

Reg No.: _____

Name: _____

APJ ABDUL KALAM TECHNOLOGICAL UNIVERSITY
EIGHTH SEMESTER B.TECH DEGREE EXAMINATION(S), OCTOBER 2019

Course Code: CS402
Course Name: DATA MINING AND WAREHOUSING

Max. Marks: 100

Duration: 3 Hours

PART A*Answer all questions, each carries 4 marks.*

Marks

- | | | |
|----|--|-----|
| 1 | How is data warehouse different from a database? How are they similar? | (4) |
| 2 | Compare star and snowflake schema dimension table. | (4) |
| 3 | Use the two methods below to normalize the following group of data:
100,200,300,500,900
i) min-max normalization by setting min=0 and max=1
ii) z-score normalization | (4) |
| 4 | Explain the attribute selection method in decision trees . | (4) |
| 5 | Distinguish between hold out method and cross validation method. | (4) |
| 6 | Explain prepruning and postpruning approaches in decision tree algorithm. | (4) |
| 7 | Differentiate between support and confidence. | (4) |
| 8 | How to compute the dissimilarity between objects described by binary variables? | (4) |
| 9 | Differentiate between Agglomerative and Divisive hierarchical clustering method. | (4) |
| 10 | Explain web content mining? | (4) |

PART B*Answer any two full questions, each carries 9 marks.*

- | | | |
|----|---|-------------------|
| 11 | The following data is given in increasing order for the attribute age:
13,15,16,16,19,20,20,21,22,22,25,25,25,25,30,33,33,35,35,35,36,40,45,46,52,70.
a) Use smoothing by bin boundaries to smooth these data, using bin depth of 3.
b) How might you determine outliers in the data?
c) What other methods are there for data smoothing? | (3)
(3)
(3) |
| 12 | Explain the following procedures for attribute subset selection
a) Stepwise forward selection
b) Stepwise backward elimination
c) A combination of forward selection and backward elimination | (3)
(3)
(3) |

- 13 a) Suppose a datawarehouse consists of three measures customer, account and branch and two measures count (number of customers in the branch) and balance. Draw the schema diagram using snowflake schema. (4)
- b) Real-world data tend to be incomplete, noisy, and inconsistent. What are the various approaches adopted to clean the data? (5)

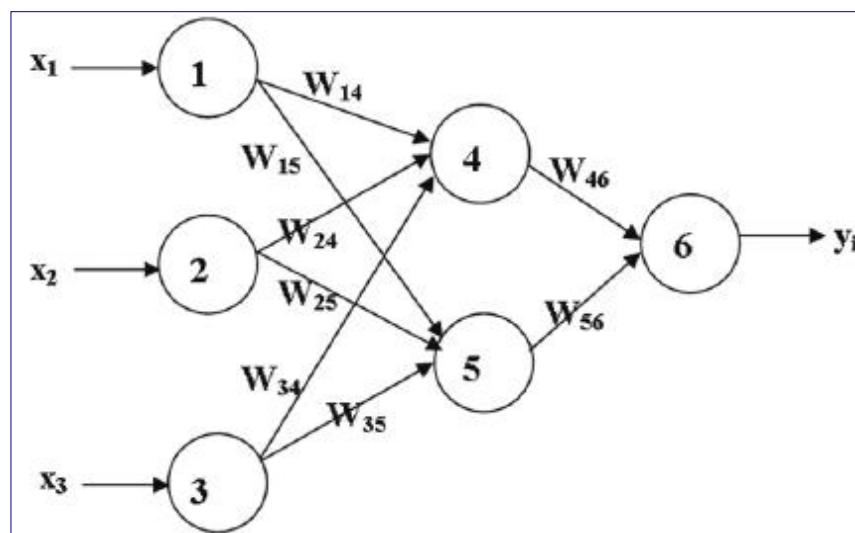
PART C

Answer any two full questions, each carries 9 marks.

- 14 Given the following data on a certain set of patients seen by a doctor, can the doctor conclude that a person having chills, fever, mild headache and without running nose has the flu?(Use Naive Bayes algorithm for prediction) (9)

chills	running nose	headache	fever	has flu
Y	N	mild	Y	N
Y	Y	no	N	Y
Y	N	strong	Y	Y
N	Y	mild	Y	Y
N	N	no	N	N
N	Y	strong	Y	Y
N	Y	strong	N	N
Y	Y	mild	Y	Y

- 15 The following figure shows a multilayer feed-forward neural network. Let the learning rate be 0.9. The initial weight and bias values of the network is given in the table below. The activation function used is the sigmoid function. (9)



X ₁	X ₂	X ₃	W ₁₄	W ₁₅	W ₂₄	W ₂₅	W ₃₄	W ₃₅	W ₄₆	W ₅₆	θ ₄	θ ₅	θ ₆
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

Show weight and bias updation with the first training sample (1,0,1) with class label 1, using backpropagation algorithm

- 16 a) Explain classification by C4.5 algorithm. (6)
 b) What is meant by Maximum Marginal Hyperplane (MMH)? (3)

PART D

Answer any two full questions, each carries 12 marks.

- 17 Consider the transaction database given below. Set minimum support count as 2 and minimum confidence threshold as 70%

Transaction ID	List of Item_Ids
T100	I1,I2,I5
T200	I2,I4
T300	I2,I3
T400	I1,I2,I4
T500	I1,I3
T600	I2,I3
T700	I1,I3
T800	I1,I2,I3,I5
T900	I1,I2,I3

- a) Find the frequent itemset using Apriori Algorithm. (8)
 b) Generate strong association rules . (4)
- 18 a) Explain DBSCAN algorithm . (8)
 b) State the pros and cons of DBSCAN method. (4)
- 19 a) Explain clustering by k-medoid algorithm. (6)
 b) Explain Apriori based frequent subgraph mining. (6)

15/11

MAY 2019

PART A

1. What is datamining related to business intelligence? (Mod 1)

Datamining is the process of converting raw data into useful information. It is an important analytics process designed to explore data. Data mining is also referred to as data/knowledge discovery. It is about understanding the data you already have to make better business decisions.

Business Intelligence as a discipline is made up of several related activities, including data mining, online analytical processing, querying and reporting. It often relies on data mining to look for solutions to problems in the data already available. It can provide insights of things that you did not know. / it can confirm/deny a hypothesis you may have. Either way it allows you to make business decisions.

16 Q. Differentiate b/w OLTP and OLAP. (Mod 1)

Basis for comparison	OLTP	OLAP
Bank	It is an online transactional system & manages database modification	It is an online data retrieving & data analysis sys.
Focus	Insert, Delete, Update info from the database	Extract data for analysing that helps in decision making
Data	OLTP and its transactions are the original source of data	Diff OLTPs database become the source of data for OLAP
Transaction	Short transaction	Long transactions
Time	Processing time of a transaction is comparatively less in OLTP	Processing time of a transaction is comparatively more in OLAP
Queries	Simple	Complex
Normalization	Tables in OLTP database are normalized (3NF)	Tables in OLAP database are normalized
Integrity	OLTP database must maintain data integrity constraint	OLAP database does not get frequently modified. Hence, data integrity is not affected.

2 3 Why do we need data transformation? What are the different ways of data transformation? (Mod II)

Data transformations (eg normalization) may be applied, where data are scaled to fall within a smaller range like 0-1. This can improve the accuracy and efficiency of mining algorithms involving distance measurements.

These techniques are not mutually exclusive, they may work together. For eg, data cleaning can involve transformations to correct wrong data such as by transforming all entries for a date field to a common format.

- * Smoothing :- Remove noise from data

- * Aggregation :- Summarization

- * Generalization :- Concept hierarchy climbing

- * Normalization :- Scaled to fall within a small, specified range

- min-max normalization

- Z-score normalization

- normalization by decimal scaling

- * Attribute / feature construction

- New attribute constructed from the given one.

→ Min-max normalization to $(\text{new}_{-}\text{min}_A, \text{new}_{-}\text{max}_A)$

$$v' = \frac{v - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new}_{-}\text{max}_A - \text{new}_{-}\text{min}_A) + \text{new}_{-}\text{min}_A$$

Eg:- Let income range \$ 12000 to \$ 98000 normalized to [0.0, 1.0]

Then \$ 73600 is mapped to $\frac{73600 - 12000}{98000 - 12000} (1-0) + 0$
= 0.716

→ Z-score normalization (μ : mean, σ : standard deviation) 3

$$v' = \frac{v - \mu_A}{\sigma_A}$$

Eg: Let $\mu = 54000$, $\sigma = 16000$. Then 73600 is mapped to
 $= \frac{73600 - 54000}{16000} = 1.225$

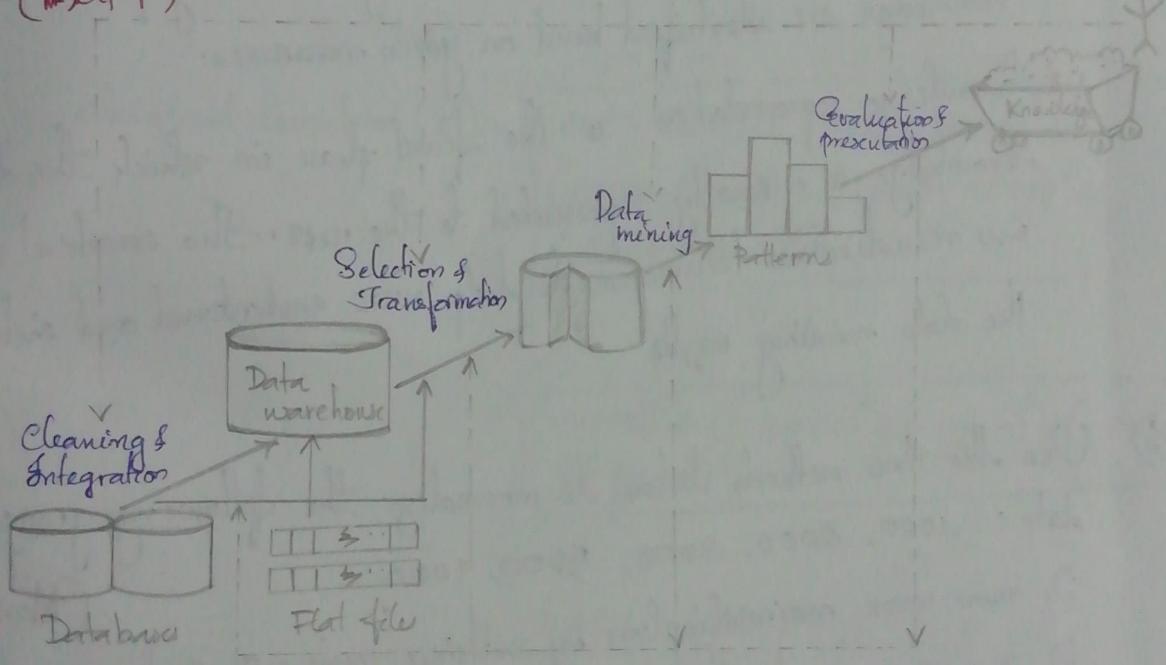
→ Normalization by decimal scaling.

$$v' = \frac{v}{10^j} \text{ where } j \text{ is the smallest integer such that } \max(|v'|) < 1$$

$j = \text{no. of digits in max no.}$

PART B

- II a) Explain various stages in knowledge discovery process with neat diagram
~~(Mod 1)~~



Knowledge Discovery in databases (KDD) is the process of finding useful information and patterns in data. Datamining is the use of algorithms to extract the information and patterns derived by the KDD process.

The iterative process consists of the following steps:

- * Data cleaning: also known as data cleaning, it is a phase ~~in~~

A) which noise data and irrelevant data are removed from the collection.

- * Data Integration:- at this stage, multiple data sources, often heterogeneous, may be combined in a common source.
- * Data selection :- At this step, the data relevant to the analysis is decided on and retrieved from the data collection.
- * Data transformation:- also known as data consolidation it is a phase in which the selected data is transformed into forms appropriate for the mining procedure.
- * Data mining :- is the crucial step in which clever techniques are applied to extract patterns potentially useful.
- * Pattern evaluation: in this step, strictly interesting patterns representing knowledge are identified based on given measures.
- * Knowledge representation:- is the final phase in which the discovered knowledge is visually represented to the user. This essential step uses visualization techniques to help users understand and interpret the data mining results.

