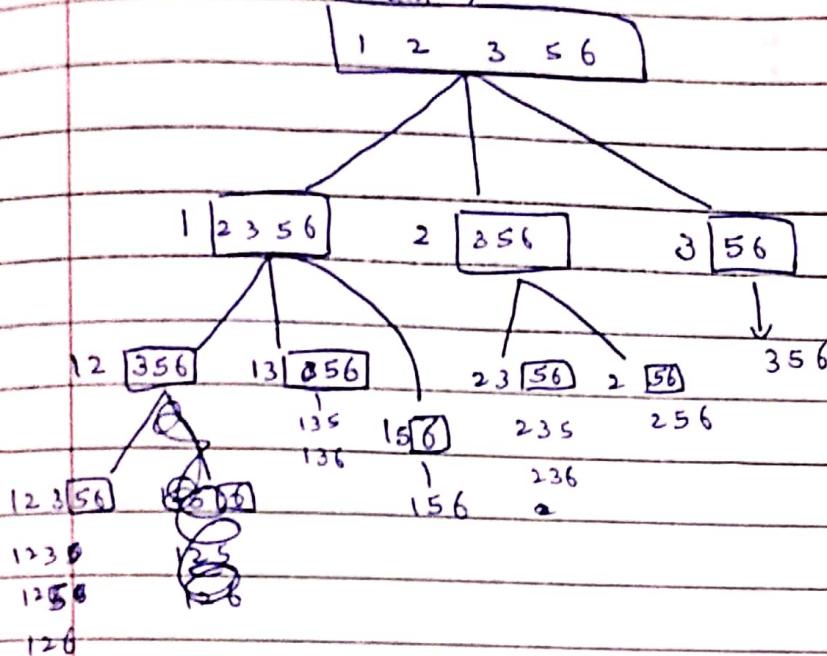
### \* Support Counting

one approach for doing this is to compare each transaction against every candidate itemset and to update the support counts of candidates contained in the transaction

this approach is computationally expensive so we consider an alternative approach in this each itemset keeps its items in increasing lexicographic order

and an itemset can be generated by specifying smallest item first followed by larger items.

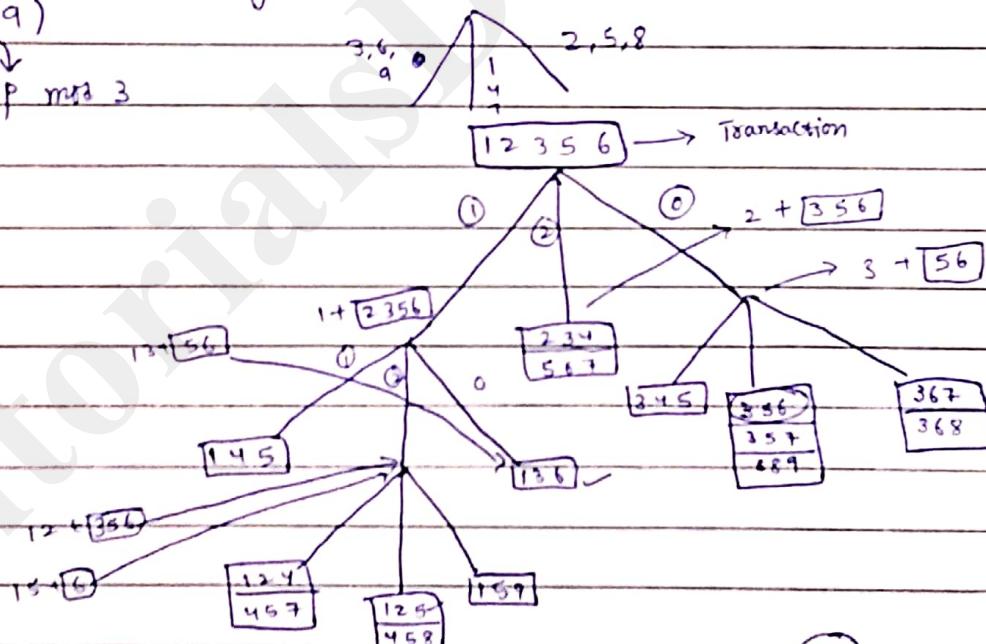
Transaction, t



Here we have to determine whether each generated 3 itemset correspond to an existing candidate itemset if it matches one of the candidates then support count of the corresponding candidate is incremented

### \* Support Counting using Hash Tree

$$(1-9) \\ \downarrow \\ h(p) = p \bmod 3$$



$$3 \rightarrow \{1, 3, 5, 6, 1, 3, 6, 1, 2, 5\}$$

In the apriori algorithm candidate itemsets are partitioned into different buckets and stored in a hash tree during support counting. Itemsets contained in each transaction are also hashed into their appropriate buckets as a result each itemset in the transaction is matched only against candidate itemsets that belong to the same bucket.

