

Data Warehousing and Data Mining

- Data Mining- Mining Knowledge from large amounts of data
- Mining of gold from rocks or sand
- Knowledge mining from DB, knowledge of extraction data pattern / analysis.

Vim
Bm

Knowledge discovery in databases

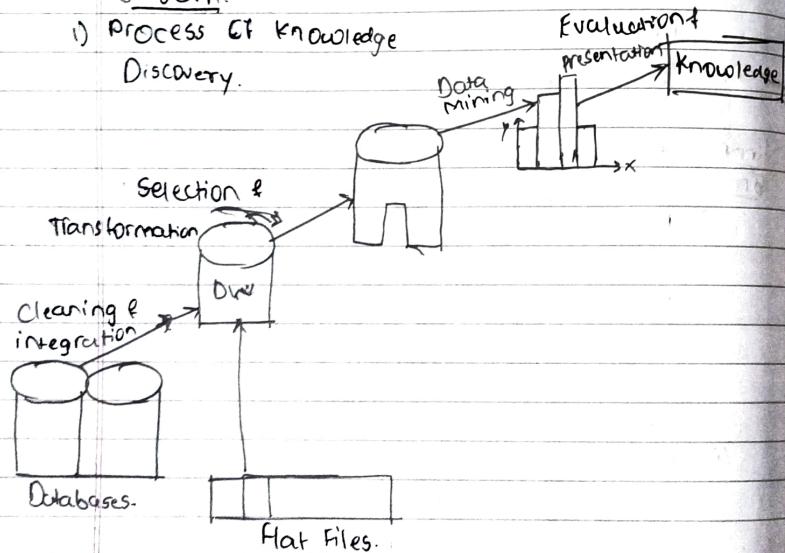
Knowledge discovery in databases (KDD)

- 1) Data Cleaning- To remove noise and inconsistent data
- 2) Data integration: where multiple data sources may be combined
- 3) Data selection: where data relevant to analysis are retrieved from databases.
- 4) Data transformation: where data transformation is done. Consolidated in 2 forms for mining by performing summarizing or aggregation approach
- 5) Data Mining: where intelligent methods are applied for extracting data patterns.

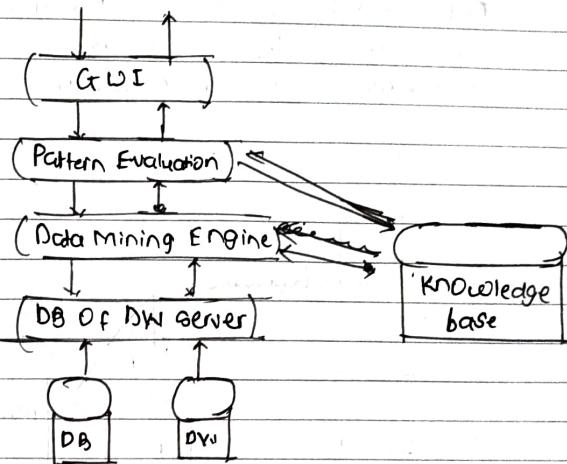
- 6) Pattern evaluation - to identify truly interesting patterns representing patterns-knowledge
 7) Knowledge presentation - where knowledge and visualization represent the mind knowledge to users

Diagram:

- i) Process of knowledge discovery.



Eg: Architecture Of typical system.



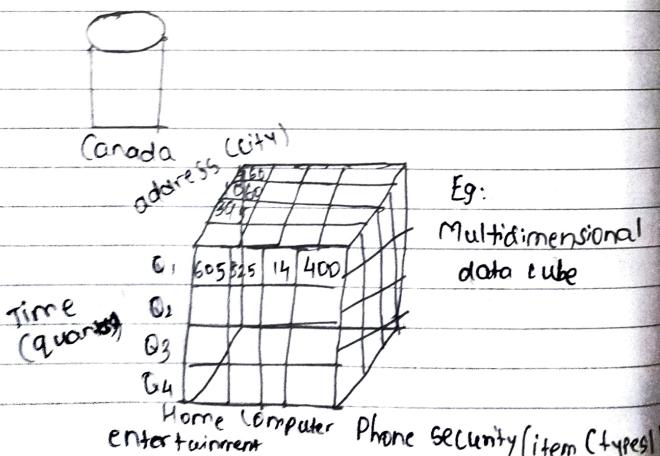
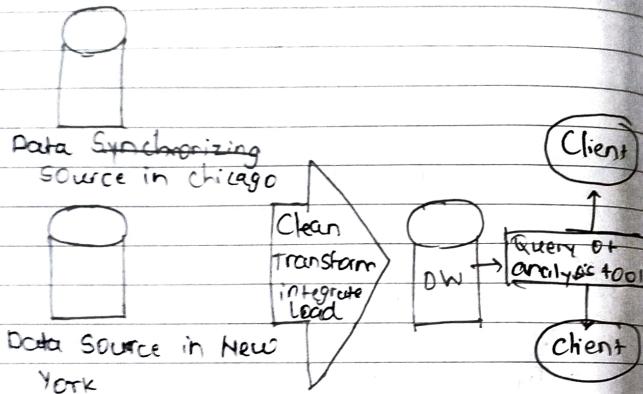
* Types of Data

