# Introduction To Data Mining

## (1.1) Motivation & Importance:

The information industry in recent years contains the huge amount of data. This huge amount of data must be transformed into useful information & knowledge. This knowledge & information used in many systems like market analysis, business management, production & analysis, science & engineering etc., This is the motivation behind the data mining technology.

The evolution of database technology contains the functions like data collection & database creation, data management (store, Retrive & transaction process) & data analysis & understanding

The database evolution is shown in below diagram.

```
┌─────────────────────────────────────────────┐
│   Data Collection & Database creation         │
│        (1960 & early)                          │
│     primitive file processing                  │
└─────────────────────────────────────────────┘
                     ↓
┌─────────────────────────────────────────────┐
│       Data Management Systems                  │
│       (1970's - early 1980's)                  │
│   Hierarchical & network database systems      │
│   Relational database systems                  │
│   Data modeling tools: Entity-relationship model etc., │
│   Indexing & Data organisation techniques:     │
│         B⁺ tree, hashing etc.,                 │
│   Query languages: SQL etc.,                   │
│   user interfaces: forms., Reports.            │
│   Transaction Management: Recovery, concurrency │
│                    -control, etc.,             │
│   OLTP systems (Online transaction processing system) │
└─────────────────────────────────────────────┘
        ↓                ↓                ↓
┌──────────────┐  ┌──────────────┐  ┌──────────────┐
│Advanced db   │  │ DW & DM.      │  │ web-base data│
│systems       │  │ (late 1980's- │  │ -base system.│
│(mid 1980's   │  │  pres         │  │ (1990's      │
│ present)     │  │  -ent)        │  │  present)    │
│Advanced data │  │ - DW & OLAP   │  │ - xml databa │
│mod-els:      │  │   tech        │  │ -se systems. │
│object-       │  │   -nology.    │  │              │
│relational etc│  │ - DM & KD     │  │              │
│Application-  │  │               │  │              │
│oriented      │  │               │  │              │
│spatial,      │  │               │  │              │
│temporal,     │  │               │  │              │
│multi-media,  │  │               │  │              │
│scientific    │  │               │  │              │
└──────────────┘  └──────────────┘  └──────────────┘
        ↓                ↓                ↓
      ┌─────────────────────────────────────┐
      │ Integrated Information systems (IIS)  │
      │          (2000 - - - -)               │
      └─────────────────────────────────────┘
```

fig 1.1 The Evolution of database technology

In the year of 1960 & Early data is stored in the form of primitive file processing. This is used as a base of the rest of the database systems.

In the year of 1970, db systems are developed. The initial db systems are Hierarchical & network, later relational db systems, data modeling techniques, indexing & data organisation & query languages are developed. After that to provide Easy understanding to user Several user interfaces, form reports & query processing & query optimization, & transaction management techniques are developed.

In the year of early 1980's OLTP (Online, transaction processing systems are developed). In the year of mid 1980's advanced db systems are developed like advanced data model, application oriented db systems are developed.

In the year of late 1980's Data ware -housing DW has been developed. It is a reposit -ry of hetrogeneous data sources organised und

a schema to support management decision.

Data Warehouse contains the OLAP (Online analytical processing) technology. OLAP contains the functions like summarization & aggregation. Later on data mining & knowledge discovery techniques are developed.

In the year of 1990 web based db systems are developed. By integrating above three systems, we get integrated information systems. This systems developed in the year of 2000.

Therefore, the information industry contains the large amount of data & also data is stored in different sources that exceeds human ability. Therefore, we need data mining tools to analyse & search for interesting patterns or knowledge. This knowledge is used by the management to take Accurate decisions.

dab
that
-en
that
of

clear
Integr
D

## 1.2 What is data mining:

Data mining refers to extracting or mining data from large databases. Some of the people think that data mining is a synonym of KDD knowledge discovery in databases & also some of the people think that data mining as a essential step in the the process of KDD. This is shown in below diagram.



Data mining as a step in the process of knowledge discove

1. Data cleaning:

The data from the large hetrogeneous databases are retrieved & noisy data or inconsistent data is removed before going to the next step.

2. Data Integration:

The data from the different sources are integrated & then this data is loaded into data warehouse

3. Data warehouse: (DWH)

It is a centralized repository it contains all the organizations data.

4. Data Selection:

The selected data is retrieved from the data warehouse.

5. Data Transformation:

Here data is transformed into a form that must be appropriate for the data mining.

6. Data Mining:

It is essential process. It contains intellige -nt techniques or methods like summarization, aggre -gation etc.

7. F

8. kr

to

archi

**7. pattern Evolution:**

Here interesting patterns are Evolutied.

**8. knowledge presentation:**

Here visualization techniques are applied to present mined knowledge to the user.

Based on this the typical data mining architecture contains the following components.



fig 1.3 : Typical data mining Architecture

databases
a is

) are
warehouse

contains

the

form

intellige

aggre

1. Database, DWH or Any other information Repository:

Here data is retrieved from any database, or data warehouse or any other information repository we apply data cleaning & integration tech-nique before data is going to database or DWH server.

2. Database or DWH Server:

It contains all the data according to the user specification.

3. knowledge base:

Here we retrieve the domain knowledge & this domain knowledge is used for searching the interesting patterns

4. Data Mining Engine:

It is a essential process. It contains technique like characterization, classification, association & cluster analysis etc.

5. pattern Evolution:

Here interesting patterns are evoluated.

user
pres
1.3

infc
be
adi
-cc
1.3.

&
ta
q
t
l

fo·
-c
er

**6. Graphical user interface:**

It provides the communication between user & data mining system through this end user can present the queries to datamining system.

## 1.3 Data Mining - on what kind of data:

Data mining can be applicable to any information repository. This information repository may be a relational databases, DWH, transactional database, advanced database systems, & Advanced database applications.

### 1.3.1 Relational Databases:

Relation database system consist of tables & also each table contains the unique name. Each table consist of set of attributes or columns & generally stores large set of data in the form of tuples or records or rows. As well as each row is uniquely accessed by using a key.

**Example:**

For example, relation database tables for ALL Electronics Company contains the tables like customer, item, employee, purchases etc.,

Customer

| Cust-Id | name | address | age | income | .... |
|---------|------|---------|-----|--------|------|
| C1 | Smith | 123 Ham-St canada | 21 | $25,000 | |

Here customer table contains the attributes like cust-Id, name, address, age, income etc.,

Item

| Item-Id | name | brand | category | type | place-made | cost |
|---------|------|-------|----------|------|------------|------|
| I3 | high res-TV | Toshiba | high resolution | TV | Japan | $2000 |

Employee

| emp_id | name | category | group | salary |
|--------|------|----------|-------|--------|
| E5 | John | home, entertainment | manager | $10,000 |

Branch

| Branch-id | name | address |
|-----------|------|---------|
| b1 | city square | 123 main St, Tomato, canada |

purchases

| trans-id | cust-id | item-id | date | time | cost |
|----------|---------|---------|------|------|------|
| T100 | C1 | I3 | 29/06/09 | 15:45 | $100 |

Item_Sold

| Trans_id | Item_id | Quantity |
|----------|---------|----------|
| T100     | I3      | 1        |

fig 1.4 Relational tables for AllElectronics database

Here some of the tables represents the relation between multiple tables like purchases, item_sold etc.,

The Relation database is a most popularly available & rich information repository in datamining In Relation databases data modelling is done by using ER-model

## 1.3.2 Data Ware house (DWH):

It is a centralized repository organised under a schema to support management decision. We load data into 'data warehouse' i.e., the data warehouse is constructed by data cleaning, data transformation, data integration & data load.

for example, typical data warehouse for All Electronics company.

---

tributes

c., f

cost

$2000

ury

000

cost

$100

Data source in chicago

Data source in Newyork

Data source in Tornato

Data Source in Vancouver

clean Transform Integrate load

Data warehouse

Query And Analysis tools.

client

client

Fig 1.5 Typical Architecture of Data Warehouse for All Electronics Company.

Here AllElectronics company has branches All around the world. Therefore data from different branches first of all cleaned, transformed, integrated & finally loaded into data Warehouse, then we apply query & analysis tools before going to present the data to the management.

Data warehouse contains the historical data e·) 5 to 10 years of data & used to support

man
is d
of
this
pre
Exc

ee ALLE
the
valu
( &
\ -0
Phi

management decision. In data warehouse data modeling is done by using schema but the physical architecture of data warehouse contains the "data cube". Through this data cube "multi-dimensional view of data is presented".

Example

Data cube for summarized sales. data of "ALL Electronics" is shown below. This data cube contains the three dimensions i.e., address (it contains city values chicago, Newyork, Toronto, Vancouver), Time (Quater values $Q_1$, $Q_2$, $Q_3$ & $Q_4$) & item (it cont -ains item typer like home entertainment, computer phones, security).

address(cities)
Chicago /440
Newyork /1560
Toronto /395
Vancouver

time(quaters) Q1  605  825  14  400
Q2
Q3
Q4

<Vancouver,
Q1, securi
ty>

home  computer | security
entertainment  phone
item (types)

(b) Drill-down on time data for Q1

address(cities) Chicago
Newyork
Toronto
Vancouver

time( ) Jan 150
feb 100
march 150

address(countries)
USA /2000
Canada /1000

time(quarters) Q1
Q2
Q3
Q4

Rollup on address

fig 1·6  multi-dimensional data cube for All Electronics

1·6 (a) : summarized sales Data

(b) : Drill Down & Roll up operations

Vancouver,
Q1, security

Here data Warehouse contains the data about all the subjects & all the organizations. Therefore, its scope is enterprise - wide. On the other hand datamart, it is a subject of data warehouse. It contains the single subject area & its scope is department - wide.

Data Warehouse provides multi-dimensional view & summarized data. Therefore, data Warehouse is well suited to OLAP (online Analytical processing).

address

OLAP contains the techniques like "drill - down" & "roll up".

In fig 1·6(b) drill - down on the time data for Q1 by monthly & also roll up on address by countries.

Electronics

### 1·3·3 Transactional Databases:

Transactional databases contains the files where each record represents a transaction. A transaction uniquely identified by the transaction-i

trans - id . These transactions are can be stored in tables & each record/row represents transaction. The fragment of transactional database for ALL Electronics is shown in below.

Sales

| Trans - id | list of item_id |
|---|---|
| T100 | $I_1$, $I_3$, $I_8$ |
| . . . . | . . . . |

fig 1.7 Fragment of Transactional database for sales at All Electronics.

Here the sales table contains the nested relation i.e., it contains the list of item-id's

1.3.4. Advanced Database systems & Advanced Database applications:

The new applications like spacial data (such as maps), Engineering & Design (such as designing buildings, Designing components, Integrating circuits), Text & multimedia (such as audio, video, images) Temporal & time related data (such as historical data & stock exchange data) & world wide web (widely distributed information repositories on the internet) etc.,

# Object - oriented database systems:

These object oriented database syste-ms mainly based on object oriented paradigm. Here each entity is treated as object. In this the data & the code related to object are enc-apsulated into single unit.

Each object contains the following.

(i) Set of variables that describes the object:

This is similar to attributes in ER-model

(ii) Set of messages to communicate with one object to another object.

(iii) Set of methods each method contains the code to implement message.

The similar objects are combined into class. This class is known as object class then each object c this class is an instance for that class.

Eg: employee class contains variables like name, address, salary etc.,

# Object - related database systems:

The object related database systems are constructed from object relational data model.

These models are extended from relational model &
also new data types has been added to handle the
complex objects. The constructors also defined to handle
this added data types in relational query languages.

In data mining system the objects
oriented database systems & object related database
systems share some similarties.

Spacial database Systems:

Spacial database systems contains the data
in the form of geographic maps & are represented
in the form of Raster format with n-dimensional
bit maps or pixel maps.

For example, 2D satellite image is represen
-ted in raster format, where each pixel represent
the rain fall in this specific area.

The maps are also represented in
Vector form i.e., where buildings, roads etc., are
represented in the form of points, polyline or
polygon.

## Text & multimedia Database Systems:

Text database systems contains long senten -s, paragraphs to specify product specification or error or bug specification or summary report spe -ification etc., These text database systems are un- structured.

Multimedia databases contains the date in the form of audio, video, images. These are used in applications like voice message systems, video on demand systems, world wide web information systems etc.;

## Temporal & Time - series Database systems:

These two database systems contains the set of data. This data continuously changes with time like stock exchange data.

## World Wide Web information systems:

These are the widely distributed information systems like yahoo or America online or online service. These are linked with objects to exchange the information through internet i.e., the end-

user can get these services through internet & also end user get attractive web pages. These web pages are un-structured i.e. doesn't contain any sche -ma or pattern.

## 1.4 Data Mining functionalities:

Data mining functionalities are used to specify the kind of patterns to be mined found in data mining tasks.

Data mining tasks mainly classified into

1. Descriptive
2. predictive

"Descriptive means it explains the general characteristics of a data in database"

"predictive means it provides conclusion on current data through prediction"

Here end user can go for the search for ultiple patterns because end user does not know

(What kind of patterns are interesting.)

Therefore, data mining system provides the multiple search facility & also provides the interesting patterns are search in various granularities. (levels)

Therefore, datamining system provides the hints to the end user to search for the interesting patterns.

The data mining functionalities mainly classified into

1.4.1 concept/class description: characterisation & Descrimination:

The data is associated with the concept or a class. For example in ALLElectronics database class of items for sale contains the computers & printers.

Similarly concept of customers include big spenders & budget spenders.

Therefore, the data associated with concept or class is called as concept/class description. Thus

descriptions are derived via

    1. Data characterization

    2. Data discrimination

    3. Both Data characterization & Data discri-mination.

## 1. Data characterization:

It is the summarization of general characteristics or features of a target class.

Data characterization can be achieved through several techniques.

for ex, Rollup operation in OLAP.

once the data characterization is completed then data can be presented in various forms "bar charts, pie charts, multi-dimensional datacubes"

## 2. Data Descrimination:

It is the comparison of target class with one or more constructing classes. This is called as the "data descrimination"

Here End user specifies the target classes & constructing classes & also End user can the data from database by using queries.

3. Data characterization & Data Descrimination:

NOTE: we have to write both (1) & (2) here

1.4.2 Association Analysis:

The association analysis is can be used to identify the association rule between each pair of attributes.

This association analysis mainly used in transactional data analysis.

The association rule $x \Rightarrow y$ i.e.,

"$A_1 \wedge A_2 \ldots \ldots A_m \longrightarrow B_1 \wedge B_2 \ldots B_n$" where

$A_i$ for $i \in \{1, \ldots m\}$ and $B_j$ for $j \in \{1, \ldots n\}$ for pair of attributes.

For Example, AllElectronics database contains the association rule

$age(x, ``20\ldots29") \wedge income(x, 20k-30k) \Rightarrow buys($

$``CD\ player")$

[support = 2% confidence = 60%]

As mentioned in the above Ex, a person whas age between 20 to 29 & whose income will be 20k

to 30k can buy or purchase a CD player at
ALLElectronics company. With support 2% & confid
-ence 60%.

(Or)

This association rule indicates that $x \Rightarrow y$ can
be interpreted as customers 2% support are 20
to 29 years of age with income 20k to 30k buys
or purchases a CD player at ALL Electronics company.

There is a probability confidence 60%
of this age group & this income can buy
CD player.

## 1.4.3 Classification & Prediction:

The classification is the process of
finding set of models or functions that describes
class i.e., the classification mainly used to identify
class labels. But many of the applications the end
user search for the data. In most of the cases
usually we search numerical data.

Therefore search for numerical data

in

1.4

the
we
an

Eas
lab
cai
-te

fig 1.8
A
-er
C

in database is called "prediction".

### 1.4.4 Cluster Analysis:

Cluster Analysis is nothing but identifying the data objects with respect to our requirement, we have to also examine the data object which are not relevant to our clan label.

By performing cluster analysis, we can easily identify similarity & dissimilarity of class label. so that we can group the similar objects & we can leave the non-similar objects which is represented in below diagram.
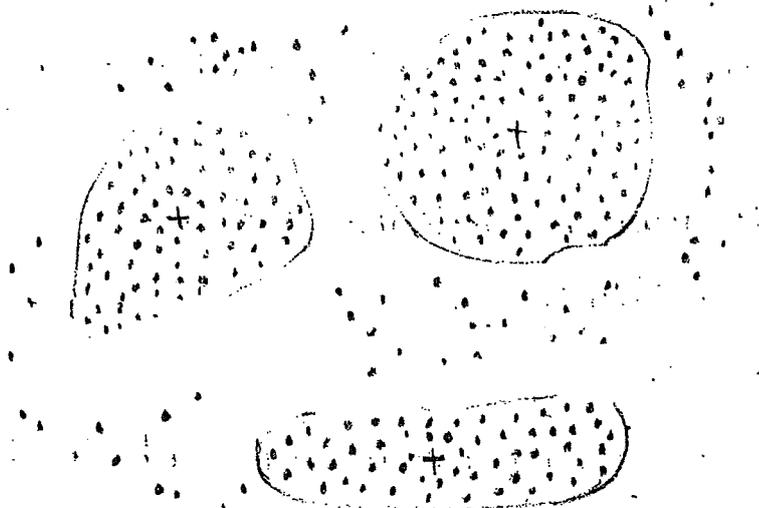
fig 1·8
A 2-D plot of customer data with respect to custor -er locations in a city, showing three data clusters Each cluster "center" is marked with a "+"

1.4.5 Evolution Analysis:

The data evolution analysis describes data objects those are changing with time.

These data objects are evoluted through charaterization, descrimination, association analysis & cluster analysis;

1.5 Are All of the patterns interesting?:

(Or)

Interestingness of a pattern:

The datamining system capable of providing 1000's of patterns. A pattern is interesting if

* That must be understandable to the human beings.

* It must be applicable for new data.

* potentially useful.

* Novel.

Interesting pattern represents the knowledge. Several objective measures for interestingness of a pattern. In these association rule is one of the

method. Therefore association rule $(x \Rightarrow y)$ then support $(x \Rightarrow y) = P(x \cup y)$, where $x \cup y$ indicates that a transaction contains both $x$ & $y$ i.e., the union of item $x$ & item $y$.

Confidence $(x \Rightarrow y) = P(y/x)$ conditional probability i.e., the probability that a transaction containing $x$ also contains $y$.

## 1.6 Classification of data mining systems:

The classification of data mining system is an integration of multiple disciplines (fields). This is shown in below

Database technology → Data mining ← Statistics

Information science → Data mining ← machine Learning

Visualisation → Data mining ← other disciplines

fig 1.9 Classification data mining systems

Here data mining system is integrated with
database Technology, Statistics, machine learning,
visualization, information science etc.,

✓ Therefore data mining system is classified
as following

1. classification According to kind of database it is
mined :

The data mining system is classified based
on database mined. But the database technology is
classified based on data model or type of data
(i.e., application oriented).

If the database Technology is classified
based on data model this may be a relational
database, transactional database, object - oriented
database, object relational or data warehouse.

If the database is classified based
on the type of data then it may be a spacial,
temporal & time series, text & multi-media etc.,

Systems.

2. classification according to the kind of knowledg
it is mined :

Data mining system is classified based on
knowledge mined i.e., we use the data mining functi
-nalities like concept or class description, association
analysis, cluster analysis & evolution analysis.

Data mining system is also classif
-ed based on knowledge mined at different levels.
or top level or primitive level, (low level or
raw level) or multiple levels.

3. Classification according to kind of techniques
utilized :

The data mining system is classified base
on techniques utilized. These techniques describes
the degree of user interaction involved.

For example query-driven systems, autonomo
-us systems, interactive systems etc., or methods
for data analysis.

**4. classification according to kinds of applications adapted :**

The datamining system is classified based on applications adapted like finance data -base system, tele communication database syste -m, share market database system etc.

**★ 1.7 Major issues of datamining :**

These are mainly classified into

1. mining methodology & user interaction issues.

2. performance issues.

3. Issues related to database types.

1. mining methodology & user interaction issues :

These specifies the kinds of knowledge mined, knowledge mined at different level use of domain knowledge & knowledge presentation

Mining knowledge from Databases :
Different users require different

kinds of knowledge. Therefore, data mining system must provide wide range of data analysis & also knowledge discovery through different data mining techniques, like data characterization, data discri-mination, association analysis, cluster analysis & evolution analysis.

## 1.2 mining knowledge at different levels :-

Here end user interact with the data mining system & use the different OLAP tech (online analytical processing) -niques like drill down & roll up, through this techniques end user mined the knowledge at different levels.

## 1.3 Background knowledge :

This background knowledge guides the discovery process if end user has a background know-edge about the database like the constraints or association rules or conditions then this data mini-g process is speed up.

## 1.5 Presentation & Visualization:

The mined knowledge is presented to the end user & this must be understanble to the end user. Therefore, we use the several visualisation techniques like trees, graphs, charts, matrix etc..,

## 1.4 Data mining query languages for datamining Tasks:

Like the relational query language, for example SQL. In SQL we present. Adhoc Query & (temporary query) get the data. similarly data mining query langu -ages has to be developed & it must support end user Adhoc datamining tasks, like data analy - sis & get the domain knowledge, understand the constraints & conditions etc..,

## 1.6 Handling noisy (or) incomplete data:

Generally databases contains the noisy or incomplete data. Therefore, we use the data mining techniques like data cleaning &

data analysis to avoid this noisy or incomplet
-ed data.

## 1.7 Pattern Evolution:

Still the data mining system uncover
thousands of patterns & also many of the discov
-everd patterns uninteresting to the user. Theref
-re, End user has to specify the measures or cons
-traints to reduce the search criteria.

## 2. Performance Issues:

These issues contains the efficiency,
scalability & parallelization of data mining
algorithms.

## 2.1 Efficiency & scalability of datamining Algorithms:

We have to retrieve the data from
large databases. Therefore, the data mining algo
-rithms must be efficient & scalable.

In datamining system knowledge
discovery, efficiency & scalability are the mai

---

ed to

the

ralization
etc.,

ning

, for
query)
rery &
y langu
support
a analy
stand the

x noisy

the

&

key terms.

## 2.2 Parallel, distributed & Incremental data mining algorithms :

End user has to access the data parallely from different data sources & also data is distributed to different data sources. These two are done by using parallel & distributed datamining algorithms. In these algorithms data is divided into partitioned then partition processed paralley & then results of this partitions are merged.

The incremental datamining algorithms are mainly used to reduce the cost of data mining process & also update the data in database without performing the search from the begining of DataBase.

## 3. Issues related to database types:

## 3.1 Handling relational & complex data:

The relational data is handled

efficiently by using relational database syste
data warehouse systems.

The complex data like spatial data, tex
& multi-media data, time series data. These
complex data is handled by using systems spac
-al db systems, text & multi-media systems,
time series db systems etc.,

## 3.2 Mining data from heterogeneous databases

Mining the data like LAN systems,
WAN systems, distributed systems & Hetrogeneous
database systems is only possible through data minin
System. This data mining system also provides &
improves the information exchange & interoperatabi
-ty in hetrogeneous databases.

# Data Warehouse And OLAP Technology for Data Mining

## Basic Concepts of data warehouse :-

(i) What is Data Warehouse

According to warehouse Inman, "A Data
-Warehouse is a Subject Oriented, integrated, non
-Volatile, Time-Variant collection of data, to support
management decisions".

## Subject Oriented:

For example, consider the retailer, these retailer contains the Database systems like Retail Database System, Catalog Database system and outlet Database Systems. These database systems individually support for different queries. But a user want to run a query in all the Sales. This is only possible through data -Warehouse. therefore Data warehouse Organizes Subject areas like customer, products, suppliers, Sales etc. Here the subject area is "Sales". This is shown in the following diagram.
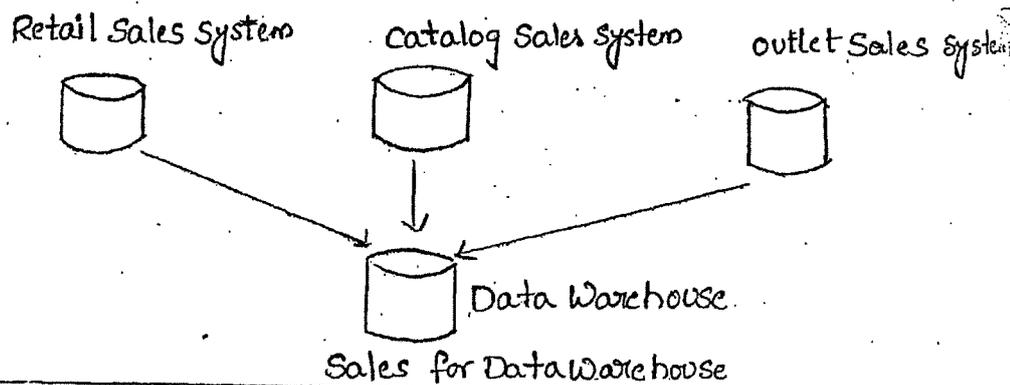
Retail Sales System    Catalog Sales System    outlet Sales System

Data Warehouse.

Sales for Data Warehouse

Fig (2.1.1) Subj-Oriented Sales Information

# Integrated :

Here, different data Sources data are integrated. Then this data loaded it into Data -warehouse. These Sources may be relational databases, flat files, Excel Sheets, Online Transaction records, etc. Before the data is loaded it into DWH, we apply data cleaning and data analysis.

In the above example, 3 data Source, data are integrated then we get the unique key. Using this unique key, we uniquely identify the data record in dataware. This is shown in the following diagram.

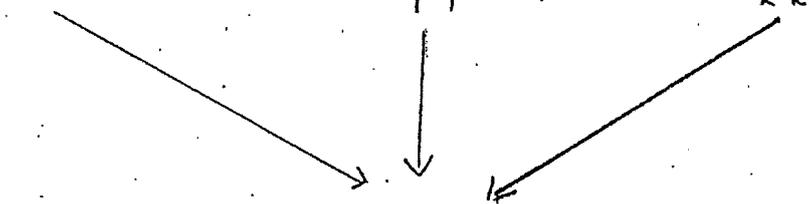Retail Sales System     Catalog Sales system     outlet Sales - System

Product - code:     Product - code:     Product - code:

    XX          YY          ZZ

Product - code for DWH:

A1

fig (2.1.2). Integrated information for Sales

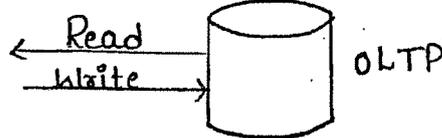## Non-Volatile:

The non-Volatile means read only. The DWH User always read the data. But the OLTP User can read the data as well as write the data. This is shown in below.
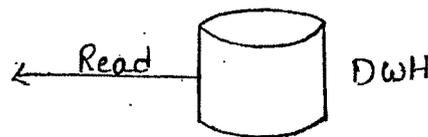


fig (2.1.3) Non-Volatile

## Time - Variant:

The DWH contains the historical data i.e; 5 to 10 years of data. But the OLTP System contains the current data.

```
                    ┌──┬──┬──┐
                    │Jul 2011 │
                    ├──┼──┼──┤
                    │  │  │  │
                    └──┴──┴──┘

┌──┬──┬──┐  ┌──┬──┬──┐  ┌──┬──┬──┐  ┌──┬──┬──┐  ┌──┬──┬──┐
│ 2007  │  │ 2008  │  │ 2009  │  │ 2010 │  │ 2011 │
├──┼──┼──┤  ├──┼──┼──┤  ├──┼──┼──┤  ├──┼──┼──┤  ├──┼──┼──┤
│  │  │  │  │  │  │  │  │  │  │  │  │  │  │  │  │  │  │  │
└──┴──┴──┘  └──┴──┴──┘  └──┴──┴──┘  └──┴──┴──┘  └──┴──┴──┘
```
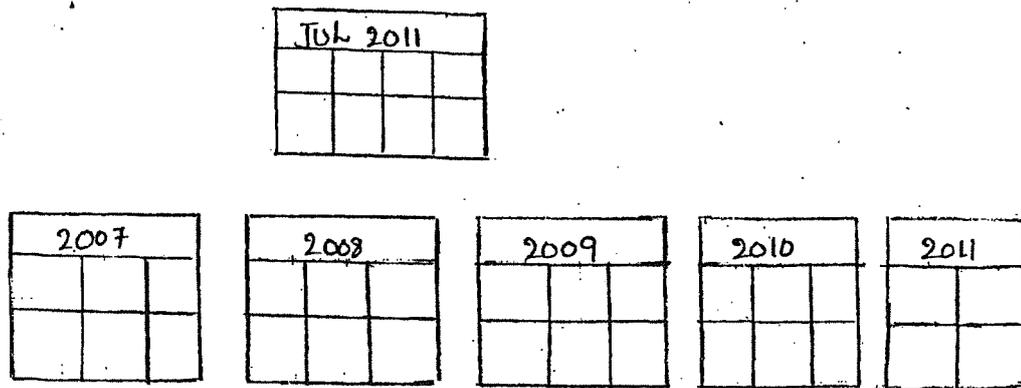
fig (2.1.4) . Time Variant

Finally we extract the data from the DWH & then we apply the different visualization Techniques like Trees, Graphs, Charts, Tables, etc. to present the data to management. The management using this data they take the accurate decision.

## 2.1.1 Differences between operational database systems & Data warehouses:-

The online operational database systems mainly used for online transactions & query processing. These include the transactions like sales, payroll, marketing, manufacturing, etc. These type of systems are called as "Online Transaction Processing" (OLTP). systems.

These systems mainly used for day to day transactions.

The DWH contains the historical data. This is used for data modelling and data analysis. This is used for decision making. This type of systems are called as "Online Analytical Processing" (OLAP) systems.

### Advantages of Data Warehouse:-

1. It makes the data permenant.
2. It makes the data accessable.
3. It identifies hidden business operations.
4. It improves the customer Relationship.
5. It provides the security.

## Differences between OLTP and OLAP :-

### 1. Users and System Orientation :

The OLTP system is customer orientation to support the Online Transaction.

The OLAP system is market orientation to support decision making.

### 2. Data :

The OLTP systems contains the current data. This data does not be used for decision making.

The OLAP system contains the historical data and supports functions like Summerization & aggregation and also data stores and manages at different levels.

### 3. Database design :

In OLTP Systems data-modelling is done by using ER-model.

In OLAP systems data modelling is done by Star (o) Snowflake scheme.

4. View: Using OLTP Systems we extract the data within a Enterprise & department.

In OLAP systems we acass the data from different organization's & different data sources.

5. Access pattern: The OLTP systems Contains the data. This data used for Online Transac -tion's..

The OLAP systems Contains the historical data This data is used for data Analysis.

DWH ⟹ OCAP
OpDBS → OLTP.

Users and System orientation

Data:
Database design

View
Acces patter.

## Comparision between OLTP and OLAP :

| Feature | OLTP | OLAP |
|---|---|---|
| Characteristic | Transactional - Processing | Informational - processing. |
| Users | Clerk, DBA, DB - Professional | HR, manager, Analyst, executive. |
| Functions | Day - to - Day Transa -ctions. | Decision making. |
| Orientation | Online Transaction - Orientation | Analysis Orientation |
| Data | current data | Historical data |
| Database Design | ER - model | star (☆) Snowflake Schema |
| View | Relational (2-D) | Multi Dimensional. |
| Units of work | Simple Transactions | Complex Queries. |
| Acess pattern | Read / write | Read only |
| number of Users | thousands | Hundreds |
| Number of data - records accessed | Tens | millions |
| Data storage | 100MB to GB | 100 GB - TB |

fig (2.1.5). Comparision between OLTP & OLAP

# Data Warehouse Modeling : Data cube & OLAP

## 2.2. Multi-Dimensional Data model:

The Data Warehouse (or) OLAP tools contains the data in the form of multi-dimensional model. i.e; It contains the Data cube.

## 2.2.1. From Tables and Excel Sheets to Data cubes:

• Using Data cube we can view the data (or) analyze the data in the form of multi-dimensional model. It is defined by Dimensions and Facts.

The Dimension is nothing but the Entity. Using this Organizations store the data. For Example, Sales DWH contains the dimensions like Time, Item, location, supplier etc. Using this dimensions we can analyze the things like monthly sales of a item, branches & locations at which the item were sold.

This Dimension information is stored in a Table. This is called as "Dimension Table". Each Dimension Table contains the Set of Attributes. For Example, item - Dimension contains the attributes

like item-Id, name, category, brand, type, etc.

The multi-Dimensional data model the entire data typically organized under the central theme like sales. This "Central Theme" is called "Fact Table".

The facts are numerical entities (or) numerical measures like dollars-sold (sales amount in dollars), units-sold (total units sold).

For example, consider the table for sales DWH for All electronic's company is shown in below. It contains the dimensional's like Time, items, and location.

| Location = "Vancouver" | | | | |
|---|---|---|---|---|
| Time (Quarter) | Item (types) | | | |
| | Home Entertainme -nt | Computers | phones | Security |
| Q1 | 600 | 700 | 500 | 400 |
| Q2 | 700 | 800 | 600 | 700 |
| Q3 | 500 | 600 | 400 | 700 |
| Q4 | 600 | 700 | 800 | 900 |

fig (2·2·1·1). Table for sales DWH, it contains Dimensions, time, item & location

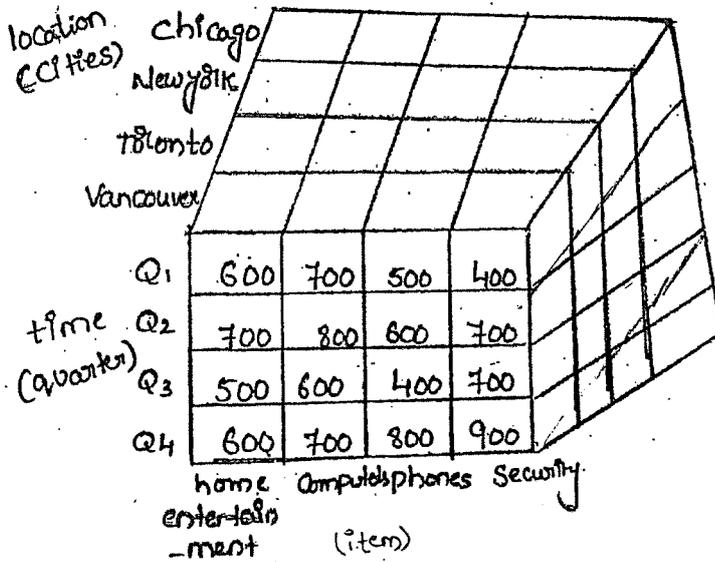The table information is represented in 3D-cube the 3D-cube is shown in below.



| location (Cities) | Chicago | NewYork | Toronto | Vancouver |
|---|---|---|---|---|

time (quarter)

| | home entertain-ment | computers/phones | | security |
|---|---|---|---|---|
| Q1 | 600 | 700 | 500 | 400 |
| Q2 | 700 | 800 | 600 | 700 |
| Q3 | 500 | 600 | 400 | 700 |
| Q4 | 600 | 700 | 800 | 900 |

home entertain-ment    computers/phones   Security

(item)

fig (2.2.1.2): A 3-D data cube for sales it contains dimensions time, items & locations.

The Datacube contains the $n$-dimensions.

The above Data cube may contains the set of dimensions. Using this set of Dimensions we can construct the lattice relation of cuboids. Each cuboid gives the different level of summerization. once the cuboids are completed we make the data cube.

This is shown in the following diagram.

Fig (2.2.1.3): A lattice of cuboids makes 4D Data cubes. It Contains Dimensions time, item, location & suppliers.

In the above cuboid, 4D cuboid is called as The "Base cuboid". it gives the lowest level of Summerization. The 0-D cuboid is called as "apex" cuboid. It gives the highest level of Summerization. It is represented by the keyword "all".

## 2.2.2. Star, Snowflake and fact constellation Schemas for multidimensional databases:-

In Relational database the data is loaded by Using entity Relationship model. This is best Suitable for the online Transactions. But the Data warehouse Contains the subject areas like customers, products, Sales, Suppliers, etc. This is Used for data analysis.

Therefore to load the data into DWH we Use multi-dimensional Data model. Such model exists (or) form of Star Schema (or) Snowflake Schema (or) fact constellation schema.

### Star Schema:

That is the most popular data model to load the data into DWH. It contains 2 tables.