Each internal node of the given hash tree uses the following hash function

$$h(p) = p \bmod 3$$

Consider a transaction  $T = \{1, 2, 3, 5, 6\}$  to update the support counts of the candidate itemsets the hash tree must be traversed in such a way that all the leaf nodes containing candidate 3 itemsets belonging to  $T$  must be visited at least once.

Item  $\{1, 2, 3\}$  is hashed to the left child of root node  $\{1\}$ , Items  $\{2, 5, 6\}$  are hashed

$$\{1, 4, 7\}$$

to the middle while items  $\{3, 6, 9\}$  is hashed to the rightmost child.

This process continues until the leaf nodes of hash tree are reached.

The candidate itemsets stored at the ~~reached~~ visited leafnodes are compared against the transaction if ~~the~~ candidate is a subset of transaction its support count is incremented

### (SNNP)

- \* The Computational Complexity of Apriori algorithm can be affected by following factors
  - 1) Support threshold:- Lowering the support threshold results in more itemsets being declared as frequent thus affecting the computational complexity of algo.
  - 2) No. of items:- As the no. of items increases more space will be needed to store the support count of items.
  - 3) No. of transactions:- As the no. of transaction increases the apriori algo makes repeated passes over the dataset.
  - 4) Average transaction width:- The max. size of frequent itemsets tend to increase as the average transaction width increases. As a result more candidate itemsets must be examined during candidate generation and support counting. This will also increase the no. of hash tree traversals performed during support counting.

### Time Complexity

- 1) Generation of Frequent 1 itemset  $\rightarrow \Theta(NW)$  time  
 $N \rightarrow$  total no. of transaction  
 $W \rightarrow$  average transaction width
- 2) Candidate Generation  $\rightarrow \sum_{k=2}^W (k-2) |C_k| < \text{cost of merging} < \sum_{k=2}^W (k-2) |C_{k-1}|^2$   
 To generate Candidate K itemsets pairs of frequent K-1 itemsets are merged to determine whether they have at least K-2 items in common. Worst case scenario when we have to merge every pair of frequent K-1 itemsets.
- 3) A hash tree is also constructed during the candidate generation to store candidate itemsets, so the cost for populating hash tree with max depth K  
 $O\left(\sum_{k=2}^W k |C_k|\right)$
- 4) The cost for looking up a candidate in hash tree is  $O(k)$  and the candidate pruning step =  $O\left(\sum_{k=2}^W k (k-2) |C_k|\right)$

- 5) Support Counting - Cost is  ~~$O(N \sum_{k=2}^W C_k)$~~   $O(N \sum_{k=2}^W C_k)$   
 $w \rightarrow$  max transaction width       $C_k \rightarrow$  cost for updating support count of a candidate K itemset in the hash tree

### \* RULE GENERATION

$$\text{frequent } k \text{ itemset rules} = 2^{k-2} \\ \rightarrow \emptyset \rightarrow X \\ \emptyset \rightarrow ]X$$

An association rule can be extracted by partitioning the itemset Y into 2 non-empty subsets X, Y-X such that  $X \rightarrow Y-X$  satisfies the confidence threshold.

Let  $X = \{1, 2, 3\}$  be a frequent itemset

$$\begin{array}{ll} \{1, 2\} \rightarrow \{3\} & \{1, 3\} \rightarrow \{2\} \\ \{1, 3\} \rightarrow \{2\} & \{2, 3\} \rightarrow \{1\} \\ \{2, 3\} \rightarrow \{1\} & \{3\} \rightarrow \{1, 2\} \end{array}$$

$$\text{Confidence} = \frac{\sigma(\{1, 2, 3\})}{\sigma(\{1, 2\})} = 2$$

Since the support counts for both the itemsets were already found during frequent itemset generation there is no need to traverse to the entire dataset again. Confidence doesn't have any monotone property but it holds the following theorem.

→ If a rule  $x \rightarrow y - x$  doesn't satisfy the confidence threshold then any rule  $x' \rightarrow y - x'$  where  $x' \subseteq x$  must not satisfy the confidence threshold as well.

$$x \rightarrow y - x \rightarrow \frac{\sigma(y)}{\sigma(x)}$$

$$S(x) = \sigma(x)$$

$$x' \rightarrow y - x' \rightarrow \frac{\sigma(y)}{\sigma(x')}$$

$$\sigma(x') \geq \sigma(x)$$

$$\text{So } S(x) \geq S(x')$$

Since  $x' \subseteq x$   $\sigma(x') \geq \sigma(x)$  will hold due to antimonotone property of support  
so  $x' \rightarrow y - x'$  can't have higher confidence than  $x \rightarrow y - x$ .