- b) Use the two methods below to normalize the following group of data : 1000, 2000, 3000, 5000, 9000. (Mod II)
- i) min max normalization by setting $\text{min} = 0$ and $\text{max} = 1$
 - ii) Z-score normalization.

D) Min max normalization. $v' = \frac{v - \text{min}}{\text{max} - \text{min}} (\text{new_max} - \text{new_min}) + \text{new_min}$

$$\text{min} = 1000$$

$$\text{max} = 9000$$

$$\text{new_min} = 0$$

$$\text{new_max} = 1$$

$$v = 1000 \Rightarrow v' = \frac{1000 - 1000}{9000 - 1000} (1 - 0) + 0 = 0$$

$$v = 2000 \Rightarrow v' = \frac{2000 - 1000}{9000 - 1000} (1-0) + 0 = 0.125$$

$$v = 3000 \Rightarrow v' = \frac{3000 - 1000}{9000 - 1000} (1-0) + 0 = 0.25$$

$$v = 5000 \Rightarrow v' = \frac{5000 - 1000}{9000 - 1000} (1-0) + 0 = 0.5$$

$$v = 9000 \Rightarrow v' = \frac{9000 - 1000}{9000 - 1000} (1-0) + 0 - 1$$

Values	1000	2000	3000	5000	9000
Normalized	0	0.125	0.25	0.5	1

ii) Z-score normalization.

$$\text{Mean } (\mu) = \frac{1000 + 2000 + 3000 + 5000 + 9000}{5}$$

$$= 4000$$

$$\begin{aligned} \text{Standard deviation } (\sigma) &= \sqrt{\frac{(1000 - 4000)^2 + (2000 - 4000)^2 + (3000 - 4000)^2 + (5000 - 4000)^2 + (9000 - 4000)^2}{5-1}} \\ &= \sqrt{\frac{(-3000)^2 + (-2000)^2 + (-1000)^2 + (-1000)^2 + (5000)^2}{2}} \\ &= \frac{6324.5582}{2} \\ &= 3162.2776 \end{aligned}$$

$$\text{Z-score normalization } v' = \frac{v - \mu}{\sigma}$$

$$v = 1000 \Rightarrow v' = \frac{1000 - 4000}{3162.2776} = -0.9486$$

$$v = 2000 \Rightarrow v' = \frac{2000 - 4000}{3162.2776} = -0.6324$$

6

$$n = 3000 \Rightarrow n' = \frac{3000 - 4000}{3162 \cdot 2776} = -0.3162$$

$$n = 5000 \Rightarrow n' = \frac{5000 - 4000}{3162 \cdot 2776} = 0.3162$$

$$n = 9000 \Rightarrow n' = \frac{9000 - 4000}{3162 \cdot 2776} = 1.5811$$

Values	1000	2000	3000	5000	9000
Normalised	-0.9486	-0.6324	-0.3162	0.3162	1.5811

- 12 Suppose that a datawarehouse for University consists of four dimension date, spectator, location and game and two measures count and charge, where charge is the fare that a spectator pays when watching a game on the given date. Spectator may be students, adults or seniors, with each category having its own charge rate.

- Draw a star schema for the data warehouse
- Starting with the basic cuboid [date, spectator, location, game] what specific OLAP operation should be performed in order to list the total charge paid by student spectators at GM_PLACE in 2010.

(Mod II)

a) date-dimension table

date_id
day
month
quarter
year

Sales fact table.

date_id
spectator_id
game_id
location_id
count
charge

spectator dimension table

spectator_id
spectator_name
status
phone
address

game-dimension table

game_id
game_name
description
producer

location-dimension table

location_id
location_name
phone#
street
city
province
country

- b) → Roll up on date from date_id to year
- Roll up on game from game_id to all
- Roll up on location from location_id to location_name
- Roll up on spectator from spectator_id to status.
- Dice with status = "Students", location name = "GM place" and year = 2001.

13 Summarize the various preprocessing activities involved in data mining. (Mod II)

There are several data preprocessing techniques:

- * Data Cleaning can be applied to remove noise and correct inconsistencies in data
- * Data Integration merged data from multiple sources into a coherent data store such as a data warehouse
- * Data Reduction can reduce any data size by, for instance aggregating, eliminating redundant features or clustering.
- * Data Transformations (eg normalization) maybe applied, where data are scaled to fall within a smaller range like 0.0 to 1.0. This can improve the accuracy and efficiency of mining algorithms involving distance measurements

These techniques are not ~~and~~ mutually exclusive; they may work together. For eg, data cleaning can involve transformations to correct wrong data, such as by transforming all entries for a ~~date~~ date field to common format.

Data processing techniques, when applied before mining can substantially improve the overall quality of the patterns mined.

8 and for the time required for mining

OCTOBER 2019

PART A

1. How is datarehouse different from a database? How are they similar? (Mod 1)

- * Database is a collection of related data that represents some elements of the real world whereas Data warehouse is an information system that stores historical and cumulative data from single or multiple sources.
- * Database is designed to record data whereas the Data warehouse is designed to analyse data.
- * Database is application-oriented - collection of data whereas Data warehouse is the subject oriented collection of data.
- * Database uses Online Transactional Processing (OLTP) whereas Data warehouse uses Online Analytical Processing (OLAP).
- * Database tables and joins are complicated because they are normalized whereas Data Warehouse tables and joins are easy because they are denormalized.
- * ER modelling techniques are used for designing Database whereas Data modelling techniques are used for designing Data Warehouse.

2. Compute star and snowflake schema dimension table (Mod II)

STAR SCHEMA	SNOWFLAKE SCHEMA
* It consists of fact table with a single table for each dimension (dimension table)	* It is a variation of star schema, which have multiple level of dimension table.
* Highly denormalized.	* Normalized (in dimension table)
* Category wise single dimension table	* Dimension table further splits into additional table
* More data dependency and redundancy	* Less
* No need of complicated joins	* Require complicated joins

3. Use the two methods below to normalize the following group of data: 100, 200, 300, 500, 900

- i) min max normalization by setting min=0 and max=1
- ii) Z score normalization (Mod II)

i) Min max normalization. $v' = \frac{v - \text{min}}{\text{max} - \text{min}} (\text{new_max} - \text{new_min}) + \text{new_min}$

$$\text{min} = 100$$

$$\text{max} = 900$$

$$\text{new_min} = 900$$

$$\text{new_max} = 1$$

$$v = 100 \Rightarrow v' = \frac{100 - 100}{900 - 100} (1 - 0) + 0 = 0$$

$$v = 200 \Rightarrow v' = \frac{200 - 100}{900 - 100} (1 - 0) + 0 = 0.125$$

$$v = 300 \Rightarrow v' = \frac{300 - 100(1-\alpha) + \alpha}{900 - 100} = 0.25$$

$$v = 500 \Rightarrow v' = \frac{500 - 100(1-\alpha) + \alpha}{900 - 100} = 0.5$$

$$v = 900 \Rightarrow v' = \frac{900 - 100(1-\alpha) + \alpha}{900 - 100} = 1$$

Values	100	200	300	500	900
Normalized	0	0.125	0.25	0.5	1

ii) Z-score normalization

$$\text{Mean } (\mu) = \frac{100 + 200 + 300 + 500 + 900}{5} = 400$$

$$\begin{aligned} \text{Standard deviation } (\sigma) &= \sqrt{\frac{(100-400)^2 + (200-400)^2 + (300-400)^2 + (500-400)^2 + (900-400)^2}{5-1}} \\ &= \sqrt{\frac{(-300)^2 + (-200)^2 + (-100)^2 + (100)^2 + (500)^2}{2}} \\ &= \frac{632.455532}{2} \\ &= 316.227766 \end{aligned}$$

$$\text{Z-score normalization} = \frac{v - \mu}{\sigma_A}$$

$$v = 100 \Rightarrow v' = \frac{100 - 400}{316.2277} = -0.9486$$

$$v = 200 \Rightarrow v' = \frac{200 - 400}{316 \cdot 2277} = -0.6324$$

$$v = 300 \Rightarrow v' = \frac{300 - 400}{316 \cdot 2277} = -0.3162$$

$$v = 500 \Rightarrow v' = \frac{500 - 400}{316 \cdot 2277} = 0.3162$$

$$v = 900 \Rightarrow v' = \frac{900 - 400}{316 \cdot 2277} = 1.5811$$

Values	100	200	300	500	900
Normalized	-0.9486	-0.6324	-0.3162	0.3162	1.5811

PART B

The following data is given in increasing order for the attribute age : 13, 15, 16, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70.

- a) Use smoothing by bin boundaries to smooth these data, using bin depth of 3.
- b) How might you determine outliers in the data?
- c) What other methods are there for data smoothing? (Mod II)
- d) Bin depth = 3; No. of attributes = 27; Each bin contains $\frac{N}{D} = \frac{27}{3} = 9$ items.

		Smoothed values	Smoothed value
Bin 1	: 13, 15, 16	13, 13, 16	13, 13, 16
Bin 2	: 16, 19, 20	16, 16, 20	16, 16, 20
Bin 3	: 20, 21, 22	20, 20, 22	20, 20, 22
Bin 4	: 22, 25, 25	22, 22, 25	22, 22, 25
Bin 5	: 25, 25, 30	25, 25, 30	25, 25, 30
Bin 6	: 33, 33, 35	33, 33, 35	33, 33, 35
Bin 7	: 35, 35, 35	35, 35, 35	35, 35, 35
Bin 8	: 36, 40, 45	36, 36, 45	36, 36, 45
Bin 9	: 46, 52, 70	46, 46, 70	46, 46, 70

- 12) Outliers in data can be identified in several ways
- By dividing the data into equi-width histograms and identifying the outlying histograms
 - By clustering the data into groups. Any data that do not fall in a group can be taken as outliers.
 - In general, fit a model to the data. Any data points that deviate significantly (based on some threshold) from the model can be considered outliers.

- c) Other methods that can be used for data smoothing include alternate forms of binning such as:
- Smoothing by bin medians
 - Smoothing by bin boundaries

Alternatively, equal width bin can be used to implement any of the forms of binning, where the internal range of values in each bin is constant.

Methods other than binning include using regression technique to smooth data by setting it to a function such as through linear or multiple regression.

Classification techniques can be used to implement concept hierarchies & that can smooth the data by rolling-up lower level concepts to higher level concepts.

12. Explain the following procedures for attribute subset selection:
- Stepwise forward selection
 - Stepwise backward elimination
 - A combination of forward selection and backward elimination.

a) Stepwise forward selection.

This procedure starts with an empty set of attributes as the minimal set. The most relevant attributes are chosen (having minimal p-value) and are added to the minimal set. In each iteration, one attribute is added to a reduced set.

b) Stepwise backward selection

Here all the attributes are considered in the initial set of attributes. In each iteration, one attribute is eliminated from the set of attributes whose p-value is higher than significance level.

c) A combination of forward selection and backward elimination.

The stepwise forward selection and backward elimination are combined so as to select the relevant attributes most efficiently. This is the most common technique which is generally used for attribute selection.