1. Fact Table ( It contains large amount of data without any duplication).

2. Dimension Tables ( one Table for each dimension

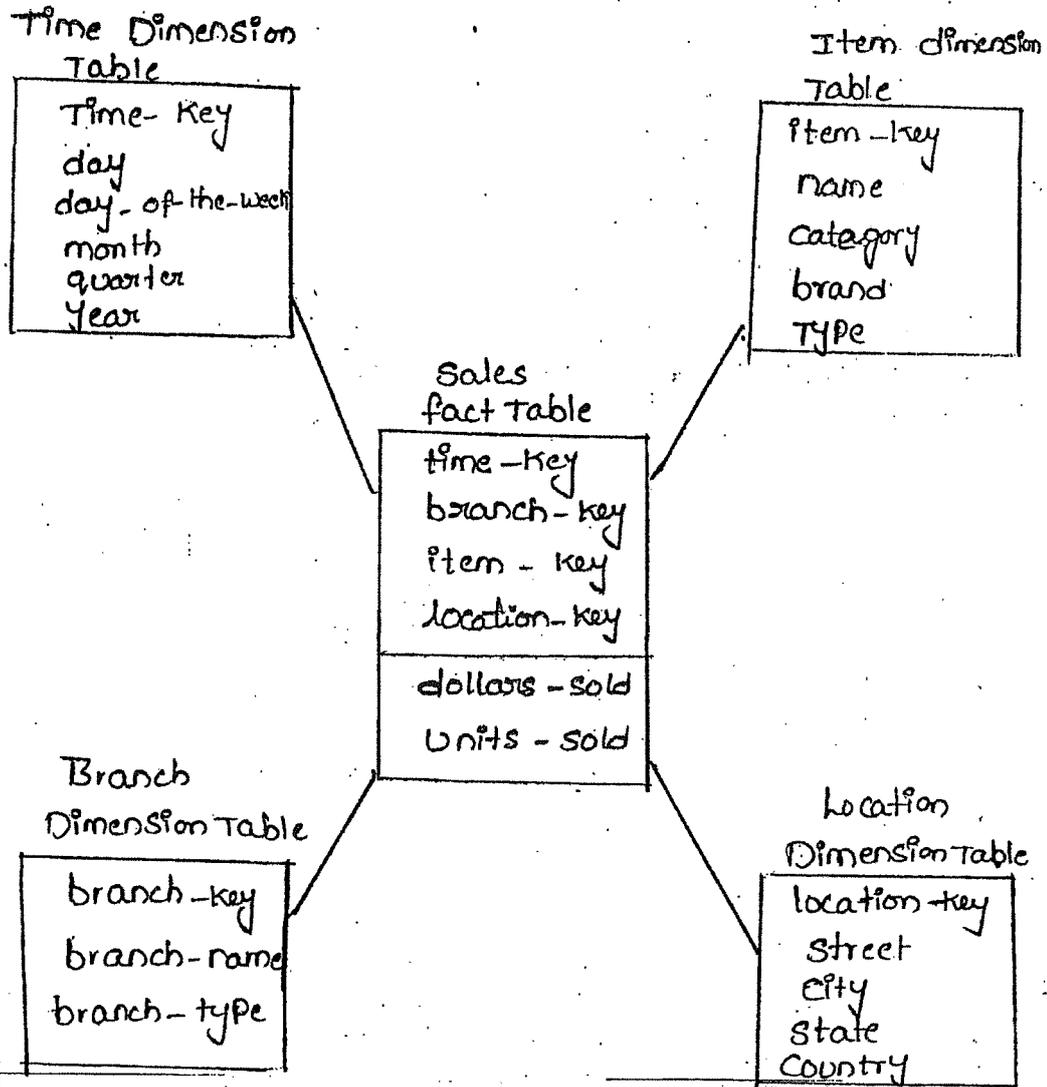This Schema look like the star. Because of that it is named as "star schema". This is shown in below.

Time Dimension Table

| Time - Key |
| --- |
| day |
| day - of-the-week |
| month |
| quarter |
| Year |

Item dimension Table

| Item - key |
| --- |
| name |
| category |
| brand |
| Type |

Sales fact Table

| time - key |
| --- |
| branch - key |
| item - key |
| location - key |
| dollars - sold |
| units - sold |

Branch Dimension Table

| branch - key |
| --- |
| branch - name |
| branch - type |

Location Dimension Table

| location - key |
| --- |
| street |
| city |
| state |
| country |

fig (2-2-2-1): Star Schema for All Electronics Sales DWH

Star Schema for fact Table contains 2 Parts.

(1) Key for Each dimension Table.

(2) Measures that is to be analyzed.
     i.e; like dollars - sold, units - sold.

# Snowflake Schema:

In Star Schema the dimensions tables are not normalized. But in snowflake schema the dimension tables are normalized. i.e; The tables are further splitted it into more tables. This is shown in below.
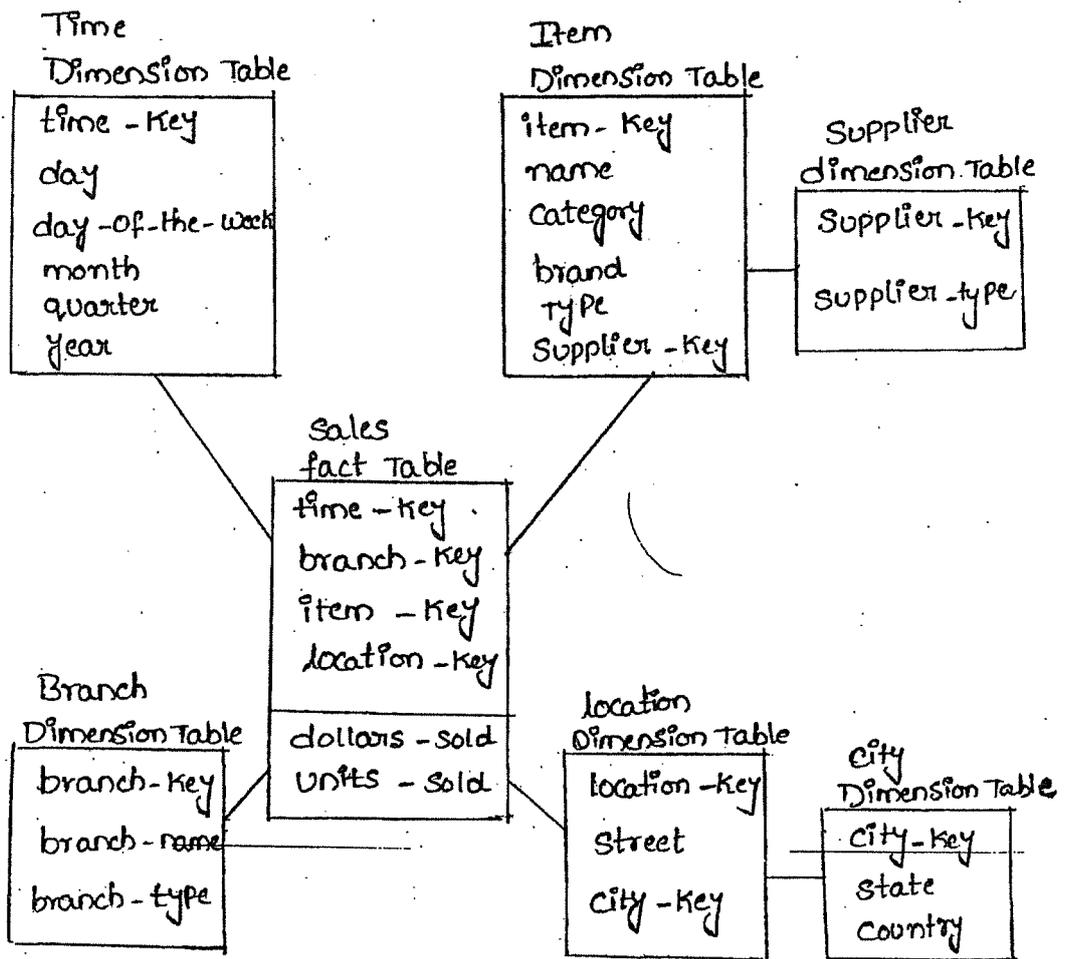
Time
Dimension Table

| time - Key |
|---|
| day |
| day -of-the- week |
| month |
| quarter |
| year |

Item
Dimension Table

| item - Key |
|---|
| name |
| Category |
| brand |
| Type |
| Supplier - Key |

Supplier
dimension Table

| Supplier - Key |
|---|
| supplier - type |

Sales
fact Table

| time - Key |
|---|
| branch - Key |
| item - Key |
| location -Key |
| dollars - Sold |
| units - Sold |

Branch
Dimension Table

| branch - Key |
|---|
| branch - name |
| branch - type |

location
Dimension Table

| location - Key |
|---|
| Street |
| city - Key |

City
Dimension Table

| city - Key |
|---|
| state |
| country |

fig (2.2.2.2) Snowflake Schema for All Electronics Sales DWH

# Fact Constellation Schema:

In some of the applications it requires the multiple fact tables to represent the Dimension Tables. Such a Schema is called as the "fact Conste -llation schema".

The DWH contains the information about the Entire the organization .... Its scope is enterprise wide. But the Data mart is subset of DWH and it contains the single subject area & it's scope is department wide. If it is a DWH (or) Data mart. It contains 2 Definition's.

## (1) Cube Definition:

Syn: define cube ∠ cube-name > [∠ dimension - list>]:
∠ measures - list>

## (2) Dimension Definition:

Syn: define dimension ∠ dimension - name> as ∠ attributes - list>

## 2.2.3. Measures:—

These are the numeric values. that is to be analyzed. we find the measure value at any time

by aggregating the data with corresponding dimensions.

The measures are mainly classified into 3.

(1) **Distributive:**

If it is a distributive measure, then it contains the function's like sum(), max(), min(), count(), etc...

(2) **Algebric:** If the measure is algebric · means it contains the function's like avg(). The average is

$$avg() = \frac{sum()}{count()}$$

Here sum() & count() are distributive measures.

(3) **Holistic:**

If the measure is holistic means it contains the functions like mean(), median(), mode().

### 2.2.4. Introduction to Concept Hierarchies:-

The concept hierarchies defines sequence of mappings from lowest level to highest level. The concept hierarchy for the location is shown in the following diagram.

location

all



countries

states

cities

```
                        ┌─────────┐
                        │   all   │
                        └─────────┘
                       /           \
            ┌──────────┐           ┌──────────┐
            │  Canada  │           │   USA    │
            └──────────┘           └──────────┘
             /        \             /         \
    ┌──────────┐  ┌──────────┐  ┌──────────┐  ┌──────────┐
    │ Colombia │..│  Antaria │  │ New York │..│  Chicago │
    └──────────┘  └──────────┘  └──────────┘  └──────────┘
         │              │            │             │
    ┌──────────┐  ┌──────────┐  ┌──────────┐  ┌──────────┐
    │Vancouver │..│  Toronto │  │ New York │..│  Chicago │
    └──────────┘  └──────────┘  └──────────┘  └──────────┘
```

fig (2.2.4.1): Concept hierarchies for location Dimension

The concept of hierarchy for location dimension contains the attributes street, city, state and country. Using this attributes we can define the concept hierarchy.

" Street $<$ city $<$ State $<$ country ".

These attributes are organized in partial order then we get the "lattice". The lattice for location Dimension.



country
State
city
street

fig (2.2.4.2). Lattice for location Dimension

We can also define the lattice for time Dimension.



fig (2.2.4.3). Lattice for Time Dimension

The concept hierarchy is also defined for grouping of values then it is called as "Set - grouping hierarchy". The Set Grouping hierarchy for price Dimension is shown in below.
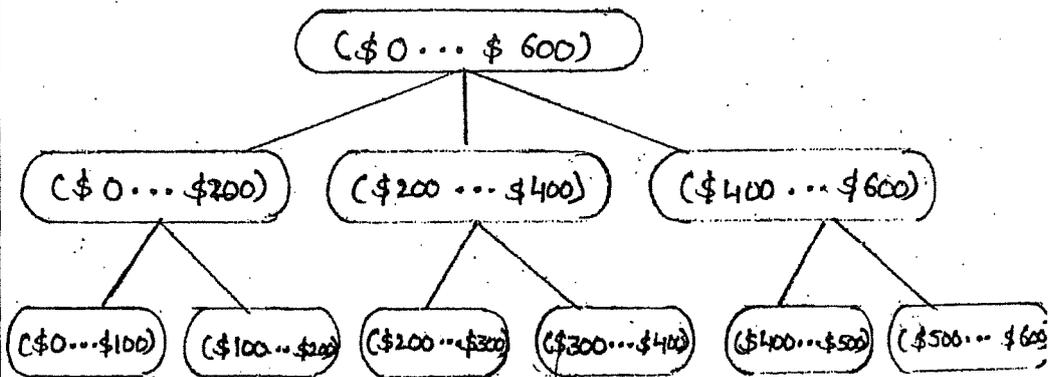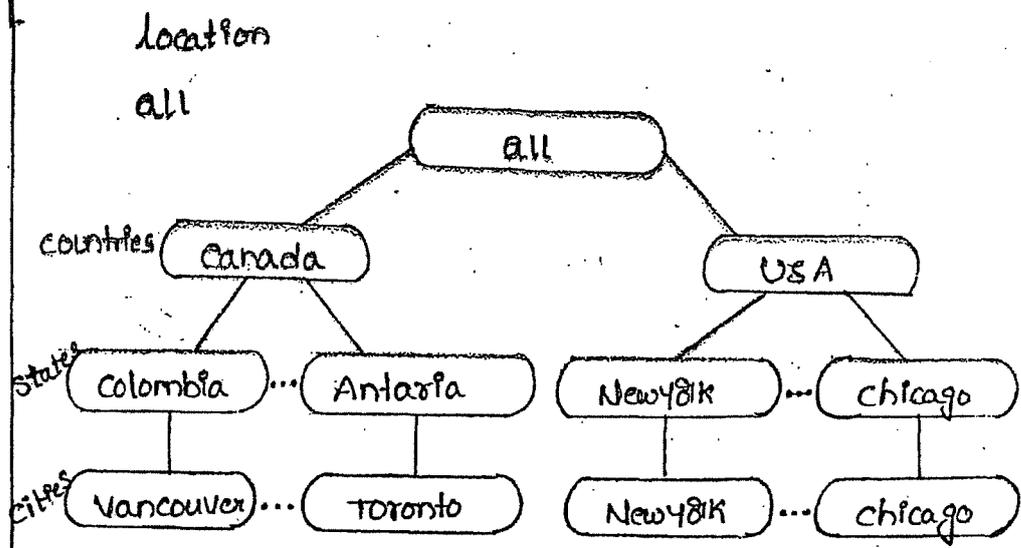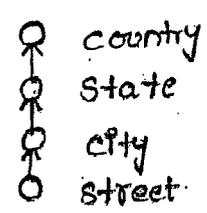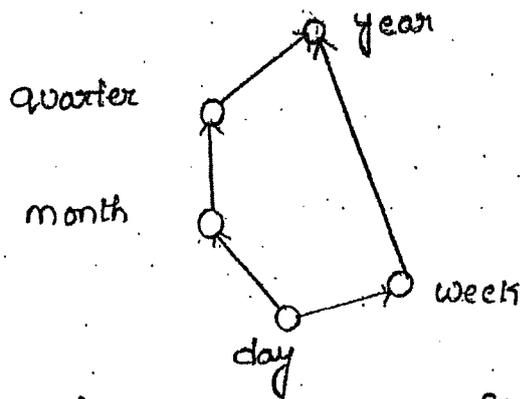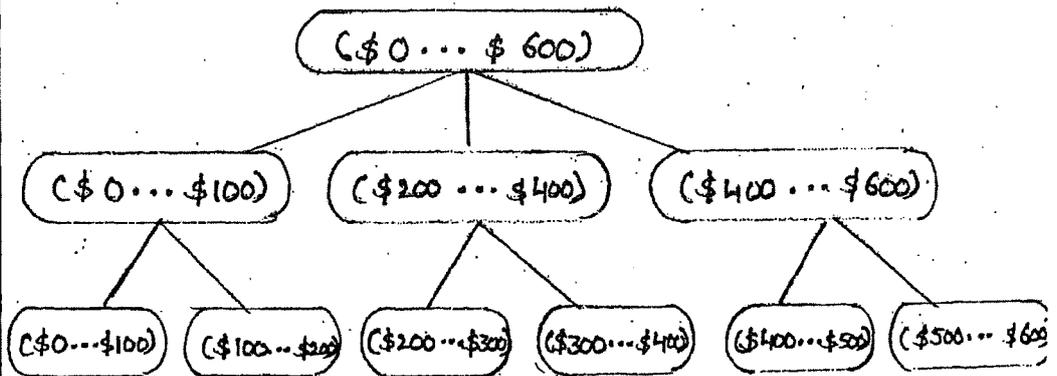


fig (2.2.4.4). Set Grouping hierarchy for price dimension

location

all



fig (2·2·4·1): Concept hierarchies for location Dimension

The concept of hierarchy for location dimension contains the attributes street, city, state and country. Using this attributes we can define the concept hierarchy.

" Street < city < state < country ".

These attributes are organized in partial order then we get the "lattice". The lattice for location Dimension.



fig (2·2·4·2). Lattice for location Dimension

we can also define the lattice for time Dimension.



fig (2·2·4·3) Lattice for Time Dimension

The concept hierarchy is also defined for grouping of values then it is called as "Set - grouping hierarchy". The Set Grouping hierarchy for price Dimension is shown in below.



fig (2·2·4·4). Set Grouping hierarchy for price dimension

## 2.2.5 OLAP operations in multi Dimensional Data Model :-

The multi Dimensional Data model allows the data is organized in multiple Dimensions and also Each dimension contains the set of levels as defined in concept hierarchy. The OLAP operations are shown in the following diagram.
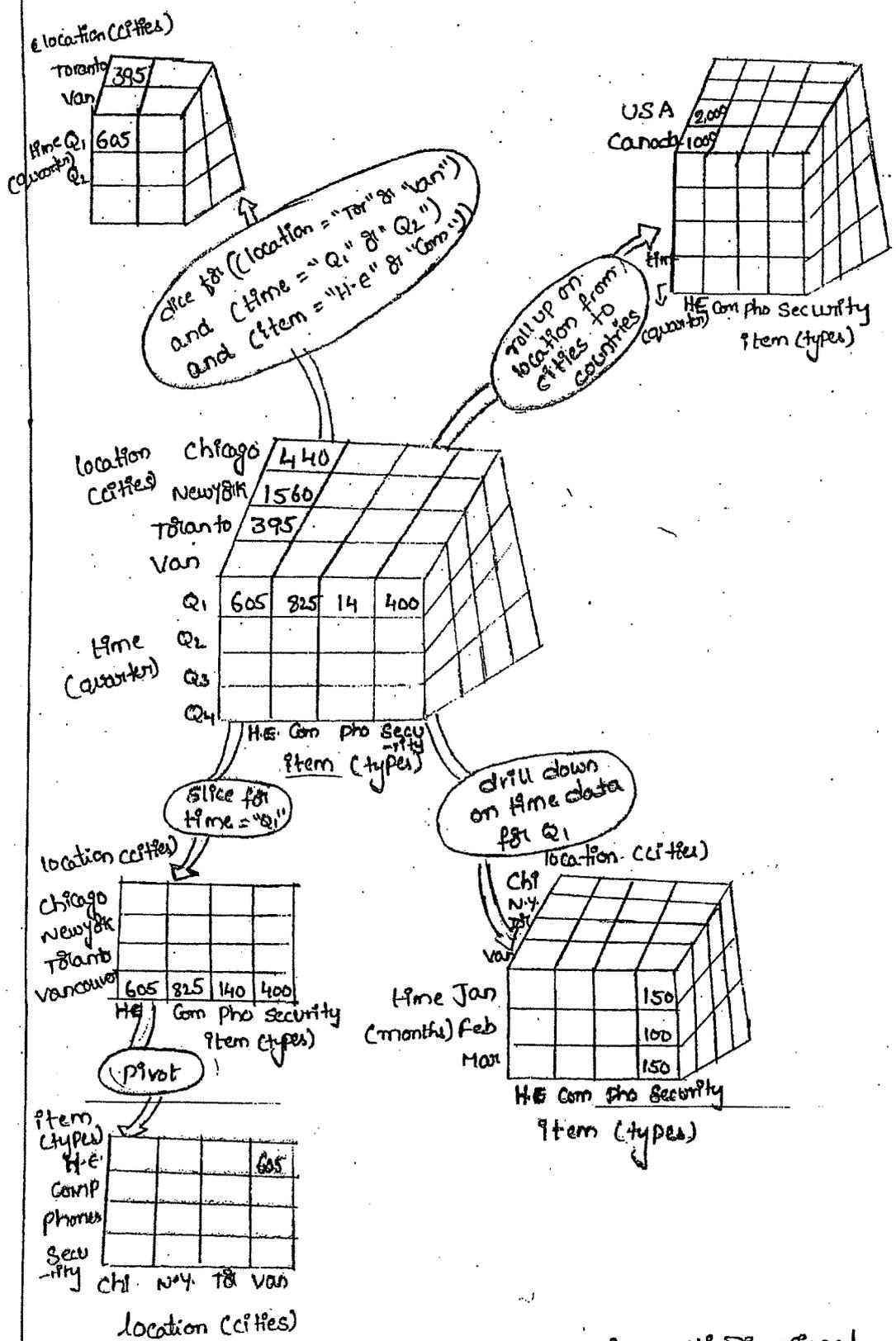
combination of data-mart is going to form

DWH

ε location (cities)

Toronto 395
Van
time Q1 605
(quarter) Q2

Slice for ((location = "Tor" & "Van")
and (time = "Q1" & "Q2")
and (item = "H-e" & "Comp"))

Roll up on
location from
cities to
countries

USA 2,000
Canada 1000

time
(country)

H.E Com pho Security
item (types)

location    Chicago  440
cities      NewYork  1560
            Toranto  395
            Van

            Q1  605  825  14  400
time        Q2
(quarter)   Q3
            Q4
                H.E  Com  pho  Secu-rity
                    item (types)

Slice for
time = "Q1"

drill down
on time data
for Q1

location (cities)

Chicago
NewYork
Toranto
Vancouver  605  825  140  400
           H.E  Com  pho  security
                item (types)

Pivot

item
(types)
H-e            605
Comp
phones
secu-rity
           Chi  N.Y  Tor  Van
           location (cities)

location (cities)

Chi
N.Y
Tor
Van

time Jan          150
(months) Feb      100
Mar               150
       H.E Com pho Security
           item (types)

fig (2·2·5·1):  OLAP operations for Multi Dimensional
                Data Model

(1) **Roll-Up:**

It is also called as Drill-Up. It specified by claimbing up a concept hierarchy for a given dimension. Here roll-up operation is specified for location (from cities to countries).

(2) **Drill-Down:** It is the reverse of roll-up operation. Here Drill-Down operation is specified for time data on $Q_1$. Drill-Down means it is the stepping down the concept hierarchy for the given dimension.

(3) **Slice and Dice:** In slice operation we select only one Dimension to get the sub cube. Here slice operation is specified for time data = $Q_1$.

In Dice operation we select the 2 or more Dimensions to create the sub cube. Here Dice operation is specified for location = "Toronto" & "Vancouver", Time = "$Q_1$" & "$Q_2$", item = "home entertainment" & "computer".

**Pivot:** It is the Visualization operation. Here we change the data acass to present the another view of data. Pivot is also called as "Rotate".

In Some of the OLAP applications additionally contains two operations.

(1) **Drill Across:** This uses the relational db requires to access the data from multiple fact tables.

2, **Drill-Through:** It also uses the relational db queries to drill lowest level of date cube. i.e; down to its back end relational table

### 2.2.6. Starnet Query Model for querying multi-dimensional Database's:

The starnet query model consists of radial lines originating from the centre and each radial line represent concept hierarchy for given dimension and also each level in concept hierarchy is called as "foot print".

The starnet query model for All Electronics database.



fig (2.2.6.1) : Starnet query model for

All Electronics Database

**2.3 Data Warehouse Architecture :-**

**2.3.1. The steps for designing and construction of DWH :-**

**Reasons for DWH Design :-**

1. The DWH provides the competative advantage. i.e., it contains the live Data.

2. Using DWH we can extend our bussiness.

3. It provides the security.

4. It improves the Customer Relationship.

5. It Reduces the cost by using Relational DB querries To Design the DWH generally we follow the 4 views.

   a. Top-Down view.
   b. Data Source view.
   c. Data ware house view.
   d. Business process view.

a. **Top-Down view :-**

In this view we select the essential data before going to the DWH Design.

b. **Data source view :-**

In this view it exposes the data stored & Managed by the operational Data bases.

c. **DWH view :-**

In this view Data is used stored in the form of fact & Dimensional Tables.

d, **Business process view** :-

Here Data is extracted from the DWH and it presents to the end user according to the opinion of the End user.

**The process of DWH Design** :-

To Design the DWH generally we follow the 3 approaches :

1. Top - Down Approach.
2. Bottom - Up Approach.
3. Combined Approach.

1. **Top - Down Approach** :-

In this approach we start from the strategic planning and Design. This approach is best suitable for the Technology is known & business problems are solved.

2. **Bottom - Up Approach** :-

In this approach we always perform the Experimentation and here the system is completed quickly.

3. **Combined Approach** :-

In this approach it contains the 2 advantages that is strategic planning from top - down design & Rapid implementation from bottom - up approach.

Generally to design the DWH we follow the any one of the 2 Methods. i.e., water fall Method (or) spiral Method.

## Waterfall Model :-

In this small model we construct systematic and structured analysis for each step and before going to the next step this is just like the waterfalling from one step to next step.

## Spiral Model :-

In this Model updation is easy i.e., it can be adapted to any new Technology or any Model. Because of this it is best suitable for DWH & Data Marks.

## The steps for Designing DWH process :-

The steps are

① choose a bussiness process to Model. i.e., If we require entire data and its Scope is enterprise wide then go for the DWH Model.

The Data scope is limitted to with in the department then go for the Data Mast-Mart Model.

(2) We have to define the level for each dimension.

(3) We have to define the dimension Tables that must be included in fact Table.

(4) We have to define the Measures that is to be analyzed like dollars-sold, units-sold.

# Three-Tier DWH Architecture :-

Query Reporting  Analysis

Top Tier : front-end Tool.
Data Mining



presentation Tools

Output

OLAP Server

Middle Tier : OLAP Server
all operations)

OLAP Server

Data clean
Data transistor
- mation).

bottom Tier :
DWH Server.

Metadata Repository

DWH

Datamart

Extract clean Trans form Integrate load

operation Databases

External sources.

fig (2.3.2) A Three - Tier DWH Architecture.

This architecture contains the 3 Tires.

1) The Bottom tier contains the DWH Server - Here Data is Extracted from the operation Databases (or) external sources like Text-file, excel-sheet online-Data Record etc. Then we apply data cleaning & Data Transformation then this entire data is integrated and finally this is loaded it into DWH. To extract the data we use the application programs. These application programs are called as gateways.

2) The Middle Tier contains the OLAP Servers. These servers are constructed directly from DWH data (or) Data Marks data. To construct these servers we use 2 Models.

    i) ROLAP Model (Relational OLAP) - ϒ 3 schemas star, snowflakes, factconstell tion

    ii) MOLAP Model (Multi-Dimensional OLAP).

        ROLAP Model, it is the Relational Db system. It maps the Multi-Dimensional data to Relational Db systems.

        MOLAP Model, These are the special servers. These serves provide the multi Dimensional view

3) The Top Tier contains the front end Tools like query and reporting tools, analysis Tools, Data Mining Tools. Using these Tools the data is presented to Management.

# Recommended Architecture for DWH :-

The Recommended architecture contains the 3 steps.

## Step 1 :

We have to define the high level enterprise Data Model. This model must be consist Realible and it must integrate the all the subject areas data.

## Step 2 :

In step 2, we design independent data marks at the same time we have to develop the enterprise DWH.

## Step 3 :

In step 3, first of all we develop the distributed data marks. These are used to integrate the different data marks data. Then finally we construct the multi-Tier architecture for DWH.



fig (2.3.2.2) A Recommended Architecture for DWH (two-Tier).

### 2.3.3 Types of OLAP Servers : ROLAP, MOLAP & HOLAP :-

The Types of OLAP servers Mainly classified into 3.

1. ROLAP (Relation OLAP).
2. MOLAP (Multidimension OLAP).
3. HOLAP (Hybrid OLAP).

These Servers provides the Multidimensional view to the end user of the data is accessed from DWH & data marks.

**1. ROLAP :- (Relation OLAP)**

: This is the extensions of relational db Management system. It is also called as Relational db Management system. It contains the advantage of greater scalability and these ROLAP servers placed in b/w DWH & front End Tools.

**2. MOLAP (Multidimensional OLAP) :-**

These Servers store the data in the form of Multidimensional arrays and provides the Multidimensional view to the end user. These ~~speca~~ Servers Contains the advantage of efficient space utilization.

**3. HOLAP (Hybrid OLAP) :-**

It is the Combination of ROLAP & MOLAP. It Contains the advantage of scalability from ROLAP & efficient space utilization from MOLAP.

## 2.4 Implementation of DWH : Using Data cubes and OLAP

The DWH contains the large amount of data. Therefore It is a critical for Data Mining system to provide efficient cube computation Techniques & query process Techniques.

### 2.4.1 Efficient computation of Data Cubes :-

In Multidimensional view cube computation extends SQL and includes "Compute cube operator". For example sales data cube for all electronics contains 3 dimensions and one measure. i.e., item product year and measure is sales in dollars. This data is analyzed by using the following

"Compute sum of sales, group by item, year"

4y "Compute sum of sales. group by year"

"Compute sum of sales group by item".

The above cube contains the total no. of Cuboids (or) group by's are $= 2^3 = 8$ i.e., { (item, product, year), (item, year), (item, product) (product, year), (item), (product), (year), ()}

This above cube is represented in the form of lattice of cuboids. This is shown in below.

() +1 is represented as all

$2^3 = 8$ (7+1) item

item

product

year

(item year)

(product, year)

(item, product)

(item, product, year)

--- 0-D (apex) cuboid

--- 1-D cuboid

--- 2-D cuboid.

--- 3-D (base cuboid).

fig (2.4.1): Lattice of cuboids makes the 3-D (Database) Datacube, the Dimensions are item, product, & year.

Compute sales groupby item(or) product oo year

Here, apex cuboid is represented by

() or call it- Specify group by is Empty total is total sum of all the sales. This is computed by using the query.

"Computed sum of total sales". Similarly one dimension operation is specified by using query

"Compute sum of sales group by item".

lly 2-D operation is specified by using query

"Compute sum of sales group by item and year.

finally 3-D operation is specified by using query.

"Compute sum of sales group by item, product and year".

The above data cube is defined by using the DMQL.

DMQL → Data Mining Query Language

Define cube sales {item, product, yearly} :
Sum (sales_in_dollars). A cube with $n$ dimension contains the Total cuboids $= 2^n$, if there is no hierarchy for each dimension. The above contain the 3 dimensions. Therefore Total cuboids $= 2^3 = 8$. But if the each dimension contains the level of hierarchy like time dimension i.e., day < month < quarter < year ; then total cuboids $= \prod_{i=1}^{n} (L_i + 1)$, where $n$ is the no. of dimensions $L_i$ means levels associated with the dimension $i$.

Here '1' specifies the high level summa-rization i.e., represented by R all , or, ( )

For example a cube contains 10 dimensions and each dimension contains the 4 levels then approximately. The total no. of cuboids is $5^{10}$. Therefore it contains the large no. of cuboids and also each cuboid requires the large amount of space. This is avoided by using concept hierarchy. — 3 level stores 2 level 1 compute stored

card → countries → street

3. The aggregates are computed from the previously computed aggregate Rather than taking from fact Table.

## MOLAP cube Computation Techniques :-

In MOLAP the data is stored in the form of Multidimensional array. Here each row is divided it into chunk. The chunk is a subcube and small enough to fit in the memory available for cube. The n-Dimensional array is divided it into n-chunks. This is called as chunking. Each chunk is stored as a object on the disk. Each chunk is identified by "chunkID + offset".

For example 3-Dimension array contains 3 dimensions. i.e., A, B & C, then the chunks are shown in below.



fig (2.4.2): 3-D array contains 3 dimensions A, B & C organized into 64 chunks.

## partial materialization :-

The Data Cube materialization mainly classified into 3.

1. No Materialization ( eg - item)
2. Full Materialization
3. partial Materialization ( no of dimension & value represent)

In those 3rd one is best option i.e., partial Materialization. partial Materialization Means selected cuboids are materialized. i.e., selected cuboids are analyzed. i.e., The selected cuboids data is represented by using Data cubes.

The partial Materialization contains the 3 steps.

1. Identify the (Subjects) of cuboids that is to be materialized.
2. Analyze cuboids during the query processing
3. update the cuboids during the data load process.

## Multi-way Array Aggregation in the Computation of Data cubes :-

To Analyze the Data ROLAP contains tuple & tables these Two are the basic data Structures in ROLAP.

## ROLAP cube Computation Techniques :

The ROLAP uses the

1. Sorting, Hashing & grouping operations on dimension attributes to recorder the data.
2. The ROLAP uses the grouping & sub queries to fastly the processing of sub querries.

The Datacube contains the following.

1. The high level summerization is represented by () or all.

2. The 2 dimension i.e., AB, AC & BC is calcolated by grouping AB, AC, BC.

3. The one dimension i.e., A, B & C is calcolated by grouping dimensions individually.

4. The 3-D (base) cuboid is represented by ABC.

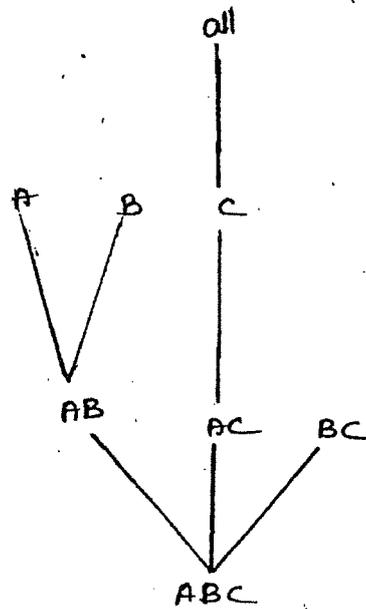**Two arrangements for Multiway array aggregation for 3-D cube** concept hierarchy of $5 \& $ values of -el.
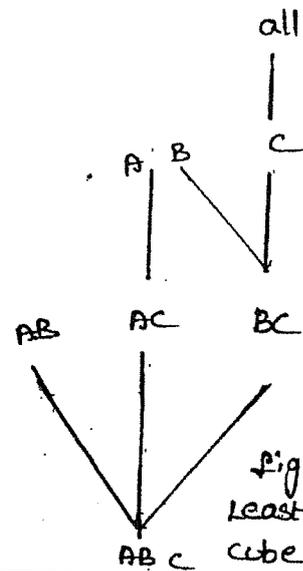


fig (a) Best · used Cube computing

fig (b).
Leastly used cube computat.

for ex, the Memory units for dimension A, B & C are 40,400 & 4000 Memory units. then largest space required for 2-Dimension. i.e., BC (for units = 400×4000 = 1600000 MU) then 2nd largest space is required for 2-D.

i.e., AC (for units = $40 \times 4000 = 160000$ MU). The smallest space is required for 2-D.

i.e., AB (for units = $40 \times 4000 = 16000$ MU).

$\therefore$ In fig (a), Minimum space required.

$= AB$ (for whole AB) $+ AC$ (for one row for AC)

$\qquad\qquad + BC$ (for one chunk for BC).

$= 40 \times 4000 + 10 \times 4000 + 100 \times 1000$

$= 16000 + 40000 + 100000 = 1,56,000$ MU.

In fig (b), Minimum space required.

$= BC$ (for whole BC) $+ AC$ (for one row for AC)

$\qquad\qquad + AB$ (one chunk for AB).

$= 400 \times 4000 + \frac{40}{10} \times 4000 + 10 \times 100$

$= 1600000 + 40000 + 1000$

$= 16,41,000$ MU.

## 2.4.2 : Indexing OLAP data :

To provide efficient data access most of the DWH systems contains index structure and materialization. The materialization is provided through data cubes. The index structure is provided by using bit map indexing and join indexing.

### Bit map indexing :-

This is the most popular data accessing Technique in OLAP. Using this we find the data quickly in data cube. This is the alternative Technique for Record-Id (RID) lists. In Bit-Map indexing. Each attribute value

consists of distinct bit vector values. If the attribute value 'v' for the given row in base table. then its value is set to '1' corresponding to the row in the bit-map index, all the remaining bits are set to '0' in the corresponding row at bit-map index.

for example consider the all electronics db.

Base Table

| RID | item | city |
|-----|------|------|
| R₁ | H | V |
| R₂ | C | V |
| R₃ | P | V |
| R₄ | S | T |
| R₅ | H | T |

Item Bit-map Index Table

| RID | H | C | P | S |
|-----|---|---|---|---|
| R₁ | 1 | 0 | 0 | 0 |
| R₂ | 0 | 1 | 0 | 0 |
| R₃ | 0 | 0 | 1 | 0 |
| R₄ | 0 | 0 | 0 | 1 |
| R₅ | 1 | 0 | 0 | 0 |

City Bit-map Index Table

| RID | V | T |
|-----|---|---|
| R₁ | 1 | 0 |
| R₂ | 1 | 0 |
| R₃ | 1 | 0 |
| R₄ | 0 | 1 |
| R₅ | 0 | 1 |

fig (2.4.2.1): OLAP index by using Bit-map Indexing

### Advantages of Bitmap indexing :-

(1) It is the easy as well as the simple compared with the indexing of hashing.