Apriori algorithm uses a levelwise approach for generating association rules.

Initially all the high confidence rules that have only 1 item in the rule consequent are extracted. Next candidate rules are generated by merging the consequent of 2 rules. If any node in the lattice has low confidence then according to the above theorem the entire subgraph spanned by the node can be pruned immediately.

Suppose the confidence of  $\{a, b, c\} \rightarrow \{a, b\}$  is low. Then all rules containing  $a$  or  $b$  or  $c$  at right are set to 0.

The only difference in rule generation from candidate generation is that we don't make additional passes over dataset to compute the confidence of the candidate rules instead we use the support counts computed during frequent itemset generation.

## CLUSTER ANALYSIS

It divides data into groups or clusters that are meaningful and useful. We do clustering for 2 purposes:-

- 1) Understanding -  
(SICB)  
a) Biology b) Info. Retrieval c) Climate d) Business
- 2) Utilities -  
a) Summarization b) Compression c) Efficiently finding nearest neighbors

Instead of applying the algo. to entire dataset it can be applied only to cluster prototypes if the objects are relatively close to the prototype of their cluster then we can use the prototypes to reduce the no. of distance computations.

Also the nearness of 2 clusters can be measured by finding the distance b/w their prototypes.

\* Clustering :- The diff. ways to group a set of objects into a set of clusters. The greater the similarity within a group and greater the difference b/w the groups the better the clustering is.

### \* DIFFERENT TYPES OF CLUSTERING (From book)

- 1) Hierarchical vs Partitional
- 2) Exclusive vs Overlapping vs Fuzzy
- 3) Complete vs Partial → in case of noise & outliers

In Fuzzy clustering every object belongs to every cluster with a membership weight that is b/w 0 and 1. There's an additional constraint that sum of the weights for each object must equal 1.

(W.P.G.I.S)

### \* DIFFERENT TYPES OF CLUSTERS

- 1) Well separated → well separated boundary b/w them.
- 2) Prototype based → the data objects will be closest to prototype of their clusters as compared to the prototype of any other cluster. Also known as centre based
- 3) Graph based → All the points in a cluster will be closer to each other w.r.t. points of any other cluster.
- 4) Density based → The clusters in which density at center is max and decreases as we move outwards
- 5) Shared property (conceptual clusters) → All objects in a cluster have same property

in case of centroid

## \* K-means clustering Algorithm ( $k$ -medoid)

It's a prototype based clustering technique which defines a prototype in terms of a centroid (the mean or group of points) it's applied to the objects in a cluster.

n-D space.

$k$  medoid defines a prototype which is the most representative point for group of points while a centroid may or may not correspond to an actual data point a medoid must be an actual data point.

### Basic K-Means Algorithm

We first choose  $k$  initial centroids where  $k$  is a user specified parameter (no. of clusters). Each point is then assigned to the closest centroid and each collection of points assigned to a centroid is a cluster. The centroid of each cluster is then updated based on the points assigned to the cluster. We repeat the assignment and update the steps until no point changes the cluster or until the centroids remain the same.

Diagram from book

### → Assigning points to the closest Centroid

To assign a point to the closest centroid Euclidean distance is most often used for the data points in Euclidean space the goal of clustering is expressed by an objective function that depends on the proximities of points to one another or to the cluster centroids for e.g. SSE (sum of squared Errors) we need to minimize the square distance of each point to its closest centroid we will prefer the set of clusters with the smallest SSE meaning that prototypes of this clustering are a better representation of the points in their cluster.

$x \rightarrow$  a data object  $C_i \rightarrow$  ith cluster  $c_i \rightarrow$  centroid of cluster  $C_i$

$m_i \rightarrow$  no. of objects in  $i$ th cluster

$m \rightarrow$  total no. of objects in dataset

$k \rightarrow$  no. of clusters

$$\boxed{SSE = \sum_{i=1}^k \sum_{x \in C_i} \text{dist}(c_i, x)^2}$$

$$c_i = \frac{1}{m} \sum_{x \in C_i} x$$

→ Suppose we have 3, 3D point the coordinates are  $(1, 2, 3), (3, 4, 5), (1, 2, 4)$   
Centroid =  $\left( \frac{1+3+1}{3}, \frac{2+4+2}{3}, \frac{3+5+4}{3} \right) = \left( \frac{5}{3}, \frac{8}{3}, 4 \right)$