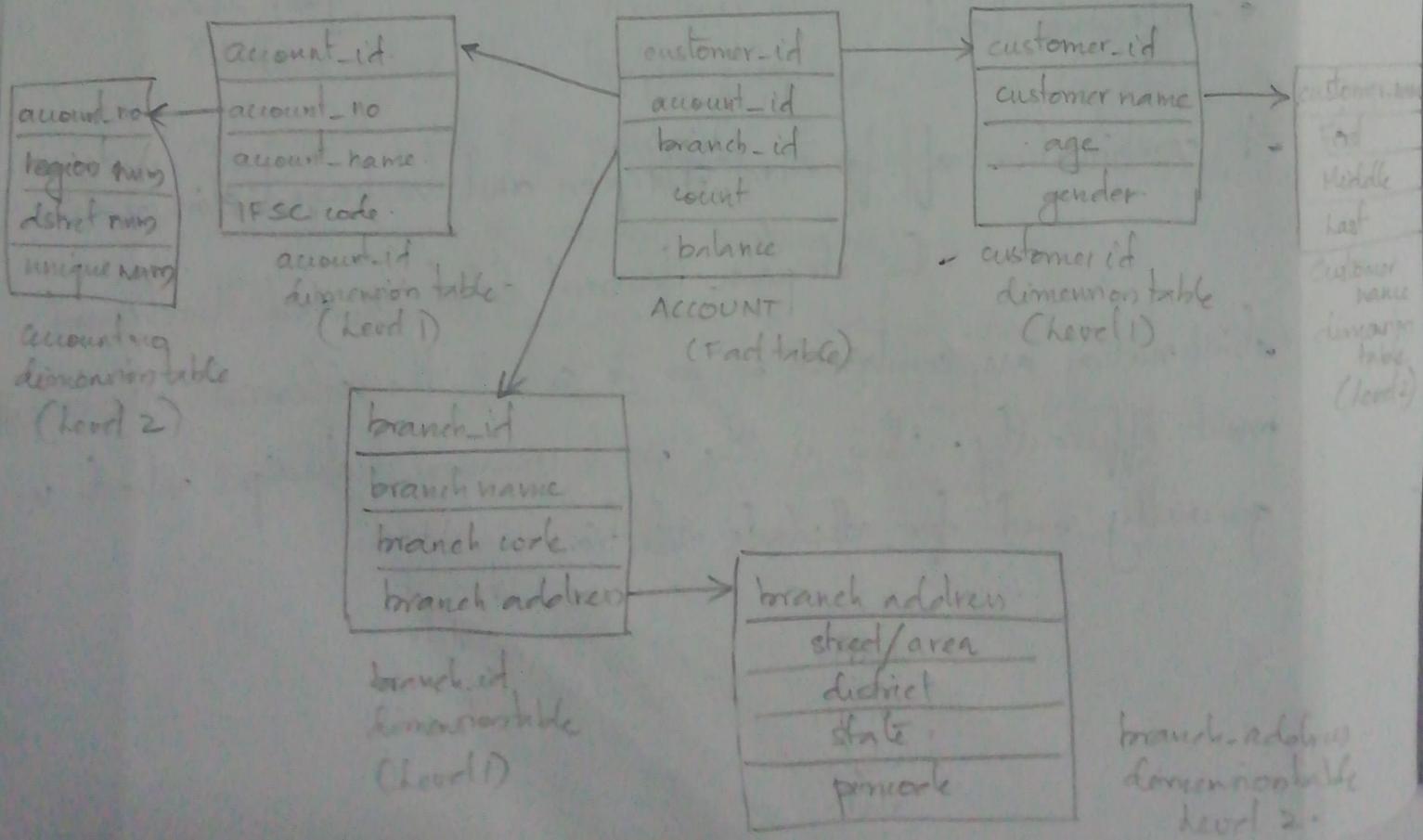
18) b) a) Suppose a dataware house consists of three measures customer, account and branch and two measures count (no: of customers in the branch) and balance. Draw the schema diagram using snowflake schema.

b) Real world entity data tend to be incomplete, noisy and inconsistent. What are the various approaches adopted to clean the data?

b) Data cleaning tasks

- Fill in missing values
- Identify outliers and smooth out noisy data
- Correct inconsistent data
- Resolve redundancy caused by data integration

a)



PART A

4 An airport security screening station wants to determine if passengers are criminals or not. To do this, the face of passengers are scanned and kept in a database. Is this a classification or prediction task? Justify. (Mod III)

Classification can be defined as the process of finding a model (or function) that describes and distinguishes data classes or concepts, for the purpose of being able to use the model to predict the class of objects whose class label is unknown. The desired model is based on the analysis of a set of training data.

The face of the passenger is scanned and its basic pattern (distance b/w eyes, size and shape of mouth, head, etc.) is identified. This pattern is compared to entries in a database to see if it matches any patterns that are associated with known offenders.

5 Where do we use linear regression? Explain linear regression. (Mod IV)

A straight line regression analysis involves a response variable y and a single predictor variable x .

$$\text{ie } y = b + cx$$

b and c are regression coefficients, specifying y intercept and slope of the line respectively.

16

→ The regression coefficients can be thought of as weights.

$$y = w_0 + w_1 x$$

→ These coefficients can be solved for using by the method of least square.

→ Let D be the training set consisting of values of predictor variable x and their associated response variable y .

The training set contains $|D|$ datapoints of the form $(x_1, y_1), (x_2, y_2), \dots, (x_{|D|}, y_{|D|})$.

The regression coefficients can be estimated by.

$$w_1 = \frac{\sum_{i=1}^{|D|} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{|D|} (x_i - \bar{x})^2}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

\bar{x} = Mean value of $x_1, x_2, \dots, x_{|D|}$

\bar{y} = Mean value of $y_1, y_2, \dots, y_{|D|}$

→ Multiple linear regression is an extension of straight line regression which involves more than one predictor variable.

→ It allows response variable y to be modelled as a linear function of n predictor variables or attributes A_1, A_2, \dots, A_n describing a tuple X

Example of multiple linear regression model based on two predictor attributes are

$$y = w_0 + w_1 x_1 + w_2 x_2$$

x_1 and $x_2 \rightarrow$ Value of attribute A₁ & A₂

The method of least squares can be extended to calculate w_0 , w_1 & w_2 .

→ The major uses for regression analysis are

- Determining the strength of predictors
- Forecasting an effect
- Trend forecasting.

Decision tree
for prediction

Can be used for prediction

Decision tree
for prediction
a change in one
more independent var

6. What is the significance of tree pruning in decision tree algorithms (Mod III)

→ Main approach to avoid overfitting is pruning

→ Technique that reduces the size of decision tree by removing sections of the tree that provide little power to classify instances.

→ Reduces complexity of the final classifier and hence improves predictive accuracy by the reduction of overfitting.

→ The pruning phase might remove redundant comparisons or remove subtrees to achieve better performance.

→ When decision trees are built, many of the branches may reflect noise or outliers in the training data.

Tree pruning methods address this problem of overfitting the data.

→ Tree pruning methods attempt to identify and remove such branches, with the goal of improving classification accuracy on unseen data.

- Decision tree suffer from repetition and replication, making them overwhelming to interpret
- Repetition occurs when an attribute is repeatedly ~~making~~ tested along a given branch of the tree
- For replication, duplicate subtrees exist within the tree
- These situations can impede the accuracy and comprehensibility of a decision tree.

PART C

- 14 Based on the following data determine the gender of a person having height 6 ft, weight 130 lbs and foot size 8 cm (Naïve Bayes algorithm)

Person	Height (feet)	Weight (lbs)	foot size (inches)
male	6.00	180	10
male	6.00	180	10
male	5.50	170	8
male	6.00	170	10
female	5.00	130	8
female	5.50	150	6
female	5.00	130	6
female	6.00	150	8

(Mod III)

Person	Height (feet)			Weight (lbs)			foot size (inches)			
	6.00	5.50	5.00	180	170	150	130	10	8	6
Male	3	1	0	2	2	0	0	3	1	0
Female	1	1	2	0	0	2	2	0	2	2

Conditional
Probability table

Person	Height (feet)	Weight (lbs)	foot size (inches)
Male	6.00	5.50	5.00
Female	3/4	1/4	0

$$P(\text{Male}) = \frac{\text{No. of male}}{\text{Total persons}} = \frac{1}{8} = \frac{1}{2}$$

19

$$P(\text{Female}) = \frac{\text{No. of female}}{\text{Total persons}} = \frac{4}{8} = \frac{1}{2}$$

Probability if the conditions satisfied in male

$$\begin{aligned} q_1 &= P(\text{6ft/Male}) \times P(130\text{lbs/Male}) \times P(8\text{in/Male}) \times P(\text{Male}) \\ &= \frac{3}{4} \times 0 \times \frac{1}{4} \times \frac{1}{2} \\ &= 0 \end{aligned}$$

$$\begin{aligned} q_2 &= P(\text{6ft/Female}) \times P(130\text{lbs/Female}) \times P(8\text{in/Female}) \times P(\text{Female}) \\ &= \frac{1}{4} \times \frac{3}{4} \times \frac{2}{4} \times \frac{1}{2} \\ &= \frac{1}{32} \\ &= 0.03125 \end{aligned}$$

$$\max(q_1, q_2) = \max(0, 0.03125) \\ = 0.03125 = q_2$$

So the gender of the person having height 6ft, weight 130lbs and foot size 8inch is Female

$$P(A|C) = \frac{P(C|A) P(A)}{P(C)}$$

$P(A|C)$ → Posterior probability - which represents the degree to which we believe a given model accurately describes the situation given the available data and all of our prior information.

$P(C|A)$ → Likelihood - which describes how well the model predicts the data

$P(A)$ → Prior probability - which describes the degree to which we believe the model accurately describes reality based on all of our prior information.

$P(C)$ → Normalizing constant - The constant that makes the posterior density integrate to one.

15 The "Restaurant A" sells burger with optional flavours : Pepper, Ginger and Chilly. Every day this week you have tried a burger (A to E) and kept a record of which you liked. Using Hamming distance, show how the 3NN classifier with majority voting would classify.

$\{ \text{pepper} = \text{false}; \text{ginger} = \text{true}; \text{chilly} = \text{true} \}$

	Pepper	Ginger	Chilly	Liked
A	false	true	true	false
B	true	false	false	true
C	false	true	true	false
D	false	true	false	true
E	true	false	false	true

(Mod IV)

The earliest local distance function known as the overlap func, returns 0 if the two values are equal and 1 otherwise

$$\text{dist}_A(x, A, q, A) = \text{def } \begin{cases} 0 & \text{if } x \cdot A = q \cdot A \\ 1 & \text{otherwise} \end{cases}$$

Hamming Distance :

Let Q_n be the $\{\text{pepper} = \text{false}, \text{ginger} = \text{true}, \text{chilly} = \text{true}\}$.

$$\text{Dist}(A, Q_n) = 1 + 0 + 0 = 1$$

$$\text{Dist}(B, Q_n) = 1 + 1 + 1 = 3$$

$$\text{Dist}(C, Q_n) = 0 + 0 + 0 = 0$$

$$\text{Dist}(D, Q_n) = 0 + 0 + 1 = 1$$

$$\text{Dist}(E, Q_n) = 1 + 1 + 1 = 3$$

true & true $\Rightarrow 0$
true & false $\Rightarrow 1$

0, 1, 1, 3, 3.
C, A, D, B, E

Arrange the distances in ascending order and the last 3 distances are taken. So A, C, D are taken.

While considering them A - False B - False D - True
So the majority voting is taken and 2 votes for false

and 1 vote for tree ~~so considering one~~

Hence No. of fabr > No of tree

So concluding that class liked = False for the condition {pepper = false; ginger = tree; chilly = true}

- 16 a) How C4.5 differs from ID3 algorithm? (Mod III)
 b) How does backpropagation algorithm works? (Mod IV)

a)	ID3	C4.5
Splitting criteria	Information Gain	Gain ratio
Attribute type	Handles only categorical value	Handles only categorical & numerical values
Missing values	Do not handle missing values	Handle missing values
Pruning strategy	No pruning is done	Error based pruning is used
Outlier detection	Susceptible to outliers	Susceptible to outliers

b) Backpropagation learns by iteratively processing a data set of training tuples, comparing the network's prediction for each tuple with the actual known target value.

The target value may be known class label of the training tuple (for classification prob) or continuous values (for numeric prediction).

For each training tuple, the weights are modified so as to minimize the mean squared error b/w the net's prediction and the actual target value.

These modifications are made in the "backwards" direction

(i.e. from output layer) through each hidden layer down to the first hidden layer (-hence the name propagation).

Although it is not guaranteed, in general the weights will eventually converge and the learning process stops.

