

Practical No.1

Aim: To build Data warehouse and explore WEKA.



Date: 01/02/23

Practical No. 1

Aim: To build Data Warehouse and Explore WEKA.

Theory:

Data warehouse :

Data warehousing is the process of constructing and using a data warehouse. A data warehouse is constructed by integrating data from multiple heterogeneous sources that support analytical reporting, structured and for ad hoc queries, and decision making.

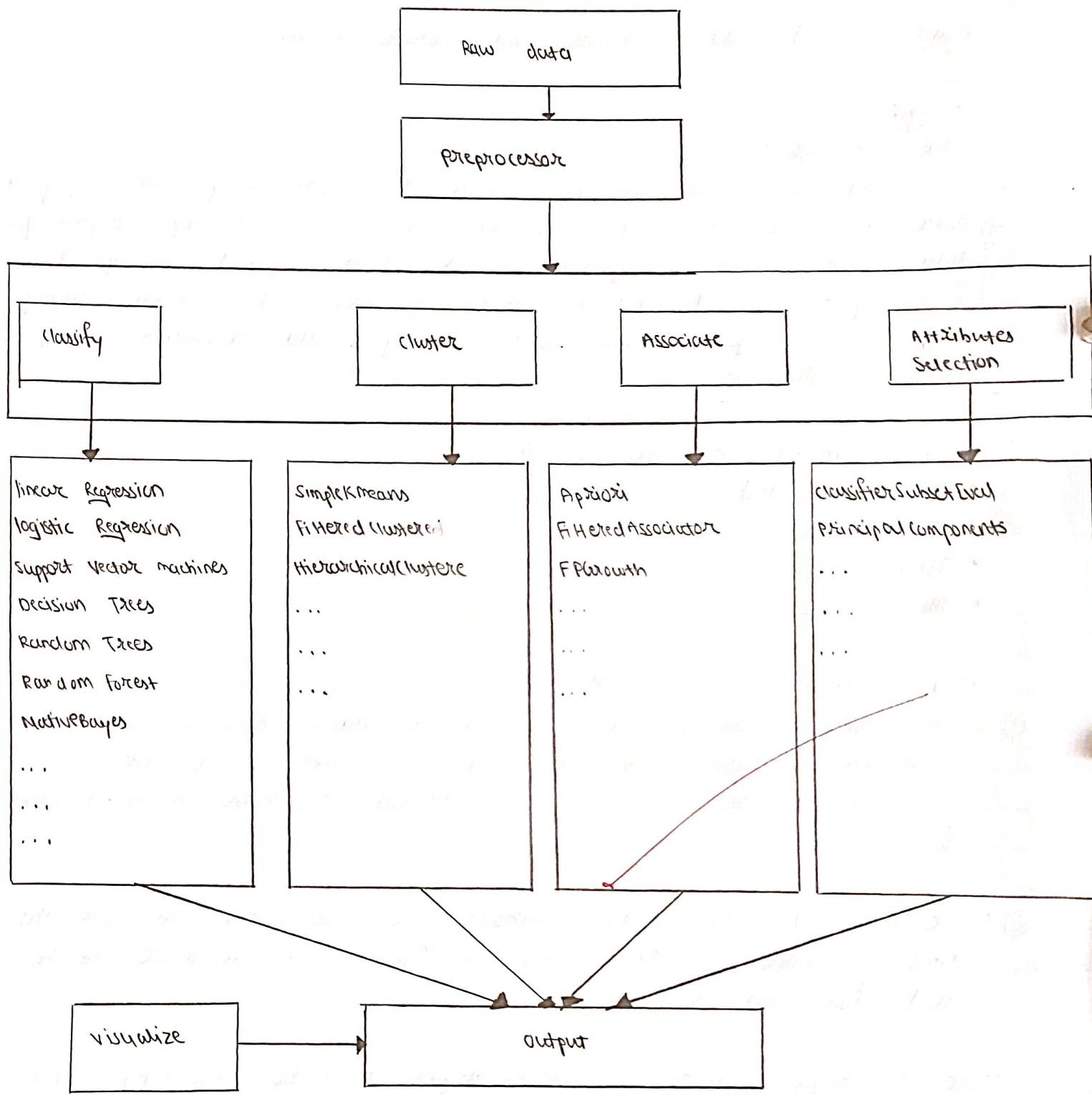
Data warehousing involves data cleaning, data integration, and data consolidations.

Data warehouse characteristics: →

- Subject-oriented
- Integrated
- Time-variant
- Non-Volatile

Need of Data warehouse : →

- ① Business user: Business users require a data warehouse to view summarized data from the past. Since these people are non-technical, the data may be present to them in an elementary form.
- ② Store historical data: Data warehouse is required to store the time variable data from the past. This input is made to be used for various purposes.
- ③ Make strategic decision: Some strategies may be depending upon the data in the data warehouse. So, data warehouse contributes to making strategic decisions.





Date :

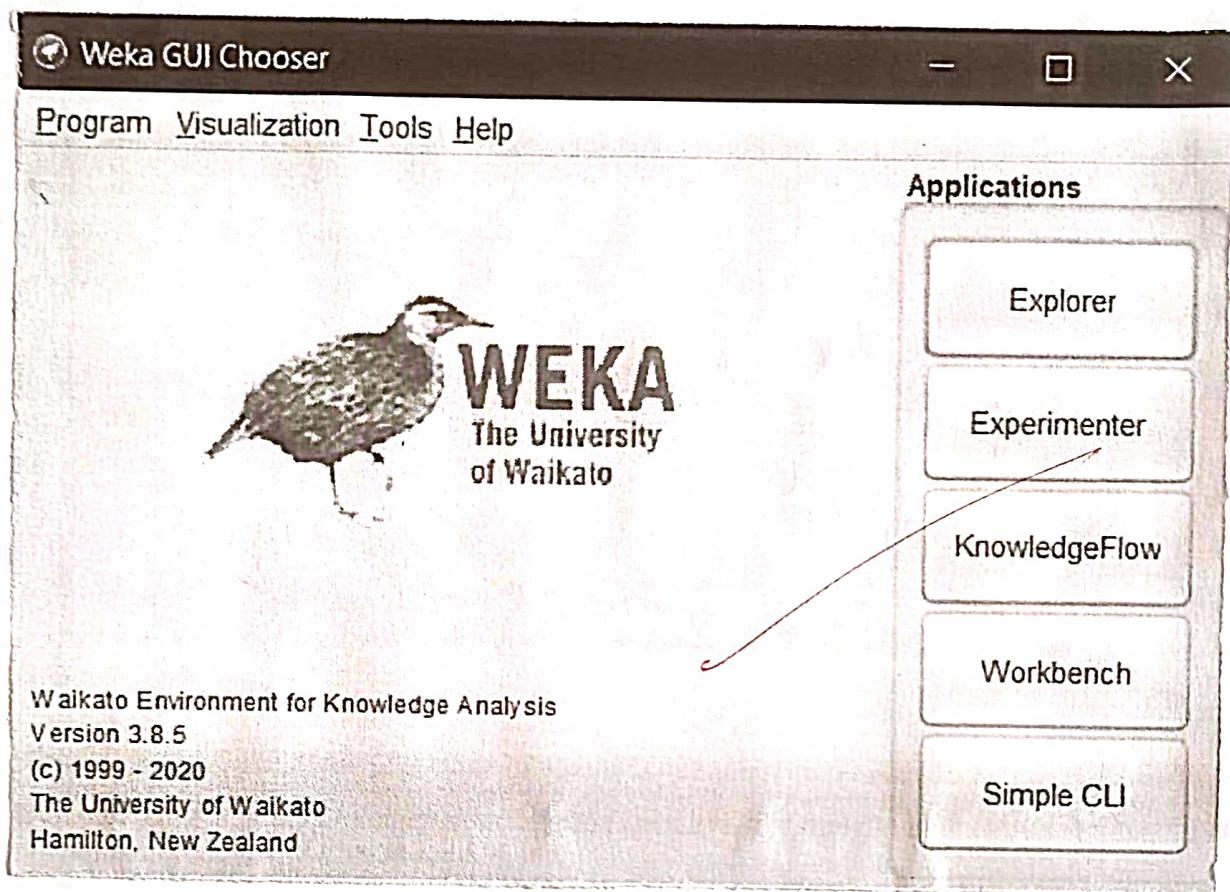
- ④ For data consistency and quality : bringing the data from different sources at a common place, the user can effectively undertake to bring the uniformity and consistency in data.
- ⑤ High response time : data warehouse has to be ready for somewhat unexpected loads and types of queries, which demands a significant degree of flexibility and quick response time.

Benefits of data warehouse : →

- ① Understand business trends and make better forecasting decisions.
- ② Data warehouses are designed to perform well enormous amount of data.
- ③ The structure of data warehouses is more accessible for end-users to navigate, understand and query.
- ④ Queries that would be complex in many normalized databases could be easier to build and maintain in data warehouses.
- ⑤ Data warehousing is an efficient method to manage demand for lots of information from lots of users.
- ⑥ Data warehousing provides the capabilities to analyze a large amount of historical data.

Explaining weka tool : →

WEKA an open source software provides tools for data preprocessing, implementation of several machine learning algorithms, and visualization tools so that you can develop machine learning techniques and apply them to real-world data mining problems. What WEKA offers is summarized in the diagram -





Date :

If you observe the beginning of the flow of the image, you will understand that there are many stages involved with big data to make it suitable for machine learning.

First, you will start with the raw data collected from the field. This data may contain several null values and irrelevant fields. You use the data preprocessing tools provided in WEKA to cleanse the data.

Then, you would save the preprocessed data in your local storage for applying ML algorithms.

Next, depending on the kind of ML model that you are trying to develop you would select one of the option such as classify, cluster or associate. The attributes selection allows the automatic selection of features to creates a reduced dataset.

WEKA provides the implementation of several algorithms. You would select an algorithm of your choice, set the desired parameters and run it on the dataset.

Then, WEKA would give you the statistical output of the model processing. It provides you a visualization tool to inspect the data.

The various models can be applied on the same dataset. You can then compare the outputs of different models and select the best that meets your purpose.

Thus, the use of WEKA results in quicker development of machine learning models on the whole.

Now that we have seen what WEKA is and what it does, in the next chapter let us learn how to install WEKA on your local computer.

WEKA is a data mining software that uses a collection of machine learning algorithms. These algorithm can be applied directly to the data or called from the Java code.



Date :

Weka is a collection of tools for:

- Regression
- Clustering
- Association
- Data pre-processing
- Classification
- Visualisation weka application interface.

There are totally five application interfaces available for weka. When we open weka, it will start the weka GUIchooser screen from where we can open the weka application interfaces Explorer, Preprocessing, Attribute Selection, Learning, Visualization, Experimenter, Testing and Evaluating machine learning algorithms. Knowledge flow visual design of KODI process simple command line A simple interface for typing commands.

Weka data formats:

Weka uses the Attribute Relation File format for data analysis by default. But listed below are some formats that weka supports, from where data can be imported:

- arff
- arff.gz
- bsi
- csv
- dat
- data
- json
- json.gz
- libsvm
- m
- names
- xarff
- xarff.gz



Date :

ARFF Format :-

- An ARFF file contains two sections - headers and data.
- The header describes the attributes types.
- The data section contains a comma separated list of data.
- An ARFF file requires the declarations of the relations.

@relation :-

This is the first line in any ARFF file, written in the header sections, followed by the relation / dataset name. The relation name must be a string and if it contains spaces, then it should be enclosed between quotes.

@attribute :-

These are declared with their names and the type or range in the header sections. Weka supports the following data types for attributes :

- Numeric
- <nominal-specification>
- String
- date
- @data - defined in the data section followed by the list of all duty segments.

Creating a Student Table with the help of data mining Tool weka:

@relation students

@attribute name {Sakil, Swaj, Mayur, Prathamesh}

@attribute rollno numeric

@attribute exp {low, medium, high}

@attribute gender {male, female}

@attribute phone numeric

Training Data

Output:

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit Save... Filter Choose None

Current relation Relation students Instances 4 Attributes

No. 1 name 2 rollno 3 exp 4 gender 5 phone

Attributes: 5 Sum of weights: 4

Selected attribute Name: name Missing: 0 (0%) Distinct: 4 Type: Nominal Unique 4 (100%)

No. Label Count Weight

1 Sahil 1 1
2 Suraj 1 1
3 Mayur 1 1
4 Pratmesh 1 1

Viewer

Relation: students

No.	1: name	2: rollno	3: exp	4: gender	5: phone
1	Sahil	157.0	high	male	80879.0
2	Suraj	162.0	medl.	male	70568.0
3	Mayur	147.0	low	female	75396.0
4	Pratmesh	151.0	high	male	80645.0

Add instance Undo OK Cancel

Status CK

Class: phone (Num)

Visualize All Log

Conclusion:

Thus, the training data table is created and WEKA tool is explored.



Date :

@data

Sachil, 157, high, male, 80879

Swaj, 162, medium, male, 70568

Mayur, 147, low, female, 75896

Prathamesh, 151, high, male, 89645

Conclusion:-

Thus the training data table is created and weka tool is explored.

Viva Questions

① Define Data Warehouse.

→ It is a system used for reporting and data analysis and is considered a key component of business intelligence. A datagwarehouse is a subject oriented, integrated, non-volatile and time variant collection of data in support of managements decision making process.

② what are the characteristics of Data warehousing?

→ There are four characteristic of data warehousing and they are:

a) Subject oriented

c) Time - variant

b) Integrated

d) Non - volatile.

③ what is weka tool and what are the significance of weka?

→ (1) weka is an open source software software.

(2) It is collection of machine learning algorithms for data mining tasks.

(3) weka contains tools for data pre-processing, classification, regression, clustering, association rules and visualization.

* Significance of weka!

① Free availability under GNU General Public License Page No. _____



Date :

- ② Portability, since it is fully implemented in Java programming language
- ③ A comprehensive collection of data preprocessing and modelling techniques.
- ④ Ease of use due to its graphical user interfaces.
- ⑤ It provides you a visualization tool to inspect the data.
- ⑥ The use of WEKA results in quicker development of ML models on the whole.

21/02/2022
R

Practical Mu. 02

Aim: To demonstrate pre-processing on provided dataset.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose None

Current relation Relation: weather.symbolic Instances: 14 Attributes: 5 Sum of weights: 14

Attributes

No.	Name
1	outlook
2	temperature
3	humidity
4	windy
5	play

Selected attribute Name: outlook Missing: 0 (0%) Distinct: 3 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	sunny	5	5
2	overcast	4	4
3	rainy	5	5

Class play (Nom) Visualize All

(Loading data from file)

Open

Look In: PRO2

dataset

File Name: weather.arff

Files of Type: Arff data files (*.arff)

Open Cancel

Status: OK

The screenshot shows the Weka Explorer interface. The 'Selected attribute' table displays the distribution of the 'outlook' attribute with three distinct values (sunny, overcast, rainy) each having a count of 5 and a weight of 5. Below this is a bar chart titled 'Class play (Nom)' showing the same distribution. A red arrow points from the 'play' column in the 'Selected attribute' table to the bars in the chart. On the left, there is an 'Open' dialog box showing the path 'PRO2\dataset' and the file name 'weather.arff'. A red arrow also points from the 'File Name' field in the dialog to the corresponding field in the Weka interface. The status bar at the bottom indicates 'OK'.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter Choose None

Current relation Relation: weather.symbolic Instances: 14 Attributes: 5 Sum of weights: 14

Attributes

No.	Name
1	outlook
2	temperature
3	humidity
4	windy
5	play

Selected attribute Name: outlook Missing: 0 (0%) Distinct: 3 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	sunny	5	5
2	overcast	4	4
3	rainy	5	5

Class: play (Nom) Visualize All

Load Instances

Enter the source URL: <https://stream.cs.stanford.edu/~gweiss>

OK Cancel

Status: OK

(Loading data from URL)

The screenshot shows the Weka Explorer interface again. The 'Selected attribute' table for 'outlook' is identical to the previous one. A red arrow points from the 'play' column in the 'Selected attribute' table to the bars in the chart. On the left, there is a 'Load Instances' dialog box with a URL field containing 'https://stream.cs.stanford.edu/~gweiss'. A red arrow points from the 'URL' field in the dialog to the corresponding field in the Weka interface. The status bar at the bottom indicates 'OK'.



Aim: To demonstrate pre-processing on provided dataset.

Theory:

The data that we collected from the field contains many unwanted things that leads to wrong analysis. For example, the data may contain null fields, it may contain columns that are irrelevant to the current analysis, and so on. Thus, the data must be preprocessed to meet the requirements of the type of analysis you are seeking. This is the done in the preprocessing module.

To demonstrate the available features in preprocessing, we can use the datasets that is provided in the installation.

Ways to load Data : →

The data can be loaded from the following sources -

- Local file system
- web
- Database.

Load data from local File System : →

Just under the machine learning tabs that you studied in the previous lesson, you would find the following three buttons -

open file ...

open URL...

Open DB...

Click on the open file ... button, a directory navigator window opens as shown in following screen -

Load data from web : →

Once you click on the open URL... button, you get a window as follows:

File | Preferences | Classify | Cluster | Associate | Select attributes | Visualize

Open file... Open URL... Open DB... Generate... Undo Redo Save...