(2) It requires the less amount of accessing time because the value is represented by using only 1 bit.

(3) It requires less amount of space.

### Join indexing :-

It is mainly used in relational db by using this we easily find the joinable toples. fot ex, Two relations are R(RID, A) and S(B, SID) join on attributes A and B. The join index record pair (RID, SID) where RID and SID are the record identifiers for relation R and S.

For example consider the sales fact table and 2D Tables i.e., location and item. The relation b/w the Sales fact table and dimension tables are shown in below.



fig(4.2.2) : Relation b/w sales fact table and 2D tables i.e., location and item.

Here, main street value in location dimension joins the tuples with T59, T234 & T883 of the sales fact tables similarly the value sony TV of item dimension joins the Tuples T59 and T449 of the sales fact table.

The join index Tables are shown in below :

join index Table
Location / Sales

| Location | Sales −key |
|---|---|
| Main _street | T59 |
| Main-street | T237 |
| Main-street | T883 |

join index Table
item / Sales

| item | Sales − item |
|---|---|
| Sony−TV | T59 |
| sony−TV | T449 |

join Index Table for Two Dimensions
Location / item / Sales.

| Location | item | Sales _key |
|----------|------|-----------|
| - - - - - | - - - - - | - - - - - |
| Main - street | Sony TV | T59 |
| - - - - - | - - - - - | - - - - - |

fig (2·4·2·3) join Index Tables.

This join index Tables mainly used to maintain the relation b/w primary key and foreign key. This is mainly used in staar schema to maintain the relation b/w foreign key of the fact table with primary key of dimension tables.

## 2·4·3 Metadata Repository :-

Metadata is nothing but data about the data. The Metadata repository contains the following steps.

1. The structure of DWH Which includes schemas views, dimensions, hierarchies and data definitions.

2. The operational repository which includes data usage and monitoring information.

3. It also includes the algorithms for summari- -zation.

4. To map data sources to data ware house which includes data cleaning, data transformatic data integration and security Techniques.

5. The bussiness repository which includes bussine policy, conditions, terms and business definiti

## 2.5 Datacube Technology :-

Here initially we write about the data cube in multidimensional data model. It is mainly classified into.

### 2.5.1 "Discovery Driven Exploration" of Data cubes:-

Here analyst or user search for the interest pattern by using OLAP operations i.e., drill-down, drill-up, slice and dice & pivot. OLAP opera -tion imp

If the discovery process is not automated then end user search for the interesting patterns manually. This is difficult. This is avoided by using the method discovery driven exploration. In this discovery driven exploration we identify the interesting pattern and it is marked. To provide the visualization to the end user.

### Discovery Driven Exploration :-

In this method previously computed measures indicates the data exceptions after that these measures treated as exception indicators and using this enduser identifies the interesting pattern easily. ∴ These measures provide the "degree of Surprise" i.e., w.r.to existing value in the cell with expected value.

These measures are calculated by using statical analysis & grouping functions. These measures mainly classified into 3.

# Chapter-3
## Data preprocessing

Database consists of massive volume of data which is collected from heterogeneous sources due to this heterogenity, real world data tends to be inconsistent and noisy. If data is inconsistent, then there is a possibility that mining process can lead to confusion which results in accurate.

## Need for preprocessing the Data:

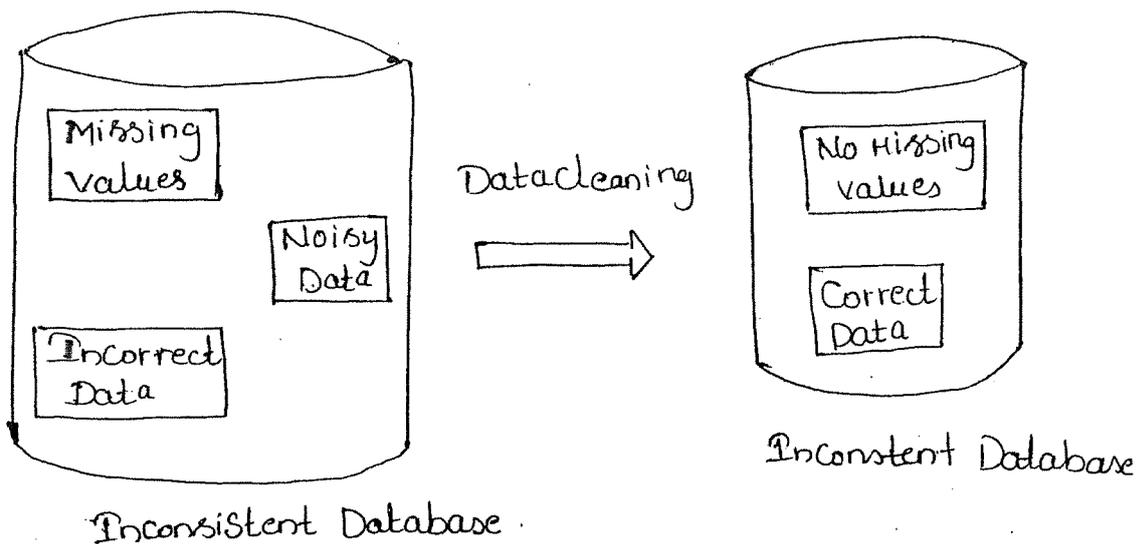Incomplete, noisy and inconsistent data is common in large real world databases and datawarehouse

Incompleted data can occur for a no. of reasons.

1. Attributes of interest may not always be available.
2. Relevant data may not be recorded due to a misunderstanding.
3. Data that is inconsistent with other recorded data might be deleted.
4. The data collection instruments used may be faulty.
5. There may have been human or computer errors occuring at data entry.
6. Errors in data transmission can also occur.

→ To overcome the above problems the following data preprocessing techniques are required.

1. Data cleaning
2. Data Integration
3. Data Transformation
4. Data Reduction.
5. Data Discretization.

1. **Data Cleaning** : when data is collected from datasources then there are chances that the data can be inconsistent, incompleted and noisy. Data cleaning is a Process of removing unnecessary and inconsistent data from the databases. The main purpose of data cleaning is to improve the quality of data by filling missing values, reconfiguring the data to make sure that data is in consistent format.



Inconsistent Database.

fig(1): Datacleaning

1.1: **Missing Values** :

    The missing values consists the following techniques and corrected befor applying DKQL.

a) **Ignore the tupple:**

This technique is simple. In this just we avoid the tupple which doesn't contain values. But it's not recommended.

b) **Fill the missing values manually:**

Here we manually fill the missing values. But this technique is not recommended for large database which contains important values.

c) **Using Global Constant values:**

Here the missing values has been filled with the global constant such as unknown as (or) infinity. But this technique not succeeded because it deviates the Datamining Process.

d) **Using attribute mean value:**

Here each missing value has been filled by mean value i.e, In All Electronics customer average income is 2800 $. Then the customer record i.e the income attribute's missing value is filled with 2800 $

e) **Use the most portable values:**

Here we find out the most portable values by using different techniques like Bayesian classification and decision Tree Induction etc.

## 1.2 Noise data:

Here the noise data has been smoothing out by comparing neighbouring values. In this technique

datavalues has been distributed into different wings.

1. Binning Method :

a) : Smoothing by Bin-by-mean

Here we find out the mean value for each Bin & missing values are replaced with that mean value.

b). Smoothing by Bin Bounding :

Here we find out min and max values for each Bin then each Bin's value is replaced with the closest value to the min or max.

Ex: price in dollars of an item is as follows.

4, 8, 15, 21, 21, 24, 25, 28, 34.

Bin 1 : 4, 8, 15

Bin2 : 21, 21, 24
Bin3 : 25, 28, 34

Smoothing by Bin

Bin 1 : 9 9 9

Bin 2 : 22 22 22

Bin 3 : 29 29 29

Smoothing by Bin Boundary value;

Bin 1 : 4 4 15

Bin 2 : 21 21 24

Bin 3 : 25 25 34

In Bin boundary we identify the data if it is near
to the min value then we replace that value with
min I or min values otherwise if it is near to the
max value then we replace that value with the
max value.

Clustering: In this we find out outliers i.e, clustering
identifies similar data objects those are placed in
one cluster & also identifies dis-similar data objects
and those are called outliers, in other words
grouping the similar data objects into 1 place is
called as clustering.



→ outliers.

→ similar data objects.

Using combination of computers & human inspection:

In this outliers can be identified through
a combination of computer & human inspection
i.e, we need to identify any algorithmic approa
-ch to find out the outliers and clusters in our data
set along with this we have to take the help of
manual procedure to identify clusters & outliers.

Inconsistent data:

Inconsistent data means the data with

duplicated values. For Eg; Item Id is used to categor-ize items in AllElectronics. If the item-id entered wrong or missing then we can say that Item-id is inconsistent data. Then by applying different Smoothing techniques Bin by mean, Bin by Boundary etc. to Convert the data into consistent data.

## Data Integration & Data Transformation:

Here the data from different data sources is Collected and placed into the DwH, but it contains different difficulties to make unique Identification or unique representation for the data of different dataSources.

Data mining requires data integration. Data integration is the merging of data from multiple data Stores into a coherent data Store as in DwH. These Sources may include multiple databases, datacubes or flat files.

Data Integration: Data integration is a process of combining data from heterogeneous data Sources Such as different data bases, flatfiles etc, to form a Single consistent data repository.



Datacube

Database

Data Integration

Data Repasi-tory.

Data Integration.

For eg: In one data source customer-id is entered as cust-id where as in another data source it entered as Customer Number.

The integration contains different problems while collecting the data from multiple tables. This can be avoided with an analysis called as correlation analysis. For eg: The Correlation analysis or relation b/w two attributes can be measured as,

$$= \frac{\sum (A-\bar{A})(B-\bar{B})}{(n-1)\ \overline{\sigma A}\ \overline{\sigma B}}$$

In the above Equation 'n' represents no of tupples $\bar{A}$ & $\bar{B}$ represents mean values of A,B $\sum A$ & $\sum B$ represents standards deviations of A,B.

we can find the mean value of A as $\bar{A} = \frac{\sum A}{n}$

Similarly we can find out the standard deviation of A as

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

If the result of correlation analysis is greater than 0 (zero) then two attributes are positively correlated ie, if we increase the value of attribute A then the value of attribute B will also increase.

If the correlation analysis result is '0' then two attributes are independent to each other.

If the Correlation analysis result is less than 0 then the 2 attributes are negatively correlated ie If the value of A increases then the value of B will decrease.

## Data Transformation:

Data Transformation is nothing but converting diff data sources into a format that must be acceptable for the datamining system.

5, 48, 99, 35, 81 $\xrightarrow{\text{Data Transformation}}$ 0.005, 0.048, 0.049, 0.035, 0.081.

### 1. Smoothing:

Here the noise data is subjected for smoothing by using different smoothing techniques, Binning, clustering etc, .

### 2. Aggregation:

Here we use some aggregated functions to perform data transformation. For eg: using daily sales we can compute monthly sales. In the same way, using monthly sales we can compute yearly or Annual sales.

### 3. General Generalization:

Here lower level value is replaced with higher level value by using concept hierarchy. For eg: The dimension location contains concept hierarchy street, city, state & country instead of placing the street value would be more generalized to place Country value; the data mining system.

### 4. Normalization:

In Normalization we force each value should in specific range. -1.0 to +1.0

### 5. Attribute Construction: In this new attribute

will be constructed by integrating attributes from different sources.

In the above 5 methods the Best one is "Normalization".

## Normalization:

It contains following techniques.

### 1. Min-Max Normalization:

Here $min_A$, $max_A$ are the minimum & maximum values of the attribute A. Then the value v of A is transformed into v'. By using new range of new $min_A$ & new $max_A$.

we compute new values i.e., v' By using below formulae.

$$v' = \frac{v - min_A}{max_A - min_A}(new\ max_A - new\ min_A) + new\ min_A$$

min max Normalization maintains same relation tip of original datavalues. For Eg: The income attribute contains the min & max values of $ 12000 & $ 98000 then the income attribute changed into new range (0,1).

By using min, using min max Normalization the original value is $ 73600 will be Transfor -med, into

$$v' = \frac{73600 - 1200}{98000 - 12000}(1-0) + 0.$$

$$v' = 0.716$$

## 2. Z-Score Normalization:

Z stands for Zero mean. So this normalization is also called as Zero mean normalization. In this attribute is find out by using mean & standard deviation value.

The value v of A will be transformed by using the formulae.

$$v' = \frac{v - \bar{A}}{\bar{\sigma A}}$$

Here $\bar{A}$ is called as mean value $\Sigma A$ is called as standard deviation. For Eg: The attribute 'A' mean & standard deviation values are $ 54000 $16,000 then we can Z-Score the Normalization $F_1$ value v $ 73600 will transformed into

$$v' = \frac{73600 - 54000}{16000}$$

$$= 1.225$$

## Normalization by decimal Scaling:

In this technique we move the decimal value based on absolute value of an attribute. By using this technique the original B of A is transformed into B' by using.

$$v' = \frac{v}{10^j}$$

where $j$ is the smallest integer based on the attribute value.

for Eg The attribute A contains values between $-986$ to $+917$. Then we find out the absolute value as $986$ by using the decimal scaling we can find out the new value $v'$ by dividing it with $1000$

i.e $J = 3$

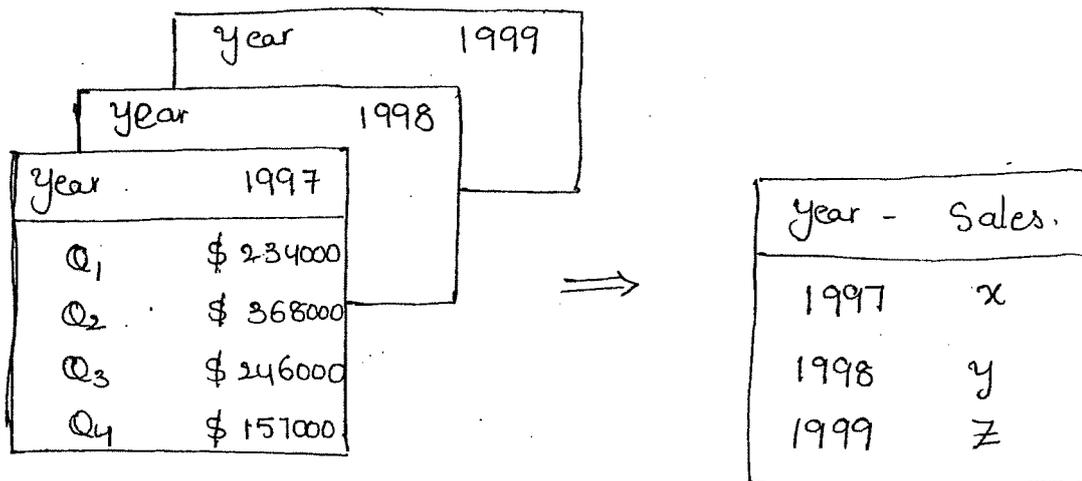$$J = 3$$

$$v' = \frac{-986}{10^3}$$

$$= -0.986$$

## Data Reduction:

Here we apply aggregation operations, Redundant attributes are removed from the data. Then we can apply data compression Technique on data.

In data reduction, we reduce the size of the data in such a way that essential features does not effect. It contains several techniques.

i) Data aggregation:

In Data aggregation we will apply some aggregation operation on the data to construct the data cube & we can also identify some clustered data to reduce the size of the data.

For Eg. In All Electronics database, the data is stored quarterly based for 3 years i.e, for the year 1997, 1998 & 1999 which is shown in the below diagram.

| year | 1999 |
| year | 1998 |

| Year | 1997 |
|------|------|
| $Q_1$ | $ 234000 |
| $Q_2$ | $ 368000 |
| $Q_3$ | $ 246000 |
| $Q_4$ | $ 157000 |

$\Rightarrow$

| Year - | Sales. |
|--------|--------|
| 1997 | $x$ |
| 1998 | $y$ |
| 1999 | $z$ |

we can apply different aggregation operations for the annual sales to represent the data in yearly format rather than storing it into Quarterly wise.

In the reduction, we can also represent the data in data cube format to represent multiple dimension which is represented in the below diagram.



fig: Sales datacube for all electronics.

ii) Dimensions Reduction: Here the redundant attribu or dimension is identified & it is removed. It contai -ns several methods to reduce dimensions. In those methods, the main aim is we need to identify th min no. of attributes to represent the dataset. Among the available techniques, the greedy method is important. In this method, it enable us to find out the best attribute or worst attribute based on decision tree analysis. It contains the following techniques.

1. Forward Selection: For Eg: Consider the following decision tree.



fig: Decision Tree data.

In forward selection we will start the identification with the empty set. Consider the data set as { A, A₂, A₃, A₄, A₅, A₆ } initial set is { } (empty set) In the order the first attribute to consider these {A₁}. The next dataset we consider is { A₁, A₄} & the next data set we consider is { A₁, A₄, A₆} we leave other attributes which are not important.

2. The Backward Eliminations:

Here we start our searching from the complete

Dataset. In each step we eliminate worst attribute & that is removed from the dataset. This process will contained until all the dataset is complete. Here the dataset is $\{A_1, A_2, A_3, A_4, A_5, A_6\}$.

So we start over search from this only. In first step, we eliminate the attribute $\{A_2\}$ So the dataset becomes.

$$\{A_1, A_3, A_4, A_5, A_6\}.$$
$$\Downarrow$$
$$\{A_1, A_4, A_5, A_6\}$$
$$\Downarrow$$
$$\{A_1, A_4, A_6\}.$$

3. Combination of forward selection & Backward elimination: It is the integration of above 2 techniques. This is basically applicable for large Dataset. In this we need to identify the best attribute as well as worst attribute based on our requirement. we will add the best attribute to the dataset at the same time we remove the worst attribute from the Dataset.

To perform all these redundant techniques we need to use decision tree induction. In decision tree induction, we perform test on attributes & we identify the results based on the classifications. Each classification is represented by "class labels" In the last example, the total dataset i.e $[A_1, A_2, A_3, A_4, A_5, A_6]$ is reduced as $\{A_1, A_4, A_6\}$ & final result

will fall into either class 1 or class 2.

4. Definition of Data Reduction :

Data reduction is a process of compressing massive volume of data into limited data set without sacri -ficing data integrity.

| Year | 1998 |
|------|------|
| Sem | Total |
| 1 | 850 |
| 2 | 800 |

| Year | 1999 |
|------|------|
| Sem | Total |
| 1 | 825 |
| 2 | 830 |

Data Reduction →

| Year | Total |
|------|-------|
| 1998 | 1650 |
| 1999 | 1655 |
| 2000 | 1600 |

Year

| Year | 2000 |
|------|------|
| Sem | Total |
| 1 | 800 |
| 2 | 800 |

iii Data Compression

It means we reduce the data in such a way that it must be in smaller size i.e, we reconstruct the data from the data in the form of compressed. If we compress the data without any loss of information that Data Compression technique is called as "Loss Less Data Compression" otherwise it is called as " Loss Data Compression".

This Data compression contains two techniques

1. Discrete wavelet Transform (DWT)
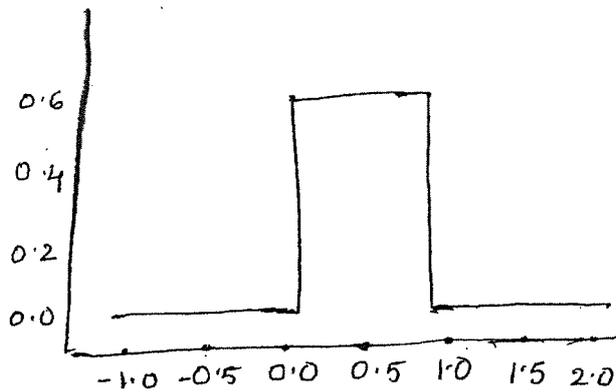2. Principle Components Analysis (PCA).

1. Discrete wavelet Transform :

In this technique, the data Vector-D is transformed into D' of wavelet co-efficient. In this technique, we transform any value into new value by considering the coefficients with specific Range. This technique is derived form DFT (Discrete Fourier Transform). This DWT contains several wavelet Transform techniques. But familiar ones are.
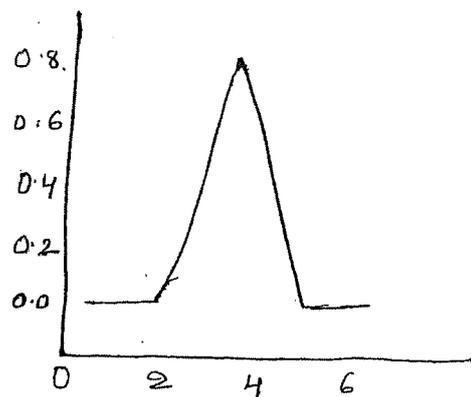
(a) Haar-2
(b) Daubenchies 4

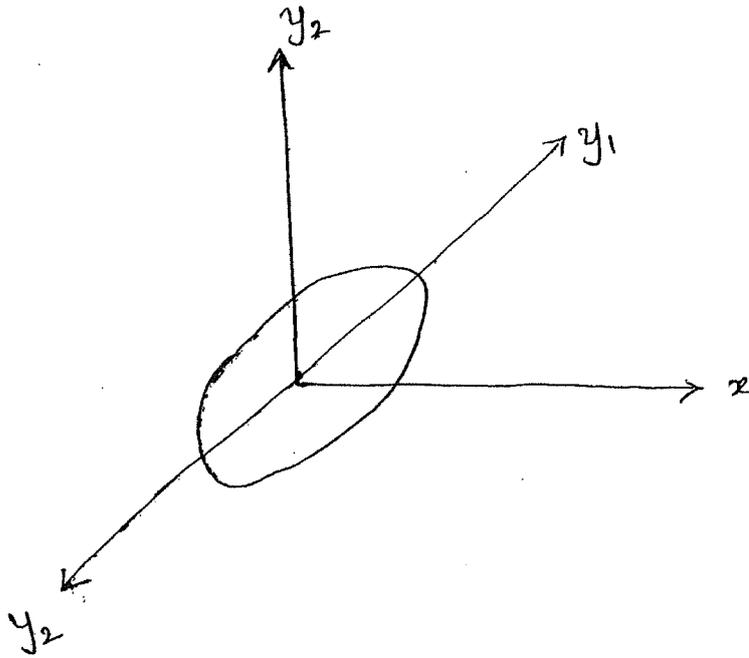which are represented in the belao diagram.



(a) Haar-2

(b) Daubechies-4

In Haar-2 technique we will consider the data based on the specific range co-efficient r. In the above diagram, the range for considerable values is from o to r. Then we consider the class labels in that specific rare by leaving other one.

Daubechies 4 also follows same technique but represent in the form of pyramid. The only difference is in Daubechies 4, the dataset length L is represented in the form of even values or in the form of L integer power of 2. Here we apply two functions to compress the data. In that initially we apply smoothing techniques such as sum & average and we findout the differences. If the dataset is too large, then we can divide the dataset consider its length as $L_2$.

The procedure will repeat until the length of the dataset becomes 2. Then we will get compressed data which can be placed into the Datamining system.

2. Principle Components Analysis:

Here we compress the data of N-tuples or N Data vectors reduced into K-dimensional orthogonal vectors. By using this technique we can represent the data into compressed vectors. This technique is called as KL (Karhunen-Loébe) technique.

$y_2$

$y_1$

$x$

$y_2$

In principle component of analysis we apply the normalisation of the dataset to represent the entire values in specific range we identify k-dimensi -onal orthogonal vectors & then we compute orthogonal vectors c with new range suchthat

$c <= k$.

Here we construct new access to represents the compressed dataset. According to the above diagram $x_1$ & $x_2$ are original axis and $y_1$ & $y$ are new axis for the original data which represents the compressed data from the original data.

In this method, final higher principle compon -ents are integrated to represent the original data.

iv Numerosity Reduction :

Here we reduce the data by selecting the small
er group of values. It is classified into.

i) parametric

ii) Nonparametric

Parametric : In parametric Reduction Technique,
the data is reduced by using

i, Linear regression

ii, Log-Linear regression Models.

→ In Linear Regression model, we find the best
line to be fit b/w the pair of attributes. If 1 value
is given then we find the another value. Here we
use the formula.

$$y = \alpha + \beta x$$

Where $\alpha, \beta$ are the regression coefficients.

$x$ is the predictor value and

$y$ is the Responsive value.

→ In Log Linear Model we predict the cell
value of base cuboid. Using this we find the high
-er level cuboids values. i.e By using low
level cuboid values we can find high level cuboid
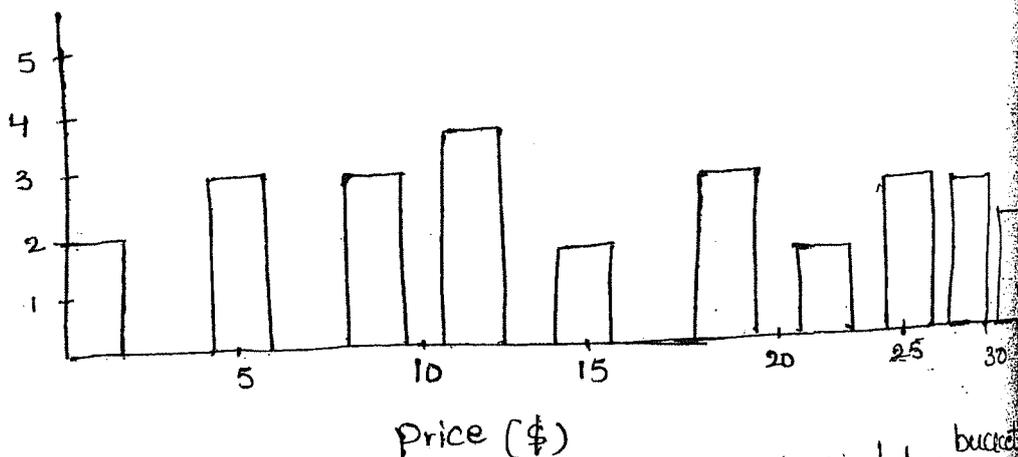values.

## Non - parametric :

Here data is reduced by using the several methods. The methods are.

1. Histograms
2. Clustering
3. Sampling

### 1. Histograms :

Here we use the several smoothing techniques. The histogram of an attribute 'A' is partitioned into disjoint parts. These are called as buckets. These buckets are represented in horizontal axis and the bucket count i.e, frequency is represented in vertical axis

For example, consider the all electronics sales data for price is specified in dollars i.e ; 1, 1, 5, 5, 8, 8, 8, 10, 10, 10, 10, 15, 15, 18, 18, 18, 21, 21, 25, 25, 25, 28, 28, 28, 30, 30. Then the histogram for the above data is show in below.
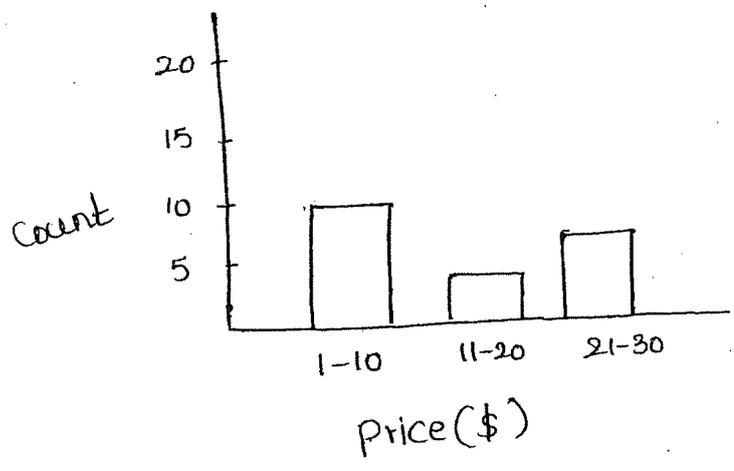


Price ($)

fig: Histogram for price according to singleton bucket

Here, each bucket contains only one value then those buckets are called as singleton buckets to partition the data we use several partitioning techniques. In those familiar one's are:

i. Equi width:

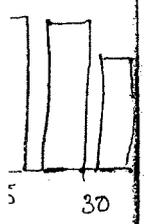The equiwidth histogram contains equal width or constant width. The equiwidth histogram for the above data.



ii. Equidepth

The equidepth histogram contains the constant frequency that is, each bucket approximately contains the equal no. of samples.

iii) v-optimal:

The v-optimal histogram contains the how variance. The histogram variance is nothing but sum of values of the bucket.
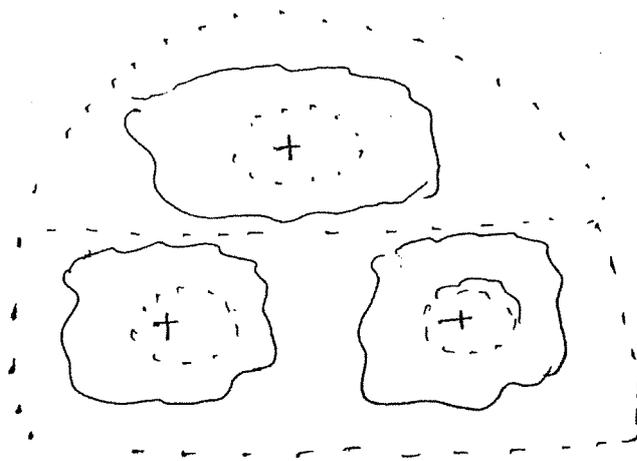
L

lled

Sales
5,5,
i,25,
the

5    30

buckets
on

**iv) Max.-Diff :**

The Max-Diff histogram contains the maximum difference of each pair of adjacent values. Here, we use the formula "$\beta - 1$", Here '$\beta$' is the maximum difference. This is specified by the end-user.
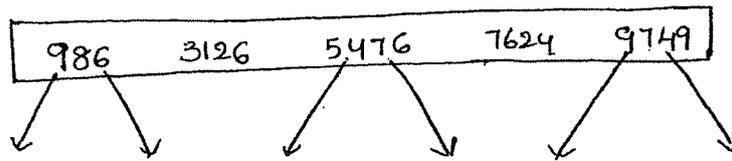
**2. Clustering :**

In clustering similar data objects are stored in one cluster and dissimilar data objects are stored in another cluster. For example, 2D customer data w.r.to customer locations in city is shown in below.



Fig (a) : 2D customer Data w.r.to custom location in city.

Here center of the cluster is marked with '+'. But in data reduction requires the original values to store the data in efficient way and access the data in efficient way. We use multidimensional index Tree structure. $B^+$ tree organization.

The B+ Tree organization is shown in below.

| 986 | 3126 | 5476 | 7624 | 9749 |
|-----|------|------|------|------|

fig(b): B+ Tree organization for data set.

For ex consider the large database. This datab-ase contains the 10,000 tuples. These are represented by using the keys 1 to 9999. We partitioned this data by using the equidepth Histogram. i.e. Each bucket approximately contains the equal no. of values. Then we get the 5 keys. These are ranging from 1 to 985, 986 to 3125, 3126 to 5475, 5476 to 7623, 7624 to 9748, 9749 to 9999.

Therefore, Each bucket approximately contains the 10,000/6 and also the keys are again divided it into subkeys.

3. Sampling :

In sampling we reduce the data i·e Large amount of data is replaced with small random sample For eg. the data set "D" and it contains the 'N' no. of tuples. Then the samplings are:

i) Simple Random Sample without Replacement (SRSWR) for size n :

Here we Select the 'n' tuples where n < N. The Probability of drawing the tuple from D is $1/N$. Here we draw the tuple but, that tuple never be rewritten.
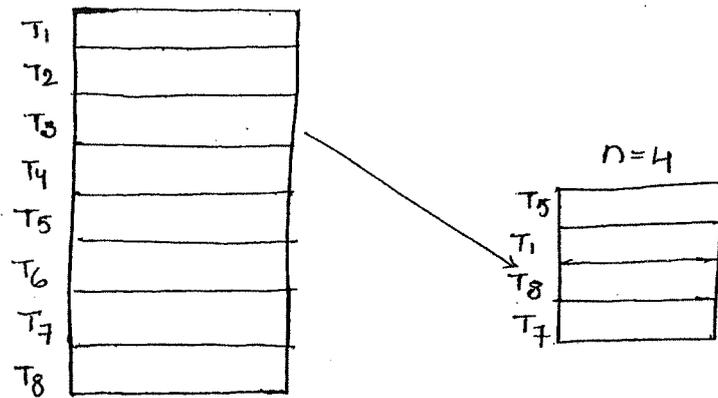
fig (a) : SRSWOR for n=4

ii) **Simple Random sample With Replacement (SRSWR).**
   **for size n :**

   Here, we access the tuple, recorded and rewritten. This technique allows the end user. The tuple $x$ is accessed again.
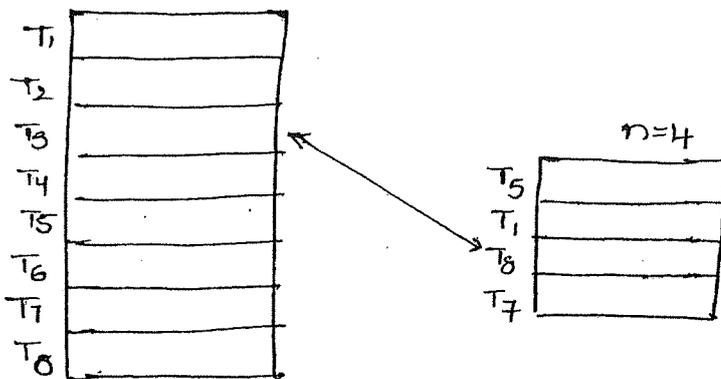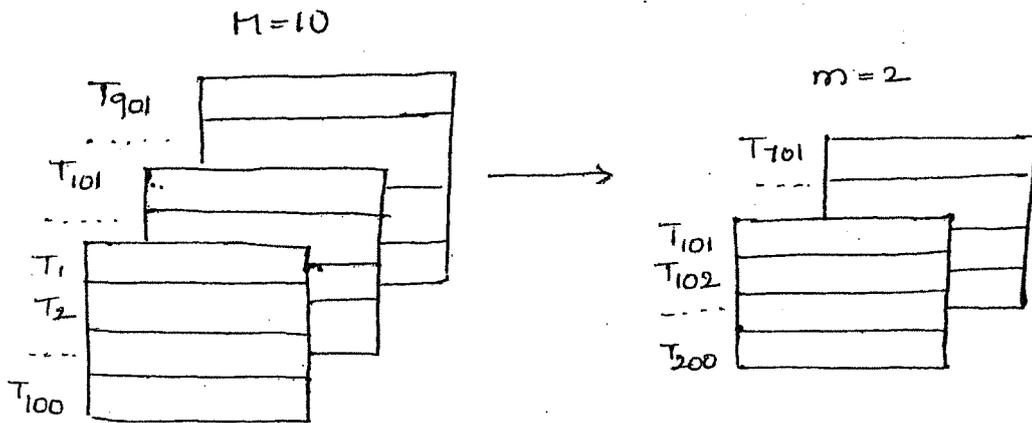


fig (b) : SRSWR for n=4.

iii) **cluster sample :**

   Here the tuples in 'D' grouped into 'M' mutually disjoint clusters. Then we select 'm' clusters. This is shown in below.
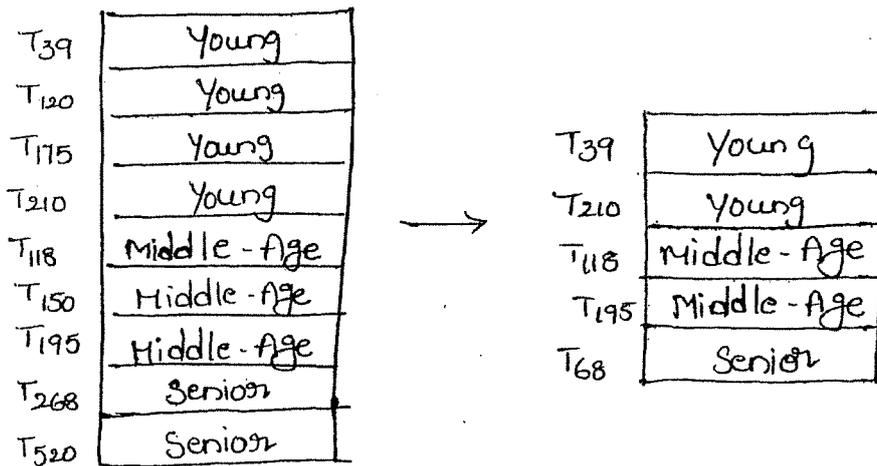
H=10



m=2

$T_{901}$
$T_{101}$
$T_1$
$T_2$
$T_{100}$

$T_{101}$
$T_{101}$
$T_{102}$
$T_{200}$

fig (c): cluster sample with m=2.

iv) **Stratified Sample:**

Here, Large tuples of 'D' partitioned into mutually disjoint parts. These parts are called as strata. Then we generate the stratified samples of D for each stratum. This is shown in below.

| | |
|---|---|
| $T_{39}$ | Young |
| $T_{120}$ | Young |
| $T_{175}$ | Young |
| $T_{210}$ | Young |
| $T_{118}$ | Middle-Age |
| $T_{150}$ | Middle-Age |
| $T_{195}$ | Middle-Age |
| $T_{268}$ | Senior |
| $T_{520}$ | Senior |

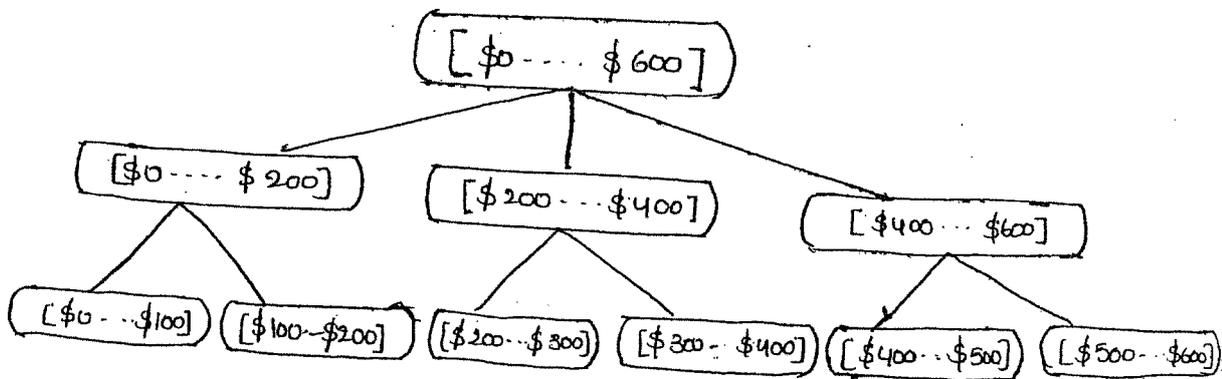| | |
|---|---|
| $T_{39}$ | Young |
| $T_{210}$ | Young |
| $T_{118}$ | Middle-Age |
| $T_{195}$ | Middle-Age |
| $T_{68}$ | Senior |

fig (d): stratified sample according to Age.

