



Dr. M.G.R.
EDUCATIONAL AND RESEARCH INSTITUTE
DEEMED TO BE UNIVERSITY

University with Graded Autonomy Status

(An ISO 21001 : 2018 Certified Institution)

Periyar E.V.R. High Road, Maduravoyal, Chennai-95. Tamilnadu, India.



RECORD NOTEBOOK

BCS18L11-DATA MINING LAB

2024– 2025 (ODD SEMESTER)

DEPARTMENT

OF

COMPUTER SCIENCE AND ENGINEERING

NAME : J P RAGHUL RAGHAVENDRA

REGISTER NO : 211191101121

COURSE : B.TECH CSE (DS & AI)

YEAR / SEM / SEC : IV / VII / C



Dr. M.G.R.
EDUCATIONAL AND RESEARCH INSTITUTE
DEEMED TO BE UNIVERSITY

University with Graded Autonomy Status

(An ISO 21001 : 2018 Certified Institution)

Periyar E.V.R. High Road, Maduravoyal, Chennai-95. Tamilnadu, India.



BONAFIDE CERTIFICATE

REGISTER NO : 211191101121

NAME OF LAB : BCS18L11-DATA MINING LAB

DEPARTMENT : COMPUTER SCIENCE AND ENGINEERING

Certified that this is the bonafide record of work done by JP RAGHUL RAGHAVENDRA of IV Year B . Tech - CSE (DS & AI), Sec- C in the **BCS18L11-DATA MINING LAB** during the year 2024-2025.

Signature of Lab-in-Charge

Signature of Head of Dept

Submitted for the Practical Examination held on _____

Internal Examiner

External Examiner

Table of Contents

EX.NO	DATE	NAME OF EXERCISE	PAGE NO	SIGN
1		INTRODUCTION TO THE WEKA MACHINE LEARNING TOOLKIT	1	
2		CLASSIFICATION INTRODUCTION TO EXPLORATORY DATA ANALYSIS USING RATTLE AN OPEN SOURCE TOOL(R)	8	
3		INTRODUCTION TO REGRESSION USING RATTLE AN OPEN SOURCE TOOL. (R)	12	
4		USING THE WEKA TOOLKIT – PART 1	16	
5		CLASSIFICATION USING THE WEKA TOOLKIT - PART-2	20	
6		PERFORMING DATA PREPROCESSING FOR DATA MINING IN WEKA	25	
7		PERFORMING CLUSTERING IN WEKA	28	
8		ASSOCIATION RULE ANALYSIS IN WEKA	33	
9		DATA MINING CASE STUDY USING THE CRISP-DM STANDARD	35	
10		DATA MINING CASE STUDY USING THE CRISP-DM STANDARD	50	

EXP:1

DATE :

INTRODUCTION TO THE WEKA MACHINE LEARNING TOOLKIT

AIM :

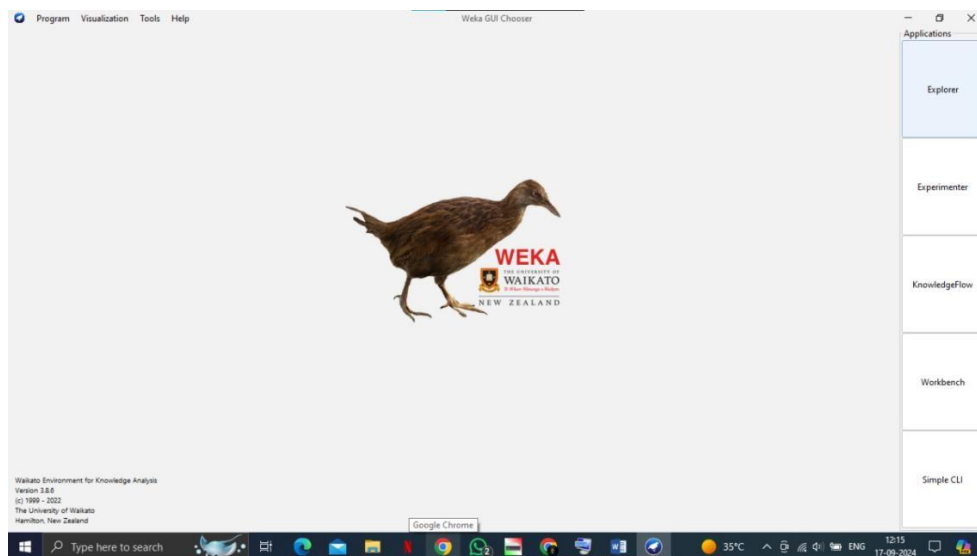
To study WEKA Tool Kit

WEKA is open source java code created by researchers at the University of Waikato in New Zealand. It provides many different machine learning algorithms, including the following classifiers:

1. Decision tree (j4.8, an extension of C4.5)
2. MLP, aka multiple layer perceptron (a type of neural net)
3. Naïve bayes
4. Rule induction algorithms such as JRip
5. Support vector machine

The GUI WEKA

The Weka GUI Chooser (class `weka.gui.GUIChooser`) provides a starting point for launching Weka's main GUI applications and supporting tools. If one prefers a MDI ("multiple document interface") appearance, then this is provided by an alternative launcher called "Main" (class `weka.gui.Main`). The GUI Chooser consists of four buttons—one for each of the four major Weka applications—and four menus.



The buttons can be used to start the following applications:

- **Explorer** : An environment for exploring data with WEKA (the rest of this documentation deals with this application in more detail).
- **Experimenter**: An environment for performing experiments and conducting statistical tests between learning schemes.
- **Knowledge Flow**: This environment supports essentially the same functions as the Explorer but with a drag-and-drop interface. One advantage is that it supports incremental learning
- **Simple CLI**: Provides a simple command-line interface that allows direct execution of WEKA commands for operating systems that do not provide their own command line interface.

The menu consists of four sections:

WEKA Explorer

The user interface

Section Tabs

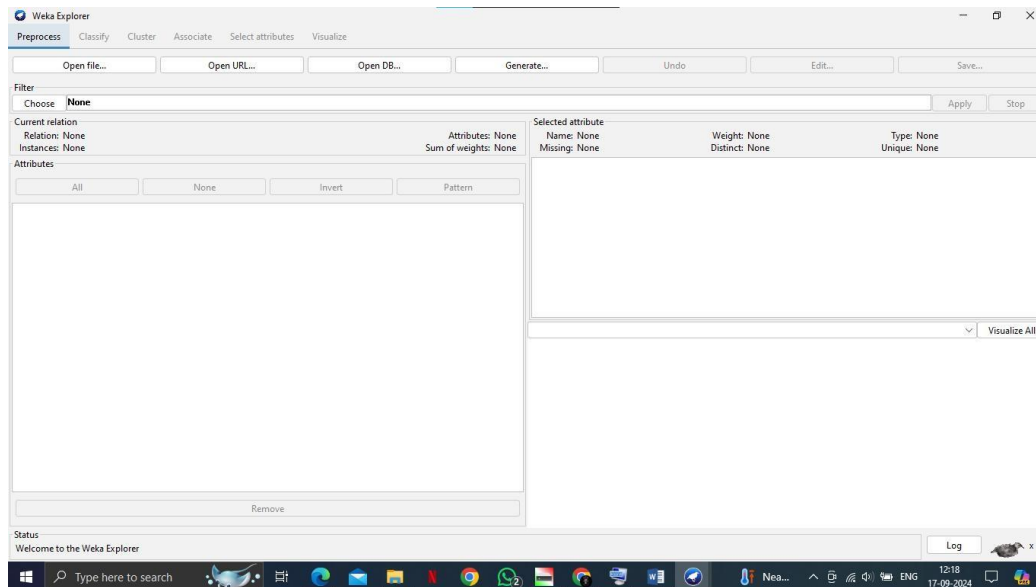
At the very top of the window, just below the title bar, is a row of tabs. When the Explorer is first started only the first tab is active; the others are grayed out. This is because it is necessary to open (and potentially pre-process) a data set before starting to explore the data.

The tabs are as follows:

1. **Preprocess**. Choose and modify the data being acted on.
2. **Classify**. Train and test learning schemes that classify or perform regression.
3. **Cluster**. Learn clusters for the data.
4. **Associate**. Learn association rules for the data.
5. **Select attributes**. Select the most relevant attributes in the data.
6. **Visualize**. View an interactive 2D plot of the data.

Once the tabs are active, clicking on them flicks between different screens, on which the respective actions can be performed. The bottom area of the window (including the status box, the log button, and the Weka bird) stays visible regardless of which section you are in. The Explorer can be easily extended with custom tabs.

1. PREPROCESSING



OPTIONS

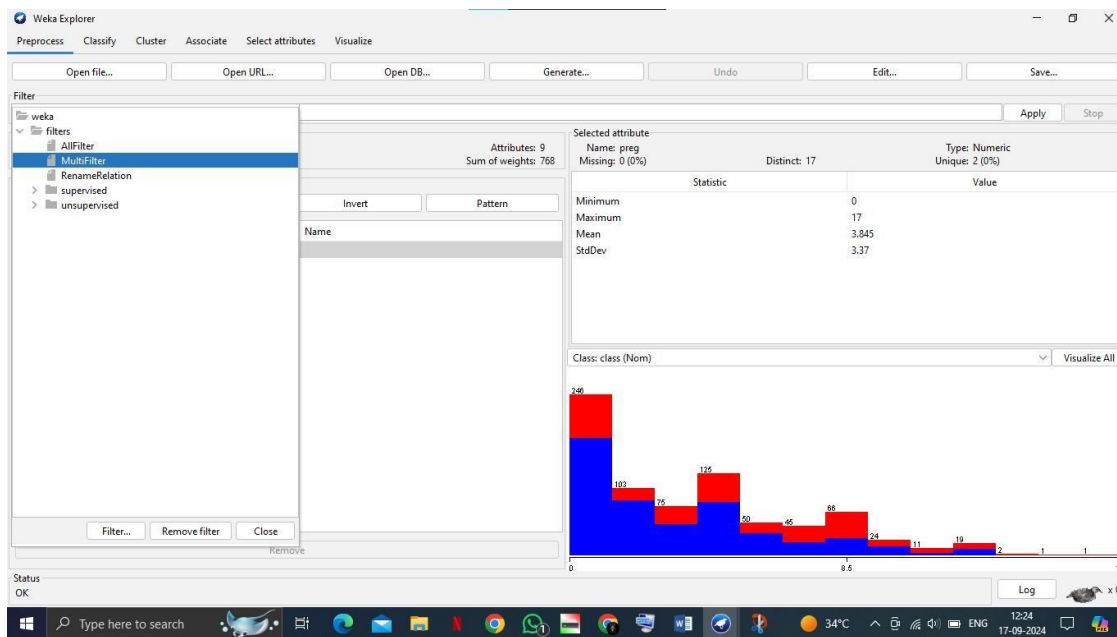
1. **All.** All boxes are ticked.
2. **None.** All boxes are cleared (unticked).
3. **Invert.** Boxes that are ticked become unticked and vice versa.
4. **Pattern.** Enables the user to select attributes based on a Perl 5 Regular Expression. E.g., `.* id` selects all attributes which name ends with `id`. Once the desired attributes have been selected, they can be removed by clicking the Remove button below the list of attributes. Note that this can be undone by clicking the Undo button, which is located next to the Edit button in the top-right corner of the Preprocess panel.

Loading Data

The first four buttons at the top of the preprocess section enable you to load data into WEKA:

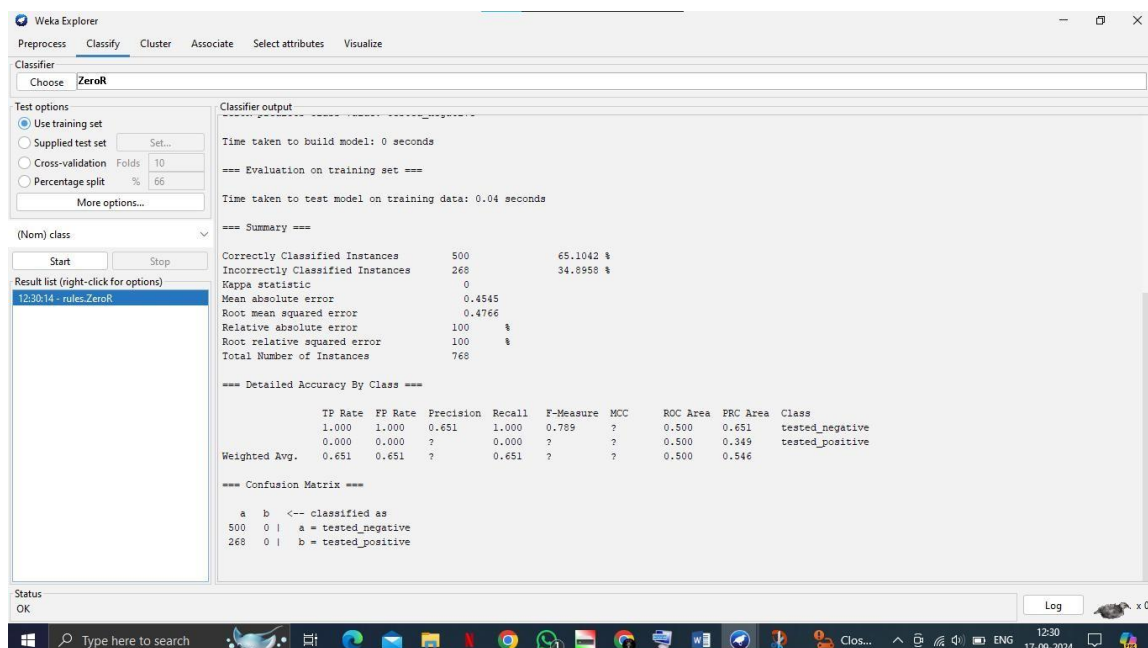
1. **Open file.** Brings up a dialog box allowing you to browse for the data file on the local file system.
2. **Open URL.** Asks for a Uniform Resource Locator address for where the data is stored.
3. **Open DB.** Reads data from a database. (Note that to make this work you might have to edit the file in `weka/experiment/DatabaseUtils.props`.)
4. **Generate.** Enables you to generate artificial data from a variety of Data Generators. Using the Open file. button you can read files in a variety of formats: WEKA's ARFF format, CSV format,

C4.5 format, or serialized Instances format. ARFF files typically have a .arff extension, CSV files a .csv extension, C4.5 files a .data and .names extension and serialized Instances objects a .bsi extension.



2 CLASSIFICATION

The preprocess section allows filters to be defined that transform the data in various ways. The Filter box is used to set up the filters that are required. At the left of the Filter box is a Choose button. By clicking this button it is possible to select one of the filters in WEKA. Once a filter has been selected, its name and options are shown in the field next to the Choose button. Clicking on



this box with the left mouse button brings up a GenericObjectEditor dialog box. A click with the right mouse button (or Alt+Shift+left click) brings up a menu where you can choose, either to display the properties in a Generic Object Editor dialog box, or to copy the current setup string to the clipboard.

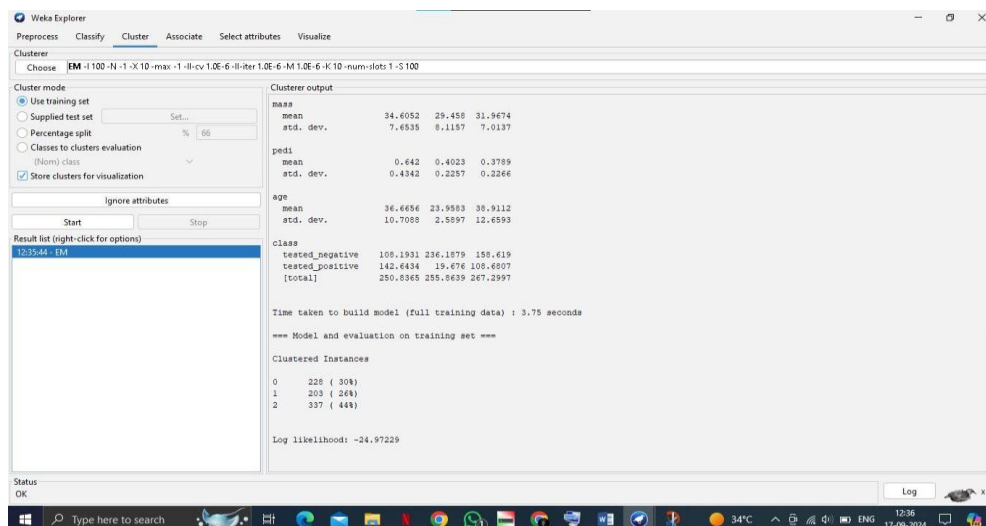
Selecting a Classifier

At the top of the classify section is the Classifier box. This box has a text field that gives the name of the currently selected classifier, and its options. Clicking on the text box with the left mouse button brings up a Generic Object Editor dialog box, just the same as for filters that you can use to configure the options of the current classifier. With a right click (or Alt+Shift+left click) you can once again copy the setup string to the clipboard or display the properties in a Generic Object Editor dialog box. The Choose button allows you to choose one of the classifiers that are available in WEKA.

3. CLUSTERING

Selecting a Clusterer

By now you will be familiar with the process of selecting and configuring objects. Clicking on the clustering scheme listed in the Clusterer box at the top of the window brings up a Generic Object Editor



dialog with which to choose a new clustering scheme.

