

14/3/23

Date		
Page No.		

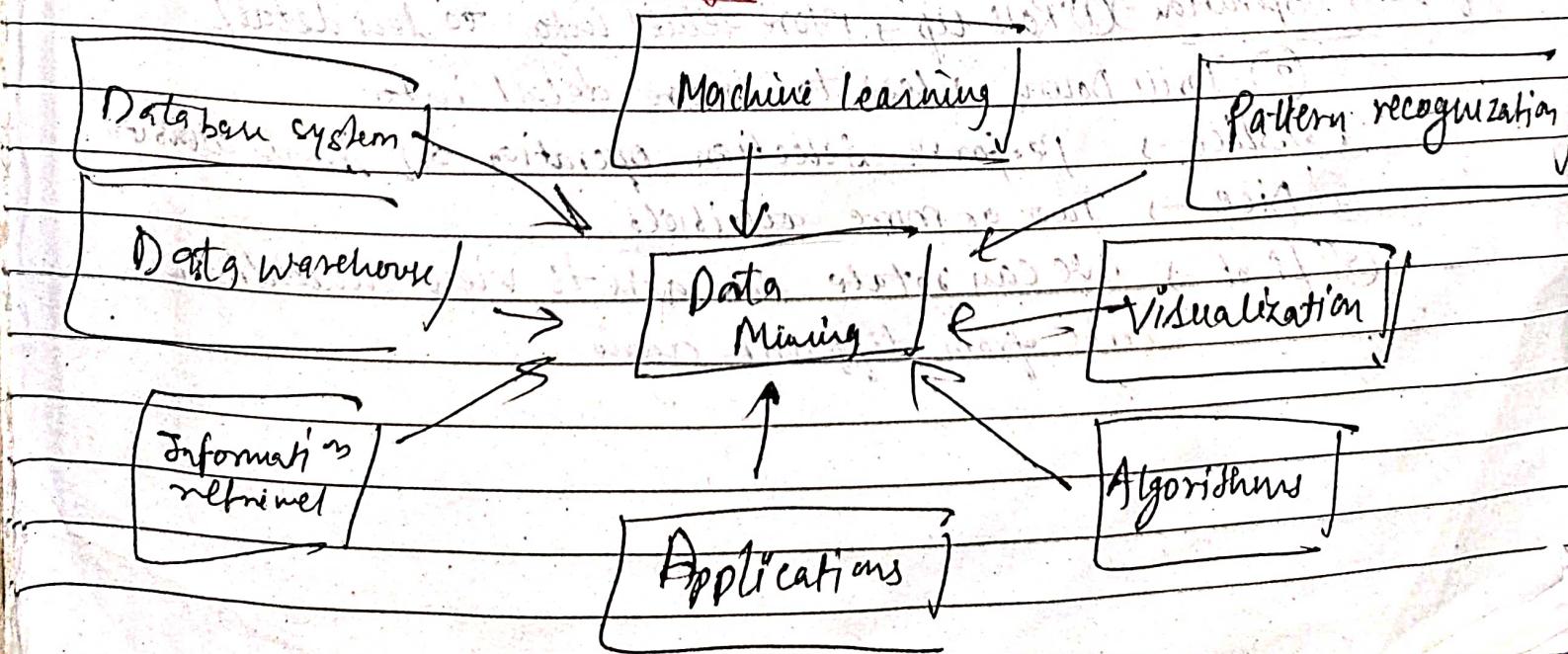
* What is data mining?

- Data mining is one of the most useful techniques that help enterpreneurs, researchers and individuals to extract valuable information from large set of data. Data mining is also called knowledge discovery in Database (KDD).
- Process of extracting useful information from larger sets of data.
- Also known as knowledge discovery / knowledge extraction / pattern analysis

(1960's) → Storing of data in form of hardcopy.

(1970's) → RDB on database. → Storing of data in the form of many tables, in which row & columns are therefore

Main purpose of Data Mining



* Data Warehouse:

Date		
Page No.		

Date		
Page No.		

- A data warehouse is a repository of heterogeneous database collected from multiple sources over a history of time, stored under a unified schema and used for analysis and decision making.
- Used for OLAP → Online analytical processing
- Data warehouse are constructed via a process of data cleaning, data integration, data transformation, data loading etc.
- Usually modeled by a multidimensional data structure, called a data cube.
- Each dimension corresponds to an attribute or a soft off attribute.
- advantage → 1) Perform various types of OLAP operation
2) View multiple data at once
- Ex of OLAP operation → Roll up → More detail info to less detail
→ Drill down → Less to more detail info
→ Slice → perform selection operation on single attribute
→ Dice → Two or more attribute
- Pivot → We can rotate the axis to view multidimensional cube from different angle
- | Item | Category | Quantity | Price |
|------|----------|----------|-------|
| Q1 | 605 | 82.5 | 14 |
| Q2 | 606 | 82.5 | 14 |
| Q3 | 607 | 82.5 | 14 |
| Q4 | 608 | 82.5 | 14 |

Date		
Page No.		