Algorithm

Input :- D, a dataset consisting of the training tuples and their associated target values
l, the learning rate
network, a multilayer feed-forward m/c

Output :- a trained neural network

Method :

1. Initialize all weights and bias in m/c
2. while terminating condition not satisfied {
 3. for each training tuple $\mathbf{x}_i \in D$ {
 4. // Propagate the inputs forward for each input layer unit j {
 5. $O_j^0 = I_j$; // o/p of an ip unit is its actual value
 6. for each hidden or o/p layer unit j {
 7. $I_j = \sum_i w_{ij} O_i + O_j$; // compute the net of unit j wrt pre-layer, i
 8. $O_j = \frac{1}{1 + e^{-I_j}}$; // compute the o/p of each unit j
 9. // Backpropagate the errors for each unit j in the o/p layer
 10. $E_{xj} = O_j(1 - O_j)(T_j - O_j)$; // compute the error

13. for each unit j in the hidden layers, from the last to the first hidden layer:
14. $E_{\text{err},j} = O_j(1 - O_j) E_{\text{err},k} w_{jk};$
 // compute the error w.r.t the next higher layer, k
15. for each w_{ij} in network {
16. $w_{ij} = (l) E_{\text{err},j} O_i;$ // weight increment
17. $w_{ij} = w_{ij} + \Delta w_{ij};$ } // weight update
18. for each bias θ_j in network {
19. $\theta_j = (l) E_{\text{err},j};$ // bias increment
20. $\theta_j = \theta_j + \Delta \theta_j;$ } // bias update
21. }

Terminating conditions: Training stops when

- all Δw_{ij} in the previous epoch are so small as to be below some specified threshold or
- The percentage of tuples misclassified in the previous epoch is below some threshold or
- A prespecified number of epochs has expired.

A. Explain the attribute selection method in decision trees (Not 11)



→ Entropy :- Entropy is a measure of the randomness in the information being processed. The higher the entropy, the harder it is to draw any conclusions from that information.

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

where $S \rightarrow$ current state

$p_i \rightarrow$ probability of an event i of state S or percentage of class i in a node of state S .

→ Information Gain : Information Gain is a Statistical property that measures how well a given attribute separates the training examples according to their target classification. Constructing a decision tree is all about finding an attribute that returns the highest information gain and the smallest entropy.

$$\text{Information Gain} = \text{Entropy}(\text{before}) - \sum_{j=1}^K \text{Entropy}(j, \text{after})$$

where "before" is the dataset before the split, K is the no. of subsets generated by the split, and (j, after) is subset j after the split.

→ Gini Index - Gini Index is a cost function used to evaluate splits in the dataset. It is calculated by subtracting the sum of the squared probabilities of each class from one.

$$\text{Gini} = 1 - \sum_{i=1}^c (p_i)^2$$

Gain index works with the categorical target variable "Success" or "Failure". It performs only Binary splits.

→ Gain Ratio : Gain Ratio overcomes the problem with information gain by taking into account the number of branches that would result before making the split. It corrects information gain by taking the intrinsic information of a split into account.

$$\text{Gain ratio} = \frac{\text{Information Gain}}{\text{Split Info}} = \frac{\text{Entropy}(\text{before}) - \sum_{j=1}^k \text{Entropy}(j) \frac{\omega_j}{\text{after}}}{\sum_{j=1}^k \omega_j \log_2 \omega_j}$$

→ Reduction in Variance - Reduction in Variance is an algorithm used for continuous target variables (regression problems). This algorithm uses the standard formula of variance to choose the split. The split with lower variance is selected as the criteria to split the population.

$$\text{Variance} = \frac{\sum_i (x_i - \bar{x})^2}{n}$$

→ Chi Square - It is one of the oldest tree classification methods. It finds out the statistical significance b/w the differences between subnodes and parent node. We measure it by the sum of squares of standardized diff b/w observed and expected frequencies of the target variable.

$$\chi^2 = \sum \frac{(O-E)^2}{E} \quad \text{where}$$

$$\begin{aligned} \chi^2 &= \text{Chi square obtained} \\ \Sigma &= \text{the sum of} \\ O &= \text{observed score} \\ E &= \text{expected score} \end{aligned}$$

26 5 Distinguish b/w hold out method and cross validation method. (Mod 11)

Cross validation is usually preferred method because it gives your model the opportunity to learn on multiple train-test splits. This gives you a better indication of how well your model will perform on unseen data. Hold out on other hand, is dependent on just one train-test split. That makes the holdout method less dependent on how the data is split into train and test sets.

The hold out method is good to use when you have a very large dataset, you're on a time crunch, or you are starting to build an empirical model in your data science project. Keep in mind that because cross-validation uses multiple-train-test splits, it takes more computational power and time to run than using the holdout method.

6 Explain prepruning and post pruning approaches in decision tree algorithm. (Mod 11)

Tree pruning is performed in order to remove anomalies in the training due to noise or outliers. The pruned trees are smaller and less complex.

Types of pruning

Post pruning - The tree is pruned by halting its construction early

Pre pruning - This approach removes a subtree from a fully grown tree

The cost complexity is measured by following two parameters

- * No. of leaves in the tree

- * Error rate of the tree.

Post Pruning

- First build the tree
- Then prune it
 - * Fully grown tree shows all attribute interactions
- Problem: Some subtrees might be due to chance effects
- Two pruning operations
 - * Subtree replacement
 - * Subtree raising.
- Possible strategies
 - * Error estimation
 - * Significance testing
 - * MDL principle.

Pre Pruning

- Based on statistical significance test.
- * Stop growing the tree when there is no statistically significant association b/w any attribute and the class at a particular node.
- * Use all available data for training and apply statistical test to estimate whether expanding / pruning a node is to produce an improvement beyond the training set.
- Most popular test - Chi squared test
- ID3 used chi square test in addition to information gain
 - * Only statistically significant attributes were allowed to be selected by information gain procedure.
- Early stopping:- Pre pruning may stop the growth process prematurely.

PART C

28

- 14 Given the following data on a certain set of patients seen by a doctor, can the doctor conclude that a person having chills, fever, mild headache and without running nose has the flu? (Use Naive Bayes algorithm for prediction)

Chills	Running nose	Headache	Fever.	Has flu
Y	N	Mild	Y	N
Y	Y	No	N	Y
Y	N	Strong	Y	Y
N	Y	Mild	Y	Y
N	N	No	N	N
N	Y	Strong	Y	Y
N	Y	Strong	N	N
Y	Y	Mild	Y	Y

(Mod III)

Has flu	Chills		Running nose		Headache		Fever		
	Yes	No	Yes	No	Mild	Strong	No	Yes	No
Yes	3	2	4	1	2	2	1	4	1
No	1	2	1	2	1	1	1	1	2

Conditional Probability Table ↴

Has flu	Chills		Running Nose		Headache		Fever		
	Yes	No	Yes	No	Mild	Strong	No	Yes	No
Yes	3/5	2/5	4/5	1/5	2/5	2/5	1/5	4/5	1/5
No	1/3	2/3	1/3	2/3	1/3	1/3	1/3	1/3	2/3

$$P(\text{Has flu} = \text{yes}) = \frac{\text{No. of Has flu} = \text{"yes}}{\text{Total no.}} = \frac{5}{8}$$

$$P(\text{Has flu} = \text{"no"}) = \frac{\text{No. of Has flu} = \text{"no"} }{\text{Total no.}} = \frac{3}{8}$$

Probability of having chills, fever, mild headache & without running nose has flu

$$\begin{aligned} q_1 &= P(\text{Chills} = \text{yes} / \text{Has flu} = \text{yes}) \times P(\text{Fever} = \text{yes} / \text{Has flu} = \text{yes}) \times \\ &\quad P(\text{Mild headache} / \text{Has flu} = \text{yes}) \times P(\text{Running nose} = \text{no} / \text{Has flu} = \text{yes}) \\ &\quad \times P(\text{Has flu} = \text{yes}) \\ &= \frac{3}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{4}{5} \times \frac{5}{8} \\ &= \frac{12}{125} \\ &= 0.096 \end{aligned}$$

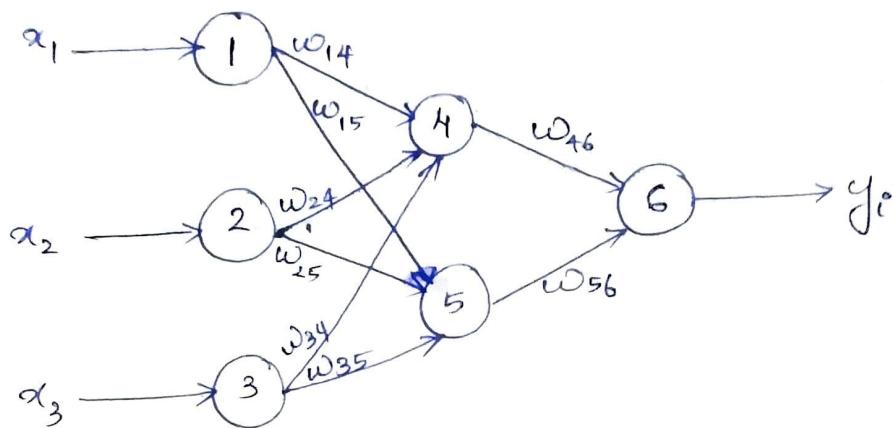
$$\begin{aligned} q_2 &= P(\text{Chills} = \text{yes} / \text{Has flu} = \text{"no"}) \times P(\text{Fever} = \text{yes} / \text{Has flu} = \text{"no"}) \times \\ &\quad P(\text{Mild headache} / \text{Has flu} = \text{"no"}) \times P(\text{Running nose} = \text{no} / \text{Has flu} = \text{"no"}) \\ &\quad \times P(\text{Has flu} = \text{"no"}) \\ &= \frac{1}{3} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} \times \frac{3}{8} \\ &= \frac{1}{54} \\ &= 0.0185 \end{aligned}$$

$$\begin{aligned} \max(q_1, q_2) &= \max(0.096, 0.0185) = \\ &= 0.096 = q_1 \end{aligned}$$

So the person has flu = Yes

30

15 The following figure shows a multilayer feed-forward neural network. Let the learning rate be 0.9. The initial weights and bias values of the network is given in the table below. The activation function used is the sigmoid function.



x_1	x_2	x_3	w_{14}	w_{15}	w_{24}	w_{25}	w_{34}	w_{35}	w_{46}	w_{56}	θ_4	θ_5
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2

Show weight and bias update with the first training sample (1, 0, 1) with class label 1, using backpropagation algorithm.
(Mod IV)

Net Input and Output Calculations

Unit j	Net Input I_j	Output, O_j
4	$x_1 w_{14} + x_2 w_{24} + x_3 w_{34} + \theta_4$ $= (1 \times 0.2) + (0 \times 0.4) + (1 \times -0.5) + (-0.4)$ $= 0.2 + 0 + -0.5 - 0.4$ $= -0.7$	$\frac{1}{1+e^{-0.7}}$ $= 0.332$
5	$x_1 w_{15} + x_2 w_{25} + x_3 w_{35} + \theta_5$ $= (1 \times -0.3) + (0 \times 0.1) + (1 \times 0.2) + (0.2)$ $= -0.3 + 0 + 0.2 + 0.2$ $= 0.1$	$\frac{1}{1+e^{-0.1}}$ $= 0.525$

$$\begin{aligned}
 & w_{46} \times O_4 + w_{56} \times O_5 + \Theta_6 \\
 & = (-0.3)(0.332) + (0.2)(0.525) + 0.1 \\
 & = -0.105
 \end{aligned}
 \quad \frac{1}{1+e^{-0.105}} = 0.474$$

Calculation of the error at each node

Out_j	$Error_j$
6	$(0.474)(1-0.474)(1-0.474) = 0.1311$
5	$(0.525)(1-0.525)(0.1311)(-0.2) = -0.0065$
4	$(0.332)(1-0.332)(0.1311)(-0.3) = -0.0087$

Calculation for Weight and Bias Updating.