Cluster Modes

The Cluster mode box is used to choose what to cluster and how to evaluate the results. The first three options are the same as for classification: Use training set, Supplied test set and Percentage split except that now the data is assigned to clusters instead of trying to predict a specific class. The fourth mode, Classes to clusters evaluation, compares how well the chosen clusters match up with a pre-assigned class

in the data. The drop-down box below this option selects the class, just as in the Classify panel. An additional option in the Cluster mode box, the Store clusters for visualization tick box, determines whether or not it will be possible to visualize the clusters once training is complete. When dealing with datasets that are so large that memory becomes a problem it may be helpful to disable this option

4. ASSOCIATING

SETTING UP : This panel contains schemes for learning association rules, and the learners are chosen and configured in the same way as the clusterers, filters, and classifiers in the other panels

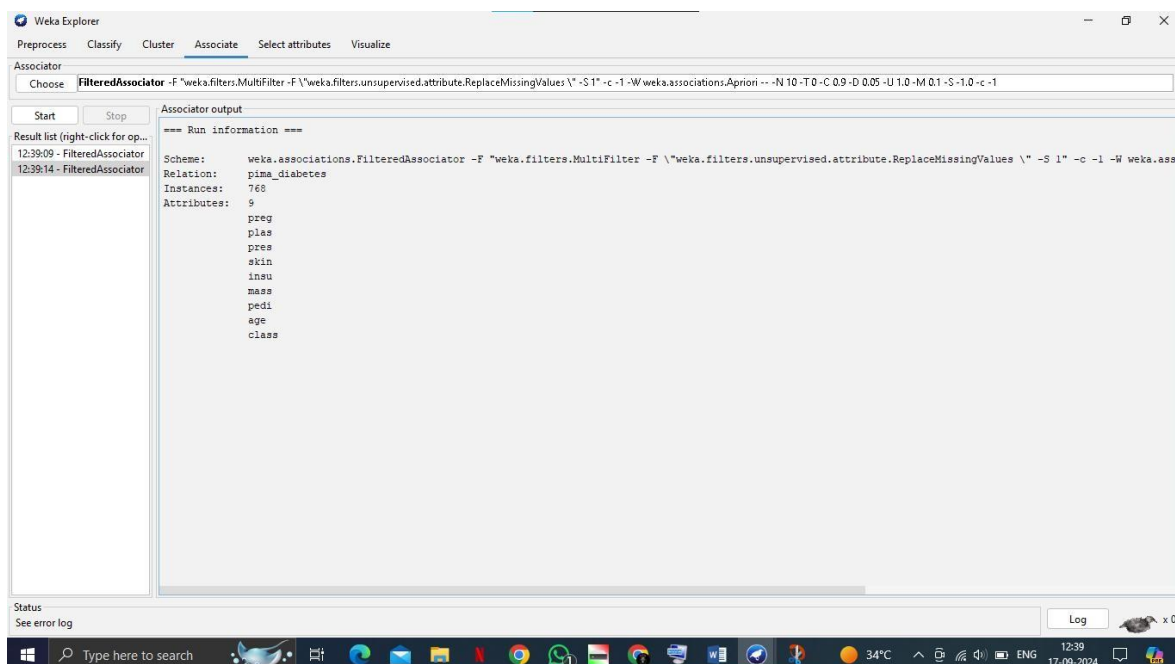
Learning Associations

Once appropriate parameters for the association rule learner have been set, click the Start button.

When complete, right-clicking on an entry in the result list allows the results to be viewed or saved.

Selecting Attributes

Searching and Evaluating Attribute selection involves searching through all possible combinations of attributes in the data to find which subset of attributes works best for prediction. To do this, two objects must be set up: an attribute evaluator and a search method. The evaluator determines what method is used to assign a worth to each subset of attributes. The search method determines what style of search is



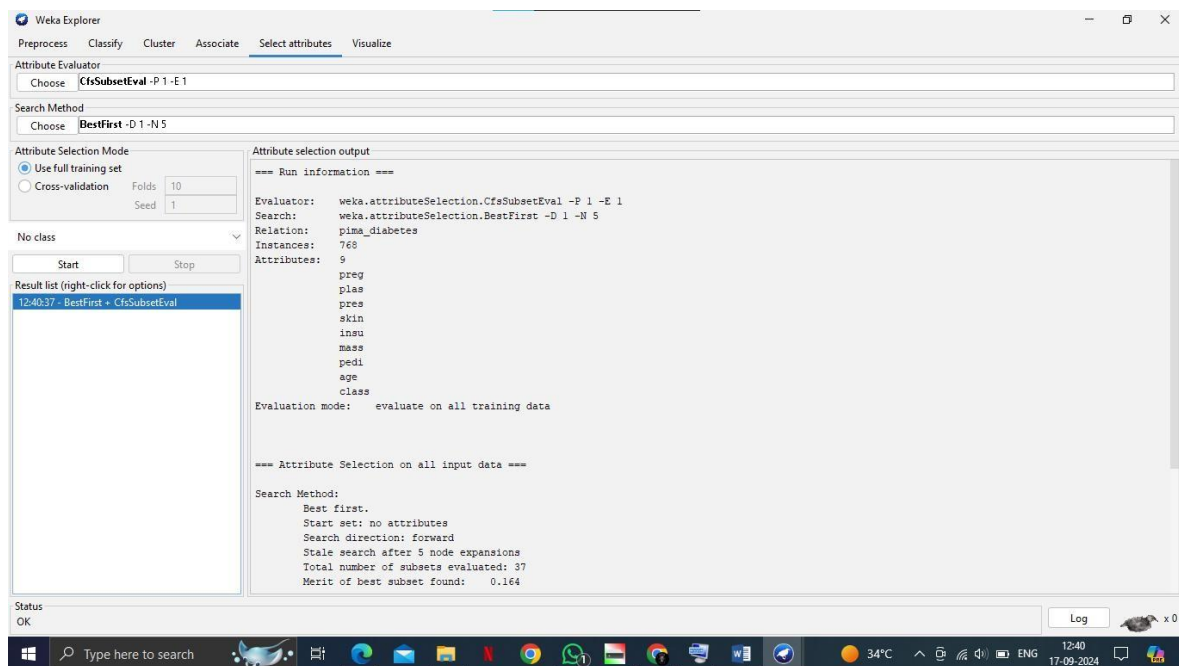
performed.

Options

The Attribute Selection Mode box has two options:

Use full training set. The worth of the attribute subset is determined using the full set of training data.

Cross-validation. The worth of the attribute subset is determined by a process of cross-validation. The Fold and Seed fields set the number of folds to use and the random seed used when shuffling the data. As with Classify, there is a drop-down box that can be used to specify which attribute to treat as the class



Result:

The WEKA tool kit is studied and output is verified

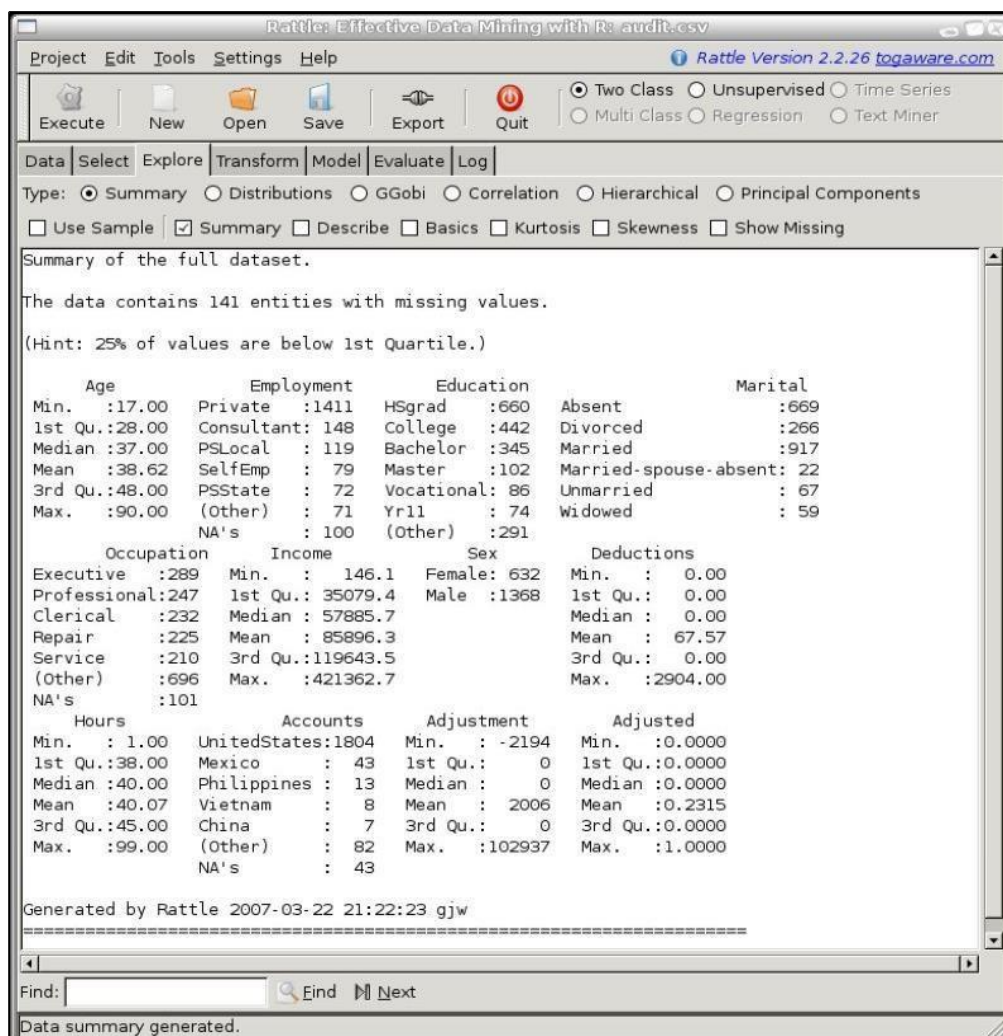
CLASSIFICATION INTRODUCTION TO EXPLORATORY DATA ANALYSIS USING RATTLE AN OPEN SOURCE TOOL(R)

AIM :

To Conduct Exploratory Data Analysis Using Rattle

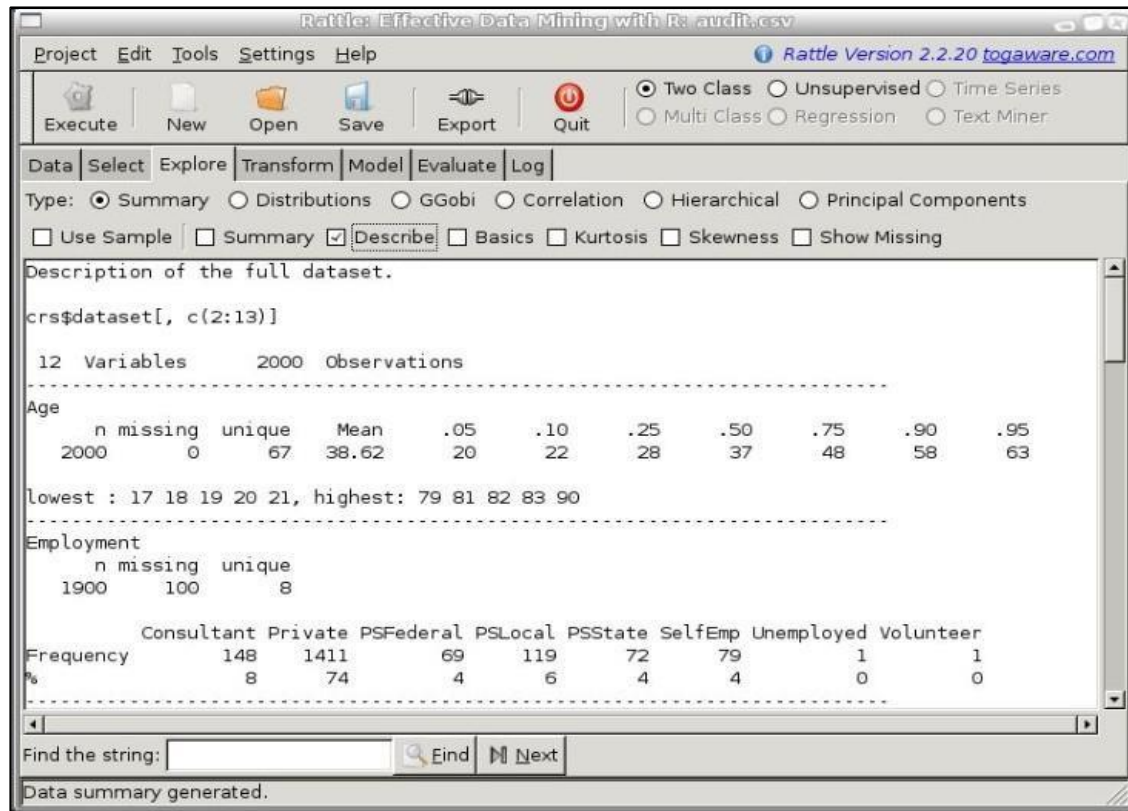
(A) SUMMARY

The Summary check box provides numerous measures for each variable, including, in the first instance, the minimum, maximum, median, mean, and the first and third quartiles. Generally, if the mean and median are significantly different then we would think that there are some entities with very large values in the data pulling the mean in one direction. It does not seem to be the case for Age but is for Income.

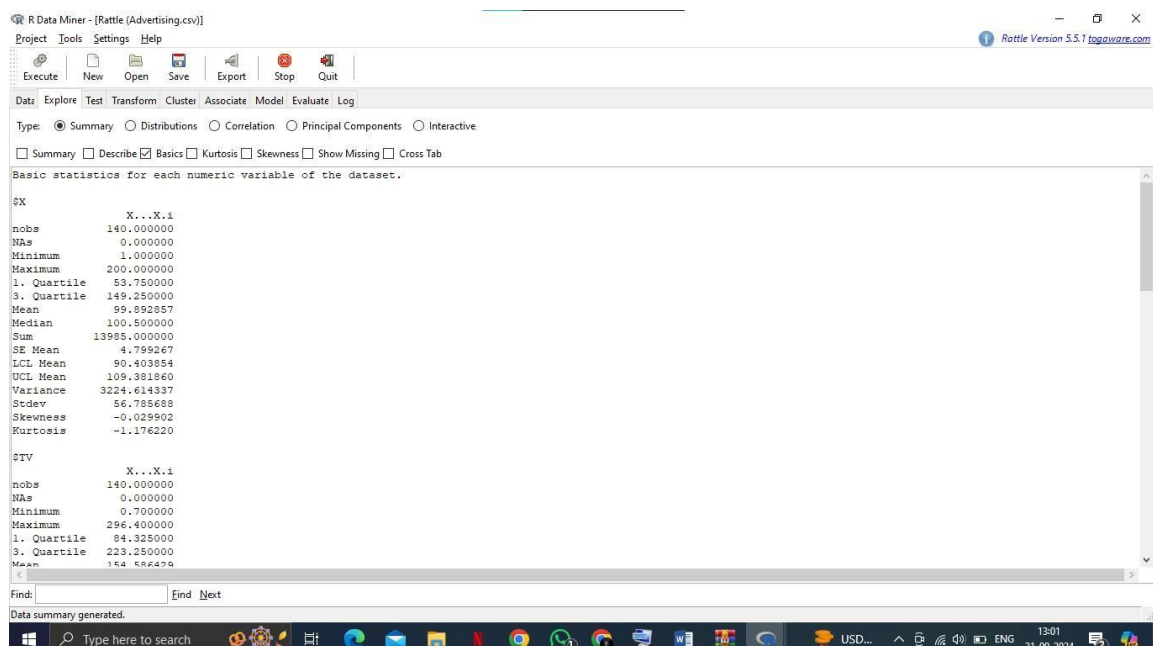


(B) DESCRIBE

The Describe check box

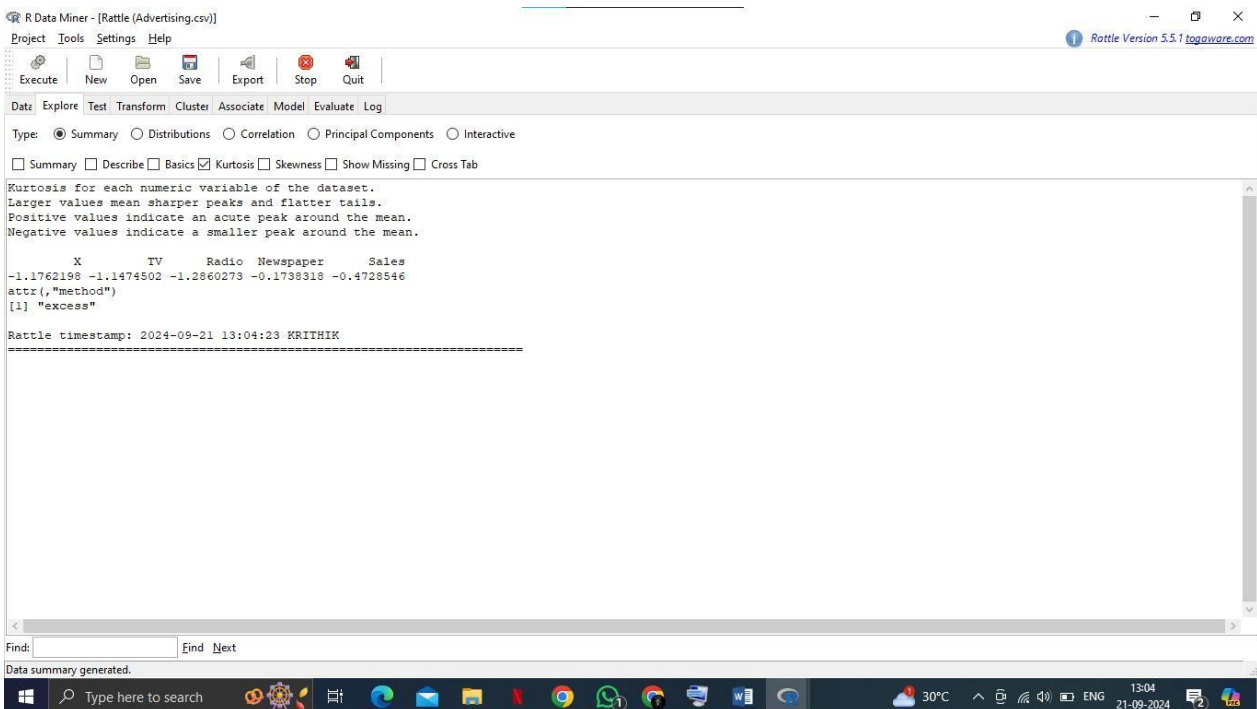


(C) BASICS



(D) KURTOSIS

The kurtosis is a measure of the nature of the peaks in the distribution of the data. A larger value for the kurtosis will indicate that the distribution has a sharper peak, as we can see in comparing the distributions of Income and Adjustment. A lower kurtosis indicates a smoother peak.



