

Data Mining

SAQ's :-

- Q) Give an example of Data Matrix and Dissimilarity Matrix?

Ans Data Matrix \Rightarrow are used to build the evolutionary trees.

For ex \Rightarrow Datamatrix code - encoding the text "wikipedia" use the color to show "datagreen", "padding yellow" "error correction red", finder & Timing Agenda "an unused orange".

\Rightarrow Data matrix symbols are made up modules arranged within a nemometer finder & timing pattern.

* Dissimilarity matrix \Rightarrow is a matrix that will express the similarity pair to pair btw two sets. It is square and symmetric in nature. \Rightarrow The diagonal members are defined as "0" meaning that the 0 is measure of dissimilarity between an element and itself.

- Q) Define Data Mining and its functionalities.

Ans Data Mining \Rightarrow The process of automatically discovering patterns, relationships, and insights from large datasets, using various techniques and tools.

Ex: Imagine a supermarket analyzing its sales data to identify patterns and relationships between products.

FUNCTIONALITIES \Rightarrow .

\Rightarrow Data can be associated with classes and concept
 \Rightarrow Association Analysis.

- ⇒ classification and prediction of data
- ⇒ clustering analysis.
- ⇒ evolution and derivation analysis.

3Q) Enumerate different technologies Used in Data Mining?

Technology Used :-

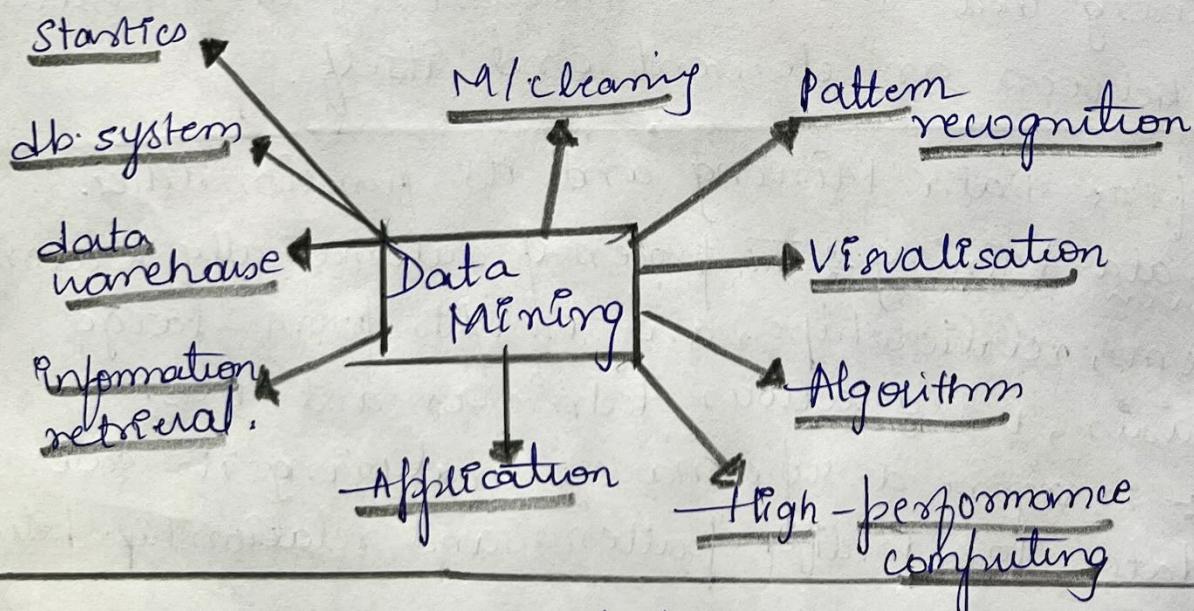
- 1) ⇒ classification according to the kinds of database mined.
- 2) ⇒ classification according to the kinds of knowledge "
- 3) ⇒ classification according to the kind of technique utilize.

* Techniques + sophisticated data mining system will often adopt multiple data mining technique or workout an objective, integrated technique which combines the merits of a few individual approaches.

Techniques Used :-

- Statistics
- Machine learning

- Database system & Data warehouse
- Information Retrieval.



4Q) Calculate the mode of following data and its modality? 2, 4, 3, 8, 12, 14, 12, 13, 15, 20, 22.

Ans . Mode:- The most common / repeated ^{most} value in a dataset.
∴ Mode = 12. Modality :- unimodal.