Date		
Page No.		

→ Steps involved in data mining process
The knowledge discovery process consists of performing following steps:

Evaluation & presentation → patterns

knowledge

Data mining → patterns

Selection

Transformation

Data warehouse

Data cleaning

Data integration

Data base

① Characterization → Summarization of the general features of a target class of data

② Association Rule Mining →

③ Data cleaning → To remove noise and inconsistent data
→ In real life we encounter the cleaned data have to check if first

④ Data integration → When multiple data sources may be combined.

⑤ Data selection → When data relevant to the analysis task are selected from the database
→ Selection of only the relevant features out of all the data

⑥ Data transformation → Where data are transformed and consolidated into forms appropriate for mining by performing summary.

⑦ Data mining → An essential process where intelligent methods are applied to extract data patterns

⑧ Pattern evaluation → To identify the truly interesting patterns representing knowledge based on interesting measures
⑨ Knowledge presentation → Where visualization & knowledge representation techniques are used to present mined knowledge to user

Date	
Page No.	

Date	
Page No.	

- ① collect the data
② then do processing

④ Predictions → predict missing or unavailable numeric data

⑤ Clustering → Unlike classification & regression, which analyze class labeled datasets, clustering analyzes data objects without consulting their class labels.

⑥ Outlier Analysis → Analysis of noisy data.

⑦ Evaluation analysis → Trying to develop the model which describes and models trends for

- objects whose behavior changes over time.

15/03/23

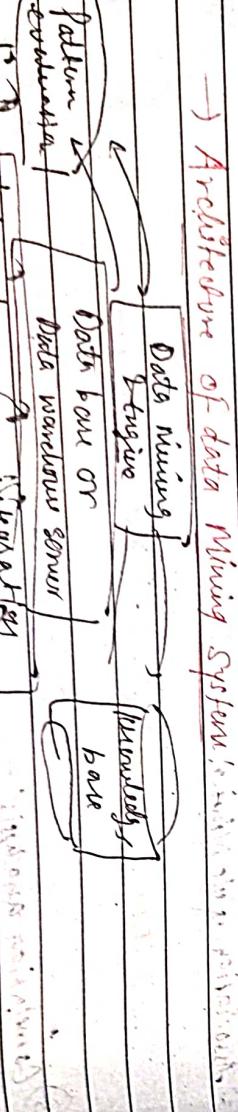
⑧ Types of databases

① Relational Database → A relational database is a collection of tables,

each of which is assigned a unique name. Each table consists of a set of attributes (columns or fields) and usually stores large set of tuples (records or rows). Each tuple in a relational table represents an object identified by a unique key & described by a set of attribute values.

② Data warehouse → already written

③ Transactional database → Each record in a transactional



→ data cleaning, integration and user interface selection.

- ④ Object relational database → All object concepts
- ⑤ Temporal database → Relational database which includes
 - related attributes
 - - Bookings (Meeting Room, Time)



- ⑥ Object relational database → All object concepts
- ⑦ Temporal database → Relational database which includes
 - related attributes
 - - Bookings (Meeting Room, Time)



16/03/23

Date	
Page No.	

(B) Sequence database → Stores sequence of ordered events with or without concrete notion of time.

Ex → Formation of DNA

→ Customer shopping sequences

(C) Time-series database → Stores sequences of values over repeated measurement of events, obtained over repeated time.

(Hourly, daily, weekly, monthly)

Ex → Stock exchange data.

(D) Text Database → Contains long sentences & paragraphs

Ex → Product Specification
→ Error report
→ Warning message.

Ex → Highly unstructured - ex - Web page
→ Semi-structured - ex - Email
↳ Structured - ex - Library catalogue

(E) Multimedia database → Stores images, audio, video.

(F) Heterogeneous database → Consist of a set of intramodular autonomous components of database.

→ VLSI → utilizes the ability to design ultra-se space.

→ It saves a huge amount of space-related data, including maps, preprocessed remote sensing, or medical imaging records, and VLSI chip design data.

Application of geographical data base

① Forestry and ecology, planning

2) Providing public service information regarding the location of telephone and electric cable, water pipes & sewage system.

3) Vehicle navigation and dispatching system.

4) Text Database: → These contain word description for NLP simple keywords but objects

→ Contains long sentences & paragraphs

5) Multimedia database: → ex - Web page

→ Error report
→ Warning message.

6) Heterogeneous database: → ex - Email

7) Structured - ex - Library catalogue

8) Multimedai database: → Stores images, audio, video.

9) VLSI → Utilizes the ability to design ultra-se space.