R Data Miner - [Rattle (Advertising.csv)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: ☒ Summary ☐ Distributions ☐ Correlation ☐ Principal Components ☐ Interactive

☐ Summary ☐ Describe ☐ Basics ☒ Kurtosis ☐ Skewness ☐ Show Missing ☐ Cross Tab

Kurtosis for each numeric variable of the dataset.
Larger values mean sharper peaks and flatter tails.
Positive values indicate an acute peak around the mean.
Negative values indicate a smaller peak around the mean.

X	TV	Radio	Newspaper	Sales
-1.1762198	-1.1474502	-1.2860273	-0.1738318	-0.4728846

```
attr(,"method")  
[1] "excess"
```

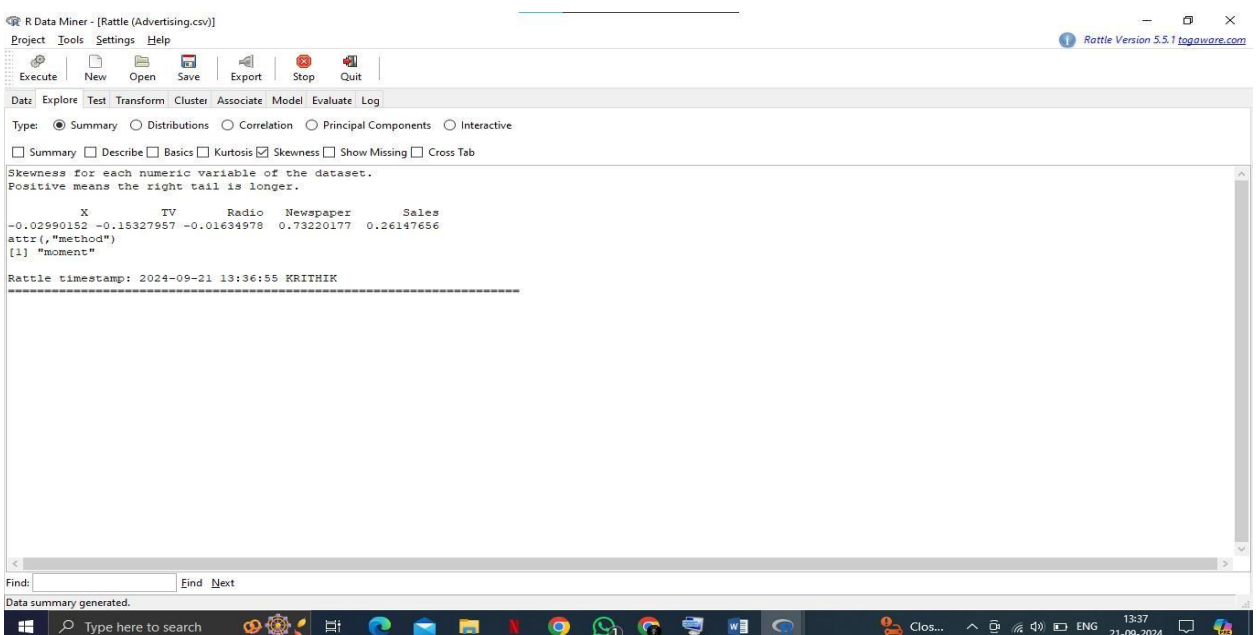
Rattle timestamp: 2024-09-21 13:04:23 KRITHIK

Find: Find Next

Data summary generated.

(E) SKEWNESS

The skewness is a measure of how asymmetrical our data is distributed. A positive skew indicates that the



R Data Miner - [Rattle (Advertising.csv)]

Project Tools Settings Help

Execute New Open Save Export Stop Quit

Data Explore Test Transform Cluster Associate Model Evaluate Log

Type: ☒ Summary ☐ Distributions ☐ Correlation ☐ Principal Components ☐ Interactive

☐ Summary ☐ Describe ☐ Basics ☐ Kurtosis ☒ Skewness ☐ Show Missing ☐ Cross Tab

Skewness for each numeric variable of the dataset.
Positive means the right tail is longer.

X	TV	Radio	Newspaper	Sales
-0.02990152	-0.15327957	-0.01634978	0.73220177	0.26147656

```
attr(,"method")  
[1] "moment"
```

Rattle timestamp: 2024-09-21 13:36:58 KRITHIK

Find: Find Next

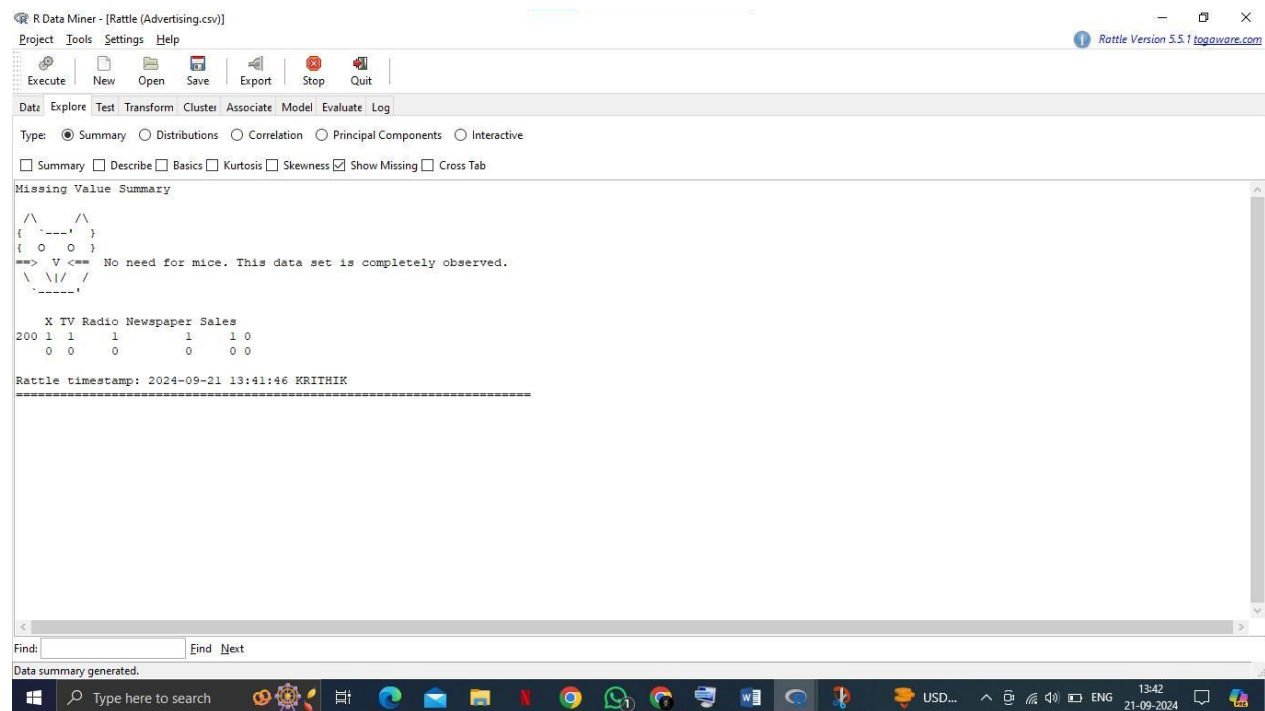
Data summary generated.

tail to the right is longer, and a negative skew that the tail to the left is longer.

(F) MISSING

Missing values present challenges to data mining. The Show Missing check button of the Summary option of the Explore tab provides a summary of missing values in our dataset. Following figure illustrates the missing value summary. Such information is useful in understanding structure in the missing values.

The missing value summary table is presented with the variables listed along the top. Each row corresponds to a pattern of missing values. A 1 indicates a value is present, whereas a 0 indicates a value is missing.



RESULT :

Classification introduction to exploratory data analysis using rattle an open source tool(r) is executed and results verified.

EXP NO. 3

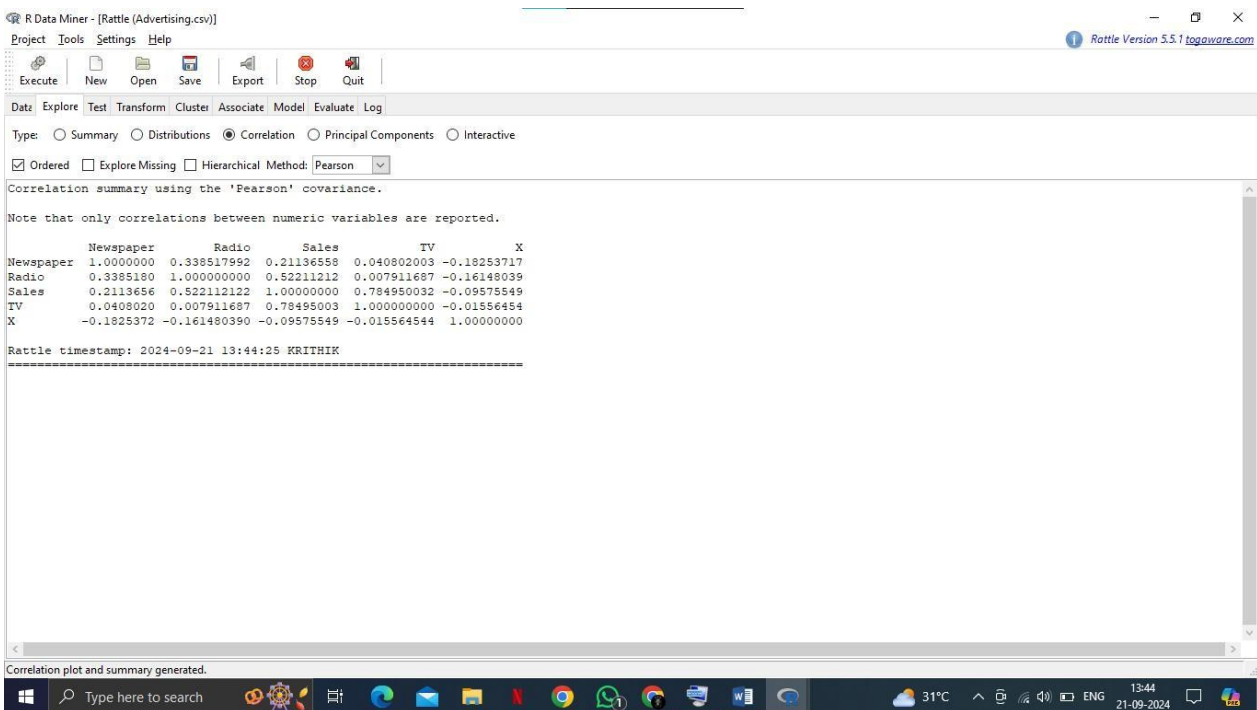
DATE :

INTRODUCTION TO REGRESSION USING RATTLE AN OPEN SOURCE TOOL. (R)

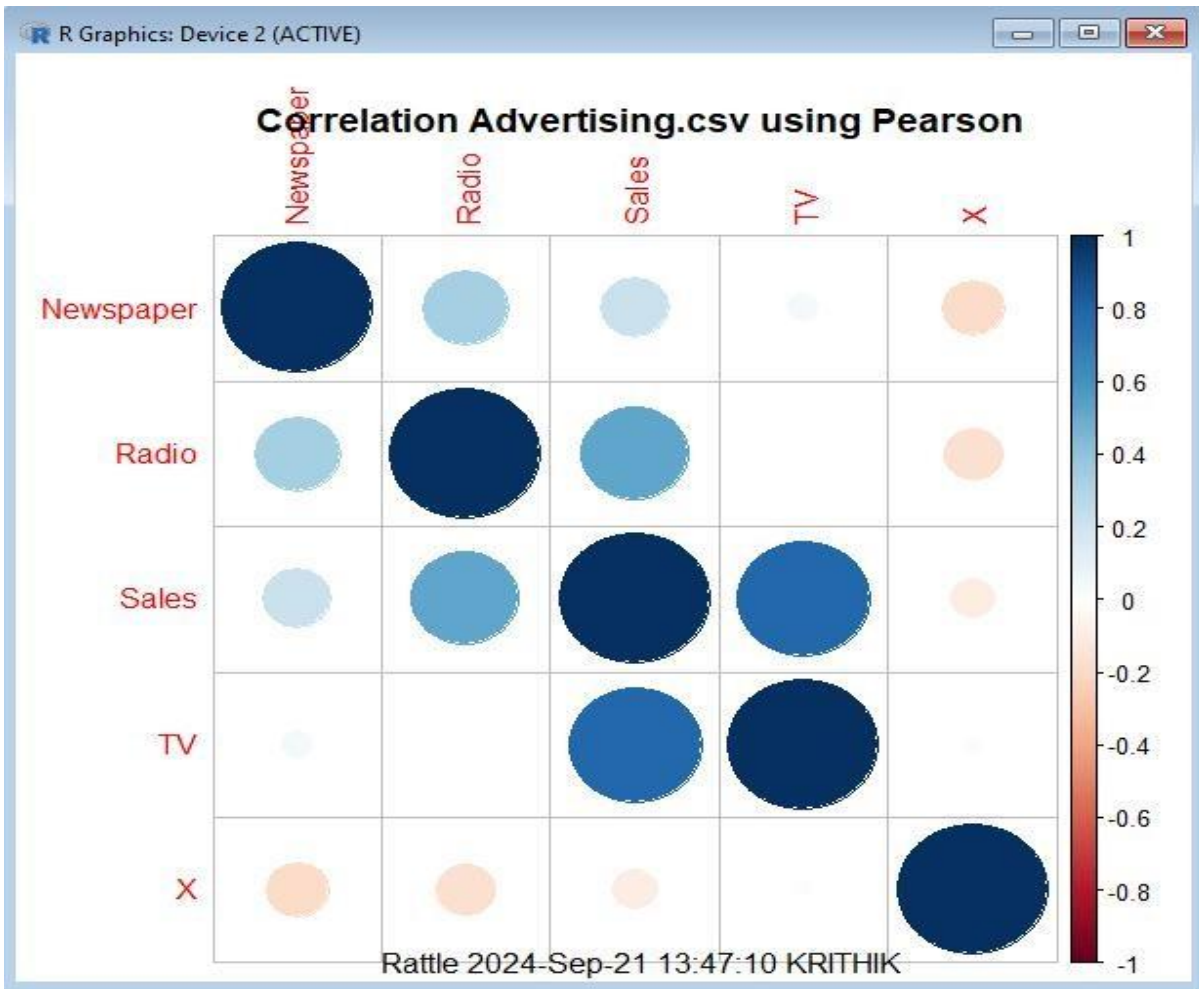
AIM :

To perform Correlation and linear regression analysis that are the most commonly used techniques for investigating the relationship between two quantitative variables.

CORRELATION ANALYSIS



A correlation plot will display correlations between the values of variables in the dataset. In addition to the usual correlation calculated between values of different variables, the correlation between missing values can be explored by checking the Explore Missing check box. The first thing to notice for this correlation plot is that only the numeric variables appear. Rattle only computes correlations between numeric variables at this time. The second thing to note about the graphic is that it is symmetric about the diagonal. The correlation between two variables is the same, irrespective of the order in which we view the two variables. The third thing to note is that the order of the variables does not correspond to the order in the dataset, but to the order of the strength of any correlations, from the least to the greatest. This is done simply to achieve a more pleasing graphic which is easier to take in.

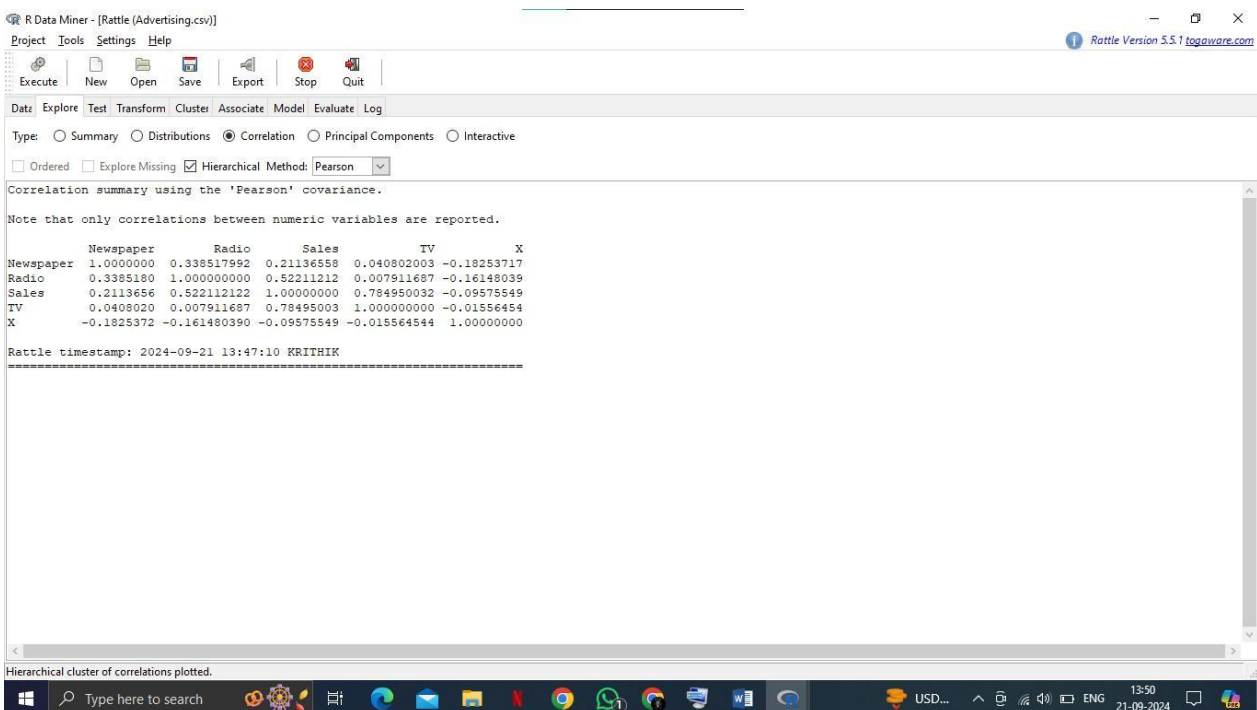


We interpret the degree of any correlation by both the shape and colour of the graphic elements. Any variable is, of course, perfectly correlated with itself, and this is reflected as the diagonal lies on the diagonal of the graphic. Where the graphic element is a perfect circle, then there is no correlation between the variables, as is the case in the correlation between Hours and Deductions although in fact there is a correlation, just a very weak one. The colours used to shade the circles give another (if perhaps redundant) clue to the strength of the correlation. The intensity of the colour is maximal for a perfect correlation, and minimal (white) if there is no correlation. Shades of red are used for negative correlations and blue for positive correlations.