LAlg's $\frac{1}{2}$

- Q) Suppose we have the following 2-D dataset.
 Consider the data as 2-data points. Given a new data point, $n = (1.4, 1.6)$ as a query, rank the database point based on similarity with the query using i) Euclidean Distance ii) Manhattan Distance iii) Supremum Distance & iv) Cosine Similarity

	A ₁	A ₂
x ₁	1.5	1.7
x ₂	2	1.9
x ₃	1.6	1.8
x ₄	1.2	1.5
x ₅	1.5	1.0

Sol. Considering the following data :-

A (x ₁)	B (y ₁)	Euclidean distance	Manhattan Distance	Minkowski distance	Supremum Distance	Cosine similarity
1.5	1.7	0.1414	0.2	0.1414	0.1	0.99999
2	1.9	0.6708	0.9	0.6708	0.6	0.99575
1.6	1.8	0.2828	0.4	0.2828	0.2	0.99997
1.2	1.5	0.2236	0.3	0.2236	0.2	0.99903
1.5	1.0	0.6083	0.7	0.6083	0.6	0.96536

let P₁ be a point with X₁, Y₁

P₂ be a point with X₂, Y₂

$$X = \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = (1.4, 1.6)$$

Let P₁ (X₁, Y₁)

P₂ (X₂, Y₂).

a) Euclidean distance $\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

\Leftrightarrow $\star [(x_2, y_2) = (1.4, 1.6)]$ — Given.

ED of $(1.5, 1.7) = \sqrt{(1.4 - 1.5)^2 + (1.6 - 1.7)^2} = 0.1414$

ED of $(2, 1.9) = \sqrt{(1.4 - 2)^2 + (1.6 - 1.9)^2} = 0.6708$

ED of $(1.6, 1.8) = \sqrt{(1.4 - 1.6)^2 + (1.6 - 1.8)^2} = 0.2828$

ED of $(1.2, 1.5) = \sqrt{(1.4 - 1.2)^2 + (1.6 - 1.5)^2} = 0.2236$

ED of $(1.5, 1.0) = \sqrt{(1.4 - 1.5)^2 + (1.6 - 1.0)^2} = 0.6083$

b) Manhattan distance \Rightarrow

$$= |x_2 - x_1| + |y_2 - y_1|$$

MD of $(1.5, 1.7) = 0.2$

MD of $(2, 1.9) = 0.9$

MD of $(1.6, 1.8) = 0.4$

MD of $(1.2, 1.5) = 0.3$

MD of $(1.5, 1.0) = 0.7$

c) Minkowski distance:

$$\Rightarrow h \sqrt{(x_2 - x_1)^h + (y_2 - y_1)^h} \quad h \rightarrow \text{height}$$

(get cancel) $\Rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad m = (1.4, 1.6)$

\Downarrow ~~so~~ equals to euclidean distance.

d) supremum distance

$$\Rightarrow \max(|x_2 - x_1|, |y_2 - y_1|)$$

SD of $(1.5, 1.7) = \max(|1.4 - 1.5|, |1.6 - 1.7|) = 0.1$

SD of $(2, 1.9) = \max(|1.4 - 2|, |1.6 - 1.9|) = 0.6$

SD of $(1.6, 1.8) = \max(|1.4 - 1.6|, |1.6 - 1.8|) = 0.2$

$$SD \text{ of } (2, 1.5) = \max(11.4 - 1.2), 1(1.6 - 1.5) | = 0.2$$

$$SD \text{ of } (1.5, 1.0) = \max(11.4 - 1.5), 1(1.6 - 1.0) | = 0.6$$

e) Cosine similarity :-

$$\frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum A_i B_i}{\sqrt{\sum A_i^2} \sqrt{\sum B_i^2}} \quad \boxed{\begin{array}{l} (x_1, y_1) \\ = (1.4, 1.6) \end{array}}$$

$\left(\because \sum A_i^2 = \sum B_i^2 \right)$

$$\text{i)} \frac{(1.4 * 1.5) + (1.6 * 1.7)}{\sqrt{(1.5)^2 + (1.7)^2} * \sqrt{(1.4)^2 + (1.6)^2}} = 0.999999$$

$$\text{ii)} \frac{(1.4 * 2.0) + (1.6 * 1.9)}{\sqrt{(2.0)^2 + (1.9)^2} * \sqrt{(1.4)^2 + (1.6)^2}} = 0.99575.$$

$$\text{iii)} \frac{(1.4 * 1.6) + (1.6 * 1.8)}{\sqrt{(1.6)^2 + (1.8)^2} * \sqrt{(1.4)^2 + (1.6)^2}} = 0.99997$$

$$\text{iv)} \frac{(1.4 * 2.0) + (1.6 * 1.5)}{\sqrt{(2.0)^2 + (1.5)^2} * \sqrt{(1.4)^2 + (1.6)^2}} = 0.99903$$

$$\text{v)} \frac{(1.4 * 1.5) + (1.6 * 1.0)}{\sqrt{(1.5)^2 + (1.0)^2} * \sqrt{(1.4)^2 + (1.6)^2}} = 0.96536$$

