



INDUS INSTITUTE OF TECHNOLOGY & ENGINEERING

DATA WAREHOUSE AND MINING

LAB PRACTICALS

Name: Ahmed Shaikh

Enrollment Number: IU19410050002

CLASS: CE-A

Semester : 7th

PRACTICAL - 1

AIM: Study Practical: Introduction to Weka



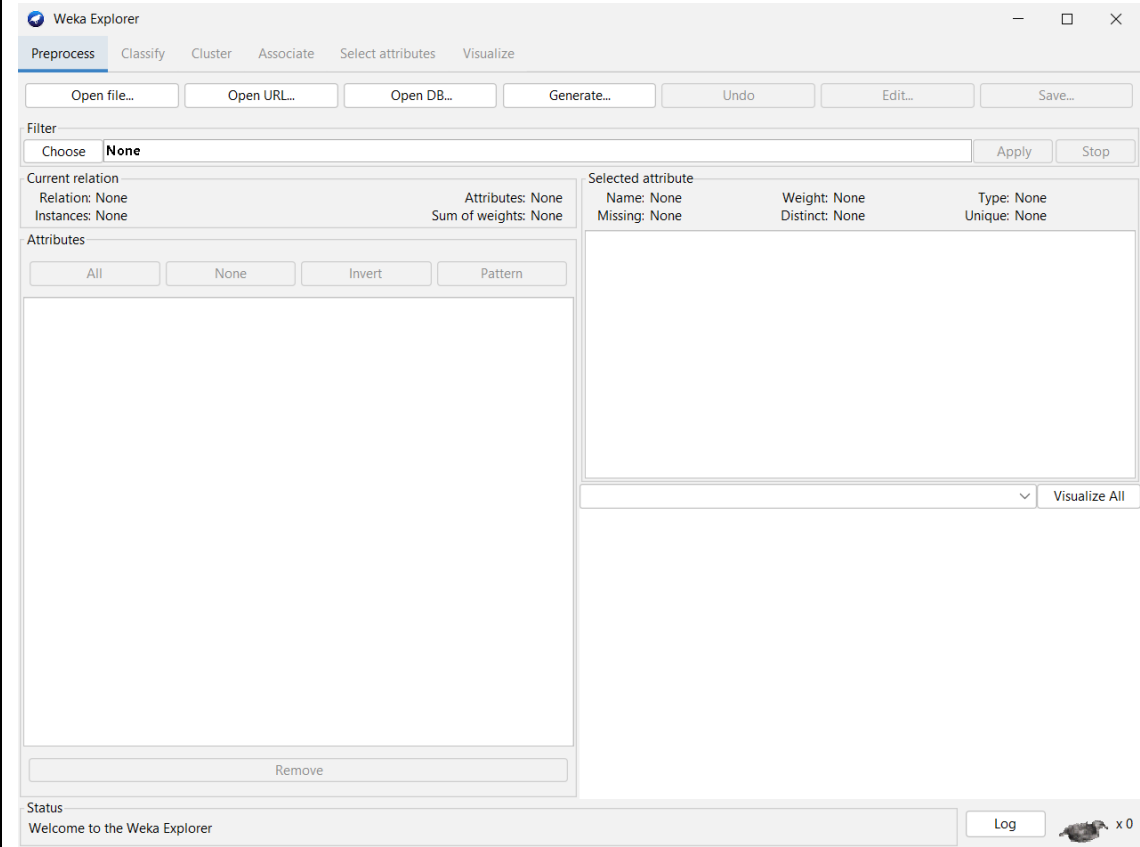
WEKA - an open source software provides tools for data preprocessing, implementation of several Machine Learning algorithms, and visualization tools so that you can develop machine learning techniques and apply them to real-world data mining problems.

Explorer :

On the top, there are several tabs as listed here –

- 1) **Preprocess:** This allows us to choose the data file.
- 2) **Classify:** This allows us to apply and experiment with different algorithms on preprocessed data files.
- 3) **Cluster:** This allows us to apply different clustering tools, which identify clusters within the data file.
- 4) **Association:** This allows us to apply association rules, which identify the association within the data.
- 5) **Select attributes:** These allow us to see the changes on the inclusion and exclusion of attributes from the experiment.
- 6) **Visualize:** This allows us to see the possible visualisation produced on the data set in a 2D format, in scatter plot and bar graph output.

The user cannot move between the different tabs until the initial preprocessing of the data set has been completed.

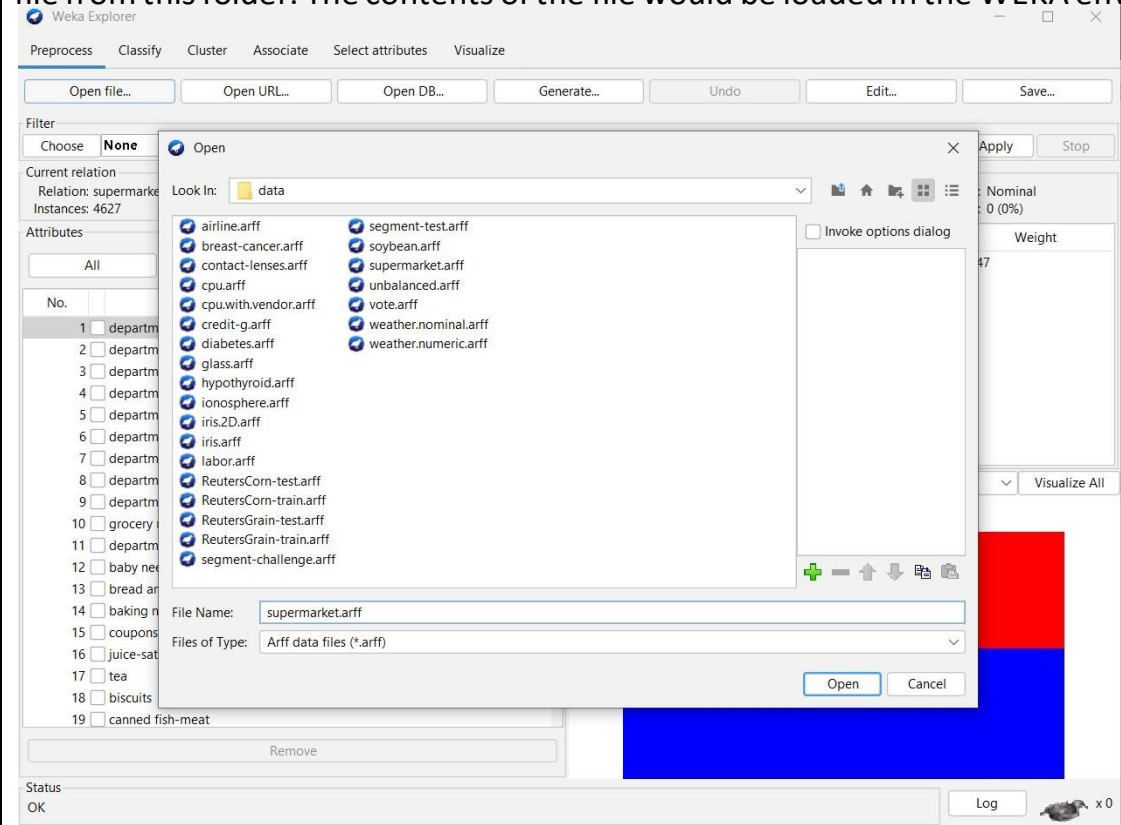


The data can be loaded from the following sources –

- Local file system
- Web
- Database

Loading Data from Local File System

Click on the Open file button. Now, navigate to the folder where your data files are stored. select any data file from this folder. The contents of the file would be loaded in the WEKA environment.