If your dataset is document dataset similarity measure is

$$\text{Total cohesion} = \sum_{i=1}^k \sum_{x \in C_i} \text{Cosine}(x, c_i)$$

### \* Choosing initial Centroids

Choosing proper I.C. is the key step of K-Means also a common approach is to choose the initial centroids randomly but the resulting clusters are often poor. So one solution is to perform multiple runs each with a diff. set of randomly chosen initial centroids and then select the set of clusters with a minimum SSE. An optimal clustering will be obtained as long as 2 initial centroids fall anywhere in a pair of clusters since the centroids will redistribute themselves 1 to each cluster. Other effective techniques employed for initialization are:-

- 1) Take a sample of points and cluster them using hierarchical clustering technique the centroids from the clusters formed in this technique are used as initial centroids it works well only if the sample size is small or no. of clusters are small.
- 2) We start with a centroid of all the points then for each successive initial centroid we select the point that's farthest away from any of the initial centroids already selected. problem This can select outliers rather than the points in dense regions  
b) It's expensive to compute the farthest point from the current set of initial centroids so the solution to these problems is  
Instead of applying this approach to all data points it's applied only to sample of points. Since outliers are rare they usually won't be selected as a part of random sample and also the points from very dense region are likely to be included in this sample.

## KNN (K-Nearest Neighbour)

$k=13$ ;  $\#(\text{sqrt}(n))$

normalize ← function( $x$ ) {  
     $\Sigma$

        return  $(x - \text{mean}(x)) / \text{sd}(x)$

}

library(class)

// class library for the kNN

$m1 <- \text{knn}(\text{iris-train}, \text{iris-test}, \text{iris-train[,5]}, k=13)$

confusionMatrix(as.factor(iris-test[,5]), m1)

### Time and Space Complexity

Space =  $O((m+k)*n)$

Time Complexity =  $O(I + K * m + n)$

I → no. of iterations  
K → no. of centroids  
m → no. of data objects  
n → no. of attributes

+

→ Handling Empty Clusters :- If no points are allocated to a cluster during an assignment step then a strategy is needed to choose a replacement centroid.

- 1) Eliminate the point that currently contributes most to the total squared error.
- 2) To choose a replacement centroid from the cluster that has the highest SSE.

⊗

2 → Handling Outliers :- When the outliers are present the SSE will be higher so it's useful to discover outliers and eliminate them beforehand.

- 1) We can keep track of the SSE contributed by each point and eliminate the point with high contributions especially over multiple runs.
- 2) We may want to eliminate small clusters since they frequently represent the group of outliers.

### Reducing the SSE with postprocessing

- 1) One way to reduce the SSE is to find more clusters the strategy is to minimize the SSE contributed by individual clusters since the total SSE is the sum of the SSE contributed by each cluster. We can change the total SSE by either splitting or merging the clusters.

Two strategies to decrease the total SSE by increasing the no. of clusters:-

- a) Split a cluster - the cluster with the largest SSE is usually chosen for splitting.
- b) Introduce a new cluster centroid - 1) A point farthest from any cluster center is chosen.  
2) We can choose randomly from all points or from the points with highest SSE.

Two Strategies that determine the no. of clusters while trying to minimize the increase in total SSE

- 1) Disperse a cluster :- We can remove the centroid that corresponds to the cluster and reassign the points to other clusters.
- 2) Merge 2 clusters :- We can merge the 2 clusters that result in the smallest increase in total SSE. The technique used for this is Centroid method and Ward's Method.

④

#### 4 → Updating the Centroids Incrementally

Instead of updating the cluster centroids after all the points have been assigned to a cluster the centroids can be updated incrementally after each assignment of a point to a cluster.

Advantages :-

- a) This strategy guarantees that empty clusters are not produced because if a cluster has ever had one point that point will be reassigned to the same cluster.
- b) Suppose we use given an arbitrary objective function to measure the goodness of the set of clusters when we predict an individual point we can compute the value of objective function for each possible cluster assignment and then choose the one that optimizes the objective.

Disadvantages :-

- a) The clusters produced may depend on the order in which the points are processed.
- b) Also incremental updates are slightly more expensive.

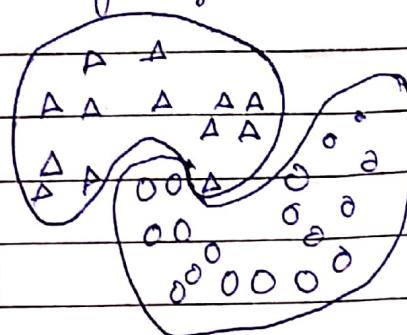
#### Strengths and Weaknesses of K-Means Algorithm

Strength :-

- 1) It's simple and can be used for a wide variety of data types.
- 2) It's quite efficient.

Weaknesses :-

- 1) It's not suitable for all types of data
- 2) It can't handle non-globular clusters or clusters of different sizes and densities.
- 3) It has trouble ~~classifying~~ clustering the data that contain outliers.
- 4) It's generally restricted to data for which we can find a center point.