Elucidate about Decision tree Induction Algorithm with an example?

(6)
9)
Ans

DECISION TREE INDUCTION

- Flowchart like tree structure
- supports in taking decisions as it classifies the data.
- It defines the rules visually in form of tree.

Types of nodes

- 1) Root node - Main question
- 2) Branch node - Intermediate processing node
- 3) leaf node - Answer

* Attribute Selection Measures :-

↳ Information Gain :-

How much it does the answer to the specific question provide.

ii) entropy :-

Measures the amount of uncertainty in the info
As $IG \uparrow$ Entropy \downarrow .

Ex:- Credit Score Rating

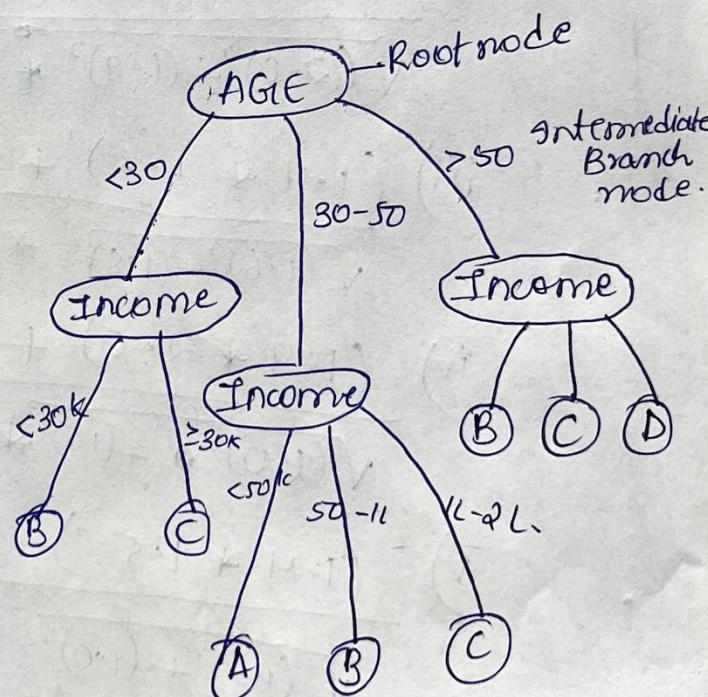
A - Avg

B - Bad

C - Good

D - excellent

} classification



* Rules that can be defined are :-

→ if age < 30 , income $< 30k$
the credit score = Bad.

→ if age < 30 , income $\geq 30k$
then credit score = Good.

Advantages of Decision Tree :-

→ easy to understand and map.

nicely to a set of production rules.

- successfully can be applied to a real problems.
- No prior assumptions about the nature of the data.
- Decision Tree are able to build the models with datasets containing numerical as well as categorical data
- DT generation are of 2 types :-

i) Tree construction :-

- At start all the training examples are at the root
- Partitioning examples recursively based on selected attributes.

ii) Tree pruning :-

- Identify the root node branches that reflect noise and out layers.
- The use of this DT is it can classify the unknown samples and attribute values can be tested from the sample against the DT.

Q) Discuss how to improve efficiency of Apriori Algorithm?

Ans) → Improving the efficiency of Apriori :-
There are several variations for improving the efficiency of Apriori Algorithm.

