

Semester: V

Time: 1 Hr. & 15 Mints

Branch: CMPN

Subject: DWM

Marks: 30

Q. 1)	Attempt any Five (2 Marks Each)	CO	BL
a)	Define the terms Data Mining & KDD.	2	1
b)	Consider a dataset, $D = \{8, 11, 12, 12, 23, 29, 29, 35, 38, 45, 46, 66\}$ . Give five number summary of this dataset.	2	3
c)	Consider two objects, $A = \{19, 22, 07, 34\}$ and $B = \{23, 12, 17, 30\}$ . Calculate Minkovski distance ( $b=3$ ) and supremum distance between these objects.	2	3
d)	Consider a dataset, $D = \{4, 8, 15, 21, 21, 24, 25, 28, 34\}$ . Perform smoothing by bin boundaries. Let there be 3 bins and are equal frequency.	2	3
e)	Normalize the following group of data: 200,300,400,600,1000 using min-max normalization by setting min=0 and max=1	2	3
f)	What are the limitations of ID3 tree induction algorithm?	3	2
g)	Define the terms: Precision & Recall with respect to evaluation of classifiers	4	1
h)	Define the terms: Support & Confidence with respect to Association Analysis	3	1

Q. 2)	Attempt any two. (5 Marks Each)																																																								
a)	Explain the architecture of a typical Data Mining system with neat diagram.	2	1																																																						
b)	Consider the following training set to calculate the Information gain for any one attribute using attribute relevance analysis.	3	3																																																						
	<table border="1"> <thead> <tr> <th>Name</th> <th>Hair</th> <th>Height</th> <th>Weight</th> <th>Dublin</th> <th>Result</th> </tr> </thead> <tbody> <tr> <td>Sarah</td> <td>Blonde</td> <td>Average</td> <td>Light</td> <td>No</td> <td>Sunburned</td> </tr> <tr> <td>Dana</td> <td>Blonde</td> <td>Tall</td> <td>Average</td> <td>Yes</td> <td>None</td> </tr> <tr> <td>Alex</td> <td>Brown</td> <td>Short</td> <td>Average</td> <td>Yes</td> <td>None</td> </tr> <tr> <td>Annie</td> <td>Blonde</td> <td>Short</td> <td>Average</td> <td>No</td> <td>Sunburned</td> </tr> <tr> <td>Emily</td> <td>Red</td> <td>Average</td> <td>Heavy</td> <td>No</td> <td>Sunburned</td> </tr> <tr> <td>Pete</td> <td>Brown</td> <td>Tall</td> <td>Heavy</td> <td>No</td> <td>None</td> </tr> <tr> <td>John</td> <td>Brown</td> <td>Average</td> <td>Heavy</td> <td>No</td> <td>None</td> </tr> <tr> <td>Katie</td> <td>Brown</td> <td>Short</td> <td>Light</td> <td>Yes</td> <td>None</td> </tr> </tbody> </table>	Name	Hair	Height	Weight	Dublin	Result	Sarah	Blonde	Average	Light	No	Sunburned	Dana	Blonde	Tall	Average	Yes	None	Alex	Brown	Short	Average	Yes	None	Annie	Blonde	Short	Average	No	Sunburned	Emily	Red	Average	Heavy	No	Sunburned	Pete	Brown	Tall	Heavy	No	None	John	Brown	Average	Heavy	No	None	Katie	Brown	Short	Light	Yes	None		
Name	Hair	Height	Weight	Dublin	Result																																																				
Sarah	Blonde	Average	Light	No	Sunburned																																																				
Dana	Blonde	Tall	Average	Yes	None																																																				
Alex	Brown	Short	Average	Yes	None																																																				
Annie	Blonde	Short	Average	No	Sunburned																																																				
Emily	Red	Average	Heavy	No	Sunburned																																																				
Pete	Brown	Tall	Heavy	No	None																																																				
John	Brown	Average	Heavy	No	None																																																				
Katie	Brown	Short	Light	Yes	None																																																				
c)	Consider the training set in Q2 B. Classify a new sample $X=\{\text{Blonde, Short, Light, No}\}$ using the Naïve Bayes Classifier.	3	3																																																						

Q 3)	Attempt any One (10 Marks Each)		
a)	Explain different steps involved in data preprocessing	2	1
b)	Explain different data mining tasks in brief.	2	1

C01	Understand data warehouse fundamentals, design dimensional model and apply OLAP operations.
C02	Understand data mining principles, perform Data preprocessing and Visualization.
C03	Identify & implement appropriate data mining algorithms to solve real world problems
C04	Compare and evaluate different data mining techniques like classification, prediction, clustering and association rule mining
C05	Describe complex information and social networks with respect to web mining.

