

* Data Warehousing and Data Mining:

* Data Warehousing: It is a coherent collection of Subject oriented, integrated, time variant and non-volatile collection of data

Subject oriented :- stores data related to particular subject

e.g:- Sales, products

Integrated :- It combines multiple homogeneous and heterogeneous data sources.

Time Variant :- collects data related to particular period

e.g:- Sales

Non-volatile :- If you want to add some data then previous data is not changed

* Data Mining:

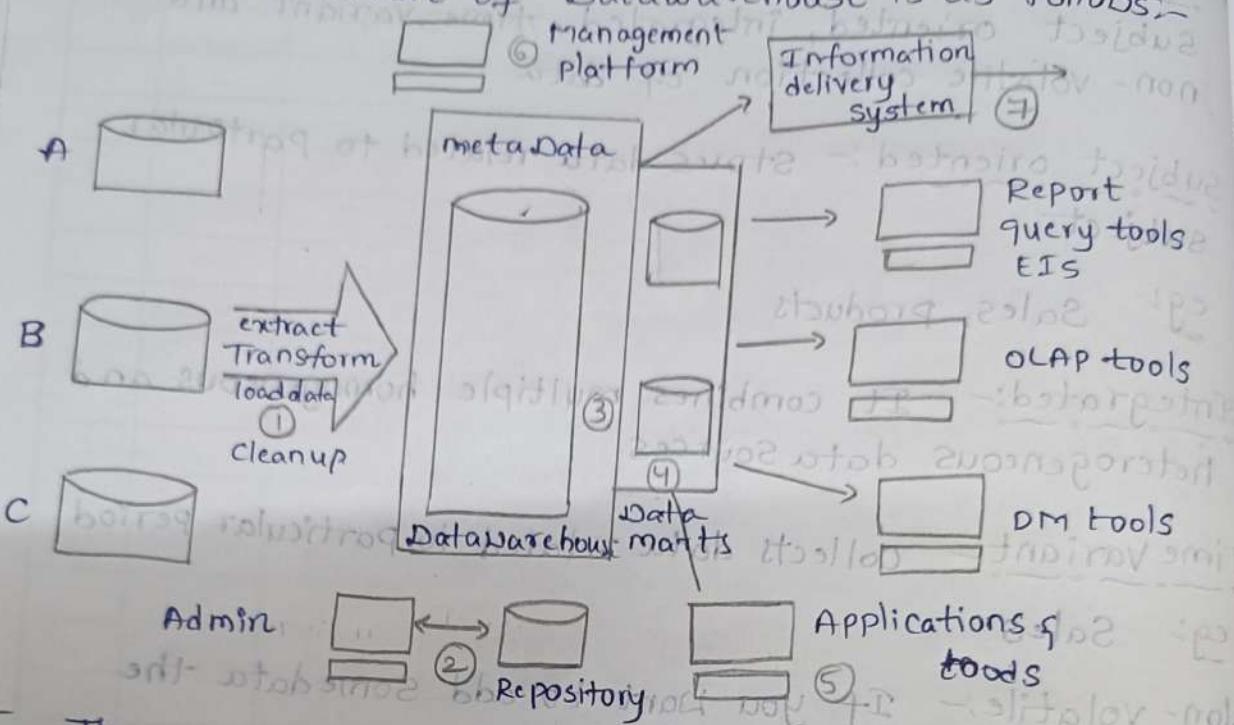
It is the process of extracting information from huge volumes of data in datawarehouse is called datamining. apply some datamining techniques to extract information from datawarehouse

* Advantages of DWDML

i) Taking better business decisions

* Architecture of Datawarehouse and its components

The Architecture of Datawarehouse is as follows:



* The source for Datawarehouse is operational databases (or) files.

- 2) The data entered into the Datawarehouse is transformed into an integrated structure and formats
- 1) Datawarehouse Database: - This is the central part of data warehousing environment. This is implemented based on RDBMS technology.
- 2) Sourcing, acquiring, cleanup tools:
 - 1) To remove unwanted data from operational database
 - 2) converting to common datatypes & attributes
 - 3) calculating summaries & derived data

4) Establishing defaults for missing data

3) meta data:- It is data about data. It is used for maintaining and managing and using the data warehouse.

4) Access tools:-

Its purpose is to provide information to the business users for decision making.

1) Dataquery and reporting tools

2) Application development tools

3) Executive information system tools

4) OLAP tools

5) Datamining tools

Dat

} tools for
accessing
data

5) DataMarts

Departmental subsets that focuses on selected subjects.