- Hash-based technique
- Transaction reduction
- Partitioning
- Sampling
- Dynamic - Itemset counting

1) Hash-based Technique :-

- hashing of itemset into corresponding buckets.
- tech is used to reduce the size of candidate k.

c_k for $k > 1$

ex:- create hash table H_2 using hash function

$$h(n, y) = (\text{order of } n) \times n + (\text{order of } y)$$

th

bucket address	0	1	2	3	4	5	6
bucket count	2	2	4	2	2	4	4
bucket content	I_1, I_9	I_1, I_5	\vdots	$\{\}$	\vdots	\vdots	\vdots

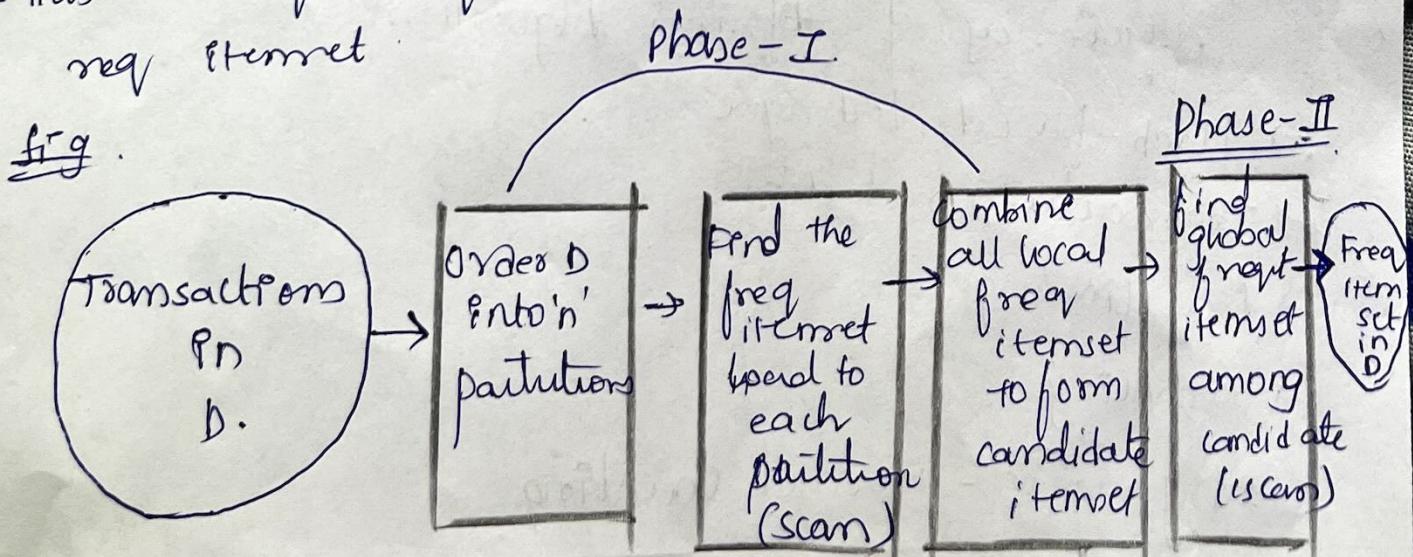
2) Transaction Reduction :-

- reducing the no. of transactions scanned in future iterations.
- At T^* that does not contain any frequent itemset cannot contain any freq. ($k+1$) itemsets
- Therefore such T^* can be removed ($f < k$).

3) Partitioning :-

Partition the data to find candidate itemsets.

- This technique requires 2 db to be scan to mine freq itemset



Using by Partitioning the Data.

Phase I :-

Using Algor $\Rightarrow D \xrightarrow{n}$ have
 non overlap partitions $\rightarrow \min_sup \times \text{no. of } T_p$
~~no. of T_p~~ \times ~~no. of T_p~~

- load freq itemset \rightarrow itemset freq in that partition
 may or may not be freq wrt the entire db 'D'.
- ∴ Any itemset that is potentially freq wrt to D must occur as a freq itemset in at least one of the partitions.

Phase-II :-

This is the second scan of D is conducted in which actual support of each candidate is assessed to determine the global freq itemset.

- Partition size and no. of partitions are set so that each partition can fit into main memory & \therefore be read only once in each phase.

4) Sampling :- mining on a subset of given data.

5) Dynamic itemset counting :-

- adding candidate itemset at different points during a scan
- partitioning into blocks starting by a starting point
- In this, new candidate itemsets can be added at any start point, unlike in Apriori, which determines new candidate itemset only immediately before each complete db scan.

→ This technique uses count so far as the lowerbound of actual count. And if the count passes min-support then it is added to freq itemset collection & can be used to generate longer candidates. This lead to fewer db scan than with Apriori for finding the freq itemsets.

8Q) Explain how to mine closed frequent itemsets?

Ans → Closed Frequent Itemset Mining (CFIM) is a technique used to discover sets of items that frequently appear together in a dataset.