Non-globular clusters

Q Given the following points  $2, 4, 10, 12, 3, 20, 30, 11, 25$  given  $K = 3$  and initial means  $\mu_1 = 2, \mu_2 = 4, \mu_3 = 6$ . Show the clusters obtained and the new means after each iteration using K-Means algorithm.

Ans 1st iteration

|               |    |                                                     |
|---------------|----|-----------------------------------------------------|
| $d(2, \mu_1)$ | 2  | Closest to $\mu_1$ , so it's assigned to cluster 1. |
| $d(2, \mu_2)$ | 4  | - - $\mu_2$ - - - - - 2.                            |
| $d(2, \mu_3)$ | 10 | - - - $\mu_3$ - - - - - 3.                          |
|               | 12 | - - - - - 3.                                        |
|               | 3  | → $\frac{108}{6}$                                   |

at the end of 1st iteration

$$\begin{aligned} C_1 &= (\cancel{2, 4}) (2, 3) & = (2 \cdot 5) = \mu_1 \\ C_2 &= (4, ) & 4 = \mu_2 \\ C_3 &= (10, 12, 20, 11, 25, 30) & \mu_3 = 18 \end{aligned}$$

In 2nd iteration

$$\begin{aligned} C_1 &= (2, 3) & \mu_1 = 2 \cdot 5 \\ C_2 &= (4, 10, 11) & \mu_2 = \frac{25}{3} = 8 \cdot 3 \\ C_3 &= (12, 20, 25, 30) & \mu_3 = 21 \cdot 75 \end{aligned}$$

In 3rd iteration

$$\begin{aligned} C_1 &= (2, 3, 4) & \mu_1 = 3 \\ C_2 &= (4, 10, 11, 12) & \mu_2 = \cancel{8 \cdot 3} 11 \\ C_3 &= (20, 25, 30) & \mu_3 = 25 \end{aligned}$$

10 Marks

+ Classified  
Page No. \_\_\_\_\_ Date \_\_\_\_\_

In 4th iteration

$$C_1 = (2, 3, 4)$$

$$\mu_1 = 3$$

$$C_2 = (10, 11, 12)$$

$$\mu_2 = 11$$

$$C_3 = (20, 25, 30)$$

$$\mu_3 = 25$$

4th iteration (~~the~~ clusters same as 3rd so final clusters are

$$C_1 = (2, 3, 4)$$

$$C_2 = (10, 11, 12)$$

$$C_3 = (20, 25, 30)$$

### \* EXHAUSTIVE HIERARCHICAL CLUSTERING

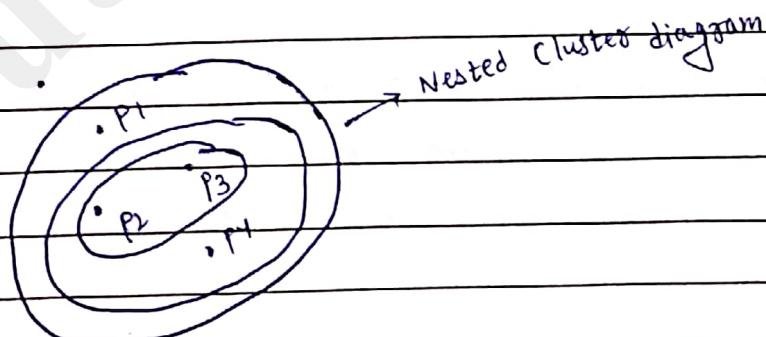
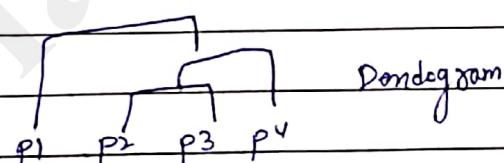
### \* AGGLOMERATIVE HIERARCHICAL CLUSTERING

There are two basic approaches for generating a hierarchical clustering:-

- 1) Agglomerative :- In this we start with the points as individual clusters and at each step merge the closest pair of clusters until one cluster remains.
- 2) Divisive :- In this we start with one all inclusive cluster and at each step split a cluster until only singleton clusters or individual points remain

A hierarchical clustering is often represented graphically using a tree like diagram called dendrogram which displays both the subcluster relationships and the order cluster-

in which the ~~sub~~clusters were merged.