Filter

Choose Remove

Current selection: Relation: weather symbolic-weights filtering generalized attributeAdd TSTR-Noposition-L2 WID weights filtering priority. Instances: 14 Attributes: 5 Sum of weights: 14

Attributes

AB None Invert Patterns

No. Name

- 1 outlook
- 2 temperature
- 3 humidity
- 4 windy
- 5 play

Selected attribute: Name: outdoor Missing: 0 (0%) Distance: 3 Type: Nominal Unique: 3 (3%)

No.	Label	Count	Weight
1	sunny	5	5
2	overcast	4	4
3	rainy	5	5

Class play (Nom)

Visualize All

Log X

File | Preferences | Classify | Cluster | Associate | Select attributes | Visualize

Open file... Open URL... Open DB... Generate... Undo Redo Save...

Filter

Choose Add TSTR-Noposition-L2 WID 0

Current selection: Relation: weather symbolic-weights filtering generalized attributeAdd TSTR-Noposition-L2 WID 0 Instances: 14 Attributes: 6 Sum of weights: 14

Attributes

AB None Invert Patterns

No. Name

- 1 outlook
- 2 temperature
- 3 humidity
- 4 windy
- 5 play
- 6 ?

Selected attribute: Name: outdoor Missing: 0 (0%) Distance: 3 Type: Nominal Unique: 3 (3%)

No.	Label	Count	Weight
1	sunny	5	5
2	overcast	4	4
3	rainy	5	5

View

Relation: weather symbolic-weights filtering generalized attributeAdd TSTR-Noposition-L2 WID 0

	1 outlook	2 temperature	3 humidity	4 windy	5 play
1	outlook	normal	normal	normal	normal
2	sunny	hot	high	FALSE	no
3	sunny	hot	high	TRUE	no
4	overcast	normal	high	FALSE	yes
5	overcast	normal	high	FALSE	yes
6	overcast	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	cool	normal	FALSE	no
9	sunny	cool	high	FALSE	yes
10	sunny	mild	normal	FALSE	yes
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	normal	normal	FALSE	yes
14	sunny	mild	high	TRUE	no

Add instance | Delete | OK | Cancel

Class play (Nom)

Visualize All

Log X

(Adding individual)



We will open the file from a public URL type the following URL in the popup box -

for example : <https://storm.cis.fordham.edu/~guess1/data-mining/weka-data/weather-nominal.arff>

You may specify any other URL where your data is stored. The Explorer will load the data from the remote site into its environment.

Understanding Data →

Let us first look at the highlighted current relation sub window. It shows the name of database that is currently loaded. You can infer two points from this sub window. There are 14 instances - the number of rows in the table. The table contains 5 attributes - the field, which are discussed in the upcoming sections.

On the left side, notice the attributes sub window that displays the staticus fields in the database.

Applying Filters →

Some of the machine learning techniques such as associations rule mining required categorical data. To illustrate the use of filters, we will use weather-numeric.arff database that contain two numeric attributes - temperature and humidity. We will convert these to nominal by applying a filters on our raw data, click on the choose button in the filter subwindow and select the following filters - weka → filters → supervised → attribute → discretize.

Click on the Apply button and examine the temperature and for humidity attribute. You will notice that these have changed from numeric to nominal types. Page No.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo EDA... Save...

Filter Choose Remove

Current relation: Relation weather symbolic-weka filters.unsupervised.attribute.AddTSTR-Normalization-C2-WIDD-weka filters.unsupervised... Attributes: 5 Instances: 14 Sum of weights: 14

Attributes

All	None	Invert	Pattern
No	Name		
1 outlook			
2 temperature			
3 humidity			
4 windy			
5 play			

Selected attribute:
Name: outlook
Missing: 0 (0%)
Distinct: 3
Type: Nominal
Unique: 0 (0%)

No.	Label	Count	Weight
1 sunny	5	5	
2 overcast	4	4	
3 rainy	5	5	

Class: play (Nom) Visualize All

(Removing attribute)

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo EDA... Save...

Filter Choose Normalize-S10-T00

Current relation: Relation weather symbolic-weka filters.unsupervised.attribute.AddTSTR-Normalization-C2-WIDD-weka filters.unsupervised... Attributes: 5 Instances: 14 Sum of weights: 14

Attributes

All	None	Invert	Pattern
No	Name		
1 outlook			
2 temperature			
3 humidity			
4 windy			
5 play			

Selected attribute:
Name: outlook
Missing: 0 (0%)
Distinct: 3
Type: Nominal
Unique: 0 (0%)

No.	Label	Count	Weight
1 sunny	5	5	
2 overcast	4	4	
3 rainy	5	5	

Class: play (Nom) Visualize All

(Normalization)



- Like this we can apply following prepossessing on the data set.
- (1) Add : To add attribute. It is present in unsupervised data. In that we can enter attribute index, type, data format, nominal label values.
 - (2) Remove : To remove attributes, select them and click on the remove button at the bottom. The selected attributes would be removed from db.
 - (3) Normalization : It is in unsupervised data. It is good technique to use when you do not know the distribution of your data. You can normalized all of the your dataset.

Result :

The pre-processing on the given data set is executed.

Viva Questions :-

- (1) Define prepossessing Technique ?
→ Preprocessing tools in WEKA are called filters. The preprocess receives data from a file, SQL db or URL. Some preprocessing are Add, remove, normalization etc..
- (2) What are data preprocessing steps?
→ Step 1 : Open weka
Step 2 : load dataset from URL, SQL database or file.
Step 3 : click on filters
Step 4 : click on supervised
Step 5 : click on attribute
Step 6 : Attribute selection or select the prepossessing technique ie. add, remove or etc..
Step 7 : click on apply.



Date :

Q) What are the 3 stages of data processing?

→ Stage 1: Loading data or dataset

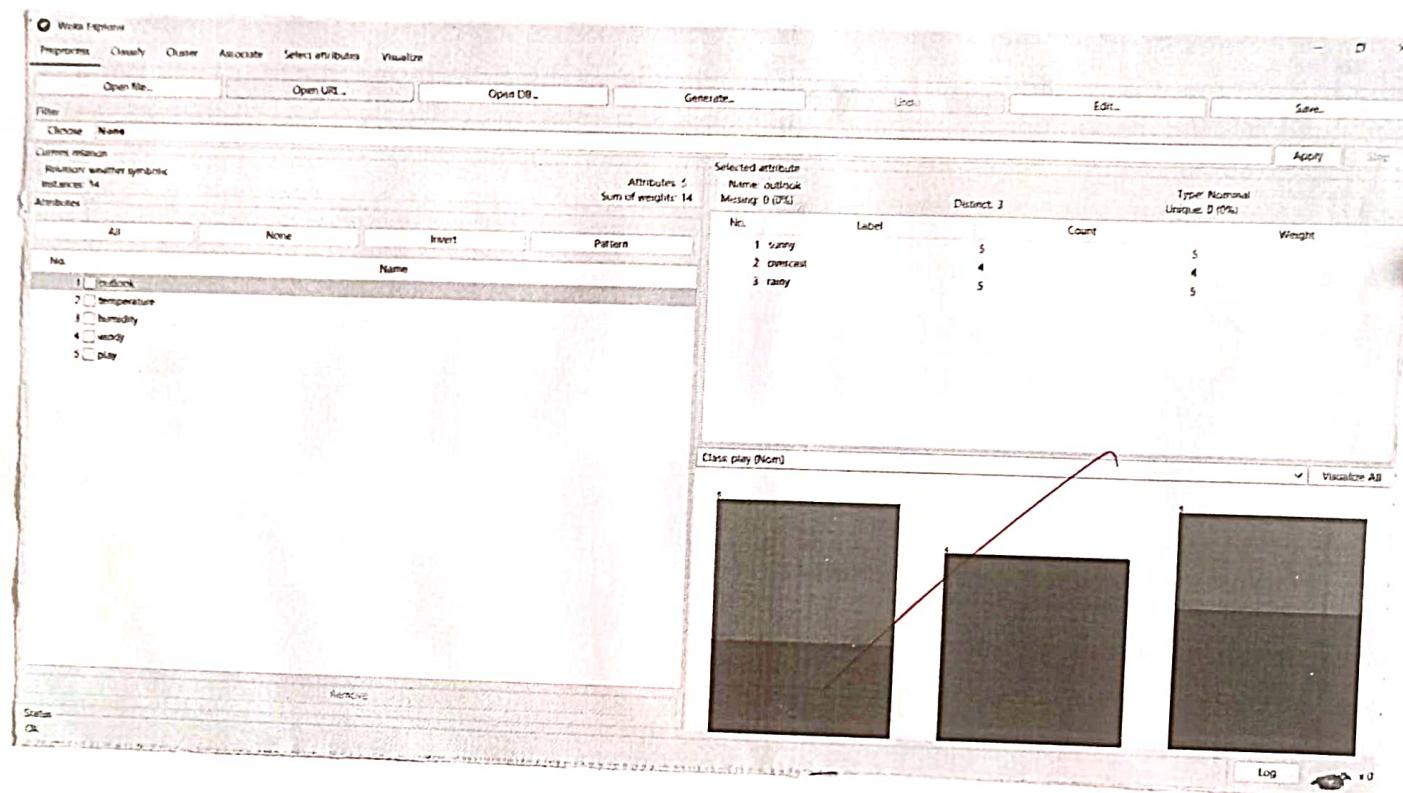
Stage 2: Understanding Data

Stage 3: Apply filters on dataset and save.

Ques
Ans

Practical No. 03

Aim: To demonstrate performing classification on data set.





Practical No. 63

Aim: To demonstrate performing classification on data sets.

Theory :

Classification:

Classification is the process for finding a model that describes the data values and concept for the purpose of prediction.

Classification in data mining is a common technique that separates data points into different classes.

It allows you to organize data set of all sorts, including complex and large datasets as well as small and simple ones. It primarily involves using algorithms that you can easily modify to improve the data quality. The algorithm establishes the link between the variables for predictions.

The algorithm you use for classification in data mining is called the classifier, and observations you make through the same are called instances.

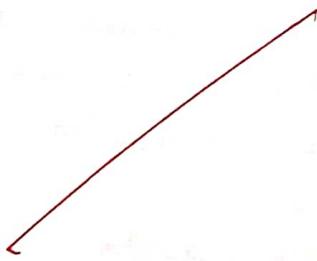
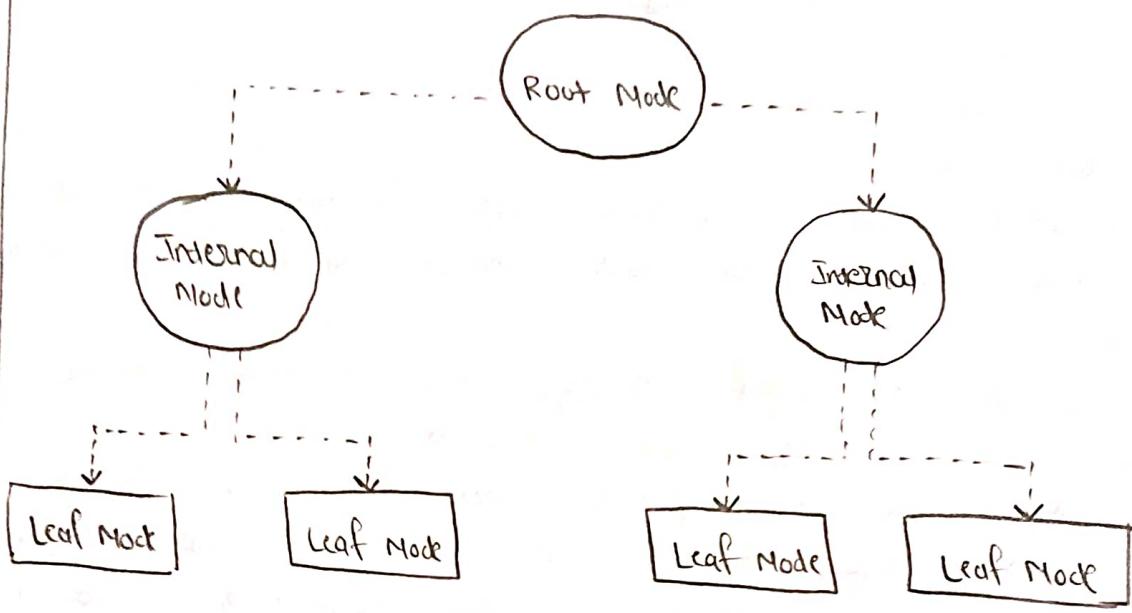
~~There are multiple type of classification algorithm, each with its unique functionality and application. All of those algorithm are used to extract data from a dataset.~~

Data Mining Algorithm for Classification:

- Decision trees
- Logistic Regression
- Naive Bayes Classification
- k-nearest neighbors
- Support Vector machine.

Decision Tree:

A decision tree is a structure that includes a root node, branches and leaf nodes. Each internal node performs a test on an





Date :

attribute, each branch denotes the outcome of a test and each leaf node holds a class label. The topmost node in the tree is the root node.

A decision tree is a classification scheme to generate a tree consisting of root node, internal nodes and external nodes. Root nodes representing the attributes. Internal nodes are also the attributes. External nodes are the class and each branch represent the values of the attributes.

Decision tree also contains set of rules for a given data set; there are two subsets in decision tree, one is a training data set and second one is a testing data set. Training data set is previously classified data. Testing data set is newly generated data.

The benefits of having a decision tree are as follows:-

- It does not require any domain knowledge.
- It is easy to comprehend.
- The learning and classification steps of a decision tree are simple and fast.
- J48 (classification and its decision tree)
- C4.5 algorithm | J48
- The C4.5 algo. is a classification algo. which provides decision trees based on information theory. It is an extension of Ross Quinlan's earlier ID3 algorithm also known in weka as J48, J standing for Java. The decision trees generated by C4.5 are used for classification, and for this reason, C4.5 is often referred to as a statistical classifier.
- The J48 implementation of the C4.5 algo. has many additional features including accounting for missing values, decision trees pruning, continuous attribute value range, derivation of rule and etc..
- In WEKA, data mining tool, J48 is an open page wise Java

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier
Choose: J48 (0.25 M 2)

Test options
 Use training set
 Selected test set
 Cross-validation Folds: 10
 Percentage split %: 10
 More options...

Result plot
 Start Stop
 Result list (right-click for options): [trees/J48](#)

Classifier output
 Weka Version: 3.7.1 (WEKA 3.7.1)
 package: weka
 1. Windy = FALSE: no (1.0)
 1. Windy = TRUE: yes (1.0)

Number of leaves: 5

Size of the tree: 3

Time taken to build model: 0 seconds

Weka Prioritized Cross Validation

— Summary —

Correctly Classified Instances	50	%
Incorrectly Classified Instances	2	%
Kappa statistic	-0.0424	
Mean absolute error	0.4147	
Root mean squared error	0.5984	
Relative absolute error	87.5	%
Root relative squared error	121.2987	%
Total Number of Instances	52	

— Detailed Accuracy By Class —

Class	TP	FP	FN	Precision	Recall	F-Measure	NCC	RND Area	PRC Area	Class
0.0	0.000	0.400	0.625	0.500	0.500	0.500	-0.043	0.633	0.718	yes
0.1	0.000	0.400	0.625	0.400	0.333	0.333	-0.043	0.433	0.537	no
Weighted Avg.	0.000	0.400	0.625	0.500	0.500	0.500	-0.043	0.633	0.650	

— Confusion Matrix —

	0.0	1.0
0.0	48	2
1.0	2	50

0.0 = **classified as**
 1.0 = **actual**
 2.0 = **error**

Status: OK **Log:** **XO:**

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier
Choose: J48 (0.25 M 2)

Test options
 Use training set
 Selected test set
 Cross-validation Folds: 10
 Percentage split %: 10
 More options...

Result plot
 Start Stop
 Result list (right-click for options): [trees/J48](#)

Classifier output
 Weka Version: 3.7.1 (WEKA 3.7.1)
 package: weka
 1. Windy = FALSE: no (1.0)
 1. Windy = TRUE: yes (1.0)

Number of leaves: 5

Size of the tree: 3

Time taken to build model: 0 seconds

Weka Prioritized Cross Validation

— Summary —

Correctly Classified Instances	50	%
Incorrectly Classified Instances	2	%
Kappa statistic	-0.0424	
Mean absolute error	0.4147	
Root mean squared error	0.5984	
Relative absolute error	87.5	%
Root relative squared error	121.2987	%
Total Number of Instances	52	

— Detailed Accuracy By Class —

Class	TP	FP	FN	Precision	Recall	F-Measure	NCC	RND Area	PRC Area	Class
0.0	0.000	0.400	0.625	0.500	0.500	0.500	-0.043	0.633	0.718	yes
0.1	0.000	0.400	0.625	0.400	0.333	0.333	-0.043	0.433	0.537	no
Weighted Avg.	0.000	0.400	0.625	0.500	0.500	0.500	-0.043	0.633	0.650	

— Confusion Matrix —

	0.0	1.0
0.0	48	2
1.0	2	50

0.0 = **classified as**
 1.0 = **actual**
 2.0 = **error**

Weka Classifier Tree Visualizer (19.12.03 - trees/J48 weather symbols)

Tree View

```

graph TD
    Root((Outlook)) -- "= sunny" --> Node1((Humidity))
    Root -- "= overcast" --> Node2((yes[yes (4.0)]))
    Root -- "= rainy" --> Node3((Wind))
    Node1 -- "= high" --> Node4((no[no (2.0)]))
    Node1 -- "= normal" --> Node5((yes[yes (2.0)]))
    Node3 -- "= TRUE" --> Node6((no[no (2.0)]))
    Node3 -- "= FALSE" --> Node7((yes[yes (2.0)]))
  
```

Status: OK **Log:** **XO:**



Date :

Implementation of the C4.5 algorithm. It allows classification via either decision trees or rules generated from them.