Branch  
CmpnTest Date  
rSemester  
VDiv.  
A

Roll No.

Student's Signature  
YH

# Test No.

mSE

Subject  
Dwm

Junior Supervisor's full signature with date :

Question No.	1	2	3	Total 20	Examiners Signature	Student's Sign After receiving the assessed answer sheet
Marks obtained						

- (a) Define the term Data mining and KDD  
 Data Mining is the process to find patterns by analysing large databases. The patterns should be
- (a) valid: hold on new data with some certainty
  - (b) novel: non-obvious to the system
  - (c) useful: should be possible to act on the item
  - (d) understandable: humans should be able to interpret the pattern

KDD: Knowledge discovery in Databases also called as Datamining is the process of discovering useful knowledge from a collection of data.

- (b) Consider a dataset,  $D = \{8, 11, 12, 12, 23, 29, 29, 35, 38, 45, 46, 66\}$   
 Give five number summary of this dataset

→ Dataset  $D = \{8, 11, 12, 12, 23, 29, 29, 35, 38, 45, 46, 66\}$

The median :- As the <sup>total</sup> number of ~~n~~ is even  

$$\frac{29+29}{2} = 29$$

First Quartile =

$$Q_1 = \{8, 11, 12, 12, 23, 29\} \\ = (12+12)/2 = 12$$

Third Quartile

$$Q_3 = \{29, 35, 38, 45, 46, 61\} \\ (38+45)/2 = 42.$$

minimum = ~~19~~ 8

maximum = 66.

So the Five number Summary is ~~{19, 12, 29, 42, 61}~~

(Q) Consider two objects

$$A = \{19, 22, 07, 34\} \text{ & } B = \{23, 12, 17, 30\}$$

Calculate Min Kowalski distance ( $h=3$ ) &  
supremum distance between these objects

Min Kowalski distance formula

$$D(x_i, x_j) = \left( \sum_{l=1}^d |x_{il} - x_{jl}|^p \right)^p$$

~~1st~~  $h=3$  i.e  $p=3$   
 $= 12.862$

The supremum distance = ~~12~~ 21

(c) Consider a dataset  $D = \{4, 8, 15, 21, 21, 24, 24, 25, 28, 34\}$ , perform smoothing by bin boundaries. Let these be 3 bins and are equal frequency.

$$\begin{aligned} \text{Bin 1 : } & \{4, 8, 15\} \\ \text{Bin 2 : } & \{21, 21, 24\} \\ \text{Bin 3 : } & \{25, 28, 34\} \end{aligned}$$

Smoothing by bin

$$\text{Bin 1 : } \{4, 8, 15\}$$

$$\text{Bin 2 : } \{21, 21, 24\}$$

$$\text{Bin 3 : } \{25, 28, 34\}$$

e) Normalize the following group of data.  
 $200, 300, 400, 600, 1000$ .

min-max normalization by setting  $\min = 200$  &  $\max = 1000$

$$v_1 = \frac{(200 - 200) \times (1-0) + 0}{1000 - 200} = 0$$

$$v_2 = 300 = \frac{(300 - 200) \times (1-0) + 0}{1000 - 200} = 0.125$$

$$v_3 = 400 = \frac{(400 - 200) \times (1-0) + 0}{1000 - 200} = 0.25$$

$$V_1 = 600$$

$$\frac{600 - 200}{(1000 - 200)} \times (1 - 0) + 0 = 0.5$$

$$V_1 = 1000$$

$$\frac{1000 - 200}{(1000 - 200)} \times (1 - 0) + 0 = 1$$

After normalization.

### (ii) Limitations of ID3

The Limitation of ID3 is overfitting which can be overcome by using Pruning method.

Pruning method consists of two types.

(1) Pre-pruning

(2) Post-pruning.

### (iii) Precision & Recall

Precision is a good measure to determine the cost of false positive emails. When the cost of false positive emails is high, a spam detection model may detect a email as not spam but detected spam. The email user might lose important emails if the precision is not high for the spam detection model.