6) Datawarehouse Admin and Management

i, Security and priority management

ii, Monitoring updates from multiple sources

iii, Data quality checks

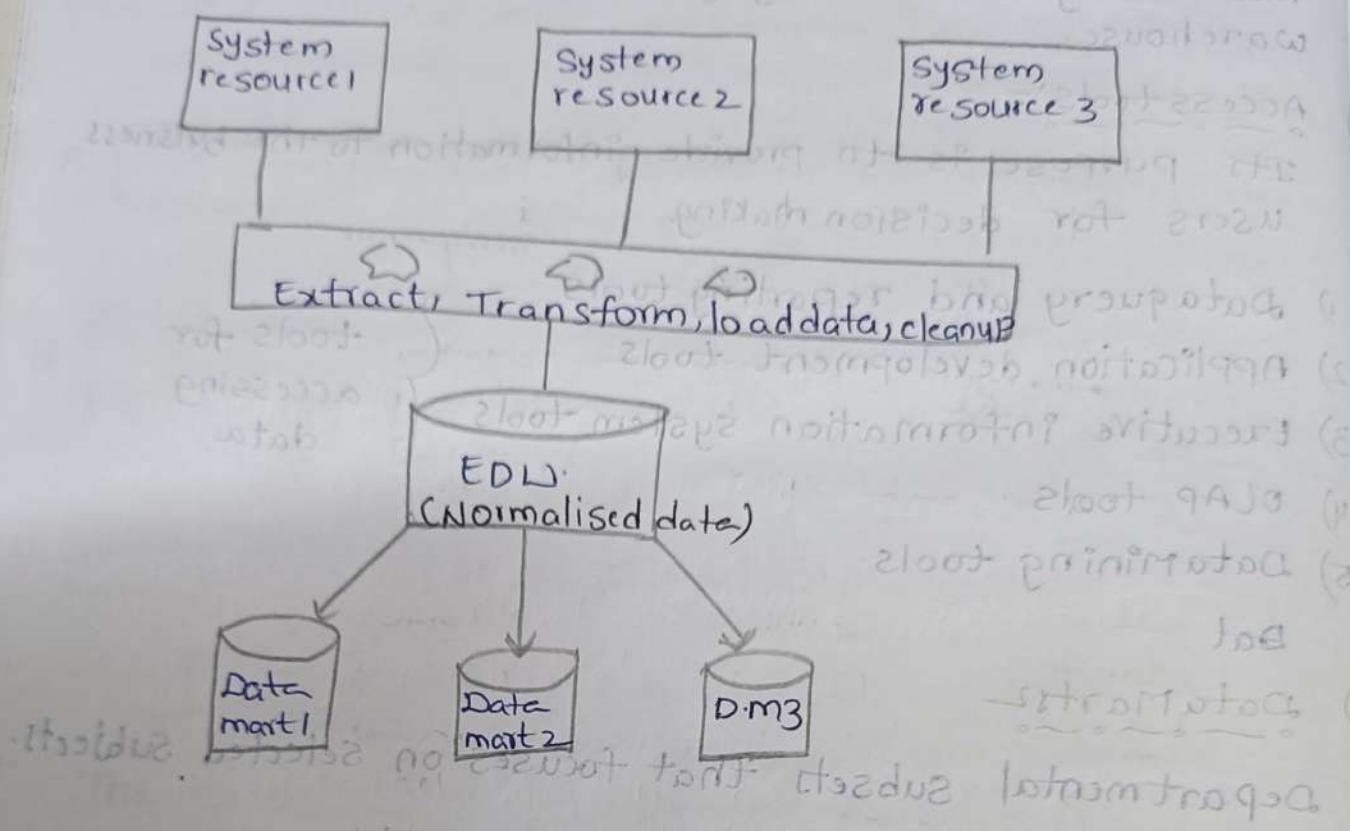
iv, Managing & updating metadata

7) Information delivery system

The Information delivery system delivers the information to the data warehouse user.

realme Shot on realme 8s 5G

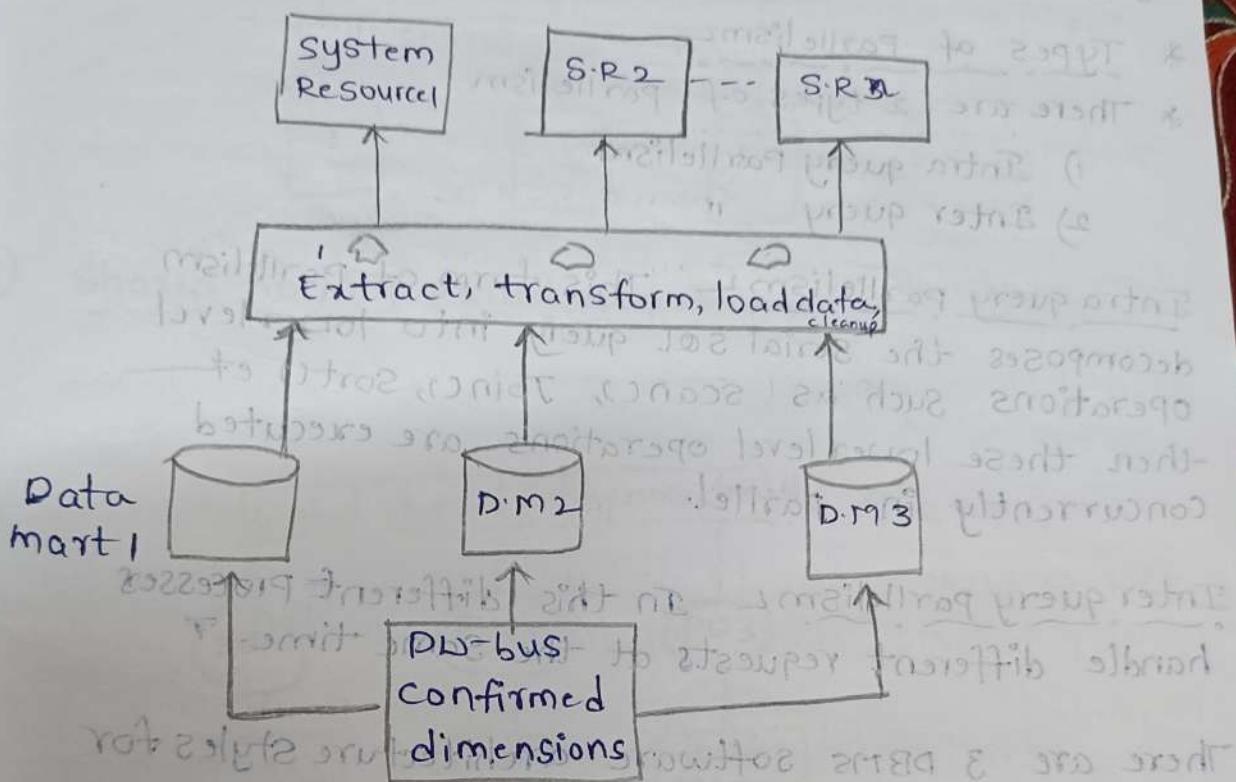
- * Building a datawarehouse:-
- * Top down approach:-



- * There are 2 reasons why organisations consider data warehousing at critical need.
 - Business factors:
Business users want to take decisions quickly & correctly using all available data.
 - Technical factors:
 - TO address the incompatibility of operational data stores
 - IT Infrastructure is changing rapidly. Its capacity is increasing & cost is decreasing. So, building a datawarehouse is easy.

* There are 2 ways to design a data warehouse

- 1) Top down approach
 - 2) Bottom up approach
- 2) Bottom up approach



* Database Architecture for parallel processing

- 1) The functions of datawarehouse are based on the relational database technology. The RDBMS technology is implemented in parallel manner.
- 2) There are 2 advantages of having parallel Relational database technology for data warehouse:-
 - 1) Linear - speed up
 - 2) Linear - scale up

Linear-Speed up It refers the ability to increase the number of processors to reduce the response time

2) Shar

Linear-Scale up It refers the ability to provide same performance on the same requests as the database size increases.

* Types of Parallelism

* There are 2 types of parallelism

1) Intra query Parallelism

2) Inter query "