Generating a decision tree from training tuples of data partition D

Algorithm : Generate decision tree

Input :

Data partition D , which is a set of training tuples and their associated class labels.

attribute-list . the set of candidate attributes.

Attribute selection method, a procedure to determine the splitting criterion that best partitions the data tuples into individual classes . This criterion includes a splitting attribute and either a splitting point or splitting subset .

Output :

A Decision Tree

Method :

Create a node N ;

if tuples in D are all of the same class, C then
return N as leaf node labeled with class C ;

if attribute-list is empty then

return N as leaf node with labeled

with majority class in D ; // majority Voting

apply attribute-selection method (D , attribute-list)

to find the best splitting-criterion;

label node N with splitting-criterion;

if splitting-attribute is discrete-valued and

multiway split allowed then // no restriction to binary trees

Result:

Thus, the classification on Data set is performed by decision tree (J48) method.



Date :

attribute-list = splitting attribute // remove splitting attribute
for each outcome j of splitting criterion

// partition the tuples and grow subtrees for each partitions
let D_j be the set of data tuples in a satisfying outcome j ;
// a partition

If D_j is empty then
attach a leaf labeled with the majority class in D to node N_j ;

else

attach the node returned by generate decision tree (D_j , attribute list) to node N_j ;

end for

return N_j

Result:

Thus the classification on data set is performed by decision tree (ID3) method.

Viva Question :

① What is classification?

→ Classification is the process for finding a model that describes the data values and concept for the purpose of prediction.

② What is the need of classification?

→ The need of classification is to accurately predict the target class for each case in the data, it allows you to organize data set of all sorts, including complex and large datasets as well as small and simple ones.



Date :

Q.3 What are the different methods of classification?

→ Logistic regression, Naive Bayes, Decision tree, k-Nearest Neighbors, Support vector machine, Bayes classifiers, function classifier, lazy classifier, meta classifier and so on.

Q.4 What are the advantages of a decision tree classifier?

- ① It is easy to comprehend
② It does not require any domain knowledge.
③ The learning and classification steps of decision tree are simple and fast
④ Less data preparation, Non-parametric, versatility, Non-linearity.

11
08/03/23
X

Practical No. 4

Nim: To demonstrate classification rule generation dataset using naive bayes algorithm.

Dataset: Iris dataset. It is a multiclass dataset in which we have to predict the species of flower based on four features: sepal length, sepal width, petal length, and petal width. The dataset contains 150 samples, each with five attributes: Sepal Length, Sepal Width, Petal Length, Petal Width, and Species. The Species attribute has three possible values: Setosa, Versicolor, and Virginica.

Naive Bayes classifier: A simple probabilistic classifier based on applying Bayes' theorem with strong (naive) independence assumptions between the features. It is a generative model that assumes all variables are continuous.

Implementation: We will implement the Naive Bayes classifier from scratch using Python. The implementation will include data loading, feature extraction, model training, and prediction.

Data Preprocessing: The first step is to load the dataset and preprocess it. This includes handling missing values, normalizing the data, and splitting it into training and testing sets.

Model Training: We will train the Naive Bayes classifier on the training dataset. This involves estimating the parameters of the model, such as the prior probabilities of each class and the conditional probabilities of each feature given a class.

Prediction: Once the model is trained, we can use it to make predictions on the testing dataset. The predicted class for each sample is determined by the class with the highest posterior probability.

Evaluation: Finally, we evaluate the performance of the classifier using various metrics such as accuracy, precision, recall, and F1 score. These metrics help us understand how well the classifier is performing and identify areas for improvement.



Aim: To demonstrate classification rule process on dataset using naive Bayes algorithm.

Theory:

Naive Bayes classifier algorithm:

- Naive Bayes algo. is a supervised learning algo., which is based on Bayes theorem and used for solving classification problems.
- It is mainly used in text classification that includes a high-dimensional training dataset.
- Naive Bayes classifier is one of the simple and most effective classification algs.
- It is a probabilistic classifier, which means it predict on the basis of the probability of an object.
- Some popular example of Naive Bayes algo. are spam filtration, sentimental analysis and classifying articles.

why is it called Naives Bayes?

The Naive Bayes algo. is comprised of two words Naive and Bayes, which can be described as:

~~Naive~~: It is called Naive because it assume that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape and taste then red, spherical, and sweet fruit is identifier or recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

~~Bayes~~: It called Bayes because it depends on the principle of Bayes theorem.

output:

(1)

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose NaiveBayes

Test options Use training set Supplied test set Cross-validation Folds 10 Percentage split % More options...

(Nom) play Start Stop Result List (right-click for options) 16.03.37-bayes.NaiveBayes

Classifier output

--- Run information ---

Scheme: weka.classifiers.bayes.NaiveBayes
Relation: weather.symbolic
Instances: 14
Attributes: 4
outlook
temperature
humidity
windy
play

Test mode: 10-fold cross-validation

--- Classifier model (full training set) ---

Naive Bayes Classifier

Attribute	Class	yes	no
(Nom) play		(0.61)	(0.38)
outlook			
sunny	yes	2.0	4.0
sunny	no	12.0	8.0
overcast	yes	5.0	1.0
rainy	yes	4.0	3.0
rainy	no	11.0	7.0
[total]		12.0	8.0
temperature			
hot	yes	3.0	3.0
hot	no	12.0	8.0
mild	yes	5.0	3.0
cool	yes	4.0	2.0
cool	no	11.0	7.0
[total]		12.0	8.0
humidity			
high	yes	4.0	5.0
high	no	11.0	7.0
normal	yes	7.0	2.0
normal	no	11.0	7.0
[total]		11.0	7.0
windy			
TRUE	yes	4.0	4.0
TRUE	no	7.0	3.0
[total]		11.0	7.0

Time taken to build model: 0 seconds

--- Stratified cross-validation ---

--- Summary ---

	Correctly Classified Instances	Incorrectly Classified Instances	%
outlook	8	6	57.1429 %
temperature	8	6	42.8571 %
humidity	8	6	42.8571 %
windy	8	6	42.8571 %
Total Number of Instances	14	14	100.0000 %

Correctly Classified Instances: 8
Incorrectly Classified Instances: 6
Kappa statistic: -0.0244
Mean absolute error: 0.4374
Root mean squared error: 0.4916
Relative absolute error: 91.8631 %
Root relative squared error: 99.6452 %
Total Number of Instances: 14

Status OK

(2)

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose NaiveBayes

Test options Use training set Supplied test set Cross-validation Folds 10 Percentage split % More options...

(Nom) play Start Stop Result List (right-click for options) 16.03.37-bayes.NaiveBayes

Classifier output

--- Run information ---

Scheme: weka.classifiers.bayes.NaiveBayes
Relation: weather.symbolic
Instances: 14
Attributes: 4
outlook
temperature
humidity
windy
play

Test mode: 10-fold cross-validation

--- Classifier model (full training set) ---

Naive Bayes Classifier

Attribute	Class	yes	no
(Nom) play		(0.61)	(0.38)
outlook			
sunny	yes	2.0	4.0
sunny	no	12.0	8.0
overcast	yes	5.0	1.0
rainy	yes	4.0	3.0
rainy	no	11.0	7.0
[total]		12.0	8.0
temperature			
hot	yes	3.0	3.0
hot	no	12.0	8.0
mild	yes	5.0	3.0
cool	yes	4.0	2.0
cool	no	11.0	7.0
[total]		12.0	8.0
humidity			
high	yes	4.0	5.0
high	no	11.0	7.0
normal	yes	7.0	2.0
normal	no	11.0	7.0
[total]		11.0	7.0
windy			
TRUE	yes	4.0	4.0
TRUE	no	7.0	3.0
[total]		11.0	7.0

Time taken to build model: 0 seconds

--- Stratified cross-validation ---

--- Summary ---

	Correctly Classified Instances	Incorrectly Classified Instances	%
outlook	8	6	57.1429 %
temperature	8	6	42.8571 %
humidity	8	6	42.8571 %
windy	8	6	42.8571 %
Total Number of Instances	14	14	100.0000 %

Correctly Classified Instances: 8
Incorrectly Classified Instances: 6
Kappa statistic: -0.0244
Mean absolute error: 0.4374
Root mean squared error: 0.4916
Relative absolute error: 91.8631 %
Root relative squared error: 99.6452 %
Total Number of Instances: 14

--- Detailed Accuracy By Class ---

	TP	P-F1	FN	FP	Precision	Recall	F-Measure	Roc Area	Class
outlook	0.75	0.75	0.25	0.25	0.75	0.75	0.75	0.579	yes
temperature	0.75	0.75	0.25	0.25	0.75	0.75	0.75	0.578	no
humidity	0.2	0.22	0.78	0.2	0.25	0.25	0.25	0.578	no
windy	0.2	0.22	0.78	0.2	0.25	0.25	0.25	0.578	no
Total	0.571	0.594	0.528	0.571	0.539	0.539	0.539	0.576	
Weighted Avg.	0.571	0.594	0.528	0.571	0.539	0.539	0.539	0.576	

--- Confusion Matrix ---

	a b	a c	b c	c a	c b
a	8	2	0	0	0
b	2	0	0	0	0
c	0	0	0	0	0
a b	8	2	0	0	0
a c	2	0	0	0	0
b c	0	0	0	0	0
c a	0	0	0	0	0
c b	0	0	0	0	0

a | b --> classified as
2 | 2 | a = yes
4 | 1 | b = no

Status OK



Date :

Bayes' theorem:

- Bayes' theorem is also known as Bayes' rule or Bayes' law, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.
- The formula for Bayes' theorem is given as:

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

where,

$P(A|B)$ is posterior probability : probability of hypothesis A on the observed event B.

$P(B|A)$ is likelihood probability : probability of the evidence given that the probability of a hypothesis true.

$P(A)$ is prior probability: probability of hypothesis before observing the evidence.

$P(B)$ is marginal probability : probability of evidence.

Working of Naive Bayes' classifier:

Working of Naive Bayes' classifier can be understood with the help of the below example:

Suppose we have dataset of weather conditions and corresponding target variable "play". So using this dataset we need to decide that whether we should play or not on a particular day according to the weather condition. So to solve this problem, we need to follow the below steps:

- Convert the given dataset into frequency tables.
- Generate likelihood table by finding the probabilities of given frequency.
- Now, use Bayes' theorem to calculate the posterior probability.

Step 1: Initially, we have to load the required dataset in the weka tool using choose file option. Here we selecting the weather - nominal dataset to execute.

Result:

Thus, the classification on dataset is performed by Naïves Bayes classification.



Date :

Step 2: Now we have to go to the classify tab on the left side and click on the choose button and select the naive Bayesian algorithm in it.

Step 3: Now to change the parameters click on the right side at the choose button and we are accepting the default values in the example.

Step 4: We choose the percentage split as our measurement method from the "Test" choices in the main panel. since we don't have a separate test data collection, we'll use the percentage split of 66 percent to get a good idea of the model's accuracy. Our dataset contains 14 examples, with 9 being used for training and 5 being used for testing.

Step 5: To generate the model, we now click 'start'. When the model is done, the evaluation statistic will appear in right panel.

Result:

Thus the classification on data set is performed by Naïves Bayes classification.

Viva Questions:

① What is a Naïve Bayes classifier?

→ A naïve Bayes classifier is an algo. that uses Bayes' theorem to classify objects. Naïve Bayes classifiers assume strong or naïve, independence between attributes of data points.

② What are the basic assumption?

→ The basic assumption in naïve Bayes is one of conditional independence between all independent variable features.



Date :

Ques.3) what are the advantages of Naive Bayes classifier?

- ① It is simple and easy to implement.
- ② It doesn't require as much training data.
- ③ It handle both discrete and continuous data.
- ④ It is fast and can be used to make real time predictions.
- ⑤ It is not sensitive to irrelevant features.

Ques.4) Name the different problem statement you can solve using Naive Baye's.

- ① A fruit may be considered to be an apple if it is red, round and about 3 inch in diameter.
- ② If the weather is sunny, then the player should play or not?

Ques.3)
Ans
A

Date :

Practical No. 5



Aim: To implement k-means algorithm.

Description:

clustering: →

clustering is the method of dividing a set of abstract objects into groups. points to keep in mind A set of data object can be viewed as a single entity when performing cluster analysis. we divide the data set into groups based on data similarity, then assign labels to the groups.

Simple k-means clustering: →

K means clustering is a simple unsupervised learning algorithm. In this, the data objects ('n') are grouped into a total 'k' clusters, with each observation belonging to the cluster with closest mean. It defines 'k' sets, one for each clusters $k \in n$. The clusters are separated by a large distance.

The data is then organized into acceptable data sets and linked to the nearest collection. If no data is pending, the first stage is more difficult to complete; in the case, an early grouping is performed. The 'k' new set must be recalculated as the barycenters of the clusters from the previous stage.

The same data set points and the nearest new sets are bound together after these 'k' new sets have been created. After that, a loop is created. The 'k' sets change their position step by step until no further changes are made as a result of this loop.

K-means clustering algo. computes the centroids and iterates until we it finds optimal centroid. It assumes that the no. of clusters are already known. It is also called flat clustering algo. The no. of clusters identified from data by algo. is represented by 'k' in k-means.

In this algo. the data points are assigned to a cluster in such a manner that the sum of the squared distances between the

Activities weka-gui-FileChooser Mar 23 4:43 PM

Preprocess Classify Cluster Associate Select attributes Visualize

Clusterer
 Choose SimpleKMeans -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10
Cluster mode
 Use training set
 Supplied test set Set
 Percentage split %
 Classes to clusters evaluation
 Store clusters for visualization
Ignore attributes
 Start Step
Result list (right-click for options) 16:43:42 - SimpleKMeans

Weka Explorer

Clusterer output

```
==== Run information ====
Schema: weka.clusterers.SimpleKMeans -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -S 10
Relation: labor-neg-data
Instances: 571
Attributes: 17
duration
wage-increase-first-year
wage-increase-second-year
wage-increase-third-year
cost-of-living-adjustment
working-hours
pension
standby-pay
shift-differential
education-allowance
statutory-holidays
vacation
longterm-disability-assistance
contribution-to-dental-plan
bereavement-assistance
contribution-to-health-plan
class

Test mode: evaluate on training data
==== Model and evaluation on training set ====
kMeans
=====
Number of iterations: 3
Within cluster sum of squared errors: 137.79496140156423
Missing values globally replaced with mean/mode
Cluster centroids:
```

Attribute	Full Data (571)	Clusters	
		0 (48)	1 (9)
duration	2.1607	2.2533	1.6657
wage-increase-first-year	3.8036	3.9834	2.8444
wage-increase-second-year	3.9717	4.0209	3.7097
wage-increase-third-year	3.9133	3.9511	3.7119
cost-of-living-adjustment	none	none	none
working-hours	39.0392	37.4541	36.5599
pension	empl.contr	empl.contr	none
standby-pay	7.4444	7.7431	5.0519
shift-differential	4.871	5.2298	2.957
education-allowance	no	no	no
statutory-holidays	11.0803	11.237	10.3333
vacation	below_average	below_average	below_average
longterm-disability-assistance	yes	yes	no
contribution-to-dental-plan	half	half	none
bereavement-assistance	yes	yes	yes
contribution-to-health-plan	full	full	none
class	good	good	bad

Status OK

Activities weka-gui-Clustering Mar 23 4:44 PM

Preprocess Classify Cluster Associate Select attributes Visualize

Clusterer SimpleKMeans - N=3 - A:meta.core.EuclideanDistance - R: first last 4 500 6 10

Cluster mode

- Use training set
- Supplied test set 481
- Percentage split %
- Classes to clusters evaluation
- Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)
16:43:42 - SimpleKMeans
16:44:29 - SimpleKMeans

Cluster output

```

@relation "HR Attrs"
@attribute duration numeric
@attribute wage-increase-first-year numeric
@attribute wage-increase-second-year numeric
@attribute wage-increase-third-year numeric
@attribute cost-of-living-adjustment numeric
@attribute working-hours numeric
@attribute pension numeric
@attribute standby-pay numeric
@attribute shift-differential numeric
@attribute education-allowance numeric
@attribute statutory-holidays numeric
@attribute vacation numeric
@attribute below-average numeric
@attribute long-term-disability-assistance yes/no
@attribute contribution-to-dental-plan half/full
@attribute bereavement-assistance yes/no
@attribute contribution-to-health-plan full
@attribute class good/bad

```

Test mode: evaluate on training data

--- Model and evaluation on training set ---

KMeans

Number of iterations: 3
 Within cluster sum of squared errors: 119.5224194214812
 Missing values globally replaced with mean/node

Cluster centroids:

Attribute	Full Data (57)	Clusters		
		0 (36)	1 (5)	2 (16)
duration	2.1137	2.2267	1.4	2.25
wage-increase-first-year	3.8136	4.4695	3.2	2.4938
wage-increase-second-year	3.9717	4.4175	4.183	2.9027
wage-increase-third-year	3.6133	4.1093	3.9133	3.4725
cost-of-living-adjustment	none	none	none	none
working-hours	38.0392	37.4756	39.2078	38.94
pension	7.4444	empl_centr	none	empl_centr
standby-pay	7.4444	7.0938	6.7556	6.4226
shift-differential	4.871	5.4776	3.1494	4.0444
education-allowance	no	no	no	no
statutory-holidays	11.0943	11.4801	10.6	10.3809
vacation	below-average	generous	below-average	below-average
long-term-disability-assistance	yes	yes	no	yes
contribution-to-dental-plan	half	half	none	half
bereavement-assistance	yes	yes	no	yes
contribution-to-health-plan	full	full	none	full
class	good	good	bad	bad

