

# Data Mining and Analytics

Information - collection of data

Data :: Unprocessed things / raw data.

Mining :: Extraction of information.

Knowledge :: Useful information converting from the raw data.

Pattern :: Generate patterns / insights from available data.

Analysis:-

Market

Basket

Analysis:-

⇒

Extract useful & meaningful information / patterns from the data.

clustering ⇒ grouping, No class / label, grouped based on similarity.

classification ⇒ labeling, putting under a label or class.

Note:-

classification is based on the label.

clustering → classification.

Prediction:- (CV-2)

Predicting with the already existing data

eg:- weather prediction.  
Stock market

Business Sector  
Education

Health care - Cancer prediction  
House rent prediction.

## Association Rule Mining (V-2)

if - then rule.

if (age group  $18 < \text{age-group} < 22$ ) {

then college student.

Pre-processing:-

⇒ Removal of unwanted data = Data cleaning

⇒ Data transformation

⇒ Data integration

⇒ Data reduction

DBMS:-

Querying ⇒ Fetching

Database vs Data warehouse

Database:-

⇒ Small amount of data

(transactional data)

⇒ Online Transaction processing

OLTP

Data warehouse:-

⇒ Large amount of data

⇒ Can't update, but more data can be included

⇒ OLAP ⇒ Online

Analytical

processing

# Machine Learning

Data mining :-

extracting previously unknown and potentially useful patterns from large amounts of data.