## Discretization & concept Hierarchy Generation:

The Discretization means we reduce the data value of an continuous attribute by dividing it into intervals. In this best method is concept Hierarchy In concept Hierarchy we collect and replaced with

low level Hierarchies. For ex, Age attribute contains the values young, middle and seniers.

Then the Age attribute is replaced with one of these 3 values. Concept Hierarchy for price is shown in below:

```
                    [$0 .... $600]
          /              |             \
   [$0 .... $200]   [$200 ... $400]   [$400 ... $600]
    /      \          /       \          /       \
[$0..$100][$100-$200][$200..$300][$300-$400][$400..$500][$500-$600]
```

fig(a): Concept Hierarchy for price.

## Discretization and Concept Hierarchy Generation for Numeric data:

We apply the Concept Hierarchy easily for Numeric data by using the method distributed attribute analysis. This method contains the 5 techniques.

1. Binning
2. Histogram Analysis
3. Cluster Analysis
4. Entropy - Based Discretization
5. Data Segmentation by natural partitioning

### 1. Binning:

Here data is partitioned and distributed it into different buckets or bins. Here we use the several Smoothing techniques.

- Smoothing by bin Mean:
Here, we find the mean Value for each bin and

that is replaced with each value in bin.

ii. Smoothing by bin Boundary :

Here we find the maximum value, minimum value. Then each bin value is replaced with closest boundary. For ex, Consider the price in dollars are 4,8,15,21,21, 24,25,28,34 Partitioned with equidepth i.e , 3.

bin 1 :  4   8   15

bin 2 :  21   21   24

bin 3 :  25   28   34

Smoothing by bin Mean :

bin 1 :  9   9   9

bin 2 :  22   22   22

bin 3 :  29   29   29

Smoothing by bin boundary :

bin 1 :  4   4   15

bin 2 :  21   21   24

bin 3 :  25   25   34
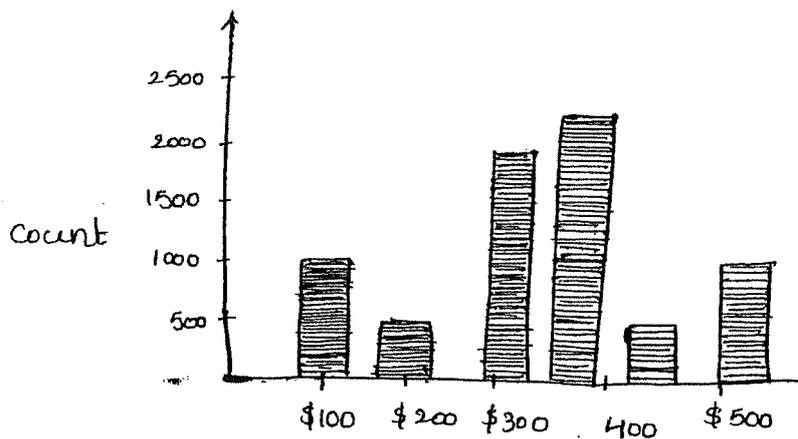
2. Histogram Analysis :

The Histogram for an attribute A is composed by Partitioning data into disjoint bucket. These buckets are represented in horizontal axis and the bucket count i.e frequency is represented in vertical axis. For ex. Consider the most frequency range for price is $300 to $350 and this data is partitioned by using Equiwidth partition i.e. Each bucket contains the
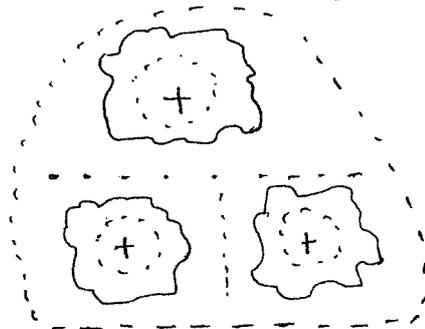
Equiwidth



fig(a): Histogram Analysis for price ($).

## 3. Cluster analysis:

In cluster Analysis, similar data objects are stored in one cluster and dissimilar data objects are stored in another cluster. The 2D customer data with respect to customer Locations in city is shown in below.



fig(a): 2D customer Data w.r to customer Locations in city.

Here center of the cluster is marked with '+' and the cluster again divided it into smaller clusters.

## 4. Entropy Based Discretization:

The information based measure is called as Entropy. This is applied recursively for the numeric data of Attribute A.

Let us consider 'm' classes. The sample 'S' contains $S_i$ samples in class $C_i$ for $i = 1, 2, 3, \cdots m$. Then expected

information for the given sample.

$$I(S_1, \cdots, S_m) = -\sum_{i=1}^{m} \frac{S_i}{S} \log_2 {}^{S_i/S} \longrightarrow \textcircled{1}$$

Then we can calculate Entropy for an attribute A as.

$$E_A = \sum_{i=1}^{v} \frac{(S_{1i} + \cdots + S_{mi})}{S} \times I(S_{1j}, \cdots, S_{mj}) \longrightarrow \textcircled{2}$$

Finally we can calculate information gain as follows

$$gain(A) = I(S_1, \cdots S_m) - E(A)$$

Based on the information gain, we can identify strongly relevant attributes and weakly relevant attributes.

5. **Data segmentation by Natural partitioning** :

Many of the users prefer the data is distributed uniformly. The uniform distribution allows to the end-user. The end-user read the data easily. For ex, Annual Salaries of a particular company ranging from

[$ 50,000 \cdots $ 60,000] rather than specifying.

[$ 51,252.50 \cdots $ 61,252.612].

In this we use the method 3-4-5 i.e, the data is partitioned into Equiwidth of 3,4 or intervales. This method contains the following steps.

Step-1: If the most significant bit is 3,6,7 or 9 then the data is partitioned it into 3 equiwidth intervals.

Step-2: If the most significant bit is 2,4 or 8 then the data is partitioned it into 4 equiwidth interval.

Step 3: If the most significant bit is 1,5 & 10 then the data is partitioned it into 5 equiwidth intervals.

For ex, consider the profit at all branches for all electronics company ranging from - $ 351,986 ... $ 4,700, 876·50 and also 5% & 95% percentiles are - $ 159,896 and $ 1,838,761. It contains the following steps.

Step-1: The given data MIN = -$ 351,986, MAX = $ 4,700,876·50. Then the 5% is treated as low and 95% Percentile value is treated as high value

Step-2: If we examine low and the high most Significant bit is million dollar bit position i.e, most Significant bit = $ 1,000,000 (msd) most significant Digit. Rounding low down to significant bit i·e million dollar bit position

we get low' = - $ 1,000,000

Similarly, Rounding high upto million dollar bit position we get high' = $ 2,000,000

$$\Rightarrow (high'-low')/msd = \frac{\$2000000 - (-1000000)}{1000000}$$
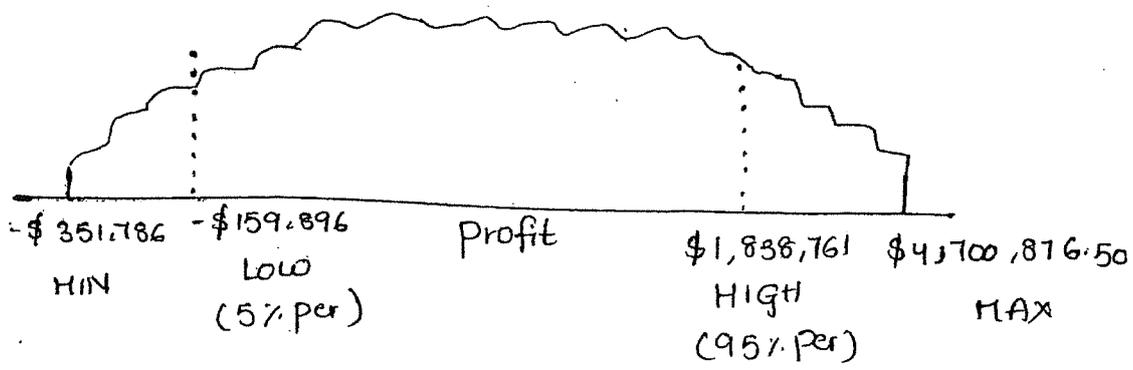
Then we get the new range. This range is partitioned into 3 equiwidth.

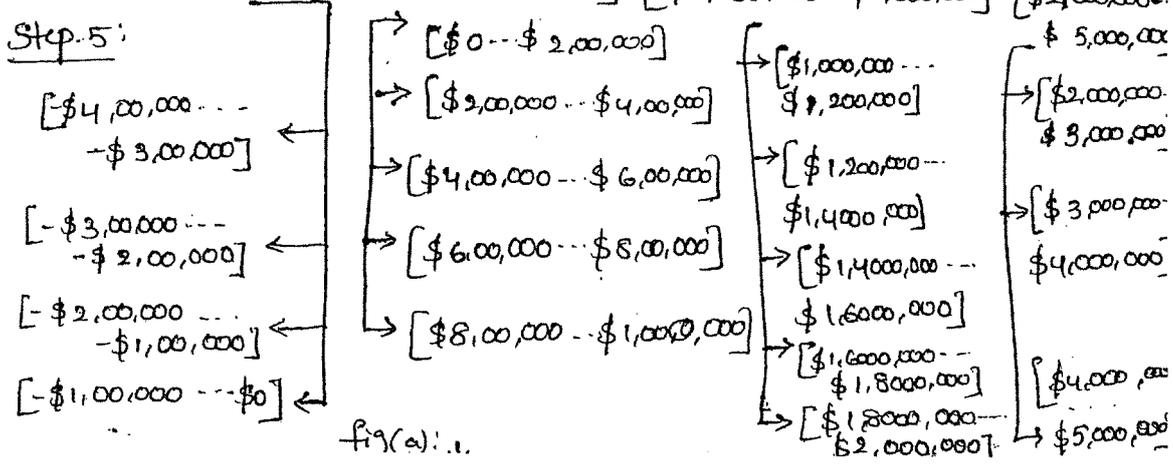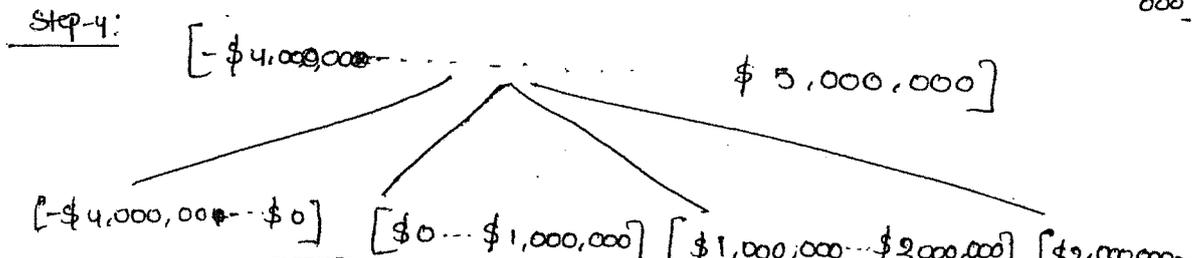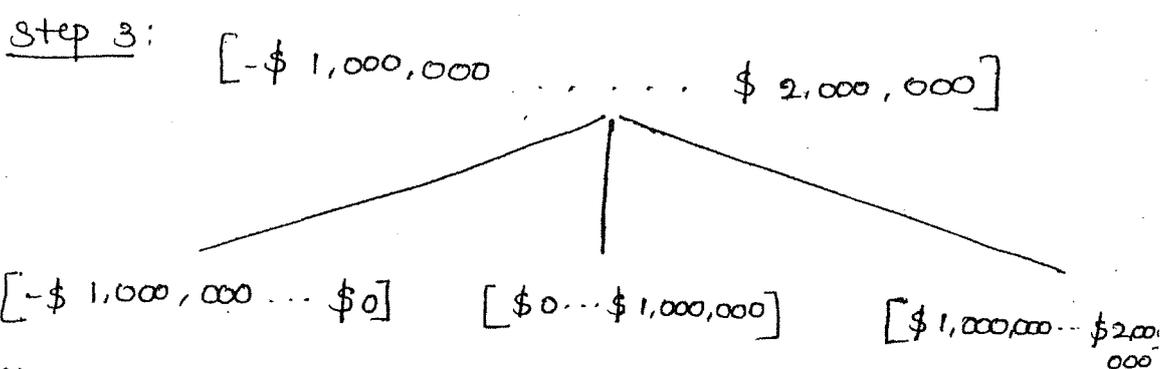Step-3: If we examine MIN and MAX then we get the MIN' = -$ 4,000,000 and MAX' = $ 5,000,000. Then we get the new range. This new range is partitioned into 4 equiwidth intervals.

Then th
is
for
6...
les
the

=
low
h value
most
, Most
ficant
1.c

Step-4: These partitions again partitioned into sub partitions to define the partition Hierarchy.

The entire process is shown below.



-$351.786    -$159.896        Profit        $1,838,761    $4,700,816.50
  MIN            Low                           HIGH           MAX
              (5% per)                        (95% per)

Step 2:     msd = $ 1,000,000
            Low' = $ 1,000,000        HIGH = $ 2,000,000.

Step 3:

[-$ 1,000,000 . . . . . . $ 2,000,000]

[-$ 1,000,000 ... $0]    [$0 ... $ 1,000,000]    [$ 1,000,000 ... $2,00,000]

Step-4:

[-$4,000,000 . . . . . . $ 5,000,000]

[-$4,000,000 ... $0]    [$0 ... $1,000,000]  [$1,000,000 ... $2,000,000]  [$2,000,000 ... $5,000,000]

Step-5:

[-$4,00,000 ... -$3,00,000]    [$0 ... $2,00,000]       [$1,000,000 ... $1,200,000]    [$2,000,000 ... $3,000,000]

[-$3,00,000 ... -$2,00,000]    [$2,00,000 ... $4,00,000]  [$1,200,000 ... $1,4000,000]   [$3,000,000 ... $4,000,000]

[-$2,00,000 ... -$1,00,000]    [$4,00,000 ... $6,00,000]  [$1,4000,000 ... $1,6000,000]  [$4,000,000 ... $5,000,000]

[-$1,00,000 ... $0]            [$6,00,000 ... $8,00,000]  [$1,6000,000 ... $1,8000,000]

                               [$8,00,000 ... $1,000,000] [$1,8000,000 ... $2,000,000]

fig(a).1.

fig(a): Concept Hierarchy for Profit based on 3-4-5 rule

## Concept Hierarchy Generation for Categorical Data:

The Categorical data means discrete. It contains the finite no: of values for ex, job category, item types etc. It is mainly classified into 2.

1. **Specification of partial ordering of attributes at the schema level by user:**

Here we specify the partial ordering of concept Hierarchy at the schema level itself. For ex, Location Dimension Contains attributes Street, city, State, country. Then the Partial.

# Chapter-4

## Getting to know your Data

### 4.1 Data objects and Attribute types

Data sets are made up of data objects. A data object represents an entity.

For ex:- In Sales DB, the objects may be customer, store items, sales;

medical DB the objects may be Patient.

Data objects are typically described by attributes. Data objects can also be referred to as samples, examples, instances, data point or objects.

If the data objects are stored in a data base they are data tuples. i.e, the rows of a DB correspond to the data objects, & columns corresponds to the attributes.

### What is an Attribute?

An attribute is a data field, representing a characteristic or feature of a data object. The nouns attribute, dimension, feature & variable are often used interchangeably in the literature. The term dimension is commonly used in DWH. In Machine learning we use the term feature, while statisticians Refer the term variable.

Observed values for a given attribute are known as observations. A set of attributes used to describe a given object is called an attribute vector.

## Types of attributes :-

The type of an attribute is determined by the set of possible values. they are

1) Nominal
2) binary
3) ordinal
4) numeric.

## 1) Nominal Attributes :-

Nominal means " relating to names". the values of a nominal attribute are symbols or names of things. Each value represents some kind of category, code(or) state. Nominal attributes are also reffered to as categorical. The values do not have any meaningful order. In computer science, the values also knowns as enumerations.

Ex:- Suppose that hair_color & marital_status are two attributes describing person objects. Possible values for hair_color are black, brown, red, gray white.

the attribute marital-status can take on the values single, married, divorced & widowed.

2, **Binary Attributes:-**

A binary att! is a nominal att! with only two categories or states: 0 (or) 1, where 0 is represent absent of an attribute, & 1 represents present of an att! Binary att! are referred to as Boolean if the two states correspond to true & false.

for ex:- Given the attribute Smoker describing a patient object, 1 indicates that the patient smokes, while 0 indicates that the patient does not.

A binary attribute is symmetric if both of its states are equally valuable and carry & carry the same weight; ie, there is no preference on which outcome should be coded as 0 or 1.

A binary attribute is Asymmetric, if the outcomes of the states are not equally important. such as the positive & negative outcomes of a medical test for HIV. By convention, we code the most important outcome, which is usually the rarest one, by 1 & the other by 0.

3) Ordinal Attributes :-

An ordinal att. is an att. with possible values
that have a meaningful order or ranking among them,
but the magnitude b/w successive values is not known

Ordinal att/s are useful for registering subjective
assessments of qualities that cannot be measured objectively.
Ordinal attributes are often used in surveys for
rating.

Ex :- In one survey, participants were asked to rate
how satisfied they were as customer. Customer
satisfication had the following ordinal categories.
0: very dissatisfied  1: some what dissatisfied
2: neutral  3: satisfied  & 4: very satisfied.

4) Numeric Attributes :-

A numeric attribute is quantitative; ie,
it is a measurable quantity, represented in integer
or real values. Numeric att's can be interval-scaled
or ratio-scaled.

Interval-Scaled Attributes

Interval-Scaled att/s are measured on a scale
of equal-size units. The values of interval-scaled

attributes have order and can be positive, 0, or -ve. In addition to providing a ranking of values, such att's allow us to compare & quantify the diff. b/w values.

Ex:

1) celsius Temperature
2) IQ (Intelligence scale
3) Time on a clock with hands.

Ratio - Scaled Attributes :-

A ratio-scaled Att: is a numeric att: with an inherent zero point. i·e, if a measurement is ratio scaled, we can speak of a value as being a multiple of another value. In addition, we can also compute the diff: b/w values, aswell as mean, median & mode.

Ex :-

1) Age, weight, Height.

2) Ruler Measurements

3) Years of education.

5). Discrete versus Continous Attributes :—

An we have organized attrs into nominal, binary, ordinal, & numeric types. There are many ways to organize att: types. The types are not mutually exclusive.

A discrete attribute :—

It has a finite or countably infinite set of values, which may or may not be represented as integers. The attri hair-color, smoker, medical-test each have a finite number of values, such as are for binary attrs and are discrete.

Continuos Attribute :—

If an attr: is not discrete, then it is continous. The terms numeric attribute & continous att: are often used interchangeably in literature. Real values are represented using a finite number of digits. Continous attributes are typically represented as floating point variables.

## 4.2. Basic Statistical Description :-

Basic statistical descriptions can be used to identify properties of the data and highlight which data values should be treated as noise or outliers. In this we discusses three areas of basic statistical descriptions. measures of central Tendency, measures of dispersion & graphic Displays.

## 4.2.1 Measuring the Central Tendency : mean, median and Mode :

A measure of central tendency is a single value that describes the way in which a group of data cluster arround a central value. In other words, it is a way to describe the centre of a data set. they there are three measures of central tendency :

1) Mean
2) Median
3) Mode.

## Mean :-

The mean is preferred measure of central tendency because it consider all the values of dataset. In order to calculate the mean, data must be numerical. You cannot use the mean for nominal data, which is data on characteristics like gender.

the mean for set of values is

$$\bar{x} = \frac{\sum_{i=1}^{N} x_i}{N} = \frac{x_1 + x_2 + \cdots + x_N}{N}$$

Ex :- Suppose we have the following values for salary, shown in increasing order:

30, 36, 47, 50, 52, 52, 56, 60, 63, 70, 70, 110

$$\bar{x} = \frac{30 + 36 + 47 + 50 + 52 + 52 + 56 + 60 + 63 + 70 + 70 + 110}{12}$$

$$= \frac{696}{12}$$

$$= 58$$

the mean salary is $58,000.

Sometimes, each value $x_i$ in a set may be associated with a weight $w_i$ for $i = 1, 2 \ldots N$. The weights reflect the significance.

$$\bar{x} = \frac{\sum\limits_{i=1}^{N} w_i x_i}{\sum\limits_{i=1}^{N} w_i} = \frac{w_1 x_1 + w_2 x_2 + \ldots w_N x_N}{w_1 + w_2 + \ldots + w_N}$$

This is called as weighted average or weighted Arithmetic mean (AM)

## Median :-

The median is the middle value for the given data, that has been arranged in order of magnitude the median

If we have 2 middle value, we can find the mean value b/w the 2 values & take it as a median.

The median is expensive to compute when we have a large number of observations.

$$\text{Median} = l + \left( \frac{\frac{n}{2} - cf}{f} \right) \times h.$$

$l =$ lowest values (lowest interval)

$\frac{n}{2} = \dfrac{\text{no. of observations}}{2}$

ef = highest frequency

ef = above the f

h = class interval.

## ModeL

A statistical term that refers to the most frequently occurring number found in a set of numbers. The mode is found by collecting and organizing the data in order to count the frequency of each result.

Ex: 8 suppose we have the following values

30, 36, 47, 50, 52, 62, 56, 60, 68, 70, 70, 70, 110

mode = 70, 52. This is bimodal.

For unimodal numeric data that are moderately skewed, we have the following expression.

$$\text{mean} - \text{mode} \approx 3 \times (\text{mean} - \text{median})$$

i.e, the mode for unimodal frequency curves that are moderately skewed can easily be approximated if the mean & median values are known.
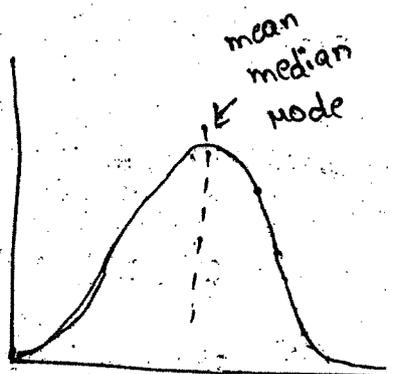
## Mid range :
~~Range~~

(the range of a set of data is the diff btw largest & smallest values)

For the above example the midrage is as follows

$$\frac{30,000 + 110,000}{2} = 70,000$$

In a unimodal frequency curve with perfect symmetric data distribution the mean, median & mode are all at the same center value.
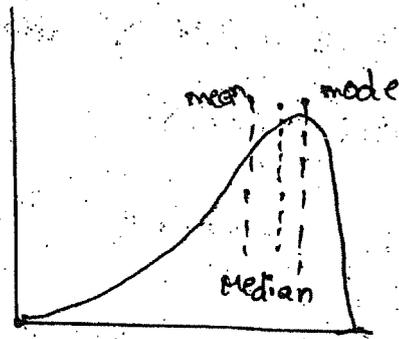
Data is most real application are not symmetric. They may instead be either "positively skewed", where the mode occurs at a value th ic, smaller than median, or negatively skewed, where the mode occurs at a value greater than the median.



(a) Symmetric data      (b) positively skewd data.

(c) Negatively skewed data.

fig(4.1) Mean, median, mode of symmetric v/s positively & negatively skewed data.

## 4.2 Measuring the Dispersion of Data :-

A measure of spread, or measure of dispersion is used to describe the variability in a sma sample or population. The measures include range, quantiles, quartiles, percentiles and the interquartile range. The five number summary, which can be displayed as a boxplot, is usefull to identifying outliers.
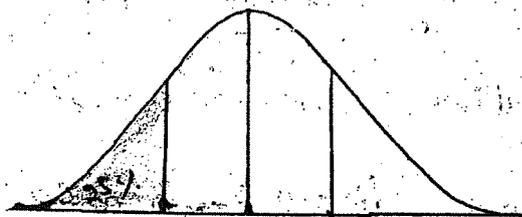
### Range, Quartiles, and Interquartiles range

let $x_1, x_2 \cdots x_N$ be a set of observations for some numeric attribute, $x$. The range of the set is the difference b/w the large & small valuel

$$\text{range} = \max(x) - \min(x).$$

Quartiles tell us about the spread of the data set by breaking the data set into quarters, just like median breaks it in half.

The First quartile($Q_1$) lies b/w the $25^{th}$ & $26^{th}$ observation, the Second quartile ($Q_2$) b/w $50^{th}$ & $51^{th}$ observation & the third quartile ($Q_3$) b/w the $75^{th}$ & $76^{th}$ observation.



$$\begin{array}{ccc} Q_1 & Q_2 & Q_3 \\ 25^{th} & \text{median} & 75^{th} \\ \text{Percentile} & & \text{Percentile}. \end{array}$$

The distance b/w the first & third Quartiles is simple measure of spread that gives the range covered by the middle half of the data. this distance is called "Interquartile range"(IQR)

$$IQR = Q_3 - Q_1.$$

**Ex:** Suppose we have have the following values for salary is

30, 36, 47, 50, 52, 52, 56, 60, 63, 70, 70, 110.

It contains 12 observations, already stored in in creasing order. thus, the quertails are for the data are 3rd, 6th & 9th values respectively.

$\therefore$ $Q_1 = 47,000$

$Q_2 = 52,000$

$Q_3 = 63000$

$IQR = 63 - 47 = 16,000.$

## Five-number Summary, Boxplots and outliers:-

The outlier is an observation that lies an abnormal distance from oa other values. A common rule of thumb for identifying suspected outliers is to single value falling at least $1.5 \times IQR$ above the 3rd Quartile or below the 1st Quartile.

The five-number summary of a distribution consists of the median $(Q_2)$, the quartiles $Q_1$ & $Q_3$, and the smallest & largest individual observations written in the order of min, $Q_1$, median, $Q_3$, max.

Boxplots are a popular way of visualizing a distribution. A boxplot incorporates the five number summary as follows:

* The ends of the box are the the quartiles, so the box lenght is interquartile range.

* The median is marked by a line within the box.

* Two lines (called whiskers) outside the box extend to the smallest & largest observations.



Fig. 3 Box plot for the unit price data for items sold at four branches of all electronics during a given time period.

# Variance and Standard Deviation:

Variance and standard deviation are measures of data dispersion. They indicate how spread out a data distribution is.

A low standard deviation means that the data observations tend to be very close to the mean, while high S.D indicates that the data are spread out over a large range of values.

The variance of N observations $x_1, x_2, \ldots x_N$ for a numeric attribute $x$ is

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i - \bar{x})^2$$

$$= \left(\frac{1}{N} \sum_{i=1}^{N} x_i^2\right) - \bar{x}^2$$

where $\bar{x}$ is the mean value of the observations.

& the S.D $\sigma$, of the observations is the square root of the variance $\sigma^2$

Ex:-   In the above observation we found $\bar{x} = 58,000$ and $N = 12$.

$$\sigma^2 = \frac{1}{12} (30^2 + 36^2 + \ldots + 110^2) - 58^2$$

$$= 379.17$$

$$\sigma = \sqrt{379.17}$$
$$\simeq 19.47$$

The basic properties of S.D $\sigma$, as a measure of spread are as follows:

1/ $\sigma$ measures spread about the mean and should be considered only when the mean is chosen as the measure of center.

2, $\sigma = 0$ only, when there is no spread, ie, when all observations have same value. otherwise $\sigma > 0$.

## 4.2.3 Graphic Displays of Basic Statistical Descriptions of Data :-

The graphic displays of basic statistical descriptions include quantile plots, quantile-quantile plots, histograms & scatter plots. Such graphs are useful for the visual inspection of data, which is useful for data preprocessing. First three of these show univariate distributions, and while scatter plots show bivariate distributions.

Quantile plot:-

A quantile plot is a simple, we cover and effective way to have a first look at a univariate

data distribution. It displays all of the data for the given attribute. then, it plots quantile information.

Ex! let $x_i$, for $i=1$ to $N$, be the data sorted in increasing order so that $x_1$ is the smallest element & $x_N$ is the largest for some ordinal or numeric attribute $x$.

A set of unit price Data for Items sold at a Branch of All Electronics

| unit price ($) | count of Items sold |
|---|---|
| 40 | 275 |
| 43 | 300 |
| 47 | 250 |
| 74 | — |
| 75 | 360 |
| 78 | 515 |
| — | 540 |
| 115 | — |
| 117 | 320 |
| 120 | 270 |
| | 350 |

Table 2.1

A quantile plot for the unit price data of Table 21.

## Quantile - Quantile plot :-

A quantile - quantile plot, or q-q plot, graphs the quantiles of one univariate distribution against the corresponding quantiles of another. It is a powerful visualization tool in that it allows the user to view whether there is a shift in going from one distribution to another.

## Histograms:-

A histogram is an accurate representation of numerical data. To construct a histogram for an attribute X, partitions the data distributions of X into buckets. Each bucket is represents only a single attribute of value/frequency pair.



A histogram for table 4.1

## Scatter plots and Data Correlation:-

A scatter plot is one of the most effective graphical methods for determining if there appere to be a relationship, pattern or trend b/w two numeric attributes.

To construct a scatter plot, each pair of value is treated as a pair of coordinates in an algebraic and plotted as points in the plane.

scatter plot for the table 2.1

The scatter plot is useful to represent the bivariate data to clusters of points & outliers, or to explore the possibility of correlation relations. Two attributes, X, and Y, are correlated if one attribute implies to other. Correlation can be +ve, -ve or null



(a)
Positive

(b)
negative.

Scatter plots can be used to find (a) positive (or) (b) negative b/w attributes.

# ★ 4.3 Data Visualization

Data visualization aims to communicate data clearly and effectively through graphical represents. Data visualization has been used extensively in many applications. f

for ex! at work for reporting, managing business operations, and tracking progress of tasks.

More popularly, we can take advantage of visualization techniques to discover data relationships that are otherwise not easily observable by looking at the raw data.

## 4.3.1 pixel oriented visualization Techniques :-

A simple way to visualize the value of dimension is use a pixel where the colors of the pixel reflect the dimension's value. For the data set of m dimensions, pixel-oriented technique creates m windows on the screen, one for each dimension.

Ex! AllElectronics maintains a customer info table which consists of income, credit limit, transaction

volume and age. we can analyze correlation b/w income and the other attributes by visualization?

we can sort all the customers in income ascending order, & used this order to layout the customer data in four visualization windows, as follows.



(a) income   (b) credit limit   (c) transaction - volume   (d) age.

Pixel oriented visualization of four att.

credit limit increases as income increases; whose income is in the middle range are more likely to purchase more from All Electronics; there is no clear correlation b/w income & age.

Filling a window by laying out the data records in a linear way may not work well for a wide window. we can lay out the data records in a space-filling curve to fill the windows. A It is a curve with a range that covers the entire n-dimensional unit hyper cube. Since the visualization

windows are 2-D, we can use any 2-D space filling curve as follows.



The circle segment technique. (a) Representing a data record in circle segments (b) laying out pixels in circle segments.

## 4.3.2 Geometric projection Visualization Techniques:-

Geometric projection techniques help users find interesting projections of multidimensional data sets. The central challenge the geometric projection techniques try to address is how to visualize a high diem dimensional space on a 2-D Display.

A Scatter plot displays 2-D data points using Cartesian coordinates. A third dimension can be added using different colors or shapes to shapes to represent diff. data points.

Visualization of a 2-D data set using a scatter plot.

where x and y are two spatial attributes and the third dim; is represented by different shapes. Through this visualization, we can see that the points of type "+" and "X".

A scatter plot matrix technique is a useful extension to the scatter plot. For an n-dimensional data set, a scatter plot matrix is an n×n grid of 2-D scatter plots that provides a visualization of each dim; with every other dim; The scatter plot matrix becomes less effective as the dimensionality increases.

## 4.3.3 :- Icon Based Visualization Techniques

Icon-based visualization techniques use small icons to represent multidimensional data values. There are two popular icon-based techniques chernoff faces and stick figures.

Chernoff faces were introduced in 1973 by statistician Herman chernoff. They display multidimensional data of up to 18 variables as a cartoon human face.



Chernoff faces. Each face represents an n-dimensional

data point $(n \leq 18)$

## 4.3.4 Hierarchical Visualization Techniques :-

The visualization techniques discussed so far focus on visualizing multiple dimensions simultaneously. However, for a large data set of high dimensionality, it would be difficult to visualize all dimensions at the same time. Hierarchial visualization techniques partition all dimensions into subsets. The subsets are visualized in a hierarchical manner.

words "worlds - within - worlds" also known as n-vision, is a representative hierarchical visualization method. Suppose we want to visualize a 6-D dataset, where the dimensions are $F, X_1 \cdots X_5$. we can first fix the values of dimensions $X_3, X_4, X_5$ to some selected values, say $C_3, C_4, C_5$. we can then visualize $F, X_1 X_2$ using 3-D plot.

## 2.3.5 Visualizing Complex Data and Relations:-

Visualization techniques were mainly for numeric data. Recently, more and more non-numeric data, such as text and social networks, have become available. for ex: many people on the web tag various objects suchas pictures, blog entries and product reviews.

L.

## 4.4 Measuring Data Similarity and Dissimilarity :-

In data mining applications, such as clustering, outliers analysis, and nearest-neighbour classification, we need ways to assess how alike or unalike objects are in comparison to one another.

for ex: a store may want to search for clusters of customer objects, resulting in group of customers with similar characteristics. (income, are of residence & age). such information can then be used for marketing.

A cluster is a collection of data objects such a cluster are similar to one the objects.

so employs clustering-based potential outliers as objects to others.

ct similarities can also be used ation schemes where a given label based on its similarity the model.

SECTION - A

Answer all questions.                    (4×15=60)

1.
   a.  What is multi - dimensional analysis? How it is
       implemented through OLAP? Explain various types
       of OLAP systems used for multi - di..
       data...

## 4.4.1 Data Matrix Versus Dissimilarity Matrix :-

We studying the central tendency, dispersion & spread of observed values for some attribute x. that are one-dimensional, i.e, described by a single attribute. We talk about objects described by multiple attributes. Therefore, we need a change in notation.

Suppose that we have n objects (ex: persons, items or courses) described by p attributes (age, height, weight or gender). The objects are $X_1 = (x_{11}, x_{12}, \cdots x_{1p})$, $X_2 = (x_{21}, x_{22}, \cdots x_{2p})$ & and so on, where $x_{ij}$ is the value for object $x_i$ of the $j^{th}$ attribute.

Main memory-based clustering and nearest-neighbour algorithms typically operate on eighter of the following two data structures.