Weight Bias	New Value	
w_{46}	$-0.3 + (0.9)(0.1311)(0.332)$	$= -0.261$
w_{56}	$-0.2 + (0.9)(0.1311)(0.525)$	$= -0.138$
w_{14}	$0.2 + (0.9)(-0.0087)(1)$	$= 0.192$
w_{15}	$-0.3 + (0.9)(-0.0065)(1)$	$= -0.306$
w_{24}	$0.4 + (0.9)(-0.0087)(0)$	$= 0.4$
w_{25}	$0.1 + (0.9)(-0.0065)(0)$	$= 0.1$
w_{34}	$-0.5 + (0.9)(-0.0087)(1)$	$= -0.508$
w_{35}	$-0.2 + (0.9)(-0.0065)(1)$	$= 0.194$
Θ_6	$0.1 + (0.9)(0.1311)$	$= 0.218$
Θ_5	$0.2 + (0.9)(-0.0065)$	$= 0.194$
O_4	$-0.4 + (0.9)(0.0087)$	$= -0.408$

- 32
- a) Explain classification by C4.5 algorithm (Mod III)
b) What is meant by Maximum Marginal Hyperplane (MMH) (Mod IV)

- a) C4.5 algorithm is used to generate a decision tree which was developed by Ross Quinlan.
- It is an extension of Quinlan's ID3 algorithm.
 - C4.5 generates decision trees which can be used for classification & therefore C4.5 is often referred to as statistical classifier.
 - It is better than ID3 algorithm because it deals with both continuous and discrete attributes and also with the missing values and pruning trees after construction.
 - C5.0 is the commercial successor of C4.5 because it is a lot faster, more memory efficient and used for building smaller decision tree. C4.5 performs by default a tree pruning process.
 - This leads to the formation of smaller trees, more simple rules and produces more intuitive interpretations.
- C4.5 follows 3 steps in tree growth.

- i) For splitting of categorical ~~values~~ attributes, C4.5 follows the similar approach to ID3 algo. Continuous attribute always generate binary splits.
- ii) Selecting attribute with highest gain ratio
- iii) These steps are applied to new tree branches and growth of the tree is stopped after checking of stop criterion.
Information gain bias the attribute with more no. of values. Thus, C4.5 uses gain ratio which is a less biased selection criterion.

Advantages of C4.5

- i) Easy to implement
- ii) Builds models that can be easily interpreted
- iii) It can deal with noise and deal with missing value attributes

Disadvantages of C4.5

- i) A small variation in data can lead to diff decision trees when using C4.5
- ii) For a small training set, C4.5 does not work well.

2) b) There are infinite number of separating lines that could be drawn and we want to find the best one which have the minimum classification error on previously unseen tuples.

If our data were 3-D (ie three attributes), we would want to find the best separating plane.

Generalizing to n dimensions, we want to find the best hyperplane. Hyperplane refers to the decision boundary that we are seeking, regardless of the number of input attributes.

An SVM approaches this problem by searching for the maximum marginal hyperplane.

In fig both hyperplanes can classify all the given data tuples. However we expect the hyperplane with the larger margin to be more accurate at classifying future data tuples than the hyperplane with the smaller margin.

This is why (during the learning or training phase) the SVM searches for the hyperplane with the largest margin, ie., the maximum marginal hyperplane (MMH). The associated margin gives the largest separation between classes.



The shortest distance from a hyperplane to one side of its margin is equal to the shortest distance from the hyperplane to the other side of its margin, where the "sides" of the margin are parallel to the hyperplane.

When dealing with MMH, this distance is in fact the shortest dist from MMH to closest training tuple of either class.

7. What are two measures used for rule interestingness? (Mark V)

The strength of a rule is measured by its support and confidence.

Support (X) is the no. of times X in the transaction divided by total no. of transaction.

Confidence of the association rule $X \rightarrow Y$ is defined as the ratio of X and Y together to the support of X .

$$\text{Support}(X) = \frac{\text{No. of item } X \text{ appears}}{\text{Total no. of transaction}} = P(X)$$

$$\text{Confidence}(X \rightarrow Y) = \frac{\text{Support of } X \text{ & } Y \text{ together}}{\text{Support of } X} = \frac{P(X \cap Y)}{P(X)} = P(Y|X)$$

Support is a useful measure because if it is low, the rule may just occur due to chance. Furthermore in a business environment a rule covering too few cases (or transactions) may not be useful because it does not make business sense to act on such a rule (not profitable).

For eg., In the database, the itemset {milk, bread, butter} has a support of $\frac{1}{5} = 0.2$ since it occurs in 20% of all transactions (1 out of 5 transactions).

Confidence thus determines the predictability of the rule. If the confidence of a rule is too low, one cannot reliably infer or predict Y from X . A rule with low predictability is of limited use.

For eg., the rule {butter, bread} \rightarrow {milk} has confidence of $\frac{0.2}{0.2} = 1.0$ in the database, which means that for 100% of transactions containing butter and bread the rule is correct.

Given two objects represented by the tuples $(22, 1, 42, 10)$ and $(20, 0, 36, 8)$. Compute the Manhattan distance b/w the two objects

$$\text{Manhattan distance} = |22-20| + |1-0| + |42-36| + |10-8|$$

$$(\text{City block distance}) = |2| + |1| + |6| + |2|$$

$$= \underline{\underline{11}} \quad \text{distance} = \sum_{i=0}^{n-1} |(x[i] - y[i])|$$

Manhattan distance - The distance b/w two points measured along axes at right angles. In a plane with p_1 at (x_1, y_1) and p_2 at (x_2, y_2) , it is $|x_1 - x_2| + |y_1 - y_2|$

9 How can density based clustering varies from other methods? (Mod VI)

- Density based method finds ~~clusters~~ clusters of arbitrary shape. It grows the clusters with as many points as possible till some threshold is met. The ϵ -neighbourhood of a point is used to find dense regions in the database.
- Density based algorithms typically regard clusters as dense regions of objects in the data space that are separated by regions of low density. The main idea of density-based approach is to find regions of high density and low density, with high-density regions being separated from low-density regions. These approaches can make it easy to discover arbitrary clusters.
- A common way to divide the high dimensional space is to density-based grid units. Units containing relatively high densities are the cluster centres and the boundaries b/w clusters fall in the regions of low-density units
- Diff type of density based methods
 - DBSCAN
 - OPTICS
 - DENCLUE

10 Differentiate web content mining and web structure mining?
 (Mod VI)

Web Content mining

- Web Content mining is the application of extracting useful information from the content of the web document.
- Web content consists of several types of data - text, image, audio, video etc
- Content data is the group of facts that a web page is designed.
- It can provide effective and interesting patterns about user needs
- Text documents are related to text mining, machine learning and natural language processing
- This mining is also known as text mining
- This type of mining performs scanning and mining of the text, images and groups of web pages according to the content of the input.

Web Structure mining

- Web Structure mining is the application of discovering structure information from the web.
- The structure of the web graph consists of web pages as nodes, and hyperlinks as edges connecting related pages.
- Structure mining basically shows the structured summary of a particular website
- It defines relationship b/w web pages linked by information or direct link connections
- To determine the connection b/w two commercial websites, Web Structure mining can be very useful.
- This type of datamining can be performed document level (contra page) ↪ hyperlink level (contra page)

PART D

- 38 17. Consider the transaction database given below. Set minimum support count as 2 and minimum confidence threshold as 70%.

Transaction ID	List of Item-IDs
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₃
T800	I ₁ , I ₂ , I ₃ , I ₅
T900	I ₁ , I ₂ , I ₃

- a) Find the frequent itemset using FP Growth Algorithm
 b) Generate strong association rules
 (Mod 1)
 (Mod 1)

- a) Minimum support count = 2
 Total no. of transactions = 9.

Step: Construct list of items & table with support count

Itemset	Support count
I ₁	6
I ₂	7
I ₃	6
I ₄	2
I ₅	2

Table 1

Table 2

Step II: Write the items in descending order of support count.

Items	Sup. count
I ₂	7
I ₁	6
I ₃	6
I ₄	2
I ₅	2

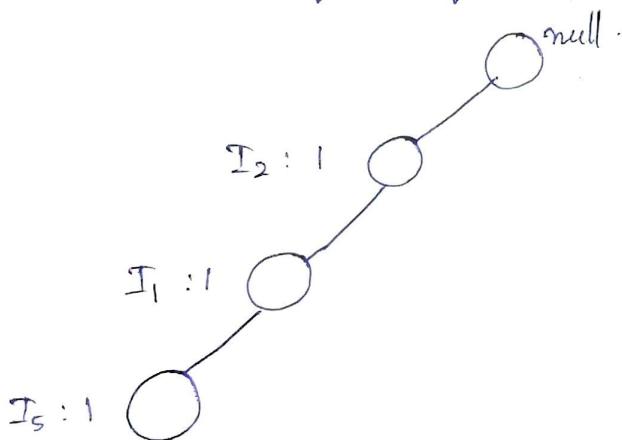
Table 3

Step III: Start drawing the FP tree with root node as null.



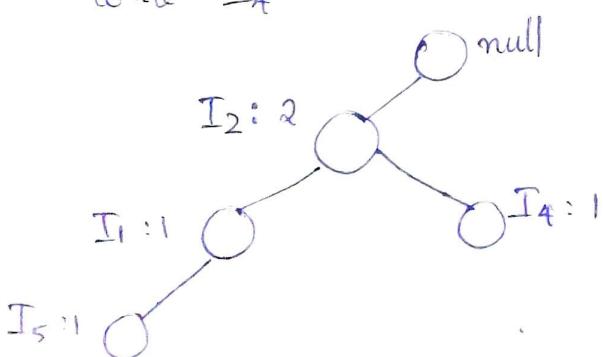
Step IV: Now consider the first transaction in table (T₁₀₀) and arrange in descending order as of sup count.
So T₁₀₀ \Rightarrow I₂(7), I₁(6), I₅(2).

Start from left node.



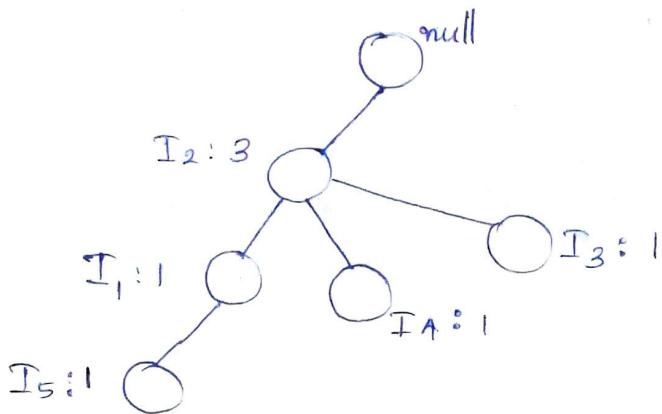
Next transaction T₂₀₀ arrange it \Rightarrow I₂, I₄

Now check. I₂ is there from the root node if there increment the count of I₂ then check whether I₄ is there after I₂ if not then add another branch from I₂ and write I₄ with count 1