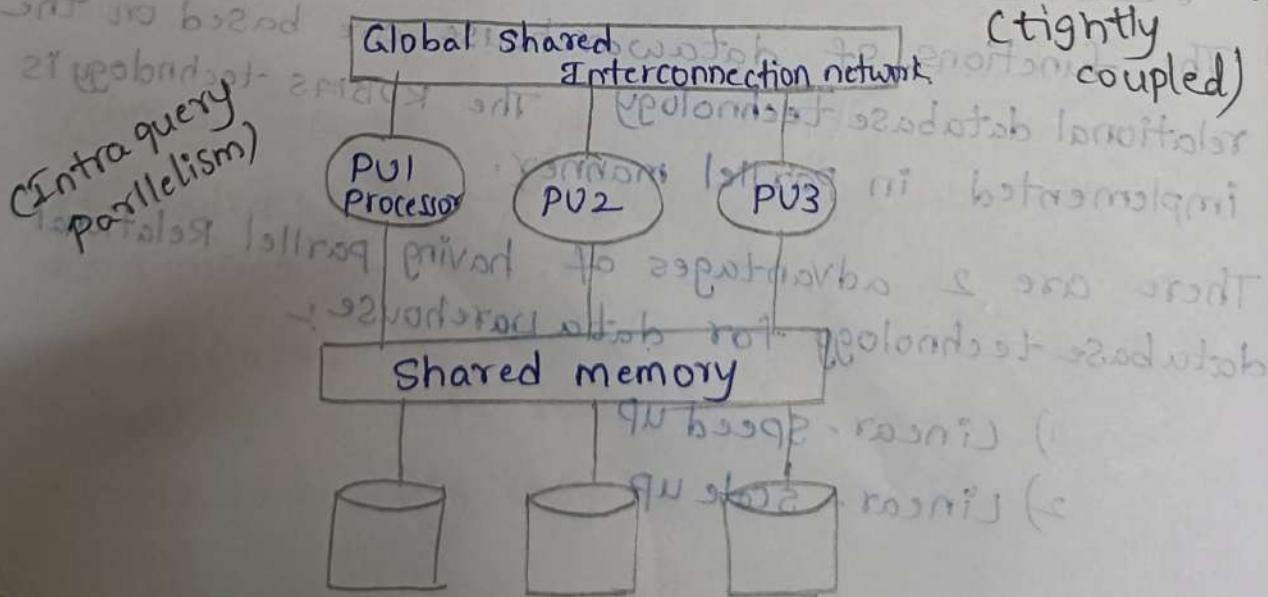
Intra query Parallelism This form of Parallelism decomposes the serial SQL query into lower level operations such as scan(), Join(), Sort() etc then these lower level operations are executed concurrently in parallel.

3)

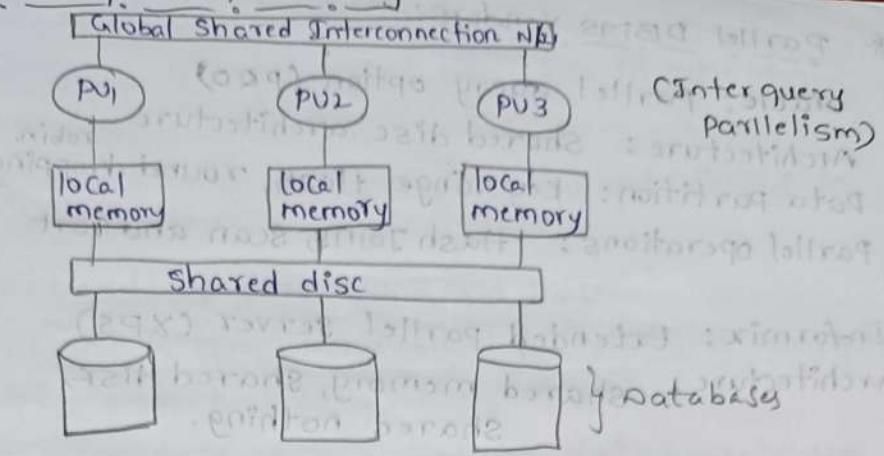
Inter query parallelism In this different processes handle different requests at the same time.

* There are 3 DBMS software architecture styles for Parallel Processing

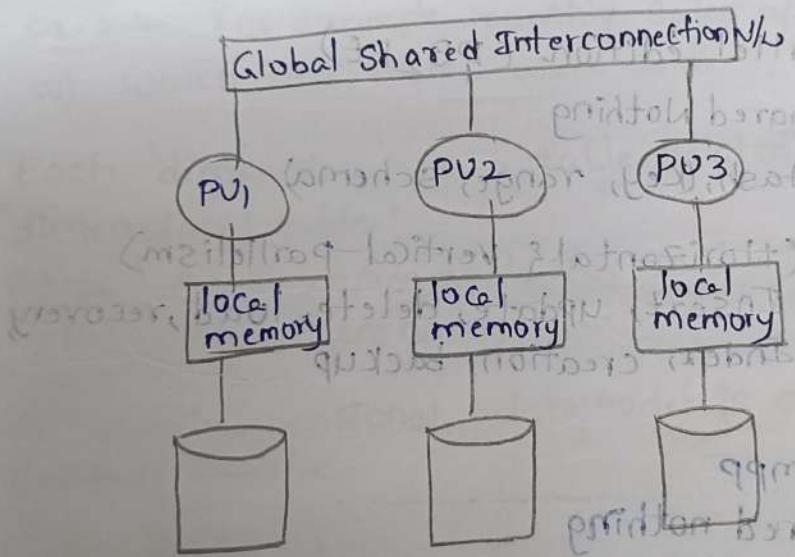
1) Shared memory / shared everything parallel processing



2) Shared disk parallel processing



3) Shared nothing Parallel processing



* Parallel DBMS Vendors:

1) oracle: Parallel query option (PQO)

Architecture: Shared disc architecture

Data partition: key range, hash, round robin dropping

Parallel operations: Hash Joins, Scan and Sort

2) Informix: Extended parallel server (XPS)

Architecture: shared memory, shared disk, shared nothing.

Data partition: round robin, hash, schema, key ranges, User defined.

Parallel operations: Insert, update, delete

3) IBM: DB 2 Parallel edition (DB₂ PE)

Architecture: Shared Nothing

Data partition: Hash, (key, range, schema)

Parallel

operations: (Horizontal & vertical parallelism)

Insert, update, delete, load, recovery, Index, creation, backup

4) Sybase: Sybase mpp

Architecture: shared nothing

Data partition: Hash, key, range, schema

Parallel operations : Horizontal & Vertical parallelism.

* Multidimensional datamodel

- 1) A multidimensional datamodel views data in the form of datacube.
- 2) A datacube enables data to be model and viewed in multiple dimensions.
- 3) It is defined by dimensions and Facts
- 4) The dimensions are perspectives (or) Entities concentrating which an organisation keeps records
- 5) ex:- A shop may create a sales datawarehouse to keep records of stores' sales for the dimension time, item and location. These dimensions allow the shop to keep track of things.

ex-2 For example monthly sales of items and locations at which the items were sold.

- 5) Each dimension has a table related to it called a "dimensional table".
ex-2 A dimension table for an item may contain the attributes item name, brand and type.
- 6) A Multidimensional datamodel is organised around a central theme.
ex-2 Sales
- 7) This theme is represented by a "Fact table". Facts are numerical measures. The fact table contains the names of facts (or) measures of related tables.

	Location	Item	Hyd	Chen	Can
Time	Q1	1700	305	305	200
	Q2		2000		50
	Q3				70
	Q4				90

8) Multidimensional datamodel is based on OLAP tool

* Database Schemas for decision support

- 1) Star Schema
- 2) Snowflake schema
- 3) Fact constellation

* Concept hierarchy

- The mapping of low level form to higher level (or)
Generalised form is called 'concept hierarchy.'
- It defines the sequence of mappings from low-level concepts to higher-level (or) Generalised concepts

* Types:-

- 1) Schema hierarchy → cuboid (collection of nodes)
- 2) Set Grouping hierarchy

Base value, Target Value has to be identified
then high level data converted to low level format

* Characteristics of OLAP systems:

Transactions/ which are performed on huge volumes of Analysis data is OLAP.