Precision is

$$\frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

(a) Define the term : support & confidence with respect to association analysis

→ A support of 1% means that 1% of all the transactions under analysis shows that computer & s/w are purchased together.

Support ( $x \Rightarrow y$ )

No of tuples that satisfies both  $x$  &  $y$   
Total no of tuples

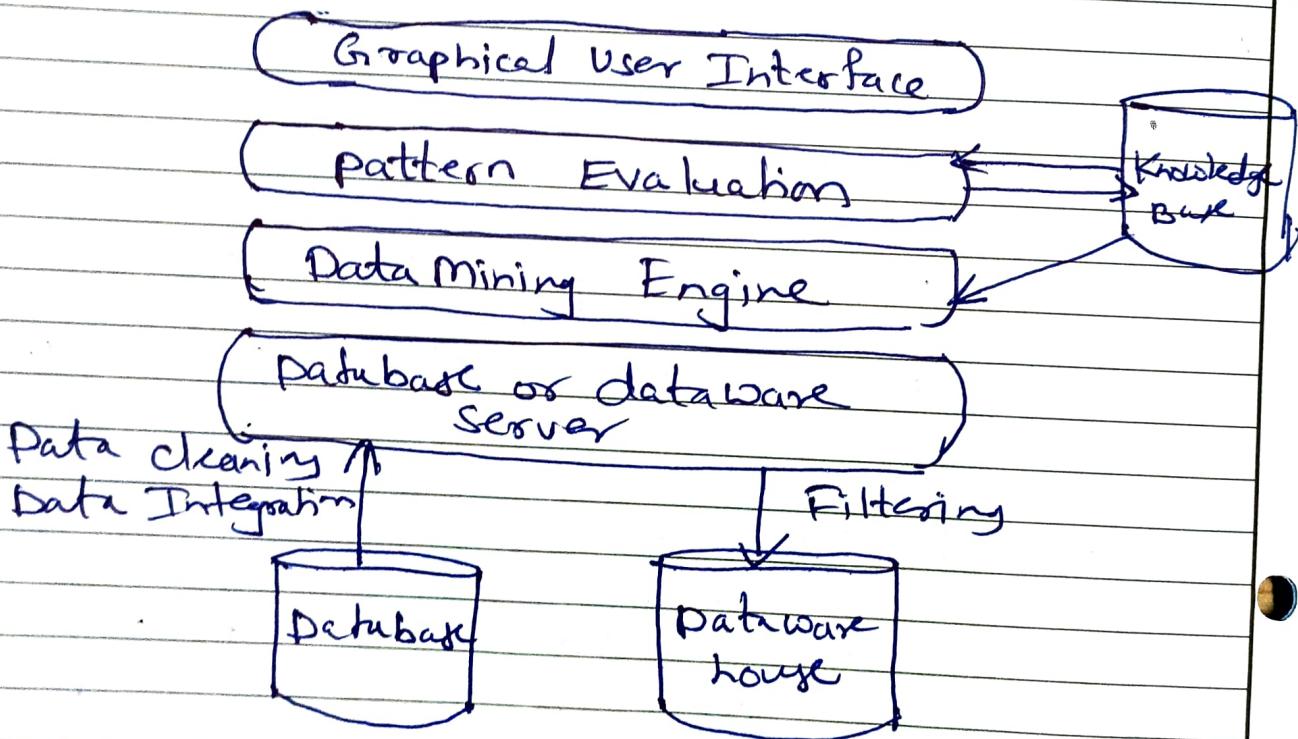
Confidence ( $x \Rightarrow y$ )

No of tuples that satisfy  $bzny$   
No of tuples that satisfy  $x$ .

(a) Explain the architecture of a typical DSS system with neat diagram  
A DB is used to extract required data warehouse by executing SQL queries. If data is selected from a data warehouse then no pre-processing is required because a data warehouse contains cleaned, transformed and integrated data.

A data mining engine is used to extract all hidden pattern using intelligent algorithms

- (3) Additional information required to form a particular Data mining task is called back ground knowledge and is available from knowledge base.
- (4) The pattern Evaluation Engine and the datamining engine are normally combined in a single component
- (5) A GUI is used to convert extracted knowledge into required form. It allows users to interact with Datamining systems.