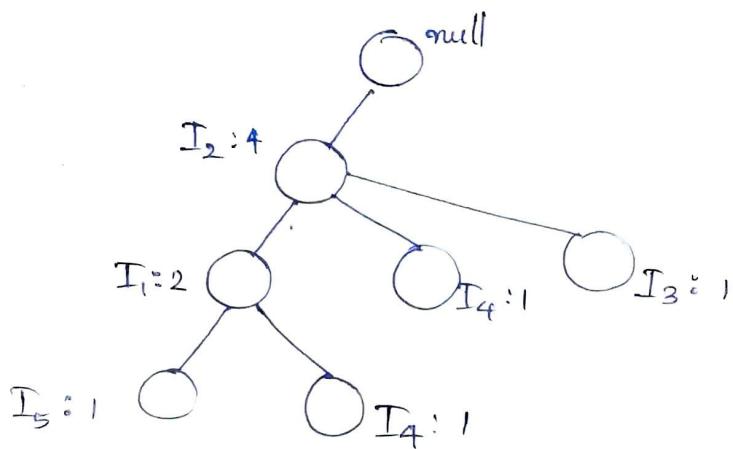


AO

Next transaction $T_{300} \Rightarrow$ Arranging $\Rightarrow I_2, I_3.$

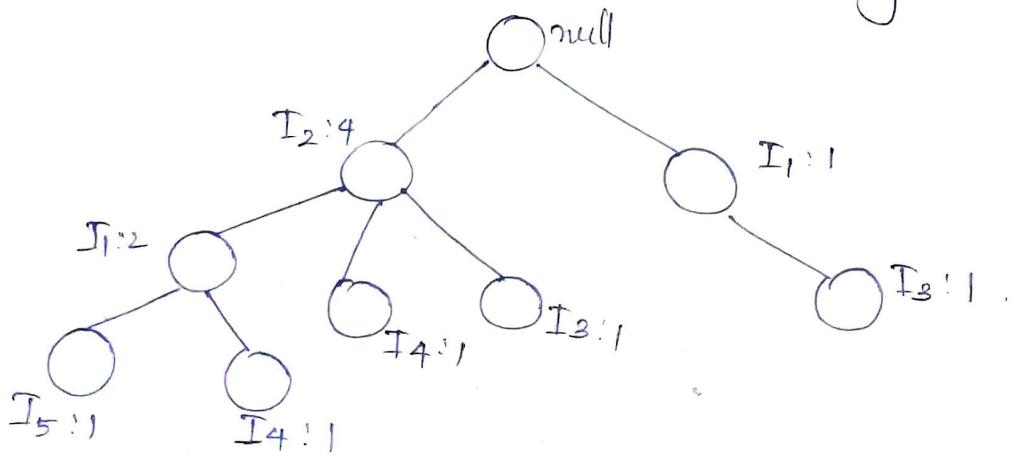


$T_{400} \Rightarrow$ Arranging $\Rightarrow I_2, I_1, I_4$

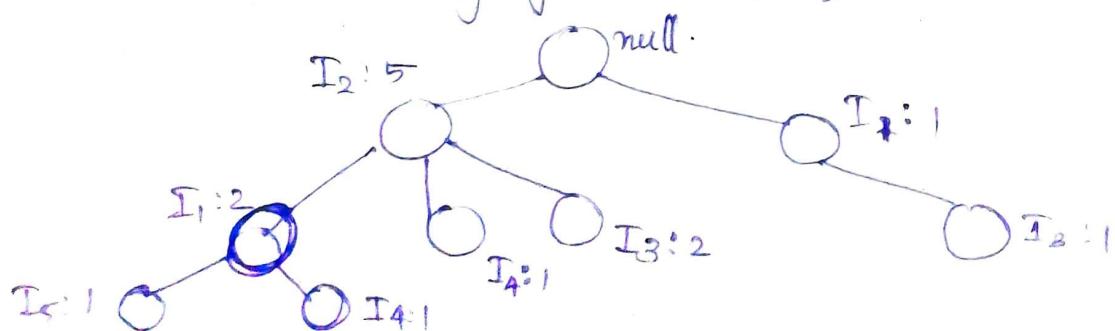


$T_{500} \Rightarrow$ Arranging $\Rightarrow I_1, I_3.$

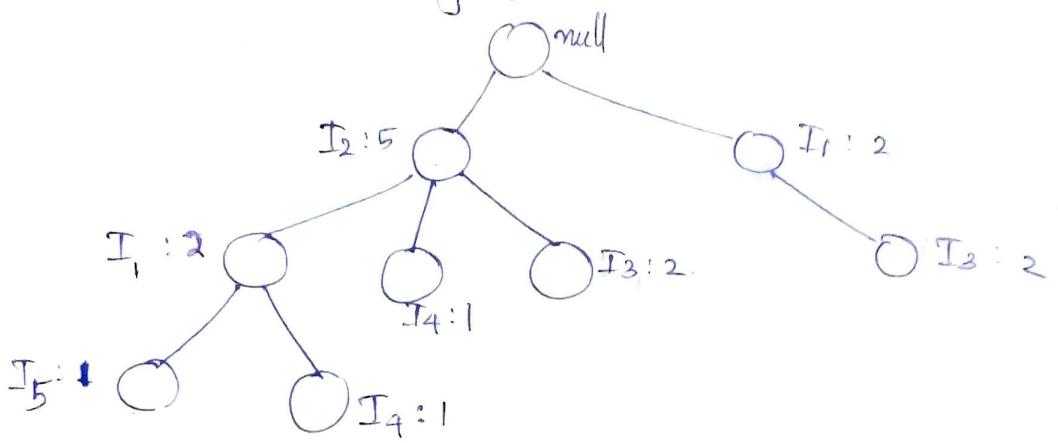
There is no I_1 from null so generate another branch.



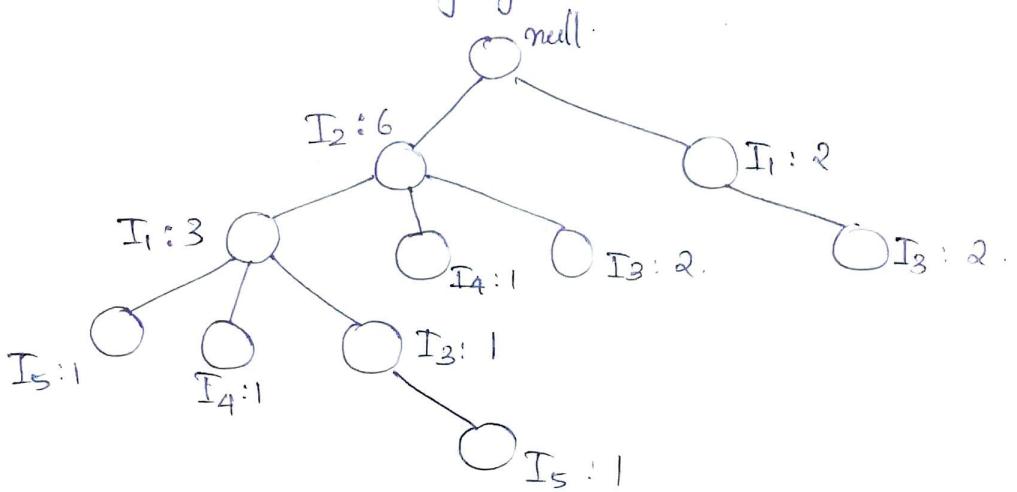
$T_{600} \Rightarrow$ Arranging $\Rightarrow I_2, I_3.$



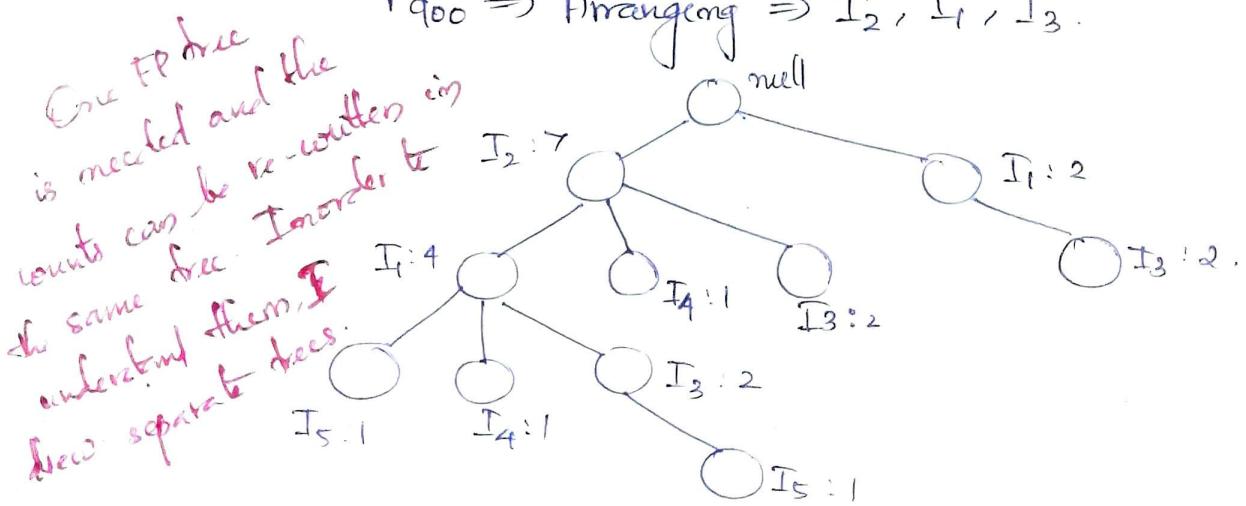
$T_{700} \Rightarrow$ Arranging $\Rightarrow I_1, I_2$



$T_{800} \Rightarrow$ Arranging $\Rightarrow I_2, I_1, I_3, I_5$



$T_{900} \Rightarrow$ Arranging $\Rightarrow I_2, I_1, I_3$



Now check the total count of transactions in the tree and compare it with the table. If the counts are same then the tree is correct.

Item	Conditional Pattern Base	Conditional FP tree
I_5	$\{\{I_2, I_1 : 1\}, \{I_2, I_1, I_3 : 1\}\}$	$\{I_2 : 2, I_1 : 2\}$
I_4	$\{\{I_2, I_1 : 1\}, \{I_2 : 1\}\}$	$\{I_2 : 2\}$
I_3	$\{\{I_2, I_1 : 2\}, \{I_2 : 2\}, \{I_1 : 2\}\}$	$\{I_2 : 4, I_1 : 2\}$
I_1	$\{\{I_2 : 4\}\}$	$\{I_2 : 4\}$

Frequent patterns generated.

$\{I_2, I_5 : 2\}, \{I_1, I_5 : 2\}, \{I_2, I_1, I_5 : 2\}$.

$\{I_2, I_4 : 2\}$.

$\{I_2, I_3 : 4\} \{I_1, I_3 : 4\} \{I_2, I_1, I_3 : 2\}$.

$\{I_2, I_1 : 4\}$.

Consider the database (Table 3) check for least sup count first.
So here I_5 / I_4 can be taken. Taking I_5 .

Check FP tree
How we traverse from root node to I_5 is checked and then writing each path and writing the respective sup count of I_5 in that path in the respective paths. There are ~~two~~ under conditional patterns base

Now in the column conditional FP tree write the count for each item. for item $I_5 \Rightarrow I_2 = 2, I_1 = 2, I_3 = 1$

But we don't write I_3 since its not satisfying min count = 2
So neglecting I_3 .

Now in freq patterns generated joining all each sets in conditional FP tree with item.

So joining I_2, I_5 with min count $\Rightarrow \{I_2, I_5 : 2\}$

" I_1, I_5 " " " $\Rightarrow \{I_1, I_5 : 2\}$.

" I_2, I_1, I_5 " " " $\Rightarrow \{I_2, I_1, I_5 : 2\}$.

Min count checked from conditional FP tree.

Similarly generate for I_4 , I_3 and I_1 , too.

Generate paths wherever the items are located.

For I_3 check for both RHTs and LHTs and

while writing the conditional FP tree write the paths in separate brackets when obtained from different LHTs and RHTs branch.

For $I_3 \Rightarrow \{I_1, I_3 : 2\}$ first but again from RHTs branch we get again $\{I_1, I_3 : 2\}$. So adding the count and combining so $\{I_1, I_3 : 4\}$

For $\{I_2, I_1, I_3 : 2\}$ since the least count = 2.

Thus frequent patterns are generated.

$\{I_2, I_5 : 2\}$, $\{I_1, I_5 : 2\}$, $\{I_2, I_1, I_5 : 2\}$

$\{I_2, I_4 : 2\}$

$\{I_2, I_3 : 4\}$, $\{I_1, I_3 : 4\}$, $\{I_2, I_1, I_3 : 4\}$

$\{I_2, I_1 : 4\}$.

Strong Association rules from items (where strong association rule satisfy both minimum support and minimum confidence)
Association rule can be generated as follows

- For each frequent itemset I_i , generate all nonempty subsets of I_i
- for every nonempty subset of I_i , output the rule

" $s \Rightarrow (l-s)$ " if $\frac{\text{support_count}(l)}{\text{support_count}(s)} \geq \text{min_conf}$,

where min-conf is the minimum confidence threshold.

confidence ($A \Rightarrow B$) = $P(B|A) = \frac{\text{support count}(A \cup B)}{\text{support count}(A)}$

3. frequent itemsets, taking I_1, I_2, I_5

Take the non empty subsets of I_1, I_2, I_5 . \emptyset set is not

I_1	\rightarrow	$I_2 \cap I_5$	6 rules generated.
I_2	\rightarrow	$I_1 \cap I_5$	
I_5	\rightarrow	$I_2 \cap I_5$	
$I_1 \cap I_2$	\rightarrow	I_5	
$I_1 \cap I_5$	\rightarrow	I_2	
$I_2 \cap I_5$	\rightarrow	I_1	
$I_1 \cap I_2 \cap I_5$	\rightarrow	\emptyset Not considering since there is no RHS.	

Now we have to check $s \rightarrow l-s \geq \text{min-conf}$.

Support count of $\{I_1, I_2, I_5\} : 2$, if it is 2.