OLTP → Performs on database

OLAP → Performs on data warehouse

OLAP → complex, OLTP → Simple tasks & easy
Big tasks

→ The analysis which is performed on previous data → OLAP

* Characteristics:

i) In FASMI characteristics of OLAP methods - the term derived from first letters of the characteristics

F - Fast

A - Analysis

S - Shared

M - Multidimensional

I - Information

Fast:- It defines which the system targeted to deliver the most feedback to the client within about 5 seconds. Taking no more than 1 sec & very few taking more than 20 sec.

Analysis:-

	<u>OLTP</u>	<u>OLAP</u>
1. characteristics	Handles a large no. of small transactions	Handles large volumes of data with complex queries
2. Query types	Simple, standardised queries	Complex queries
3. operations	Based on insert, update, delete commands	Based on select command to aggregate data for reporting
4. Response time	milliseconds	Based on the query type, seconds, minutes (or) hours depending on amount of data to process
5. Design	Industry-specific such as retail, manufacturing (or) Banking	Subject-specific such as sales, inventory (or) marketing
6. Source	Transactions	aggregated data from Transactions.
7. purpose	Controls run essential business operations in real-time	plan, solve problem, support decisions, discover hidden insights
8. Space requirements	Generally, small until historical data is achieved	Generally large, due to aggregating large data sets.

9. Backups Recovery	Regular backups required to achieve Business continuity & meet legal & Governance requirements	Lost data can be reloaded from OLTP database as needed in view of regular backups
10. productivity.	OLTP increases the productivity of end users	Increases productivity of business managers, Data analysts, executives
11. dataview	day-day business transactions	Multidimensional view of enterprise data
12. database design	Normalised databases for efficiency	Denormalised databases for analysis

* Typical operations of OLAP Systems:-

- 1) OLAP stands for online analytical processing server
- 2) It is a software technology that allows users to analyse information from multiple database systems at the same time
- 3) It is based on multidimensional data model and allows the user to query on multidimensional data
- 4) OLAP databases are divided into one or more cubes and these cubes are known as "hypercubes"
- 5) There are 5 basic analytical operations that can be performed on OLAP cube:-

3) Drill down and roll up

From high level to low level

To add more detail information

In this operation, the less detailed data is converted into highly detailed data.

Q1	1000	Loc	B	P
Q2	200	Loc	A	M
Q3	100	Loc	K	M
Q4	50	Loc	J	T

(No. of dimensions increased)

less detailed data
high level
dimensions to be
increased to high level

Loc	B	P
Q1	500	200
Q2	200	100
Q3	100	50

2) Roll up It is just opposite of the drilldown operation. It performs aggregation on OLAP cube. It can be done by climbing up in the concept hierarchy, reducing the dimensions.

Q1	1000	Loc	B	P
Q2	200	Loc	A	M
Q3	100	Loc	K	M
Q4	50	Loc	J	T

low level data

Loc	B	P
Q1	500	200
Q2	200	100
Q3	100	50

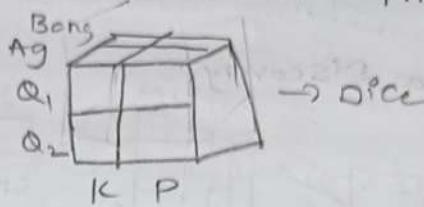
high level data

3) Dice It selects a subcube from the OLAP cube by selecting 2 (or) more dimensions. In the cube given in the Overview section a subcube is selected by selecting following dimensions with criteria.

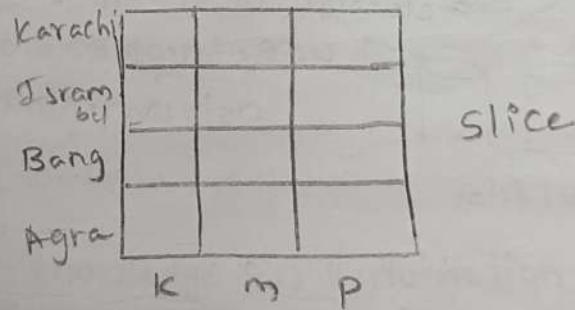
Location = "Agra" (or) "Bangalore"

Time = " Q_1 " / " Q_2 "

Item = "Keyboard" (or) "Printer".

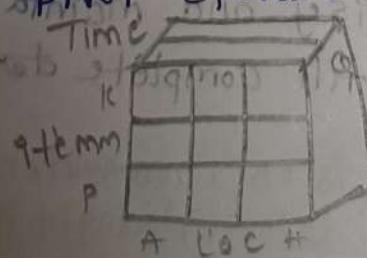


4) Slice Selects a single dimension from the OLAP cube which results a new subcube creation. In the cube given in the overview section, slice is performed on the "dimension time = Q_1 ".



5) Pivot It is also known as rotation operation as it rotates the current view to get a new view of the representation.

* In the subcube obtain after the slice operation, performing pivot operation gives a new cube of it.