The Architecture of a typical DM system

(b) For given data set calculate the IG for any one attribute using attribute relevance analysis

$$X = \{ \text{Blonde, short, light, No} \}$$

The IG for the entire dataset using ID3 algorithm.

As there are ~~2~~<sup>None</sup> classes ~~2~~<sup>2</sup> son bar = 3

$$\text{Info}(D) = -\frac{5}{8} \log_2 \frac{5}{8} - \frac{3}{8} \log_2 \frac{3}{8}$$

\$ For attribute

$$\text{Info}_{\text{hair}}(D) = \left[ \frac{3}{8} \times \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) \right] +$$

$$\begin{aligned}
 & \text{Hair} \\
 & \quad \text{Red} \\
 & \quad \text{Blonde, Brown} \quad (1) \\
 & \quad (3) \quad (4) \\
 & \quad \text{No} \\
 & \quad \text{son} \quad (1) \quad 0 \quad (4) \\
 & \quad \left[ \frac{4}{8} \times \left( -\frac{4}{4} \log_2 \frac{4}{4} \right) + 0 \right] + \\
 & \quad \left[ \frac{1}{8} \times \left( -\frac{1}{1} \log_2 \frac{1}{1} \right) \right]
 \end{aligned}$$

$$IG_{\text{hair}} = \text{Info}(D) - \text{Info}_{\text{hair}}(D)$$

$$= 0.454$$

(Q2) classify a new sample  $X = \{ \text{Blonde, Short, Light, No} \}$  using Naive Bayes

$\Rightarrow$  Class =

$$C_1 = \text{Sunburned} = 3$$

$$C_2 = \text{None} = 5$$

$$P(\text{Result} = \text{Sunburned}) = 3/8$$

$$P(\text{Result} = \text{None}) = 5/8$$

For  $X = \text{Blonde, Short, Light, no}$

~~Prob~~ Result

$$P(\text{Hair} = \text{Blonde} | \text{Result} = \text{Sunburned}) = \frac{2}{3}$$

$$P(\text{Height} = \text{Short} | \text{Result} = \text{Sunburned}) = \frac{1}{3}$$

$$P(\text{Weight} = \text{Light} | \text{Result} = \text{Sunburned}) = \frac{1}{3}$$

$$P(\text{Dublin} = \text{No} | \text{Result} = \text{Sunburned}) = \frac{3}{3} = 1$$

$$P(\text{hair} = \text{Blonde} | \text{Result} = \text{None}) = \frac{1}{5}$$

$$P(\text{Height} = \text{Short} | \text{Result} = \text{None}) = \frac{2}{5}$$

$$P(\text{Weight} = \text{Light} | \text{Result} = \text{None}) = \frac{1}{5}$$

$$P(\text{Dublin} = \text{No} | \text{Result} = \text{None}) = \frac{2}{5}$$

$$P(X | \text{Result} = \text{Sunburned}) = \frac{3}{8} \times \frac{2}{3} \times \frac{1}{3} \times \frac{1}{3} = 0.0277$$

$$P(X | \text{Result} = \text{None}) = \frac{5}{8} \times \frac{1}{3} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5} = 0.004$$

as probability of  $C_1 = \text{sunburned}$  is greater than probability of  $C_2 = \text{None}$

$\therefore$  Unknown tuple ' $X'$  belongs to class  $C_1 = \text{Sunburned}$ .

Branch	Test Date	Semester	Div.	Roll No.	Student's Signature
CMPN		V	A		<u>KY</u>

IA Test No. <u>mst</u>	Subject <u>Dwm</u>
---------------------------	-----------------------

Junior Supervisor's full signature with date :	Question No.	1	2	3	Total 20	Examiners Signature	Student's Sign After receiving the assessed answer sheet
	Marks obtained						

- (Q3.A) Explain different steps involved in Preprocessing
- (1) Data cleaning : Data in the real world is dirty as there is lots of potentially incorrect data.  
eg occupation = " " missing data  
noisy :- containing noise, errors or outliers.  
eg salary = 10' (an error)  
Age = 42, BD = 03/02/2010.  
To handle missing value.
- (2) Ignore the tuple
- (2) Use the attribute mean to fill the missing value
- (3) predict the missing value by using a learning algorithm.
- (4) Identify the outliers & smooth out noisy data using following techniques
- (1) Binning  
(2) Clustering  
(3) Regression  
(4) correct inconsistent data.