By selecting the Explore Missing check box you can obtain a correlation plot that will show any correlations between the missing values of variables. This is particularly useful to understand how missing values in one variable are related to missing values in another. We notice immediately that only three variables are included in this correlation plot. Rattle has identified that the other variables in fact have no missing values, and so there is no point including them in the plot. We also notice that a categorical

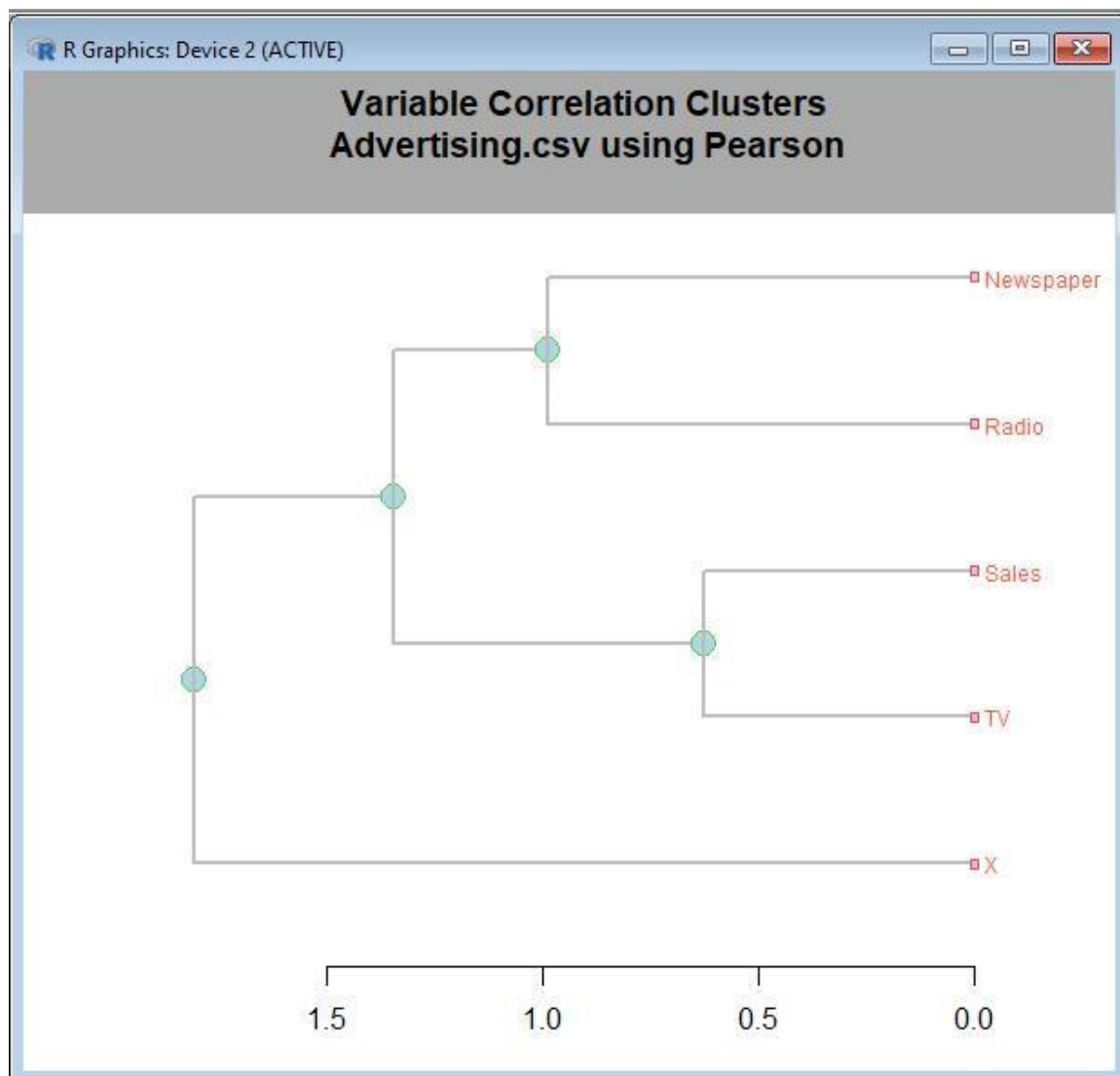
variable, Accounts, is included in the plot even though it was not included in the usual correlation plot. In this case we can obtain a correlation for categorical variables since we only measure missing and presence of a value, which is easily interpreted as numeric.

The graphic shows us that Employment and Occupation are highly correlated in their presence of missing values. That is, when Employment has a missing value, so does Occupation, and vice versa, at least in general. The actual correlation is 0.995 (which can be read from the Rattle text view window), which is very close to 1. On the other hand, there is no (in fact very little at 0.013) correlation between Accounts and the other two variables, with regard missing values. It is important to note that the correlations showing missing values may be based on very small samples, and this information is included in the text view of the Rattle window. For example, in this example we can see that there are only 100, 101, and 43 missing values, respectively, for each of the three variables having any missing values. This corresponds to approximately 5%, 5%, and 2% of the entities, respectively, having missing values for these variables.



Rattle uses the default R correlation calculation known as Pearson's correlation, a common measure of correlation.

HIERARCHICAL CORRELATION



RESULT :

Introduction to regression using rattle an open source tool is executed and results are verified

EXP NO. 4

DATE :

USING THE WEKA TOOLKIT – PART 1

AIM :

Demonstration of classification rule process on dataset student.arff using j48 algorithm

This experiment illustrates the use of j-48 classifier in weka. The sample data set used in this experiment is “student” data available at arff format. This document assumes that appropriate data pre processing has been performed.

STEPS INVOLVED IN THIS EXPERIMENT:

Step-1: We begin the experiment by loading the data (student.arff) into weka.

Step2: Next we select the “classify” tab and click “choose” button to select the “j48” classifier.

Step3: Now we specify the various parameters. These can be specified by clicking in the text box to the right of the choose button. In this example, we accept the default values. The default version does perform some pruning but does not perform error pruning.

Step4: Under the “text” options in the main panel. We select the 10-fold cross validation as our evaluation approach. Since we don’t have separate evaluation data set, this is necessary to get a reasonable idea of accuracy of generated model.

Step-5: We now click “start” to generate the model. The ASCII version of the tree as well as evaluation statistic will appear in the right panel when the model construction is complete.

Step-6: Note that the classification accuracy of model is about 69%. This indicates that we may find more work. (Either in preprocessing or in selecting current parameters for the classification)

Step-7: Now weka also lets us view a graphical version of the classification tree. This can be done by right clicking the last result set and selecting “visualize tree” from the pop-up menu.

Step-8: We will use our model to classify the new instances.

Step-9: In the main panel under “text” options click the “supplied test set” radio button and then click the “set” button. This will pop-up a window which will allow you to open the file containing test instances.

Dataset student .arff

@relation student

@attribute age {<30,30-40,>40} @attribute income {low, medium, high} @attribute student {yes, no}

@attribute credit-rating {fair, excellent}

@attribute buyspc {yes, no}

@data

%

<30, high, no, fair, no

<30, high, no, excellent, no

30-40, high, no, fair, yes

>40, medium, no, fair, yes

>40, low, yes, fair, yes

>40, low, yes, excellent, no

30-40, low, yes, excellent, yes

<30, medium, no, fair, no

<30, low, yes, fair, no

>40, medium, yes, fair, yes

<30, medium, yes, excellent, yes

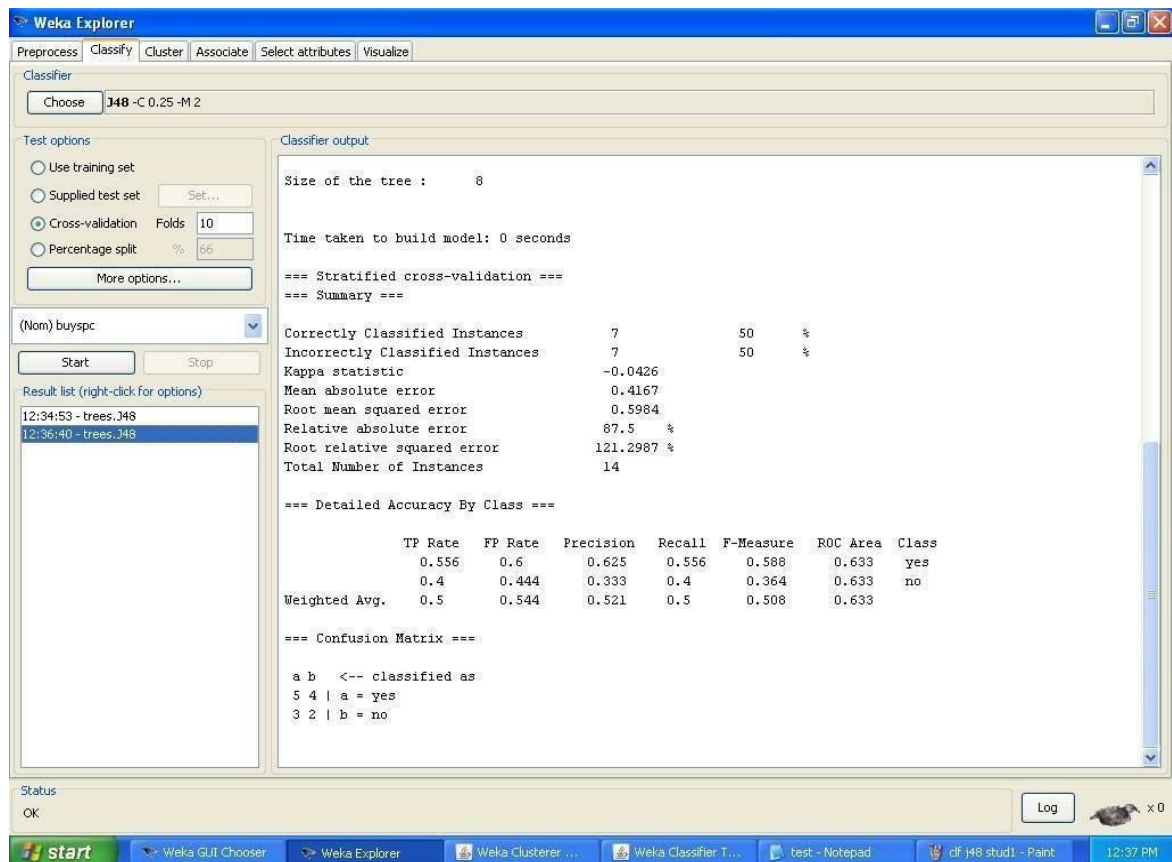
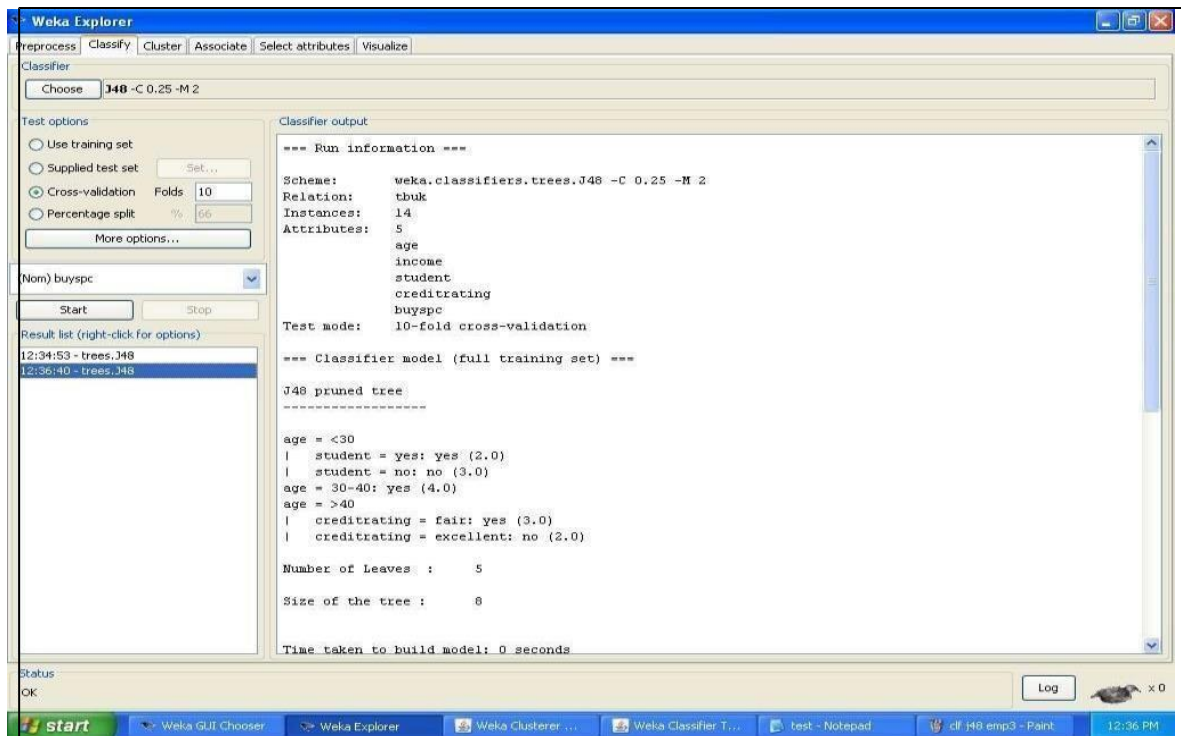
30-40, medium, no, excellent, yes

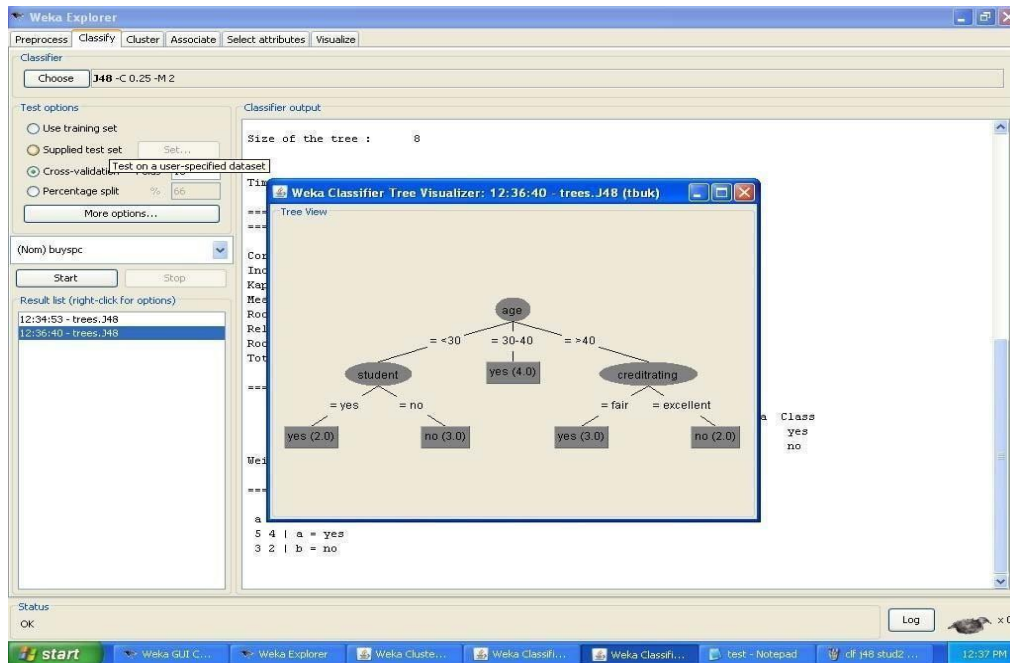
30-40, high, yes, fair, yes

>40, medium, no, excellent, no

%

The following screenshot shows the classification rules that were generated when j48 algorithm is applied on the given dataset.





RESULT :

Demonstration of classification rule process on dataset student.arff using j48 algorithm is executed and results are verified

EXP NO. 5

DATE :

CLASSIFICATION USING THE WEKA TOOLKIT – PART 2

AIM :

Demonstration of classification algorithm NaiveBayes using Iris.arff dataset in Weka Tool.

STEPS INVOLVED IN THIS EXPERIMENT:

LOAD DATASET:

- Open Weka.
- Go to the "Explorer" section and click on "Open File."
- Load the Iris.arff dataset.

SELECT CLASSIFIER:

- Navigate to the "Classify" tab.
- In the "Classifier" section, choose the **NaiveBayes** classifier from the tree.