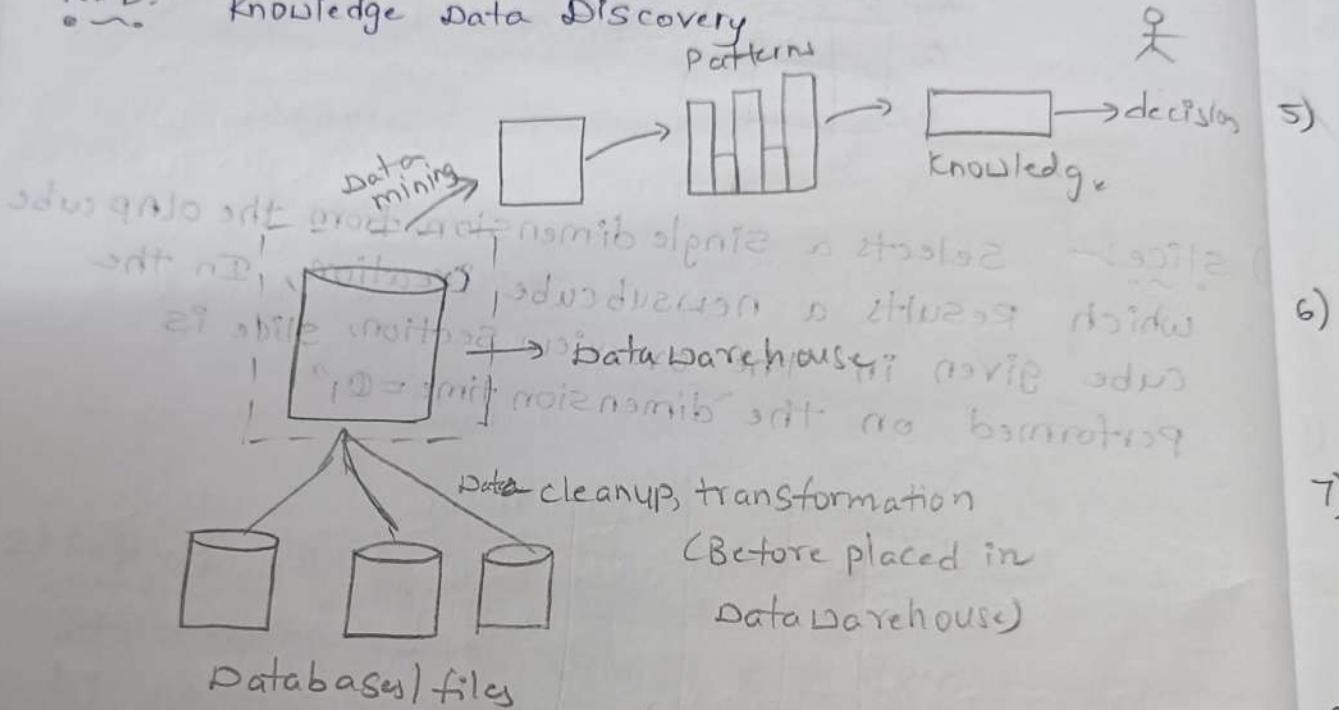
2. Introduction to Data Mining

- 1) Data mining is a process used by organisations to extract specific data from huge databases to solve business problems.
- 2) The main aim of data mining turns raw data into useful information.

KDD:-

knowledge data discovery

patterns



1) Data Selection

Select, task relevant data from different operational databases (or) flat files

2) Data cleaning

In order to maintain quality data in the warehouse for the purpose of extracting correct information & then by removing noisy and maintaining data consistency and to fill complete details data.

(Remove noise, inconsistent, incompleteness in the data before placing it in the warehouse)

To noise, inconsistency, incomplete

3) Data transformation

To convert the selected data to required format related to our task.

4) Data warehouse

The collection of raw data from different heterogeneous (or) homogeneous databases (or) flat files related to Subject

5) Datamining:-

Extraction of Information (or) knowledge from huge volumes of Raw data

6) Patterns:-

The extracted information is reported to the datawarehouse user by using interesting patterns like Graphs, charts etc

7) Knowledge (or) Information

After processing of raw data what we can get is called 'knowledge'

8) Decision

Based on knowledge (or) Information the datawarehouse user/ decision maker to take correct and quality decision regarding about his organisation.

* Data mining:-

- * Datamining is the process by which organisation detects patterns in data for insights relevant to their Business needs.
- * It is essential for Both Business intelligence in data science
- * There are many datamining techniques:-
 - 1) Association rules (or) mining
 - 2) classification
 - 3) Prediction
 - 4) Clustering
 - 5) Regression
 - 6) Artificial neutral network
 - 7) Outlayer detection
 - 8) Genetic algorithm
- * Applications of Datamining:-
 - 1) The Amount of data collected is said to be almost double of every year.
 - 2) An extracting data (or) seeking knowledge from massive data, datamining techniques are used.
 - 3) In almost all datamining is used in almost all places where a large amount of data is stored & processed.
 - 1) Scientific Analysis
 - 2) Inclusion detection
 - 3) Business transaction
 - 4) Market Basket Analysis
 - 5) Education
 - 6) Research

- 7) Health care & Insurance
- 8) Transportation

* Issues (or) problems in datamining:-

1) Datamining is not a easy task as the algorithms used is very complex & it is not available at 1 place
2) It needs to be integrated from various heterogeneous data sources.

3) These factors also create some issues.

4) The major issues are:-

i, Mining methodology & user interaction

ii, Performance

(iii) Diverse data drivers issue

1) Mining methodology & user interaction issues:-

i, Mining different kinds of knowledge from different databases

ii, Interactive mining of knowledge at multiple levels of abstraction

iii, Incorporation of background knowledge

iv, datamining query languages ad hoc

v, presentation & visualisation of datamining results

vi, Handling noisy & incomplete data

vii, Pattern Evaluation

2) Performance issues:-

Efficiency & scalability of algorithms can be used
Parallel, distributed & incremental mining algorithms

3) Diverse datatypes issue

Handling of relational & complex types of data
Mining information from heterogenous databases and
global information system

* Data objects & Attribute types

1) Object → Real time thing

2) Attribute Property/ characteristics of an object

ex:- Attributes of student :- sid, Sname, phno, add

* Types of Attributes

1) Nominal attribute Represents the names of the objects

It also consider as a categorical attribute

ex:- Attribute colour values

color

Blue, red, green

Fruit

Banana, mango

2) Binary attribute always represents 2 values

Yes / affected / unaffected (or) True / false.

* These are further divided into

1) Symetric type :-

2) unsymetric type :-

Symetric Equal importance

Attribute values

Gender Male / female

Sametype

- 2) Asymmetric Result pass, fail
- 3) ordinal attributes Represents the attributes in a particular order. i.e either character values/numeric values etc.

eg :- Attribute Value

Gender

Grade

- 4) Numerical attributes Represents the Numerical value in finite no. of values.

eg :- Attribute Value

Age

12, 20, 40

1) Ratio scaled : multiply/divide

2) Interval scaled (equal interval)

e.g calendar

- 5) discrete attributes stores the infinite no. of values.

eg :- Attribute Value

Weight 60 40 35

Height 5.4 5.6 5.2

* Statistical Description of Information

Numerical attributes These are quantitative which can be measurable represented as integer (or) real values

- There 2 types of numeric attributes

1) Interval scaled attribute

2) Ratio scaled attribute

Interval scaled These attributes will have values which follows some order. We can add/subtract to get next value in the sequence

ex :- Temp, calendar

realme (1-31) Shot on realme 8s 5G

2) Ratio Scaled These attributes will have the values. We can multiply the previous value to get next value in a sequence.

e.g. 10k is twice greater than 5k

5k 10k 20k 40k

3) Discrete attribute: These attributes have finite or countably infinite set of values.

e.g.

Attribute	Values
ZIP code	522403, 522017
Profession	Teacher, doctor

4) Continuous attribute: These attributes have infinite no. of values. These values are float type.

e.g.: Values b/w 2 & 3 are infinite

Attribute	Values
Height	5.4, 5.7, 5.8

* Statistical description of Data

The statistical functions are useful to identify the properties of data and also used to find noisy, outliers.