- i) Relational Databases: Database Management Systems consists of collection of interrelated data called as database. A relational database is a collection of tables each of which is assigned a unique name. Each table consists the set of attributes and stores large set of tuples.

Customer

Cust Id	Name	Address	Age	Income
101	AAA	Canada	21	\$2700.

2) Data Warehouse: It is modeled by multi-dimensional structure where each dimension corresponds to an attribute in the schema.



3) Transactional databases: Transactional database consist of a file where each record represents transaction. A transaction includes transaction identity no which is unique, and list of items making up the transaction. it further contains other info such as transaction Id, no branch, Sales etc.

Sales

Trans-Id	List Of Items-Id
T100	T1, T3, T8
T200	T2, T4, T6
;	;
))

4) Advanced database Applications: Database Applications include

- i) Handling spatial data
- ii) Engineering design data
- iii) HyperText ~~multimedia data~~
- iv) Time related data
- v) WWW
- vi) Multimedia data

5) Object Oriented Databases:

It is based on OOPS where each entity is considered as objects. Objects can be individual employee, customer or item. Data and code are related to an object encapsulated in single unit.

- Spatial DB
- Temporal DB and Transferring DB
- Text DB and Multimedia DB
- Heterogeneous DB
- WWW

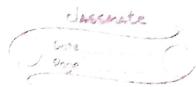
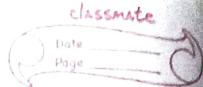
12 marks DM Functionalities - Kinds of patterns can be mined

- Data Mining functionalities
 - 1) Descriptive
 - 2) Predictive

1) Concept / class Description : Characterization

2) Association Analysis.

a) It is the discovery of association rules showing attribute value condition that occur frequently together in a given set of data



a) $\text{age}(x, "20\text{-}29") \wedge \text{income}(x, 20k..29k) \Rightarrow \text{buys}(x, "CD\ player")$ [Support = 21, Confidence = 60%]

b) $\text{contains}(T, "Computer") \Rightarrow \text{contains}(T, "Software")$ [Support = 1%, Confidence = 50%]

3) Classification and Prediction

- Distinguish data classes / concepts.
- Training data.
- (IF-THEN) Rules

Eg: Good Response, Mild Response and Non-response

- All electronics company.

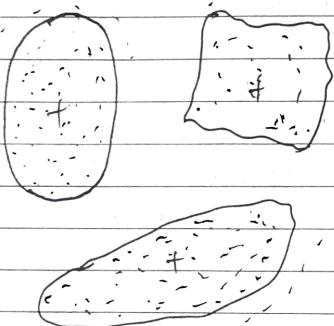
4) Clustering Analysis

- Analysis data objects without consulting a known class label
- Class labels are not present in training data

Eg.	Question	Label
1	What is Cluster analysis	Prescriptive
2	Water content in bottom	Numeric
3	Where ACT is located	Location

- * Unsupervised → without label
- * Supervised → with label

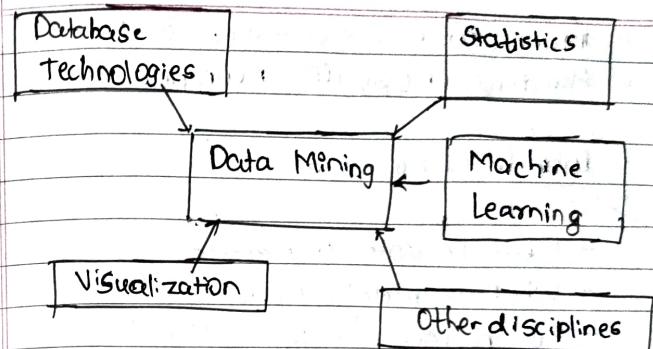
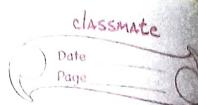
Defn: The Objects are clustered or grouped based on the principle of maximizing intra-class similarity and minimizing the inter-class similarity. The clusters of objects are formed so that objects within a cluster have high similarity in comparison with one another and dissimilarity in other clusters.