Loading Data from Web

Once you click on the Open URL button, you can see a window as follows –

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Open file...Open URL...Open DB...Generate...UndoEdit...Save...

Filter

ChooseNone

ApplyStop

Current relation

Relation: None

Instances: None

Attributes: None

Sum of weights: None

Selected attribute

Name: None

Weight: None

Type: None

Missing: None

Distinct: None

Unique: None


Attributes

AllNoneInvertPattern

Remove

Status

Welcome to the Weka Explorer

Logx0

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Open file...Open URL...Open DB...Generate...UndoEdit...Save...

Filter

ChooseNone

ApplyStop

Current relation

Relation: weather

Instances: 14

Attributes: 5

Sum of weights: 14

Selected attribute

Name: outlook

Missing: 0 (0%)

Distinct: 3

Type: Nominal

Unique: 0 (0%)

Attributes

AllNoneInvertPattern

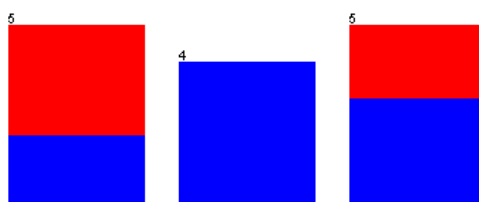
No.	Name
1	<input type="checkbox"/> outlook
2	<input type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input type="checkbox"/> play

Remove

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


Class: play (Nom)

Visualize All



Status

OK

Logx0

Current relation

Relation: weather

Instances: 14

Attributes: 5

Sum of weights: 14

Attributes

AllNoneInvertPattern

No.	Name
1	<input type="checkbox"/> outlook
2	<input checked="" type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input type="checkbox"/> play

Remove

Selected attribute

Name: temperature

Missing: 0 (0%)

Distinct: 12

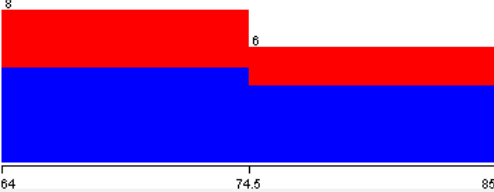
Type: Numeric

Unique: 10 (71%)

Statistic	Value
Minimum	64
Maximum	85
Mean	73.571
StdDev	6.572

Class: play (Nom)


Visualize All



Status

OK

Log

 x 0

Current relation

Relation: weather

Instances: 14

Attributes: 5

Sum of weights: 14

Attributes

AllNoneInvertPattern

No.	Name
1	<input type="checkbox"/> outlook
2	<input type="checkbox"/> temperature
3	<input type="checkbox"/> humidity
4	<input type="checkbox"/> windy
5	<input checked="" type="checkbox"/> play

Remove

Selected attribute

Name: play

Missing: 0 (0%)

Distinct: 2

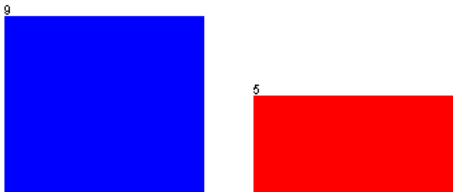
Type: Nominal

Unique: 0 (0%)

No.	Label	Count	Weight
1	yes	9	9
2	no	5	5

Class: play (Nom)


Visualize All



Status

OK

Log

 x 0

Loading Data from DB

Connection

Once you click on the Open DB button, you can see a window as follows –

Query

Execute

Clear

History...

max. rows 100

Result

Close

Close all

Re-use query

Optimal width

Info

Clear

Copy

Generate sparse data

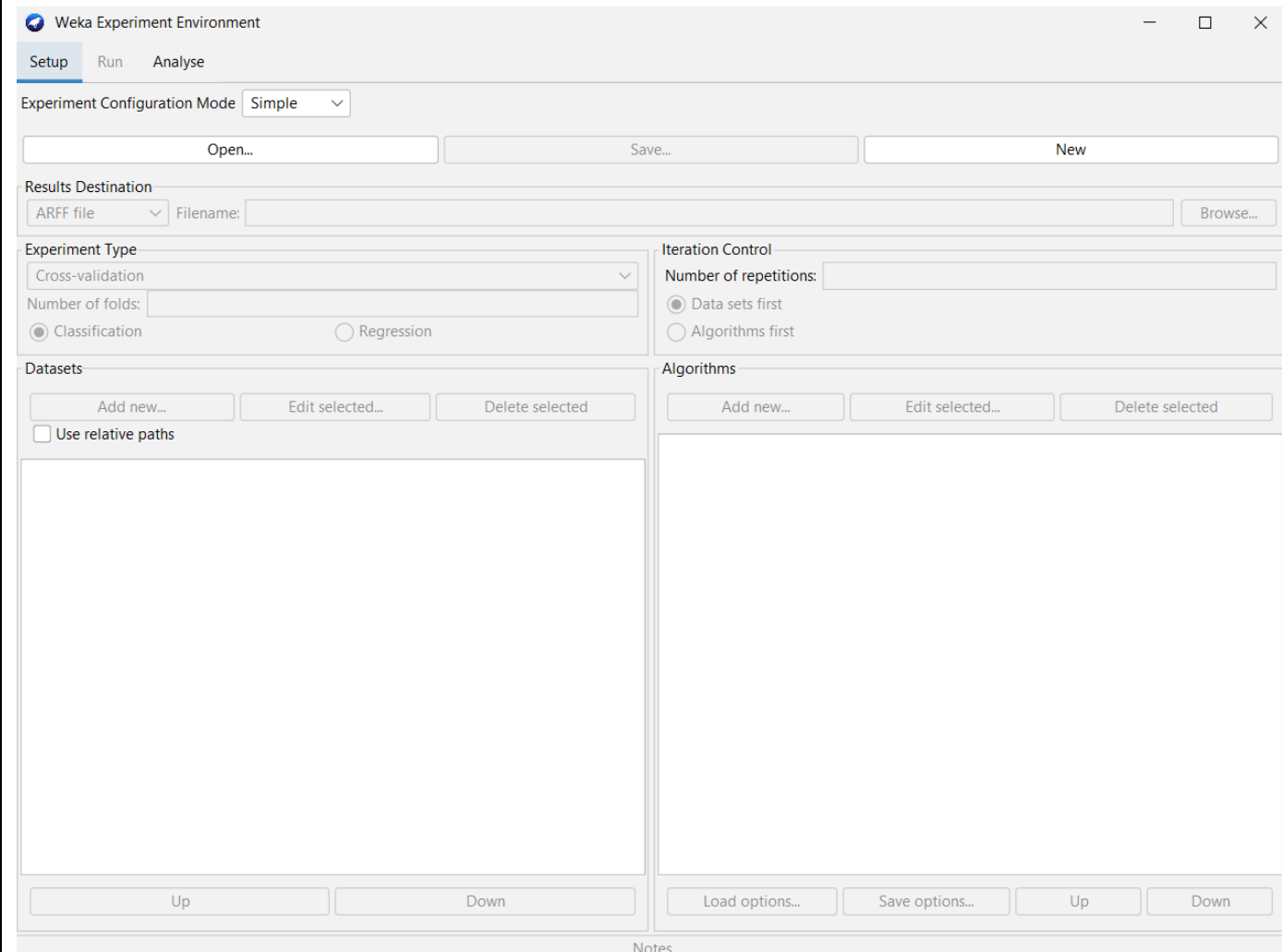
OK

Cancel

Set the connection string to your database, set up the query for data selection, process the query and load the selected records in WEKA.

