

Data Mining is also known as KDD

Date
Page

Batch

MOD 1

~~Def~~

DATA MINING

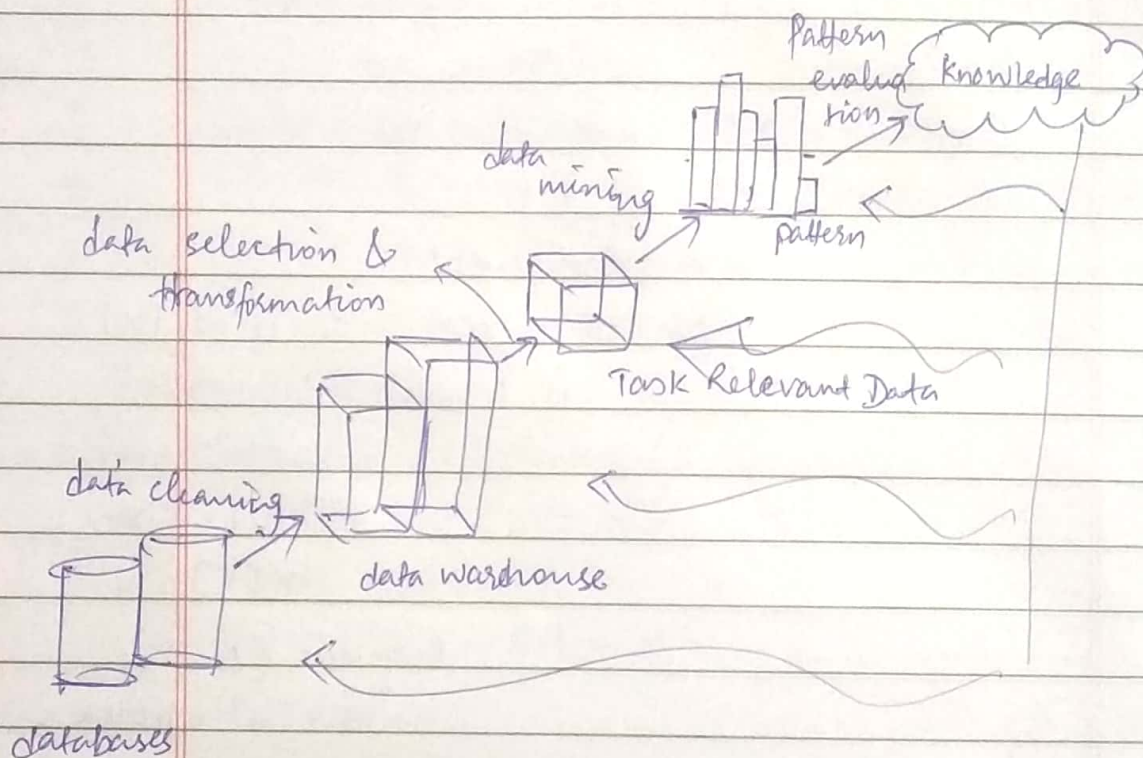
process of

#

Data Mining is the extraction of interesting (non-trivial, implicit, previously unknown, and potentially useful) patterns or knowledge from huge amount of data.

#

Knowledge Discovery in Database (KDD) process -



Need of DM on what data -

- RDBMS, Transactional data
- data streams & sensor data
- time series data, sequence data
- structure data, graphs, social networks
- multimedia dB
- text dB
- www.

DM Functions -

1. Generalization - Info integration & data warehouse construction includes data cleaning, transformation, integration & multidimensional data model such as generalise, summarise and contrast data characteristics.

2. Association & Correlation Analysis - defining & identifying frequent patterns through association & correlation & how to use this pattern for classification, clustering & other applications.