Knowledge Discovery in Database (KDD) :-

knowledge extraction  
data / pattern analysis

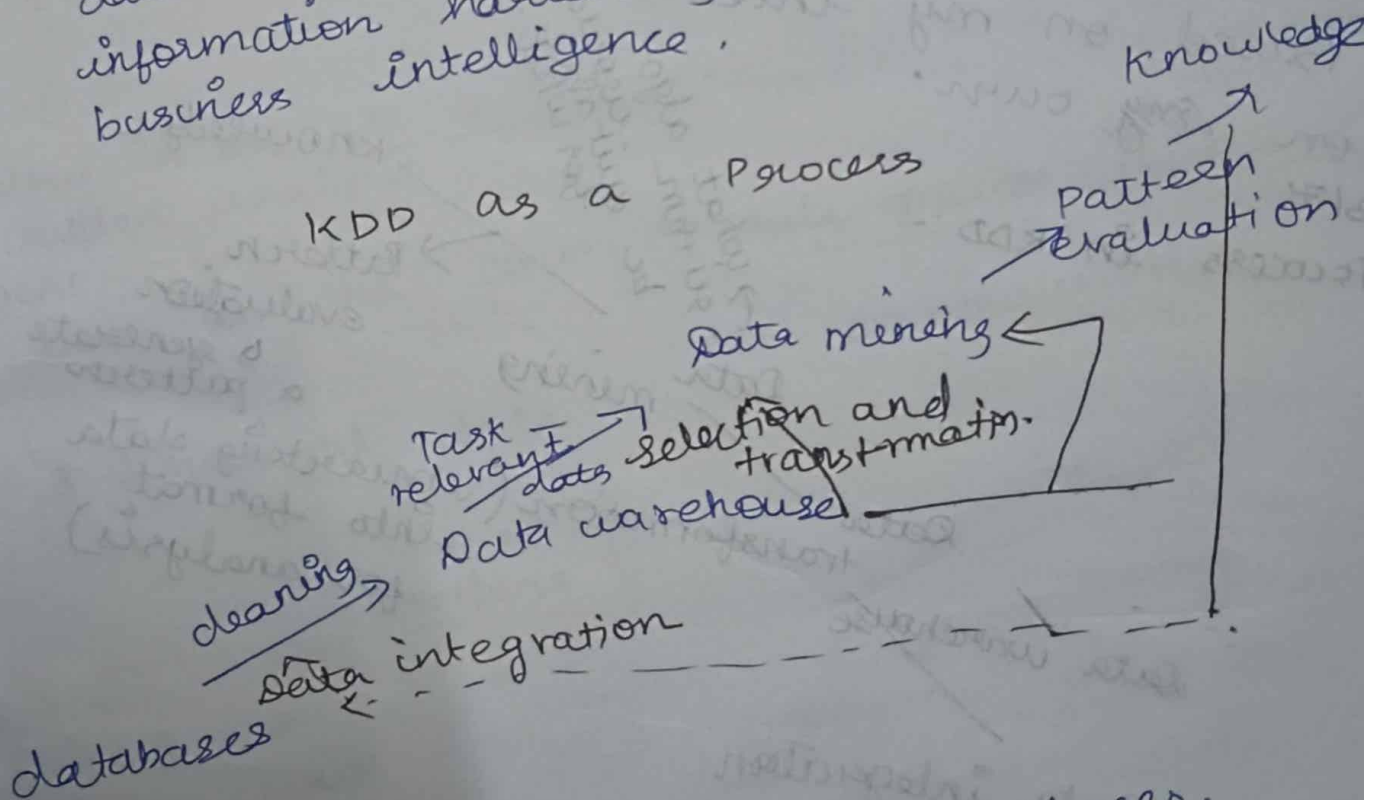
data archeology

data dredging

information harvesting

business intelligence.

KDD as a Process



Data cleaning :-  
Dealing with missing values.

Data lake  $\Rightarrow$  source of data  
large amt, structured / unstructured,  
place where raw data is stored.



Unstructured data :-

Video / audio files.

Semi-structured.

Excel + video files

Supervised :- labelled data

eg :- Classification

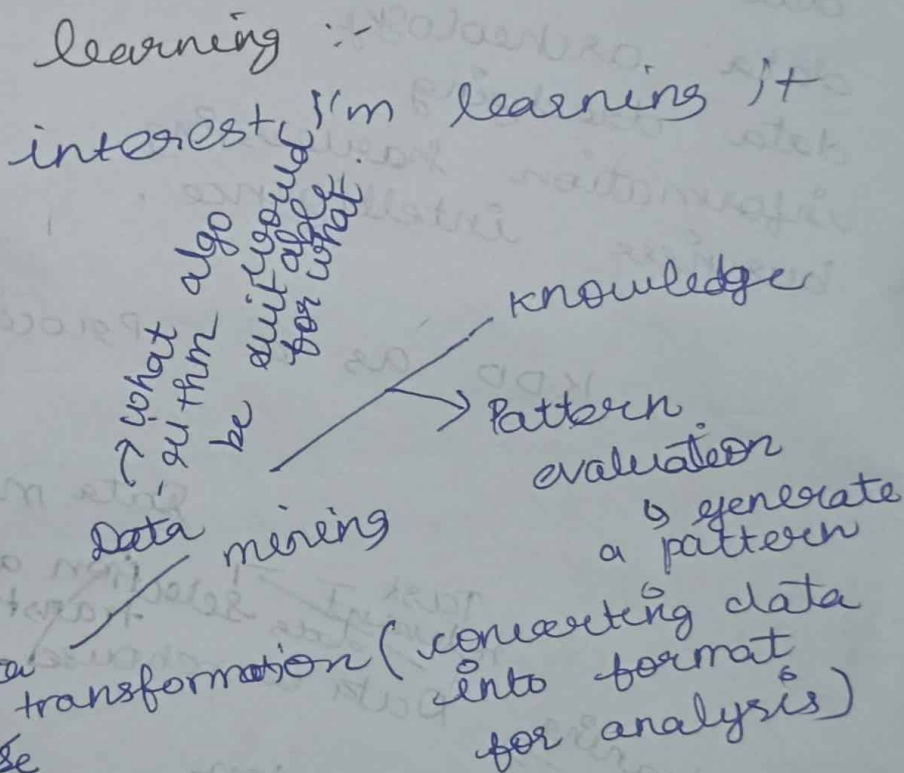
Unsupervised :- Un-labelled data

eg :- Clustering

Re-inforcement learning :-

Based on my  
own.

Process in KDD :-



Data warehouse

Data integration

Data cleaning

Data sources

Flat files DB warehouse

# Extraction of knowledge $\Rightarrow$ Data Mining.

Technologies related to DM :-

- \* DBMS
- \* ML
- \* OLAP
- \* Statistics

DBMS  $\Rightarrow$  software that's used to manage the data

DB  $\Rightarrow$  storage space of the data.

tuples  $\Rightarrow$  row      column  $\Rightarrow$  attributes

Note :-  
DM is used when it involves more than 2 or 3 join queries.

ML :-  
used in DM for better accuracy and prediction

Outlier analysis

OLAP :-

Data in the form of data ~~queries~~ cubes.  
 $\downarrow$   
3D structure  
3-dimensional data.

