

# Data Mining (Trouble-Free)

## Intro to Data Mining, Types of Data

→ is defined as procedure of extracting info. from huge sets of data.

— also known as mining knowledge from data.

What types of data can be mined :-

3 types - database (DB), data warehouse, transactional data.

## ① Database Data (RDBMS) - (Relational DBMS)

→ set of tables — has some features

↑  
tuples      attributes

While mining databases, we can search for trends / data patterns.

ex- ① Analysing customer data to predict the credit risks of new customers based on prev. data.

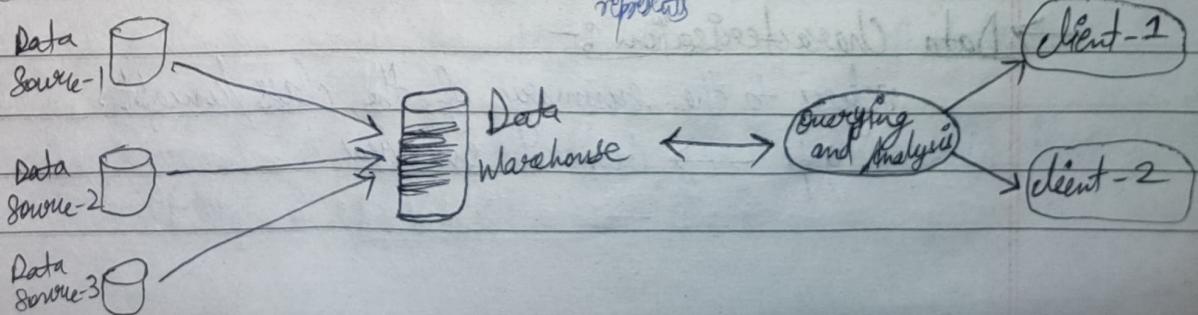
## ② Analysing sales data — (any deviations)

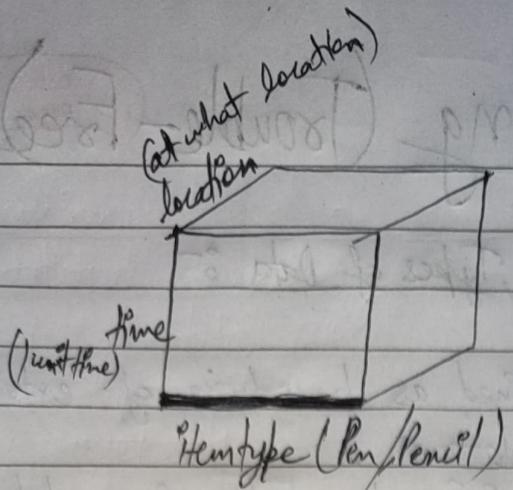
Like or less  
(dip)

② DataWarehouse data :-

Collection of data integrated from diff. sources with querying & decision making on data.

In data warehouse, data is stored in multidimensional structure (datacube) where each dimension is each attribute.





### ③ Transactional databases :-

Each record is called as transaction.

(sales, flight booking, user click on web page.)

Transaction has transaction ID, list of other items making transaction from transaction DB, we can mine frequent patterns.

### → other types of data :-

(stock market) (continuous data) (image) (IC design)

Sequence data, data streams, spatial data, engineering design data, hypertext, multimedia, web data, etc.

### ④ #2 Data Mining Functionalities :-

#### ① Concept / class definitions :- (descriptions)

data is always associated with class/concept.

descriptions can be done in 2 ways.

Data characterisation

Data  
Discrimination

#### → Data Characterisation :-

refers to the summary of the class/concept.

O/p  $\longleftrightarrow$  General Overview

## → Data Discretization :-

↳ Compares the common features of the two classes.

Then we noted down the changing values, and give off → bar charts / curves, etc.

## ② Mining frequent patterns, associations & correlations :-

### Frequent Patterns :-

↳ Things which are found most commonly in data.

Frequent Itemsets (data items / data objects)

Frequent Subsequence

Frequent Substructure

relationship analysis

### Association Analysis :-

It's a way of identifying the relation b/w various item.

Ex:- Used to determine sales of items that are frequently purchased together.

### Correlation Analysis :-

— mathematical technique.

— shows how strongly pair of attributes are related together.

Ex:-

Tall people, tend to have ~~more~~ more weight.