⑤ Outliers Analysis.

- Don't comply with the general behaviour

→ Classification of Data Mining System.



- 1) Classification according to databases mined.
- 2) Classification according to knowledge mined
↳ Based, characterization, etc
- 3) Classification according to Techniques Mined
↳ ML, Datawarehouse, Visualization.
- 4) Classification according to Applications adopted
↳ Finance, Stock Market, Weather Forecasting.
- ★ Major Issues in Data Mining
 - 1) Mining different kinds of knowledge in DB.
 - 2) Interactive mining of knowledge at multiple levels of abstraction.
 - 3) Data Mining query languages.

- 4) Present And visualization of DM results.
 5) Handling Noisy Or incomplete data.

Data Processing

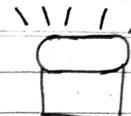
Learnings

- Need of data preprocessing
- What is a data processing
- To improve the quality of data
- Data pre-processing techniques.

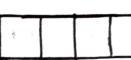
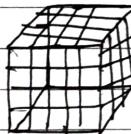
Data Cleaning, Data Integration, Data Transformation,
 and Data Reduction.

- 1) Data Cleaning: is used to remove noise and correct inconsistent data.
- 2) Data integration: Merges data from coherent data multiple sources into coherent data to store such as a Datawarehouse.
- 3) Data Transformation: Normalization may be applied to improve efficiency & accuracy of mining algos.
- 4) Data Reduction: Can reduce the data size Aggregating, eliminating redundant features or cluster.

Data
Cleaning



Data
Integration



Data transformation - $2.32, 100, 59 \rightarrow 0.02, 0.32, 0.100, 0.59$

Data Reduction

	A ₁	A ₂	A ₃	...	A ₁₂₆
T ₁					
T ₂					
T ₃					
T ₄					
T ₅					
T ₂₀₀₀					

	A ₁	A ₂	...	A ₁₅
T ₄				
T ₅				

Fig Forms Of Data Preprocessing.

Data Cleaning

i) Data Cleaning attends to fill the missing values, identifies the outliers and correct inconsistencies in data and smooth out the noise.

i) Missing Values

- Many tuples have no recorded value for several attributes such as cust-name

- Methods to fill missing values.

i) Ignore the tuples

ii) Use a global constant (unknown label -∞)

iii) Use attribute mean to fill the missing values.

iv) Use the most probable values for missing value. (Using other customer attribute in data set consider decision tree to predict the missing value.)

v) Fill missing values manually.

ii) Noisy Data

→ It's a random error or variance in a measured variable

a) - Data smoothing technique

i) Binning

Sorted data for price

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

Partition into (equidepth) bins:

Bin1: 4, 8, 15

Bin2: 21, 21, 24

Bin3: 25, 28, 34

Smoothing by bin mean

Bin 1: $\frac{4+8+15}{3} = 9$, 9, 9.

2: 22, 22, 22

3: 29, 29, 29.

Smoothing by bin boundaries

Bin 1: 4, 4, 15

2: 21, 21, 24

3: 25, 25, 34

b) Clustering

Outliers can be detected by clustering where similar values are organized into groups or clusters. Values that fall out of set of clusters considered as outliers.



c) Regression

- Data can be smoothed by fitting data to a function.

Linear Regression finds best line to fit two variables, One variable can be used to predict other. Multiple Linear regression where more than 2 variables are involved.

→ Data Integration and Termination

Issue

Data Integration

Issues to Consider data integration

- 1) Schema integration → Entity identification
- 2) Redundancy
- 3) Detection and resolution of data value conflicts.

Data transformation: The data transformed or consolidated into form appropriate for mining. It involves:

D) Smoothing: Which work to remove noise from data.

2) Aggregation: Where summary or aggregation is applied to the data.

Generalization: Generalization of data where low level data replaced by higher level concepts through use of hierarchy.

Eg: Street → City or Country

Age → young, middle age and senior

Normalization: Where attribute data are scaled to fall within a specified range ~~as to~~

0.2 - 1

Attribute Construction: Where new attributes are constructed and added from given set of attributes.

* Data Reduction:

- To obtain a reduced representation of the data set
- Mention integrity of original data.
- 1) Data cube Aggregation
- 2) Dimension Reduction
- 3) Data Compression.