1) Data matrix (object-by-attribute structure) :—

This is also called as object-by-attribute structure. This structure stores the n data objects in form of relational table, or $n \times p$ matrix ( n objects × p attributes).

$$\begin{bmatrix} x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ x_{n1} & \cdots & x_{nf} & \cdots & x_{np} \end{bmatrix}$$

2) Dissimilarity matrix :—

This is also called as object-by-object structure. This structure stores a collection of proximities that are available for all pairs of n objects. It is also represented as n×n table.

$$\begin{bmatrix} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ \vdots & \vdots & \vdots & \\ d(n,1) & d(n,2) & \cdots & 0 \end{bmatrix}$$

where $d(i,j)$ is the measured dissimilarity (or) difference b/w objects $i$ & $j$. $d(i,j)$ is a non-ve number that is close to 0 when object $i$ & $j$ are highly similar or near each other.

$d(i,i) = 0$ i.e, the difference b/w an object and itself is 0.

Measure of similarity can be expressed as a function of measures of dissimilarity.

Ex :— for nominal data

$$sim(i,j) = 1 - d(i,j)$$

where $sim(i,j)$ is the similarity b/w objects $i$ and $j$.

A data matrix is made up of two entities they are rows (for objects) and columns (for attributes) therefore, the data matrix is often called as "two-mode" matrix. The dissimilarity matrix contains one kind of entity & so is called as "one-mode" matrix.

## 4.4.2 Proximity measures for Nominal Attributes:—

A nominal attribute can take on two or more states. for example: map-color is a nominal attribute that may have say five states: red, yellow, green, pink and blue.

let the no. of states of a nominal attribute be M. the states can be denoted by letters, symbols or a set of integers, such as $1, 2, \ldots M$.

### dissimilarity b/w objects described by nominal attributes

The dissimilarity b/w two objects i and j can be computed based on the ratio of mismatches:

$$d(i,j) = \frac{p - m}{p}$$

where m is the no. of matches i.e, no. of attributes for which i & j are in the same state. and p is the total no. of attributes describing the objects.

Ex: Suppose that we have the sample data of as follows, expect that only the object-identifier and the attribute test-1 are available, where test-1 is nomial.

Table 4.2. A simple Data Table Containing Attributes of mixed type.

| Object Identifier | test-1 (nominal) | test-2 (ordinal) | test-3 (numeric) |
|---|---|---|---|
| 1 | codeA | excellent | 45 |
| 2 | code B | fair | 22 |
| 3 | code c | good | 64 |
| 4 | code A | excellent | 28 |

lets compute the dissimilarity matrix.

$$\begin{bmatrix} 0 & & & \\ d(2,1) & 0 & & \\ d(3,1) & d(3,2) & 0 & \\ d(4,1) & d(4,2) & d(4,3) & 0 \end{bmatrix}$$

we set p=1, so that d(i,j) evaluates to 0 if the i & j are match, & = 1 if the objects differ. Thus we get

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

we see all the objects are dissimilar except objects 1 and 4 i.e (d(4,1) = 0.

Alternatively, similarity can be computed as

$$\text{sim}(i,j) = 1 - d(i,j) = \frac{m}{p}$$

## 4.4.8 :- Proximity Measure for Binary Attributes:-

A binary attribute has only one of two states: 0 or 1 where 0 means that the attribute is absent, and 1 means that is present. Treating binary attributes as if they are numeric can be misleading. Therefore, methods to specific to binary data are necessary for computing dissimilarity.

One approach involves computing a dissimilarity matric from the given binary data. If all binary attrs are thought of a having the same weight, we have the 2x2 contingency table as follows

|  |  | object j | | |
|---|---|---|---|---|
|  |  | 1 | 0 | sum |
| object i | 1 | q | r | q+r |
|  | 0 | s | t | s+t |
|  | sum | q+s | r+t | p |

contingency Table for Binary Attributes.

where '$q$' is the number of attributes that equal to 1 for both objects $i$ & $j$, '$r$' is the no. of attributes that equal '1' for object $i$, but equal '0' for object $j$, $s$ is the number of attributes that equal '1' for object '$i$' but equal 0 for object $j$, and '$t$' is the no. of attributes that equal '0' for both objects $i$ & $j$. The total no. of attributes is $P$

$$P = q + r + s + t.$$

For Symmetric binary attributes, each state is equally valuable. Dissimilarity that is based on symmetric binary attributes is called symmetric binary dissimilarity. If objects i and j are described by symmetric binary attributes, then the dissimilarity b/w $i$ & $j$ is

$$d(i,j) = \frac{r+s}{q+r+s+t}$$

For asymmetric binary attributes, the two states are not equally important, given two asymmetric binary attributes, the agreement of two 1's is then considered more significant than that of two 0's. Therefore, such binary attributes are often considered "monary" (having one state). The dissimilarity based on these 0's is called asymmetric binary dissimilarity.

$$d(i,j) = \frac{r+s}{q+r+s}$$

where the no of negative matches t, is considered unimportant and is thus ignored.

when both Symmetric & asymmetric binary attributes occur in the same data set, the mixed attributes approach can be applied.

Ex⌐

Suppose the patient record table contains the attribute name, gender, fever, cough, test-1, test-2, test-3, & test-4.

where name is an object identifier, gender is a symmetric attribute, and the remaining attributes are asymmetric binary.

| name | gender | fever | cough | test-1 | test-2 | test-3 | test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Jim | M | Y | Y | N | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

4.4 Relational Table where patients are described by Binary Attributes.

for Asymmetric attribute values, the values Y and P be set to 1, and the value N (no or negative) be set to '0'.

The distance b/w each pair of the three patients Jack, Mary & Jim is

$$d(Jack, Jim) = \frac{1+1}{1+1+1} = 0.67$$

$$d(Jack, Mary) = \frac{0+1}{2+0+1} = 0.33$$

$$d(Jim, Mary) = \frac{1+2}{1+1+2} = 0.75.$$

These measurements suggest that Jim & Mary are unlikely to have a similar disease. and Jack and Mary are the most likely to have a similar disease.

## 4.4.4 Dissimilarity of Numeric Data: Minkowski Distance

We describe distance measures that are commonly use for computing the dissimilarity of objects by numeric attributes. These measures include the Euclidean, Manhattan and minkowski distances.

In some cases, the data are normilized before applying distance calculations. This involves transforming the data to fall within a smaller or common range such as [-1, 1] or [0.0, 1.0].

The most popular distance measure is "Euclidean distance". let $i = (x_{i1}, x_{i2} \cdots x_{ip})$ and $j = (x_{j1}, x_{j2} \cdots x_{jp})$ be two objects described by $P$ numeric attributes. The Euclidean distance b/w objects $i$ & $j$ is

$$d(i,j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \cdots + (x_{ip} - x_{jp})^2}$$

ex: let $x_1 = (1,2)$ & $x_2 = (3,5)$, the Euclidean distance b/w two objects is

$$d(x_1, x_2) = \sqrt{(1-3)^2 + (2-5)^2} = \sqrt{2^2 + 3^2} = 3.61.$$

Another well-known measure is the Manhattan distance or city block distance, name so because it is the distance in blocks b/w any two points in a city.

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \cdots + |x_{ip} - x_{jp}|.$$

ex:
Manhattan distance for the above example

$$d(x_1, x_2) = |1-3| + |2-5| = 2+3 = 5.$$

Both the Euclidean and the Manhattan distance satisfy the following mathematical properties.

1) Non-negativity: $d(i,j) \geq 0$. : Distance is a non-ve number.

2) Identity of indiscernibles : $d(i,i) = 0$ the distance of an object to itself is '0'.

3) Symmetry: $d(i,j) = d(j,i)$ : Distance is a symmetric function.

4) Triangle inequality: $d(i,j) \leq d(i,k) + d(k,j)$.

"Minkowski distance" is a generalization of the Euclidean and Manhattan distances.

$$d(i,j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \cdots + |x_{ip} - x_{jp}|^h}$$

where h is a real number such that $h \geq 1$. It represent the manhattan distance when $h=1$ and Euclidean distance when $h=2$.

The supremum distance is a generalization of the Minkowski distance for $h \to \infty$. To compute it, we find the att. $f$ that gives the maximum difference in values b/w two objects.

$$d(i,j) = \lim_{h \to \infty} \left( \sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{1/h}$$

$$= \max_{f}^{p} |x_{if} - x_{jf}|$$

$x_2 = (3,5)$

Euclidean distance
$= (2^2 + 3^2)^{1/2} = 3.61$

Manhattan distance
$= 2 + 3 = 5$

Supremum distance
$= 5 - 2 = 3$.

Euclidean, Manhattan, supremum distance between two objects

## 4.4.5 Proximity Measure for Ordinal Attributes :—

The values of an ordinal attributes have a meaningful order or ranking. The treatment of ordinal attributes is quite similar to that of numeric attributes when computing dissimilarity b/w two objects.

Suppose $f$ is an attribute from the set of ordinal attributes describing $n$ objects. The dissimilarity computation w.r.to $f$ involves the following steps:

1) The value of $f$ for $i^{th}$ object is $x_{if}$, and $f$ has $M_f$ ordered states, representing the ranking $1, \ldots M_f$. replace each $x_{if}$ by its corresponding rank $r_{if} \in \{1 \ldots M_f\}$

2) We perform data normalization by replacing the rank $r_{if}$ of the $i^{th}$ object in the $f^{th}$ attribute by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

3) dissimilarity can be computed using any of the distance measures described in numeric attributes, using $z_{if}$ to represent the $f$ value for $i^{th}$ object.

# 4.4.6 Dissimilarity for Attributes of Mixed Types:-

One approach is to group each type of attribute together, Performing separate data mining analysis for each type eg! clustering. This is feasible if these analyses derive compatible results.

Suppose that the data set contains p attributes of mixed type. The dissimilarity blw i & j objects is defined as

$$d(i,j) = \frac{\sum_{f=1}^{p} \delta_{ij}^{(f)} \cdot d_{ij}^{(f)}}{\sum_{f=1}^{p} \delta_{ij}^{f}},$$

where the indicator $\delta_{ij}^{(f)} = 0$ if either

1) $x_{if}$ or $x_{jf}$ is missing. (or)

2) $x_{ij} = x_{jf} = 0$ and attribute 'f' is asymmetric binary.

Other wise, $\delta_{ij}^{(f)} = 1$.

The contribution of attribute f to the dissimilarity blw i & j i.e, $d_{ij}^{(f)}$ is computed dependent on its type:

* If f is numeric $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{max_h x_{hf} - min_h x_{hf}}$,

where h runs over all non missing objects for att. f.

* If f is nominal or binary: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$; otherwise $d_{ij}^{(f)} = 1$.

\* If $f$ is ordinal: compute the ranks $r_{if}$ and $z_{if} = \dfrac{r_{if} - 1}{M_f - 1}$, and treat $z_{if}$ as numeric.

## 4.4.7 Cosine Similarity :-

A document can be represented by thousands of attributes, each recording the frequency of a particular word or phrase in the document. Thus, each document is an object represented by a "term-frequency vector".

Cosine Similarity is a measure of similarity that can be used to compare document (or) give a ranking of documents w.r. to a given vector of query words. Let $x$ & $y$ be two vectors for comparison. Using cosine measure as a similarity function, we have

$$sim(x, y) = \frac{x \cdot y}{\|x\| \|y\|}.$$

where $\|x\|$ is the Euclidean norm of vector $x = (x_1, x_2 \ldots x_p)$ defined as $\sqrt{x_1^2 + x_2^2 + \ldots + x_p^2}$. $\|y\|$ $\|y\|$ is the Euclidean norm of vector $y$.

Ex :

- suppose that $x$ & $y$ are the first two term-frequency vectors in the below table. i.e $x = (5, 0, 3, 0, 2, 90, 2, 0, 0)$ & $y = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$

## Document Vector or Term-Frequency Vector

| Document | team | coach | hockey | baseball | soccer | penalty | score |
|---|---|---|---|---|---|---|---|
| Document 1 | 5 | 0 | 3 | 0 | 2 | 0 | 0 |
| Document 2 | 3 | 0 | 2 | 0 | 1 | 1 | 0 |
| Document 3 | 0 | 7 | 0 | 2 | 1 | 0 | 0 |
| Document 4 | 0 | 1 | 0 | 0 | 1 | 2 | 2 |

| win | loss | season |
|---|---|---|
| 2 | 0 | 0 |
| 1 | 0 | 1 |
| 3 | 0 | 0 |
| 0 | 3 | 0 |

compute the cosine similarity b/w 2 vectors $x$ & $y$

$$x^t \cdot y = 5 \times 3 + 0 \times 0 + 3 \times 2 + 0 \times 0 + 2 \times 1 + 0 \times 1 + 0 \times 0 + 2 \times 1$$
$$+ 0 \times 0 + 0 \times 1 = 25$$

$$\|x\| = \sqrt{5^2 + 0^2 + 3^2 + 0^2 + 2^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2} = 6.48$$

$$\|y\| = \sqrt{3^2 + 0^2 + 2^2 + 0^2 + 1^2 + 1^2 + 0^2 + 1^2 + 0^2 + 1^2} = 4.12$$

$$\text{sim}(x, y) = \frac{25}{6.48 \times 4.12}$$
$$= 0.94$$

Therefore, if we were using the cosine similarity measure to compare these documents, they would be considered quite similar.

# * Architecture of Data mining System:

To get the efficient architecture of DM System. The DM System must be integrated with db (or) DWH. If the DM System integrated with db or DWH then, we have to consider the following coupling schemas.

① No coupling
② loosly coupling
③ Semitight Coupling
④ Tight Coupling.

① No Coupling :- Here, DM System doesnot integrate with db or DWH. These types of systems contains the following drawbacks.

① If the DM System doesn't integrate with the db then it doesn't provide any Performance, Scaliability etc...;

② If the DM System doesn't integrate with the DWH then it doesnot provide data cleaning, data transformation, Data integration etc...;

These Systems extract the data from the files because these Systems

doesn't contain any data structures or any algorithms & the results are send to file.

∴ No coupling means poor design.

2. Loosely Coupling:

Here, DM system extract the essential features from db or DWH. These systems extract the data from centralized depository & the results are send to any file or db or DWH. These systems far better than the no of coupling systems.

3. Semitight Coupling:

Here, DM Systems integrate with db or DWH to get the DM primitives. i.e., store, retrieve, histogram analysis, join indexing etc..,

4. Tight Coupling:

Here, DM Systems fully integrated with db or DWH to get all the features & to provide the integrated information environment. This is the desired architecture. finally using this tight coupling we design the typical architecture of DM.

fig (a): Typical Architecture of DM System

The typical architecture of DM System contains the following components.

① DB, DWH (or) Any information Repository:-
   Here, we extract the data from any db (or) DWH or any information Repository then we apply the data cleaning, data transformation, data integration & finally data load. This data is loaded it into db & DWH Server.

② DB (or) DWH Server:-
   It contains the data according to user Specification.

③ Knowledge Base :-

Here, end user can get the knowledge from the different data sources. If the end user of the knowledge then the DM process is Simplified.

④ Data mining engine :-

Here, we apply the DM functions.

⑤ Pattern evolutions :-

Here, interesting patterns are evaluated.

⑥ GUI :-

It provides the communication b/w end user & DM System. i.e.., through this end user Present the queries to the DM System, and also discovered Patterns are visualized to the end user by using different visualization Techniques.

Concept Description : characteristization & comparison :-

The DM is mainly classified into two

1. Descriptive mining

2. Predictive mining

The descriptive mining analyses the collected data set in the form of summeriz-ation. i.e.., it explains the general charact-eristics of data.

The Predictive mining analyses the data in order to construct one or more models using this models we can predict the behaviour of a new data set.

The descriptive mining is called as Concept description.

what is Concept description :-

The concept is nothing but the collec-tion of data. For eg.., frequently-buys, graduate-students etc...., The concept description summerizes and explains the general characteristics of data. The Concept description is also called as class description. It contains the two types.

1. class characterization

2. Class Comparison

The class characterization explains the characteristics of collected data.

The class Comparison Compares the two or more collections of data.

# Comparison b/w OLAP & Concept Description:

**a) Aggrigate functions vs Complex datas.**

The DWH & OLAP Systems the data is stored in the form of multi dimensional data model. i.e.., Data cubes. The data cubes contains the dimensions & measures. The measures are analysed by using the aggregate functions, i.e.., Sum(), Avg() etc.., But the Concept description in DB handle the complex data types like non numeric, text, Special, multi media etc...,

**b) Manual vs automations.**

The DWH and OLAP Systems we have to specify the dimensions & also we have to specify the OLAP operations. But the Concept description in db we select the dimensions & also we select the data at multiple levels.

## Data Generalization & Summerization Based characterizatione:-

For eg..., the data is store in large db. then, this data generalization provides the data abstraction from low

Conceptual level to high conceptual

The data generalization mainly classified into two.

(1) Data cube approach - we have to see the 2<sup>nd</sup> chapter.

(2) Attribute Oriented Induction (AOI) Approach.

(a) Attribute Oriented Induction (AOI) Approach:

This is introduced in the year of 1989. This is the best Approach for data generalization & Summerization based characterization. It contains the two steps.

1. Select Task Relevant data.
2. Data Generalization.

1. Select Task Relevant data:

For eg... end user analyse the characteristics of graduate students from the Big university - DB & the attributes are name, gender, major, Birth-Place, Birth-date, Residence & Phone number.

Note: The data cube approach is based on DWH Orientation. But the AOI approach based on Relational Database.

... on, this is specified in DMQL

Syntax :-

```
use Big-university - DB
mine characteristics as "science-students"
in relevance to name, gender, major,
birth-place , birth date
residence, phone #
from student
where status in "graduates"
```

This is transformed into relation query

Syntax :.

```
use Big-university - DB
select name, gender, major, birth-place, birth-date,
residence, phone #.
from student
where status in {"dica", "MBA", "MA", "Ccom"}
```

Once this query is executed, we get the table. This table is called as Task relevant data or initial working relation. This is shown in below.

| name | gender | major | birth-place | birth-date | residence | phone # |
|------|--------|-------|-------------|-----------|-----------|---------|
| John | cl | Science | Hyd, AP India | 12/07/76 | 123.Bhil Hyd | 234576 |
| Lora | F | Engineering | vancouver BC, Canada | 10/08/70 | Mainst vanco -uver | 977786 |

Data Generalisations:

Again it contains two steps.
① Attribute Removal
② Attribute Generalization.

① Attribute Removal :-
It contains the two rules.

Rule 1:- If the attribute as large no of distinct values and no hierarchy is defined then remove that attribute from the relation.

Rule 2:- If the attribute as high concept hierarchy is defined then remove the low concept hierarcy & replaced with high concept hierarchy.

② Attribute Generalization :-
If the attribute as the large no of distinct values & concept hierarchy is defined then we select the generalization threshold value. After applying the AOI rules in Table 1.1 we get the following steps.

① name :- If contains the large no of distinct values & no concept hierarchy is defined. ∴ Remove that attribute from the working relation.

② Gender :- It contains 2 values and it remains the same.

③ major :- It contains 3 values i.e., Science engineering, & business. It remains

the Same.

④ birth-Place :

It's Concept hierarchy is defined city, state & country then this is Replace with birth country

⑤ birth-date :

This attribute is transformed into age & age is transformed into age-range.

⑥ Residence : It's Concept hierarchy is defined Street to city. This attribute is replaced with residence city.

⑦ Phone # : It contains the large no of distinct values & no concept hierarchy is defined.

∴ Remove it from the Relation.

Here, Similar tuples are merged & specified by using the count attribute

Then, finally we get the table 1.2

| gender | major | birth-country | age-range | Residence city | count |
|---|---|---|---|---|---|
| M | Science | India | 20...30 | Hyd | 20 |
| F | engineering | Canada | 30...40 | vancouver | 30 |

Table (1.2)  After applying AOI Rules, the final Relational Table.

* **Efficient Implementation of attribute oriented Induction:**

For efficient implementation of AOI we use the following algorithm.

**Algorithm:- Attribute · oriented · Induction.**

i/p :- i) DB, a Relational db
   ii) DMQueries , A DM Query.
   iii) A-list, list of Attributes
   iv) Gen(ai), It is the Generalization
      operators · for each ai.
   v) A-thersh · value (ai), It is the
      Generalization Threshold value for
      each ai.

o/p :-   P ← Prime · Relation

methods :- W ← task- relevant-data (DB-
                     DM Query).

Here 'W' is the initial working relatio

① P ← prepare · for · generalization (W)
  (a) If the attribute ai contains the larg
  no of distinct values & no concept
  hierarchy is defined then remove
  that attribute from the working relati

  (b) If the attribute 'ai' high level conce
  hierarchy is defined then remove
  the low level concept hierarchy
  & replace with high level concept
  hierarchy.

| ty | count |
|----|-------|
|    | 20    |
| er | 30    |

P← Generalization(b):

(a) If the attribute contains the large no of distinct values & concept hierarchy is defined. then we select the Threshold range.

(b) If the generalization contains the similar tuples these are merged by using the count attribute.

Simple AOI Algorithm.

→ * Presentation of Derived Generalization:

we apply the AOI on Relational DB then we get the generalization Relation. This is Presented at end user. by using several visualization techniques i.e., cross Table, Barcharts, Piecharts etc.,

for eg., Consider the sales table for all electronics is shown in below.

| location | item | Sales (in dollars) | count (in thousands) |
|---|---|---|---|
| Europe | TV | 15 | 200 |
| Asia | TV | 12 | 350 |
| North America | TV | 28 | 400 |
| Europe | Computer | 120 | 1000 |
| Asia | ,, | 150 | 1200 |
| North America | ,, | 250 | 1800 |

Sales of all electronics

This table information a represented in cross table.

| location item | TV | | Computer | |
|---|---|---|---|---|
| | Sales | Count | Sales | Count |
| Europe | 15 | 200 | 120 | 1000 |
| Asia | 12 | 350 | 130 | 1200 |
| North America | 18 | 400 | 250 | 1800 |

fig(b) . cross Table.



☐ Europe
◼ Asia
▨ North America.

fig(c) Bar chart .

Pie chart for TV sales:



26% Europe
Asia 22%
52% North America

Pie chart for Computer Sales:



Europe 25%
25% Asia
50% North America

# * Analytical Characterization :-

using this attribute relevance we find the weekly relevant & irrelevant attributes. These are removed from concept description. using this we also find the most relevant attributes these are included in concept description.

why attribute relevance is acquired :-

The DWH & OLAP systems Contains the drawbacks of enduser has to Specify the dimensions & also he has to Specify the high Conceptual level. This is Specified by using stmt " generalize dimension location to the country level".

The enduser does not know the attrib-ute acquired to achieve the DM task then he Specify all the attributes by using the stmt in relevance to *". But, this does not give the accurate data -:. To identify mostly relevant attributes & delete weekly relevant attributes. we acquire the attribute relevance analysis.

## methods for Attribute Relevance Analysis:-

In this method we integrate the information gain analysis with

Finally, we find the information gain for each attribute & then the highest information gain attribute is added in concept description & remove the lowest information gain attribute from the Concept description.

For eg.., Consider the training samples 's' & also each training sample class-label must be node. To identify the class label we use one attribute for eg.., Consider the attribute Status using this attribute we find the whether he is graduate student or under graduate student.

Let us consider the 'm' classes the sample 's' contains 'si' samples in class ci for $i = 1, 2, \ldots m$ then the sample belong to class 'ci' with the probability $\frac{s_i}{s}$ where s - Total no of samples then the expected information needed to classify the given sample.

$$I(S_1, \ldots S_m) = -\sum_{i=1}^{m} \frac{s_i}{s} \log_2 \frac{s_i/s}{} \quad \text{——} \quad ①$$

Let us consider the attribute partition values $\{a_1, a_2, \ldots a_v\}$ then th

attribute partition the given set 'S' into $\{S_1, \dots S_v\}$ then Let us consider the $S_j$ falling $c_{ij}$ then entrophy of A.

$$E(A) = \sum_{j=1}^{v} \frac{(S_{ij} + \dots + S_{mj})}{S} \times I(S_{ij}, \dots, S_{mj}) \quad \text{(2)}$$

finally, information $gain(A) = I(S_1, \dots S_m) - E(A)$

This method contains the following steps.

1. **Data Collection:-**

we collect the data for target class & contrasting class. The target class Contains the data that is to be characterised. The contrasti -ng the class Contains the comparitive data ~~that is to be characterised. The C~~

2. **Apply AOI Rules :-**

Here, we apply the AOI Rules.

3. **Remove irrelevant & weakly Relevant Attributes**
It contains three steps.

Step 1:- Compute the expected information gain

Step 2:- Compute the entrophy value.

Step 3:- Compute the information gain. The highest information gain attributes is added in Concept description.

**Example for class characterisation:-**

It contains the following steps

Step 1:- collect the data for target class. i.e...

graduate students. This is shown in below.

| gender | major | birth-country | age-range | count |
|--------|-------|---------------|-----------|-------|
| M | Science | CANADA | 20.......25 | 24 |
| F | Science | CANADA | 20.......25 | 30 |
| M | Science | USA | 20.......25 | 30 |
| M | Engineer | CANADA | 25.......30 | 20 |
| F | Engineer | USA | 25.......30 | 16 |

fig(a) Data for target class i.e.., graduate students.

Step 2:- Collect the data for Contrasting class i.e...., under graduate students.

| gender | major | birth-country | age-range | Count. |
|--------|-------|---------------|-----------|--------|
| M | Science | CANADA | 20.......25 | 20 |
| F | Science | USA | 20.......25 | 22 |
| M | Engineer | CANADA | 25.......30 | 20 |
| F | Engineer | USA | 25.......30 | 26 |
| M | Business | CANADA | 30.......35 | 20 |
| F | Business | USA | 30.......35 | 22 |

fig(b) Data for Contrasting class i.e.., under graduate students.

we already apply the AOI Rules & Remove the name, phone # because they contains the large no of distinct values & no concept hierarchy is defined.

The graduate student class is represented by '$s_1$' & under graduate student class is represented by '$s_2$'. The '$s_1$' contains the 120 tuples & '$s_2$' contains the 130 tuples. ∴ expected information needed $I(s_1, s_2)$

$$I(s_1, s_2) = -\frac{120}{250} \log_2 {120}/{250} - \frac{130}{250} \log_2 {130}/{250}$$

$$= 0.99$$

Then consider the attribute major, for example.., major = "Science".

$s_{11} = 84$, $s_{21} = 42$ then total − 126

expected information needed

$$I(s_{11}, s_{21}) = -\frac{84}{126} \log_2 {84}/{126} - \frac{42}{126} \log_2 {42}/{126}$$

$$= 0.78$$

major = "Engineer".

$s_{12} = 36$, $s_{22} = 46$ then total 82

expected information needed

$$I(s_{12}, s_{22}) = -\frac{36}{82} \log_2 {36}/{82} - \frac{46}{82} \log_2 {46}/{82}$$

major = "Business"

$s_{13} = 0$, $s_{23} = 42$ then

expected information needed.

$$I(s_{13}, s_{23}) = 0$$

Then, entropy value for major

$$E(major) = \frac{126}{250} \times I(s_{11}, s_{21}) + \frac{82}{250} I(s_{12}, s_{22}) +$$

$$\frac{42}{250} \, \mathcal{I}(s_{13}, s_{23})$$

$$= \frac{126}{250} \times 0.78 + \frac{82}{250} \times 0.98 + \frac{42}{250} \times 0$$

$$= 0.77$$

∴ information gain for major.

$$G(major) = \mathcal{I}(s_1, s_2) - E(major)$$

$$= 0.99 - 0.77 = 0.22.$$

Similarly, we can find the information gain for other attributes.

i.e....
$$G(gender) = 0.003$$
$$G(birth-country) = 0.004$$
$$G(age-range) = 0.59$$

for eg...., Consider the attribute - relevance Threshold is "0.1" then less than & equal to this value. Consider as the weakly relevant attributes.

∴ The attributes major & age-range is consider as the most relevant attributes these are included in Concept description. Similarly, the attributes gender & birth-country are the weakly relevant attributes these are removed from Concept description.

* Class Comparison : Discrimination b/w different class :-

The end user always have interest or mining the data from one or more Contrasting

classes. This is called as class Comparison or class Discrimination.

Rather than extracting the data from single class . i.e...., class characterization. The target class & Contrasting class must have the similarities . for eg..., Person & item does not Comparable . But the sales of last 3years all Comparable & also graduate Students vs under graduate Students are Comparable.

* Methods for class Comparison:

It contains the following steps .

Step1:- collect the data for target class & Contrasting class through DB queries .

Step2: Apply the attribute relevance analysis .

Step3: Apply the Generalization . i.e..., we Select the generalization threshold range for Specific attribute . Then this is applied on target class & Contrasting class to get the Prime relations.

Step 4:- The derived class Comparisons are Presented to end user by using several visualization techniques.

Step1: for eg..., consider the big-university DB & the attributes are name, gender.

major, birth-place , birth-date, residence & phone # to get the

use Big-university-DB

Mine Comparison as "grade-vs-undergrad students"

in relevance to name, gender, major, birth-place, birthdate, residence, phone #
for "graduate-students"
where status in "graduate"
versus "under graduate-students"
where status in "undergraduate"
analyse count %
from student.

Once the query is executed, we get the data for target class. i.e..., graduate students.

| Name | gender | major | birth-place | birth-date | residence | phone |
|------|--------|-------|-------------|------------|-----------|-------|
| Mahi | M | Science | Hyd, AP INDIA | 12-07-76 | st-street Hyd | 234 5 |
| Lee | F | Engineering | Vancouver ABC, CANADA | 12-07-70 | AB-St, Vancouver | 599-8 |

fig(a) Data for Target class i.e..., graduate Studer

| Name | gender | major | birth-place | birth-date | residence | Phone# |
|------|--------|-------|-------------|------------|-----------|--------|
| Ram | M | Science | N.Y.ABC, USA | 12-07-80 | Clain St, N.Y | 123456 |
| Anu | F | busine-ss | Vancouver, XYE, CANADA | 10-09-77 | SR-St, Vancouver | 234-576 |

fig(b) Data for Contrasting class i.e..., undergraduate students

**Step2:** Apply the attribute relevance analysis after applying the attribute relevance analysis the attribute name, gender, birth-place, residence & Phone # are removed from the relational table.

**Step3:** Apply the generalization. i.e., the attribute birth-place is transformed into age-range. The Prime relations are shown in below table.

| major | age-range | count % |
|-------|-----------|---------|
| Science | 20.......25 | 5.32% |
| Engineering | 25........30 | 4.12% |

fig(a) Prime Relation for graduate Students.

| major | age-range | count % |
|-------|-----------|---------|
| Science | 20........25 | 3.32% |
| Business | 30........35 | 2.14% |

fig(b) Prime Relation for undergraduate Students.

**Step4:**

These Prime relations are presented to enduser by using visualization techniques. In visualization technique we analyzed the

measure is count %. for eg..., 5.32% of students with the age 20....25 & major in Science.

## Presentation of class Comparisons:-

~~To Present the class Comparisons we~~ use the method. This method is called as Statistics discrimination (or) d. weight value.

for eg..., Let's $s_A$ be the any tuple & the target class $c_j$ then d-weight value of $s_A$ in $c_j$ is defined as the tuples containing target class by using total no of tuples.

$$d\text{-weight} = Count(s_A \in c_j) \Big/ \sum_{i=1}^{m} count(s_A \in c_i)$$

Here, m contains the total no of tuples i.e..., target class & contrasing class. for eg..., Consider the data for graduate students and under graduate students.

| Status | major | age.range | count |
|---|---|---|---|
| graduate | Science | 20......25 | 210 |
| undergraduate | Engineering | 25......30 | 90 |

fig(a): Data for graduate & undergraduate students

The d. weight value for target class i.e..., graduate students.

$$d\text{-weight} = \frac{\overset{7}{210}}{\underset{10}{300}} = 0.7 = 70\%$$

d-weight value for contrasting class.
i.e., undergraduate Students

$$d\text{-weight} = 9\emptyset/30\emptyset = 0.3 = 30\%$$

finally, if the student in major in science & with age 20....25 then, there is a Probability is that he is the 70% of graduate students.

* Class Descriptions:-

The class description means it contains the 2. i.e., class characterization & class Comparison. for eg.., Consider the data Shown below·

| Location | TV | Computer |
|----------|-----|----------|
| Europe | 90 | 150 |
| North America | 210 | 520 |

fig(a) : class table data for 2 locations

Then d.weight value i.e., TV

Europe (d.weight) $= 90/300 = 0.3 = 30\%$

North America (d.weight) $= 210/300 = 0.7 = 70\%$

d-weight value i.e., Computer

Europe (d.weight) $= 150/67\emptyset = 0.22 = 22\%$

North America (d.weight) $= 520/67\emptyset = 0.78 = 78\%$

These d-weight values are shown in below.

| Location | Item | TV | | Computer | |
|---|---|---|---|---|---|
| | | count | d.weight | count | d.weight |
| Europe | | 90 | 30% | 150 | 22% |
| North America | | 210 | 70% | 520 | 78% |

fig(b) d.weight value for fig(a)

from the fig(b) we can conclude that the north America Sales are very good ✓

* Mining Statistical class Description from the large Database :-

Upto now we analyze the characterstics of a relational DB by using the Predefined aggregate functions like sum(), count() etc But many of the applications end user want to analyze the characteristics of a large DB by using the measures mean, median, mode. These 3 measures are called as "Central tendency".

* Measuring the Central tendency :-

The most popular method to find the centre value of a given set is mean. for eg...
Data items $x_1, \ldots, x_n$ then mean $= \frac{1}{n} \sum_{i=1}^{n} x_i$

This mean is equal to avg in SQL.

$$\text{average} = \frac{\text{sum}}{\text{count}}$$

Sometimes the 'a' contains the bytes - i.e... $w_i$ for $i=1$ to $n$. Then, weighted arithematic mean.

$$= \frac{\sum\limits_{i=1}^{n} x_i w_i}{\sum\limits_{i=1}^{m} w_i}$$

we can find the median value by using a formula:

$$\boxed{\text{median} = L_1 + \left(\frac{N/2 - (\Sigma f)_1}{f_{median}}\right) c}$$

where '$L_1$' is the lowest class boundary
    '$c$' is the class interval.
    $(\Sigma f)_1$ is the cummulative frequency $< f_{median}$
    $f_{median}$ is the median frequency
    '$N$' is the no of samples.

Eg:- for eg.., Consider the below table.

| class Interval '$c$' | frequency $f_{median}$ | Cummulative frequency $(\Sigma f)_1$ |
|---|---|---|
| 1 - 3 | 2 | 2 |
| 4 - 6 | 3 | 5 |
| 7 - 9 | 5 | 10 |

## Q. Mining Association Rules in large Data Bases

**Association Rule Mining:-**

The association rule mining searches the interesting relationships between items in the given data set.

An typical Example of association rule of Market Basket Analysis.

**Market Basket Analysis:-**

This Market basket analysis analyse buying behaviour of customer by using association rul between the items.

For Eg:- Consider the three customers, then Mao Basket Analysis is shown below.

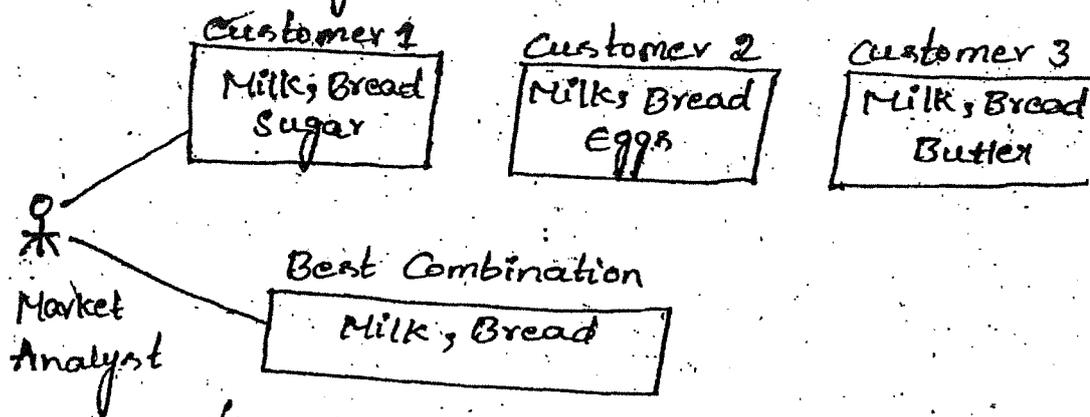| Customer 1 | Customer 2 | Customer 3 |
|---|---|---|
| Milk, Bread Sugar | Milk, Bread Eggs | Milk, Bread Butler |

Market Analyst

Best Combination

| Milk, Bread |
|---|

fig : Market Basket Analysis.

Here Market Analyst identifies the custome: Purchases milk & also Bread. Therefore this mi and Bread placed together further more increase sales.

Therefore, as well as using the results Market Basket Analysis the ll...

This market Basket Analysis uses the association rule.

For Eg: Customer purchases a Computer and also Purchases financial Mgmt Software. Then the association rule,

Computer $\Rightarrow$ financial - Management - Software

[Support 2% ; Confidence = 60%]

Here Confidence 60% means 60% of customers Purchase Computers & also purchases financial Mgmt software & support is 2%.