1<sup>st</sup> attribute

2<sup>nd</sup> attribute

## ③ Classification and regression for predictive analysis :-

### Classification :-

— process of finding a model that distinguishes data items

— decision tree is used for classification.

most appropriate  
for discrete

prediction of data

## Regression :-

Statistical methodology that is used for numeric prediction of missing data.

target data

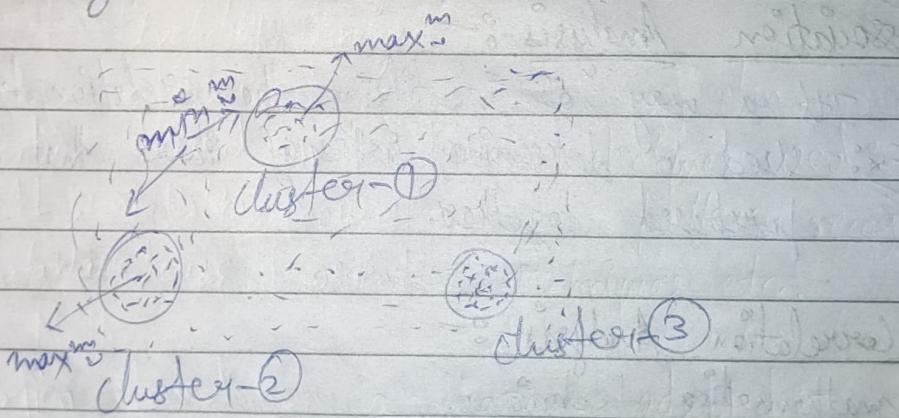
done based on Prev. data.

e.g. 2, 3, 4, 5, 6, 7

Predicted data = 6

## (4) cluster Analysis :-

The data items are clustered based on the principle of maximising the intra-class similarity and minimising the inter-class similarity.



⇒ Analysis of these clusters are cluster Analysis.

## (5) Outlier Analysis :- (anomaly mining)

e.g. {2, 4, 6, 7, 8, 10, 12, 3}

↑ outlier (not even)

Among the data items in a db, there may be some items which do not follow the general behaviour of data. Those data items → outliers (noise/exception).

## ⑦ #3. Interestingness of Patterns - What are Interesting Patterns?

### \* Interestingness of patterns :-

In a data mining system, everyday millions of data patterns are generated. Among all these patterns generated, how many are really interesting?

(means useful for user)

Actually, a small fraction of patterns generated would be of interest to any given user. This raises 3 questions :-

① What makes patterns interesting? :-

- a pattern is interesting if it's easy to understand by people
- valid on new/old test data
- potentially useful.

② Can data mining system generate all of the interesting patterns?

→ refers to completeness of a data mining (dm) system

In reality, it is not possible for a dm system to generate all interesting patterns.

③ Can dm systems generate only interesting patterns?

→ refers to optimisation of a dm system

Remember, generating only interesting patterns  $\Rightarrow$  challenging.

If only interesting patterns are generated, it becomes easy and efficient for the user.

(time is saved).

## ⑧ #4. Classification of Data Mining Systems :-

Why classification needed?

dm  $\rightarrow$  everywhere & anywhere (so that we can give what user actually wants/heeds)

Data Mining (dm) systems are classified based on several criteria :-

# ① Classification based on mined databases :-

based on type of database that is been mined.

- Relational
- Transactional
- Object-relational
- Datawarehouse

# ② Classification based on type of knowledge mined :-

- characterization
- Dissemination
- Association & Correlation Analysis
- classification.
- prediction
- outlier Analysis
- Evolution Analysis

# ③ Classification based on kinds of techniques used :-

- ML, statistics, neural networks, pattern recognition, data warehouse oriented techniques etc.

# ④ Classification based on applications adapted :-

- Finance
- Telecommunications
- DNA
- Stock Markets
- Email
- etc

Apart from above 4 classifications we have a general classification

DB technology

Statistics

Data Visualization

Data Mining

Info. Science

other disciplines

④ #5. Data Mining Task Primitives :-  
 A datamining task is represented in form of a dm query, and is defined in terms of dm task primitives.

this allows the user to interactively communicate with the dm system.

5 dm task primitives :-

also called as relevant attributes

- ① Set of task relevant data to be mined. (e.g. supermarket  $\rightarrow$  fruits)
- ② Specifies the kind of knowledge to be mined. (i.e., functionalities)
- ③ The background knowledge to be used in discovering discovery
- ④ The interestingness measures and thresholds for pattern evaluation
- ⑤ The expected representation for visualizing the discovery patterns to user.

whether in terms of rules / tables / patterns / charts / graph

⑤ #6. Integration of Data Mining System with A Database or DataWarehouse System

Integration / association / combining

If there is no integration — no communication with db.

dm  $\oplus$  db/dw

$\Rightarrow$  We have a total of 4 ~~int~~ integration schemes :-

① No Coupling :-

Dm system will not use any information function. i.e., there is no communication with db.

for this it comm.  $\uparrow$  with other storage methods.  
 (like file system)

## ② Loose Coupling :- (≈ 15%)

- will use some of the functionalities (only upto some extent)
- better than no coupling.
- suitable for small data sets.