Experimenter :

The Experimenter option available in Weka enables the user to perform some experiments on the data set by choosing different algorithms and analysing the output. It has the following components.



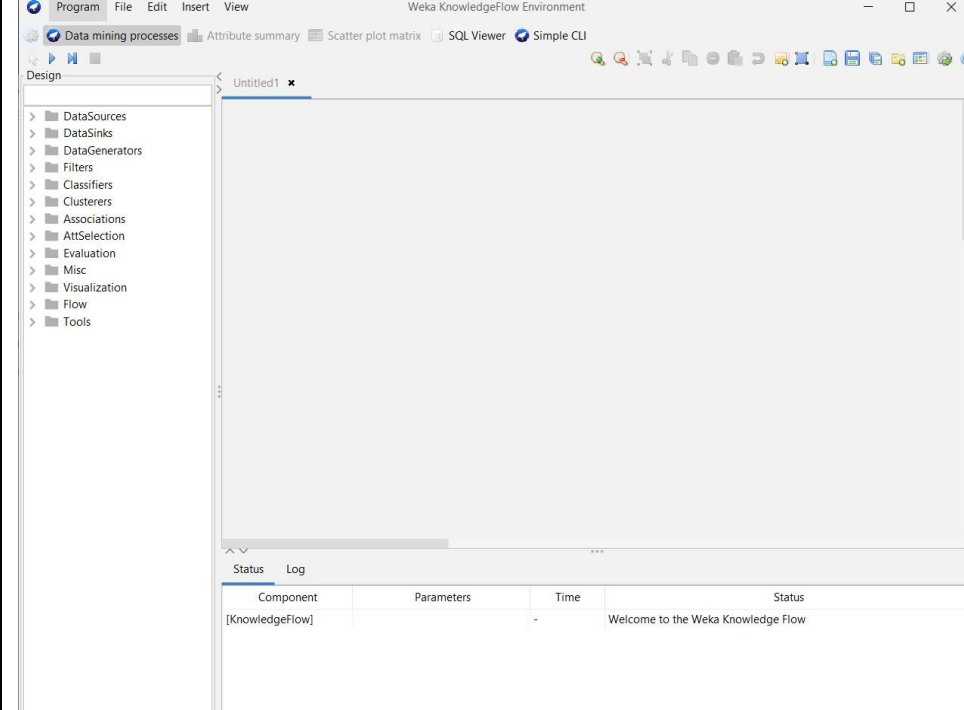
Setup: The first one is to set up the data sets, algorithms output destination, etc. Figure 4 shows an example of comparing the J4.8 decision tree with ZeroR on the IRIS data set. We can add more data sets and compare the outcome using more algorithms, if required.

Run: You can use this tab to run the experiment.

Analyse: This tab can be used to analyse the result.

Knowledge Flow:

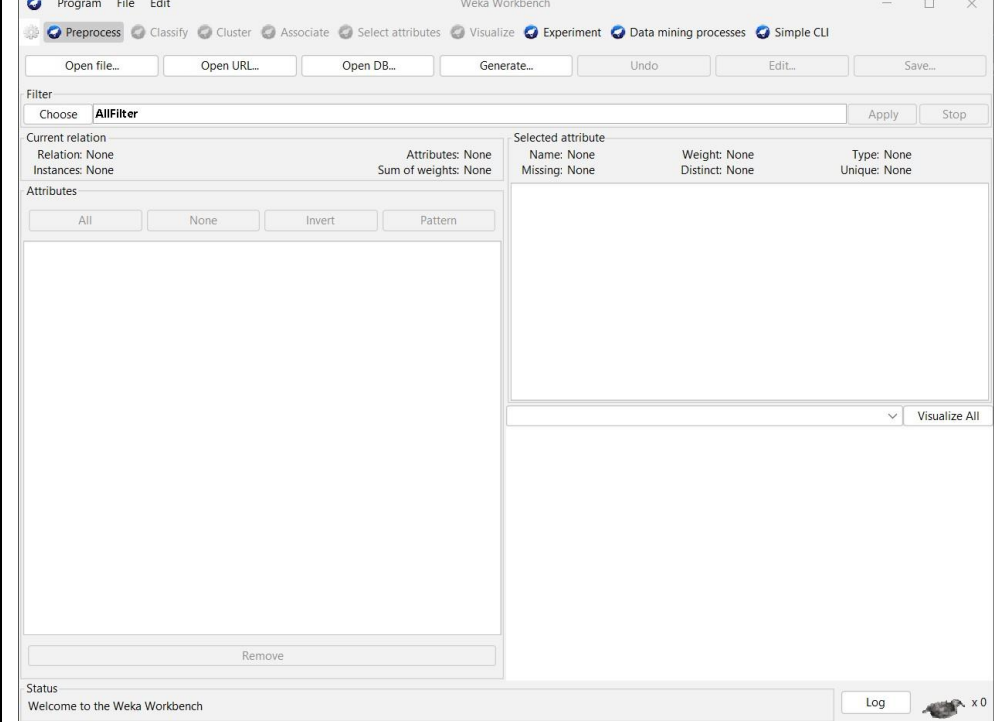
The Knowledge Flow interface is an alternative to the Explorer. You lay out filters, classifiers, evaluators, and visualizers interactively on a 2D canvas and connect them together with different kinds of connector. Data and classification models flow through the diagram.



Workbench :