## Basic Concepts :-

Consider Set of items, $I = \{i_1, i_2 - - - i_m\}$ & task relevant data Set 'D'. we get this task relevant data Set by using Set of transactions. Each transaction is represented by $T_i$  $T \subseteq I$.

As well as Each transaction is uniquely identified by using TID [Transaction ID].

For Eg, item ACI, then $A \Rightarrow B$, $A \subset I$, $B \subset I$ & $A \cap B = \phi$

Then association rule $A \Rightarrow B$ with Support 'S', where 'S' is the percentage of transactions in 'D' Contains A & B i.e.,

$$D (A \cup B)$$

The association rule $A \Rightarrow B$ has Confidence 'C', where 'C' is the Percentage of transactions in 'D' Contains A & also Contains B i.e.,

$$P(B/A)$$

$\therefore$ Support $(A \Rightarrow B) = P(A \cup B)$ i.e, Transactio Contains 'A' & 'B'.

Confidence $(A \Rightarrow B) = P(B/A)$ i.e, Transaction contains 'A' & also contains 'B'. Here, Support & confidence is represented in the form of percentage & ranges in betwe 0% to 100%.

## Association Rule Mining : A Road Map :-

The market-Basket analysis i.e, one o the association rule. The association rule classified based on the following criteries.

i, **Based on the type of value used in rule:-**

The association rule does not contain any item or attribute or dimension, then that association rule is called as boolean associ rule.

Ex:- Computer $\Rightarrow$ financial-management-software

If the association rule is described b/w quantitative values, then that association rul is called as quantitative association rule.

Ex:- $age(x, "20---29") \wedge income(x, "30k---39k"$ $\Rightarrow buys(x, "VCR")$

ii, **Based on the type of dimensions used in the**

For Ex:- The association is defined for one

dimension then that association rule is called as single dimension association rule.

Ex: buys ($X$ ; "Computer") $\Rightarrow$ buys ($X$ , "financial - Mgnt - S/w)

iii) Based on the different types of Abstractions in the Rule:-

Here the association rule contains the different types of data abstractions. i.es using this association rule we abstract data at different levels.

Here the association rule contains the different types of data levels.

Ex:- age ($x$ ; "30 --- 39") $\Rightarrow$ buys ("Computer")

age ($x$ ; "30 ----39") $\Rightarrow$ buys ("laptops")

## **Mining Single dimensional Boolean Association Rule From Transactional DbS:-

Here we use the 'Apriori' algorithm.

Apriori Algorithm:- Finding frequency Itemsets by using candidate Sets.

1) The Apriori Algorithm is used to find frequent Item Sets.

2) These are used in Boolean association rule.

3) Apriori means recursive approach, i.e., level-wise-search.

4) First of all we find frequent Itemset. This is denoted by 'L'. Then using 'L₁' we find $L_2$ i.e., frequent 2 - itemset.

5) This is used to find , $L_3$ and so on, until we can not get the no frequent set. i.e., the Apriori algorithm uses prior knowledge of frequent item set and for each $L_k$ we have to scan entire databas i.e., Đ

To simplify the Apriori algorithm, it uses apriori property. Using this Apriori Property, we reduce search length.

For Ex, item set is represented by 'I'. & it does not satisfying the minimum support threshold i.e., Sup-min, then I is not a frequent itemset, then

$$P(I) < Sup\text{-}min$$

If we add another item 'A' then it results

$I \cup A$,

It is also not a frequent itemset

$$P(P(I)) < Sup\text{-}min$$

This apriori property contains the term steps.

1) The Join Step:-

To calculate $L_k$, a candidate k itemset are joined, $L_{k-1}$ to itself. The candiate set is represented by $C_k$. Consider $L_1$ & $L_2$ item

in $\ell_{k-1}$. In general $x_i[j]$ means $j$ value of $\ell$, then we join $\ell_{k-1} \bowtie \ell_{k-1}$, it requires the first $k-2$ items are equal. Therefore, to join $\ell_1$ and $\ell_2$ of $\ell_{k-1}$ and if

$$\left[\ell_1[1] = \ell_2[1]\right) \wedge \left(\ell_1[2] = \ell_2[2]\right) \wedge \cdots \cdots \wedge \left(\ell_1[k-2] = \ell_2[k-2]\right]$$

$$\wedge \left(\ell_1[k-1] < \ell_2 \cdot [k-1]\right]$$

The last conditions

$\ell_1[k-1] < \ell_2[k-1]$ is used to avoid the duplicate values.

2) The **Prune Step**:-

The prune step, all non-empty subset of frequent items is also frequent.

Ex:- Consider transactional database of 'all Electronics'. This is represented by '$D$' & it contains of transactions.

$\therefore$ we use the Aprori algorithm to find frequent item & have minimum count is '2'

| TID | List of item-IDs |
|---|---|
| $T_{100}$ | $I_1, I_2, I_5$ |
| $T_{200}$ | $I_2, I_4$ |
| $T_{300}$ | $I_2, I_3$ |
| $T_{400}$ | $I_1, I_2, I_4$ |
| $T_{500}$ | $I_1, I_3$ |
| $T_{600}$ | $I_2, I_3$ |
| $T_{700}$ | $I_1, I_3$ |
| $T_{800}$ | $I_1, I_2, I_3, I_5$ |

Scan $D$ for Count of candidates $\rightarrow$

$C_1$

| item set | Support Count |
|---|---|
| $\{I_1\}$ | 6 |
| $\{I_2\}$ | 7 |
| $\{I_3\}$ | 6 |
| $\{I_4\}$ | 2 |
| $\{I_5\}$ | 2 |

fig : Transaction ... ... ... ... ...

Compare Candidate
Sup-count with
min sup-count
→

### $d_1$

| Item set | Sup-Count |
|---|---|
| $\{I_1\}$ | 6 |
| $\{I_2\}$ | 7 |
| $\{I_3\}$ | 6 |
| $\{I_4\}$ | 2 |
| $\{I_5\}$ | 2 |

### $C_2 = d_1 \bowtie d_1$

Generate
$C_2$
Candidate
Using $d_1$
→

| Itemset | Sup-count |
|---|---|
| $\{I_1, I_2\}$ | 4 |
| $\{I_1, I_3\}$ | 4 |
| $\{I_1, I_4\}$ | 1 |
| $\{I_1, I_5\}$ | 2 |
| $\{I_2, I_3\}$ | 4 |
| $\{I_2, I_4\}$ | 2 |
| $\{I_2, I_5\}$ | 2 |
| $\{I_3, I_4\}$ | 0 |
| $\{I_3, I_5\}$ | 1 |
| $\{I_4, I_5\}$ | 0 |

Compare
Candidate
Sup-count
with min
Sup-count
→

### $d_2$

| Itemset | Sup-Co |
|---|---|
| $\{I_1, I_2\}$ | 4 |
| $\{I_1, I_3\}$ | 4 |
| $\{I_1, I_5\}$ | 2 |
| $\{I_2, I_3\}$ | 4 |
| $\{I_2, I_4\}$ | 2 |
| $\{I_2, I_5\}$ | 2 |

### $C_3 = d_2 \bowtie d_2$

Generate
$C_3$ Candidate
Using $d_2$

| Itemset | Sup-count |
|---|---|
| $\{I_1, I_2, I_3\}$ | 2 |
| $\{I_1, I_2, I_5\}$ | 2 |

Compare
Candidate
C·C sup-count
with min
Sup-count
→

### $d_3$

| itemset | S c |
|---|---|
| $\{I_1, I_2, I_3\}$ | |
| $\{I_1, I_2, I_5\}$ | |

Using Apriori algorithm, we find Candidate d
& most frequently item set with min Support Count

Steps for finding candidate item set of frequent items.

1) This algorithm scans all the transcations of database to find total count of each itemset. This itemset is called as Candidate 1 - itemset, $c_1$.

2) Suppose, Consider minimum support threshold is $2 \, (2/9 = 0.22 = 22\%)$

3) Candidate Support Count is Compared with minimum Support Count. Then we get frequent itemset $\alpha_i$.

4) To find $c_2$, the $\alpha_1$ joined itself i.e.,

Join : $c_2 : \alpha_1 \bowtie \alpha_1$

Then we find $\alpha_2$ by Comparing Candidate Support Count with minimum Support-count.

5) Join $c_3 : \alpha_2 \bowtie \alpha_2$

$\{ \{I_1, I_2\}, \{I_1, I_3\}, \{I_1, I_5\}, \{I_2, I_3\}, \{I_2, I_4\}, \{I_2, I_5\} \} \bowtie \{ \{I_1, I_2\}, \{I_1, I_3\}, \{I_1, I_5\}, \{I_2, I_3\}, \{I_2, I_4\}, \{I_2, I_5\} \}$

$= \{ \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\}, \{I_1, I_3, I_5\}, \{I_2, I_3, I_4\}, \{I_2, I_3, I_5\}, \{I_2, I_4, I_5\} \}$

# Prune Step:-

All non-Empty frequent subsets are also frequent.

1. 2-Item sets for $\{I_1, I_2, I_3\}$ are $\{I_1, I_2\}, \{I_1, I_3\}$ and $\{I_2, I_3\}$. Then $\{I_1, I_2, I_3\}$ Add to $c_3$ because all 2-item sets present in $d_2$.

2. 2-Item sets for $\{I_1, I_2, I_5\}$ are $\{I_1, I_2\}, \{I_2, I_5\}$ and $\{I_1, I_5\}$ there all 2-item set present in $d_2$. Then $\{I_1, I_3, I_5\}$ add to $c_3$.

3. 2-Item Sets for $\{I_1, I_3, I_5\}$ are $\{I_1, I_3\}, \{I_3, I_5\}$ and $\{I_1, I_5\}$ Then $\{I_3, I_5\}$ is n in $d_2$ then $\{I_1, I_3, I_5\}$ is Removed from $c$.

4. 2-Item sets for $\{I_2, I_3, I_4\}$ are $\{I_2, I_3\}, \{I_3,$ and $\{I_2, I_4\}$ then $\{I_3, I_4\}$ is not in $d_2$. Therefore, it is not a frequent set. So $\{I_2, I_3\}$ is removed from $c_3$.

5. 2-Item sets for $\{I_2, I_3, I_5\}$ are $\{I_2, I_3\}$ $\{I_3, I_5\}$ and $\{I_2, I_5\}$ Then $\{I_2, I_3, I_5$ is removed from $c_3$.

6. 2-Item sets for $\{I_2, I_4, I_5\}$ are $\{I_2, I_4\}$, $\{I_4, I_5\}$ and $\{I_2, I_5\}$ then $\{I_4, I_5\}$ is in $d_2$. Therefore it is not frequent set. Then $\{I_2, I_4, I_5\}$ is removed from $c_3$.

7. Then $C_3 = \{\{I_1, I_2, I_3\}, \{I_1, I_2, I_5\}\}$

8. Again we scan all the transcations to find the count value for 3-itemsets. Then we get $L_3$.

| Item set | Sup-Count |
|---|---|
| $\{I_1, I_2, I_3\}$ | 2 |
| $\{I_1, I_2, I_5\}$ | 2 |

## Apriori Algorithm :-

Find frequent set using recursive level-wise approach based on Candidate item-set.

Input : Database $D$, OP-transcation; minimum Support threshold; min-sup

Output : $L$, frequent itemset.

Method : $L_1$ = first - frequent - 1 - itemset;

$\quad$ for $(k=2 ; L_{k-1} \neq \phi ; k++)$

$\quad \{$

$\quad\quad C_k$ = apriori-gen $(L_{k-1}, min-sup)$;

$\quad\quad$ for Each transaction $t \in D$

$\quad\quad \{$

$\quad\quad\quad C_k$ = Subset (count

$\quad\quad\quad$ for Each $c \in C_k$

$\quad\quad\quad\quad$ c. Count ++;

$\quad\quad \}$

$\quad\quad L_k = \{ c \in C_k / c. count \geq min-sup\}$

```
        }
        return L;
    }

Procedure apriori-gen ($L_{K-1}$: frequent $k-1$ itemset
                min-sup: mim· support threshold)
    {
    Consider $l_1$, $l_2$, $l_1 \in L_{K-1}$, $l_2 \in L_{K-1}$
    if $l_1[i] = l_2[i] \wedge ---- \wedge (l_1[k-2] = l_2[k-2])$
        $(l_1[k-1] < l_2[k-1])$ then $C = l_1 \bowtie l_2$;
    if has-infrequent-subset $(C_K, L_{K-1})$ the
    delete $c$;
    else
        add $c$ to $C_K$;
    return $C_K$;
    }

Procedure has-infrequent-subset $(c$ : candidate
            item set, $L_{K-1}$ : frequent $k$-itemse
    {
    if $C \not\subseteq L_{K-1}$ then
        return TRUE
    else
        return FALSE
    }
```

Alg:-

Apriori Algorithm for finding frequent iter
i·e, used in boolean association rule.

Generating Association rules from frequent itemset

Using Apriori alg, we find frequent item sets then we derive Strong association rules

The Strong association rule means it must include Support Count & Confidence Count. Then we use the following Equation for the Confidence to derive Strong association rules where Conditional Probability is represented in the form of the Support Count.

$$Confidence \ (A \Rightarrow B) = P(B/A) = \frac{Support-Count \ (A \cup B)}{Support-Count \ (A)}$$

where Support-Count $(A \cup B)$ represents transactions Containing A & B.

Support-Count(A) represents transaction Containing 'A'.

Using this we derive 2 associations rules for Each frequent item sets $l$, must Contain all non-Empty Subsets of '$l$'.

For Each Subset 's' of '$l$', the o/p rule "$s \Rightarrow (l-s)$" if $\frac{Support-Count(l)}{Support-Count(s)} \geq min-Confidence$

For ex, frequent itemset,

$l = \{ I_1, I_2, I_5 \}$ then Subsets are

$\{ I_1, I_2 \} \{ I_1, I_5 \} \{ I_2, I_5 \} \{ I_1 \} \{ I_2 \} \ & \ \{ I_5 \}$

Then the association rules are

$I_1 \cap I_2 \Rightarrow I_5$ , Confidence $= 2/4 = 50 \%$. —— ①

$I_1 \cap I_5 \Rightarrow I_2$ , Confidence $= 2/2 = 100 \%$. —— ②

$I_2 \cap I_5 \Rightarrow I_1$ , Confidence $= 2/2 = 100 \%$. —— ③

$I_1 \Rightarrow I_2 \cap I_5$ , Confidence $= 2/6 = 33 \%$. —— ④

$I_2 \Rightarrow I_1 \cap I_5$ , Confidence $= 2/7 = 28 \%$. —— ⑤

$I_3 \Rightarrow I_1 \cap I_2$ , Confidence $= 2/2 = 100 \%$. —— ⑥

Consider min confidence is $70 \%$ , then the o/p is ② , ③ , ⑥ only [we get] strong associati

## Improve the Efficiency of Apriori :-

we use following techniques to improve the Efficiency of Apriori alg.

### 1. Hash based Technique :-

using this , we reduce the size of candidate -k itemset. we create both hash table for 'H by using the hash function.

$$H(x,y) = ((\text{order of } x) \times 10 + (\text{order of } y)) \bmod 7$$

we consider transaction db 6-1 then

Hash table $H_2$ for Candidate-2 itemset

| Hash address | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| Hash Count | 2 | 2 | 4 | 2 | 2 | 4 | 4 |
| Hash Contents | {I₄,I₄} {I₂,I₅} | {I₁,I₅} {I₁,I₅} | {I₂,I₃} {I₂,I₃} {I₂,I₃} {I₂,I₃} | {I₂,I₄} {I₂,I₄} | {I₁,I₅} {I₂,I₅} | {I₁,I₂} {I₁,I₂} {I₁,I₂} {I₁,I₂} | {I₁ {I₁ {I₁ {I₁ |

$\{ I_1, I_2, I_5 \}$

$h(I_1, I_2) = 10 + 2 \% 7 = 5$

$h(I_1, I_5) = 10 + 5 \% 7 = 1$

$h(I_2, I_5) = 20 + 5 \% 7 = 4$

$h(I_2, I_4) = 20 + 4 \% 7 = 3$

$h(I_2, I_3) = 20 + 3 \% 7 = 2$

$h(I_1, I_4) = 10 + 4 \% 7 = 0$

$\{ I_1, I_2, I_4 \} \Rightarrow h(I_1, I_4) = 10 + 4 \% 7 = 0$

$h(I_1, I_3) = 10 + 3 \% 7 = 6$

$\{ I_1, I_2, I_3 I_5 \}$

$h(I_3, I_5) = 30 + 5 \% 7 = 8$

Consider min-Support Count is '3'. Then the addresses 0, 1, 3, 4 does not be considered for $c_2$.

## 2. Transaction Reduction:-

- If the transaction does not Contain frequent k-itemset, then it also does not Contain the frequent k+1 itemset then those transactions must be removed.

## 3. Partitioning:-

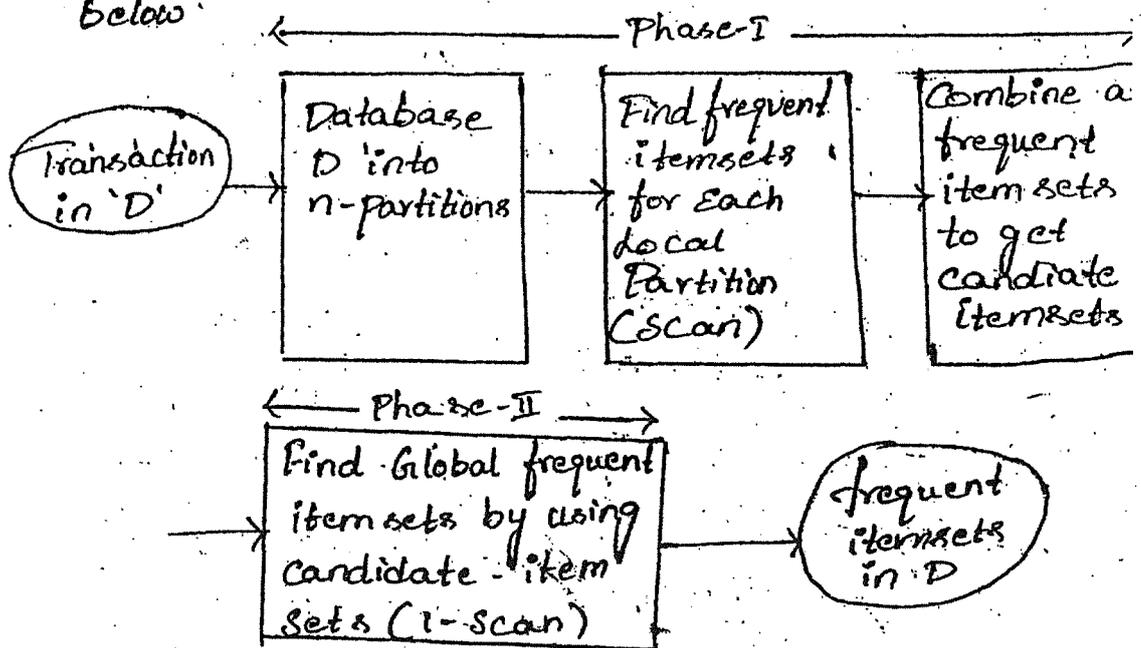In this, it requires only 2 database Scans & also it Consists 2 phases.

Phase-I: It Contains following steps.

1. Divide db 'D' into 'n' Partitions.

2. Find frequent itemset for Each local partition

3. Combine these  ~~~~
   Candidate item sets.                                        V

Phase-II
1. Find global frequent item set - This is Shown
below·



← ———————————— Phase-I ————————————→

Transaction in 'D' → Database D into n-partitions → Find frequent itemsets for Each local Partition (Scan) → Combine a frequent item sets to get candidate Itemsets

← —— Phase-II ——→

Find Global frequent item sets by using candidate - item Sets (1-Scan) → frequent itemsets in D

Partition data Method·

4. Sampling :-
     Consider the random Sample 's' in database
Then we find frequent itemsets in 's' rather
than 'D'.
     The Size of 's' in Such a way that it mu
be fit in primary memory i·e·, in single Sca
all the transactions of 's' avoided·
     Here we are finding frequent itemset in '
Therefore, we may miss Some frequent itemsets
Therefore we Consider lower Support threshold
rather than min·Support Threshold·

# Generate Frequent itemsets without using Candidate itemsets (or) FP growth Method

In Apriori algorithm, we find Candidate itemsets. Using this Candidate itemsets, we find frequent itemsets. But the Apriori Algorithm Contains the following drawbacks.

i) It generates large No. of Candidate itemsets For Ex, $10^4$ frequent - itemset. Then apriori alg generates approximately $10^7$ Candidate - 2 itemsets.

ii) It requires the large No. of database Scans to identify the value of candidate itemsets.

To avoid these drawbacks, we go for the frequent Pattern - growth Method.

This is simply Called as FP-growth method. It follows the Divide - and - Conquer technique. It contains three steps.

1. Compress the db representing itemsets into frequent pattern - tree or FP-tree. following the itemset information.

2. Transform Compressed db into conditional DB.

3. Generate frequent Patterns.

Consider the transaction db fig:

| TID | list of item - IDs |
|-----|--------------------|
| $T_{100}$ | $I_1, I_2, I_5$ |
| $T_{200}$ | $I_2, I_4$ |
| $T_{300}$ | $I_2, I_3$ |
| $T_{400}$ | $I_1, I_2, I_4$ |
| $T_{500}$ | $I_1, I_3$ |
| $T_{600}$ | $I_2, I_3$ |
| $T_{700}$ | $I_1, I_3$ |
| $T_{800}$ | $I_1, I_2, I_3, I_5$ |
| $T_{900}$ | $I_1, I_2, I_3$ |

After the first scan of db, we get L, it contains the item set & support-count & we represent this in descending order of support count & it is represented by 'L'

$$L = [I_2 : 7, I_1 : 6, I_3 : 6, I_4 : 2, I_5 : 2]$$

Consider support-count is '2', Then <u>FP-tree</u>

<u>Construction is</u>

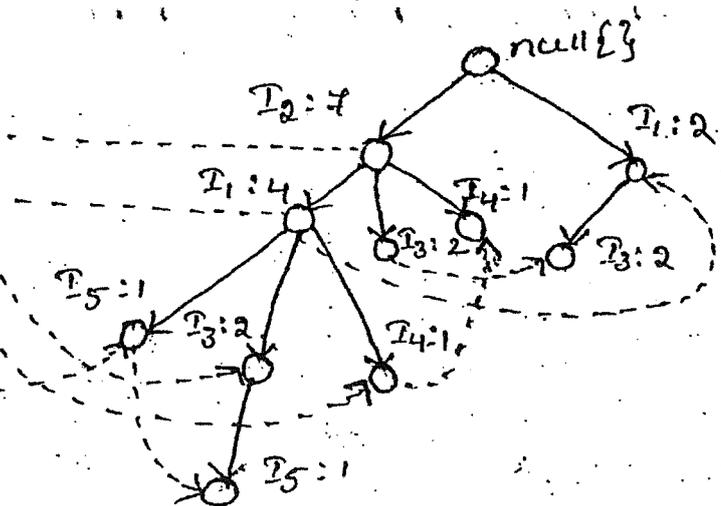after the second scan of db, the first transaction is

"$T_{100} : I_1, I_2, I_5$".

& these items are arranged based on the order of 'L'. Then we get $(I_2, I_1, I_5)$ & o these items are occured in first time, the $(I_2 : 1)$ is treated as the root & $(I_1 : 1)$ is attached to $I_2$ & $(I_5 : 1)$ is attached to $I_1$ & then second transaction.

$T_{200} : I_2, I_4$

then, $I_2$ is already root, so increment the count Value of '$I_2$' & ($I_4 : 1$) is attached to $I_2$. This is shown below.

✓ FP-Tree Construction Start with the root, this root is labelled as null {}

| Item Id | Support Count | Node LINK |
|---------|---------------|-----------|
| $I_2$   | 7             |           |
| $I_1$   | 6             |           |
| $I_3$   | 6             |           |
| $I_4$   | 2             |           |
| $I_5$   | 2             |           |



✓ FP-tree based on Transaction Database.

Construction of Conditional db by Mining FP-Tree :-

Start with the item '$I_5$' it contains the 2 branches. These branches are identified by verifying Node link. Then the paths are

$$I_2 \ I_1 \ I_5 : 1 \qquad I_2 \ I_1 \ I_3 \ I_5 : 1$$

Then $I_5$ is considered as suffix. Then we

1. Conditional patterns how i.e, as link paths of $I_5$ are ($I_2 \ I_1 : 1$) and ($I_2 I_1 I_3 : 1$)

Then Conditional FP-tree Contains only one path, then we merge above 2 paths & we get

( I2, x, -1 ) ...

We cannot include $I_3$ because its support count is '1' & it is less than the support count. Then this single path generates the frequent itemsets by concatinating suffix.

∴ Frequent patterns are

$I_2 I_5 : 2$    $I_1 I_5 : 2$

This is shown below.

| Item | Conditional pattern Base | Conditional FP-tree | Frequent Pattern Generation |
|------|--------------------------|---------------------|------------------------------|
| $I_5$ | $\{\langle I_2 I_1 : 1 \rangle \langle I_2 I_1 I_3 : 1 \rangle\}$ | $\langle I_2 : 2, I_1 : 2 \rangle$ | $I_2 I_5 : 2$  $I_1 I_5 : 2$  $I_2 I_1 I_5 : 2$ |
| $I_4$ | $\{\langle I_2 I_1 : 1 \rangle, \langle I_2 : 1 \rangle\}$ | $\langle I_2 : 2 \rangle$ | $I_2 I_4 : 2$ |
| $I_3$ | $\{\langle I_2, I_1 : 2 \rangle, \langle I_2 : 2 \rangle \langle I_1 : 2 \rangle\}$ | $\langle I_2 : 4, I_1 : 2 \rangle, \langle I_1 : 2 \rangle$ | $I_2 I_3 : 4, I_1 I_3 : 2, I_2 I_1 I_3 : 2$ |
| $I_1$ | $\{\langle I_2 : 4 \rangle\}$ | $\langle I_2 : 4 \rangle$ | $I_2 I_1 : 4$ |

Conditional db by Mining FP-tree

Finally, using this FP-growth method large Pattern is divided into smaller pattern by using suffix & then suffix is concatenated to get the frequency patterns.

This alg. reduces the search cost & also this is more faster than the apriori algorithm.

## ICEBERG QUERIES:-

The Apriori alg is used to improve the efficiency of ICE-BERG Query. This ice-berg query is mainly used in DM preferably Market

The ICE-BERG Query computes aggregation function over an attribute, to find out which aggregate values over the specified threshold.

Consider the relation 'R' with attributes $a-1, a-2, ---- a-n$, & $b$ and aggregate function $agg-f$. Then the ice-berg query looks like the following.

Select $R \cdot a-1, R \cdot a-2, ---- R \cdot a-n, agg-f(R \cdot b)$
from relation R.

group by $R \cdot a-1, R \cdot a-2 --- R \cdot a-n$
having $agg-f(b) >= threshold$.

The above ICE-BERG query accepts the large no of i/p tuples, but it gives very small No. of o/p tuples, because the o/p tuples only displays, that satisfy the threshold values.

# MINING MULTILEVEL ASSOCIATION RULES FROM TRANSACTION DATABASES:-

Mining multilevel association rules means rules are included at different levels of abstraction.
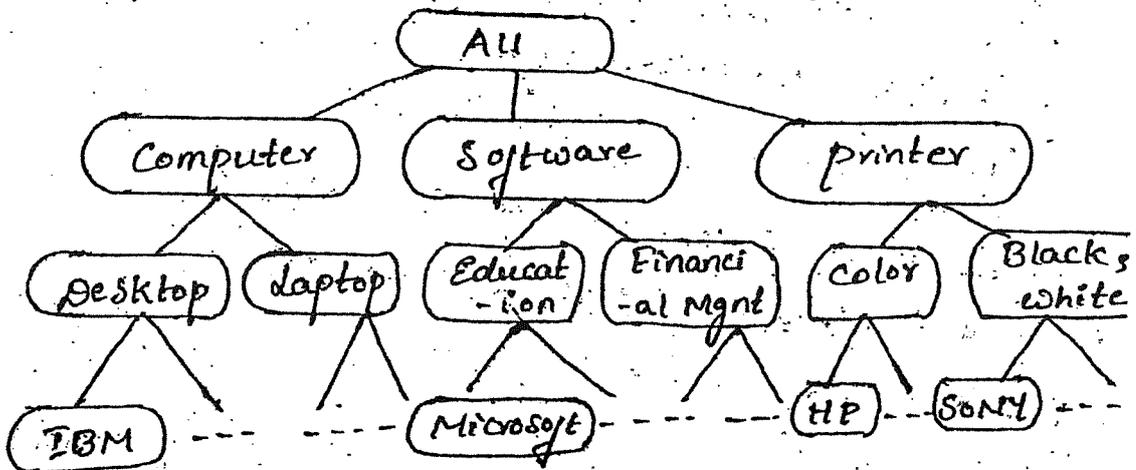
## Multilevel Association Rules:-

It is difficult to define strong association rules at low level or primitive level.

Consider the task-relevant data set as specified below.

| TID | Items Purchased |
|-----|-----------------|
| T1 | IBM, Desktop Computer, sony B/w printer, |
| T2 | Microsoft Education S/w, HP Color printe |
| T3 | Microsoft financial Mgmt S/w, |
| T4 | IBM laptop Computer |
| ⋮ | ⋮ |

Task Relevant Data set

The above Task relevant data set is represented in Concept hierarchy as.



The Concept hierarchy defines set of mapping from low level to high level.

Here it Contains 4 levels. i.e, levels 0, 1, 2 &
The top level is level '0' i.e., represented by using keyword 'all'.

The level '1' Contains Computer, S/w & printer.

The level 2 Contains Desktop Computer, laptop Computer & so on.

The level 3 Contains IBM Desktop Computer & so on.

the lowest level is level 0 & in this level ....
difficult to identify interesting patterns i.e.,
Consider "IBM Desktop Computer " & "Sony B&W
Printer " occur in very few.

∴ Define Generalization like "IBM Desktop Computers"
to " Computer " & "Sony B&w printer " to "Printer".

∴ The Combination { Computer , Printer } as many
People frequently purchased.

∴ using this multilevel Concept hierarchies , we
Easily define the interesting patterns.

## Approaches to Mining Association Rules :-

Generally we follow the top down approach ,
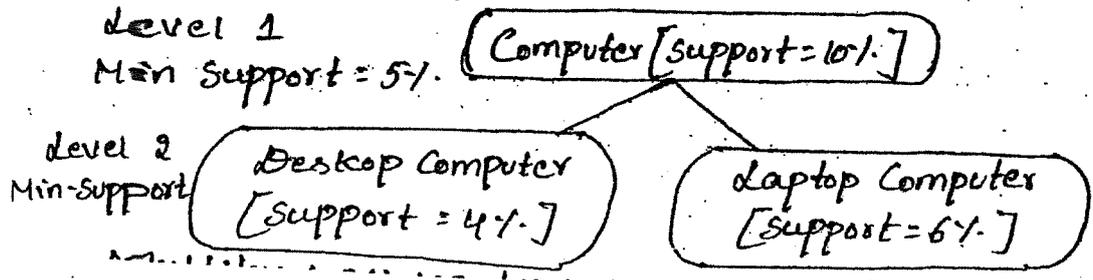-i.e, with level 1 & find frequent item
There are no frequent item Sets. Then we go for level
2 & so on.

we use the following approaches for mining
Association rules.

## Uniform Support :-

Here we define Same minimum support .
Threshold for all the levels.

For Ex, Min. Support Threshold is 5%. , then

Level 1
Min Support = 5%.  ( Computer [Support = 10%.] )

Level 2
Min-Support  ( Desktop Computer
[Support = 4%.] )    ( Laptop Computer
[Support = 6%.] )

Therefore, comp—

frequent patterns of desktop computer or not.

This approach Contains the advantages of I. Searching process is Simplified i.e, here we verify only those items satisfy Min-Support II. Here Enduser has to specify only one Min Support Threshold value.
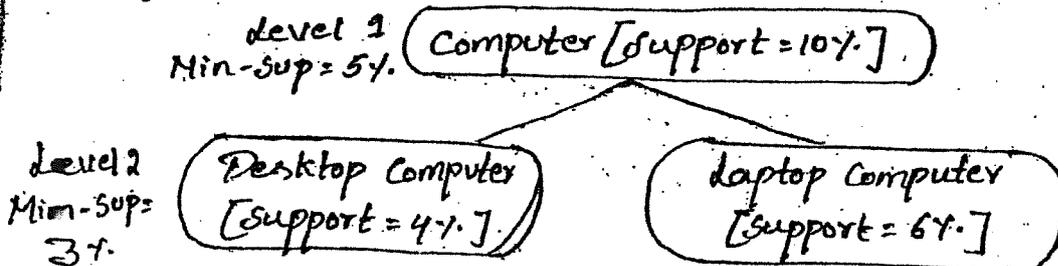
This technique Contains the drawbacks of

i) we already see the low level item is not frequer. as of upper level item. Therefore, we Cannot define Min-Support Threshold--

ii) if we specify very large support threshold tha the low level important patterns may be missed

iii) If we define very small support threshold then it Contains the large No-of unimportant patterr

ii) **Reduced Support:-**

Here Each level has its own minimum support threshold. Generally higher levels Contains the big: values than the low levels. Consider Min-Support threshold for level 1 & 2 are 5% & 3%. This is shown in below.

level 1
Min-Sup=5%.   ( Computer [Support=10%.] )

Level 2
Min-Sup=
3%.   ( Desktop Computer [Support=4%.] )     ( Laptop Computer [Support=6%.] )

Multilevel Mining with reduced support.

These Computer, Laptop Computer, ... 
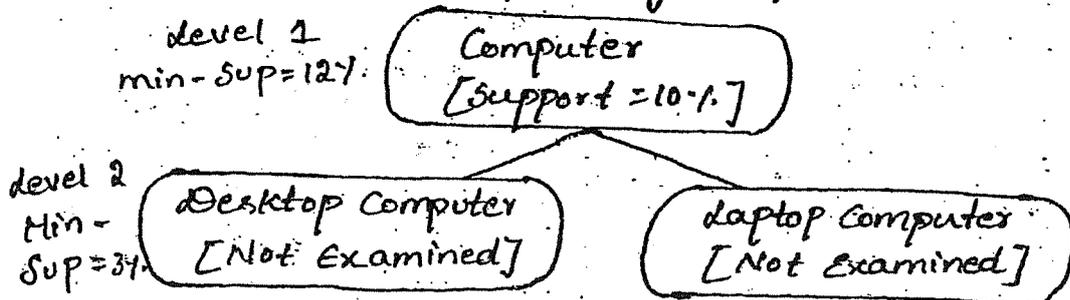are frequent patterns.

iii) Level-by-Level Independent :-

   Here each node is examined independently i.e., regardless of its parent is frequent or not.

iv) Level Cross Filtering by Single item :-

   Here 'item with' $i^{th}$ level is examined if & only iff its parent $(i-1)^{th}$ level is frequent.

   Consider the following Diagram.

Level 1
min-Sup=12%

Computer
[Support = 10%]

Level 2
Min-Sup=3%

Desktop Computer
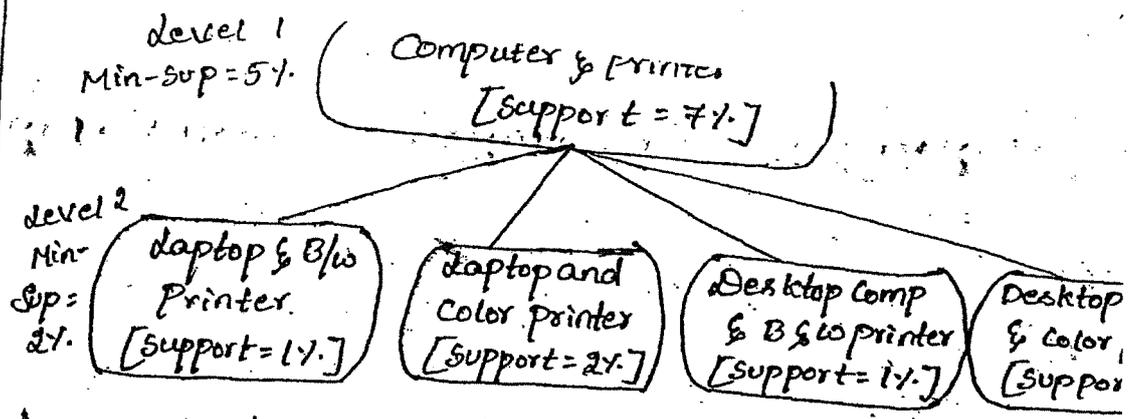[Not Examined]

Laptop Computer
[Not Examined]

Multilevel Mining with reduced support using level cross filtering by single item.