SET EVALUATION:

- Under the "Test options," select **10-fold cross-validation** to evaluate the model.

GENERATE MODEL:

- Click the "Start" button to create the model and view the results in the output window.

CHECK ACCURACY:

- In the results, note that the classification accuracy should be around **69%**.

VISUALIZE MARGIN CURVE:

- After the model finishes, right-click the result in the "Result list" and choose **Visualize margin curve** to see the plot.

EVALUATION METRICS:

- Review the precision, recall, ROC Area, and Confusion Matrix, which will be displayed in the output results.

DATA SET IRIS.ARFF:

Attributes list:

1. Sepal Length
2. Sepal Width
3. Petal Length
4. Petal width
5. Class

MODEL DATASET

S.NO	Sepal Length	Sepal Width	Petal Length	Petal width	Class
1	5.1	3.5	1.4	0.2	Iris -setosa
2	4.9	3.0	1.4	0.2	Iris -setosa
3	4.7	3.2	1.3	0.2	Iris -setosa
4	5.0	3.4	1.5	0.2	<i>Iris -setosa</i>
5	4.4	2.9	1.4	0.2	Iris -setosa
6	7.0	3.2	4.7	1.4	Iris-versicolor
7	5.9	3.0	4.2	1.5	Iris-versicolor
8	5.0	2.0	3.5	1.0	Iris-versicolor
9	6.1	2.9	4.7	1.4	Iris-versicolor
10	5.5	2.9	3.6	1.3	Iris-versicolor
11	6.5	3.0	5.8	2.2	Iris - virginica
12	6.3	3.4	5.6	2.4	Iris - virginica
13	5.7	3.1	5.1	1.9	Iris - virginica
14	5.9	3.0	5.1	5.4	Iris - virginica
15	6.2	3.4	5.4	2.3	Iris - virginica

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 %: 69

More options...

(Nom) class: **Start** **Stop**

Result list (right-click for options):

- 10:45:10 - misc.InputMappedClassifier
- 10:45:38 - bayes.NaiveBayes**

Classifier output:

```

=== Run information ===

Scheme:      weka.classifiers.bayes.NaiveBayes
Relation:    iris
Instances:   150
Attributes:  5
              sepalength
              sepalwidth
              petallength
              petalwidth
              class
Test mode:   split 69.0% train, remainder test

=== Classifier model (full training set) ===

Naive Bayes Classifier

Attribute      Class
              Iris-setosa Iris-versicolor Iris-virginica
              (0.33)      (0.33)      (0.33)
=====
sepalength
mean           4.9913      5.9379      6.5795
std. dev.      0.355      0.5042      0.6353
weight sum     50         50         50
precision      0.1059      0.1059      0.1059

sepalwidth
mean           3.4015      2.7687      2.9629
std. dev.      0.3925      0.3038      0.3088
weight sum     50         50         50
precision      0.1091      0.1091      0.1091

```

Status: OK **Log** x 0

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 %: 66

More options...

(Nom) class: **Start** **Stop**

Result list (right-click for options):

- 10:28:52 - trees.RandomTree
- 10:29:00 - trees.DecisionStump
- 10:29:54 - bayes.NaiveBayes
- 10:30:05 - bayes.NaiveBayesMultinomial
- 10:31:16 - trees.RandomTree
- 10:31:32 - bayes.NaiveBayes
- 10:32:20 - rules.OneR
- 10:32:30 - trees.DecisionStump
- 10:34:03 - bayes.NaiveBayes**

Classifier output:

```

petallength
mean           1.4694      4.2452      5.5516
std. dev.      0.1782      0.4712      0.5529
weight sum     50         50         50
precision      0.1405      0.1405      0.1405

petalwidth
mean           0.2743      1.3097      2.0343
std. dev.      0.1096      0.1915      0.2646
weight sum     50         50         50
precision      0.1143      0.1143      0.1143

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

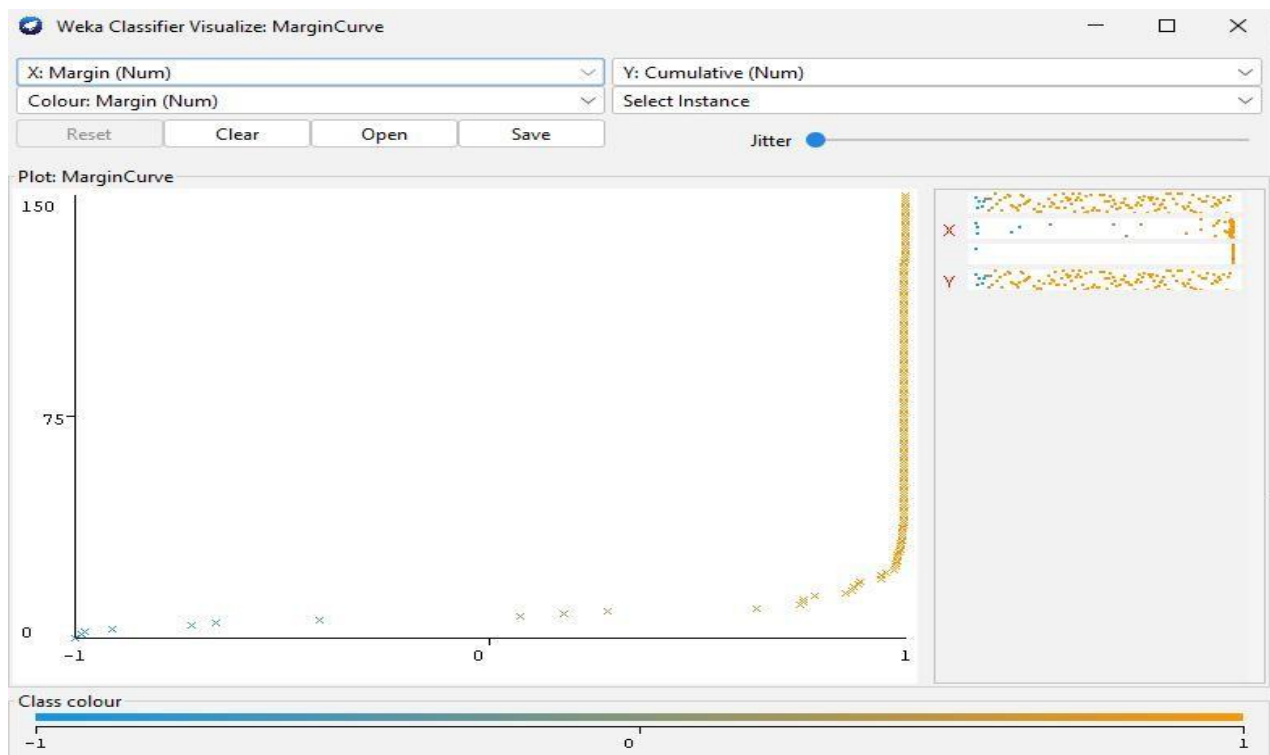
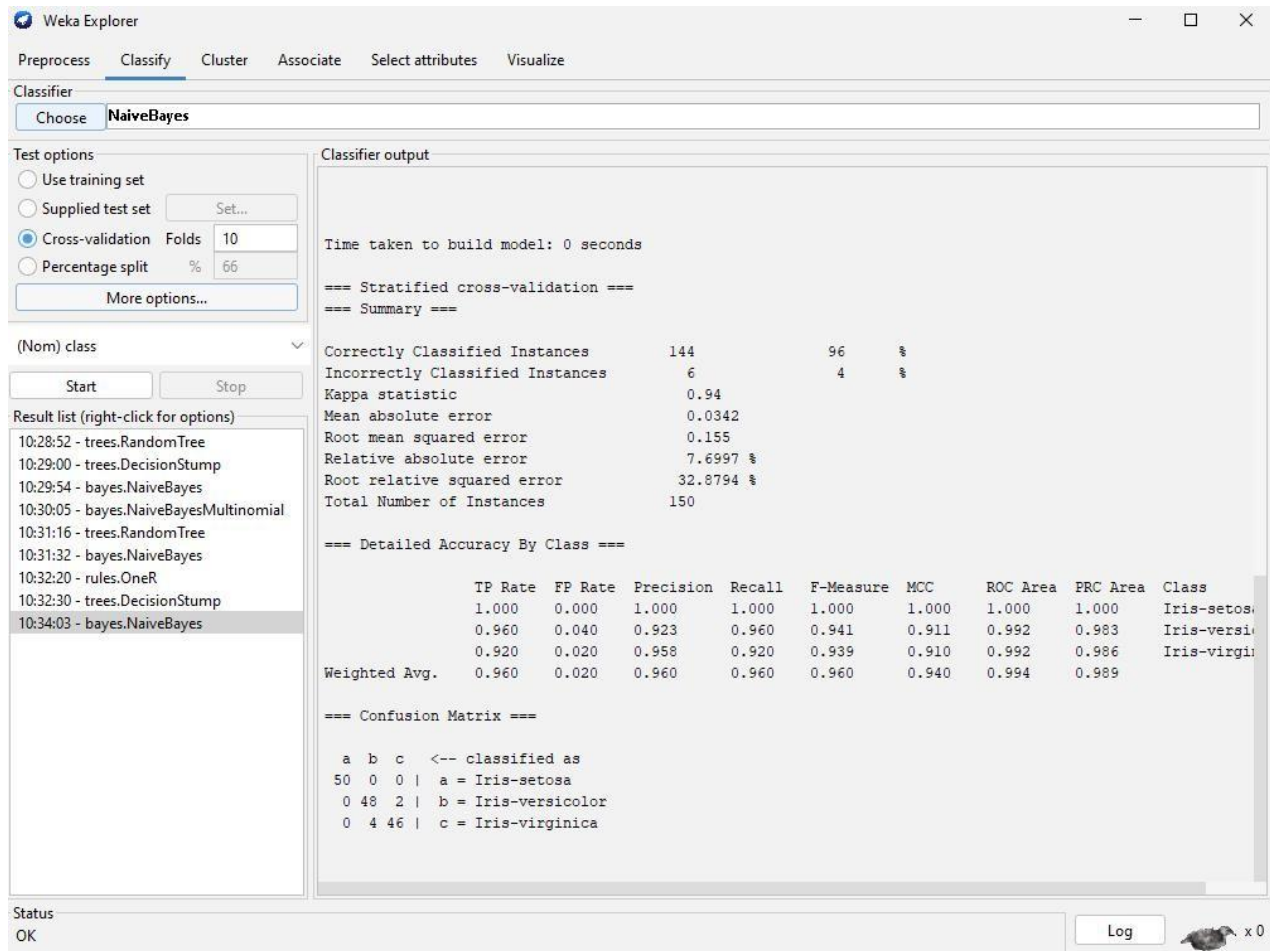
Correctly Classified Instances      144      96      %
Incorrectly Classified Instances     6       4      %
Kappa statistic                     0.94
Mean absolute error                  0.0342
Root mean squared error              0.155
Relative absolute error              7.6997 %
Root relative squared error          32.8794 %
Total Number of Instances           150

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall  F-Measure  MCC   ROC Area  PRC Area  Class
      1.000    0.000    1.000     1.000    1.000     1.000  1.000    1.000    Iris-setosa

```

Status: OK **Log** x 0



RESULT :

Demonstration of classification rule process on dataset student.arff using id3 algorithm is executed and results are verified

PERFORMING DATA PREPROCESSING FOR DATA MINING IN WEKA**AIM :**

Demonstration of preprocessing on dataset student.arff

This experiment illustrates some of the basic data preprocessing operations that can be performed using WEKA-Explorer. The sample dataset used for this example is the student data available in arff format.

STEPS INVOLVED IN THIS EXPERIMENT:

Step1: Loading the data. We can load the dataset into weka by clicking on open button in preprocessing interface and selecting the appropriate file.

Step2: Once the data is loaded, weka will recognize the attributes and during the scan of the data weka will compute some basic strategies on each attribute. The left panel in the above figure shows the list of recognized attributes while the top panel indicates the names of the base relation or table and the current working relation (which are same initially).

Step3: Clicking on an attribute in the left panel will show the basic statistics on the attributes for the categorical attributes the frequency of each attribute value is shown, while for continuous attributes we can obtain min, max, mean, standard deviation and deviation etc.,

Step4: The visualization in the right button panel in the form of cross-tabulation across two attributes.

Note: we can select another attribute using the dropdown list.

Step5: Selecting or filtering attributes

Removing an attribute-When we need to remove an attribute, we can do this by using the attribute filters in weka. In the filter model panel, click on choose button. This will show a popup window with a list of available filters.

Scroll down the list and select the “weka.filters.unsupervised.attribute.remove” filters.

Step 6: a) Next click the textbox immediately to the right of the choose button. In the resulting dialog box enter the index of the attribute to be filtered out.

Make sure that invert selection option is set to false. Click OK now in the filter box. You will see “Remove-R-7”.

Click the apply button to apply filter to this data. This will remove the attribute and create new working relation.

Save the new working relation as an arff file by clicking save button on the top(button)panel.(student.arff)

Discretization

1. Sometimes association rule mining can only be performed on categorical data. This requires performing discretization on numeric or continuous attributes.

In the following example let us discretize age attribute.

Let us divide the values of age attribute into three bins(intervals).

First load the dataset into weka(student.arff)

Select the age attribute.

Activate filter-dialog box and select “WEKA.filters.unsupervised.attribute.discretize” from the list.

To change the defaults for the filters, click on the box immediately to the right of the choose button.

We enter the index for the attribute to be discretized. In this case the attribute is age. So we must enter ‘1’ corresponding to the age attribute.

Enter ‘3’ as the number of bins. Leave the remaining field values as they are.

Click OK button.

Click apply in the filter panel. This will result in a new working relation with the selected attribute partition into 3 bins.

Save the new working relation in a file called student-data-discretized.arff

DATASET STUDENT .ARFF

@relation student

@attribute age {<30,30-40,>40} @attribute income {low, medium, high} @attribute student {yes, no}

@attribute credit-rating {fair, excellent}

@attribute buyspc {yes, no}

@data %

<30, high, no, fair, no

<30, high, no, excellent, no

30-40, high, no, fair, yes

>40, medium, no, fair, yes

>40, low, yes, fair, yes

>40, low, yes, excellent, no

30-40, low, yes, excellent, yes

<30, medium, no, fair, no

<30, low, yes, fair, no

>40, medium, yes, fair, yes

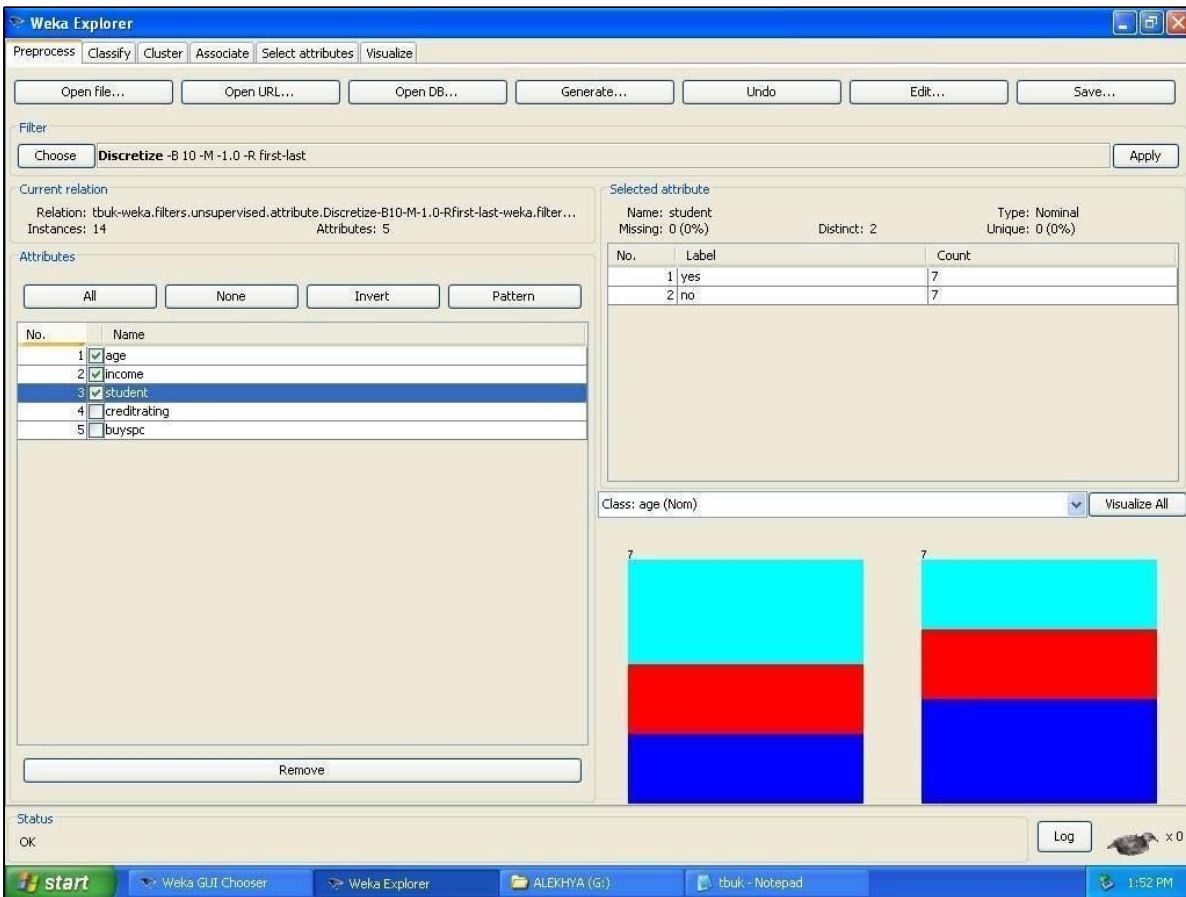
<30, medium, yes, excellent, yes

30-40, medium, no, excellent, yes

30-40, high, yes, fair, yes

>40, medium, no, excellent, no %

The following screenshot shows the effect of discretization.