Data Reduction

Data Cube Reggregation

Year = 1999	
Year = 1998	
Year = 1997	
Quarter	Sales
Q1	224
Q2	408
Q3	350
Q4	158

Year	Sales
1997	156
1998	258
1999	359

Dimension Reduction: Dataset 100's of attributes many of which are irrelevant to mining task or redundant. Dimensionality reduces the size of dataset by removing such attributes or dimensions from it.

Forward

Selection

Initial Attribute set
[A₁, A₂, A₃, A₄, A₅, A₆]

Initial reduced set
 $\{ \cdot \}$
 $\rightarrow \{ A_1 \}$

$\rightarrow \{ A_1, A_4 \}$

Reduced $\rightarrow \{ A_1, A_2, A_6 \}$

Backward

Selection

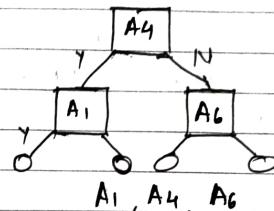
Initial Attribute set
[A₁, A₂, A₃, A₅, A₆, A₇, A₈]

$\Rightarrow \{ A_1, A_3, A_4, A_5, A_6 \}$
 $\Rightarrow \{ A_1, A_4, A_5, A_6 \}$

Reduced $\rightarrow \{ A_1, A_4, A_6 \}$

Decision Tree

Induction



Data Compression: Data encoding or transformation are applied to obtain reduced or compressed representation of data, if the original data is reconstructed from compressed data without any loss of info is called lossless

Unit - 2Data Warehouse

→ Provides tools for business executives to systematically organize, understand and use their data to make strategic decisions.

6m ★ Key features of datawarehouse

1) Subject Oriented: It is organized around various subjects such as customers, products, suppliers and sales.

2) Integrated: Data warehouse is constructed by integrating multiple heterogeneous sources.

3) Time variant - Data are stored to a historical perspective (past 5-10 yrs)

4) Non Volatile: Data warehouse is a physically separate store of data transformed from application data found in the operational environment.

Difference betwⁿ Operational DB systems and datawarehouse.

- OLTP → Cover day to day operations
Eg: Banking, inventory, purchasing
- DW → Data analysis and decision making.

*12m	OLTP	DLAB
Characteristics:		
Features:		
Features:		
1) Characteristics	Operational processing	Informational processing
2) Orientation	Transaction	Analysis
3) Users	Clerk, DBA, DBP	Knowledge worker Eg manager, executive or analyst.
4) Function	Day to day operations	Long term informational requirements.
5) Database Design	ER based	Star or Snowflake, Subject Oriented

6)  Access

Read and write

Mostly Read

7) NO of
records
accessed

Tens (10)

Millions.

8) NO of
users

Thousands

Hundreds

9) DB size

100Mb- 1GB

100GB- 1TB.

2-D view of sales data

Location = USA

item (type)

time	home	computer	phone	security
.....	rent			
Q1	605	825	14	400
Q2				
Q3				
Q4				

- Data warehousing

Multidimensional Data Model

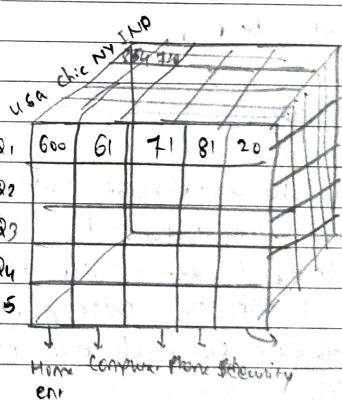
- View data in form of datacube
- Data cube \rightarrow Data to be modelled and viewed in multiple dimensions.

- Dimension

Defined by dimensions by facts.

Eg: Dimensional \rightarrow time, item, branch

Facts \rightarrow Dollars Sold, Units-Sold



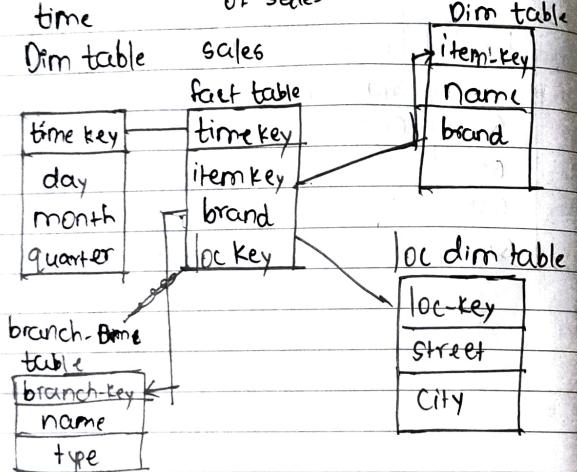
3-D Datacube pre representation:

SNOW

12m-15m Star, starflake and fact constellations
schemas for multidimensional DB

1) Star schema

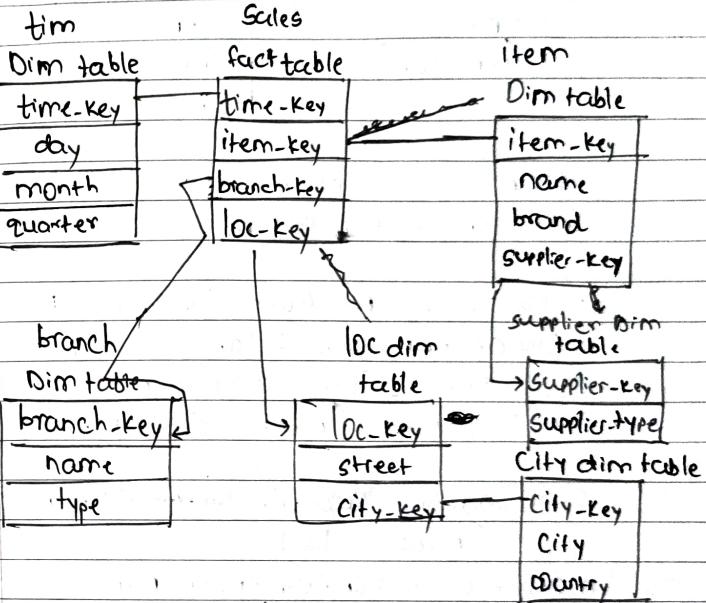
* Star schema for data warehouse
of sales



The data warehouse contains

- A large central table (fact table) containing the bulk of table data with no redundancy
 - A set of smaller attendant tables (dimension table) one for each dimension.
- 2) Snowflake Schema: It is the variant of Star schema table where some dimension tables are normalized thereby further splitting

the fact data into additional tables.



3) Fact Constellational Schema

It requires multiple fact tables to share dimensional tables. This kind of schema can be viewed as a collection of stars. And hence it's called galaxy schema or fact constellation schema.

Sales
fact table

time-key
item-key

Shipping
fact table

item-key
Shipping Key
time-key

Examples for defining star, snowflake and fact constellation Schema

- SQL based DM Query language - DML
- DW defined using 2 language primitives
 - One for Cube definition
 - One for dimension definition

Syntax: define cube <CubeName> <dimension lists>
(measure list)

2) define dimension <dimension names> as
(Attribute list)

Eg: Star Schema

- 1) define cube Sales-star [time, item, branch/
: dollars-sold = sum (sales-in-dollars), units-sold
= count (*)]
- 2) define dimension time as (time-key, day, month, yr)
define dimension item as (item-key, name, branch)

define dimension branch as (branch-key, name)

Eg: Snorflex schema

star
define cubes snorflexsales-snorflex [time, item, branch]
: dollars-sold = sum (sales-in-dollars), units-sold =
count (*)

define dimension time as (time-key, day, month, yr)

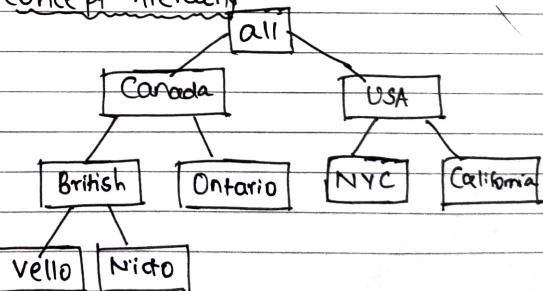
define dimension location as (location-key, day, month)

define dimension city as (city-name, cname,
(Country))

* define dimension location as (location-key, street,
city-key, state, country)

* define dimension item as (item-key, name,
type, supplier (supplier-key, type)).

* Concept Hierarchy



Country Defⁿ → Concept Hierarchy defines a sequence of mapping from low level concepts to higher level. The mapping forms a concept hierarchy for the dimension locⁿ mapping a set of low level concepts (cities) to higher level concepts (countries)

```
graph TD; Street((Street)) --> Store((Store)); Store --> City((City)); City --> State((State)); State --> Country((Country))
```