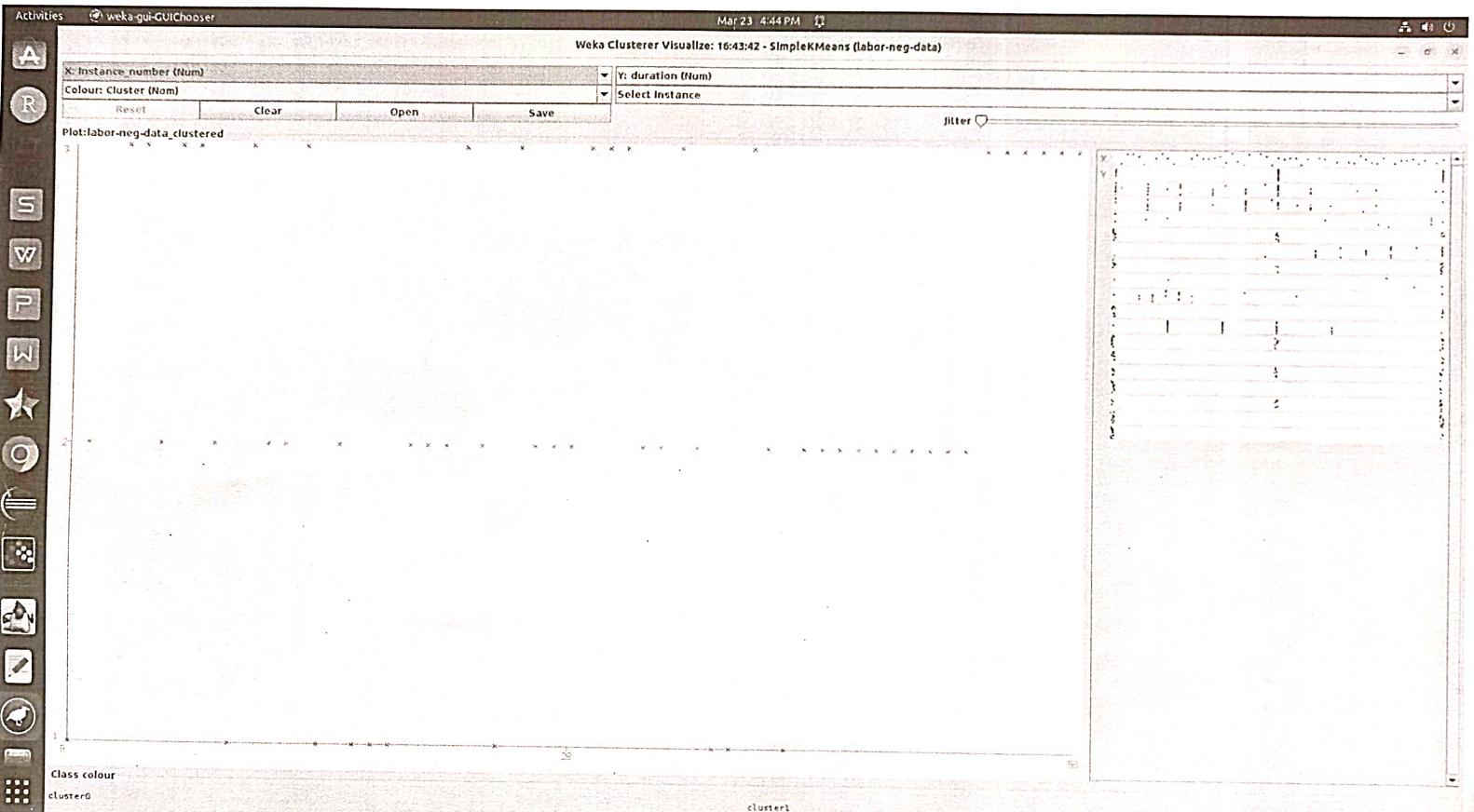
Time taken to build model (full training data) : 0 seconds

--- Model and evaluation on training set ---

Clustered Instances

	0 (62%)	1 (9%)	2 (29%)
0	36	5	16

Status OK Log





Date :

data points and centroid would be minimum. It is to be understood that less variation within the clusters will leads to more similar data points within same clusters.

Working of k-means Algorithm: →

We can understand the working of k-means clustering algo. with the following steps:

- Step 1: First, we need to specify the number of clusters, k , need to be generated by this algorithm.
- Step 2: Next, randomly select k data points and assign each data point to a cluster.
- Step 3: Now it will compute the cluster centroid.
- Step 4: Next, keep iterating the following until we find optimal centroid which is the assignment of data point to the clusters that are not changing any more -

Result:

This program has been successfully executed.

Viva Questions:

- ① What is the purpose of k-means algorithm?
→ It is used to find group which have not been explicitly labeled in the data. It is unsupervised learning algo. which groups the unlabeled datasets into different clusters.
- ② How do we use k-means clustering algo. in weka?
→ Step 1: open weka and load dataset.
Step 2: Find cluster tab in explorer and choose button to execute clustering and select simple-k-means algo.
Step 3: Then to the right of choose icon, press text button and enter three for number of clusters and leave the seed value alone.



Date :

Step 4 : The choice to use training set is selected and then 'start' button is pressed.

Step 5 : The centroid of each cluster is shown in result window and right click the result set on result , selecting to visualize cluster assignment from list column.

③ what does k-mean in k-mean algorithm?

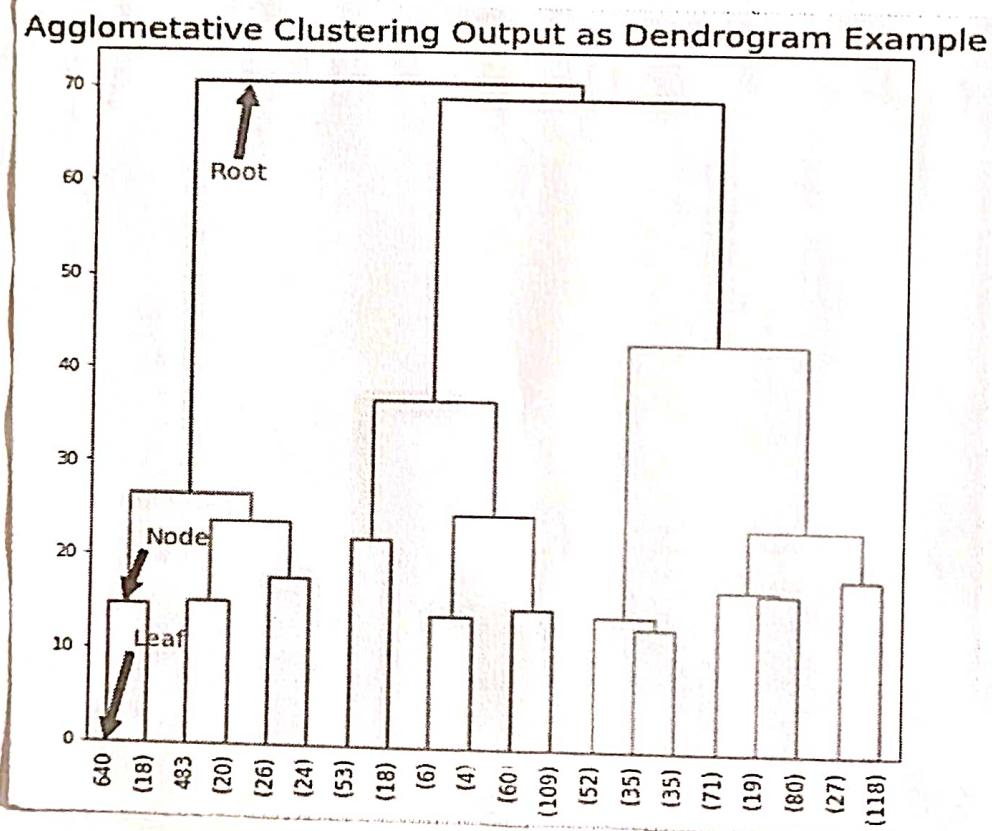
→ The number of clusters found from duty by the method is denoted by the letter 'k' in k-means.

④ what type of algorithm is k-means?

→ It is simple unsupervised learning algorithm . It is widely used centroid-based clustering algorithm.

Practical No.6

Aim: To demonstrate hierarchical clustering on given dataset.





Aim: To demonstrate hierarchical clustering on given data set.

Theory:

Clustering is one of the most common exploratory data analysis techniques used to get an intuition about the structure of the data. It can be defined as the task of identifying subgroups in the data such that data point in the same subgroup (cluster) are very similar while data point in different clusters are very different. In other words, we try to find homogeneous subgroups within the data such that data points in each cluster are as similar as possible w.r.t. to the similarity measure such as Euclidean-based distance or correlation based distance. The decision of which similarity measure to use is application-specific.

Hierarchical clustering: →

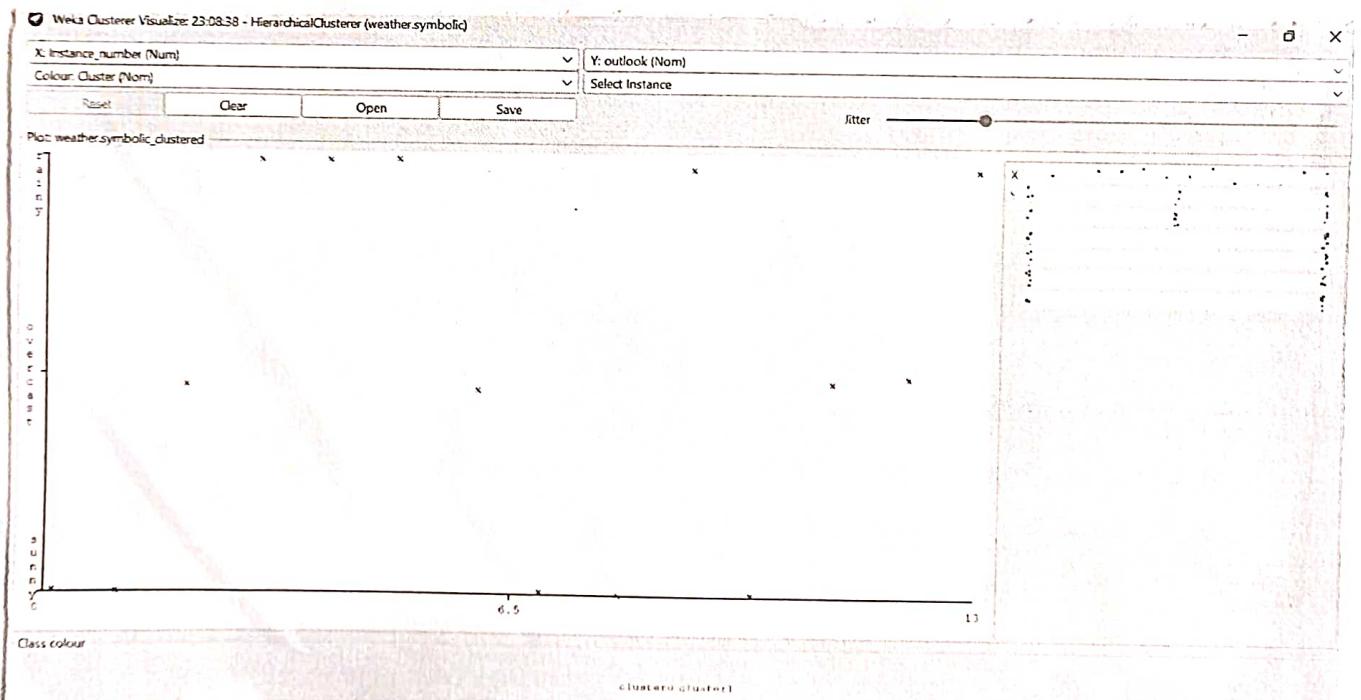
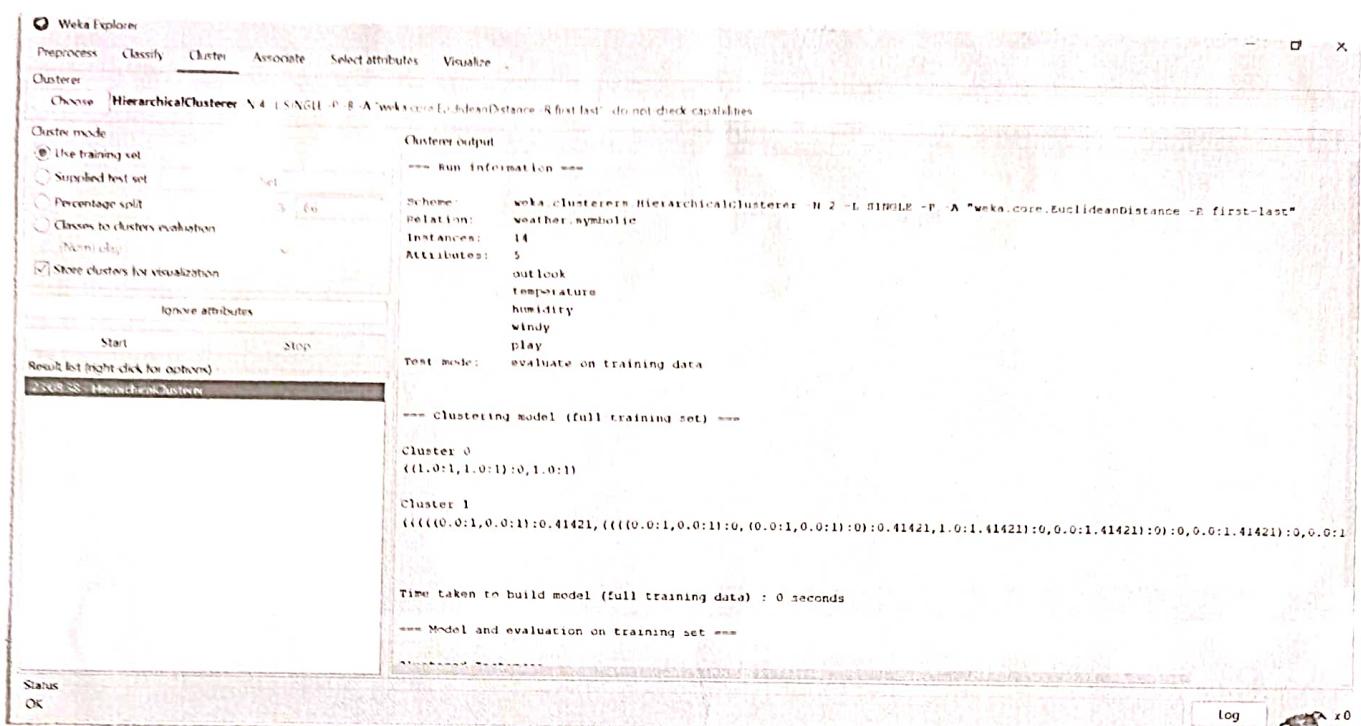
It is also known as Hierarchical cluster analysis or HCA, is an unsupervised clustering approach that includes forming group with a top-to-bottom order.

The hierarchical clustering techniques has two approaches:

- ① Agglomerative: It is bottom up approach, in which the algo starts with taking all data points as single clusters and merging them until one cluster is left.
- ② Divisive: It is the reverse of agglomerative algo. as it is a top-down approach.

Agglomerative Hierarchical clustering: →

Agglomerative clustering or bottom up clustering, essentially started from an individual cluster (each data point is considered as an individual cluster, also called leaf), then every cluster calculates their distance with each other page No. The



Conclusion: Thus, we have implemented hierarchical algo. for clustering successfully.



Date :

two clusters with the shortest distance with each other would merge creating what we called node. Newly formed clusters once again calculating the members of their cluster distance with another cluster outside of their cluster. The process is repeated until all the data points assigned to one cluster called root. The result is a tree-based representation of the objects called dendrogram.

Conclusion :

Thus we have implemented hierarchical clustering algorithm for clustering successfully.

Viva Voce :

① What is hierarchical clustering algorithm?

→ Hierarchical clustering, also known as hierarchical cluster analysis or HCA, is an unsupervised clustering approach that include forming groups with a top-to-down (bottom) order.

② What are the various types of hierarchical clustering?

→ There are two types of hierarchical clustering, they are:

a) Agglomerative

b) Divisive.

③ What is dendrogram in HCA?

→ Dendrogram - is a diagram that shows the hierarchical relationship between objects.

④ Explain the different linkage method used in HCA?

→ ① Single linkage: →

Single linkage returns minimum distance b/w two point, where each points belong to two different clusters.



Date :

② complete linkage : →

It returns the maximum distance betⁿ each data point.

③ Average linkage : →

It returns the average of distances betⁿ all pairs of data points.

④ centroid linkage : →

It returns the distance betⁿ centroid of clusters.

Practical 7

Aim: Implement Apriori Algorithm.