## \* Agglomerative Algo

- 1) Compute proximity matrix.  $O(n^2)$
- 2) repeat  $\xrightarrow{\text{in sorted}} O(n \log n)$   $O(n^3)$
- 3) Merge the closest 2 clusters  $O(m^2)$
- 4) update the proximity matrix to reflect the proximity b/w new cluster and original  $O(m)$
- 5) until only one cluster remains

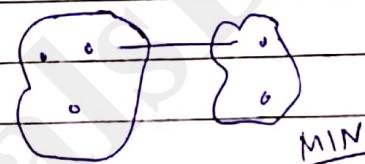
$$\text{Time} = O(n^2) + O(n^3) = O(n^3)$$

If the distance from each cluster to all other clusters are stored as a sorted list it's possible to reduce the cost of finding 2 closest clusters.

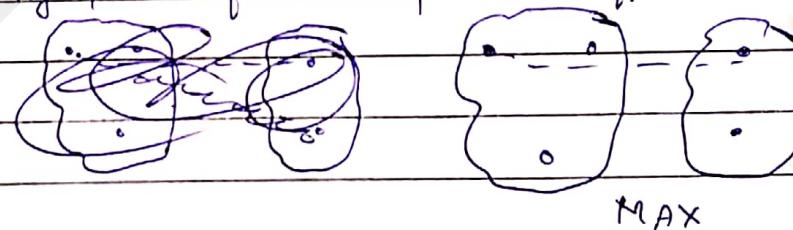
$$\text{So here Time} = O(n^2) + O(n^2 \log n) = O(n^2 \log n)$$

Defining Proximity b/w the Clusters → The key operation is the computation of proximity b/w 2 clusters many agglomerative hierarchical clustering tech such as min, max and group average are used.

MIN defines the cluster proximity as the proximity b/w the closest 2 points that are in different clusters.



MAX takes the proximity b/w the farthest 2 points in different clusters



GROUP AVERAGE define cluster proximity to be the average pairwise proximity of all pairs of points from different clusters. When using centroids the cluster proximity is commonly defined as proximity b/w cluster centroids.

|    | x    | y    |  | P1 | 0.9  | 0.63 |      | P1 | 0.0  | P2   | P3   | P4   | P5   | P6   |
|----|------|------|--|----|------|------|------|----|------|------|------|------|------|------|
| P2 | 0.22 | 0.38 |  | P2 | 0.24 | 0.22 | 0.31 | P2 | 0.24 | 0.0  | 0.15 | 0.2  | 0.29 | 0.23 |
| P3 | 0.35 | 0.32 |  | P3 | 0.29 | 0.0  | 0.2  | P3 | 0.22 | 0.15 | 0.0  | 0.15 | 0.18 | 0.11 |
| P4 | 0.26 | 0.19 |  | P4 | 0.37 | 0.2  | 0.15 | P4 | 0.37 | 0.2  | 0.15 | 0.0  | 0.21 | 0.12 |
| P5 | 0.08 | 0.41 |  | P5 | 0.39 | 0.14 | 0.28 | P5 | 0.39 | 0.11 | 0.05 | 0.0  | 0.19 | 0.19 |
| P6 | 0.45 | 0.30 |  | P6 | 0.23 | 0.25 | 0.11 | P6 | 0.23 | 0.11 | 0.05 | 0.0  | 0.11 | 0.0  |

\* single link (MIN) is good at handling non-elliptical shapes but is sensitive to noise and outliers. The height at which 2 clusters are merged in the dendrogram reflects the distance of <sup>the</sup> 2 clusters  
 $(3,6), (2,5)$

$$\begin{aligned} \text{dist}(\{3,6\}, 1) &= \min(\text{dist}(3,1), \text{dist}(6,1)) \\ &= \min(0.22, 0.23) \\ &= 0.22 \end{aligned}$$

$$\begin{aligned} \text{dist}(\{3,6\}, 2) &= 0.15 \\ \text{dist}(\{3,6\}, 4) &= 0.15 \quad \text{dist}(\{3,6\}, 5) = 0.28 \end{aligned}$$

~~Step 1~~  $\{3,6\}, \{2\}$

$$\begin{aligned} \text{dist}(\{3,6\}, \{2\}, 1) &= 0.22 \quad \text{Next cluster } \{2,5\} \\ \text{dist}(\{3,6\}, \{2\}, 4) &= 0.15 \quad \text{dist}(2,5) = 0.14 \\ \text{dist}(\{3,6\}, \{2\}, 5) &= \end{aligned}$$