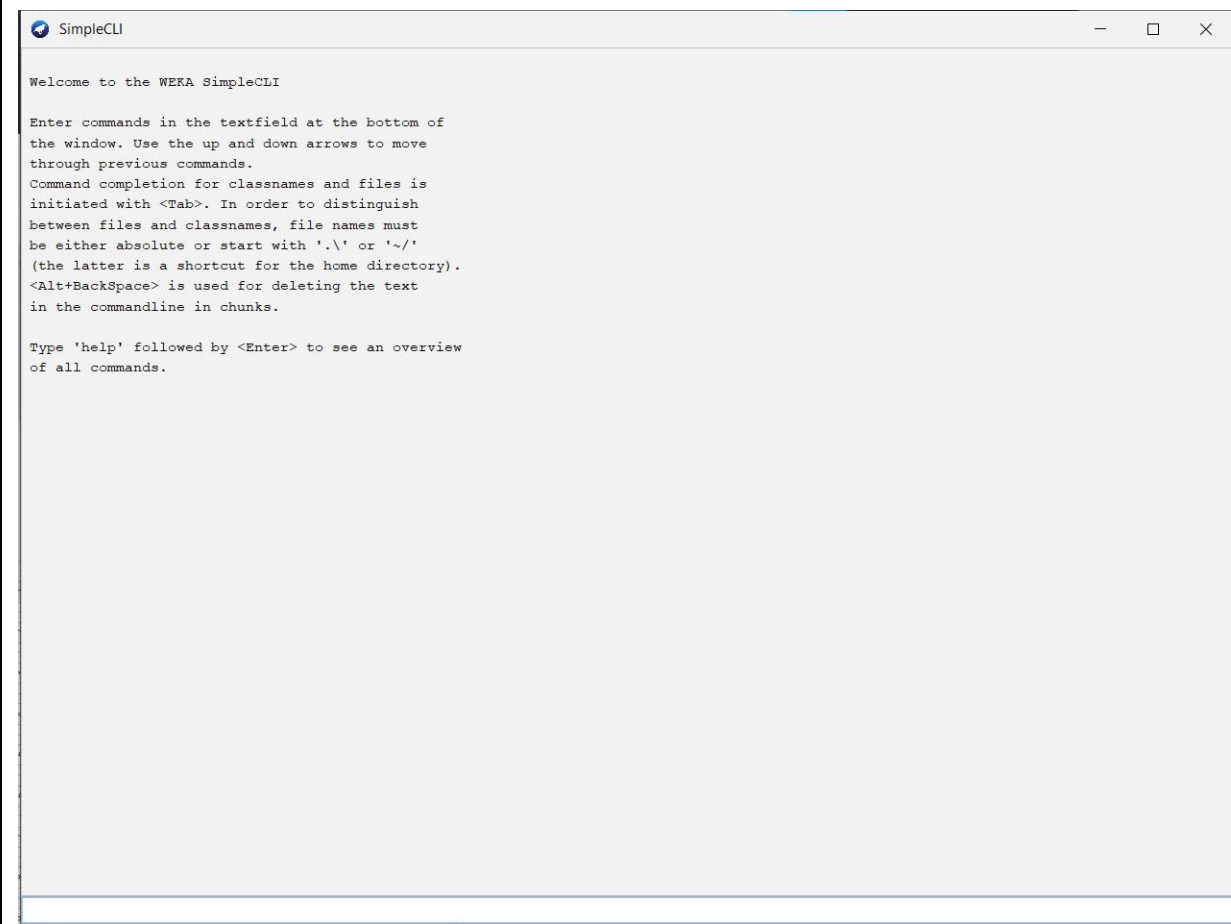
The Weka Workbench is an environment that combines all of the GUI interfaces into a single interface. It is useful if you find yourself jumping a lot between two or more different interfaces, such as between the Explorer and the Experiment

Environment.



Simple CLI:

Weka can be used from a simple Command Line Interface (CLI).This is powerful because you can write shellscripts to use the full API from command line calls with parameters, allowing you to build models, run experiments and make predictions without a graphical user interface.The SimpleCLI provides an environment where you can quickly and easily experiment with the Weka command line interface commands.



PRACTICAL - 2

AIM: Study Practical: Introduction to RStudio

1. Statistical Operations

CODE:

```

1 x <- c(2,7,3,4.2,2,2,54,-21,8,-5,2)
3 #Mean
4 mean(x)
5
6 #Median
7 median(x)
8
9 #Minimum
10 min(x)
11
12 #Maximum
13 max(x)
14
15 #Standard Deviation
16 sd(x)
17
18 #Length
19 length(x)
20
21 #Quantile
22 quantile(x)
23
24 #Range
25 range(x)
26
27 #Interquantile range
28 IQR(x)
29
30 #Variance
31 var(x)
32
33 #Coefficient of Variance
34 cv <- sd(x) / mean(x) * 100

```

OUTPUT

```

: [1] 5.290909
: [1] 2
[1] -21
[1] 54
[1] 17.97139
[1] 11
   0%   25%   50%   75%  100%
-21.0   2.0   2.0   5.6  54.0
[1] -21  54
[1] 3.6
[1] 322.9709

```

2. Charts

a. Density Plot

CODE:

OUTPUT:

b. Histogram

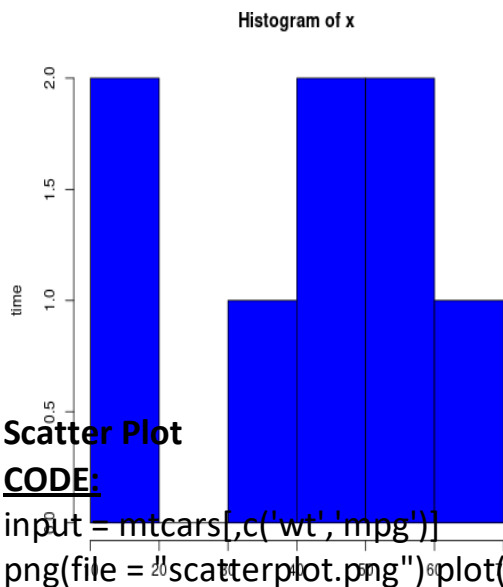
CODE:

```
x=c(10,12,54,67,45,34,56,43)
```

```
png(file = "histogram.png")
```

```
hist(x,xlab="freq",ylab="time",col="blue",border="black",breaks=5)dev.off()
```

OUTPUT:



c. Scatter Plot

CODE:

```
input = mtcars[,c('wt','mpg')]
```

```
png(file = "scatterplot.png")plot(x
```

```
= input$wt,y = input$mpg,
```

```
  xlab = "Weight",
```

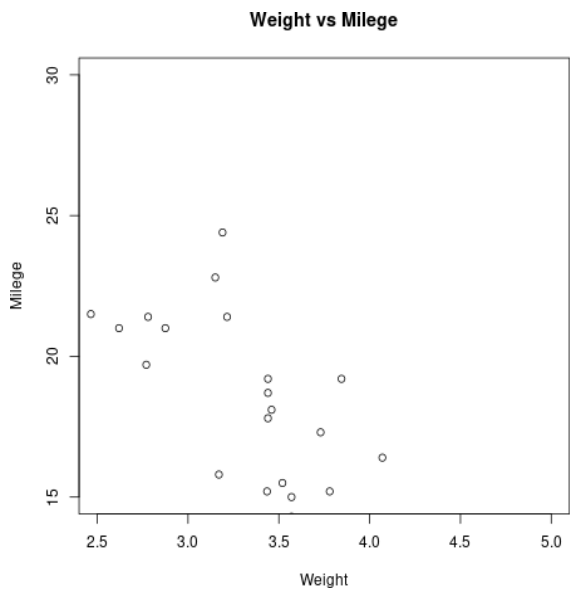
```
  ylab = "Milege",
```

```
  xlim = c(2.5,5),
```

```
  ylim = c(15,30),
```

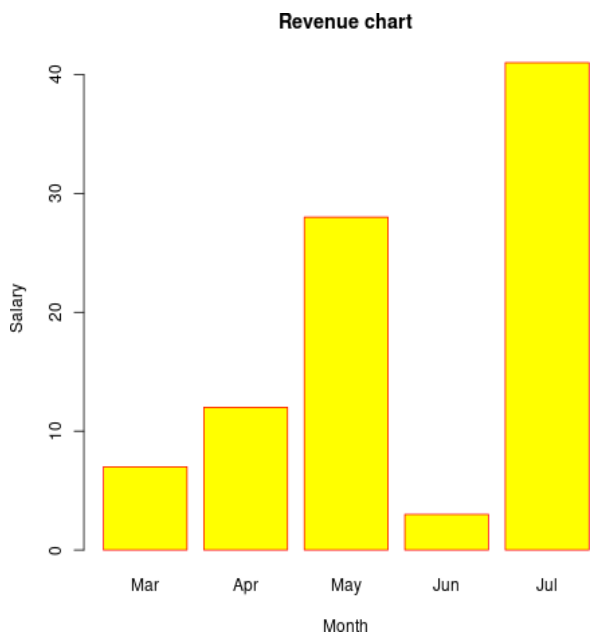
```
  main = "Weight vs Milege" )
```

OUTPUT:



d. Bar Chart

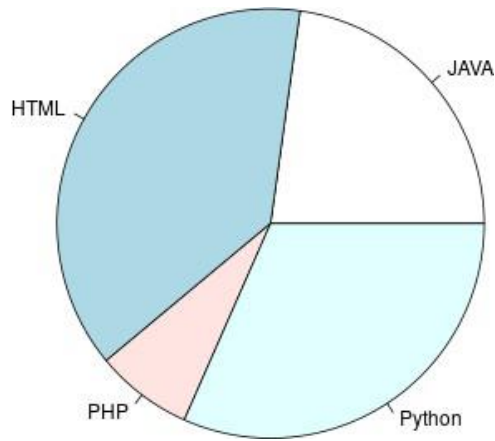
```
CODE:  
H = c(7,12,28,3,41)  
M = c("Mar","Apr","May","Jun","Jul") png(file  
= "barchart_months_revenue.png")  
barplot(H,names.arg=M,xlab="Month",ylab="Salary",col="yellow",main="Revenue  
chart",border="red")
```



e. Pie Chart

```
CODE:  
x = c(31, 52, 10, 43)  
labels = c("JAVA", "HTML", "PHP", "Python")png(file =  
"language.png")  
pie(x,labels)
```

OUTPUT:

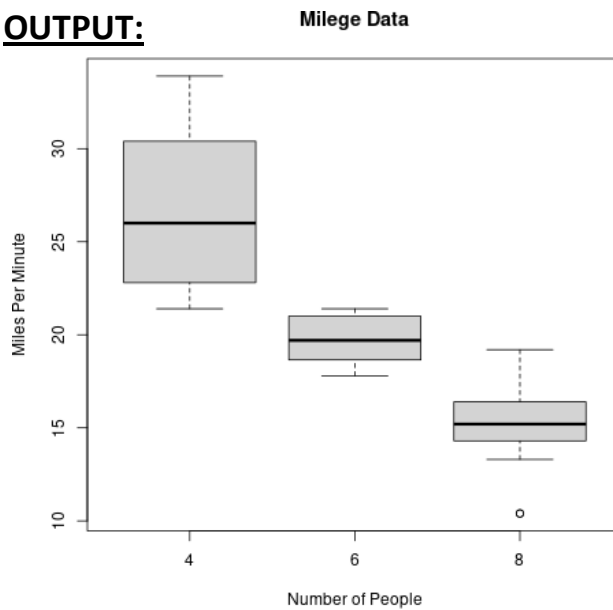


f. Box Plot

CODE:

```
png(file = "boxplot.png")
boxplot(mpg ~ cyl, data = mtcars, xlab = "Number of People", ylab =
  "Miles Per Minute", main = "Mileage Data")
```

OUTPUT:



PRACTICAL - 3

AIM: To perform Linear Regression in RStudio

Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable.

$y=ax+b$
y response variable
x predictor variable
a & b are constants called as coefficients

lm() function `lm(formula=y~x, data)`
predict() function `predict(object,newdata)`

Problem 1: Develop the equation of simple regression line to predict sales(y) from advertising(x) using the given data

CODE:

```
x = c(12.5,3.7,21.6,60.6,37.6,6.1,16.8,41.2)
y = c(148,55,338,994,541,89,126,379)
y.lm = lm(y~x)
coeffs = coefficients(y.lm)
coeffs
newdata = data.frame(x=50)
predict(y.lm,newdata)
print(summary(y.lm))
summary(y.lm)$r.squared
#plotting
png(file = "scatterplot.png")
plot(x,y)
```

OUTPUT:

```
(Intercept)      x
-45.149       15.14840
712.2706

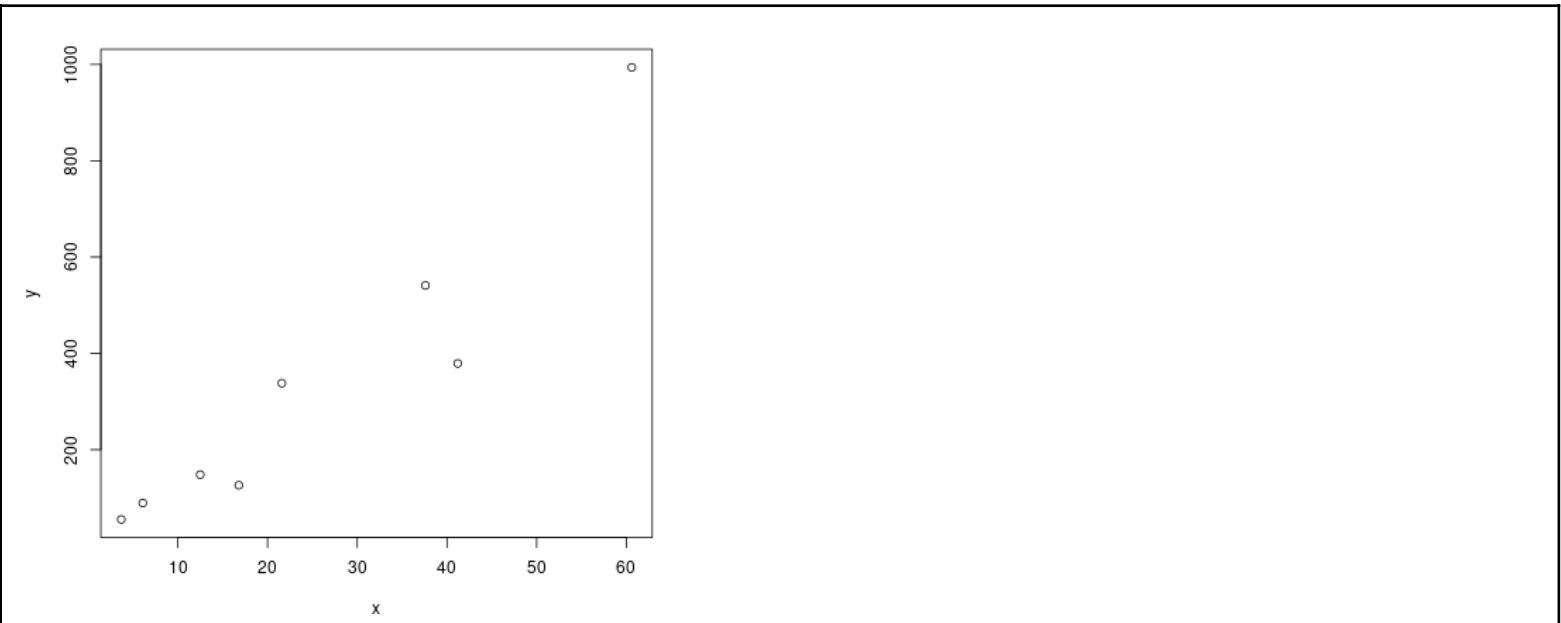
Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-199.97  -17.99   29.16   47.06  121.16

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -45.149     63.654  -0.709  0.504761
x             15.148       2.047   7.402  0.000312 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 107 on 6 degrees of freedom
Multiple R-squared:  0.9013,    Adjusted R-squared:  0.8848
F-statistic: 54.79 on 1 and 6 DF,  p-value: 0.0003123

[1] 0.9012992
```



Problem 2: Use a computer to develop the equation of the regression model for the following data. Comment on the regression coefficients. determine the predicted value of y for x1=33, x2=29,x3=13