RESULT :

Performing data preprocessing for data mining in weka is executed and results are verified

PERFORMING CLUSTERING IN WEKA**AIM :**

Demonstration of clustering rule process on dataset student.arff using simple k-means

This experiment illustrates the use of simple k-mean clustering with Weka explorer. The sample data set used for this example is based on the student data available in ARFF format. This document assumes that appropriate pre processing has been performed. This is student dataset includes 14 instances.

STEPS INVOLVED IN THIS EXPERIMENT :

Step 1: Run the Weka explorer and load the data file student.arff in preprocessing interface.

Step 2: In order to perform clustering select the 'cluster' tab in the explorer and click on the choose button. This step results in a dropdown list of available clustering algorithms.

Step 3 : In this case we select 'simple k-means'.

Step 4: Next click in text button to the right of the choose button to get popup window shown in the screenshots. In this window we enter six on the number of clusters and we leave the value of the seed on as it is. The seed value is used in generating a random number which is used for making the internal assignments of instances of clusters.

Step 5 : Once of the option have been specified. We run the clustering algorithm there we must make sure that they are in the 'cluster mode' panel. The use of training set option is selected and then we click 'start' button. This process and resulting window are shown in the following screenshots.

Step 6 : The result window shows the centroid of each cluster as well as statistics on the number and the percent of instances assigned to different clusters. Here clusters centroid are means vectors for each clusters. This clusters can be used to characterized the cluster.

Step 7: Another way of understanding characteristics of each cluster through visualization ,we can do this, try right clicking the result set on the result. List panel and selecting the visualize cluster assignments.

Interpretation of the above visualization

From the above visualization, we can understand the distribution of age and instance number in each cluster. For instance, for each cluster is dominated by age. In this case by changing the color dimension to other attributes we can see their distribution with in each of the cluster.

Step 8: We can assure that resulting dataset which included each instance along with its assign cluster. To do so we click the save button in the visualization window and save the result student k- mean .The top portion of this file is shown in the following figure.

DATASET STUDENT .ARFF

```
@relation student
@attribute age {<30,30-40,>40}
@attribute income {low,medium,high}
@attribute student {yes,no}
@attribute credit-rating {fair,excellent}
@attribute buyspc {yes,no}
@data %
<30, high, no, fair, no
<30, high, no, excellent, no
30-40, high, no, fair, yes
>40, medium, no, fair, yes
>40, low, yes, fair, yes
>40, low, yes, excellent, no
30-40, low, yes, excellent, yes
<30, medium, no, fair, no
<30, low, yes, fair, no
>40, medium, yes, fair, yes
<30, medium, yes, excellent, yes
30-40, medium, no, excellent, yes
30-40, high, yes, fair, yes
>40, medium, no, excellent, no %
```

The following screenshot shows the clustering rules that were generated when simple k-means algorithm is applied on the given dataset.

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer
Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode
☒ Use training set
☐ Supplied test set Set...
☐ Percentage split % 66
☐ Classes to clusters evaluation (Nom) Class
☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)
14:04:32 - SimpleKMeans

Status OK

Clusterer output

```

=== Run information ===

Scheme:      weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10
Relation:    breast-cancer
Instances:   286
Attributes:  10
age
menopause
tumor-size
inv-nodes
node-caps
deg-malig
breast
breast-quad
irradiat
Class

Test mode:   evaluate on training data

=== Clustering model (full training set) ===

KMeans
=====

Number of iterations: 3
Within cluster sum of squared errors: 1177.0

Initial starting points (random):
Cluster 0: 50-59,premeno,10-14,0-2,no,2,right,left_up,no,no-recurrence-events
Cluster 1: 40-49,premeno,15-19,0-2,yes,3,right,left_up,no,no-recurrence-events
  
```

Log x 0

Weka Explorer

Preprocess Classify **Cluster** Associate Select attributes Visualize

Clusterer
Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

Cluster mode
☒ Use training set
☐ Supplied test set Set...
☐ Percentage split % 66
☐ Classes to clusters evaluation (Nom) Class
☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)
14:04:32 - SimpleKMeans

Status OK

Clusterer output

```

Cluster 1: 40-49,premeno,15-19,0-2,yes,3,right,left_up,no,recurrence-events

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute      Full Data      Cluster#      0      1
                (286.0)      (225.0)      (61.0)
=====
age            50-59          50-59          40-49
menopause      premeno        premeno        premeno
tumor-size     30-34          25-29          30-34
inv-nodes      0-2            0-2            0-2
node-caps      no             no             yes
deg-malig      2              2              3
breast         left           left           left
breast-quad    left_low       left_low       left_low
irradiat       no             no             no
Class          no-recurrence-events no-recurrence-events recurrence-events

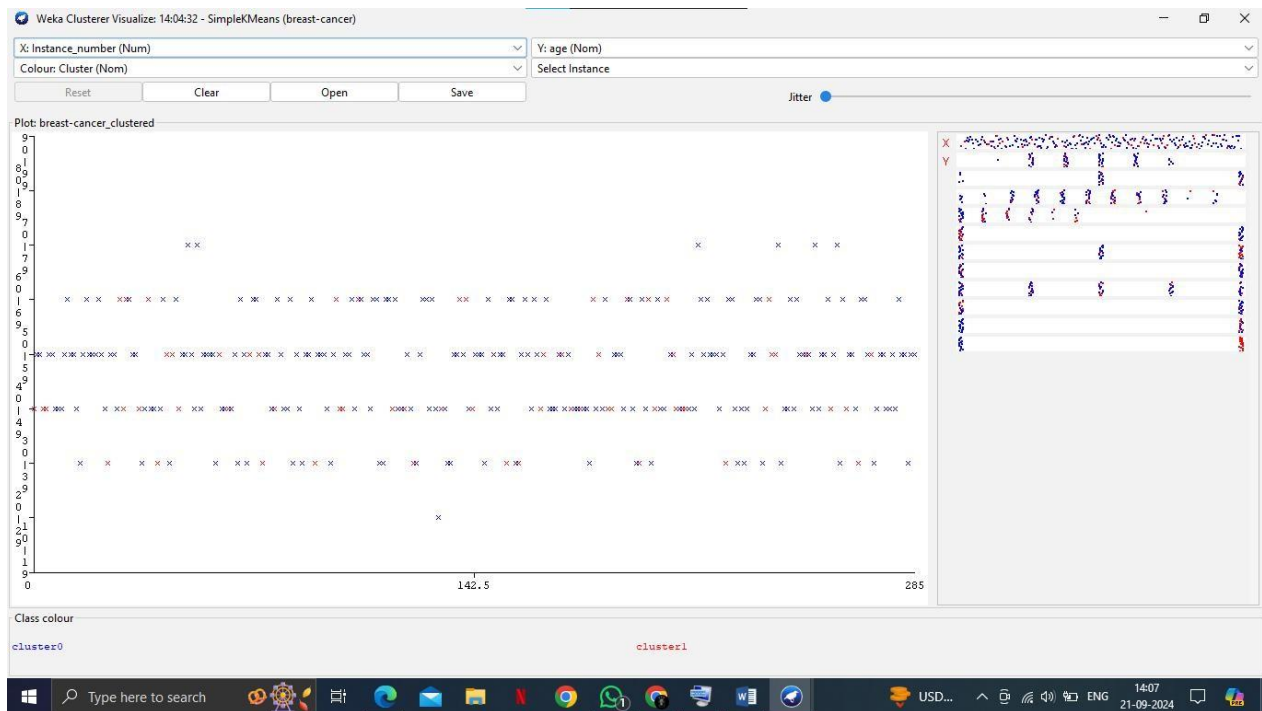
Time taken to build model (full training data) : 0.04 seconds

=== Model and evaluation on training set ===

Clustered Instances

0      225 ( 79%)
1       61 ( 21%)
  
```

Log x 0



RESULT :

Demonstration of clustering rule process on dataset student.arff using simple k-means is executed and results are verified

EXP NO. 8

DATE :

ASSOCIATION RULE ANALYSIS IN WEKA

AIM :

Demonstration of Association rule process on dataset contactlenses.arff using apriori algorithm

This experiment illustrates some of the basic elements of association rule mining using WEKA. The sample dataset used for this example is contactlenses.arff

STEPS INVOLVED IN THIS EXPERIMENT :

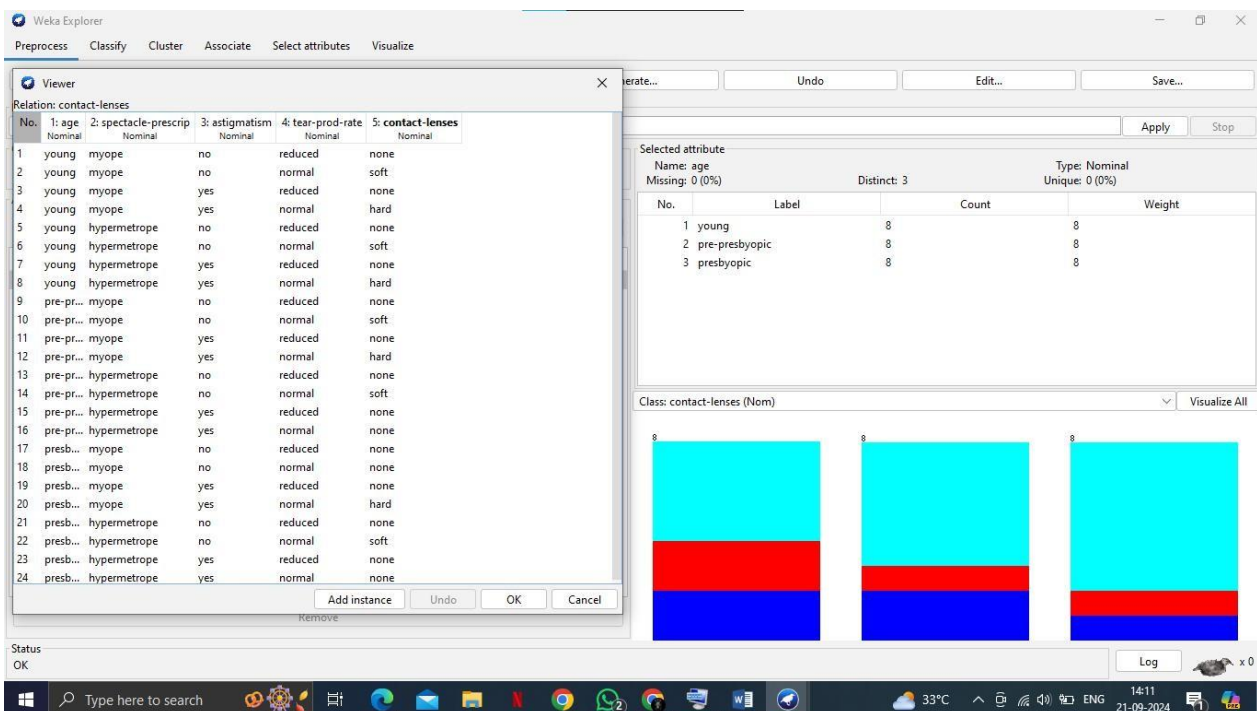
Step1: Open the data file in Weka Explorer. It is presumed that the required data fields have been discretized. In this example it is age attribute.

Step2: Clicking on the associate tab will bring up the interface for association rule algorithm.

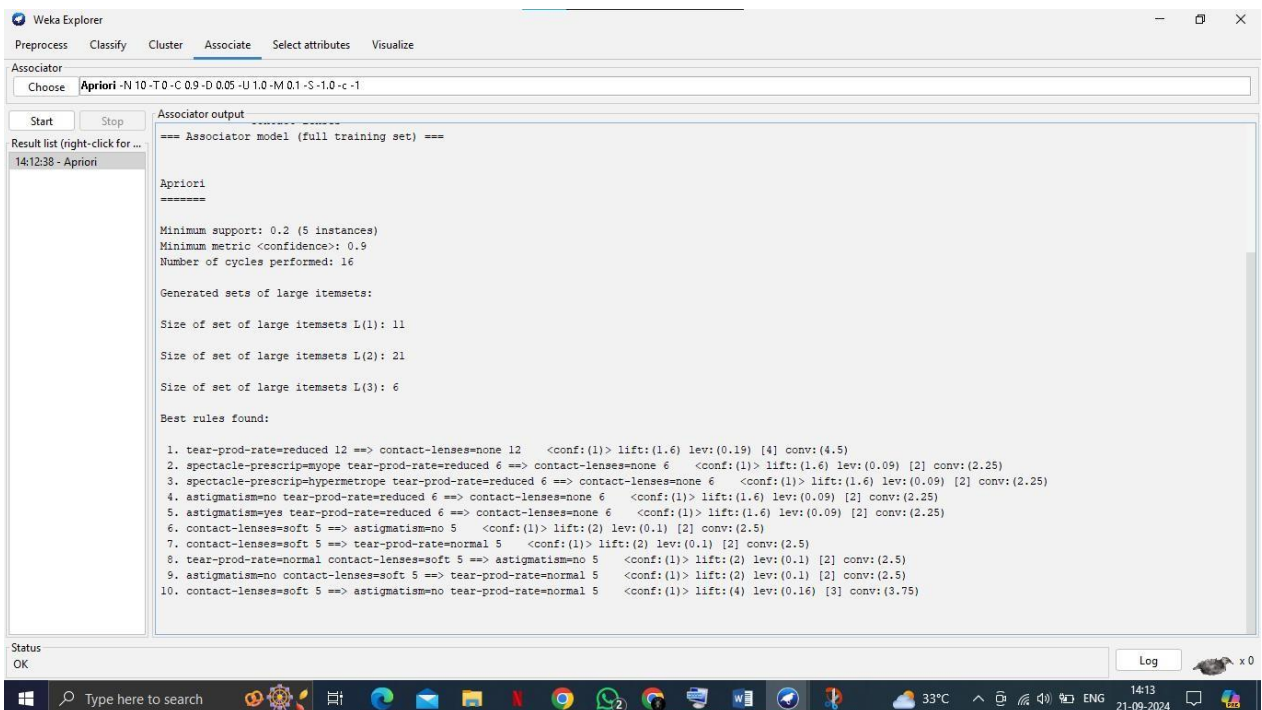
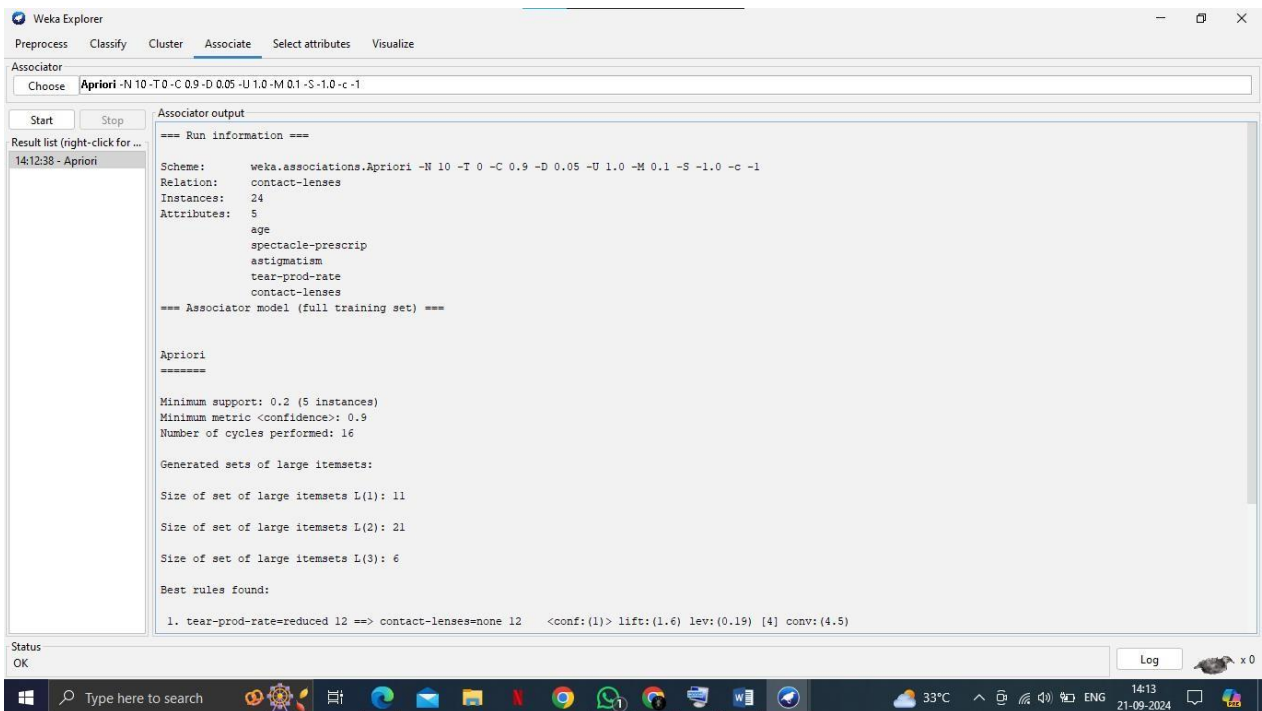
Step3: We will use apriori algorithm. This is the default algorithm.

Step4: Inorder to change the parameters for the run (example support, confidence etc) we click on the text box immediately to the right of the choose button.

Dataset contactlenses.arff



The following screenshot shows the association rules that were generated when apriori algorithm is applied on the given dataset.