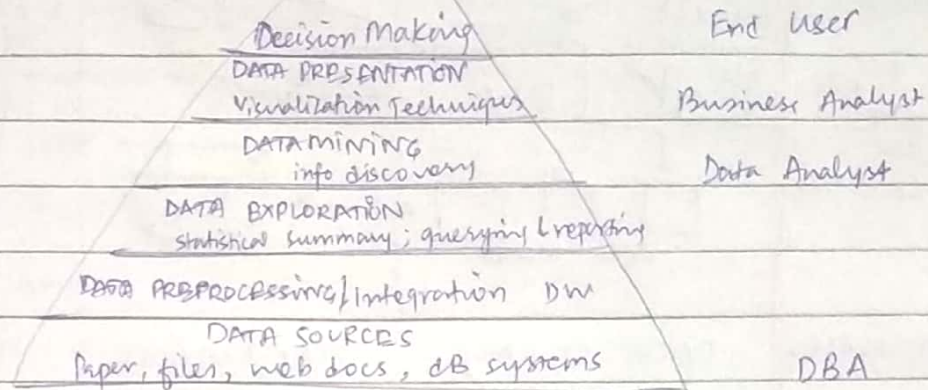
3. Classification - Construct models/fⁿ based on some training using typical methods like decision trees, neural networks, rule based classification, pattern-based class, logistics regression, etc.

4. Cluster Analysis - Group data to form new categories based on principle: maximizing intra-class similarity & minimizing interclass similarity.

5. Outlier Analysis - Outlier is a data object that does not comply with general behavior of data and they are found as a by product of clustering or regression analysis.

Business Intelligence (BI) -

Increasing potential to support business decisions.



Data Warehouse (DW) - it is a subject oriented, integrated,

USAGE

time-variant & nonvolatile collection of data in support of

① Info process

management's decision making process. the process of constructing

② Analytical processing

data warehouse is data warehousing.

③ DM

→ Types -

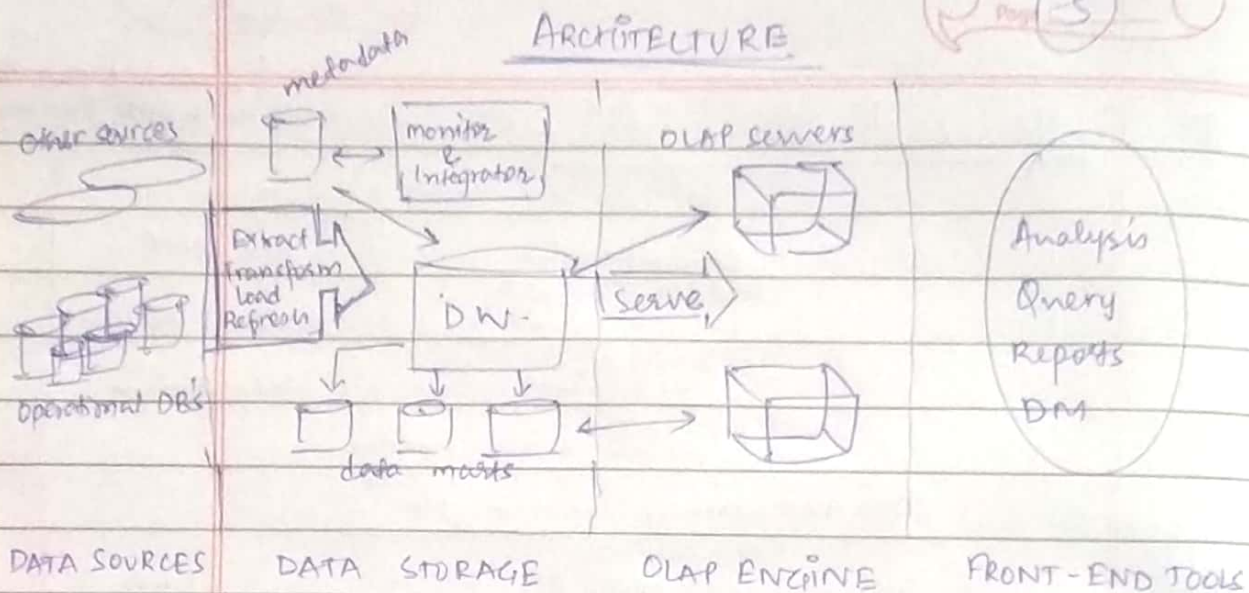
1. Subject oriented - organized around major subjects, such as product, cost or sales. Provides a simple & concise view around a particular issue by excluding data that are not useful.