CODE:

```
x1 = c(21,43,56,19,29,34,40,32,16,18,27,31)
x2 = c(6,25,42,27,20,45,33,14,4,31,12,3)
x3 = c(5,8,25,9,12,21,14,11,7,16,10,8)
y = c(114,94,87,98,101,85,94,107,119,93,108,117)
y.lm = lm(y~x1+x2+x3)
coeffs = coefficients(y.lm)
coeffs
newdata = data.frame(x1=33, x2=29,x3=13)
predict(y.lm,newdata)
print(summary(y.lm))
summary(y.lm)$r.squared
```

OUTPUT:

```
(Intercept)      x1      x2      x3
118.55951024  -0.07940245  -0.88428115   0.37690982

1
95.1949

Call:
lm(formula = y ~ x1 + x2 + x3)

Residuals:
    Min       1Q   Median       3Q      Max
-2.7481 -1.6934  0.5343  1.1214  2.6097

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 118.55951    1.85798   63.811 4.05e-12 ***
x1          -0.07940     0.06848   -1.159   0.280
x2          -0.88428     0.08631  -10.245 7.08e-06 ***
x3           0.37691     0.21973    1.715   0.125
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.134 on 8 degrees of freedom
Multiple R-squared:  0.975,    Adjusted R-squared:  0.9656
F-statistic: 103.8 on 3 and 8 DF,  p-value: 9.582e-07

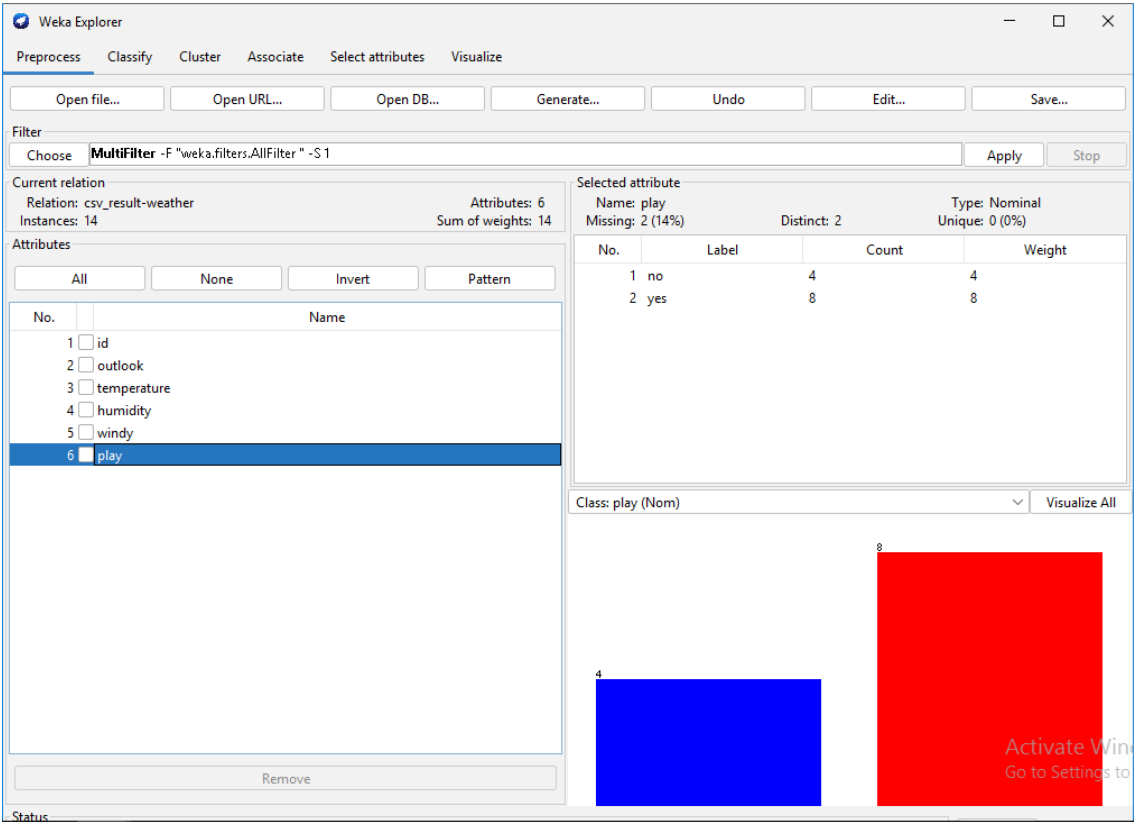
[1] 0.9749568
```

PRACTICAL - 4

AIM: Implementation of Preprocessing or Data Cleaning(Missing Values) in Weka.

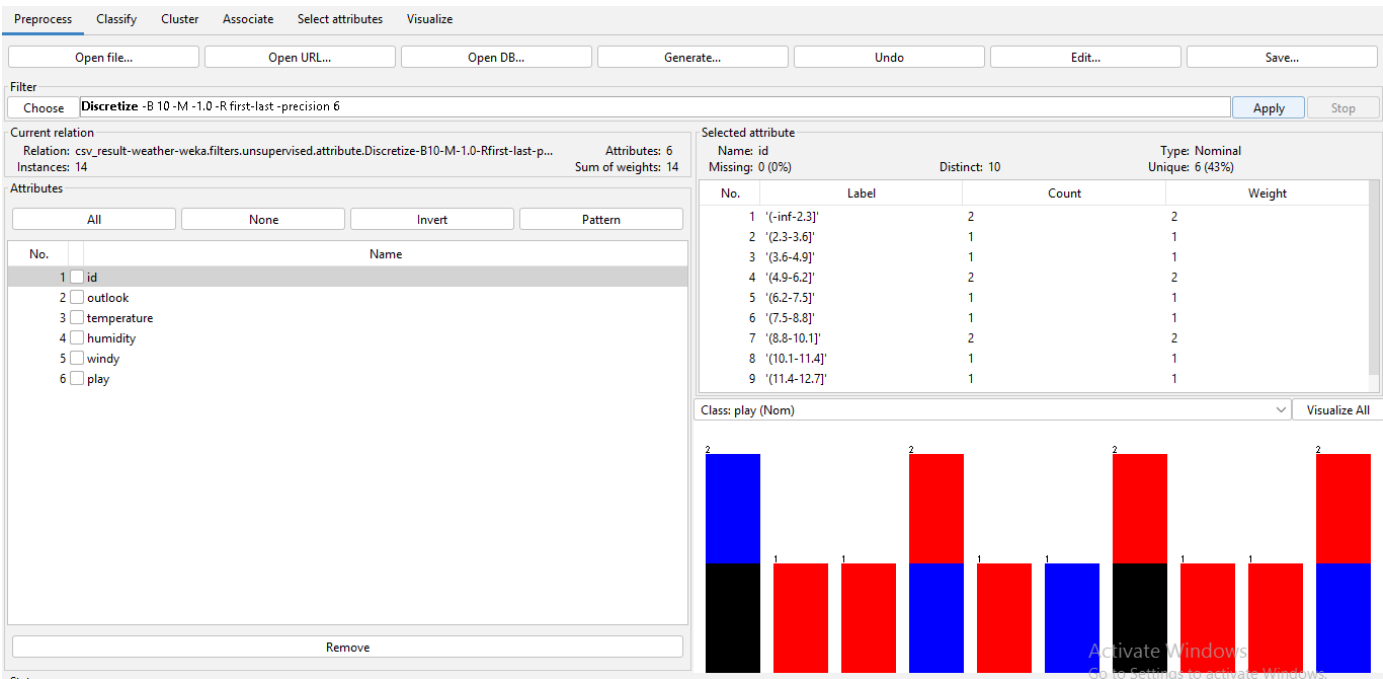
1) REPLACE MISSING VALUES:

Replaces all missing values for nominal and numeric attributes in a dataset with range.