2) Data Integration:-  
Combine data from multiple sources into a coherent store.  
It merge the data from multiple heterogeneous data source into a coherent data store. Data Integration may involve inconsistent data and therefore needs data cleaning.  
Redundant data occur when integration of multiple database that may be able to be detected by ~~co~~ correlation analysis and covariance analysis.

3) Data Reduction:-  
obtain a reduced representation of the data set that is much smaller in volume, but yet produce the same (or almost the same) analytical result  
Different data reduction strategies are

- (1) Dimensionality reduction: eg remove unimportant attributes. It can be done by using Wavelet transforms
- (2) Principal Component Analysis (PCA)
- (3) Feature subset selection or feature creation
- (4) Parametric data reduction: Using Linear regression and Log Linear model

Parametric methods assume the data fits only the parameters and discard the data dimensionality reduction to avoid the curse of dimensionality; help eliminate irrelevant features and reduce

noise. It also allow easier visualization.

(4) Data Transformation: A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values

Different methods are  
(i) smoothing, remove noise from data

(ii) Attribute / feature construction  
new attribute constructed from the given ones

(iii) Aggregation: Summarization, data cube construction

(iv) Normalization: scaled to fall within a smaller specified range

methods of Normalization are  
(a) min-max normalization  
(b) Z-score normalization

(v) normalization by decimal scaling

(5) Discretization: Concept hierarchy climbing.

Different methods ~~can~~ can be applied for recursively

(i) Binning: Top-down split, unsupervised

(ii) Histogram analysis: Top down split or bottom up merge

(iii) Clustering analysis:

(ii) Decision tree Analysis

(v) Correlation (e.g.  $\chi^2$  analysis); - Bottom up merged.

Q3) Explain different data mining task in brief

(ii) Classification and Regression for predictive Analysis

→ classification is the process of finding a model that describes and distinguishes data classes or concepts. The model are derived based on the analysis of a set of training data. i.e. Data objects of which class labels are known). The model is used to predict the class label of objects for which the class label is unknown.  
Eg:- Classify countries based on Climate

methods used are decision trees, classification, SVM, naive Bayes,

Regression Analysis:- is a statistical methodology that is most often used for numerical prediction. Regression models continuous valued function i.e regression is used to predict missing or unavailable numerical data values.

### 3) Clustering Analysis:

Clustering analyzes data objects without consulting class labels.

In many cases, class labels data may simply not exist at the beginning. Clustering can be used to generate class labels for a group of data. The objects are clustered or grouped based on the principle of maximizing the intraclass similarity and minimizing the interclass similarity.

i.e. clusters of objects are similar within a cluster and dissimilar to objects in other cluster

### Association Analysis:-

Association analysis is useful for discovering interesting relationships hidden in large datasets. The uncover relationship can be represented in the form of association rules or set of frequent items.

Applications of association Analysis  
Market Basket Analysis

$\text{buys}(X, \text{'computer'}) \Rightarrow \text{buy}('X', 'S')$

To measure support and confidence are considered.

support  $(x \Rightarrow y)$  = No of tuples that satisfy both  $x \wedge y$  / Total No of tuples

Confidence ( $x \Rightarrow y$ ) =

$$\frac{\text{No of tuples that satisfy both } x \text{ and } y}{\text{No of tuples that satisfy } x}$$

(1) Single dimensional association rule:

Association rule that contain a single predicate are referred to single dimensional association rule.

$$\text{buys}(x, \text{'computer'}) \Rightarrow \text{buys}(x, \text{'s/w'})$$

(2) Multi dimensional association rule:

The association rule that contain multiple predicates / attributes are referred to multidimensional association rule.

$$\text{e.g. } \text{age}(x, \text{'20-29'}) \wedge \text{income}(x, \text{'40k-44k'}) \\ \Rightarrow \text{buys}(x, \text{'laptop})$$

(3) Outlier Analysis..

A dataset may contain objects that do not comply with the general behaviour or model of the data. These data objects are outliers. Many applications discard outliers as noise or exception.

However, in some applications (e.g. fraud detection) the rare events can be more interesting than the more regularly occurring ones. The

analysis may & uncover fraudulent usage of credit cards by detecting purchases of unusually large amount for a given account number.