Transaction ID	onion	Potato	Burger	Milk	Beer
t1	1	1	1	0	0
t2	0	1	1	1	0
t3	0	0	0	1	1
t4	1	1	0	1	0
t5	1	1	1	0	1

weka-gui-GUIChooser Apr 5, 4:47 PM

Preprocess Classify Cluster Associate Select attributes Visualize

Associate
Choose: Apriori -N 10 -T 0 -C 0.9 -G 0.005 -U 1.0 -M 0.1 -S 1.0 -c 1

Start Stop
Result list (right) 15:47:09 - Apriori

Associate output
Run information

Scheme: weka.associations.Apriori -N 10 -T 0 -C 0.9 -G 0.005 -U 1.0 -M 0.1 -S 1.0 -c 1

Relation: sunspots

Instances: 365

Attributes:

- date
- plant-stand
- precip
- temp
- hail
- crop-hist
- area-damaged
- severity
- seed-ini
- germination
- plant-growth
- leaf-shade
- leafspots-halo
- leafspots-warg
- leafspot-size
- leaf-thread
- leaf-wilt
- leaf-wild
- size
- lodging
- stem-cankers
- stem-lesion
- fruiting-bodies
- external-decay
- mycelium
- int-discolor
- silverspot
- fruit-pods
- fruit-spots
- seed
- solid-growth
- solid-discolor
- seed-size
- spore-fall
- roots
- class

Associate model (full training set)

Apriori

Minima support: 0.8 (365 instances)
 Minima metric confidence: 0.9
 Number of cycles performed: 4

Generated sets of large itemsets:
 Size of set of large itemsets L(1): 6
 Size of set of large itemsets L(2): 6

Status OK



Date :

Practical NO.7

Aim: Implement Apriori Algorithm.

Theory:

In general association rule mining can be viewed as a two-step process:

- (1) find all frequent itemset: By definition each of these itemsets will occur at least as frequently as a predetermined minimum support count min sup.
- (2) generate strong association rules from the frequency itemsets: By definition these rules must satisfy minimum support and minimum confidence.

Let $I = \{i_1, i_2, i_3, \dots, i_n\}$ be a set of n attributes (called items) and $D = \{t_1, t_2, \dots, t_m\}$ be the set of transaction. It is called a database. Every transaction t_i in D has a unique transaction ID, and it consists of subsets of itemset in I .

A rule can be defined as an implication, $x \rightarrow y$ where x and y are subsets of I ($x, y \subseteq I$) and they have no element in common i.e. $x \cap y = \emptyset$, x and y are the antecedent and the consequent of the rule, respectively.

Let's take an easy example from the supermarket sphere. The example that we are considering is quite small and in practical situation, datasets contain millions or billions of transactions. The set of itemsets, $I = \{\text{Onion}, \text{Burger}, \text{Potato}, \text{Milk}, \text{Beer}\}$ and a database consisting of six transactions. Each transaction is a tuple of 0's and 1's where 0 represents the absence of an item and 1 the presence.

An example for a rule in this scenario would be $\{\text{Onion}, \text{Potato}\} \Rightarrow \{\text{Burger}\}$, which means that if onion and potato are bought, customers also buy a burger.

There are multiple rules possible even from a very small db, so in order to select the interesting ones, we use constraints on various measures of interest and significance. We will look at some

Weka Gui Chooser

Apr 5 14:48 PM

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Clusterer

Choose HierarchicalClusterer -N 4 -L COMPLETE -P -A "weka.core.EuclideanDistance" -R first-last*

Cluster mode

Use training set

Supplied test set

Percentage split

Classes to clusters evaluation

Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

15:42:44 - HierarchicalClusterer
15:42:27 - HierarchicalClusterer
15:42:38 - HierarchicalClusterer

Cluster output

Run information

Scheme: weka.clusterers.HierarchicalClusterer -N 4 -L COMPLETE -P -A "weka.core.EuclideanDistance" -R first-last*

Relation: pima_diabetes

Instances: 768

Attributes: 9

prep

plat

pres

skin

trou

mass

pedi

age

class

Test model/evaluate on training data

Model and evaluation on training set

Cluster 0

((((1.0 0.16243 1.0 0.16243) 0.13765 1.0 0.30008) 0.18289, (((1.0 0.22269 1.0 0.23269) 0.02193 1.0 0.25462) 0.03784, ((1.0 0.12563 1.0 0.16563) 0.11713 1.0 0.22076) 0.0595) 0.02351, ((1.0 0.22279 1.0 0.21279))

Cluster 1

((1.0 0.16243 1.0 0.16243) 0.13765 1.0 0.30008, ((1.0 0.22269 1.0 0.23269) 0.02193 1.0 0.25462) 0.03784, ((1.0 0.12563 1.0 0.16563) 0.11713 1.0 0.22076) 0.0595, ((1.0 0.16243 1.0 0.16243) 0.13765 1.0 0.30008, ((1.0 0.22269 1.0 0.23269) 0.02193 1.0 0.25462) 0.03784, ((1.0 0.12563 1.0 0.16563) 0.11713 1.0 0.22076) 0.0595)

Cluster 2

((1.0 0.59137 1.0 0.59137) 0.0818, ((1.0 0.26259 1.0 0.26259) 0.29203 1.0 0.55571) 0.11745) 0.57401, (((((1.0 0.14403 1.0 0.14403) 0.14403 1.0 0.14403) 0.14403 1.0 0.14403) 0.14403 1.0 0.14403) 0.14403 1.0 0.14403) 0.14403 1.0 0.14403)

Cluster 3

((1.0 0.45825 1.0 0.45825) 0.19461, ((1.0 0.27058 1.0 0.27058) 0.38218) 0.78809, (((1.0 0.0.22864 0.0.22864) 0.19051, ((1.0 0.0.18243 0.0.18243) 0.05164 0.0.0.23407) 0.12345, ((1.0 0.0.18292 0.0.18292) 0.05818 0.0.0.23407) 0.12345))

Time taken to build model (full training data) = 0.32 seconds

Model and evaluation on training set

Clustered Instances

0 205 (29%)
1 495 (63%)
2 19 (3%)
3 19 (3%)

Weka Classifier Tree Visualizer: 15:42:33 - Hierarch...

Status OK

Log



Date :

of these useful measures such as support, confidence, lift and conviction.

Conclusion:

Hence implementation of Apriori algorithm for association rule mining is studied.

Viva Questions:

(1) Define association rule mining?

→ Association rule mining finds interesting associations and relationships among large sets of data items. This rule shows how frequently a itemset occurs in a transaction. A typical example is market Based Analysis. market based analysis is one of the key technique used by large relation to show associations betⁿ items. It allows retailers to identify relationship betⁿ the items that people buy together frequently.

(2) Define apriori algorithm?

→ The Apriori alg. is used for data mining frequent itemsets by devising association rules from a transactional database. The parameters 'support' and 'confidence' are used. Support refers to items' frequency of occurrence; confidence is conditional probability.

(3) What is meant by frequent itemset mining?

→ Frequent pattern mining is a data mining subject with the objective of extracting frequent itemsets from a db. Frequent itemsets play an essential role in many data mining tasks and are related to interesting patterns in data, such as Association Rules.



Date :

(4) Define support and confidence?

→ The number of transactions that include items in the {x} and {y} parts of the rule as a percentage of the total no. of transaction. It is a measure of how frequently the collection of items occur together as a percentage of all transactions.

Practical No. 8

Aim: To implement Fp Growth algorithm.

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... View Edit Save...

Filter Choose None

Current relation

Relation supermarket
Instances 4627

Attributes 217
Sum of weights 4627

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> department1
2	<input type="checkbox"/> department2
3	<input type="checkbox"/> department3
4	<input type="checkbox"/> department4
5	<input type="checkbox"/> department5
6	<input type="checkbox"/> department6
7	<input type="checkbox"/> department7
8	<input type="checkbox"/> department8
9	<input type="checkbox"/> department9
10	<input type="checkbox"/> grocery misc
11	<input type="checkbox"/> department11
12	<input type="checkbox"/> baby needs
13	<input type="checkbox"/> bread and cake
14	<input type="checkbox"/> baking needs
15	<input type="checkbox"/> coupons
16	<input type="checkbox"/> juice-set-cord-ms
17	<input type="checkbox"/> tea
18	<input type="checkbox"/> beer

Remove

Status

OK

Selected attribute

Name	department1	Missing	3580 (77%)	Distinct	1	Type	Nominal
No.	Label	Count	Weight				
1	t	1047	1047.0				

Class: total (Nom)

Visualize All

Log x 0



Aim: To implement FP Growth algorithm.

Theory:

Frequent pattern Growth Algorithm:

This algo. is an improvement to the Apriori method. A FP is generated without the need for candidate generation. FP growth algo. represents the db in the form of tree called a FP tree or FP tree. This tree structure will maintain the association b/w the itemsets. The db is fragmented using one frequent item. This fragmented part is called 'pattern fragment'. The itemset of these fragmented pattern are analyzed. Thus with this method, the search for frequent item sets is reduced comparatively.

FP tree:

FP tree is a tree like structure that is made with the initial items of the database. The purpose of FP tree is to mine the most frequent pattern. Each node of FP tree represent an item of the interest itemset. The root node represents null while the lower nodes represent the itemsets. The association of the nodes with the lower nodes, that is the itemsets with the other itemsets are maintained while forming the tree.

Frequent pattern Algorithm steps:

- ① The first step is to scan the db to find the occurrences of the items in the database. This step is the same as the first step of Apriori. The count of 1-itemsets in the db is called support count or frequency of 1-itemset.
- ② The second step is to construct the FP tree. For this, create the root of the tree. The root is represented by null.

Analyst View

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate Undo Edit... Save

Filter

Choose NominalToBinary -R first-last Apply

Current relation

Relation: supermarket-weka filters unsupervised attribute No. 217 Attributes: 217 Instances: 4627 Sum of weights: 4627

Selected attribute

Name: department1	Type: Numeric
Missing: 3580 (77%)	Distinct: 1
Unique: 0 (0%)	
Statistic	Value
Minimum	0
Maximum	0
Mean	0
StdDev	0

Attributes

All None Invert Pattern

No.	Name
1	department1
2	department2
3	department3
4	department4
5	department5
6	department6
7	department7
8	department8
9	department9
10	grocery misc
11	department11
12	baby needs
13	bread and cake
14	baking needs
15	coupons
16	juice-sat-cord-ms
17	tea
18	...

Remove

Status

OK Log

Class: total (nom) Visualize All

1047



- (3) The next step is to scan the database again and examine the transaction. Examine the first transaction and findout the itemset in it. The itemset with the max count is taken at top, the next itemset with lower count and so on. It means that the branch of the tree is constructed with transaction item sets in descending order of count.
- (4) The next transaction in the database is examined. The item sets are ordered in descending order of count. If any itemset of this transaction is already present in another branch, then this transaction branch would share a common prefix to the root. This means that the common itemset is linked to the new node of another itemset in this transaction.
- (5) Also, the count of the itemset is incremented as it occurs in the transaction. Both the common node and new node count is increased by 1 as they are created and linked all to transaction.
- (6) The next step is to mine the created FP tree. The lowest node is examined first along with the links to the lowest node. It represents the frequency pattern length 1. From this, traverse path in the FP tree. This path and paths are called a conditional pattern base. Conditional pattern base is a sub-database consisting of prefix paths in the FP tree occurring with the lowest node.
- (7) Construct a conditional FP tree which is formed by count of itemsets in the path.
- (8) FP are generated from conditional FP tree.

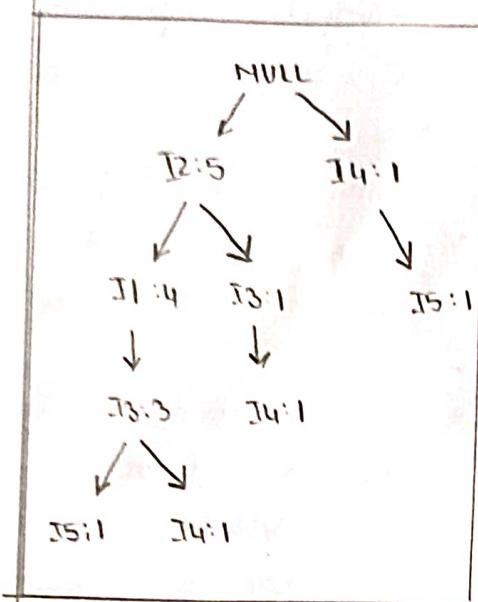
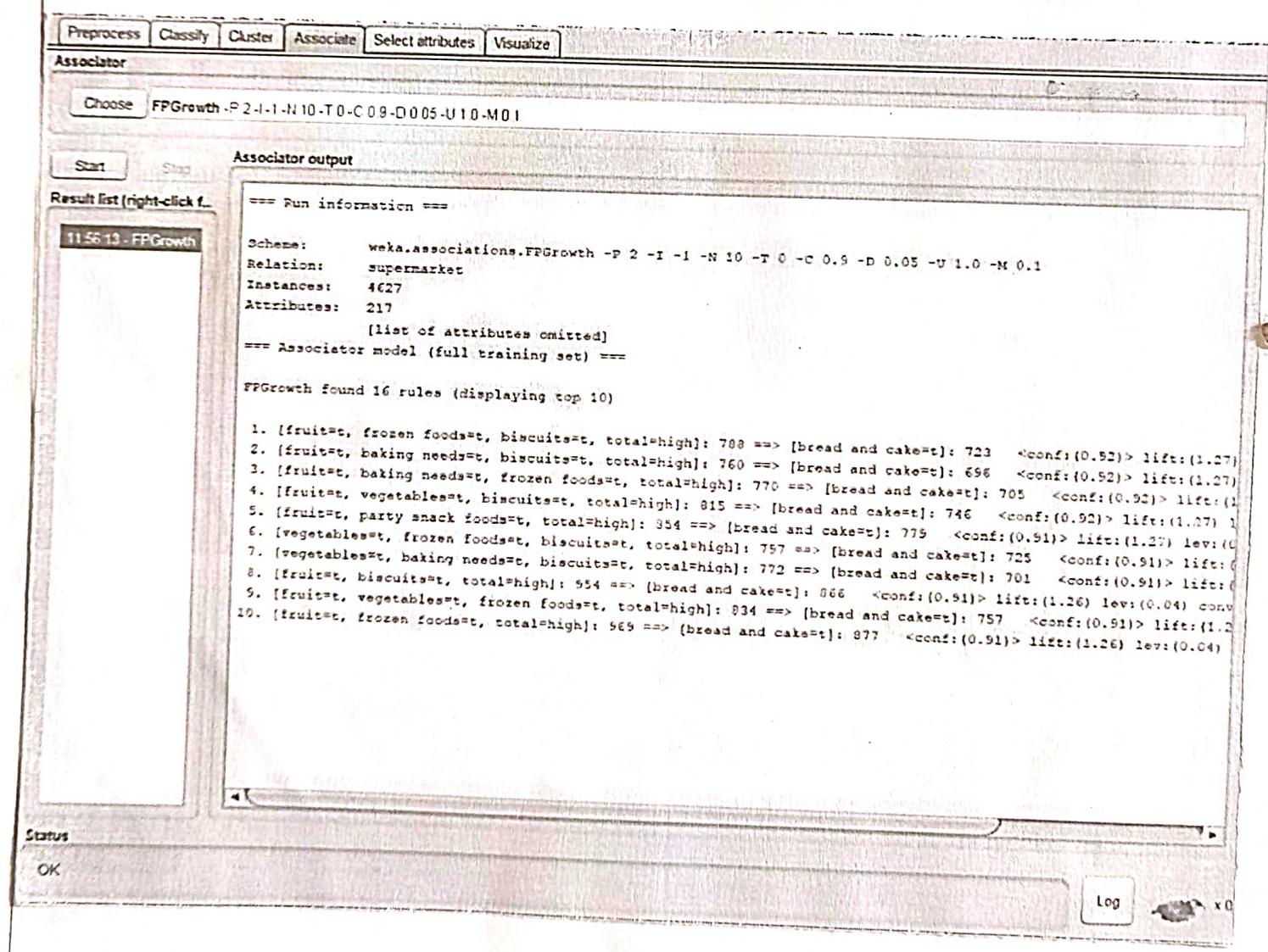


Fig : Build FP Tree





Date :

Example :

Support threshold = 50%. Confidence = 60%.

Transaction	list of items
T1	J1, J2, J3
T2	J2, J3, J4
T3	J4, J5
T4	J1, J2, J4
T5	J1, J2, J3, J5
T6	J1, J2, J3, J4

Solution:

$$\text{Support threshold} = 50\% \Rightarrow 0.5 * 6 = 3 \\ \therefore \text{min.sup} = 3$$

1. Count each item :

item	count
J1	4
J2	5
J3	4
J4	4
J5	2

⑥

2. Sort the itemset in descending order :

item	count
J2	5
J1	4
J3	4
J4	4
J5	2

3. Build FP tree :

① Considering the root node null

② The first scan of transaction T1 : J1, J2, J3 contain page No. Items



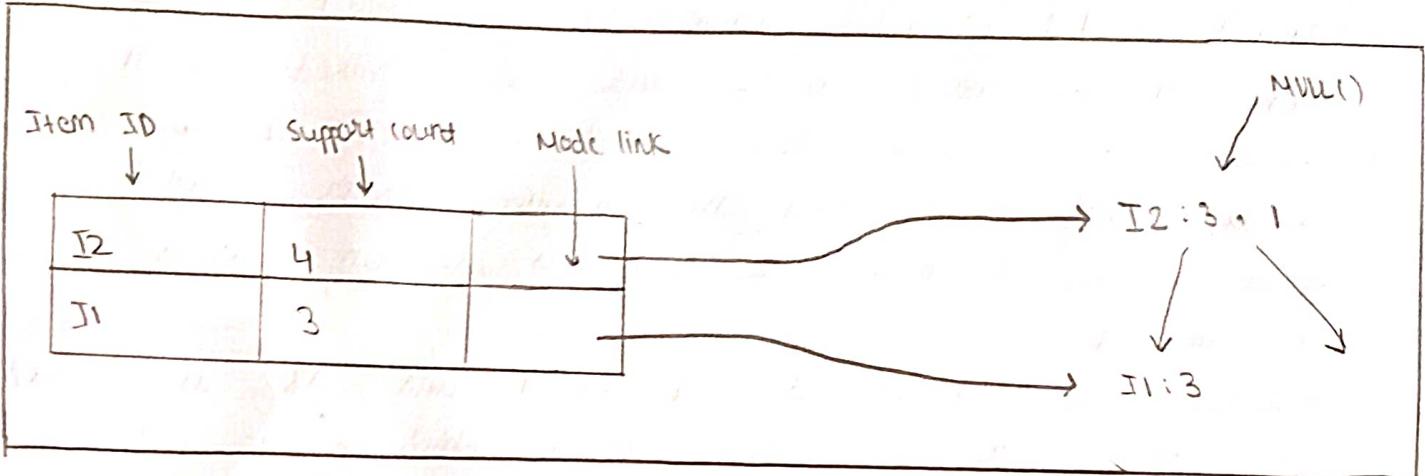
Date:

$\{J_1:1\}, \{J_2:1\}, \{J_3:1\}$ where J_2 is linked as a child to root, J_1 is linked to J_2 and J_3 is linked to J_1 .