   Here Desktop Computer & Laptop Computer is not examined because its parent 'computer' is not a frequent.

v) Level Cross Filtering by k-item Set :-

   Here k-item set is examined [k-itemset with $i^{th}$ level] if & only iff parent k-itemset of $(i-1)^{th}$ level is frequent.

   Consider the below diagram.

Level 1
Min-Sup = 5%.

Computer & Printer
[Support = 7%.]

Level 2
Min-Sup = 2%.

Laptop & B/w Printer
[Support = 1%.]

Laptop and Color Printer
[Support = 2%.]

Desktop Comp & B & w printer
[Support = 1%.]

Desktop & Color
[Suppor...

Multilevel 'Mining' with Reduced Support using level-cross filtering by k-item set, here k=2

Here { Laptop Computer , B/w printer }

{ Laptop Computer ; Color   "   }

{    "         "    ; B & w   "   }

{    "         "    ; Color   "   } are Exami

because 'Computer' & 'Printer' are frequent.

## Checking redundant Multilevel Association Ru

using Concept hierarchy we define association rules at different levels of abstraction. This is called as Multilevel association rule. when we are deriving multilevel association rules we may contain redundant rules. when we are deriving multilevel association rules we may contain redundant rules by "ancestor" property Exist different items.

For Ex Consider the concept hierarchy fig 6.4 contains the ancestor property of 'Desktop Compu is the ancestor of 'IBM'.

01

For Ex, Consider the following association rules.

Desktop Computer ⇒ B&w printer ——— ①

IBM   "     "    ⇒ B&w printer ——— ②

here rule ② does'nt give any ad

let 'R₁' is the ancestor of 'R₂' then 'R₁' replaces 'R₂'.

## Mining MultiDimensional Association rules from relational DB & DWH (or) Pattern Evolution Method:-

Here we mine association rules more than one dimension or predicates.

## Multi Dimensional Association Rules:-

For Ex, Consider the 'ALL ELECTRONICS' DB. The Boolean Association rule is i.e., it doesnt contain any dimension.

IBM Desktop Computer ⇒ Sony B&w printer.

Using this we define single dimension association rule as

buys $(x, $ "IBM Desktop Computer") ⇒
buys $(x, $ "Sony B&w printer).

Here we are using one predicate "buys" we can also define multidimensional Association rule as

age $(x, $ "20 --- 29") ∧ occupation $(x, $ "student")
⇒ buys $(x, $ "laptop").

Here, it contains occupation & buys.

Here No predication is, Repeated. Then it is called as Inter Dimensional Association Rule.

If a predicate is repeated, then Called as hybrid Dimensional Association. This is shown below.

$$age(x, ``20--29") \land buys(x, ``laptop") \Rightarrow buys(x, ``Bgw$$

Generally, attributes in database quantitative.

Here 'Categorical' means it Contains finite No. values & this values does Contain any ordering Quantitative means it Contains any numerica values & this values are ordered.
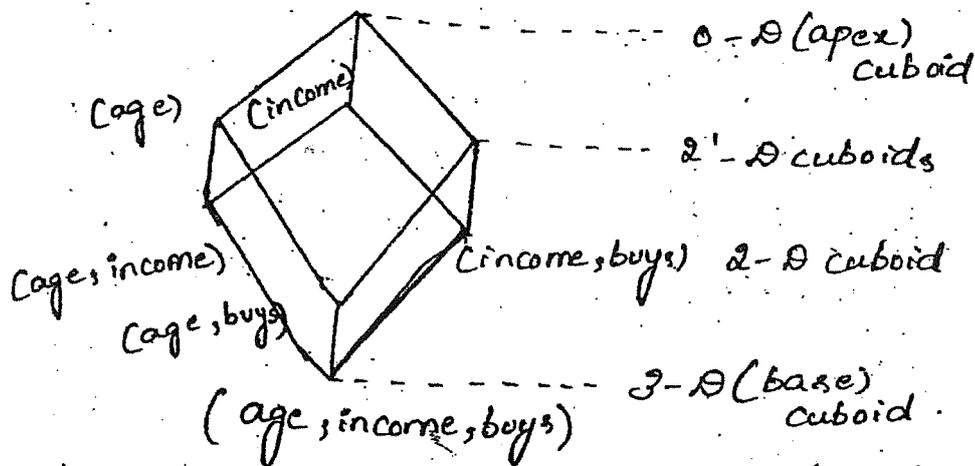
## Mining Multidimensional Association using Stati Discretization of Quantitative utes :-

Discretization is performed before the D. Consider the Concept hierarchy for in is replaced with numerical range of num values such as "20k --- 40k", "41k --- so on.

Then the discretization is static & pre-define This range of values Consider as Categories. We find the values to fall in Each Category Then we find frequent item set. using this freq item set, we derive association rules.

the task relevant data cube. Data cube may Contains the multi dimensions & also using the dataCube, we Easily derive multidimensional association rules.

The data cube Consists of lattice of cuboids. These are multidimensional data Structures & these Contains task relevant data & also Contains aggregated & grouping information. This is Shown below.



lattice of cuboids make 3-D Data cube

Here the base cuboid aggregates task relevant data by age, income, buys.

2-D cuboid i.e, (age, income), it aggregates task relevant data age & income.

1-D cuboid Specifies task relevant data of 1 dimension.

The apex cuboid specifies total No. of transcations in task relevant data.

# Mining Quantitative Association

Quantitative association rules are multi dimensional association rules but they contain numerical attributes. These are derived from discretization.

In Quantitative association rules, the left side it contains 2 Quantitative attributes & in right side it contains Categorical attributes.

These is Shown below:

$$A_{quan_1} \wedge A_{quan_2} \Rightarrow A_{cat}$$

Here, $A_{quan_1}$, $A_{quan_2}$ are quantitative Numerical range values

'$A_{cat}$' is the categorical attribute task relevant data.

For Ex, 2-D Quantitative association rule is

$$age(x, "30 --- 39") \wedge income(x, "29k---30k) \Rightarrow buys(x, "TV")$$

These quantitative association rules derived from ARCS (Association Rules Clustering System)

This ARCS approach contains the following Steps.

1) Binning:- The Quantitave attributes generally contain wide range of values. To Smooth the values, we use following binning methods.

a) Equiwidth:- Here it contains interval size Each bin is Equal.

u) Equidepth - here it contains equal no. of values in each bin.

2). Find Frequent Set :-

To find frequent set, we construct 2-dimensional array containing count values for each category. These values are scanned to get frequent set. These frequent set also satisfy min. support & min. confidence.

3) Clustering association rules:-

For Ex, consider the association rules below.

age $(x, 34)$ ∧ income $(x, "31K---40K")$ ⟹ buys $(x, "TV")$ —①
age $(x, 35)$ ∧ income $(x, "31K---40K")$ ⟹ buys $(x, "TV")$ —②
age $(x, 34)$ ∧ income $(x, "41K---50K")$ ⟹ buys $(x, "TV")$ —③
age $(x, 35)$ ∧ income $(x, "41K---50K")$ ⟹ buys $(x, "TV")$ —④
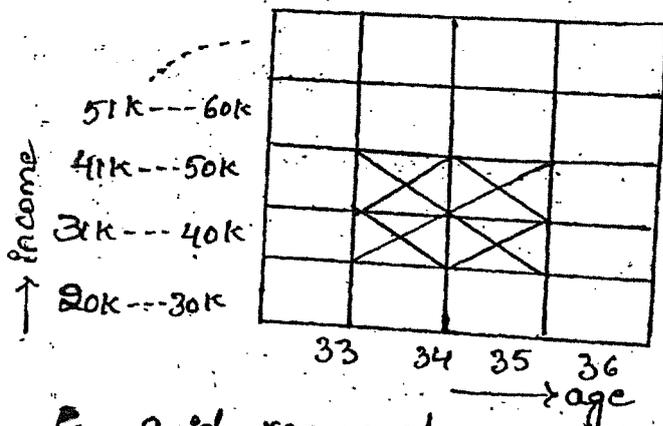
These rules are represented in 2D Grid as below



fig:- grid represents customer who buys "TV"

The above 4 association rules are simply represented in.

age $(x, "34---35")$ ∧ income $(x, "30K---50K")$ ⟹ buys $(x, "TV")$.

# Mining Distance based ~~~~~~~~

For Ex., the data for the 'price' partitioned by using equiwidth & equidepth & this is compared with distance-based partition & this is shown below.

| Prices($) | Equiwidth (width $10) | Equidepth (depth = 2) | Distance-based |
|---|---|---|---|
| 7 | [0,10] | [7,20] | [7,7] |
| 20 | [11,20] | [22,50] | [20,22] |
| 22 | [21,30] | [51,53] | [50,53] |
| 50 | [31,40] | | |
| 51 | [41,50] | | |
| 53 | [51,60] | | |

fig: Data is partitioned by using Equiwidth & Equidepth.

Here, we examined the above table. The distance based partition is the efficient c bcoz it groups the values in such a way th are closer & also closer to specified interval.

Ex. [20,22]

If u consider, Equidepth partition, it group 2 distinct values i.e., [22,50]

If u consider, Equiwidth partition, it conta equal size of intervals but some of these inter does'nt contain any values. i.e, [31,40]

∴ Distance based partition is the best one be this distance-based association contains th Drawback. i.e.,

attribute.

For Ex, Consider the association rule,

item-type $(x, \text{"electronic"}) \wedge$ manufacturer $(x, \text{"Apple"})$

$$\Rightarrow \text{Price}(x, \$200)$$

Here, 'x' is the item. In reality, we always prefer 'Apple' Electronic items are approximately '$200' rather than Exactly '$200'

∴ To define approximate values of an attribute & also derive distance based association rules, we use 2 phase algorithm.

Phase-I :-

Here we define interval for Each cluster.

Phase-II

Each Cluster is verified to identify frequent set.

Using this frequent Set, we derive association rules.

From Association Mining to Correlation-Analysis :-

Even the Strong association rules also uninteresting & gives the wrong information. Therefore we require additional measures that Contains Statistical analysis.

∴ we go for the Correlation-Analysis.

Strong Association Rules are not interesting :-

Consider the "All Electronics" database. we analyze 10,000 transaction. In those 10,000

transactions , 0000 ..........

Purchase the Computer games , 7,500 are videos.
and 4000 transactions Contains the both.

Consider Support count is 30%. & Confidence
is 60%.

buys (x , "Computer games") => buys (x ,"videos")
[Support = 40%. , confidence = 66%.]

$$Support (A => B) = \frac{\#-tuples-Containing-A \& -B}{Total-\#-of-tuples}$$

$$Confidence (A => B) = \frac{\#-tuples-Containing-A-\&-B}{\#-of-tuples-Containing}$$

$$= \frac{4000}{6000} = 2/3 = 0.66 = 66\%.$$

∴ Rule 1 is strong association rule , but the
Probability of purchasing videos is

$$\frac{7500}{10,000} = 0.75 = 75\%.$$

∴ It is greater than 66%.

So , Computer games & videos are inversely
associated - i·e, customer purchase the Computer
games that decreases the videos.

From Association Analysis to Correlation Analy.

The above association rule A => B is uninter
-ting & it gives the wrong information. So we
for the Correlation Analysis.

correlation Analysis ~~~~~ ~~~~~ ~~ y ~ is represented as

$$Corr_{A,B} = \frac{P(A \cup B)}{P(A) P(B)} \quad —— ②$$

If the result is $>1$, then 'A' & 'B' are +ve Correlated. i.e., if we increase 'A', 'B' also increases.

If the result is less than '1', then 'A' & 'B' are -vely Correlated i.e., if we increase 'A' that decreases 'B'.

If the result is Equal to '1' then 'A' & 'B' are independent. i.e., in between 'A' & 'B' there is no Correlation.

Consider the above Ex, then

Computer games = 6000

Videos = 7500

Computer games & Videos = 4000

The probability of purchasing 'computer games' is $P(\{games\}) = \frac{6000}{10000} = 0.60$

This is probability of purchasing 'videos' is $P(\{videos\}) = \frac{7500}{10000} = 0.75$

The Probability of purchasing both 'Computer games' and 'videos' is

$P(\{games, videos\}) = \frac{4000}{10000} = 0.40$

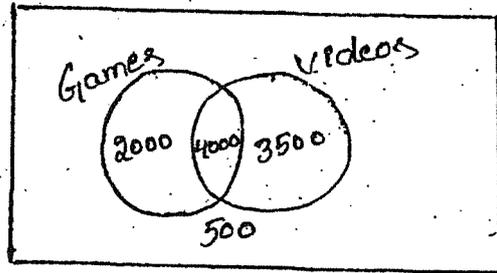∴ Then Correlation between Computer games &

Videos is

$$\text{Corr}_{\text{games, Videos}} = \frac{P(\{\text{games, videos}\})}{P(\{\text{games}\}) \times P(\{\text{videos}\})}$$

$$= \frac{0.40}{0.60 \times 0.75}$$

$$= 0.89$$

∴ Therefore the Computer games & videos are -vely Correlated.

The transactions in DB are Summarized by Contingency Table. This is shown below.



| | Games | Games‾ | Erow |
|---|---|---|---|
| Videos | 4000 | 3500 | 7500 |
| Videos‾ | 2000 | 500 | 2500 |
| Ecol | 6000 | 4000 | 10000 |

fig: 2x2 Contingency Table, Summarizing Transcation Purchased

Here 'games‾' represents dont Contain Computer games

Videos‾ represents dont Contain Computer vid

# Constraint Based Association Mining :-

Many of the DM association rules are uncovered by DM System & also discovered rules are unimportant to the End user. Therefore we go for the Constraint based association Mining.

In Constraint based association Mining we mine association rules based on Constraints these Constraints are specified by End user.

The Constraints are.

**i) Knowledge Constraints:-**

These Constraints specifies the kind of knowledge to be mined.

**ii) Data Constraints :-**

The Constraints specifies the task relevant data sets.

**iii) level Constraints:-**

These Constraints specifies the total No. of levels in concept hierarchy.

**iv) Threshold Constraints :-**

These Constraints specifies the support & confidence values.

**v) Rule Constraints :-**

These constraints specifies what type of association rules are to be mined & also specifies the no. of predicates in each rule.

# Meta Rule for Mining Association Rules:-

Many of the rules are uninteresting & the wrong information. But in Meta rules allows the user syntactic representation of association rules.

Sometimes these meta rules may contain th Constraints to improve the Efficiency of DM System.

The meta rule looks like the following.

$$P_1(x, y) \wedge P_2(x, z) \Rightarrow buys(x, \text{"Educational } S/w")$$—

Here, $P_1$ & $P_2$ are predicates

$$X \longrightarrow Customer$$

$$Y, Z \longrightarrow \text{Numeric Values related to Predicates } P_1 \text{ & } P_2.$$

Whenever we issue this meta rule, the DM System Searches for the rule that is similar to specified meta rule. Then the DM System g the similarity & return the following association rule.

$$age(x, "30---39") \wedge income(x, "30k---39k") \Rightarrow$$

$$buys(x, \text{"Educational } S/w")——②$$

This meta rule is mainly used for guiding the DM process.

For Ex we want to mine inter-dimensional association rules. Then the meta rule is Specified as

$$P_1 \land P_2 \land \text{---} \land P_n \Rightarrow Q_1 \land Q_2 \land \text{---} \land Q_n$$

Here in this meta rule, we are using total

Predicates $P = RHS$

& also $P_i$ where $i = 1 \text{---} \& $

$Q_j$ where $j = 1 \text{----} n$

## Additional Constraint Rules:-

1. This additional Constraint rules Contains Set of relations. These relations are specified by the user by using aggregate functions. These rules are mainly used in multidimensional association rule Mining.

Consider the "All Electronics" multidimensional Sales db. This Contains the following relations.

Sales (cust-Name, item-Name, trans-id)

lives (cust-Name, City, Street)

item (item-Name, price)

transaction (trans-id, day, month, year)

Here lives, item & transaction tables are dimensional tables & fact table is Sales. Here we specify the relation by using 3.

i.e., Cust-Name, item-Name, trans-id.

Ex. we want to mine association query i.e, find the Sales of cheap items (Sum of price < $100) & also find Sales of Expensive items (min of Price $500) for Vancouver

Customer in 1999.

  This query is specified in DMQL as
mine association as lives(C, "Vancover")
from Sales
where S. year = 1999 and T. year = 1999 group
by C.
having Sum(I. price) and min(I. price) will

Support threshold = 1 %.
with Confidence threshold = 50 %.

# 7. Classification & Prediction
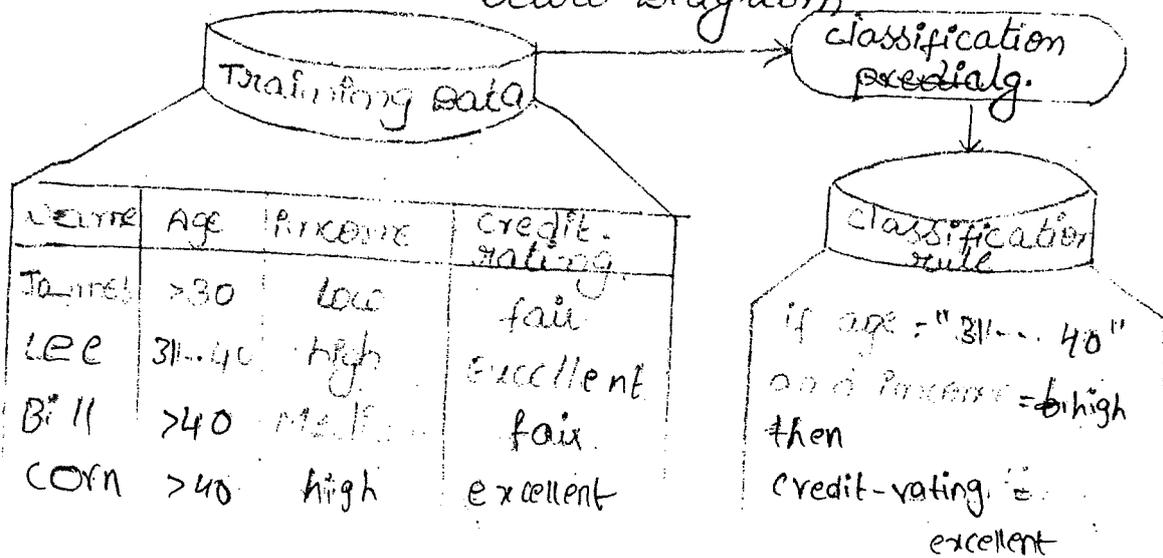
7) What is classification & what is prediction:

The data " contains 2 steps.

i) construct the model by using predefined dataset : This model uses the tuples in the db & also each tuple must fall in class. This is determined by using attribute. This attribute is called as class label attribute. Then the training data is collected & analyzed by using classification algorithm i.e, learned knowledge is represented into classification rule.

ii) classification:

In this, classification, first verify the accuracy of the " rule. its accuracy is ok. Then it will be applied for new data. This is shown in below diagram.

| Name | Age | Income | Credit rating |
|------|-----|--------|---------------|
| Jaime | >30 | low | fair |
| Lee | 31..40 | high | excellent |
| Bill | >40 | med | fair |
| Corn | >40 | high | excellent |

Classification predialg.

Classification rule

if age = "31... 40"
and Income = high
then
Credit-rating = excellent

rule

Test data

New data

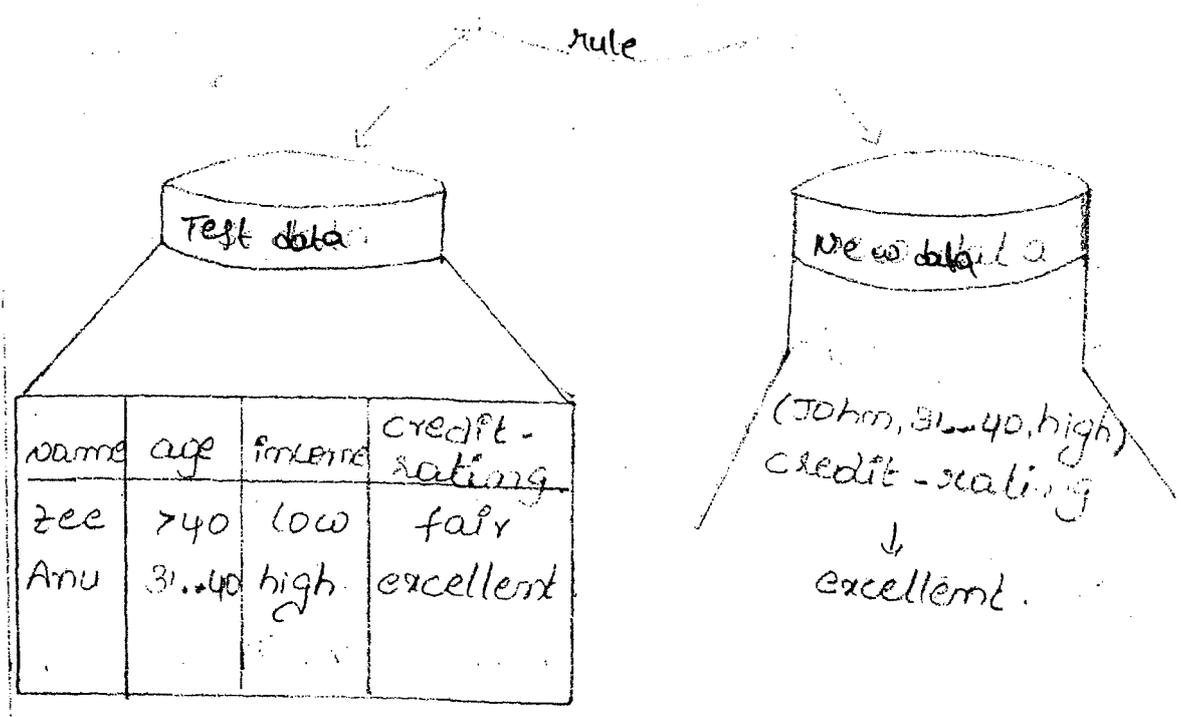| name | age | interest | credit-rating |
|------|-----|----------|---------------|
| zee | >40 | low | fair |
| Anu | 31..40 | high | excellent |

(John, 31-40, high) credit-rating
↓
excellent.

fig 71 The Data classification.

a) Learning:

Training data is analyzed by using classification algorithm & here class label attr is credit-rating. Then learning data is transformed into classification rule.

b) classification:

The classification rule is applied on to data, to find its accuracy. If it is accurate i.e., applied for new data.

In fig 71 (a) it analyzes the existing data of customers. & using this class label of ne customers is predicted.

The prediction is different from classification. The prediction allows construct the model & use the model to find the class label of unlabelled sample. The prediction uses the classification & regression.

1. **classification :**

using classification we predict the discrete values.

2. **Regression :**

using this, we predict continuous values.

## 7.2 Issues Regarding classification & Predictions

we use the following issues.

### 7.2.1. Preparing the data for classification & prediction :

Here, we initially examine the data to improve the efficiency, accuracy, scalability of classification & prediction.

i) **Data cleaning :**

Here noisy or missing data values are smoothed by using several smoothing techniques.

ii) **Relevant Analysis :**

Here, irrelevant attributes are removed from the data before applying classification & prediction.

iii) Data Transformation:

we already know higher levels contains the more frequent items rather than the low level. This is achieved by using concept hierarchy.

. The data transformation specifies generaliz-tion levels to concept hierarchy.

## 7.2.2 Comparing classification & Prediction Methods:

The classification & prediction methods are compared by using following methods.

### 1. predictive Accuracy:

This specifies how we use the model to predict class label of new data. As well as it also specifies the prediction must be accurate.

### 2. Speed:

It specifies, to predict the class label it must take the less no of comparisons.

### 3. Robustness:

This specifies that we predict the class label, if the model contains some noisy data or missing values.

### 4. Scalability:

It specifies model is constructed from large db.

### 5. Interpretability

It specifies level of understanding &

. also specifies how we interpret the model

7.3. classification by Decision Tree Induction:

Decision Tree means flow chart like tree structure. The internal nodes specifies test on attribute. Th

The branch is the outcome of the test result. The leaf node represents class. The top node of the tree is considered as a root node. The decision tree for concept hierarchy buys-computer is shown below.

age ?
<=30     31..40.     >40.
Student ?     yes.     credit-rating?
no.     yes          excellent.     fair
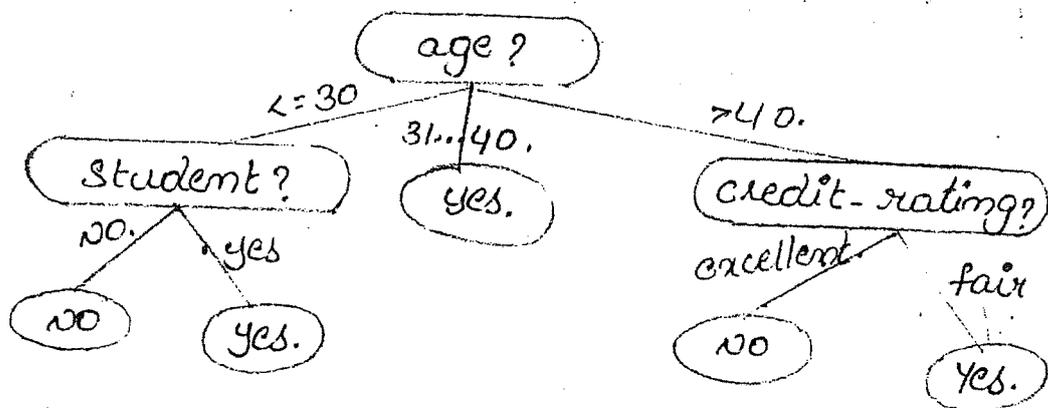NO     yes.                NO     Yes.

fig: 7.3. Decision Tree for Concept buys-computer indicates whether the customer buys computer or not.

The internal nodes represents test on attribute. The leaf node - class ( buys-computer = no or buys-computer = yes).

This decision tree contains the adv. of we classify unknown sample & also this decision tree directly converted into association

i.e basic decision tree alg

Algorithm: Generate - decision - tree.
Input: The training samples.
Output: A decision tree.
Method:

1. create a node 'N'.
2. if all the samples are same class 'c' then
3. return 'N' as a leaf labelled with class 'c'
4. if attribute - list is empty
5. return N as a leaf labelled with most common class.
6. select test-attribute from attribute-list with highest inf. gain.
7. Label node 'N'
8. Each known value $a_i$ of test-attribute
9. grow branch from node 'N'.
10. Let 'Si' be the samples in test-attribute $= a_i$.
11. if Si is empty then
12. attach the leaf labelled with most common class.
13. else attach node & recall generate-decision-tree.

## 7.3.1. Decision Tree Induction:

The above decision tree alg. follows the top down approach & divide & conquer approach.

The alg is explained by using the follow...

- Create a model (step 1)
- if the samples are same class, then return all the class labels related to that class (step 2,3)
- otherwise select the test-attribute by using entrophy based inf. gain.
- For each known value of test-attribute, create branch
- if the samples of test-attribute are empty, then return most common class labels.
  otherwise again recall the decision tree alg.

## Attribute Selection:

Let 'S' be the training samples. It contains the 'm'-classes. Let '$S_i$' be any sample in 'S' with class label '$c_i$' for $i = 1 \ldots m$. then expected inf. is needed to classify the given sample

$$I(S_1, \ldots S_m) = -\sum_{i=1}^{m} P_i \log_2(P_i)$$

Here '$P_i$' is the probability of given sample fall in class label '$c_i$'. It is calculated by $\frac{S_i}{S}$.

consider an attribute 'A' with values $\{a_1 a_2 \ldots a_v\}$. This attrib. 'A' partition the sample 'S' into subsets $\{S_1, S_2 \ldots S_v\}$

Let '$S_j$' contains '$S_{ij}$' samples of class '$c_i$'. then entrophy of 'A'

$$E(A) = \sum_{j=1}^{} \frac{S_{ij} + \cdots + S_{mj}}{S} \ I(S_{ij} \cdots , S_{mj})$$

∴ The inf. gain through this partition 'A' as

$$G(A) = I(S_1 \cdots S_m) - E(A).$$

The attribute with the highest inf. gain is treated as the most relevant attrib. that is used to construct the decision tree.

24|9

Ex. consider the data tuples from 'All Electronics' database.

Table: 7.3. Data tuples from All Electronics DataBase.

| R.I.D. | age | income | Student | credit-rating | class: buys-computer |
|---|---|---|---|---|---|
| 1. | <=30 | low | no | fair | no |
| 2. | <=30 | high | no | excellent | no |
| 3. | 31...40 | low | no | fair | yes |
| 4. | >40 | medium | yes | excellent | no |
| 5. | >40 | high | no | fair | yes |
| 6. | >40 | low | no | fair | yes |
| 7. | 31..40 | high | yes | fair | yes |
| 8. | 31..40 | low | no | excellent | yes |
| 9. | 31..40 | medium | yes | fair | yes |
| 10. | <=30 | " | yes | fair | yes |
| 11 | <=30 | low | " | excellent | yes |
| 12. | <=30 | high | no | fair | no |
| 13. | >40 | low | no | excellent | no |
| 14 | >40 | high | yes | fair | yes |

Here class label attribute is 'buys-computer'. It contains 2 values {yes, no}. The class '$c_1$' is represented by samples of 'yes', the class '$c_2$' is " " " " 'no'.

The above table contains '9' samples of 'yes' & '5' samples of 'no'. Therefore, the expected inf. is needed to classify the given sample.

$$I(S_1, S_2) = -\frac{9}{14} \log_2 9/14 - \frac{5}{14} \log_2 5/14$$

$$= 0.94.$$

Then we have to calculate the entropy value for each attribute starting with attribute 'age'. This attribute contains the different branches & fall in 2 classes {yes or no}.

for age : "<=30"

$S_{11} = 2$. $S_{21} = 3$.     ⇒ total 5

$$I(S_{11}, S_{21}) = -\frac{2}{5} \log_2 2/5 - \frac{3}{5} \log_2 3/5$$

$$= 0.97.$$

for age = "31...40"

$S_{12} = 4$, $S_{22} = 0$.       ⇒ total ⇒ 4

$$I(S_{12}, S_{21}) = 0.$$

for age = ">40"          Total ⇒ 5

$S_{13} = 3$. $S_{23} = 2$.

$$I(S_{13}, S_{23}) = -\frac{3}{5} \log_2 3/5 - \frac{2}{5} \log_2 2/5$$

$$0.97.$$

∴ Entrophy of age is

$$E(age) = \frac{5}{14} I(S_{11}, S_{21}) + \frac{4}{14} I(S_{12}, S_{22}) + \frac{5}{14} I(S_{13}, S_{23})$$

$$= \frac{5}{14}(0.97) + \frac{4}{14}(0) + \frac{5}{14}(0.97)$$

$$= \frac{10}{14}(0.97)$$

$$= \frac{9.7}{14}$$

$$= 0.69.$$

$$Gain(age) = I(S_1, S_2) - E(age)$$

$$= 0.94 - 0.69$$

$$\boxed{= 0.25.}$$

Similarly, we can calculate the inf. gain for income is

$$gain(income) = 0.02$$
$$gain(student) = 0.15$$
$$gain(credit-rating) = 0.04$$

→ gain (class: buys - computer):

Therefore, age contains the highest inf. gain. it is the most relevant attribute. Therefore, it is selected as the test attribute. Using this we construct the decision tree.
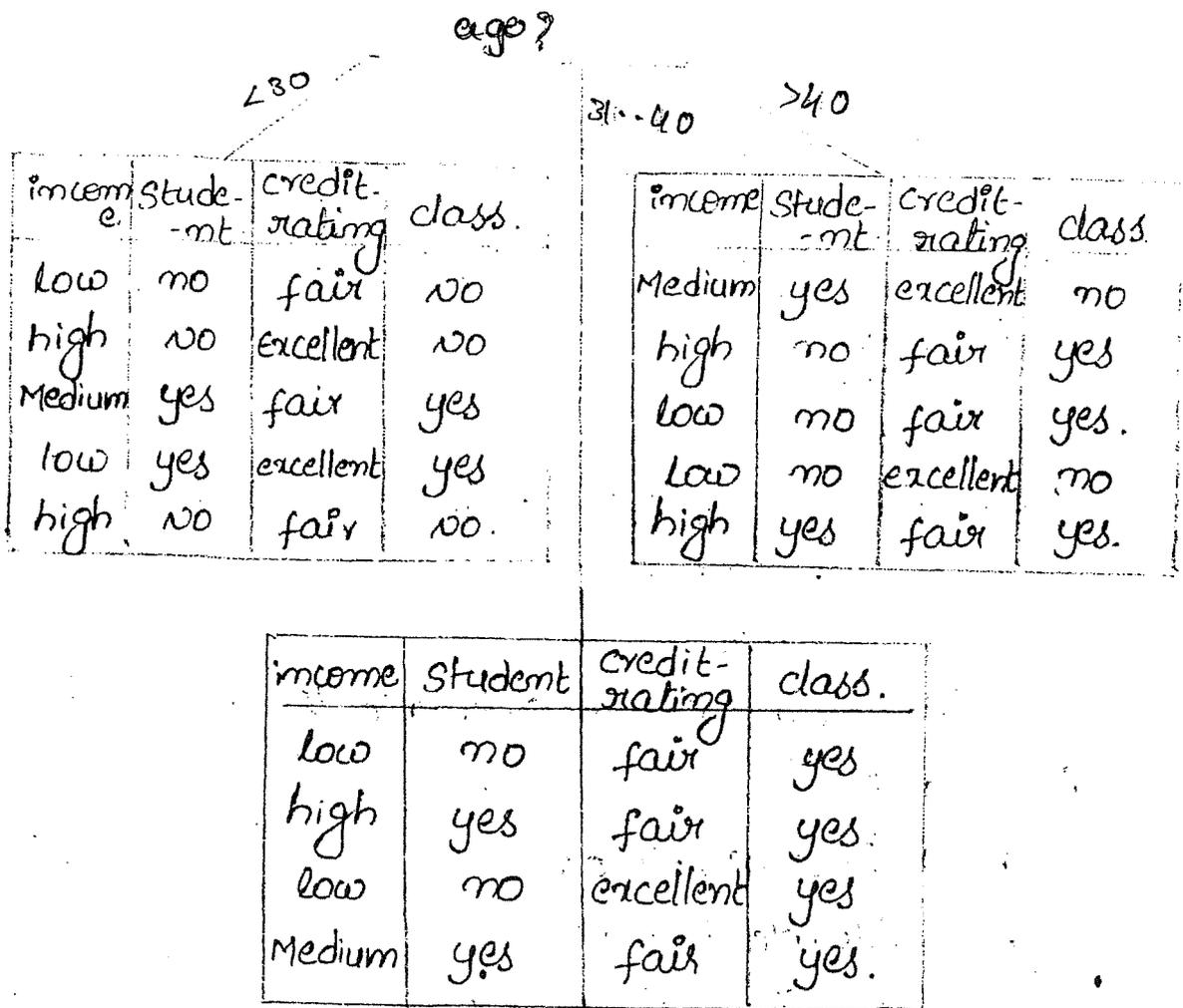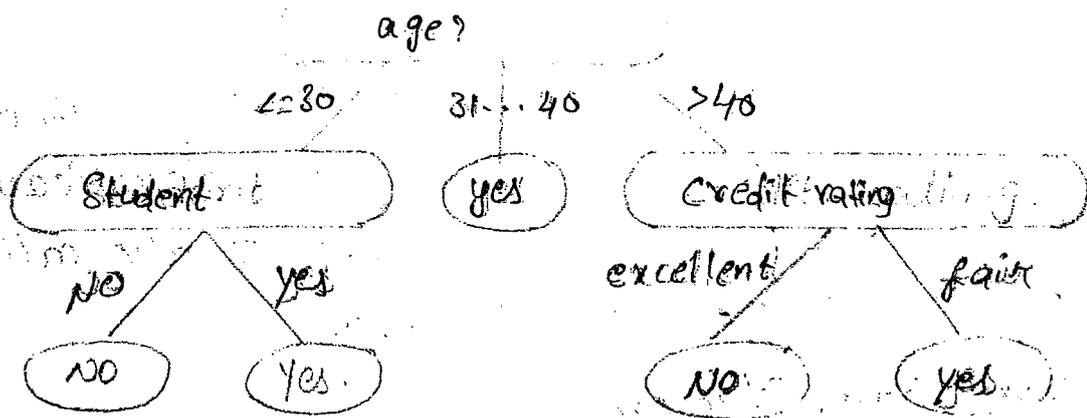
age?