$$\text{dist}(\{2,5\}, 1) = 0.24$$

$$\text{dist}(\{2,5\}, \{3,6\}) = 0.15$$

$$\text{dist}(\{2,5\}, 4) = 0.2$$

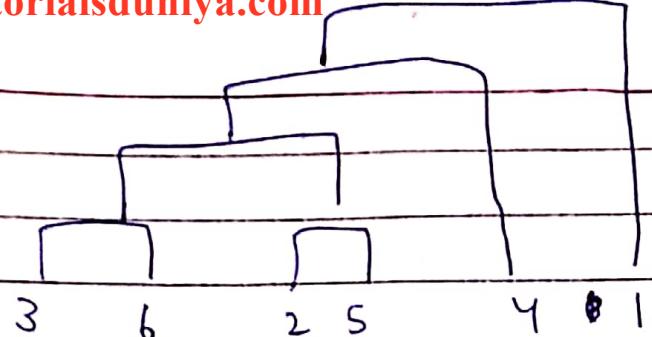
Next cluster  $\{3,6, \{2,5\}\}$

$$\text{dist}(\{3,6, \{2,5\}\}, 1) = 0.22$$

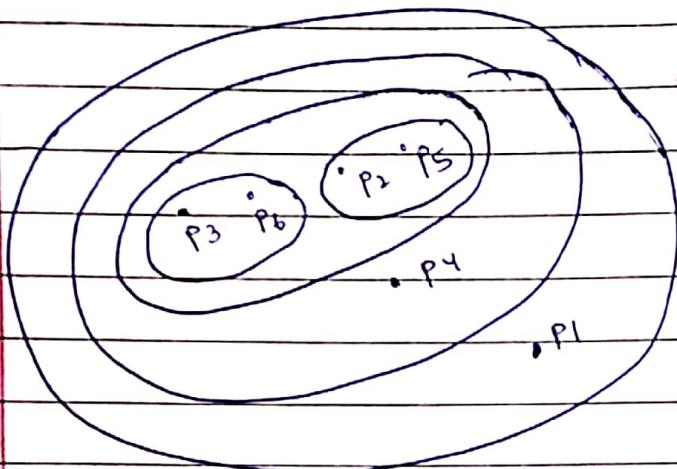
$$\text{dist}(\{3,6, \{2,5\}\}, 4) = 0.15$$

Next cluster =  $\{3,6, \{2,5\}, 4\}$

and then  $\{3,6, \{2,5\}, 4, 1\}$



dendrogram



### Bayesian

```
install.packages("e1071")
```

```
library(e1071)
```

```
model <- naivbayes(Species ~ ., data = iris_train)
```

```
pred <- predict(model, iris_test)
```

```
view(pred)
```

```
confusionMatrix(pred, iris_test$Species)
```

### K-Means

```
iris1 <- iris[, 1:4]
```

```
library(Cluster)
```

```
results <- kmeans(iris1, 3)
```

```
results
```

```
table(
```

Complete Link (MAX)(1)  $\{3, 6\}$ (2)  $\{2, 5\}$ 

$$(3) \text{dist}(\{3, 6\}, 1) = \max(\text{dist}(3, 1), \text{dist}(6, 1)) = \max(0.22, 0.23) = 0.23$$

$$\text{dist}(\{3, 6\}, \{2, 5\}) = \max(0.15, 0.28, 0.25, 0.39) = 0.39$$

$$\text{dist}(\{3, 6\}, 4) = \max(0.15, 0.22) = 0.22$$

(4)  $\{3, 6, 4\}$ 

$$\{\{3, 6\}, 4\} = 0.39$$

$$\{\{3, 6\}, \{4\}, \{1\}\} = 0.37$$

$$\cancel{\{\{3, 6\}, \{4\}\}} + \{\{2, 5\}, 1\} = 0.34$$

(4)  $\{2, 5, 1\}$ (5)  $\{3, 6, 4, 2, 5, 1\}$ 

(LAM)

\* KEY ISSUES IN HIERARCHICAL CLUSTERING

- 1) Lack of a global objective function :- Agglomerative H.C. can't be viewed as globally optimising an objective function. It uses various criteria to decide locally at each step which clusters should be merged. These approaches don't have problem in choosing the initial points.
- 2) Ability to handle different cluster sizes :- one aspect of AHC is how to treat the relative sizes of the pairs of clusters that are merged there are 2 approaches a) Weighted - Which treats all the clusters equally i.e. treating the clusters of unequal size equally gives different weights to the points in different clusters  
b) UnWeighted - Which takes the no. of points in each cluster into account i.e. it gives points in the different clusters the same weight.
- 3) Merging decisions are final :- AHC algorithms tend to make good local decisions about combining the 2 clusters since they use the info about the pairwise similarity of all points only a decision is made to merge 2 clusters it can't be undone at a later time which prevents it from becoming a global

optimization criteria some techs attempt to overcome the limitation by using a partitional clustering technique such as k-Means to create many small clusters and then perform HC using these small clusters as starting point.

### Strengths (Advantages)

- 1) These algs are generally used as they form the clusters using a hierarchy and it's considered a ~~not~~ very good clustering approach.
- 2) These algs can produce better quality clusters

### Disadvantages

- 1) HC algs are expensive in terms of their computational and storage requirement.
- 2) The fact that all merges are final can be troublesome for noisy high dimensional data.