- (3) $J_2:J_2, J_3, J_4$ contain J_2, J_3 and J_4 where J_2 is linked to root, J_3 is linked to J_2 and J_4 is linked to J_3 . But this branch would share 1 node J_2 node as common as it is already used in J_1 .
- (4) Increment the count of J_2 by 1 and J_3 is linked as a child to J_2 , J_4 is linked as a child to J_3 . The count is $\{J_2:2\}, \{J_3:1\}, \{J_4:1\}$
- (5) $J_3: J_4, J_5$. similarly, a new branch with J_5 is linked to J_4 as a child is created.
- (6) $J_4: J_1, J_2, J_3$. The sequence will be J_2, J_1 , and J_4 . J_2 is already linked to root node, hence it will be incremented by 1.
- (7) $J_5: J_1, J_2, J_3, J_4$. The sequence will be J_2, J_1, J_3 & J_5
Thus, $\{J_2:4\}, \{J_1:3\}, \{J_3:2\}, \{J_5:1\}$.
- (8) $J_6: J_1, J_2, J_3, J_4$. The sequence will be J_2, J_1, J_3, J_4 .
Thus, $\{J_2:5\}, \{J_1:4\}, \{J_3:3\}, \{J_4:1\}$

4. mining of FP-tree is summarized below:

- (1) The lowest node item 15 is not considered as it does not have a min support count, hence it is deleted.
- (2) The next lower node in J_4 . J_4 occurs in 2 branches $\{J_2, J_1, J_3: J_4:1\}$, $\{J_2, J_3, J_4: 1\}$. Therefore considering 14 as suffix the prefix paths will be $\{J_2, J_1, J_3: 1\}$, $\{J_2, J_3: 1\}$. This forms the conditional pattern base.
- (3) The conditional pattern base is considered a transaction db, on FP tree is constructed. This will contain $\{J_2:2, J_3:2\}$
 J_1 is not considered as it does not meet the main support count.
- (4) This path will generate all combinations of FP :
 $\{J_2: J_4: 2\}, \{J_3, J_4: 2\}, \{J_2, J_3, J_4: 2\}$
- (5) For J_3 , the prefix path would be $\{J_2: J_1: 3\}, \{J_2: 1\}$ this will generate a 2 node FP-tree $\{J_2: 4, J_1: 3\}$ and FP_{Page No.} generated $\{J_2, J_3: 4\}, \{J_1: J_3: 3\}, \{J_2, J_1, J_3: 3\}$.





Date :

- ⑥ For T_1 , the prefix path would be : $\{T_2:4\}$ this will generate a single node FP-tree : $\{T_2:4\}$ and FP tree generate $\{T_2, T_1:4\}$

Item	conditional pattern base	conditional FP-tree	frequent pattern generated
T_4	$\{T_2, T_1, T_3:1\}$ $\downarrow \{T_2, T_3:1\}$	$\{T_2:2, T_3:2\}$	$\{T_2, T_4:2\}, \{T_3, T_4:2\},$ $\{T_2, T_3, T_4:2\}$
T_3	$\{T_2, T_1:3\}, \{T_2:1\}$	$\{T_2:4, T_1:3\}$	$\{T_2, T_3:4\}, \{T_1, T_3:3\}$ $\{T_2, T_1, T_3:3\}$
T_2	$\{T_1:4\}$	$\{T_2:4\}$	$\{T_2, T_1:4\}$

Advantages:

- ① This algo. needs to scan the db only twice when compared to Apriori which scans the transaction for each iteration.
- ② The pairing of items is not done in this algo. and this makes it faster.
- ③ The db is stored in a compact version in memory.
- ④ It is efficient and scalable for mining both long and short frequent patterns.

Inadvantages:

- ① FP trees is more cumbersome and difficult to build than Apriori.
- ② It may be expensive.
- ③ When the db is large, the algo. may not fit in the shared memory.



Date:

Viva Voice: →

① Define FP-Growth algorithm?

→ This algo. is an improvement to the Apriori method. A FP is generated without the need for candidate generation. FP growth algo. represents the db in the form of tree called a FP tree.

② How to construct FP tree?

→ To put it simply, an FP tree is a compressed representation of the input data. It is constructed by reading the dataset one transaction at one time and mapping each transaction onto a path in the FP-tree structure. At different transactions can have the same items, their paths may overlap.

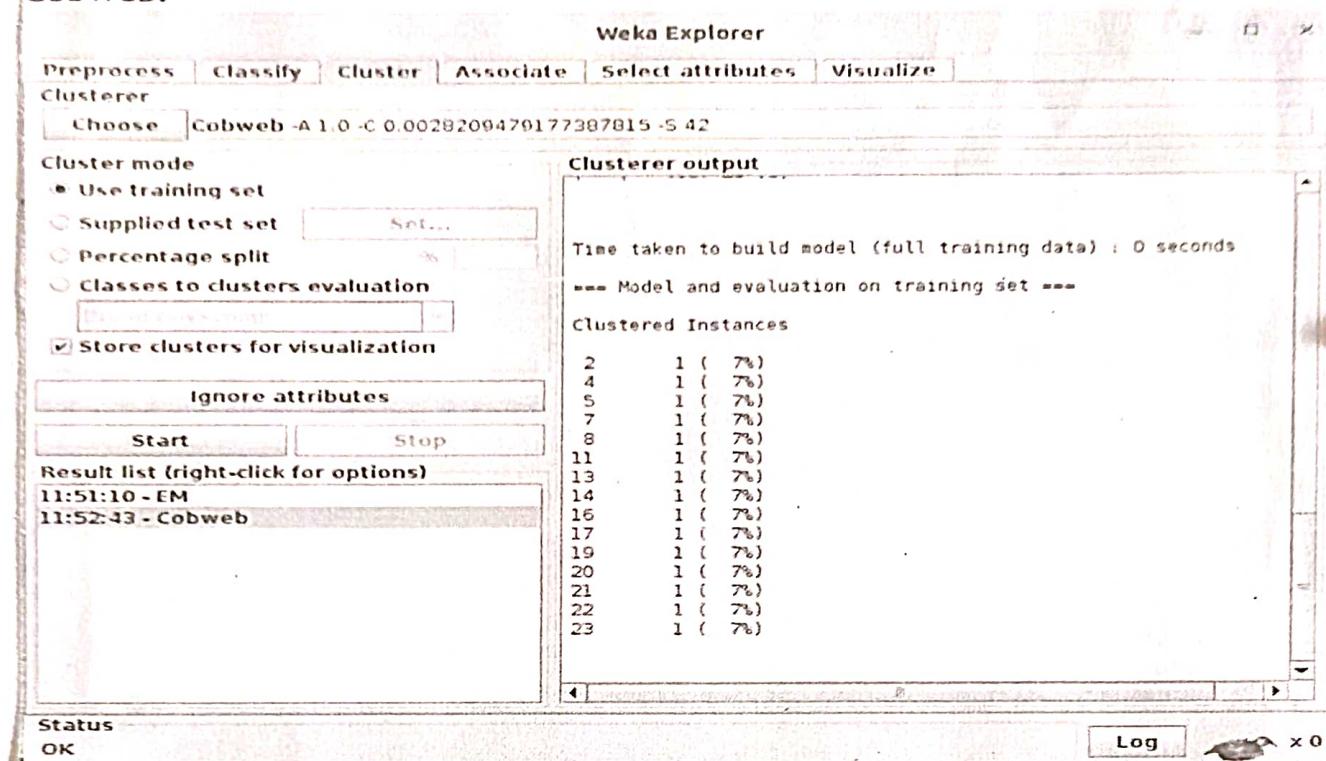
③ Define frequent pattern?

→ FP mining (aka Association rule mining) is an analytical process that finds frequent patterns, association or causal structures from data sets found in various kinds of databases such as relational db, transactional db and other databases repositories.

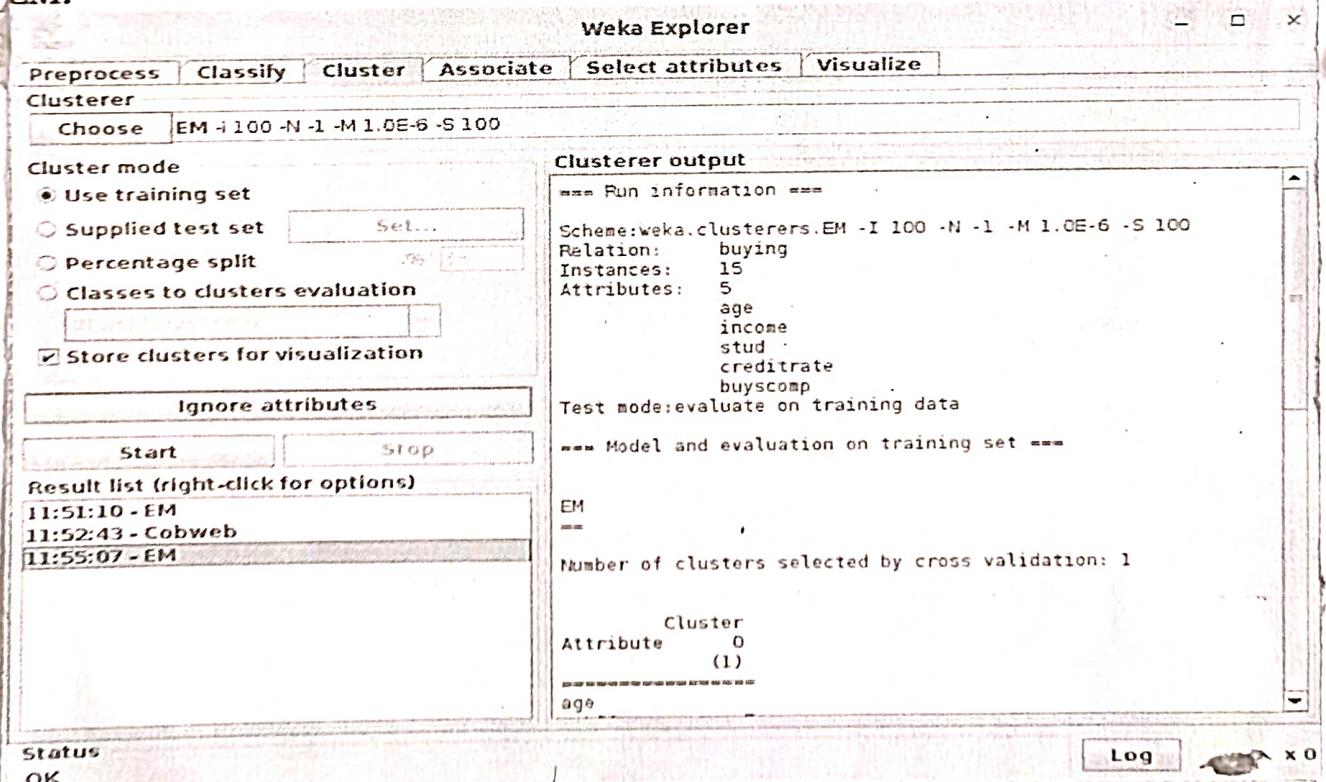
Practical No.9

Aim: Apply Cobweb, EM, further first algorithm on Banking dataset

Cobweb:



EM:





Date :

Practical No. 9

Qn: Apply (obweb, EM, Farthest First) algorithm on Banking data set.

Theory:

OBWEB: OBWEB is an incremental system for hierarchical conceptual clustering. OBWEB was invented by professor Douglas H. Fisher.

OBWEB incrementally organizes observations into a classification tree. Each node in a classification tree represents a class and is labeled by a probabilistic concept that summarizes the attribute-value distribution of object classified under the node. This classification tree can be used to predict missing attribute or the class of a new object.

There are four basic operation OBWEB employs in building the classification tree, which operation is selected depends on the category utility of the classification achieved by applying it. The operations are:

Merging two nodes:

Merging two nodes means replacing them by node whose children is union of the original nodes set of children and which summarizes the attribute-value distribution of all objects classified under them.

Splitting a node:

A node is split by replacing it with its children

Inserting a new node:

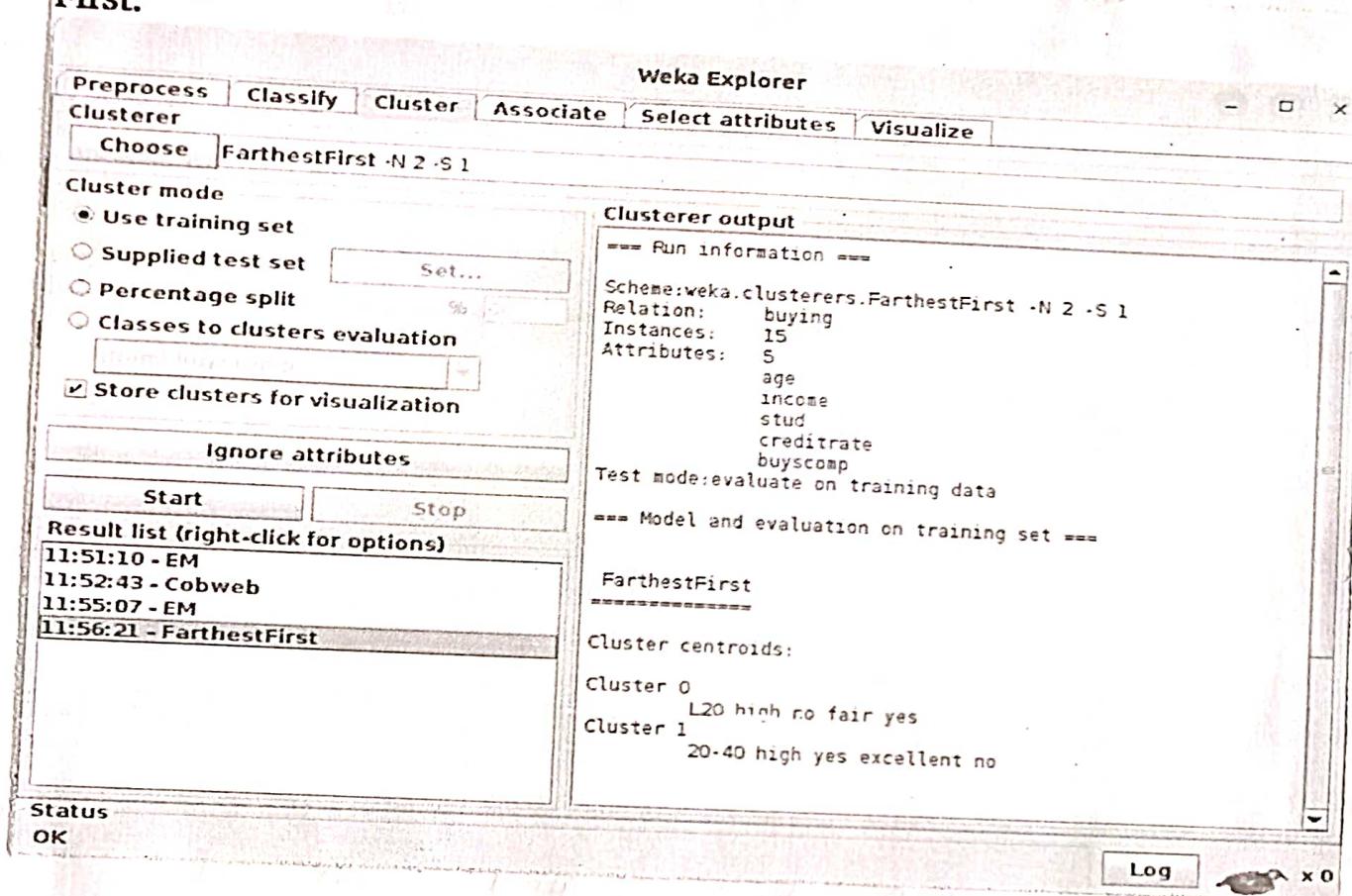
A node is created corresponding to the object being insert into the tree.

Passing an object down the hierarchy:

Effectively calling the OBWEB algorithm on the object and subtree rooted in the node.