10) It saves a huge amount of space-related data, including maps, preprocessed remote sensing, or medical imaging records, and VLSI chip design data.

11) Application of geographical data base

12) Forestry and ecology, planning

Date	
Page No.	

Date		
Page No.		

Date		
Page No.		

Q1/03/23 Data warehouse → It is a subject-oriented, integrated, time variant and non-volatile collection of data in support of management decision making process.

Difference between OLTP vs DSS:

OLTP → Online Transaction Processing

OLAP → Online Analytical Processing

Q) These are the online operational systems processing database systems

Warehouse system

② Online transaction and Data analysis and decision

③ Energy processing working

④ Customer oriented Market oriented

⑤ Manages current data Historical data

⑥ ER datamodel Star schema or snowflake schema

⑦ CR datamodel Galaxy schema

⑧ 100 TUPP - Some DB

⑨ 100GB - TB

⑩ User → manager predictive and data analysis

Database, professor

Star schema:-

The schema contains a central fact table with a bulk of data with no repetition at the centre.

→ Set of smaller dimension tables one for each dimension

Q) There is a datawarehouse consist of 3 dimension 1 fact doctor, patient & two measured → count, charge, draw, star schema for this

Hospital Doctor Dimension table

Time Dimension table

Patient Dimension table

Doctor Dimension table

Branch Dimension table

Age Dimension table

Diagnosis Dimension table

Change Dimension table

Diagnosis Dimension table

A) Different operations we can perform on multidimensional data cube:-

① Roll up → Performes aggregation on a data cube.

② Drill down → It navigates from less detailed to more detailed data.

③ Slice → Performes a selection on one dimension.

④ Dice → Selection on two or more dimension.

⑤ Pivot → Rotates the data axis in view.

Date		
Page No.		

Date		
Page No.		

★ Snowflake :-

The snowflake schema is a variant of the star schema model, where some dimension tables are further normalized.

Address-key

Address dimension table

- ① Easy to maintain
- ② Saves memory space
- ③ Reduces Repetition

- ④ Consists of the four dimensions [student, course, teacher, city]
- ⑤ Contains a snowflake schema. Δ Δ Δ Δ

Course ID	Course fee	Course name
State	City	Country

Course duration	Student key	University	Student age
Semester	Student key	University	Student age

Instructor key	Semester key	Subject	Address
State	Semester	Subject	Address

Subject	avg. grade	Instructor key	Address
City	avg. grade	Instructor key	Address

Subject	avg. grade	Instructor key	Address
Country	avg. grade	Instructor key	Address

Subject	avg. grade	Instructor key	Address
State	avg. grade	Instructor key	Address

Subject	avg. grade	Instructor key	Address
City	avg. grade	Instructor key	Address

Subject	avg. grade	Instructor key	Address
Country	avg. grade	Instructor key	Address

Course Dimension

Mark

Professor

Student Dimension

Address Dimension

- ① Roll up in time from day to year
- ② Slice for the time 2004
- ③ Roll up from doctor from individual patient to all

- ④ Suppose that a date warehouse for Big University

- ⑤ consists of the four dimensions [student, course, teacher, city]

- ⑥ Semester & Instructor and two measure [course fee, contact no.]

- ⑦ Contains a snowflake schema. Δ Δ Δ Δ

- ⑧ To year) Should you perform in order to list the average grade of CS course for each Big University Student?

- ⑨

Difference b/w Database and Databrowsing

Database System

Date

- (1) Roll up column from Course ID to department
- (2) Then Roll up on Student from Student-ID to university "you"
- (3) perform the dice operation
- (4) All on course, student with department CS L

(5) Drill-down on Student from university to Student name

(6) Draw star schema. [date, spectator, location, game]

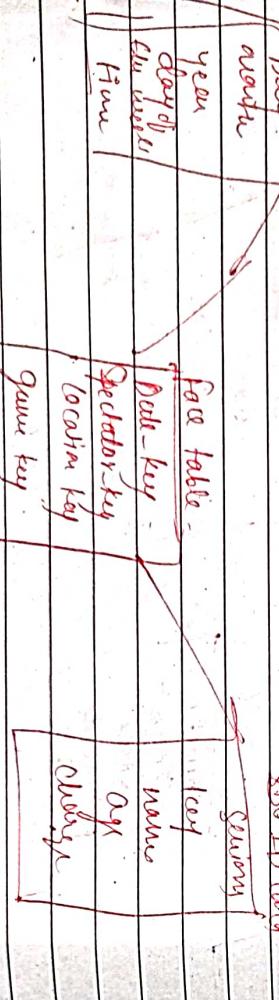
Spectator
dimension table

(7) Data is updated when transaction occurs.

Data is update of selected process.

Stu ID and

dimension table



fact table
dimensions

location key
city
state
street
country

name
genre
count