## ③ Semi-Tight Coupling :- (> 50%)

- linked to the db.
- also, some of the dm primitives are also implemented in db.

## ④ Tight Coupling :- (≈ 100%)

- dm system is completely linked to db.
- most efficient among all.

The db system is fully integrated in such a way that it becomes part of the dm system.

- efficient & optimised implementation of dm.

## # 7. Major Issues in Data Mining :-

### → Mining Methodology & User-Interface/Interaction issues :-

- Mining diff. kinds of knowledge from databases.
- interactive mining of knowledge at multiple levels of abstraction.
- Handling Noisy/Incomplete data.
- Pattern evaluation.

### → Performance Issues :-

- efficiency & scalability of dm algorithms.
- Parallel, distributed and incremental mining algorithms.

## → Diverse Data Types Issues :-

- Handling of relational and complex types of data.
- Mining information from heterogeneous db's and global info. systems.

## ⑦ #8. Data Preprocessing in Data Mining :- (-4 steps)

→ The process of transforming raw data into an understandable format.

→ 4 major tasks :-

- ① Data Cleaning
- ② Data Integration
- ③ Data Reduction
- ④ Data Transformation

\* Data Cleaning :- cleaning of missing data noisy data

Process of removal of incorrect, incomplete, inaccurate data, also replaces missing data.

### ① Handling missing values :-

In place of missing values, we can replace with "NA"

with mean value (Normal data)

with median value (non-normal data)

- sometimes replaced with most probable values.
- missing values can be filled in 2 ways.

manual

small

automate

more efficient

large datasets

② Handling noisy data :- ~~most useful~~ ~~error~~ ~~data~~ ~~available~~  
noisy data  $\rightarrow$  inconsistent/error data.

methods to handle noisy data :-

① Binning :-

(data p-) firstly, data is sorted. Then sorted data is stored in bins.

3 methods to handle data in bins :-

— smoothing by bin mean ( $\frac{a+b+c+d}{4}$ )

— smoothing by bin median (1, 2, 3, 4, 5)

— smoothing by bin boundary. (min & max)

② Regression :-

Numerical prediction of data.

③ clustering :-

— similar data items are grouped at one place.

— dissimilar items — outside the cluster.

\* Data Integration :-

Multiple heterogeneous sources of data are combined into single dataset.

There are 2 types of data integration :-

① Tight Coupling :-

Data is combined together into a physical location.

(A) + (B)  $\rightarrow$  C

$\leftarrow \times$  (no return)

② Loose Coupling :-

Only an interface is created and data is combined through the i/f (interface) and also accessed through o/f.

— Data remains in actual database only.

## \* Data Reduction :-

Volume of data is reduced to make analysis easier.

### methods for data reduction :-

#### ① Dimensionality reduction :-

reduces no. of i/p variables in the dataset,  
because large i/p variables  $\rightarrow$  poor performance.

#### ② Data Cube aggregation :-

Data is combined to construct a database

(Redundant, noisy data is removed)

(means duplicate records will be removed)

#### ③ Attribute Subset Selection :- (in Table columns, ~~rows~~)

Highly relevant attributes should be used, others should be discarded/removed.

#### ④ Numerosity Reduction :-

Here, we store only model of data instead of entire data. sample

## \* Data Transformation :-

Data is transformed into an appropriate form suitable for ~~mining~~ process.

### 4 methods for data transformation :-

① Normalisation:- Done in order to scale the data values in specified range (-1.0 to 1.0 or from 0 to 1)

#### ② Attribute Selection:-

New attributes are created using older ones.

### ③ Discretization :-

Raw values are replaced by interval levels.

### ④ Concept Hierarchy generation :-

Attributes are converted from low level to high level.

Ex:- city → Country

Low level      High level

## ⑤ #9. Frequent Patterns - Example, Market Basket Analysis :-

### \* Frequent Patterns :-

The patterns that appear frequently in a dataset.

↓  
(include frequent data items, sequences, substructures)

eg. computer      computer, mouse,      graph  
computer, mouse,      tree  
Keyboard

e.g. Milk and Bread

both bought together → frequent pattern

### \* Market-Basket-Analysis :-

Process of analysing customer buying habits, by finding the associations b/w the different items that a customer will place in their baskets.

— mainly useful for sellers.

### → Strategies used :-

1. Placing them together

2. Placing them at 2-diff. ends.

— This analysis will help sellers to plan their shelf space for increased sales.