MAIN STEPS :-

i) Frequent itemset Generation :-
→ Use Apriori algorithm or similar to generate frequent itemsets (F1)
Ex :- $\{A, B\}$ (support = 4), $\{B, C\}$ (support = 3),
 $\{A, C\}$ (support = 2).

ii) Closed itemset Generation :-
→ Use PRUNE to generate closed itemsets (C1)
from F1
Ex :- $\{A, B, C\}$ (support = 2) is a closed itemset.

iii) Item Merging :-
→ combine itemsets with the same support into a single itemset
Ex :- $\{A, B\}$ and $\{B, C\}$ are merged.
into $\{A, B, C\}$ (support = 3).

iv) Subitemset after skipping pruning :-

→ Remove itemsets that are subsets of other itemsets with the same support.

Ex :- $\{A, B\}$ is removed because it's a subset of $\{A, B, C\}$ with the same support.

v) CDFS (closed frequent itemset dataset) Mining :-

- Output :- closed frequent itemsets (CFI) with their support values.

- Ex :- $\{A, B, C\}$ (support = 2) is a closed frequent itemset.

Ex :- Dataset :-

transaction ID	Items
01	A, B, C
02	A, B, D
03	B, C, F
04	A, C, F
05	B, C, G

⇒ By applying CFIM2 to this dataset, we can discover closed frequent itemsets like $\{A, B, C\}$ (support = 2), which indicates that items A, B, and C often appear together in transactions.

(a) Suppose that the data for analysis includes the attributes age (are in increasing order) then

22, 25, 17, 19, 33, 64, 23, 17, 20, 18.

- Find 1) Mean 2) Median 3) Mode / Modality
4) 5 number summary 5) Box plot
6) Scatter plot between two variables.

Given data :-

22, 25, 17, 19, 33, 64, 23, 17, 20, 18.

Step 1 :- Arrange them in their increasing order

17, 17, 18, 19, 20, 22, 23, 25, 33, 64.

i) Mean = $\frac{\text{sum of observations}}{\text{total no. of observations}}$

$$\bar{x} = \Sigma x / n$$

$$\Rightarrow \frac{17 + 17 + 18 + 19 + 20 + 22 + 23 + 25 + 33 + 64}{10}$$

$$= \frac{228}{10} = 22.8$$

$$\therefore \text{Mean } (\bar{x}) = 22.8$$

ii) Median

17, 17, 18, 19, 20 | 22, 23, 25, 33, 64.
middle terms.

$$\text{Median} = \frac{20 + 22}{2} = \frac{42}{2} = 21 = Q_2$$

iii) min = 17 $Q_1 \Rightarrow 18$
 $Q_2 \Rightarrow 21$ $Q_3 \Rightarrow 25$ \rightarrow 5-Number Summary

max = 33.

iv) Mode :-
the most common / most repeated value
in a dataset.

\therefore 17 has occurred 2 times

$$\therefore \text{Mode} = 17$$

Modality :- Unimodal.

check outlier

$$1) Q_3 + (1.5 * IQR)$$

higher outlier

$$\begin{aligned}IQR &= Q_3 - Q_1 \\&= 25 - 18 = 7\end{aligned}$$

$$\begin{aligned}\Rightarrow 25 + (1.5 * 7) \\&= 25 + 10.5 \\&= 35.5.\end{aligned}$$

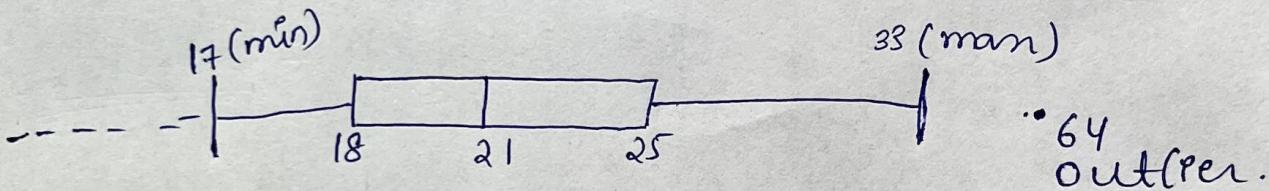
lower outlier

$$\begin{aligned}Q_1 - (1.5 * IQR) \\&\Rightarrow 18 - (1.5 * 7)\end{aligned}$$

$$\Rightarrow 7.5.$$

64-outlier.

higher outlier $\rightarrow 35.5$.



Bon plot.

vi) Scatter plot btw two variable.

leads needs.
B.