Farthest

First:





Date:

The LOBWEB algorithm:

LOBWEB (root, record):

Input: A LOBWEB node root, an instance to insert record

if root has no children then

children := {copy (root)}

newcategory (record) // adds child with record's feature value

insert (record, root) // updates root's statistics.

else

insert (record, root)

for child in root's children do

calculate category utility for merge (record, child),
see best 1, best 2 children w. best U.

end for

if newcategory (record) yields best (U) then

newcategory (record)

else if merge (best 1, best 2) yields best (U) then

merge (best 1; best 2)

LOBWEB (root, record)

else if split (best 1) yields best (U) then

split (best 1)

LOBWEB (root, record)

else

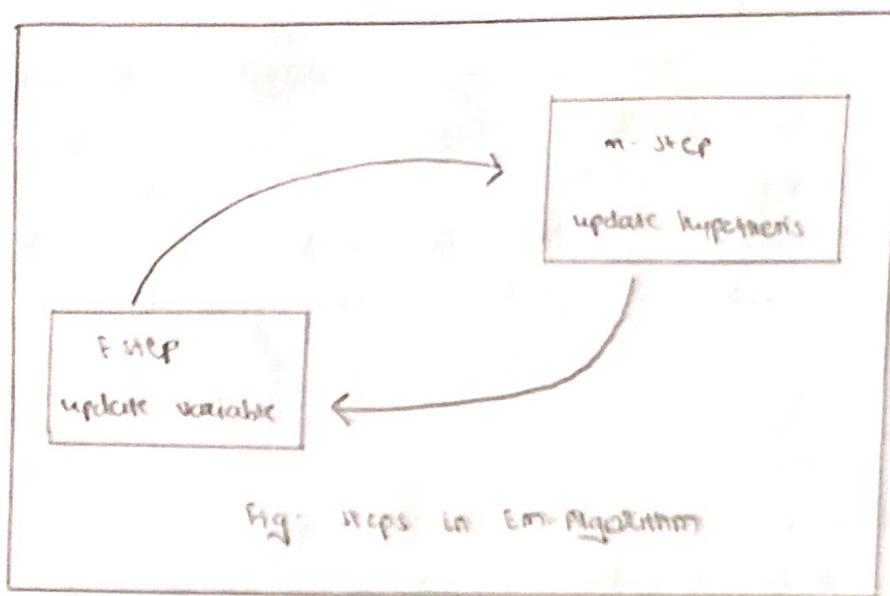
LOBWEB (best 1, record)

end if

end

EM algorithm:

The Expectation - maximization (EM) algorithm is defined as the combination of various unsupervised machine learning algorithms, which is used to determine the local maximum likelihood estimates (MLE) or maximum a posteriori estimate (MAP) for unobservable variables in statistical models. It is also referred to as the latent variable model.





The EM algo. is the combination of various unsupervised ML algo. such as K-means clustering algo. The other model is used to optimize the parameters of the models so that it can explain the data more clearly. The second mode is known as the maximization-step or m-step.

- Expectation step (E-step) : It involves the estimation of all missing values in the dataset so that after completing this step, there should not be any missing value.
- Maximization step (m-step) : This step involves the use of estimated data in the E-step and updating the parameter.
- Repeat E-step and m-step until the convergence of the values occurs.

Farthest First Algorithm:

A farthest - first traversal is a sequence of point in a compact metric space, with each point appearing at most once. The first point of the sequence may be in the space. The first point of the or each point p after the first must have the maximum possible distance from a point to a set is defined as the minimum of the pairwise distances to set points in the set.

Farthest - point traversals may be characterized by the following properties. Fix a number k , and consider the prefix formed by the first k points of the farthest point of the farthest - first traversal of any metric space. Let t be the distance b/w the final point of the prefix and other point in the prefix. Then the subset has the following two properties :

e of Practical

- All pairs of the selected points are at distance at least r from each other, and
- All points of the metric space are at distance at most r from the subset.

working as described here, it also defines initial seeds and then on basis of "k" no. of cluster which we need to know prior. In farthest first it takes - point p_i then choose next point p_i which is at maximum distance. p_i is the centroid and $p_1, p_2 \dots p_n$ are points w object of dataset belongs to cluster from equation.

$$\min \{ \max \text{ dist}(p_i, p_1), \max \text{ dist}(p_i, p_2) \dots \}$$

Farthest first actually solves problem of k-center and it is very efficient for large set of data. It takes centroid arbitrary and distance of one centroid from other is maximum figure show cluster assignment using farthest - first.

Result:

This program has been successfully executed.

Viva Question:

① what is the use of K-Means clustering?

→ It is machine learning algo. that is used to generate predictions based on data.

Teacher's Signature

Name of Practical

(2) What is EM algorithm used for?

→ It is used to find local maximum likelihood parameters of a statistical model in cases where the equation cannot be solved directly.

(3) What are the steps of EM algorithm?

→ There are two steps of EM algorithm, they are

a) Expectation Step (E-step): →

It is used to estimate the missing data in dataset

b) Maximization step (M-step): →

It is used to update the parameter after the complete data is generated in E-step.

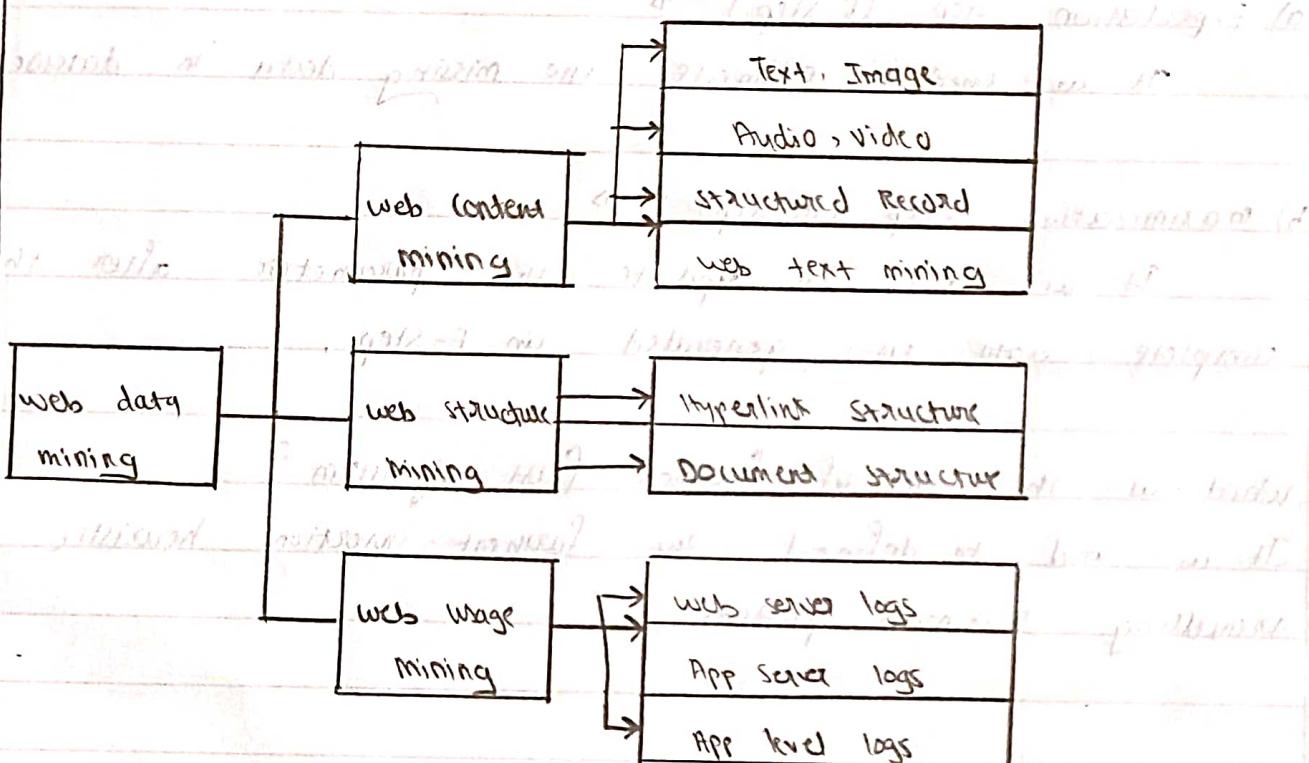
(4) What is the use of Farthest First algorithm?

→ It is used to defined the farthest insertion heuristic for travelling salesman problem.

Teacher's Signature

Practical No. 10

Aim: To compare web data mining techniques (tools and algorithms)





Aim: To compare web data mining techniques, tools and algorithms.

Theory:

Web mining is one of the types of techniques used in data mining. The main purpose of web mining is to automatically extract info. from the web. Info. over the internet is huge and increasing with passage of time due to which size of data bases are also growing. It could be unstructured data, multimedia, table, tag.

Web mining is actually an area of data mining related to the information available on internet. It is concept of extracting informative data available on web pages over the internet. user use different search to fetch their required data from the internet. that informative and user needed data is discovered through mining technique called web mining. Usage, content, are included in information gathered through web mining.

Web mining is sub categorized into three types as shown in figure.

- a) Web content mining.
- b) Web structure mining
- c) Web usage mining.

A) Web content mining: →

Content mining is a process of web mining in which needed informative data is extracted from web sites. Content includes audio, video, text, document etc.. web contents are designed to deliver data to users in the form of text, list, images, video and tables. web content over a last few decades the amount of web pages (HTML) increases to billions and still continues to grow.

web mining categories	Techniques	Tools	Algorithms
web content mining	<ul style="list-style-type: none"> - unstructured data mining - structured data mining - semi-structured data mining - multimedia data mining 	<ul style="list-style-type: none"> - screen scaper - mudenda - Automation Anywhere - web content extractor - web info extractor - rapid miner 	<ul style="list-style-type: none"> - decision tree - Naïves Bayes - support vector machine - Neural Networks
web structure mining	<ul style="list-style-type: none"> - hard-coded classification - hard-coded cluster analysis - linked type - link strength - link cardinality 	<ul style="list-style-type: none"> - google PR checker - link viewer 	<ul style="list-style-type: none"> - Page Rank algorithm - HITS algorithm - weighted page Rank algorithm - distance Rank algorithm - Eigen Ranker algorithm - weighted page content Rank algorithm - web page ranking using links attribute - Query dependent Ranking algorithm
web usage mining	<ul style="list-style-type: none"> - Data preprocessing : data cleaning user and session identification - pattern discovery : statistical analysis association Rule clustering classification sequential pattern 	<ul style="list-style-type: none"> - Data preprocessing Tools: Data Preparator Sumatra TT Lisp Miner Speed Tracer - pattern discovery Tools: SEWEBAR - CMS i-miner Argonaut MiDas 	<ul style="list-style-type: none"> - Association Rule : Apriori Alg. Maxi-min forward References markov chains Fp growth prefix span - clustering : self organized maps graph partitioning ant based technique K-mean alg. Fuzzy C-mean alg.

Name of Practical

B) web structure mining →

web mining techniques are very useful to discover knowledgeable data from web. Structure mining is one of core techniques of web mining which deals with hypertexts structure. Structure mining basically shows the structured summary of the websites. Structure mining analyzes hypertexts of the websites to collect informative data and sort out in categories like similarities and relationship. Structure analysis is also called as link-mining.

c) web usage mining Techniques: →

Following three techniques are described in detail with their sub approaches we in web usage mining. Each technique performs different tasks in a hierarchy.

- Data preprocessing.
- pattern discovery.
- pattern Analysis.

Conclusion:

Thus, the web data mining techniques, tools and algorithm are studied and compared.

Viva Voce:

- (1) what is web mining? How does it differ from regular data mining or text mining?
- web mining is the process of using data mining technique & algo. to extract info. directly from web from documents, services, hypertexts and servers logs.

Teacher's Signature

Web usage mining	- pattern analysis knowledge query mechanism OLAP	- pattern Analysis Tools Webalizer MAVIZ WebVIP2 Web miner Scratchpad	- classification: Decision trees Naïves Bayesian classifier k-nearest neighbour classification Support vector machines
	Intelligent Agents	web miner Scratchpad	- sequential pattern MDAS algorithm (mining internet data for click sequence)

Usage mining techniques comparison:

usage mining technique	method used	Data gathering	Data store	Advantages and disadvantages	Important Algorithm.
Data Preprocessing	Web status codes	- Data logs - website - user login info. - web access log - caches - cookies & etc.	- Web logs	- convert raw data to understandable format LF and extended CLF for recording	- Apriori algo. - FP-growth.
pattern discovery	- frequent median, mode used to show length, recently accessed, view time of pages	- Filtered data from preprocessing section	- session log	- Extract useful info. from discovered pattern correlation	- k-means with genetic algo. - fuzzy k-mean algo.
pattern analysis	- Roll up - Drill Down Up	- pattern analysis	Both Web logs	- irrelevant data and pattern areas separated in pre processing stage	- SQL language - OLAP

Practical

web mining is a subset of data mining that involves processing the data related to the web.

② what are the three main areas of web mining?

- ① web content mining
- ② web structure mining
- ③ web usage mining.

③ what is web content mining?

- web content mining is the browsing and mining of text, images, graphs of web pages to decide the relevance of the content to the search query.

④ what are three web usage mining techniques?

- ① data preprocessing
- ② pattern discovery
- ③ pattern analysis.

Practical No. 11

Aim: To install HADOOP.

Step 1: Download Hadoop from official website. I have chosen Hadoop 2.7.1 for my system. It is a 64 bit version of Hadoop. It has hadoop-2.7.1-bin.tgz file which contains all the required files.

Step 2: Extract the hadoop-2.7.1-bin.tgz file in the hadoop folder.

Step 3: Set environment variable HADOOP_HOME to the path of hadoop folder.

Step 4: Set environment variable HADOOP_CONF_DIR to the path of etc/hadoop folder.

Step 5: Set environment variable HADOOP_LOG_DIR to the path of logs folder.

Step 6: Set environment variable HADOOP_PID_DIR to the path of pid folder.

Step 7: Set environment variable HADOOP_SECURE_DNS_ENABLED to false.

Step 8: Set environment variable HADOOP_HEAPSIZE to 1024.

Step 9: Set environment variable HADOOP_USER_NAME to your user name.

Step 10: Set environment variable HADOOP_IDENT_STRING to your user name.

Practical No. 11

of Practical

Aim: To install Hadoop (Advance Topic).

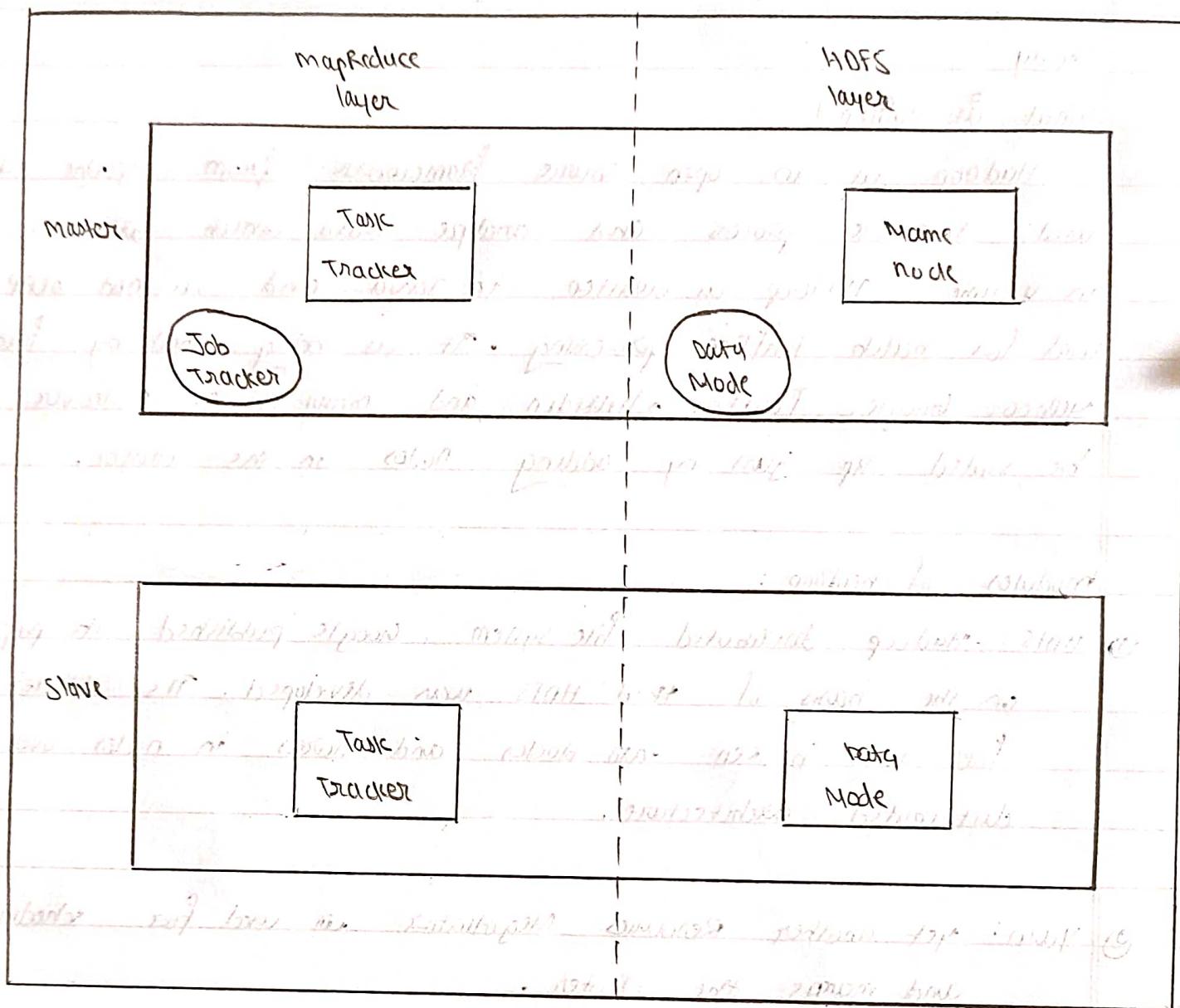
Theory:

What is Hadoop?