— frequent patterns are represented by association rules.

Ex:- Computer and Anti-virus.

computer  $\Rightarrow$  anti-virus software  
purchased purchased

[support = 2%, confidence = 6%]

out of all purchases 2% bought computer & anti-virus together.

$\rightarrow$  Support :-

Identifies how frequently a rule is applied to given dataset.

$$S(P \rightarrow Q) = \frac{\sigma(P \cup Q)}{N}, N = \text{total transactions}$$

$\rightarrow$  Confidence :-

It defines frequent occurrence of items of Q in transactions of P.

$$C(P \rightarrow Q) = P(Q|P)$$

P tends to Q

## #10. Mining Methods - APRIORI algorithm with Example :-

\* Mining Methods :-

$\rightarrow$  Apriori algorithm

$\rightarrow$  FP Growth Algorithm

$\rightarrow$  Apriori Algorithm :- Given by R. Agrawal & R. Srikant  
Shows how objects are associated with each other

objective :-

To generate an association.

$$\{(2,3), (2,5), (5,2), (5,4)\} = \text{length 2}$$

Ex8Min<sup>m</sup> support = 50%Threshold confidence = 70%  
(or min<sup>m</sup>)

| TID | Items   |
|-----|---------|
| 100 | 1 3 4   |
| 200 | 2 3 5   |
| 300 | 1 2 3 5 |
| 400 | 2 5     |

C1

| Itemset   | Support | min. support                             |
|-----------|---------|--|
| {1,3}     | 2       | $2/4 = 50\%$                             |
| {2,3}     | 3       | $3/4 = 75\%$                             |
| {3,5}     | 3       | $3/4 = 75\%$                             |
| {4,5}     | 1       | $1/4 = 25\% < 50\%$                      |
| {1,2,3,5} | 3       | $3/4 = 75\% \quad (\text{min. support})$ |

Itemset (new) = {1, 2, 3, 5}  $\Rightarrow$  form Pairs: (1, 2), (1, 3), (1, 5), (2, 3), (2, 5), (3, 5)C2

| Itemset | Support | min. support |
|---------|---------|--------------|
| (1, 2)  | 1       | $1/4 = 25\%$ |
| (1, 3)  | 2       | $2/4 = 50\%$ |
| (1, 5)  | 1       | $1/4 = 25\%$ |
| (2, 3)  | 2       | $2/4 = 50\%$ |
| (2, 5)  | 3       | $3/4 = 75\%$ |
| (3, 5)  | 2       | $2/4 = 50\%$ |

new itemset = {(1, 3), (2, 3), (2, 5), (3, 5)}

$\Rightarrow$  form triplets :-  $(1, 2, 3)$ ,  $(1, 3, 5)$ ,  $(2, 3, 5)$

| <u>C3</u>   | itemset | support      | min support |
|-------------|---------|--------------|-------------|
| $(1, 2, 3)$ | 1       | $1/4 = 25\%$ |             |
| $(1, 3, 5)$ | 1       | $1/4 = 25\%$ |             |
| $(2, 3, 5)$ | 2       | $2/4 = 50\%$ |             |

new itemset =  $\{(2, 3, 5)\}$

$\Rightarrow$  Now let's calculate support & confidence.

$$\boxed{\text{Confidence} = \frac{\text{Support}(A \cup B)}{\text{Support}(A)}}.$$

Using  $(2, 3, 5)$  we can generate association rules.

| <u>Rules</u>                 | <u>Support</u> | <u>Confidence</u>   |
|------------------------------|----------------|---|
| $(2 \wedge 3) \rightarrow 5$ | 2              | $2/2 = 100\%$ <small>strong &amp; perfect association rules</small> |
| $(3 \wedge 5) \rightarrow 2$ | 2              | $2/2 = 100\%$   |
| $(2 \wedge 5) \rightarrow 3$ | 2              | $2/3 = 66\%$ <small>among all.</small>                              |
| $2 \rightarrow (3 \wedge 5)$ | 2              | $2/3 = 66\%$  |
| $5 \rightarrow (2 \wedge 3)$ | 2              | $2/3 = 66\%$  |
| $3 \rightarrow (2 \wedge 5)$ | 2              | $2/3 = 66\%$  |

$\bullet$   $\frac{(2 \wedge 3) \rightarrow 5}{A \quad B}$  confidence =  $\frac{S((2 \wedge 3) \cup 5)}{S(2 \wedge 3)} = \frac{2}{2} = 100\%$

$\bullet$   $\frac{2 \rightarrow (3 \wedge 5)}{A \quad B}$  confidence =  $\frac{S(2 \cup (3 \wedge 5))}{S(2)} = \frac{2}{3} = 66\%$

∴  $(2^3) \rightarrow 5$  and  $(3^5) \rightarrow 2$  are association rules.