RESULT :

Demonstration of Association rule process on dataset contactlenses.arff using apriori algorithm is executed and results are verified

DATA MINING CASE STUDY USING THE CRISP-DM STANDARD**AIM :**

CRISP- DM CASE STUDY – Early Prediction Of Student Success : Mining Students Enrolment Data

This case study explores the socio-demographic variables (age, gender, ethnicity, education, work status, and disability) and study environment (course programme and course block), that may influence persistence or dropout of students at the Open Polytechnic of New Zealand. This case study examines to what extent these factors, i.e. enrolment data help us in pre-identifying successful and unsuccessful students. The data stored in the Open Polytechnic student management system covers over 450 students who enrolled to *71150 Information Systems* course was used to perform a quantitative analysis of study outcome. Based on a data mining techniques (such as feature selection and classification trees), the most important factors for student success and a profile of the typical successful and unsuccessful students are identified. The empirical results show the following: (i) the most important factors separating successful from unsuccessful students are: ethnicity, course programme and course block; (ii) among classification tree growing methods Classification and Regression Tree (CART) was the most successful in growing the tree with an overall percentage of correct classification of 60.5%; and (iii) both the risk estimated by the cross-validation and the gain diagram suggests that all trees, based only on enrolment data are not quite good in separating successful from unsuccessful students. The implications of these results for academic and administrative staff are discussed.

INTRODUCTION

Increasing student retention or persistence is a long term goal in all academic institutions. The consequences of student attrition are significant for students, academic and administrative staff. The importance of this issue for students is obvious: school leavers are more likely to earn less than those who graduated. Since one of the criteria for government funding in the tertiary education environment in New Zealand is the level of retention rate, both academic and administrative staff are under pressure to come up with strategies that could increase retention rates on their courses and programmes. The most vulnerable students to low student retention at all institutions of higher education are the first-year students, who are at greatest risk of dropping out in the first term or semester of study or not completing their programme/ degree on time.

Therefore most retention studies address the retention of first-year students. Consequently, the early identification of vulnerable students who are prone to drop their courses is crucial for the success of any retention strategy. This would allow educational institutions to undertake timely and pro-active measures.

Once identified, these ‘at-risk’ students can be then targeted with academic and administrative support to increase their chance of staying on the course.

The main objective of this study is to explore factors that may impact the study outcome in the *Information Systems* course at the Open Polytechnic. The *Information Systems* course is a core course for those majoring in IT and for most students an entry point, i.e. the first course they are taking with the Open Polytechnic. This issue have not been examined so far for Open Polytechnic and this case study attempts to fill the gap.

More specifically the enrolment data were used to achieve the following objectives:

Build models for early prediction of study outcome using the student enrolment data • Evaluate the models using cross-validation and misclassification errors to decide which model outperforms other models in term of classification accuracy Present results which can be easily understood by the users (students, academic and administrative staff)

At the time of enrolment in the Open Polytechnic of New Zealand, the only information. i.e. variables we have about students are those contained in their enrolment forms. The question we are trying to address in this case study is whether we can use the enrolment data alone to predict study outcome for newly enrolled student.

Framework for Data Mining Process

Framework for data mining applications is based on the CRISP-DM Model created by a consortium of NCR, SPSS, and Daimler-Benz companies. The modified version of the CRISP- DM model is presented on Figure 1, following the project through the general life cycle from business and data understanding, data preparation, modeling, evaluation and deployment. The feedback from deployment to data and business understanding illustrates the iterative nature of a data mining process.

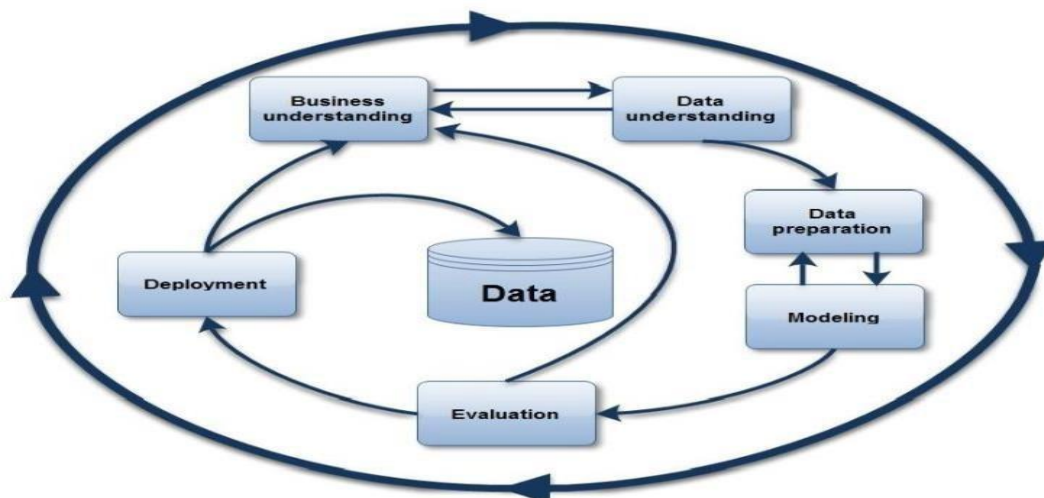


Figure 1: Modified CRISP-DM Model Version 1

Business Understanding

The business understanding phase begins with the setting up of goals for the data mining project. In this paper that would be an increasing understanding of the pre-enrolment factors that may prevent students from successfully complete the course.

Because we are planning to increase completion rate on the *Information Systems* course understanding its students population and patterns in the pre-enrolment data becomes necessary before we start developing a predictive model. In this phase we are coming up with the following questions: what is the profile of a student who successfully completes this course? Can the successful vs. unsuccessful student be distinguished in terms of demographic features (such as gender, age or ethnic origin) or study environment (such as course programme, faculty or course block)? Depending on the answers to these questions we consider the methods and approaches that can be adopted to increase the completion rate.

Data Understanding

The scope of our research in terms of data used is limited by the data available in the Open Polytechnic Student Management System (known as Integrator) and the enrolment form used for collecting data from newly enrolled students. It is important to have a full understanding of the nature of the data and how it was collected and entered before proceeding further. In this phase an initial data exploration using a pivot table was also conducted to get some insight in the data.

The enrolment form asks students to enter the following information: demographic (gender, date of birth, ethnic origins, disability, and work status), academic (secondary school qualification, course programme, course faculty, and course block), and contact details. Once the data from the enrolment form is entered into Integrator and the Enrolment Section process the application, the enrolment date is recorded and student becomes enrolled on the course.

Data Preparation

In this phase the data are put into a form suitable for the modeling phase. If required some selected variables are combined, transformed or used to create a new variables. For example, enrolment date and the course block start date were used to generate a variable labeled as “early

enrolment”. Any data excluded from the data set is documented and their removal explained. For example, only a few students enrolled to Bachelor of Arts programme were on the *Information Systems* course. They were removed from further analysis for the reason explained later in Data and

Methodology section. Data are cleaned for any duplication of records. For example, in case of *Information Systems* course, the course code changed in the past. If student enrolled in the time when the

change in the course code happened and then re-enrolled on the same course, two records exist in the data set for the same student and the same course, but under two different course codes. In this case data for this student were merged into one, single record. The dependent variable “study outcome” with two possible outcomes (labeled as *Pass* and *Fail*) indicates whether students successfully completed the course or failed the course due to voluntary transfer/ withdrawal or academic withdrawal or simply because not fulfilling course pass requirements.

Modeling

In this phase we choose and ran models on the training data set. Then we decided whether a suitable model for the data set was found that was acceptable from both analytical and managerial standpoint. In this phase we decided to use classification tree models with four different tree growing criteria.

Evaluation

The final models from the previous phase are then applied on a testing, i.e. a validation data set with the aim of assessing their predictive accuracy and consistency with the results obtained for the training data set. This phase involves an iterative process of fitting different versions of models to training and testing data set, each time evaluating their predictive performance.

Deployment

Once we decided on the final model we can apply it to current data not used during the modeling and evaluation phase. This process is known as *scoring*. The model results are used to address the issues identified in the business understanding phase. The results should be presented in the user friendly format and prepared for use by administrative staff. The final model should retain the

highest predictive accuracy and if it is to be continuously used it should be regularly updated, particularly if some organizational changes occur or if new factors are brought in. For example, if the new information about financial support is added to the student record or enrolment form, then the model should consider the new factor that might be relevant for study outcome.

DATA AND METHODOLOGY

Data preparation

Variables definition and their domains are presented in Table 1.

From the initial dataset all students granted cross-credit or credit were excluded because they didn't actually study this course. The courses they have previously completed were recognized and credited for this course. We also removed 6 students enrolled under Bachelor of Arts program. Such a small number of students in this category do not allow generalization of the results for this particular category. The total number of data was reduced to 453.

We needed to clarify the definition of study outcome that we used in our analysis. We considered only two possible outcomes, labeled as: *Pass* and *Fail*. Students labeled *Pass* successfully completed the course. Students labeled *Fail* transferred or withdrew from the course voluntarily or academic registry withdrew them for not completing the in-course assessments. Those students who stayed on the course until the end of the course but scored less than the course pass mark were also allocated into this category. Because of the data set size (only 453 students) numeric continuous variable such as age was converted into a categorical variable with only three age groups. In Ethnicity we combined Maori and Pacific Island students because of two reasons: they were found to be no different in preliminary bivariate analysis and combined together constitute a small proportion of the data (less than 10%). Combining them into one ethnic group helps with model parsimony. The Secondary school variable combines all students with no secondary school up to NCEA Level 2 on the New Zealand National Qualifications Framework into one group.

Though the software applications skills (spreadsheet and database in particular) are very important on this course we do not take into account the different skill levels present on the *Information Systems* course. Since the students enrolling on this course have different backgrounds and levels of interest in computing, we would expect that skill level has also a significant impact on the study outcome. However, the information about their Office skills level is not available in the moment of enrolment, so not included this factor in our analysis.

Variable	Description (Domain)
Student Demographics	
Gender	Student gender (binary: female or male)
Age	Student's age (numeric: 1 – under 30, 2 – 30 to 40 or 3 – over 40)
Ethnicity	Student's ethnic group (nominal: Pakeha, Maori & Pacific Islanders or Others)
Disability	Student has a disability (binary: yes or no)
Secondary school	Student's highest level of achievement from a secondary school on the New Zealand National Qualifications Framework (nominal: No secondary qualification, NCEA1, NCEA2, University entrance, NCEA3, Overseas or Other)
Work status	Student is working (binary: yes or no)
Early enrolment	Student enrolled for the first time in the course before start of the course (binary: yes or no)

Study Environment	
Course programme	Programme (nominal: Bachelor of Business or Bachelor of Applied Science)
Course block	Semester in which a course is offered (Semester 1, Semester 2 or Semester 3)
Study outcome	Study outcome (nominal: Pass – successful completion, Fail – unsuccessful completion includes also withdrawals, academic withdrawals and transfers)

Table 1: Description of variables and their domains

Methodology

Three types of data mining approaches were conducted in this study. The first approach is descriptive which is concerned with the nature of the dataset such as the frequency table and the relationship between the attributes obtained using cross tabulation analysis (contingency tables). In addition, feature selection is conducted to determine the importance of the prediction variables for modeling study outcome. The third type of data mining approach, i.e. predictive data mining is conducted by using four different classification trees. Finally, a comparison between these classification tree models was conducted to determine the best model for the dataset. We decided to use the classification tree models because of some advantages they may have over traditional statistical models such as logistic regression and discriminant analysis traditionally used in retention studies. First, they can handle a large number of predictor variables, far more than the logistic regression and discriminant analysis would allow. Secondly, the classification tree models are non-parametric and can capture nonlinear relationships and complex interactions between predictors and dependent variable.

Results and Discussion

Before growing the classification trees we summarized the variables by categories and by study outcome, i.e. whether students passed or failed the course. Feature selection was used to rank the variables by their importance for further analysis. Finally, the classification tree results for four different growing methods are presented.

Summary Statistics

As part of the data understanding phase we carried out the cross-tabulation for each variable and the study outcome after preparing and cleaning the data. The Table 2 reports the results. Based on the results shown majority of *Information Systems* students are female (over 63%). However, percentage of female students who successfully complete the course are higher (65%) which suggests that female students are more

likely to pass the course than their male counterpart. When it comes to age over 68% of students are above 30 with the age group between 30 and 40 being majority. This age group is also more likely to fail the course because their percentage of students who failed the course in this age group (39.7%) is higher than their overall participation in the student population (38.6%).

Table 2: Descriptive statistics (percentage) – Study outcome (453 students)

Variable	Domain	Count	Total	Pass	Fail
Gender	Female	286	63.1	65.0	61.5
	Male	167	36.9	35.0	38.5
Age	Under 30	136	30.0	30.4	29.7
	Between 30 and 40	175	38.6	37.4	39.7
	Above 40	142	31.3	32.2	30.5
Disability	Yes	19	4.2	3.3	5.0
	No	434	95.8	96.7	95.0
Ethnicity	Pakeha	318	70.2	75.7	65.3
	Maori & Pacific Islanders	41	9.1	2.8	14.6
	Others	94	20.8	21.5	20.1
Secondary school	No secondary school / NCEA Level 1 or 2	183	40.4	36.9	43.5
	University Entrance / NCEA Level 3	163	36.0	38.3	33.9
	Overseas Qualification or Other	107	23.6	24.8	22.6
Work status	Yes	351	77.5	78.0	77.0
	No	102	22.5	22.0	23.0
Early enrolment	Yes	317	70.0	72.4	67.8
	No	136	30.0	27.6	32.2
Course programme	Bachelor of Business	305	67.3	73.8	61.5
	Bachelor of Applied Sciences	148	32.7	26.2	38.5
Course block	First semester	139	30.7	31.3	30.1
	Second semester	201	44.4	48.6	40.6
	Third semester	113	24.9	20.1	29.3

Disability was shown to be a disadvantage for *Information Systems* students. Students with it are more likely to fail than those without it. There are huge differences in percentage of students who successfully

completed the course depending on their ethnic origin. Though Maori and Pacific Islanders make 9.1% of all students on this course their participation is significantly lower in the

“Pass” sub population (i.e. 2.8%) and higher in the “Fail” sub population (14.6%). Based on these results we can say that students with this ethnic origin are identified as students “at- risk”. Further methods of data mining will confirm this statement.

A substantial number of students (over 40%) don’t have secondary school qualification higher than NCEA Level 2 on the New Zealand National Qualification Framework and they are more vulnerable than the other two categories in this variable. Over three-fourths of *Information Systems* students are working and studying at the same time. Though the difference between those who work and those who do not is not statistically significant, it is interesting to note that the students who are working are more likely to pass the course than those not working.

We are using early enrolment as a proxy for motivation and good time management skills. Students who are motivated and are planning their study in advance will also enroll well before the enrolment closing date. The opposite category (late comers) participates with 30% in the total number of students, but these students are more likely to fail the course. Their participation in the “Fail” subpopulation increased from 30% to 32.2%.

One third of students on this course enrolled on the Bachelor of Applied Sciences program. They are more likely to fail the course when compared with students enrolled on the Bachelor of Business program. Finally, students studying this course in the summer semester are more likely to fail than those studying in the second and first semester.

Feature Selection

The number of predictor variables is not so large and we don’t have to select the subset of variables for further analysis which is the main purpose of applying feature selection to data. However, feature selection could be also used as a pre-processor for predictive data mining to rank predictors according to the strength of their relationship with dependent or outcome variable. During the feature selection process no specific form of relationship, neither linear nor nonlinear, is assumed. The outcome of the feature selection would be a rank list of predictors according to their importance for further analysis of the dependent variable with the other methods for regression and classification.

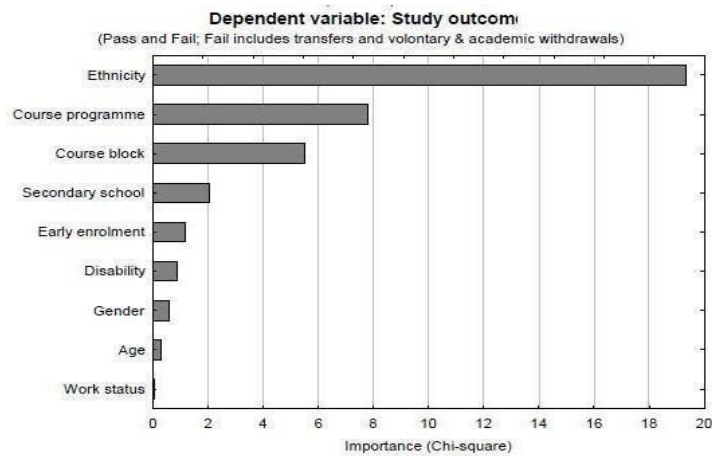


Figure 2: Importance plot for predictors

Results of feature selection are presented in Figure 2 on the importance plot and also in Table 3. The top three predictors for the study outcome are: ethnic origin of students, course programme they are enrolled on and course block, i.e. semester in which they are study.

Table 3: Best predictors for dependent variable

Variable	Chi-square	<i>P</i> -value
Ethnicity	19.35	0.00006
Course programme	7.80	0.00523
Course block	5.51	0.06354
Secondary school	2.06	0.35748
Early enrolment	1.16	0.28131
Disability	0.86	0.35363
Gender	0.58	0.44774
Age	0.28	0.86750
Work status	0.07	0.78940

From Table 3, *P*-values from the last column only the first three chi-square values are significant at 10% level. Though the results of the feature selection suggested continuing analysis with only the subset of predictors, which includes ethnicity, course programme and course block, we have included all available predictors in our classification tree analysis. We follow an advice given in Luan & Zhao (2006) who suggested that even though some variables may have little significance to the overall prediction outcome, they can be essential to a specific record.

Classification Trees

The objective of an analysis based on a classification tree is to identify factors that contribute the most to separation of successful from unsuccessful students. When the classification tree is formed we can calculate the probability of each student being successful. Once the classification tree is formed, it could be used in the new data set to predict the study outcome for newly en-rolled students.