Hadoop is an open source framework from Apache and is used to store process and analyze data which are very huge in volume. Hadoop is written in Java and is not MapReduce. It is used for batch offline processing. It is being used by Facebook, Yahoo, Google, Twitter, LinkedIn and many more. Moreover it can be scaled up just by adding nodes in the cluster.

modules of hadoop:

- ① HDFS: Hadoop distributed file system. Google published its paper HDFS on the basis of that HDFS was developed. It states that the files will be broken into blocks and stored in nodes over the distributed architecture.
- ② YARN: yet another Resources Negotiator is used for scheduling and manage the cluster.
- ③ map reduce: This is framework which help Java programs to do the parallel computation on data using key value pair. The map takes input data and convert it into a data set which can be computed in key value pair. The output of map task is consumed by reduce task and then the out of reducer gives the desired result.



Name of Practical

- ④ Hadoop common: These java libraries are used to start Hadoop and are used by other Hadoop modules.

Hadoop Architecture:

The hadoop architecture is a package of the file system, mapreduce engine and the HDFS. The mapReduce engine can be mapreduce / MR1 or YARN / MR2.

A Hadoop cluster consists of a single master and multiple slave nodes. The master node includes Job trackers, Task trackers, NameNode, datanodes whereas the slave node include Datanode and tasktrackers.

Advantages of Hadoop:

- **Fast:** In HDFS the data distributed over the cluster and are mapped which helps in faster retrieval. It also able to process tb of data in minutes and petabytes in hours.
- **Scalable:** Hadoop cluster can be extended by just adding nodes in clusters.
- **Cost effective:** Hadoop is open source and uses commodity hardware to store data so it really cost effective as compared to traditional relational DBMS.
- **Resilient to failure:** HDFS has the property with which it can replicate data over the network, so if one node is down or some other network failure happens, then hadoop takes the other copy of data and use it.

Teacher's Signature

Practical

Introduction to R and R studio:

R is a programming language and s/w environment for statistical computing and graphics supported by the R foundation for statistical computing. R is an integrated suite of s/w facilities for data manipulation, calculation and graphical display. It includes

The term 'environment' is intended to characterize it as a fully planned and coherent system, rather than an incremental accumulation of very specific and inflexible tools, as is frequently the case with other data analysis software. Many users think of R as a statistic system. It is preferable to think of it as an environment within which statistical techniques are implemented. R can be extended via packages.

The R language is widely used among statisticians and data miners for developing statistical software and data analysis.

R studio is an integrated development environment (IDE) for R. It includes a console, syntax-highlighting editor that supports direct code execution as well as tool for plotting, history, debugging and workspace management.

R studio available in open source and commercial editions and runs on the desktop or in a browser connected to R studio server or R studio server pro. R studio is available in two editions R studio desktop, where the program is run locally as a regular desktop application, and R studio server, which allows accessing R studio using a web browser while it is running on a remote linux server. R studio is written in the C++ programming languages and uses the QT framework for its graphical user interface.

daneshva@MoeVB:/S R

```
R version 3.3.1 (2016-06-21) -- "Bug in Your Hair"
Copyright (C) 2016 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)
```

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

of Practical

The R hadoop methods are the collection of packages . It contains three packages ie. rmr, rhbase and rhdbs.

The rmr package:

For hadoop framework , the rmr package provides mapreduce functionality by executing the mapping and Reducing codes in R.

The rhbase package:

This package provides R database management capability with integration with HBASE.

The rhdbs package:

This package provides file management capabilities by integrating with HDFS .

* Installation of R, RStudio and packages for RHadoop:

To install the latest version of R package , CRAN repository should be added to the system . We the following code for this purpose :

```
~$ sudo sh -c 'echo "deb http://cran.cnr.berkeley.edu/bin/linux/ubuntu xenial/" >> /etc/apt/sources.list'
```

The following commands to install the complete R system and compile R package from the source .

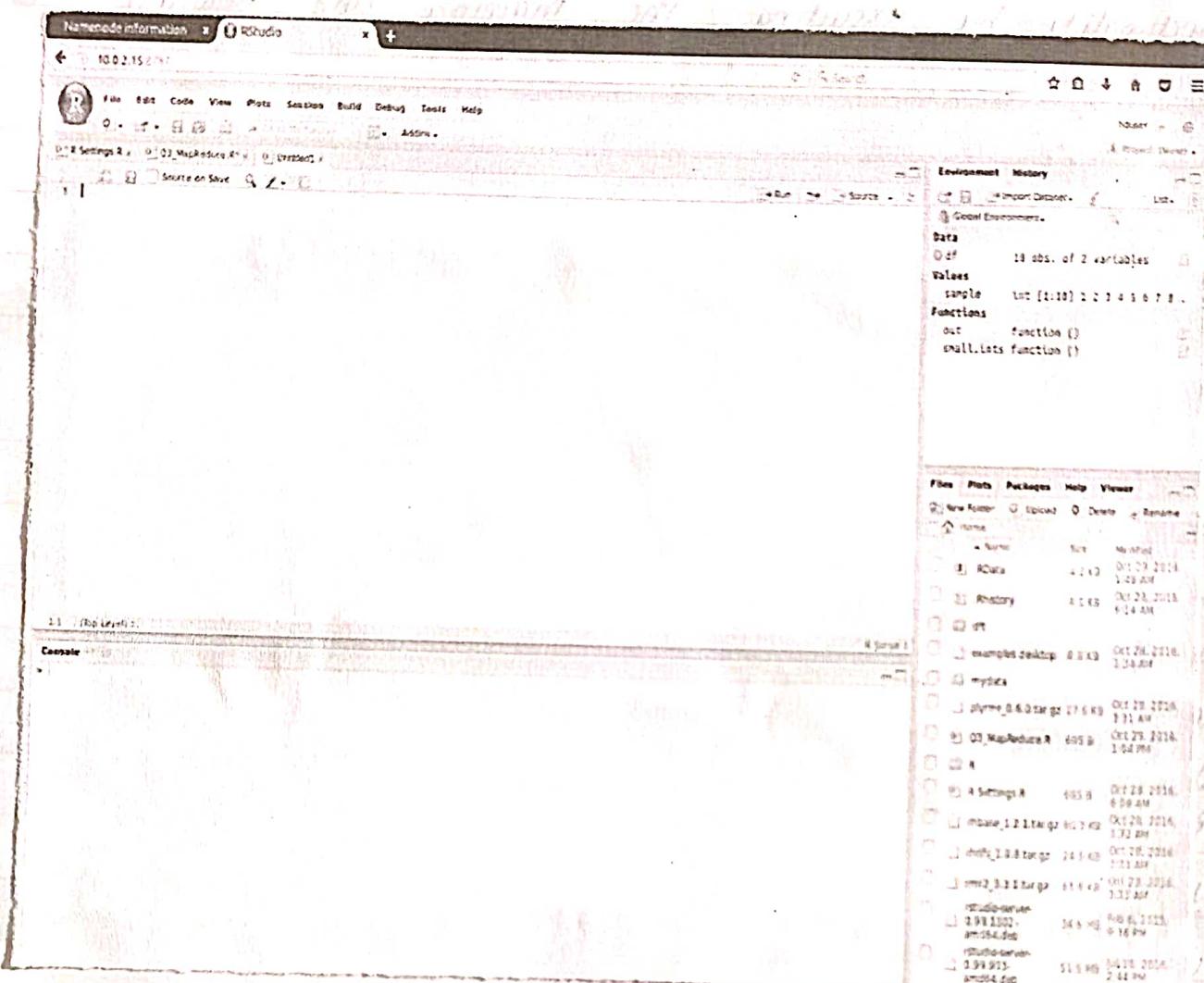
```
~$ sudo apt-get update
```

```
~$ sudo apt-get install r-base
```

```
~$ sudo apt-get install r-base-dev
```

24) In what fall program did you take part in?

Namn och information Rödöna +
← 10.0.2.15:2707



me of Practical

At this point the installation is complete and you can run R with following command.

To quit R, type q() and need to save or not as prompted
~ \$ R

To run R studio server , use following commands :

```
~ $ sudo apt-get install gdebi-core
~ $ wget https://download2.rstudio.org/rstudio-server-1.0.141-amd64.deb
~ $ sudo gdebi rstudio-server-1.0.141-amd64.deb
```

The following command is useful for identification of your VM IP address :

```
~ $ ifconfig
```

RHadoop helps in an integration interaction of R with Hadoop.

RHadoop is collection of R packages that enables R to use "MapReduce" data management. These packages are :

- ① rhdfs
- ② rmr
- ③ rJava

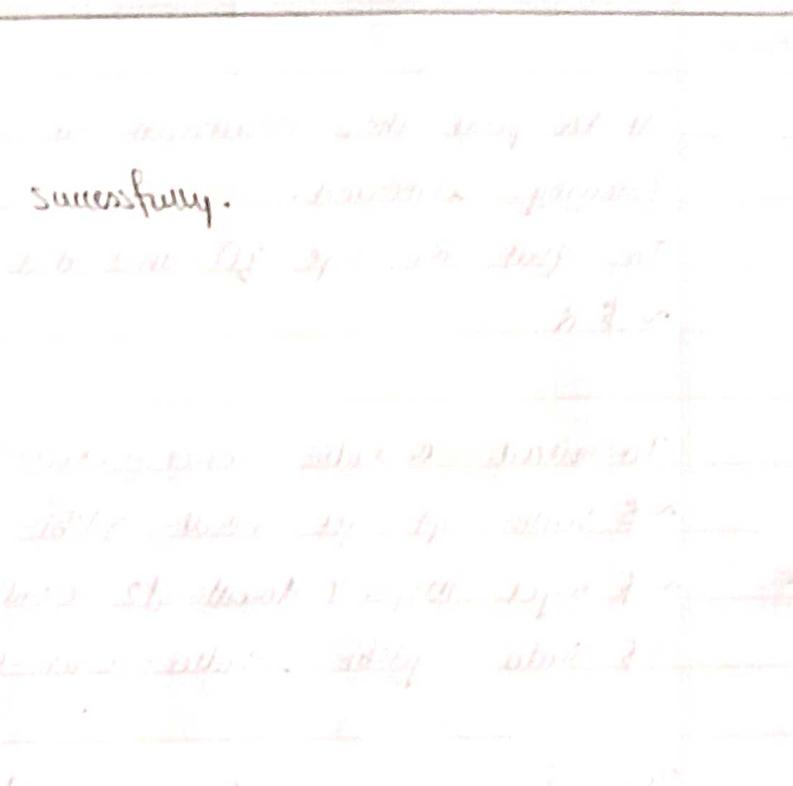
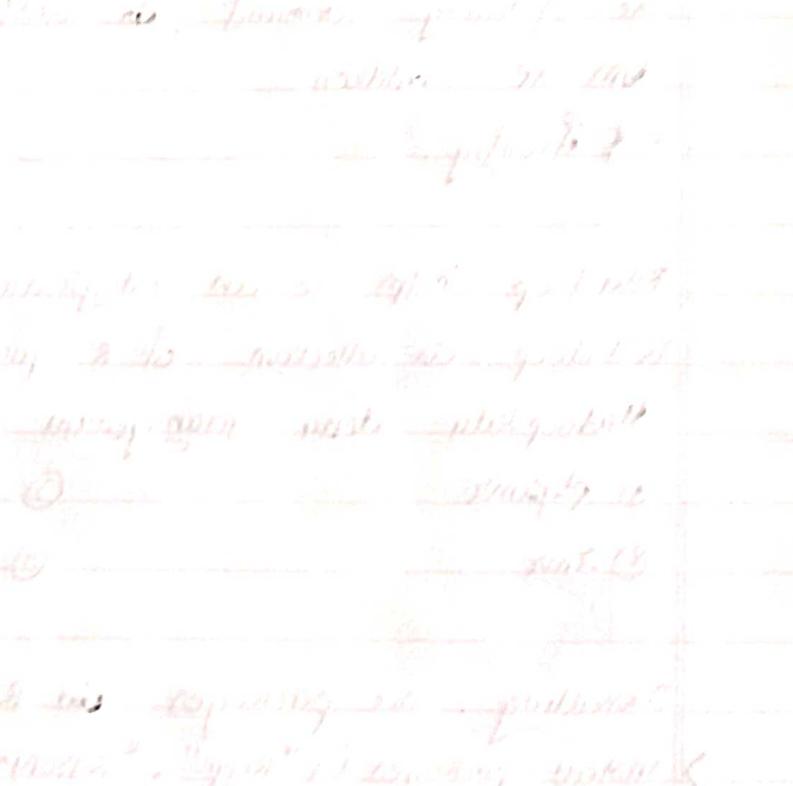
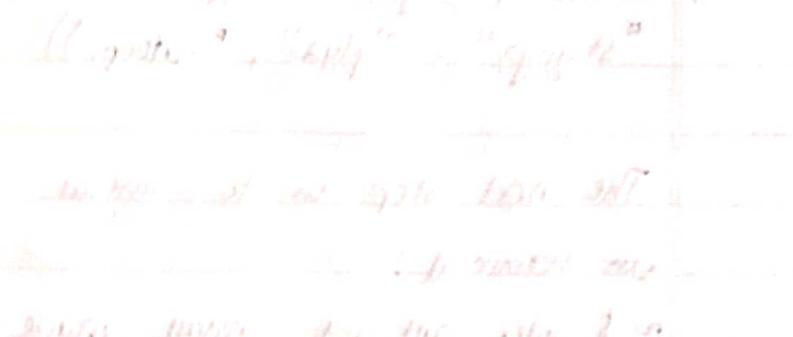
Installing the packages in R environment:

```
> install.packages(c("Rcpp", "RJSONIO", "bitops", "digest", "functional",
"stringr", "plyr", "catools"))
```

The next step is to install Java through following commands in terminal :

```
~ $ sudo apt-get install oracle-java9-installer
~ $ sudo apt-get install openjdk-9-jdk
```

Teacher's Signature

Conclusion: In conclusion, We can conclude that every bit of the code
thus, module is installed successfully. ~~but this application~~ now we
are running our application in the local host. If you see the log file
you will find that there is no error in the log file. So, it is successful.
So, we can say that the application is installed successfully. Now, we have to
check whether the application is working or not. So, we open the browser
and enter the URL which is <http://127.0.0.1:5000>.
After entering the URL, we get the following output:

The screenshot shows a white page with a large red "404" at the top. Below it, the text "Not Found" is displayed in a large, bold, black font. Underneath that, it says "The requested URL was not found on this server." At the bottom, there is a link to "View more details".
So, we can say that the application is installed successfully. Now, we have to
check whether the application is working or not. So, we open the browser
and enter the URL which is <http://127.0.0.1:5000>.
After entering the URL, we get the following output:

The screenshot shows a white page with a large green "Success" message in the center. Below it, it says "The application is running successfully." At the bottom, there is a link to "View more details".
So, we can say that the application is installed successfully. Now, we have to
check whether the application is working or not. So, we open the browser
and enter the URL which is <http://127.0.0.1:5000>.
After entering the URL, we get the following output:

The screenshot shows a white page with a large green "Success" message in the center. Below it, it says "The application is running successfully." At the bottom, there is a link to "View more details".

Name of Practical

```
~$ sudo apt-get install liblzma-dev  
~$ sudo apt-get install x-cran-zipr
```

To install the downloaded packages, move to the directory that the packages are downloaded in. In the following command:

```
~$ cd /  
~$ cd home/lwer/Downloads  
~$ Hadoop rpm="tar.gz | tar | bin | hadoop"  
~$ sudo R (mp) INSTALL phymer 0.6.0.tar.gz  
~$ sudo R (mp) INSTALL zmr2_3.3.1.tar.gz  
~$ sudo R (mp) INSTALL zhadfs_1.0.8.tar.gz  
~$ sudo R (mp) INSTALL zhtbase_1.2.1.tar.gz
```

Now, the Hadoop is ready to use.

Conclusion:

Thus Hadoop is installed successfully.

Viva questions:

a) Explain big data and list the characteristics?

→ Big data is a larger, complex set of data acquired from diverse, new, and old sources of data.

Characteristics of big data:

a) Variety: Variety of big data refers to structured, unstructured and semistructured data that is gathered from multiple resources.

b) Velocity: Velocity essentially refers to the speed at which data is being created in real time.

c) Volume: It indicates large volume of data that is being generated

Teacher's Signature

Name of Practical

on the daily basis.

Q) Explain Hadoop. List the core components of Hadoop.

→ Hadoop is an open source framework from apache and is used to store process and analyse data which are very huge in volume.

The core components of hadoop are HDFS, YARN, MapReduce and Hadoop Common.

Q) What is RStudio?

→ RStudio is an integrated development environment (IDE) for R. It includes a console, syntax highlighting editor,

Q) What is the need of R-Studio?

→ It is needed for data analysis to import, access, transform, explore, plot and model data

Teacher's Signature