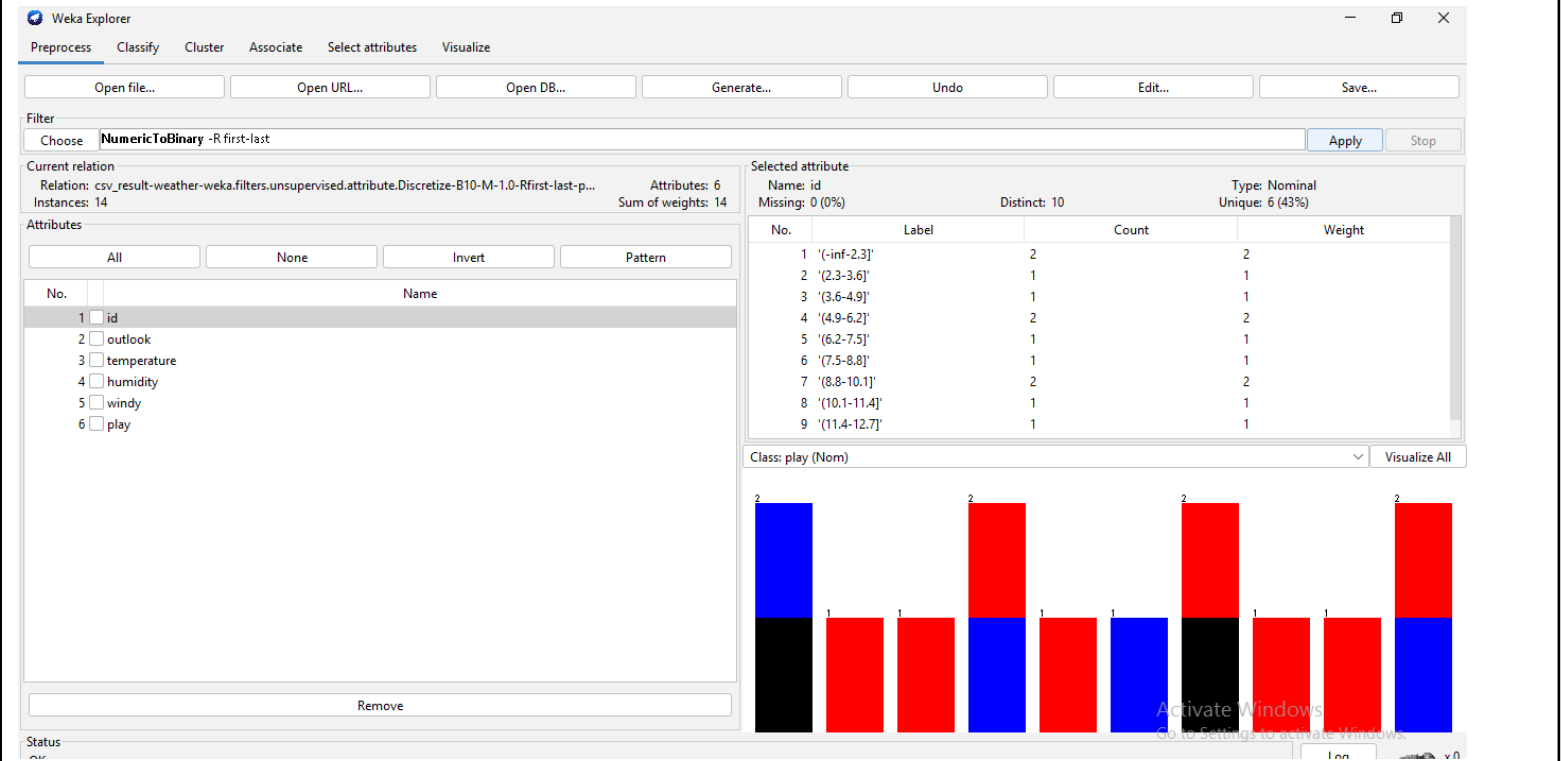
2) DISCRETIZE:

An instance filter that discretizes a range of numeric attributes in the dataset into nominal attributes



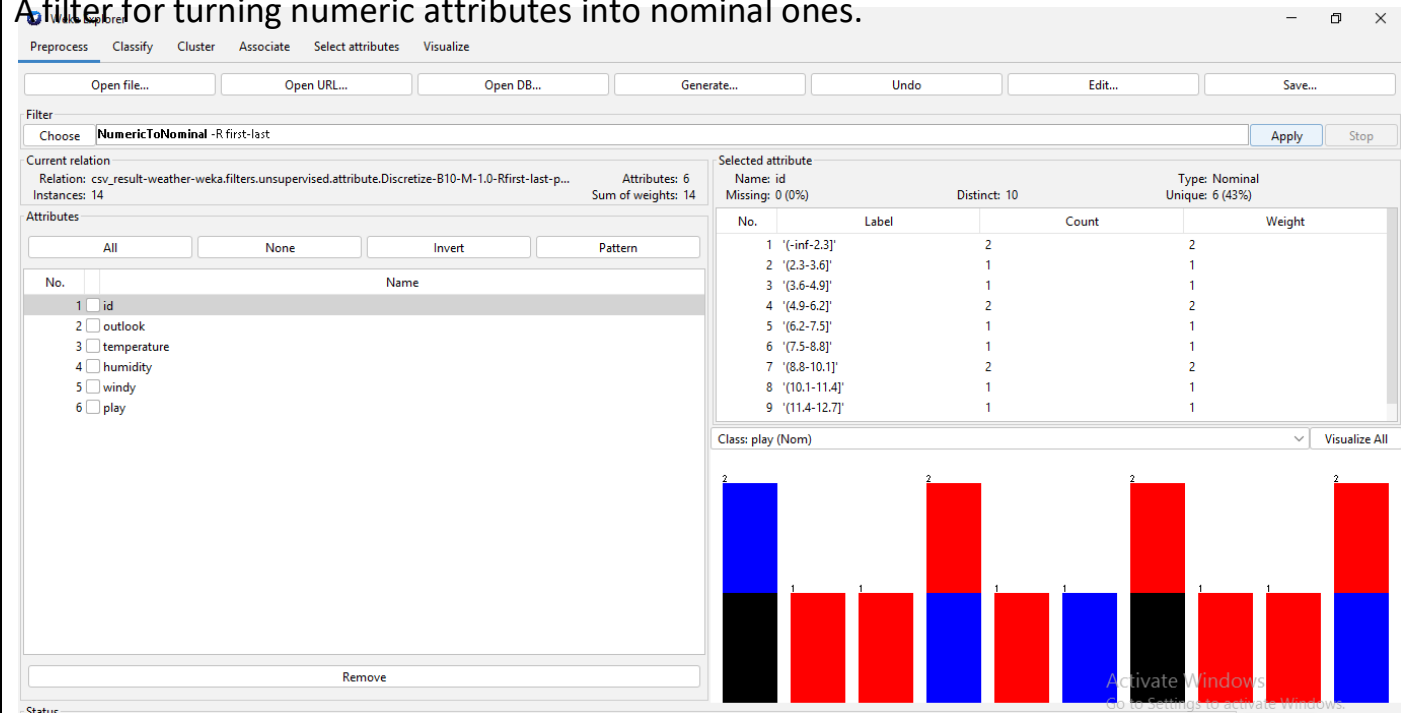
3) NUMERIC TO BINARY:

Converts all numeric attributes into binary attributes (apart from the class attribute, if set): if the value of the numeric attribute is exactly zero, the value of the new attribute will be zeros.