## DBSCAN CLUSTERING ALGO (Density Based Scan)

It isolates the regions of high density that are separated from one another by regions of low density. DBSCAN is a density based alg and is based on centre based approach in this approach density is estimated for a particular point in the dataset by counting the no. of points within a specified radius EPS of that point. This will also include the point itself. The method is simple to implement but the density of any point will depend on the specified radius. The point can be classified as:-

- a) Core point :- It's present <sup>at</sup> interior of a dense region. A point called Core point if the no. of points within a given neighbourhood & specified by the distance parameter EPS exceeds a certain threshold min points (It's also a user specified param).
- b) Border point :- It's located on edge of a dense region. It's not a core point but fall within the neighbours of a core point or core points.
- c) Noise point :- is present in a sparsely occupied region. A noise point is any point that's neither a core point nor border point.

Any 2 core points that are close enough within distance EPS of one another are put in the same cluster.

Any border point that is in the neighbourhood of a core point is put in the

Same cluster as core point.

Noise points are discarded.

Algo from book

\* Time Complexity =  ~~$O(n^2)$~~   $O(m \times \text{time to find pts in Eps neighbourhood})$   
Worst case =  $O(m^2)$

To reduce the time complexity we can use a special type of data structure called kd trees that allow efficient retrieval of points within a given distance of a specified point.

Kd trees - K-dimensional tree just like binary search tree where data in each node is a K-dimensional point in space. It's used for organizing points in K-dimensional space.

Time complexity now =  $O(m \log m)$

\* Space requirement =  $O(m)$  because we only store the points, the cluster label and identification of each point as core, border or noise.

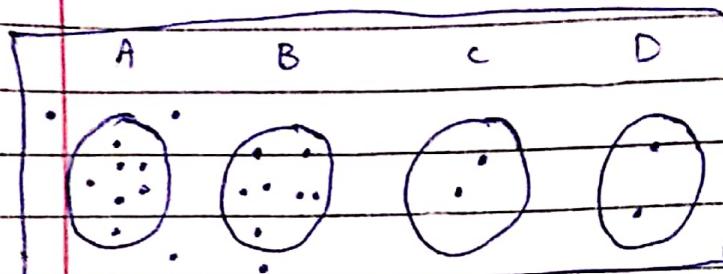
For the points that are not in the cluster the K-distance will be large.

K-distance is distance b/w current point and kth point.

If we compute the K-distance for all the points for some  $k$  sort them in increasing order and then plot the sorted value we'll see a sharp change at a value of K distance that corresponds to a suitable value of EPS.

If we select this distance as EPS parameter and take the value of K as min points parameter then points for which the K-distance is less than EPS will be labelled as core points while others are labelled as noise or border points.

## \* CLUSTERS OF VARYING DENSITY



The noise around a pair of denser clusters A and B has the same density as the clusters C and D if  $\epsilon$  is threshold is low enough <sup>that</sup> DBSCAN finds C and C and D as clusters then the

noise points surrounding A and B will also become part of a cluster.

If  $\epsilon$  threshold is high enough such that the noise points surrounding A and B are marked as noise then points in clusters C and D will also be treated as noise.

### Advantages of DBSCAN

- 1) It's relatively resistant to noise and can handle the clusters of different shapes and sizes.
- 2) It can find many clusters that couldn't be found using K-Means.

### Disadvantages :-

- 1) DBSCAN algo doesn't work when clusters have varying densities.
- 2) It's not suitable for high dimensional data bcz it's difficult to define the density for such points.
- 3) It's quite expensive as the computation of nearest neighbour requires computing all the pairwise proximities.

# **TutorialsDuniya.com**

Download FREE Computer Science Notes, Programs, Projects, Books PDF for any university student of BCA, MCA, B.Sc, B.Tech CSE, M.Sc, M.Tech at <https://www.tutorialsduniya.com>

- Algorithms Notes
- Artificial Intelligence
- Android Programming
- C & C++ Programming
- Combinatorial Optimization
- Computer Graphics
- Computer Networks
- Computer System Architecture
- DBMS & SQL Notes
- Data Analysis & Visualization
- Data Mining
- Data Science
- Data Structures
- Deep Learning
- Digital Image Processing
- Discrete Mathematics
- Information Security
- Internet Technologies
- Java Programming
- JavaScript & jQuery
- Machine Learning
- Microprocessor
- Operating System
- Operational Research
- PHP Notes
- Python Programming
- R Programming
- Software Engineering
- System Programming
- Theory of Computation
- Unix Network Programming
- Web Design & Development

**Please Share these Notes with your Friends as well**

**facebook**

**WhatsApp** 

**twitter** 

**Telegram** 