Now check $\frac{\text{support-count}(l)}{\text{support-count}(s)} \geq \text{min-conf}$.

Rule.

Rule.	Confidence.
1) $I_1 \rightarrow I_2 \cap I_5$	$2/6 = 33.3\%$
2) $I_2 \rightarrow I_1 \cap I_5$	$2/4 = 50\%$
3) $I_5 \rightarrow I_1 \cap I_2$	$2/2 = 100\%$
4) $I_1 \cap I_2 \rightarrow I_5$	$2/4 = 50\%$
5) $I_1 \cap I_5 \rightarrow I_2$	$2/2 = 100\%$
6) $I_2 \cap I_5 \rightarrow I_1$	$2/2 = 100\%$

Min-confidence = 70%. So rules greater than 70% are.

$I_1 \cap I_5 \rightarrow I_2$

$I_2 \cap I_5 \rightarrow I_1$

$I_5 \rightarrow I_1 \cap I_2$

18 a) Explain BIRCH Clustering Method (Mod VI)

b) What are the advantages of BIRCH compared to other methods (Mod VI)

a) BIRCH - Balanced Iterative Reducing and Clustering using Hierarchies.

- It is a scalable clustering method
- Designed for very large datasets.
- Only one scan of data is necessary.
- It is based on the notation of EF (Clustering Feature) a CF tree.
- EF tree is a height balanced tree that stores the clustering features for a hierarchical clustering.
- Cluster of data points is represented by a triple of members (N, LS, SS) where

N - no. of items in the sub cluster (zeroth moment of cluster)

LS - Linear sum of the points $\sum_{i=1}^n x_i$ (1st moment of cluster)

SS - sum of the squared of the points $\sum_{i=1}^n x_i^2$ (2nd moment of cluster)

→ A CF tree structure is given as below:

2 parameters

* Each non-leaf node has at most B entries

* Branching Factor

max no. of children per non-leaf node

* Each leaf node has at most L CF centers which satisfies threshold T , a maximum diameter of each cluster.

* threshold

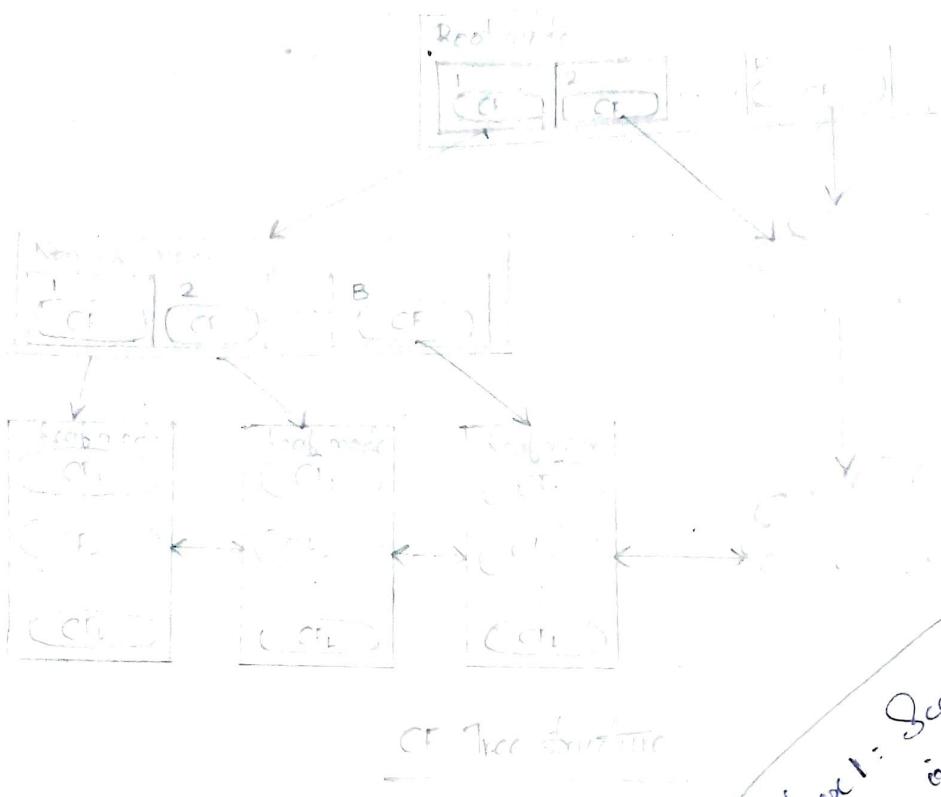
max diameter of subcluster stored at leaf node of the tree

* P (page size in bytes) is the maximum size of a node

* Compact : each ~~node~~ leaf node is a subcluster, not a data point.

3 points $(2, 5), (3, 2), (4, 3)$ in a cluster C_1 .

$$CF = \langle 3, ((2+3+4, 5+2+3), (2^2+3^2+4^2, 5^2+2^2+3^2)) \rangle = \langle 3, (9, 10), (29, 36) \rangle$$



Basic algorithm

→ Phase I : Load data into mlg

Scan DB and load data into mlg by building a CF tree.
If mlg is exhausted rebuild the tree from the leaf node.

→ Phase II : Condense data (optional)

Resize the dataset by building a smaller CF tree

Remove more outliers

Condensing is optional

→ Phase III : Global clustering

Use existing clustering algorithms (eg KMEANS, Hc)
on CF entries

→ Phase IV : Cluster refining (optional)

Refining is optional

Solves the problem with CF trees where same valued data points may be assigned to diff. leaf entries.

CF tree is dynamically built as obj are inserted

→ Obj is inserted into closest leaf entry → If diameter of the Subcluster stored in leaf node after insertion is larger than threshold value, then leaf & possibly other nodes split

→ Size of CF tree is determined threshold

→ If size of tree is larger reduce threshold

→ Similar to conversion & node split in B+ tree complexity. $O(n)$, no. of obj to be clustered

Phase I : Scans database to build initial CF tree
Phase II : Applies a clustering algorithm to the leaf nodes of CF tree,
* removes some clusters
* groups some clusters to form one cluster

Data



Phase 1: Load data internally by building a CF tree

Initial CF tree

Phase 2: Condense data by building a smaller CF tree
(optional)

Smaller CF tree



Phase 3: Global clustering

Good cluster



Phase 4: Cluster refinement (optional)

Better cluster



D) Advantages of BIRCH method

- Finds a good clustering with a single scan and improve the quality with a few additional scans.
- It is local in that each clustering decision is made without scanning all data points and currently existing clusters.
- It ~~never~~ exploits the observation that data space is not usually uniformly occupied and not every data point is equally important.
- It makes full use of available memory to derive the finest possible sub-clusters while minimizing I/O costs.
- It is also an incremental method that does not require the whole dataset in advance.

Disadvantages:- Handles only numeric data

Q8. 19. a) Explain K-means partition algorithm. What is the disadvantage of K-means? (Mod VI)

b) Term frequency matrix given in the table shows frequency of terms per document. Calculate the TF-IDF value for the term T_4 in document 3.

Document/ Term	T_1	T_2	T_3	T_4	T_5	T_6
D1	5	9	4	0	5	6
D2	0	8	5	3	10	8
D3	3	5	6	6	5	0
D4	4	6	7	8	4	4

(Mod VI)

a) K-means algorithm takes the input parameter k , and partitions a set of n objects into k clusters so that the resulting intracluster similarity is high but the intercluster similarity is low.

Cluster similarity is measured in regard to the mean of the objects in a cluster which can be viewed as the cluster's centroid or centre of gravity.

Algorithm

Input :- k - The no: of clusters

D - Dataset containing n objects

Output :-

A set of arbitrary clusters

- i) Arbitrarily choose k objects from D as the initial cluster centres
- ii) Repeat
- iii) (Re) Assign each object to the cluster to which the object is similar, based on the mean value of the objects in the cluster
- iv) Update the cluster means; calculate the mean value of the objects for each cluster.
- v) Until no change.

First it randomly select k of the objects, each of which initially represents a cluster mean or center.

For each of the remaining objects, an object is assigned to the cluster to which it is the most similar, based on the distance b/w the object and cluster mean.

It then computes the mean for each cluster.

This process iterates until the exterior function

Typically square error exterior is used

$$E = \sum_{i=1}^k \sum_{P \in C_i} |P - m_i|^2$$

P - represents an element in cluster C_i

m_i - mean of cluster C_i

Computational complexity

$O(nkt)$

n - no. of objects

k - no. of clusters

t - no. of iterations

Disadvantages of K-means

- Cannot be applied for data containing categorical attribute
- Users need to specify the no of clusters in advance
- Not suitable for discovering clusters with non convex shapes or clusters of very different size
- It is sensitive to noise and outlier datapoints

19 b) Term Frequency (TF) - Measures how frequently a term occurs in a document

$$TF(d, t) = \begin{cases} 0 & , \text{ if } freq(d, t) = 0 \\ 1 + \log(1 + \log(freq(d, t))) & , \text{ otherwise} \end{cases}$$

$t \rightarrow$ terms

$d \rightarrow$ document given

Inverse Document Frequency (IDF) - Measures how important a term is. It represents the scaling factor or importance of t .

$$IDF(t) = \log \frac{1 + |d|}{|d_t|}$$

where d is the document collection

d_t is the set of documents containing the term t

$$TF-IDF(d, t) = TF(d, t) \times IDF(t)$$

$$\begin{aligned}
 TF(d_3, t_4) &= 1 + \log(1 + \log(\text{freq}(d_3, t_4))) & 51 \\
 &= 1 + \log(1 + \log(6)) \\
 &= 1.249968
 \end{aligned}$$

$$\begin{aligned}
 IDF(t_4) &= \log \frac{1 + |d|}{|d_t|} \\
 &= \log \frac{1+4}{3} \\
 &= \log \frac{5}{3} \\
 &= 0.22184
 \end{aligned}$$

$d \rightarrow$ total no. of document = 4
 D_1, D_2, D_3, D_4
 $d_t \rightarrow$ set of documents containing t_4
 $\{d_2, d_3, d_4\}$
 $\therefore d_t = 3$

$$\begin{aligned}
 TF-IDF(d_3, t_4) &= TF(d_3, t_4) \times IDF(t_4) \\
 &= 1.249 \times 0.221 \\
 &= \underline{\underline{0.2760}}
 \end{aligned}$$

OCTOBER 2019

PART A

7. Differentiate b/w support and confidence (Mod V)

May 2019 PART A Qn 7

8. How to compute the dissimilarity b/w objects described by binary variables? (Mod V)

→ The distance between objects may be calculated based on a contingency table.

→ A binary attribute is symmetric if both of its states are equally ~~valuable~~ valuable

→ In that case, using the simple matching coefficient can assess dissimilarity b/w 2 objects:

$$d(x_i, x_j) = \frac{q + s}{q + r + s + t}$$

where,

q = no: of attributes that equal 1 for both obj

t = no: of attributes that equal 0 for both obj

r and s = no: of attributes that are unequal for both obj.

→ A binary attribute is unimetric, because if states are not equally important (usually the outcome is considered more important)

→ In this case, the denominator ignores the unimportant negative matches (t). This is called Jaccard coefficient

$$d(x_i, x_j) = \frac{r + s}{q + r + s}$$

a. Differentiate b/w Agglomerative and Divergent hierarchical clustering method. (Mod VI)

Agglomerative Hierarchical Clustering

This bottom-up strategy starts by placing each object in its own cluster and then merges these atomic clusters into larger and larger clusters, until all of the objects are in a cluster or until certain termination conditions are satisfied.

Most hierarchical clustering methods belong to this category.

They differ only in their definition of inter-cluster similarity.

Divisive Hierarchical Clustering

This top-down strategy does the reverse of agglomerative hierarchical clustering by starting with all objects in one cluster. It subdivides the cluster into smaller and smaller pieces, until each object forms a cluster on its own or until it satisfies certain termination conditions, such as a desired number of clusters is obtained or the distance b/w the two closest clusters is above a certain threshold distance.

More complex, more efficient and more accurate than that of agglomerative hierarchical clustering.

10 Explain Web Content mining? (Mod VI)

May 2019 PART A Qn 10.