2. Integrated - constructed by integrating multiple, heterogeneous data source, also applying data cleaning & integration techniques.

3. Time Variant - provides info from a historical perspective

③ Data contains an element of time, implicitly or explicitly.

4. Non-volatile - A physically separate store of data, transformed from operational environment where no updates can be made - only initial loading of data & access of data is allowed.



ETL (Extraction, Transform & Loading) -

- **Data Extraction** - get data from multiple, heterogeneous, & external sources
- **Data Cleaning** - detect errors in data & rectify them when possible.
- **Data Transformation** - convert data from legacy or host format to warehouse format
- **Load** - sort, summarize, consolidate, compute views, check integrity & build partitions.
- **Refresh** - propagate the updates from data sources to warehouse.

Relation b/w BI & DW-

BI tells us what happened, or is happening right now in your business - it describes the situation to you, and BI make use of data stored in DW and lets you apply chosen metrics to potentially huge, unstructured ^{used} data set, DM, OLAP, and reporting as well as business performance monitoring, predictive & prescriptive analytics.

OLAP - it is an online system that reports to multidimensional analytical queries like forecasting, reporting, etc

OLTP - is a system that manages transaction-oriented applications or internet like ATM.



OLTP

1. its an online transactional system & manage db modifications.
2. main focus is insert, update, delete info from db.
3. OLTP & its transactions are original source of data.
4. OLTP has short transactions,
5. Processing time of transactions is ↓
6. simple queries
7. tables in OLTP db are normalized
8. db size - 100 MB - GB
9. users - IT professionals
10. db design - application oriented

OLAP

- its an online data retrieving & data analysis system
- main focus is extract data for analyzing that helps in decision making.
- diff OLTP's db becomes the source of data.
- OLAP has long transactions
- ↑
- complex queries
- not normalized.
- 100 GB - TB
- knowledge workers
- subject oriented.

Modelling of DW (Schemas) -

1. Star - fact table in the middle connected to a set of dimension tables
2. Snowflake - a refinement of star schema where some dimensional hierarchy is normalized into a set of smaller dimension tables.
3. ^(FC)
Fact Constellations - multiple fact tables share dimension tables, viewed as a collection of stars, therefore called galaxy schema or FC.

OLAP models/operations -

1. Rollup / drill up - summarize data by climbing up hierarchy or by dimension reduction
2. Roll down / drill down - reverse of roll up: converting data from higher level to lower level or introducing new dimensions or more
3. Slice & dice - One dimension is selected and a new cube is created
4. Pivot / Rotate - you rotate the data axes to provide a substitute presentation of data.

Types of OLAP -

1. Multidimensional OLAP (MOLAP) - it implements operation in multidimensional data
2. Relational OLAP (ROLAP) - its an extended RDBMS along with multidimensional data mapping to perform the standard relational operation.
3. Hybrid OLAP (HOLAP) - in this, the aggregated tools are stored in multidimensional db while detailed data is stored in RDB. This offers efficiency of both ROLAP & MOLAP.
4. Web OLAP (WOLAP) - its in which OLAP system is accessible via web browser.
5. Mobile OLAP - it's one in which OLAP system is accessible via mobile phones.

#	ROLAP	MOLAP
1.	Relational OLAP	Multidimensional online Analytical processing
2.	data is stored and fetched from main Dm.	data is stored and fetched from the propriatary db MDPBs.
3.	data is stored in form of relational tables.	data is stored in large multidimensional array made of data cubes.
4.	large data volumes	Limited data summaries is kept in MDPBs.
5.	Slow access	Faster access
6.	& creates a multidimensional view of data dynamically.	already stores the static multidimension view of data in MDPBs.