1) Measuring central tendency:

Mean

Median

Mode

realme

Shot on realme 8s 5G

2022.09.06 22:09

1) Mean Let us consider the marks of 6 students in a particular subject as $x_1, x_2, x_3, x_4, x_5, x_6$ the value will be resp 6, 7, 4, 3, 4, 7 respectively. Then, the mean

$$\text{mean}(\bar{x}) = \frac{\sum_{i=1}^n x_i}{N} = \frac{\text{Sum of values}}{\text{no. of values}}$$

$$\text{Here } N = 6 = \frac{31.0}{6} = 5.16$$

2) Median It is the middle value among all ordered values

ex: To calculate median the values should be in arranged order (ascending order)

Consider the above example

~~6 7 4 3 4 7~~

unsorted order 3 4 4 6 7 7 9

$$\text{Odd } \frac{N}{2} = \frac{7}{2} = 3.5 = 4$$

NE $\frac{N+1}{2}$

even: 3 4 4 6 7 9

$$\text{Median} = \frac{\frac{N}{2} + \left(\frac{N}{2} + 1\right)}{2} = \frac{\left(\frac{6}{2} + \left(\frac{6}{2} + 1\right)\right)}{2} = \frac{3(4)}{2} = \frac{7}{2} = 3.5$$

$$\frac{4+6}{2} = \frac{10}{2} = 5$$

Note i) for odd no. of values

$$\text{Median} = \frac{N}{2} \text{ observation of median}$$

ii) For even no. of values

$$\text{Median} = \frac{\frac{N}{2} + \left(\frac{N}{2} + 1\right)}{2} \text{ observation will be median}$$

3) Mode :- The repeated order in a sequence is called "mode".

1) Unimodal :- No value is Repeated

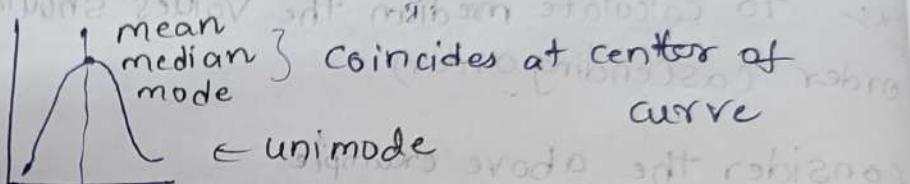
2) bimodal :- 1 value is Repeated

3) Trimodal :- 2 values are Repeated

Ex :- 3 4 4 6 7 7

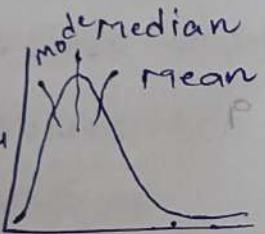
mode = 4, 7 (Trimodal)

Unimodal :- In case of symmetric data



Bimodal :-

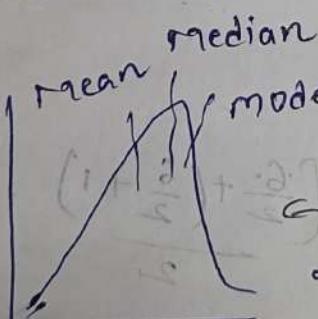
Data distribu
at left side



means mode

→ Positively skewed data /
asymmetric data

* Trimode :-



means mode

→ Negatively skewed data /
D.P. at right side

* Data in Real time applications is not in symmetric format. It may be either positively skewed data or negatively skewed data.

* Measuring dispersion of data

1) Ranger The average of min & max value

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

2) Quantile It represents half of the Data distribution

3) Quartile:- It represents $\frac{1}{4}$ th of the Data distribution

4) Percentile:- It represents 100th part of the Data distribution

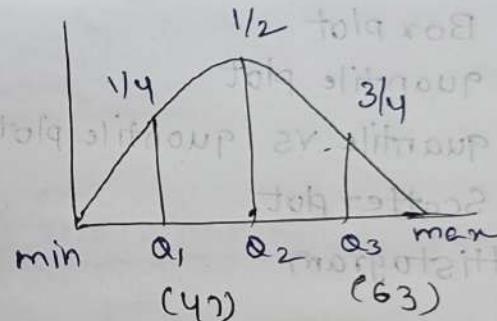
5) Interquartile range (IQR):- It is the difference b/w

$\frac{3}{4}$ th & $\frac{1}{4}$ th of the Data distribution i.e.

$$IQR = Q_3 - Q_1$$

$$\therefore IQR = 63 - 47$$

$$= 16$$



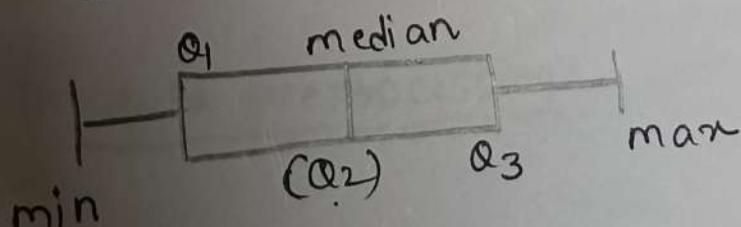
6) 5 number summary

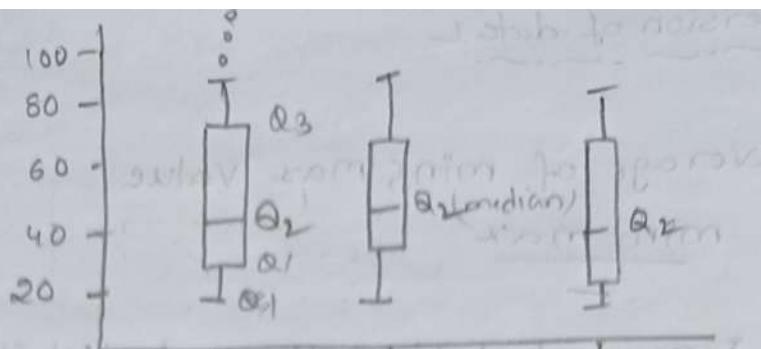
It consists of min, Q_1 , median(Q_2), Q_3 , max.

The outlier can be identified by single out value which falls either below $1.5 \times IQR$ (or) the above

$$Q_1 < 1.5 \times IQR > Q_3 \rightarrow \text{outliers}$$

Box-plot:- The 5 number summary is represented in the following Boxplot

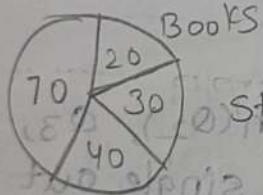




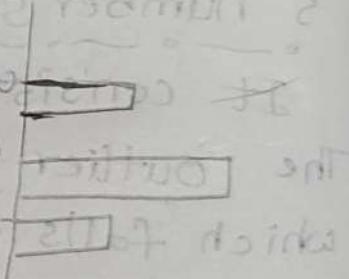
* Display of Information

- 1) Pie chart
- 2) Bar chart
- 3) Box plot
- 4) quantile plot
- 5) quantile vs quantile plot
- 6) Scatter plot.
- 7) Histogram