The classification trees for study outcome are given in Figure 2 and 4. In each tree node the number of successful students (4th line, last column) and unsuccessful students (3rd line, last column) is given, as well as the percentages for each category (2nd column) and relative and absolute size of the node (5th line). The variable names above the nodes are the predictors that provided the best split for the node according to the classification and regression tree-style exhaustive search for univariate splits method. This method looks at all possible splits for each predictor variable at each node. The search stops when the split with the largest improvement in goodness of fit, based on the Gini measure of node impurity, is found. Immediately above the nodes are categories which describe these nodes. Note that all available predictor variables in the dataset were included in the classification tree analysis in spite of their insignificance detected in the feature selection section.

CHAID, exhaustive CHAID and QUEST

Three classification tree growing methods, namely: CHAID, exhaustive CHAID and QUEST generated exactly the same tree structure presented in Figure 3. It shows that only 2 variables were used to construct the tree: (1) ethnicity and (2) course program. All the other student

demographics variables were used but not included in the final model. We could change the stopping criteria to allow further growing of the tree, but that would result in nodes with just a few students. In the most extreme case we can continue splitting the tree until we create a terminal node for every student. However, we would get a model, i.e. classification tree that fits data better, but with more likely poor performance when used on a new data set. This phenomenon is known as overfitting the tree. The largest successful group (i.e. students who successfully completed the course) consists of 274 (60.5%) students (Node 3). Ethnic origin of students in this group is either Pakeha or other ethnic groups (excluding Maori and Pacific Islands students). Students in this group opted for the Bachelor of Business programme. The largest unsuccessful group (i.e. students who were unsuccessful) contains 138 students (30.5% of all participants) (Node 4). They belong to either Pakeha or other ethnic groups (excluding Maori and Pacific Islands students). The next largest group considered also as unsuccessful students, contains 41, i.e. 9.1% of all students, where 75.4% of them are unsuccessful (Node 1). They are described as Maori and Pacific Islands students.

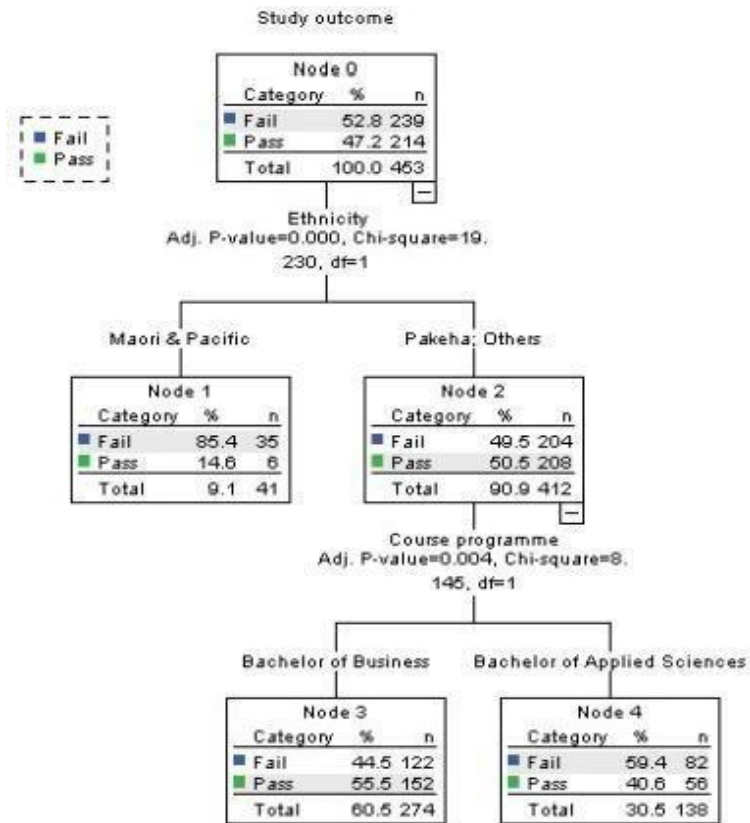


Figure 3: CHAID, exhaustive CHAID and QUEST classification tree

The overall percentage of correct classification for the study outcome is only 59.4% (Table 4).

This percentage of correct classification was achieved with 2 variables only.

The cross-validation estimate of the risk is 0.406 indicates that the category predicted by the model (successful or unsuccessful student) is wrong for 40.6% of the cases. So the risk of misclassifying a student is approximately 41%. This result is consistent with the results in the CHAID classification matrix (Table 4). The Overall percentage shows that the model only classified correctly 59% of students.

Table 4: CHAID classification matrix Observed

			<u>redicted</u>		
			Fail	Pass	Percent correct
Fail	117	122	49.0%		
Pass	62	152	71.0%		
<u>Overall percentage</u>	<u>65.4%</u>	<u>55.5%</u>	<u>59.4%</u>		

With large numbers of false positives (122) and few false negatives (62), the CHAID model is in itself poor at identifying an unsuccessful student (positive predictive value is only 49%). It will, however, pick up 65.4% of all unsuccessful students (the sensitivity). The predictive values, which take into account the prevalence of failing the course, are generally more important in determining the usefulness of a prediction model. The negative predictive value was of more concern to the course because the objective was to minimize the probability of being in error when deciding that a student is not at risk for not completing the course. However the CHAID model, as a classification tool, will pick-up with high probability successful students (negative predictive value is 71%) and correctly identifies 55.5% of those who pass the course (the specificity).

The classification matrix indicates also another problem with the model. For unsuccessful students it predicts failure for only 49% of them, which means that 51% of unsuccessful students are inaccurately classified with the successful students. Practical consequence of this misclassification is that these students would not receive additional learning support provided to the students “at risk”, simply because they will be classified by the model among successful students. This feature of the model is more critical than misclassification of the successful students among unsuccessful students (29% of successful students belong to this category). In this case these students may receive additional learning support even though they don’t need it. One option to increase percentage of correctly classified unsuccessful students is to change the misclassification cost matrix. With this option there is always a trade-off between increasing the percentage of correct classification of unsuccessful students and decreasing percentage of correct classification for successful students as well as decreasing the percentage of overall correct classification. In this case the increased cost for misclassification of unsuccessful to the successful group of students decreased significantly both remaining percentage of correct classification which was not compensated in equivalent increase in the initial 49%.

CART

Figure 5 shows the CART classification tree for study outcome. It shows that only three variables were used to construct the tree: (1) ethnicity, (2) course programme and (3) course block.

The largest successful group (i.e. students who successfully completed the course) consists of 215(47.5%) students (Node 5). The ethnic origin of students in this group is either Pakeha or other ethnic groups (excluding Maori and Pacific Islands students). Students in this group enrolled on the Bachelor of Business programme in either Semester 1 or Semester 2. The largest unsuccessful group (i.e. students who were unsuccessful) contains 138 students (30.5% of all participants)(Node 4). They belong to either Pakeha or other ethnic groups (excluding Maori and Pacific Islands students). The next largest group

considered also as unsuccessful students, contains 41, i.e 9.1% of all students, where 75.4% of them are unsuccessful (Node 1).

They are described as Maori and Pacific Islands students.

Table 5: CART classification matrix

			redicted		
			Fail	Pass	Percent correct
Fail	149	90	62.3%		
Pass	89	125	58.4%		
<u>Overall percentage</u>			<u>62.6%</u>	<u>58.1%</u>	<u>60.5%</u>

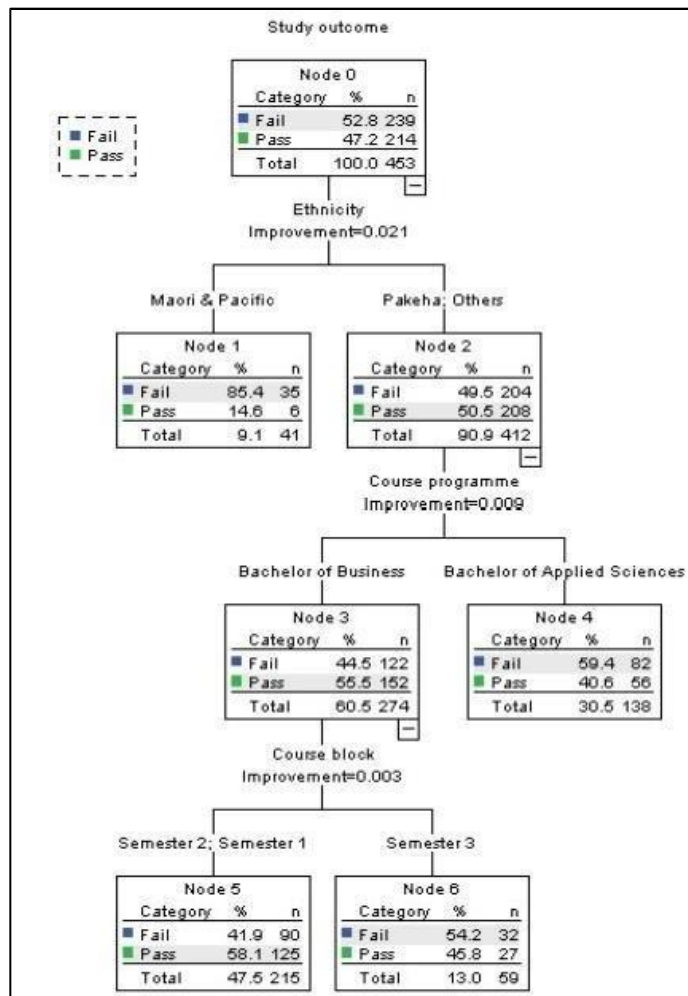


Figure 5: CART classification tree

he cross-validation estimate of the risk is 0.446 indicates that the category predicted by the model (successful or unsuccessful student) is wrong for 44.6% of the cases. The CART classification matrix

(Table 5) shows that model classify correctly 61% of students. This is slight increase in comparison to the CHAID model. The numbers of false positives (90) for the CART model decreases and therefore increasing the positive predictive value to 62.3%. In other words it will work better than the CHAID model at identifying an unsuccessful student. The price paid for increasing accuracy is reflected in decreasing sensitivity. The CART model will pick up 62.6% of all unsuccessful students (CHAID model 65.4%). At the same time the specificity will increase to 58.1% (CHAID model 55.5%).

Concluding Remarks

This study examines the background information from enrolment data that impacts upon the study outcome of Information Systems students at the Open Polytechnic. Based on results from feature selection (Figure 2 and Table 3), the CHAID tree presented in Table 4 and Figures 3 and 4, the CART trees presented in Table 5 and Figures 5 and 6 it was found that the most important factors that help separate successful from unsuccessful students are ethnicity, course programme and course block. Demographic data such as gender and age though significantly related to the study outcome, according to the feature selection result, were not used in the classification trees. Unfortunately the classification accuracy from the classification trees was not very high. In the case of the CHAID tree the overall classification accuracy was 59.4% and in the case of the CART tree slightly higher at 60.5%. This would suggest that the background information (gender, age, ethnicity, disability, secondary school, work status, and early enrolment) gathered during the enrolment process, does not contain sufficient information for an accurately separation of successful and unsuccessful students. Classifying students based on pre-enrollment information and the rules presented for each node would allow the administrative and academic staff to identify students who would be “at risk” of dropping the course even before they start with their study. Then the student support systems, such as orientation, advising, and mentoring programs, could be used to positively impact the academic successes of such students. This study is limited in three main ways that future research can perhaps address. Firstly, this research is based on background information only. Leaving out other important factors (academic achievement, number of courses completed, motivation, financial aids, etc.) that may affect study outcome, could distort results obtained with classification trees. For example, including the assignment mark after the submission of the first course assignment would probably improve predictive accuracy of the models. To improve the model, more attributes could be included to obtain prediction models with lower misclassification errors. However, the model in this case would not be a tool for pre-enrolment, i.e. early identification of ‘at-risk’ students. Secondly, we used a dichotomous variable for the study outcome with only two categories: pass and fail. However, splitting the fail category into those who stayed on the course but eventually failed the course and those who voluntary transfer or were withdrawn from the course would probably provide better profiling for each of the three categories of study outcomes. The only problem we might have with the three categories for study outcome is a low prediction accurate as a result of relatively small data set

for the course. Thirdly, from a methodological point of view an alternative to a classification tree should be considered. The prime candidates to be used with this data set are logistic regression and neural networks.

RESULT :

Analysis done on the above case study , examined the background information from enrolment data that impacts upon the study outcome of Information Systems students at the Open Polytec

DATA MINING CASE STUDY USING THE CRISP-DM STANDARD

AIM :

CRISP-DM CASE STUDY - Analysing Automobile Warranty Claims: Example Of The CRISP– DM Industry Standard Process In Action

Quality assurance continues to be a priority for automobile manufacturers, including Daimler

Chrysler. Jochen Hipp of the University of Tubingen, Germany, and Guido Lindner of

DaimlerChrysler AG, Germany, investigated patterns in the warranty claims for DaimlerChrysler automobiles.

1. Business Understanding Phase

DaimlerChrysler's objectives are to reduce costs associated with warranty claims and improve customer satisfaction. Through conversations with plant engineers, who are the technical experts in vehicle manufacturing, the researchers are able to formulate specific business problems, such as the following:

Are there interdependencies among warranty claims?

Are past warranty claims associated with similar claims in the future?

Is there an association between a certain type of claim and a particular garage?

The plan is to apply appropriate data mining techniques to try to uncover these and other possible associations.

Data Understanding Phase

The researchers make use of DaimlerChrysler's Quality Information System (QUIS), which contains information on over 7 million vehicles and is about 40 gigabytes in size. QUIS contains production details about how and where a particular vehicle was constructed, including an average of 30 or more sales codes for each vehicle. QUIS also includes warranty claim information, which the garage supplies, in the form of one of more than 5000 possible potential causes.

The researchers stressed the fact that the database was entirely unintelligible to domain non-experts: "So experts from different departments had to be located and consulted; in brief a task that turned out to be rather costly." They emphasise that analysts should not underestimate the importance, difficulty and

potential cost of this early phase of the data mining process, and that shortcuts here may lead to expensive reiterations of the process downstream.

Data Preparation Phase

The researchers found that although relational, the QUIS database had limited SQL access. They needed to select the cases and variables of interest manually, and then manually derive new variables that could be used for the modelling phase. For example, the variable *number of days from selling date until first claim* had to be derived from the appropriate date attributes.

They then turned to proprietary data mining software, which had been used at DaimlerChrysler on earlier projects. Here they ran into a common roadblock – that the data format requirements varied from algorithm to algorithm. The result was further exhaustive pre-processing of the data, to transform the attributes into a form usable for model algorithms. The researchers mention that the data preparation phase took much longer than they had planned.

Modelling Phase

Since the overall business problem from phase 1 was to investigate dependence among the warranty claims, the researchers chose to apply the following techniques: (1) Bayesian networks and (2) Association rules. Bayesian networks model uncertainty by explicitly representing the conditional dependencies among various components, thus providing a graphical visualisation of the dependency relationships among the components.

As such, Bayesian networks represent a natural choice for modelling dependence among warranty claims. Association rules are also a natural way to investigate dependence among warranty claims since the confidence measure represents a type of conditional probability, similar to Bayesian networks.

The details of the results are confidential, but we can get a general idea of the type of dependencies uncovered by the models. One insight the researchers uncovered was that a particular combination of construction specifications doubles the probability of encountering an automobile electrical cable problem. Daimler Chrysler engineers have begun to investigate how this combination of factors can result in an increase in cable problems.

The researchers investigated whether certain garages had more warranty claims of a certain type than did other garages. Their association rule results showed that, indeed, the confidence levels for the rule “If garage X, then cable problem,” varied considerably from garage to garage. They state that further investigation is warranted to reveal the reasons for the disparity.

Evaluation Phase

The researchers were disappointed that the support for sequential-type association rules was relatively small, thus precluding generalisation of the results, in their opinion. Overall, in fact, the researchers state: “In fact, we did not find any rule that our domain experts would judge as interesting, at least at first sight.” According to this criterion, then, the models were found to be lacking in effectiveness and to fall short of the objectives set for them in the business understanding phase. To account for this, the researchers point to the “legacy” structure of the database, for which automobile parts were categorised by garages and factories for historic or technical reasons and not designed for data mining. They suggest adapting and redesigning the database to make it more amenable to knowledge discovery.

Deployment Phase

The researchers have identified the foregoing project as a pilot project, and as such, do not intend to deploy any large-scale models from this first iteration. After the pilot project, however, they have applied the lessons learned from this project, with the goal of integrating their methods with the existing information technology environment at DaimlerChrysler. To further support the original goal of lowering claims costs, they intend to develop an intranet offering mining capability of QUIS for all corporate employees.

DISCUSSION :

Lessons drawn from this case study.

First, the general impression one draws is that uncovering hidden nuggets of knowledge in databases is a rocky road. In nearly every phase, the researchers ran into unexpected roadblocks and difficulties. This tells us that actually applying data mining for the first time in a company requires asking people to do something new and different, which is not always welcome. Therefore, if they expect results, corporate management must be 100% supportive of new data mining initiatives.

Another lesson to draw is that intense human participation and supervision is required at every stage of the data mining process. For example, the algorithms require specific data formats, which may require substantial pre-processing. Regardless of what some software vendor advertisements may claim, you can’t just purchase some data mining software, install it, sit back and watch it solve all your problems. Data mining is not magic. Without skilled human supervision, blind use of data mining software will only provide you with the wrong answer to the wrong question applied to the wrong type of data. The wrong analysis is worse than no analysis, since it leads to policy recommendations that will probably turn out to be expensive failures.

Finally, from this case study we can draw the lesson that there is no guarantee of positive results when mining data for actionable knowledge, any more than when one is mining for gold. Data mining is not a panacea for solving business problems. But used properly, by people who understand the models involved, the data requirements and the overall project objectives, data mining can indeed provide actionable and highly profitable results.

RESULT :

CRISP-DM CASE STUDY - Analysing Automobile Warranty Claims: Example Of The CRISP-DM

Industry Standard Process In action is done and output verified