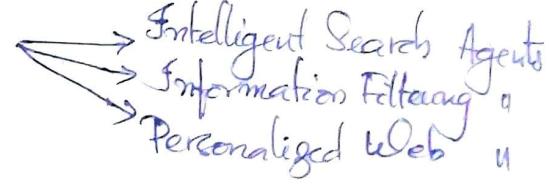
The two ~~spcl~~ approaches to web content mining

i) Agent Based Approach

Usually involves intelligent agents that can act autonomously or semi autonomously

If usually relies on autonomous agents (or bots) that can parse through the web, identify websites that are relevant and then collect information from them.

These agents can be of 3 types

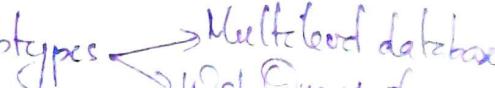


ii) Data-based Approach

Used to organize semi-structured data present on the internet into more structured data sources

Once this done, standard data querying techniques can be used to harness the data.

The approach can have two subtypes



PART D

- Q 14. Consider the transaction database given below. Set minimum support count as 2 and minimum confidence threshold as 70%.

Transactions ID	List of Item IDs
T100	I ₁ , I ₂ , I ₅
T200	I ₂ , I ₄
T300	I ₂ , I ₃
T400	I ₁ , I ₂ , I ₄
T500	I ₁ , I ₃
T600	I ₂ , I ₃
T700	I ₁ , I ₃
T800	I ₁ , I ₂ , I ₃ , I ₅
T900	I ₁ , I ₂ , I ₃

- a) Find frequent itemset using Apriori Algorithm (Mod V)
 b) Generate strong association rules. (Mod V)

- i) Apriori employs an iterative approach known as a level wise search where k itemsets are used to explore (k+1) itemsets.
 Apriori property - All non empty subsets of a frequent itemset must also be frequent.
 A two-step process is followed.

1. The join step : To find L_k, a set of candidate k-itemsets is generated by joining L_{k-1} with itself. The set of candidates is denoted by C_k.
2. The prune step : C_k is a superset of L_{k-1}, i.e., its members may or may not be frequent, but all of the frequent k-itemsets are included in C_k. To reduce the size of C_k the Apriori property is used.

Minimum support count = 2

Step 1: Scan Database D for count of each candidate.

C ₁	Itemset	Supcount
Creating 1-frequent item	I ₁	6
	I ₂	7
	I ₃	6
	I ₄	2
	I ₅	2
	I ₆	

L ₁	Itemset	Supcount
Compare candidate support count with minimum support count	I ₁	6
	I ₂	7
	I ₃	6
	I ₄	2
	I ₅	2

Creating 2-frequent itemset, C₂, created from L₁.

Generating C ₂ from L ₁	I ₁	I ₂	Scan D for each candidate's count.	C ₂
	I ₁	I ₃		I ₁ I ₂ 4
	I ₁	I ₄		I ₁ I ₃ 4
	I ₁	I ₅		I ₁ I ₄ 1
	I ₂	I ₃		I ₁ I ₅ 2
	I ₂	I ₄		I ₂ I ₃ 4
	I ₂	I ₅		I ₂ I ₄ 2
	I ₃	I ₄		I ₂ I ₅ 2
	I ₃	I ₅		I ₃ I ₄ 0

L ₂	Itemset	Supcount
Compare candidate support count with minimum sup-count	I ₁ I ₂	4
	I ₁ I ₃	4
	I ₁ I ₄	1
	I ₁ I ₅	2
	I ₂ I ₃	4
	I ₂ I ₄	2

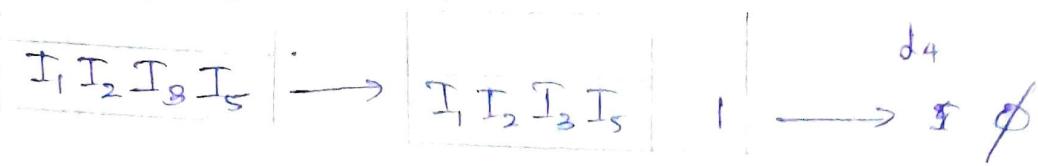
Creating 3-frequent itemset C₃ from L₂.

Generating C ₃ from L ₂	I ₁ I ₂ I ₃	Scan D for count of each candidate	C ₃
	I ₁ I ₂ I ₅		I ₁ I ₂ I ₃ 2
	I ₁ I ₂ I ₄		I ₁ I ₂ I ₅ 2
	I ₁ I ₃ I ₅		I ₁ I ₂ I ₄ 1
	I ₂ I ₃ I ₄		I ₁ I ₃ I ₅ 1
	I ₂ I ₃ I ₅		I ₂ I ₃ I ₄ 0

L ₃	Itemset	Supcount
Compare candidate support count with minimum sup-count	I ₁ I ₂ I ₃	2
	I ₁ I ₂ I ₅	2

Generating L_4 with 4-frequent items

96



L_4 is empty so L^n is $L_{4-1} = L_3$.

So freq. itemsets are $\{I_1 I_2 I_3\}$ & $\{I_1 I_2 I_5\}$.

So if $I_1 I_2 I_3$ is a frequent itemset then all its non-empty subsets must also be frequent, while checking we can see,

I_1 ; I_2 ; I_3 ; $I_1 \cap I_2$; $I_1 \cap I_3$; $I_2 \cap I_3$;

$I_1 \cap I_2 \cap I_3$ is also frequent. By $I_1 I_2 I_5$.

b)

Generate strong association rules from itemsets (where strong association rules satisfy both min support and min confidence) association rules can be generated as follows:

- For each frequent itemset l , generate all nonempty subsets of l .
- For every nonempty subset s of l , output the rule " $s \Rightarrow (l-s)$ "
if $\frac{\text{support-count}(l)}{\text{support-count}(s)} \geq \text{min-conf}$, where min-conf is the minimum confidence threshold.

$$\text{confidence } (A \Rightarrow B) = P(B|A) = \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)}$$

Taking any one from $\{I_1 I_2 I_3\}$ & $\{I_1 I_2 I_5\}$.

Taking $I_1 I_2 I_5$.

Generating all non empty subsets from $I_1 I_2 I_5$.

I_1 , I_2 , I_5 , $I_1 \cap I_2$, $I_1 \cap I_5$,

$I_2 \cap I_5$, $I_1 I_2 I_5$

Support count of $l = \{I_1 I_2 I_5\} = 2$ (from L_3)

Rule	3	$I_1 \rightarrow I_2 \cap I_5$	Confidence $2/6 = 33.3\%$
1	I_1	$\rightarrow I_2 \cap I_5$	$2/4 = 50\%$
2	I_2	$\rightarrow I_1 \cap I_5$	$2/4 = 50\%$
3	I_5	$\rightarrow I_1 \cap I_2$	$2/2 = 100\%$
4	$I_1 \cap I_2$	$\rightarrow I_5$	$2/4 = 50\%$
5	$I_1 \cap I_5$	$\rightarrow I_2$	$2/2 = 100\%$
6	$I_2 \cap I_5$	$\rightarrow I_1$	$2/2 = 100\%$
	$I_1 \cap I_2 \cap I_5$	\emptyset	So not considering since it empty in RHS.

Min. conf = 70% (Given in ques).

Here Rule 3, 5, 6 has confidence > 70%.

So strong association rules are:

$$I_5 \rightarrow I_1 \cap I_2$$

$$I_1 \cap I_5 \rightarrow I_2$$

$$I_2 \cap I_5 \rightarrow I_1$$

18 a) Explain DBSCAN algorithm (Mod VI)

b) State the pros and cons of DBSCAN method. (Mod VI)

a) A density-based clustering method based on connected regions with sufficiently high density — DBSCAN.

DBSCAN is a density-based clustering algorithm. The algorithm grows regions with sufficiently high density into clusters, and discovers clusters of arbitrary shape in spatial databases with noise. It defines a cluster as a maximal set of density-connected points.

The basic idea of density based clustering involve a number of new definition \Rightarrow The neighborhood within within a radius ϵ of a given object is called **ϵ -neighborhood of the object**.

\rightarrow If the ϵ -neighborhood of an object contains atleast a minimum number, MinPts, of objects, then the object is called a **core object**.

\rightarrow Given a set of objects, D , we say that an object ' p ' is **directly density reachable** from object ' q ' if p is within the ϵ -neighborhood of q , and ' q ' is a core object.

\rightarrow An object p is **density reachable** from ~~from~~ object q wrt ϵ and MinPts in a set of objects D , if there is a chain of objects P_1, P_2, \dots, P_n where $P_1 = q$ and $P_n = p$ such that P_{i+1} is directly density reachable from P_i wrt ϵ and MinPts, for $1 \leq i \leq n$, $P_i \in D$.

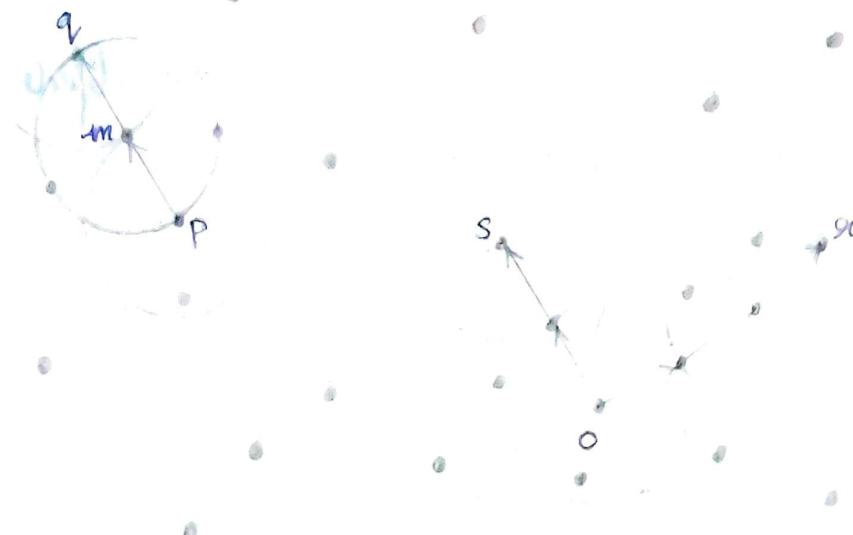
\rightarrow An object ' p ' is **density connected** to object ' q ' wrt ϵ and MinPts in a set of objects, D , if there is an object $o \in D$ such that both p and q are density reachable from o wrt ϵ and MinPts.

Density reachable is the transitive closure of direct density reachability and this relationship is asymmetric.

Only core objects are mutually density reachable.

Density connectivity, however, is a symmetric relation.

Consider the fig for a given ϵ represented by the radius of the circles, and let MinPts = 3.



Density reachability & Density connectivity

- Of the labeled points M, P, Q, and R. are core objects since each is in an ϵ -neighborhood containing at least 3 points.
- M is directly reachable from P and Q is directly density-reachable from M.
- Based on the previous observation Q is (indirectly) density-reachable from P. However P is not density reachable from Q. Similarly R and S are density reachable from O.
- O, Q and S are all density connected.

A density cluster is a set of density-connected objects that is maximal with respect to density reachability. Every object not contained in any cluster is considered to be noise.

How does DBSCAN find clusters?

DBSCAN checks the neighborhood of each point in the database. If the ϵ -neighborhood of a point p contains more than N_Pts, a new cluster with p as core obj is created. DBSCAN then iteratively collects directly density-reachable objects from these core objects, which may involve the merge of a few density-reachable clusters. The process terminates when no new point can be added to any cluster.

Algorithms :

I/p : D : a data set containing m objects

ϵ : the radius parameter

Min Pts : the neighborhood density threshold

O/p : A set of density based clusters.

Method

- 1) Mark all objects as unvisited;
- 2) do
- 3) randomly select an unvisited object p
- 4) mark p as visited;
- 5) if the ϵ -neighborhood of p has atleast Min Pts objects
- 6) create a new cluster C and add p to C
- 7) let N be the set of objects in the ϵ -neighborhood of p
- 8) for each point p' in N
- 9) if p' is unvisited
- 10) mark p' as visited
- 11) if the ϵ -neighborhood of p' has atleast Min Pts points
- 12) add those points to N ;
- 13) if p' is not yet a member of any cluster
- 14) add p' to C
- 15) end for
- 16) output C
- 17) else mark p as noise
- 18) until no object is unvisited.