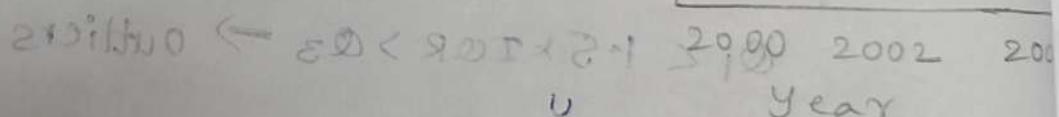
Piechart



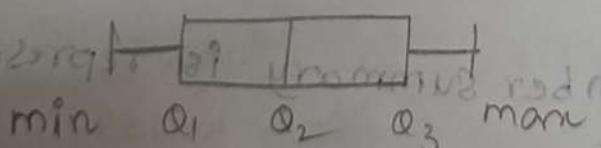
Bar chart



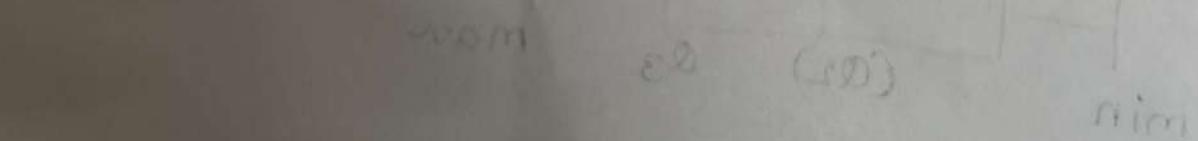
Boxplot



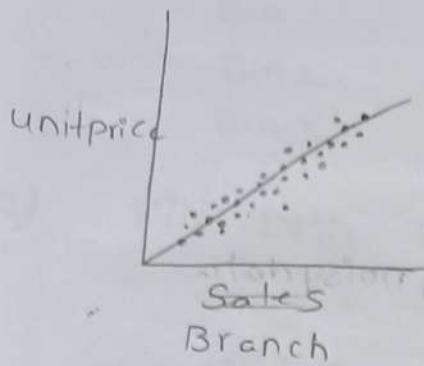
Boxplot for Year wise sales



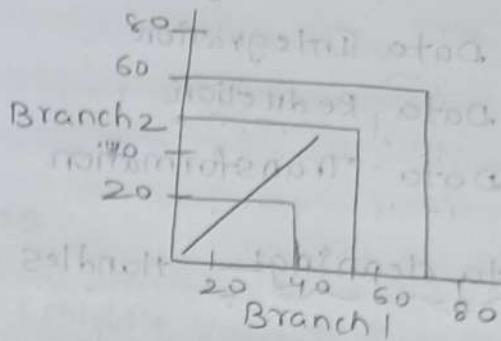
Boxplot for Year wise sales



4) quantile plot

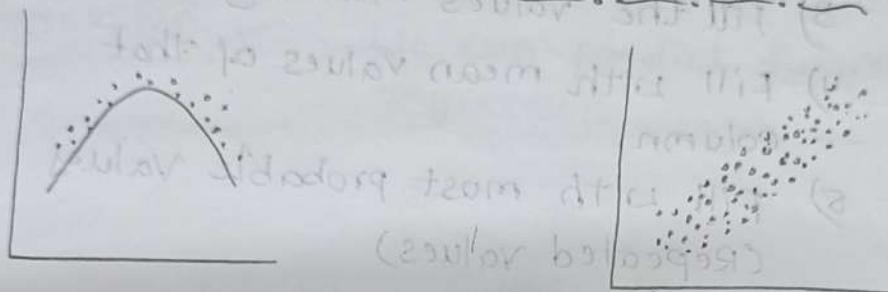


5) quantile vs quantile plot

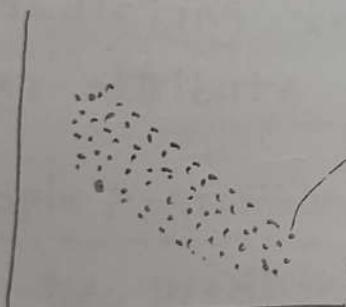


6) scatterplot

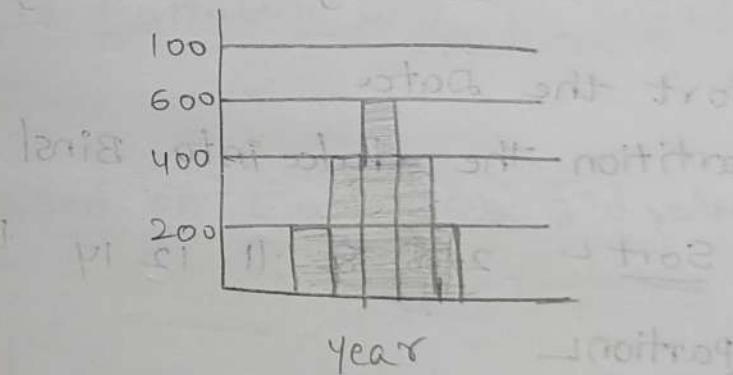
positive scatterplot



negative scatterplot



7) histogram



* Data preprocessing :-
Converting raw data into useful information is preprocessing

* steps in data preprocessing

- 1) Data cleaning
- 2) Data Integration
- 3) Data Reduction
- 4) Data Transformation

i) Data cleaning Handles missing values, noisy data

ii) Missing values
i) Ignore the tuple

ii) use Global constant (Unknown, -∞)

iii) Fill the values manually

iv) Fill with mean values of that column.

v) Fill with most probable values
(Repeated values)

ST

ii) Handling noisy data (or) errors

Data Smoothing (or) Binning

- 1) Sort the data
- 2) Partition the data into Bins/ Buckets
- 3) Sort 2 5 8 11 12 14 16 19 21

Partition

Bin 1 :- 2 5 8

Bin 2 :- 11 12 14

Bin 3 :- 16 19 21

} if any error

is detected

If any error is detected in bins then do the following:
occurred

1) Fill with the mean values

Bin 1

Bin 2

Bin 3

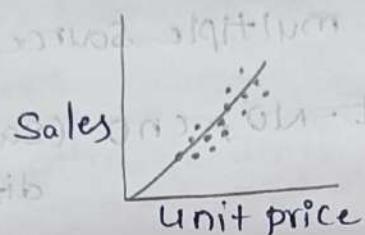
2) Fill with Boundary values

Bin 1 : 2, 2, 8 (middle value is obtained based

Bin 2 : 6, 11, 14 (on which boundary is it

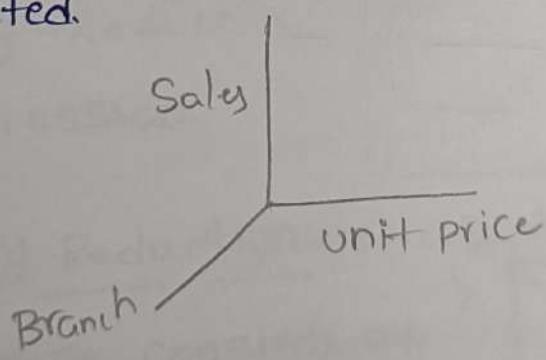
Bin 3 : 16, 21, 21 nearer to)

3) Regression :- We can predict future value based on regression



* Predicting the value of 1 attribute with the help of other attribute

Multiple Regression Based on 2 attributes 3rd value can be predicted.