<30       31..40       >40

| income | stude-mt | credit-rating | class |
|--------|----------|---------------|-------|
| low | no | fair | no |
| high | no | excellent | no |
| Medium | yes | fair | yes |
| low | yes | excellent | yes |
| high | no | fair | no |

| income | stude-mt | credit-rating | class |
|--------|----------|---------------|-------|
| Medium | yes | excellent | no |
| high | no | fair | yes |
| low | no | fair | yes |
| Low | no | excellent | no |
| high | yes | fair | yes |

| income | Student | credit-rating | class |
|--------|---------|---------------|-------|
| low | no | fair | yes |
| high | yes | fair | yes |
| low | no | excellent | yes |
| Medium | yes | fair | yes |

fig. 7.3. The 'age' attribute is the highest inf. gain. ∴ it is the root node & branches grow, with 'age' values.

Here the age 31-40 contains same samples i.e., class 's'.

The final decision tree based on the alg. is shown below.

fig: Decision tree for concept buys-computer.

age?

≤30          31...40          >40

Student          yes          Credit rating

NO / yes

NO          Yes

excellent / fair

NO          yes

## 7.3.2 Extracting classification rules from decision trees:

The knowledge, represented in decision tree extracted & represented in the form of <u>classifi-cation IF-THEN rule</u>. One rule is derived for one path starting from root node to leaf node

Consider the decision tree fig.7.8.

The classification rules are,

IF age = "≤30" and Student = "no" THEN
     buys - computer = "no".

IF age = "≤30" and Student = "yes" THEN
     buys - computer = "yes".

IF age = "31...40" THEN buys - computer = "yes".

IF age = ">40" and credit - rating = "excellent"
     THEN buys - computer = "no"

IF age = ">40" and credit - rating = "fair"
     THEN buys - computer = "yes".

# 7.4. Bayesian Classification.

It is used to predict the class label of unknown sample. It is based on Bayes Theorem & then we apply simple Bayesian classification. The simple Bayesian classification is also called as naive Bayesian classification.

## 7.4.1. Bayes Theorem:

Let 'x' be an unknown sample & let 'H' be the hypothesis. Using this 'H', we predict the class label of 'x'. Finally we determine P(H/x). This is called as posterior probability. i.e., probability of 'H' on condition x.

For ex, consider the data samples as 'fruits'. These are described by using 'colour' & 'shape'. Let 'x' is the red colour & round shape. Then 'H' gives, x is apple. Then:

P(H/x) gives the confidence of x is Apple because its colour is red & it is round shape.

To determine P(H/x) it requires the additional imf. i.e., P(H) gives the prior probability of H; P(x) gives the prior probability of 'x' and P(x/H) i.e., P(x) on condition H.

Therefore using this Bayes Theorem, we define the

$$P(H/x) = \frac{P(x/H) \cdot P(H)}{P(x)} \quad \text{——} ①$$

## 7.4.2 Naive Bayesian classification:

"Naive Bayesian classification" or simple Bayesian classification contains the following steps.

### Step 1:

Consider the 'm' class i.e., $c_1, c_2 \ldots c_m$. Let 'x' be the unknown sample. Then according to naive Bayesian classification, the unknown sample 'x' is assigned to class '$c_i$' if & only if $P(c_i/x) > P(c_j/x)$ for $1 \leq j \leq m$; $j \neq i$.

Thus we maximize probability of $P(c_i/x)$. But according to Bayes Theorem,

$$P(c_i/x) = \frac{P(x/c_i) \, P(c_i)}{P(x)} \quad \text{---} \quad (2)$$

### Step 2:

The $P(x)$ is same for all the classes. Then we maximize $P(x/c_i) \cdot P(c_i)$.

For ex, prior probabilities are not known then all the classes are equal i.e.,

$$P(c_1) = P(c_2) = \ldots = P(c_m).$$

Then we maximize $P(x/c_i)$ otherwise, we maximize $P(x/c_i) \, P(c_i)$

$$P(c_i) = S_i/S,$$

$S_i$ = class $c_i$

$S$ = Total no. of samples.

Step 3:

If there is no dependent relationship among the attributes. Then

$$P(x/c_i) = \prod_{k=1}^{n} P(x_k/c_i)$$

Here, we find the $P(x_1/c_i), P(x_2/c_i) \ldots P(x_m/c_i)$

a) If $A_k$ is categorical, then

$$P(x_k/c_i) = \frac{S_{ik}}{S_i}$$

where $S_{ik}$ is the training samples of class '$c_i$' having value $x_k$ for $A_k$.

b) If $A_k$ is continuous-valued, then we use the gaussian distribution.

$$\therefore P(x_k/c_i) = g(x_k, \mu_{c_i}, \sigma_{c_i})$$

$$= \frac{1}{\sqrt{2\pi}\,\sigma_{c_i}} e^{-(x_k - \mu_{c_i})^2 / 2\sigma_{c_i}^2}$$

where $\mu_{c_i}, \sigma_{c_i}$ are the mean & standard deviation of $A_k$.

Step 4:

The unknown sample 'x' is classified by determining $P(x/c_i) \cdot P(c_i)$ for each class '$c_i$'.

we assign unknown sample 'x' to class '$c_i$' if & only if $P(x/c_i) \cdot P(c_i) > P(x/c_j) \cdot P(c_j)$

for $1 \le j \le m, j \ne i$.

**Example:**

Consider the table 7.3.1 in that the data samples are represented by attributes age, income, student, credit_rating. The class label is buys-computer.

Buys-computer contains 2 classes {yes, no}. The class $c_1$ is represented by buys-computer = "yes"

" " $c_2$ " " " = "no"

The unknown sample $x$,

$x = $ (age = "$<=30$", income = "medium", student = "yes" credit_rating = "fair").

We have to maximize $P(x/c_i) \cdot P(c_i)$ for $i = 1, 2$.

The prior probabilities are calculated by

$p$(buys-computer = "yes") = $9/14$ = 0.64.

$p$(buys-computer = "no") = $5/14$ = 0.36.

To calculate $P(x/c_i)$ for $i = 1, 2$., we first of all calculate the conditional probabilities.

$p$(age = "$<=30$" / buys-computer = "yes") = $\frac{2}{9}$ = 0.22.

$p$(age = "$<=30$" / buys-computer = "no") = $\frac{3}{5}$ = 0.6

$P$(income = "medium" / buys-computer = "yes") = $\frac{2}{9}$ = 0.22

$P$(income = "medium" / buys-computer = "no") = $\frac{1}{5}$ = 0.2.

$p$(student = "yes" / buys-computer = "yes") = $\frac{5}{9}$ = 0.56.

$P$(student = "yes" / buys-computer = "no") = $\frac{1}{5}$ = 0.2.

$P(\text{credit-rating} = \text{"fair"} \mid \text{buys-computer} = \text{"yes"})$
$$= 7/9 = 0.78.$$

$P(\text{credit-rating} = \text{"fair"} \mid \text{buys-computer} = \text{"no"})$
$$= 2/5 = 0.4.$$

$\therefore P(x \mid \text{buys-computer} = \text{"yes"}) = 0.22 \times 0.22 \times 0.56 \times 0.78$
$$= 0.02.$$

$P(x \mid \text{buys-computer} = \text{"no"}) = 0.6 \times 0.2 \times 0.2 \times 0.4$
$$= 0.0096.$$

$\underset{P(x\mid c_i)}{P(x \mid \text{buys-computer} = \text{"yes"})} \cdot \underset{P(c_i)}{P(\text{buys-computer} = \text{"yes"})}$
$$= 0.02 \times 0.64$$
$$= 0.0128.$$

$P(x \mid \text{buys-computer} = \text{"no"}) \cdot P(\text{buys-computer} = \text{"no"})$
$$= 0.0096 \times 0.36$$
$$= 0.003456.$$

Therefore,

$P(x \mid \text{buys-computer} = \text{"yes"}) \cdot P(\text{buys-computer} = \text{"yes"})$
$> P(x \mid \text{buys-computer} = \text{"no"}) \cdot P(\text{buys-computer} = \text{"no"})$

i.e., $c_1 > c_2$.

Then, unknown sample 'x' is assigned to buys-computer = "yes".

✓

classification by Backpropagation :

The Back propagation based on the neural o/w Learning algorithm.

Here neural means set of input & output lines are connected & also each line assigns the weighted value. Then in learning process, the o/w learns by Modi--fying the weights of each connection. Then this is used to identify correct class labels of i/p samples.

The back propagation mainly applies on multilayer Feed-Forward neural o/w. & then we define the o/w topology. Finally, we apply the Back propagation.

7.5.1. Multilayer Feed-Forward neural o/w :
This " " " " " "
is shown below.



Fig 7.5. A multi-layer Feed-Forward neural o/w

The training sample $x = \{x_1, x_2 \cdots x_i\}$ form the i/p Layer.

each pair of adjacent layers. It contains the weighted connection.

For ea. $w_{ij}$ is the weight from i to j.

The multilayer means, it contains only one hidden layer. Feed-Forward means, the weighted connections never come back.

The o/ps of I/p layer are the I/ps of hidden layer. & so on.

## 7.5.2 Defining a o/w Topology:

Here end user has to define the o/w topology by specifying no. of units in i/p layer, no. of units in hidden layer, no. of units in o/p layer.

## 7.5.3 Back Propagation:

In this, the samples are repeated continuously. For each data sample, we modify the weights to minimize the error b/w o/w prediction & well-known class label.

The back propagation contains the following steps.

i. Initialize the weights:

Here the weighted connection b/w each pair of adjacent layers, we assign small random value.

for ex, this value ranging from
(-1 to 1 (or) -0.5 to 0.5)

ii) Forward the selected inputs:

In this step, we find net i/p & o/p for each unit & also each layer.

For ex. unit is 'j' in hidden or o/p layer. then its net i/p is calculated by using a formula.

$$\boxed{I_j = \sum_i w_{ij}\, O_i + \theta_j} \quad\text{——}\textcircled{1}.$$

Here $w_{ij}$ — is weight from i to j

$O_i$ — o/p of unit 'i'

$\theta_j$ — o/p of Bias. using this we differentiate the diff. i/p units.

$I_j$ — net i/p for unit j.

These i/p s are used by hidden or o/p layers. Then these layers apply activation function. This is shown in below.



fig 7.5.1 Activation function

Let unit $j$ & its net I/p $I_j$. Then $O_j$ is the o/p for unit $j$.

$$O_j = \frac{1}{1 + e^{-I_j}} \quad — ②$$

iii) Propagate error ~~forward~~ Backward:

To propagate the error ~~forward~~ Backward, we ~~update~~ modify the weights & Bias values.

For unit $j$ the error '$Err_j$' for o/p layer is calculated by using a formula

$$Err_j = O_j(1-O_j)(T_j - O_j) \quad — ③$$

Here $O_j$ — o/p of unit $j$.

$T_j$ — o/p of true value for known class label.

The error at hidden layer for unit $j$ is calculated by using the formula,

$$Err_j = O_j(1-O_j) \sum_k Err_k W_{jk} \quad — ④.$$

The weights are updated by using the formula.

$$\Delta W_{ij} = (l) Err_j O_i$$
$$W_{ij} = W_{ij} + \Delta W_{ij} \quad — ⑤$$

Here $\Delta W_{ij}$ — change of weight $W_{ij}$.

$l$ — linear rate. it is a constant value & ranging f/w $0.0$ to $1.0$.

IIly Bias are updated by using the formula

$$\Delta \theta_j = (l) Err_j$$

$$\theta_j = \theta_j + \Delta \theta_j \quad \text{——} \quad \text{⑥}$$

Ex: consider the following multilayer feed forward neural $o/\omega$



fig. 7.6. Multilayer Feed-Forward neural $o/\omega$.

The training sample,

$x = (1,0,1)$ whose class label is '1' & linear rate is 0.9. The initial i/ps, weights & Bias values are as shown in below.

| $x_1$ | $x_2$ | $x_3$ | $\omega_{14}$ | $\omega_{15}$ | $\omega_{24}$ | $\omega_{25}$ | $\omega_{34}$ | $\omega_{35}$ | $\omega_{46}$ | $\omega_{56}$ | $\theta_4$ | $\theta_5$ | $\theta_6$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0.2 | -0.3 | 0.4 | 0.1 | -0.5 | 0.2 | -0.3 | -0.2 | -0.4 | 0.2 | 0.1 |

Table: 7.7. Initial I/p weight & Bias values.

## calculation of net I/p & o/p values:

| unit j | Net Input, $I_j$ | Output, $O_j$ |
|---|---|---|
| 4 | $0.2+0-0.5-0.4=-0.7$ | $1/(1+e^{0.7})=0.332$ |
| 5 | $-0.3+0+0.2+0.2=0.1$ | $1/(1+e^{-0.1})=0.525$ |
| 6 | $(-0.3)(0.332) - (0.2)(0.525) = -0.105$ | $1/(1+e^{0.105}) = 0.474$ |

## calculation of Error:

| unit $j$ | Error |
|---|---|
| 6 | $(0.474)(1-0.474)(1-0.474) = 0.1311$ |
| 5 | $(0.525)(1-0.525)(0.1311)(-0.2) = -0.0065$ |
| 4 | $(0.332)(1-0.332)(0.1311)(-0.3) = -0.0087$ |

## calculation of weights & Bias updates:

| weight/bias | new value |
|---|---|
| $w_{46}$ | $(-0.3) + (0.9)(0.1311)(0.332) = -0.260$ |
| $w_{56}$ | $(-0.2) + (0.9)(0.1311)(0.525) = -0.138$ |
| $w_{14}$ | $0.2 + (0.9)(-0.0087)(1) = 0.192$ |
| $w_{15}$ | $-0.3 + (0.9)(-0.0065)(1) = -0.3058$ |
| $w_{24}$ | $0.4 + (0.9)(-0.0087)(0) = 0.4$ |
| $w_{25}$ | $0.1 + (0.9)(-0.0086)(0) = 0.1$ |
| $w_{34}$ | $-0.5 + (0.9)(-0.0087)(1) = 0.5078$ |
| $w_{35}$ | $0.2 + (0.9)(-0.0065)(1) = 0.194$ |
| $\theta_6$ | $0.1 + (0.9)(0.1311) = 0.2179$ |
| $\theta_5$ | $0.2 + (0.9)(-0.0065) = 0.194$ |
| $\theta_6$ | $-0.4 + (0.9)(-0.0087) = -0.407$ |

## Termination condition:

$\Delta w_{ij}$ becomes the small i.e., it is less than the specified Threshold range.

# 8. Cluster Analysis.

## 8.1. What is cluster Analysis?

The process of storing similar objects in one cluster & dissimilar objects in another cluster. The data clusters are shown in below.



fig. 8.22-⊕ customer Data w.r. to customer locations in city. Here Three data clusters are specified & centre of cluster is marked with '+'.

The following are the critical requirements for cluster in data Mining.

**(1) Scalability:**

The clustering algorithms are efficient for low data, but DM system contains large volume of data. to handle this large volume of data, the clustering analysis alg.s must be highly scalable i.e., efficient.

# 8.2 Partitioning Methods.

The simplest and most fundamental version of cluster analysis is partitioning, which organizes the object of a set into several exclusive groups or clusters.

Given a data set $D$, of $n$ objects and $K$, the no. of clusters to form, a partitioning algorithm organizes the objects into $K$ partitions $K \leq n$, where each partition represents a cluster.

## Partitioning methods :-

1) K-means : A centroid - Based Technique.
2) K- medoids: A Representative object - Based Technique.

## 1) K-means : A centroid - Based Technique :-

K-means clustering intends to partition $n$ objects into $K$ clusters in which each object belongs to the cluster with the nearest mean. This method produces exactly $K$ different clusters of greatest possible distinction.

The best number of cluster $K$ leading to the greatest separation is not known as a prior and must be computed from the data. The objective of k-means clustering is to minimize total intra - cluster variance or the squared error function;

———

$$E = \sum_{i=1}^{K} \sum_{p \in c_i} dist(p, c_i)$$

(top-left margin)

$$E = \sum_{i=1}^{K} \sum_{p \in c_i} dist(p, c_i)^2,$$

where $E$ is the sum of the squared error.

$P$ is the point in space representing a given object.

$c_i$ is the centroid of cluster $c_i$.

## Algorithm:

K-means algorithm for partitioning, where each cluster is represented by the mean value of the objects in the cluster.

### Input:

K: the no. of clusters,

D: a data set containing n objects.

### Output:-

A set of K clusters.

### Method:-

1) arbitarily choose K objects from D as the initial cluster centers;

2) repeat

3) assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;

4) update the cluster means, i.e., calculate the mean value of the objects for each cluster;

5) until no change;

(a) Initial clustering.



(b) Iterate.



(c) Final clustering.

clustering of a set of objects using the k-means

Consider a set of objects located in 2-D space, as depicted in fig (a). Let $K=3$, i.e., the objects to be partitioned into 3 clusters.

According to the algorithm, we arbitrarily choose 3 objects as 3 initial cluster centers, cluster centers are marked by '+'. Each object is assigned to a cluster based on the cluster center to which it is the nearest.

Next the cluster centers are updated i.e, the mean value of each cluster is recalculated based on the current object in the cluster.

$$E = \sum_{i=1}^{K} \sum P_i \otimes |P_i - O_i|$$

using this new cluster centers, the objects are re distributed to the cluster based on which cluster center is the nearest. The process of iteratively re assigning object to clusters to improve the Partitioning is referred to as iterative relocation.

The time complexity of the k-means algorithm is $O(nkt)$, where $n$ is the total no. of objects, $k$ is the no. of clusters, & $t$ is the no. of iterations. Therefore, the method is relatively scalable & efficient in processing large data sets.

## 2, K-Medoids: A Representative Object-Based Technique:-

The k-means algorithm is sentive to outliers because such objects are far away from the majority of the data, they can dramitically distroy the mean value of the cluster. This affects the assignment of other objects to clusters. This effect is particularly due to the use of the squared error function.

The k-medoids is an iterative clustering algorithm which iterates until each representative object is the mediad or most centrally located object of its cluster.

$$E = \sum_{}^{K} \sum |P - O_i|$$

The most common realisation of k-medoid clustering is Partitioning Around Medoids (PAM)

## K-Medoids Algorithm (PAM) :—

**Input :**

    K: the no. of clusters

    D: a dataset containing n objects

**Output :—**

    A set of K clusters

**Method :—**

1) Arbitrary choose K objects from D as representative objects.

2) Repeat

3) Assign each remaining object to the cluster with the nearest representative object

4) For each representative object $O_j$

5) Randomly select a non representative object $O_{random}$.

6) Compute the total cost $S$ of swapping representative object $O_j$ with $O_{random}$.

7) if $S < 0$ then replace $O_j$ with $O_{random}$

8) Until no change.

$$O(K(n-K)^2).$$

The complexity of each iteration is $O(K(n-K)^2)$. For large values of n and K, such computation becomes very costly.

—

Advantages

## Advantages:—

* K-Medoids method is more robust than k-Means in Presence of noise and outliers.

## Disadvantages:—

* K-medoids is more costly that the k-Means Method

* Like k-means, k-medoids requires the user to specify k

* It does not Scale well for large data Sets.

## Example:—

Data objects

| | $A_1$ | $A_2$ |
|---|---|---|
| $O_1$ | 2 | 6 |
| $O_2$ | 3 | 4 |
| $O_3$ | 3 | 8 |
| $O_4$ | 4 | 7 |
| $O_5$ | 6 | 2 |
| $O_6$ | 6 | 4 |
| $O_7$ | 7 | 3 |
| $O_8$ | 7 | 4 |
| $O_9$ | 8 | 5 |
| $O_{10}$ | 7 | 6 |



Goal: Create two clusters
choose randomly two medoids

$O_2 = (3,4)$
$O_8 = (7,4)$

* The absolute error criterion (for the set of medoids (O7,O8))

* Assign each object to the closest representative object

* using L1 Metric (Manhattan), we form the following clusters

$$Cluster1 = \{O_1, O_2, O_3, O_4\}$$
$$cluster2 = \{O_5, O_6, O_7, O_8, O_9, O_{10}\}$$

* Compute the absolute error criterion [for the set of Medoids $(O_2, O_8)$]

$$E = \sum_{i=1}^{k} \sum_{P \in C_i} |P - O_i| = |O_1 - O_2| + |O_3 - O_2| + |O_4 - O_2| +$$
$$|O_5 - O_8| + |O_6 - O_8| + |O_7 - O_8| + |O_9 - O_8| +$$
$$|O_{10} - O_8|.$$

* The absolute error criterion [for the set of Medoids $(O_2, O_8)$]

$$E = (3+4+4) + (3+1+1+2+2) = 20$$

* choose a random object $O_7$

* Swap $O_8$ and $O_7$

* Compute the absolute error criterion [for the set of medoids $(O_2, O_7)$]

$$E = (3+4+4) + (2+2+1+3+3) = 22.$$



* Compute the cost function

Absolute error $[O_2, O_7]$ - Absolute error $[O_1, O_8]$

$$S = 22 - 20$$

$S > 0 \Rightarrow$ it is a bad idea to replace $O_8$ by $O_7$

## 8.3 Hierarchical Methods:—

A hierarchical clustering method creates a hierarchical structure from data objects. The hierarchical clustering methods can be classified into two types. They are,

1) Agglomerative (bottom-up)
2) Divisive (top-down).

## Agglomerative Hierarchical Clustering:—

This involves merging of data objects. It initiates with each object forming its own group/cluster. For every pair of cluster, some value of dissimilarity is computed,

then the clusters are merged into larger and larger clusters until all the clusters are merged into one largest cluster or until a termination condition is reached. The merging of clusters is carried-out based on the Euclidean distance b/w any two objects from different clusters. The user can set the termination criteria by fixing the desired no. of clusters.

## Divisive Hierarchical clustering:-

This clustering method is also known as top-down method. This involves division of objects cluster into smaller parts. It initiates with all the objects in the same cluster. This single cluster is splitted into smaller clusers, until each object is one cluster or until a termination condition is reached.

Example for AGNES and DIANA :— AGNES (Agglomerative Nesting)

A dataset of Seven objects {p, q, r, s, t, u, v}



First AGNES on Data objects {p, q, r, s, t, u, v}

connec

alg.

with

dens

of

No'

be



1, (DIANA) Divisive Analysis on Data Objects $\{P, Q, r, s, t, v\}$

In 'AGNES', all the objects are placed into individual clusters ($c_1 - c_7$), from the second step onwards, the objects are merged into larger objects based on the minimum Euclidean distance b/w any two objects from diff: clusters. This process is continues until a single cluster containing all the objects is obtained.

In DIANA, all the data objects are first placed into a single cluster ($C_{1234567}$). This single cluster is then divided into smaller clusters based on maximum Euclidean distance b/w closest neighbouring objects in the cluster. This process continues until each object is placed in an individual clu

## 8.3.2 Distance Measures in Algorithmic Methods :-

whether using an agglomerative method or a divisive method, a core need is to measure the distance b/w 2 clusters, where each cluster is generally a set of objects.

Four widely used measures for distance b/w clusters are as follows.

minimum distance : $dist_{min}(c_i, c_j) = \min\limits_{p \in c_i,\ p' \in c_j} \{|p - p'|\}$

Maximum distance : $dist_{max}(c_i, c_j) = \max\limits_{p \in c_i,\ p' \in c_j} \{|p - p'|\}$

mean distance : $dist_{mean}(c_i, c_j) = |m_i - m_j|$

Average distance : $dist_{avg}(c_i, c_j) = \dfrac{1}{n_j \cdot n_i} \sum\limits_{p \in c_i,\ p' \in c_j} |p - p'|$ .

where $|p - p'|$ is the distance b/w 2 objects or points

$m_i$ is the mean for cluster $c_i$,

$n_i$ is the no. of objects in $c_i$

they are also known as linkage measures.

when an algorithm uses the minimum distance, $d_{min}(c_i, c_j)$ to measure the distance b/w clusters, it is called as "nearest - neighbor clustering algorithm".

An agglomerative hierarchical clustering algorithm uses the minimum distance measure is also

Called a "minimal spanning tree algorithm". The minimal spaming tree is the one with the least sum of edge weights.

when an algorithm uses the maximum distance $d_{max}(c_i, c_j)$, to measure the distance b/w clusters. It is called as "farther-neighbor clustering algorithm"
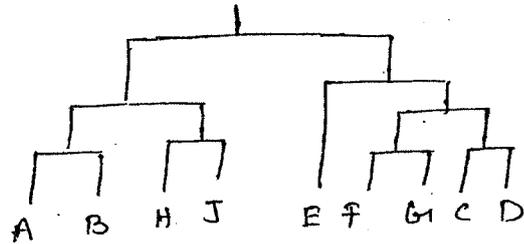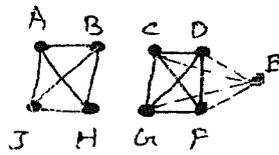
The use of mean or average distance is a compromise b/w the minimum & maximum distance and overcomes the outliers sensitivity problem.

Example :-

let us apply the hierarchical clustering to the data set { A, B, C, D, E, F, G, H, J}



clustering using single linkage

8·3

clustering using complete linkage.

### 8.3.3 BIRCH : Multiphase Hierarchical clustering using clustering Feature Trees

BIRCH (Balanced Iterative Reducing and clustering using Hierarchies) method integrates both hierarchical clustering methods, in initial stage and iterative partitioning clustering method at final stages. In other words, the o/p generated from hierarchical methods serve as I/p or preprocessing step for iterative partitioning.

Advantages of BIRCH when compared to Hierarchical Methods

1) BIRCH method improves the Scalability of clustering.

2) It has the ability to undo the mistakes which have occured in previous phases.

In BIRCH method, for an input data points, the following three factors are calculated.

1) Centroid (cu)

It determines the center of the cluster. It is

Calculated using the formula, $G_1 = \dfrac{\sum_{i=1}^{n} z_i}{n}$

2) **Radius (r)**

It determines the average distances of the I/p point from the centroid.

$$r = \sqrt{\dfrac{\sum_{i=1}^{n} (z_i - G_1)^2}{n(n-1)}}$$

3) **Diameter (D)**

It determines the average pair wise distance in a cluster.

$$D = \sqrt{\dfrac{\sum_{i,j=1}^{n} (z_i - z_j)^2}{n(n-1)}}$$

The radius and diameter factor determines the level of association of the cluster around the centroid.

**clustering feature and clustering feature Tree :-**

In BIRCH method, clustering feature (CF) is defined as a three-dimensional vector. which contains the info, about the clusters of objects that have been discarded or compressed in a summarized form.

The 'CF' for 'n' data points spread over d-dimensions can be represented as,

$$CF = (n, Lsum, Ssum)$$

where

$n \rightarrow$ no. of data points (or) objects.

$Lsum \rightarrow$ linear sum of data points,

$Ssum \rightarrow$ Square sum of data points.

clustering Feature Tree (CFT) is height-balanced tree which is used to represent clustering features of individual clusters in an hierarchical fashion. The non leaf nodes maintain the info. about the sums of the clustering features of their child nodes and there by summarizing the clustering info. of these child nodes. The size of a clustering feature tree is dependent on two factors.

1/ <u>Branching factor (B)</u>

This factor decides the maximum no. of child nodes for a non leaf node.

2/ <u>Threshold (T)</u>

Threshold decides the maximum diameter that a sub cluster i.e, a collection of non-leaf and its child nodes.

# Example of CF Tree

Node ←

Root level
(No parent for these nodes)

First level non-leaf Node ←

Leaf Node

CF₁ | CF₂ | CF₃ | CF₄ →

CF₁₁ | CF₁₂

CF₃₁

Second level ← CF₁₁₁

CF₄₁₁ →  First level Node

CF₄₁₁ | CF₄₁₂ → Second level Node.

BIRCH method is used to generate efficient clusters using those resources that are available. This method applies multiphase clustering technique which generates good clustering with minimum intervention from I/o and minimum amount of main memory for storage.

## 8.3.4  Chameleon : A Hierarchical clustering Algorithm using Dynamic Modeling :—

This method uses dynamic modeling. It basically constructs a sparse k-nearest neighbor graph, then partitions the graph into pieces and then clusters the pieces together. produces more natural clusters than DBSCAN. uses density measurements to determine the k-nearest neighbor.

Hures:-

dapt to the characteristics of the data set to find the
natural clusters.

use a dynamic model to measure the similarity
two clusters.

main property is the relative closeness and relative
connectivity of the cluster.

two clusters are combined if the resulting cluster
wares certain properties with the constituent cluster

The merging scheme preserves self-similarity

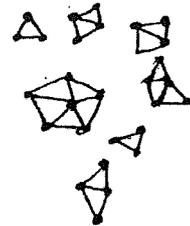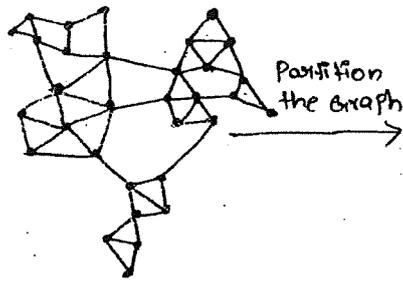one of the areas of application is spatial data.

k-nearest-neighbor graph.

level
parent for
se nodes)

1st level
node

Second level
Node.

? efficient clusters
lable. This method
ue which generates
tervention from I/o
mory for storage.

stering Algorithm

modeling. It basically
a graph. then partitions
clusters the pieces
oval clusters than DBSCAN-

determine the k-nearest

nstruct
spane
raph

merge
partition

Hierarchical clustering Based on k-nearest neighbors

mic Modeling.



Partition
the graph

Final clusters

Chameleon Steps :—

Preprocessing :—

Represent the data by a graph

1, Given a set of points, construct the $k$-nearest neighbor ($k$-NN) graph to capture the relationship b/w a point and its $k$-nearest neighbours.

2, concept of neighborhood is captured dynamically.

Phase 1 :—

use a multilevel graph partitioning algorithm on the graph to find a large no. of clusters of well-connected vertices.

Each cluster should contain mostly points from one "true" cluster i.e, is a sub-cluster of a "real" cluster.

Phase 2 :—

use hierarchical agglomerative clustering to merge Sub-clusters.

1, Two clusters are combined if the resulting cluster shares certain properties with the constituent clusters.

2, There are two key properties used to model cluster similarity.

* Relative Interconnectivity :—

Absolute interconnectivity of two clusters normalized by the internal connectivity of the clusters.

\* <u>Relative closeness :-</u>

Absolute closeness of two clusters normalited by the internal closeness of the clusters.

<u>Density-Based Methods :-</u>

Density based clustering methods have been developed to determine clusters with arbitary space. clustering based on density, such as denisity-connected points.

<u>Major Features of Density-Based Methods :-</u>

1) Discover clusters of arbitrary shape.

2) Handle noise.

3) One scan.

4) Need density parameters as termination condition.

There are three different methods of density-based clustering. they are.

1) DBSCAN
2) OPTICS
3) DENCLUE

<u>8.4.1 DBSCAN :-</u>

Density-Based Spatial Clustering of Application with Noise (DBSCAN) is a density-based clustering algorithm. In DBSCAN a cluster is a set of maximum density

connected points DBSCAN finds arbitrary-shaped clustering alg. and finds arbitrary-shaped clusters in spatial DBS with noise and expands regions having enough high density into clusters. To better understand the concept of DBSCAN method we first need to know some definitions.

## Noise

Noise is defined as the set of objects in N which do not belong to any cluster.

### 1, $\epsilon$-neighborhood of an objects:-

The $\epsilon$-neighborhood of an object is defined as for a given positive radius, the neighborhood of an object must be within the radius($\epsilon$).

### 2, Core object:-

An object is said to be a core object, if the $\epsilon$-neighborhood of an object contains at least a minimum no. of objects (Minpnts).

### 3, Directly-Density Reachable:-

An object 'e' belonging to a set of objects, N is directly density-reachable from an object 'f', if 'f' is a core object and e is within $\epsilon$-neighborhood of f.

4) **Density - Reachable :-**

Given a neighborhood $\varepsilon$ and minipnts in a group of objects $N$ and when there is a series of objects $e_1, e_2, \ldots e_n$ where $e_1 = f$ and $e_n = e$ $\Rightarrow$ $e_{a+1}$ is directly density reachable from $e_n$ with respect to $\varepsilon$ and minipnts, for $1 \le a \le n$, $e_a \in N$ then an object $e$ is density reachable from object $f$.

5) **Density - Connected :-**
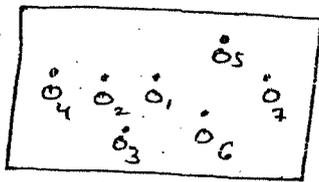
Given $\varepsilon$ and minipnts in a group of objects $N$ and when there is an object 'g' belonging to $N$ $\Rightarrow$ both $e$ & $f$ are density - reachable from 'g' with respect to $\varepsilon$ and minipnts. then an object 'e' is density - connected to object $f$.

6) **Density clusters:-**

A cluster with respect to $\varepsilon$ and minipnts is a non-empty subset of $N$ containing maximum no. of density connected objects with respect to density - reachability.

**Example :—**

let us assume that minipnts = 5 and $\varepsilon = 2 cms$.



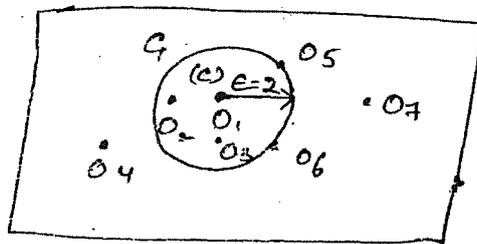objects before clustering.

**Steps for finding clusters :-**

1) Select an unclassified object $O_1$. The objects that are in the neighborhood of $O_1$ along with the radius $\varepsilon$ are,

$$N_E(O_1) = \{O_1, O_2, O_3, O_5, O_6\}$$

since, $N_E(O_1) \geq$ minpnts i.e, $5 \geq 5$

$O_1$ is considered as a core object $(c)$

Thus, these set of points together form a cluster, which is assigned a cluster-id $(C_1)$. This is shown in below.



$O_1$ Neighbouring objects

2, Now, select the object $O_2$ from the set. The neighbouring objects of $O_2$ are

$$N_E(O_2) = \{O_1, O_2, O_3, O_4\}$$

Since $N_E(O_1) <$ minpts i.e, $4 < 5$, the object $O_2$ is considered as a 'Noise' object $(N)$ and therefore no cluster is formed.

3, select the object $O_3, O_4$, and $O_5$ from set. The neighbouring object of $O_3, O_4, O_5$ are

$$N_E(O_3) = \{O_1, O_2, O_3, O_6\}, \quad N_E(O_4) = \{O_2, O_4\}$$

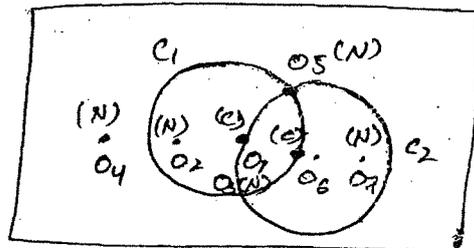$$N_E(O_5) = \{O_1, O_5, O_6, O_7\}$$

Since $N_E(O_3) <$ minpnts, $N_E(O_4) <$ minpnts, $N_E(O_5) <$ minpnts

so, these are consider as Nosie objects $(N)$. and therefore no cluster is formed.

**e)** Finally, select the object $O_6$ from the set. The neighboring objects of $O_6$, are

$$N_e(O_6) = \{O_1, O_3, O_5, O_6, O_7\}$$

Since, $N_e(O_6) \geq minpnts$ i.e, $5 \geq 5$ the object $O_6$ is considered as core object. thus, these set of points form a cluster, & assigned a cluster-Id ($C_2$).



$O_1$ & $O_6$ objects with their neighbors forming clusters $C_1$ & $C_2$.

The resultant DBSCAN clustering in which there are two clusters as,

$$C_1 = \{O_1, O_2, O_3, O_5, O_6\} \text{ &}$$
$$C_2 = \{O_1, O_3, O_5, O_6, O_7\}$$

where $O_1, O_6$ are core objects and $O_2, O_3, O_4, O_5, O_7$ are noise objects.

## 8.4.2 OPTICS :—

### Algorithm :—

DBSCAN : a density-based clustering algorithm.

**Input :**
- D : a data set containing n objects,
- $\epsilon$ : the radius parameter &
- MinPts : the neighborhood density threshold.

Output:- A set of density-based clusters;

Method:-

1. Mark all objects as unvisited;
do
    randomly select an unvisited object $p$;
    mark $p$ as visited;
    if the $\epsilon$-neighborhood of $p$ has at least minpnts objects
      Create a new cluster C, and add $p$ to c;
      let N be the set of objects in the $\epsilon$-neighborhood
                          of $P$;
      for each point $p'$ in N
        if $p'$ is unvisited
          mark $p'$ as visited;
          if the $\epsilon$-neighborhood of $p'$ has at least
                     minpnts points
         add those points to N;
        if $p'$ is not yet a member of any cluster,
              add $p'$ to c;
    end for
    Output c;
  else mark $p$ as noise;
until no object is unvisited;

Classification according to the
Database mined.
Clustr accos to the
Knowle mined.
class accord to the
techniques utilized.
class accor to the
applications adopted.