## #11. Mining Methods - FP growth algorithm with examples

### \* FP Growth Algorithm :-

Frequent Pattern

→ it's an efficient & scalable method for mining the complete set of FP using a tree structure for storing information about FP called FP tree.

ex :- min. support = 30%

| Trans id | Items         | all items are A, B, C, D, E |
|----------|---------------|-----------------------------|
| 1        | E, A, D, B    |                             |
| 2        | D, A, E, C, B |                             |
| 3        | C, A, B, E    | (2, 2, 2) (prior)           |
| 4        | B, A, D       |                             |
| 5        | D             |                             |
| 6        | D, B          |                             |
| 7        | A, D, E       |                             |
| 8        | B, C          |                             |

To write priorities :-

- more freq. → more priority
- same freq. → FCFS priority

- list out all the priorities

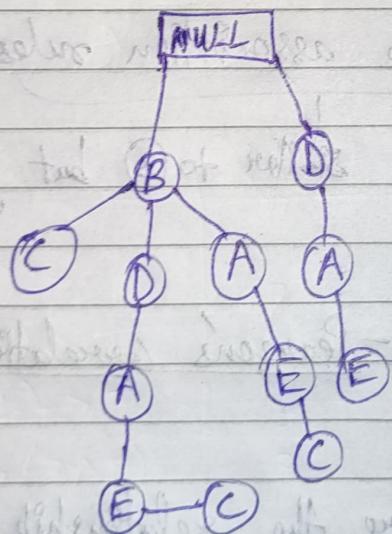
| Itemset | Freq. | Priority |
|---------|-------|----------|
| A       | 5     | 3        |
| B       | 6     | 1        |
| C       | 3     | 5        |
| D       | 6     | 2        |
| E       | 4     | 4        |

B D A E C

— order items according to their priority

| Trans ID | Items     | Ordered Items |
|----------|-----------|---------------|
| 1        | E A D B   | B D A E       |
| 2        | D A E C B | B D A E C     |
| 3        | C A B E   | B A E C       |
| 4        | B A D     | B D A         |
| 5        | D         | D             |
| 6        | D B       | B D           |
| 7        | A D E     | D A E         |
| 8        | B C       | B C           |

— construct the FP tree



B - 1, 2, 3, 4, 5, 6  
 D - 1, 2, 3, 4  
 A - 1, 2, 3  
 E - 1, 2  
 C - 1

A - 1

E - 1

C - 1

D - 1, 2

A - 1

E - 1

C - 1

$$\text{freq of } D = 4 + 2 \\ = 6$$

(you can check in previous page)  
 freq of D is same or not.  
 for correctness.

## #12. Mining ~~Various~~ various kinds of Association Rules :-

### \* Mining Various kinds of Association Rules :- (4 types)

#### ① Mining Multi-level Association rules :-

- using uniform support for all levels.
- using reduced minimum support at lower levels.
- using item or group based minimum support.

#### ② Mining ~~multidimensional~~ Association Rules from relational-db or relational dw.

each attribute considered as a predicate  
e.g. - buy(x), sell(x) etc.

#### ③ mining multi-dimensional association rules using static ~~discretization~~ discretisation of quantitative attributes.

numerical

#### ④ mining quantitative association rules.

intervals  
manually

similar to ③ but by dynamic,  
using binning

## #13. Correlation Analysis - Pearson's Correlation Coefficient :-

### \* Correlation Analysis :-

↳ used to measure the relationship b/w two variables.

$$\gamma_{A,B} = \frac{\sum (A - A')(B - B')}{(n-1) \cdot \sigma_A \cdot \sigma_B}$$

gamma

$A', B'$  = mean of A and B.  
 $\sigma_A, \sigma_B$  = standard deviation of A & B.

Karle Pearson Correlation Coefficient

$A', B'$  = mean of A and B.

$\sigma_A, \sigma_B$  = standard deviation of A & B.

$\uparrow \uparrow \rightarrow$  the correlation  
 $\uparrow \downarrow \rightarrow$  - no correlation  
 $\downarrow \uparrow \rightarrow$  or  $\downarrow \downarrow \rightarrow$  no correlation

↳ No correlation  
if no proportion  
is found

$r$  can have  
3 variables  $(0, -1, +1)$   
(or values)

$r = +1 \rightarrow$  Perfect positive correlation

$r = 0 \rightarrow$  No correlation (no dependence)

$r = -1 \rightarrow$  Perfect negative correlation

Ex:-