* Outlier Analysis

1) clustering :- Grouping of objects based on similar characteristics.

With the help of clustering we can identify the outliers

2) Data Integration:-

It is the process of collecting data from multiple data sources.

Issues:-

- 1) Schema Integration & object matching
- 2) Redundancy & correlation Analysis
- 3) Data value conflict and Resolution

1) Schema Integration

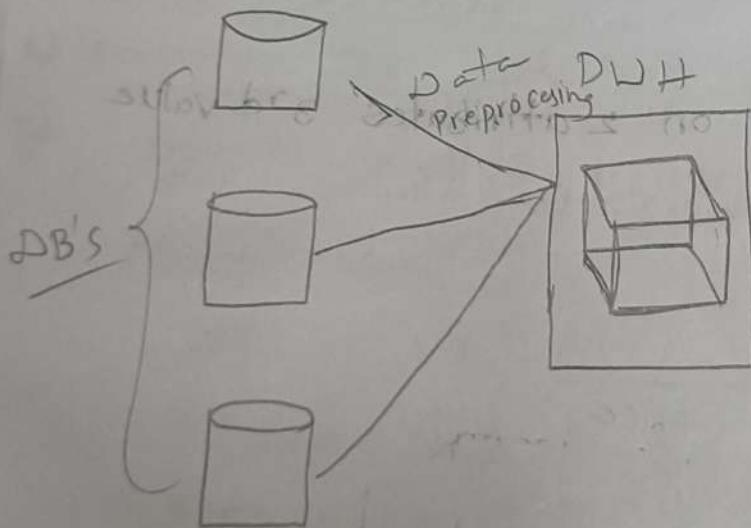
It occurs due to mismatch in attribute names which are collected from multiple sources

ex:- cust-ID, cust-NO, cno (same data but with diff names)

2) Redundancy:-

Same type of data is available.

Redundancy also occurs due to derived attribute



ex Name, DOB, Age,

DM knowledge

Object matching

- e.g. Type of currency $\rightarrow \$, \text{Rs}, \text{pounds}$
- * When merging is done either \$ changes to Rs or vice versa. Mismatching of data occurs in case of object matching

Correlation analysis

describing Relation b/w normal attribute and the derived attribute

$$\text{ex: } \text{Sname} \rightarrow \text{std}$$

$$\text{Dname} \rightarrow \text{Dno}$$

Sid depends on Sname

Dno " Dname

[Relation b/w 2 attributes]

Data value conflict & resolution

- converting different units into single unit

diff attribute units

circumstances CBSE ICSE

Data Reduction

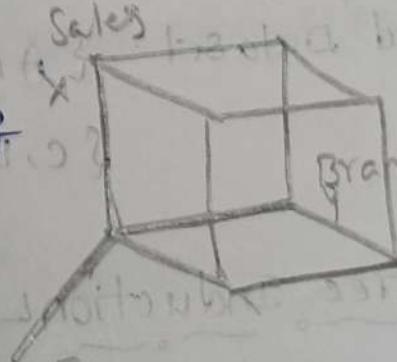
Reducing the size of data before loading into DW

Steps in Data Reduction

- Dimensionality Reduction
- Numerosity Reduction
- Data compression

Dimensionality Reduction

- Data warehouse consists of Datacubes



- Reducing the Irrelevant data in DW (or)

Removing the unused attributes from database

Steps in Dimensionality Reduction

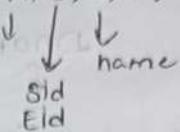
i) Stepwise forward selection

ii) Stepwise Backward elimination

iii) Decision tree Induction

i) Stepwise forward selection

* Original dataset: {A, B, C, D, E, F}



ii) Reduced dataset:

{}

{B}

{B, C}

{B, C, E}

ii) Stepwise Backward elimination

i) Opposite to the forward selection

ii) Identify most irrelevant attributes from dataset. Remove them

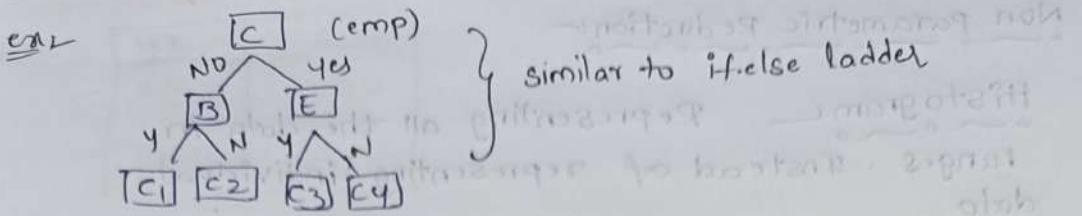
Original dataset: {A, B, C, D, E}

Reduced dataset: {A, B, C, D}

{C, B, D}

iii) Decision tree Induction

Represents in tree format
most relevant attribute becomes root node, 2nd, 3rd relevant attributes become children to root node.



2) Numerosity Reduction

- 1) Parametric Reduction
- 2) Non parametric Reduction

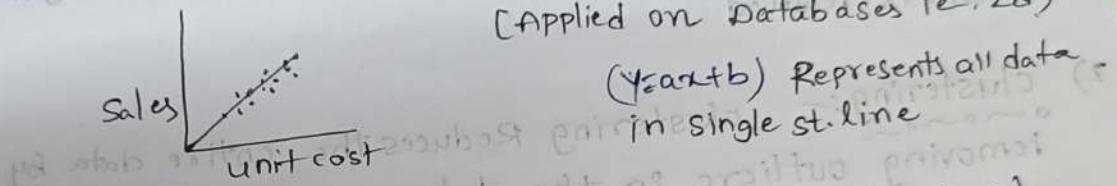
Linear Regression

multiple "

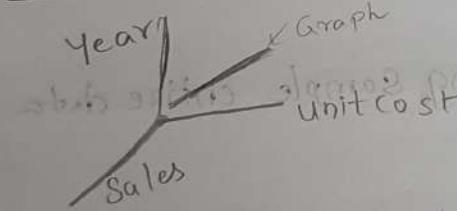
- histograms
- clustering
- sampling
- Datacube aggregation



- 1) Linear Regression We can Represent entire data within a single Straight line (b₂ predictor variable)
(Applied on Databases i.e. 2d)



- 2) multiple Regression Applied on Datacubes i.e. 3d



(more than 2 attributes)

* Entire data is reduced and represented in the straight line.

	1	2	3
1	T	F	F
2	F	T	F

	1	2	3	4
1	0.25	0.25	0.25	0.25
2	0.25	0.25	0.25	0.25
3	0.25	0.25	0.25	0.25