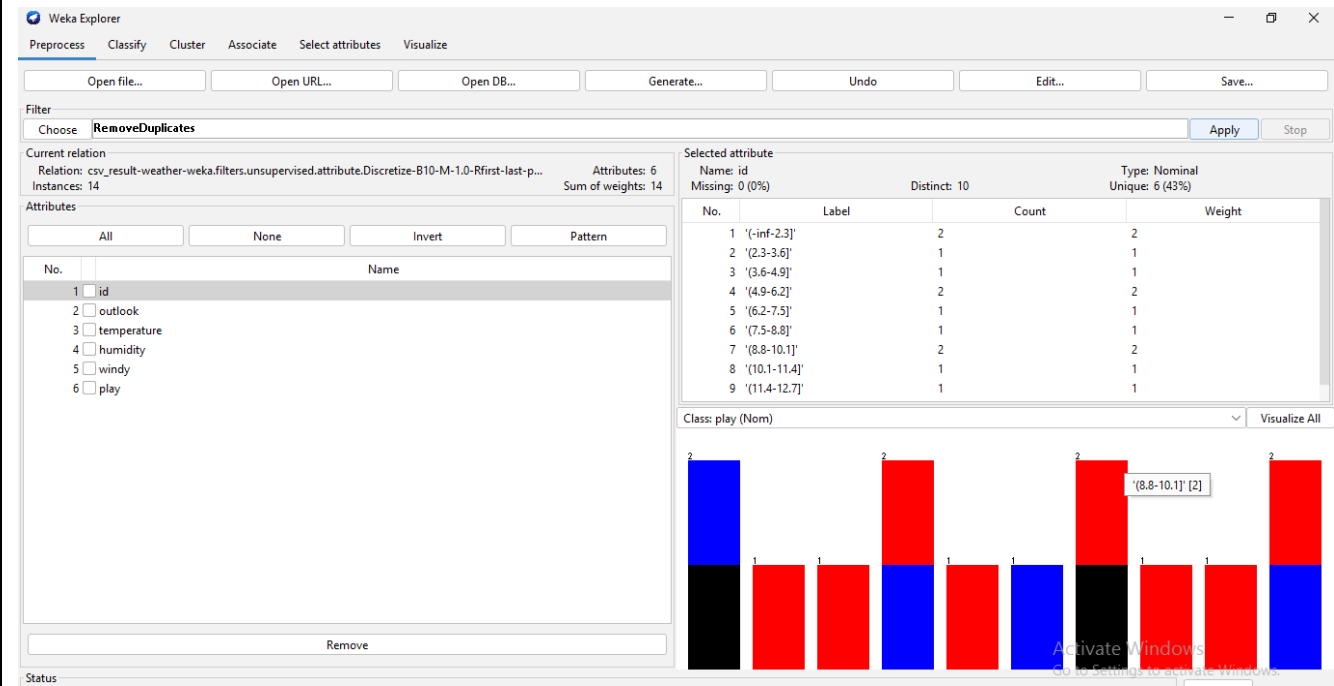


4) NUMERIC TO NOMINAL:

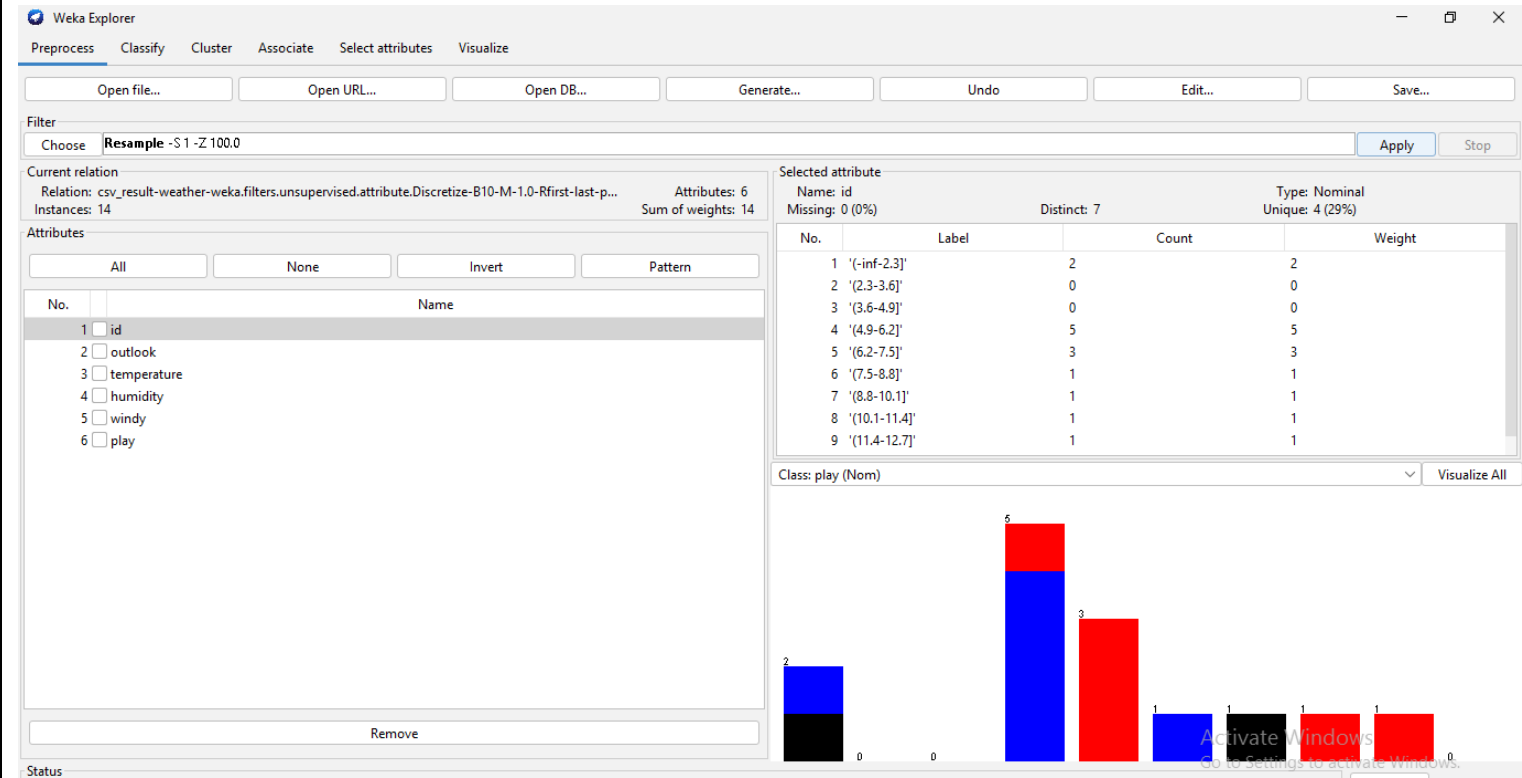
A filter for turning numeric attributes into nominal ones.



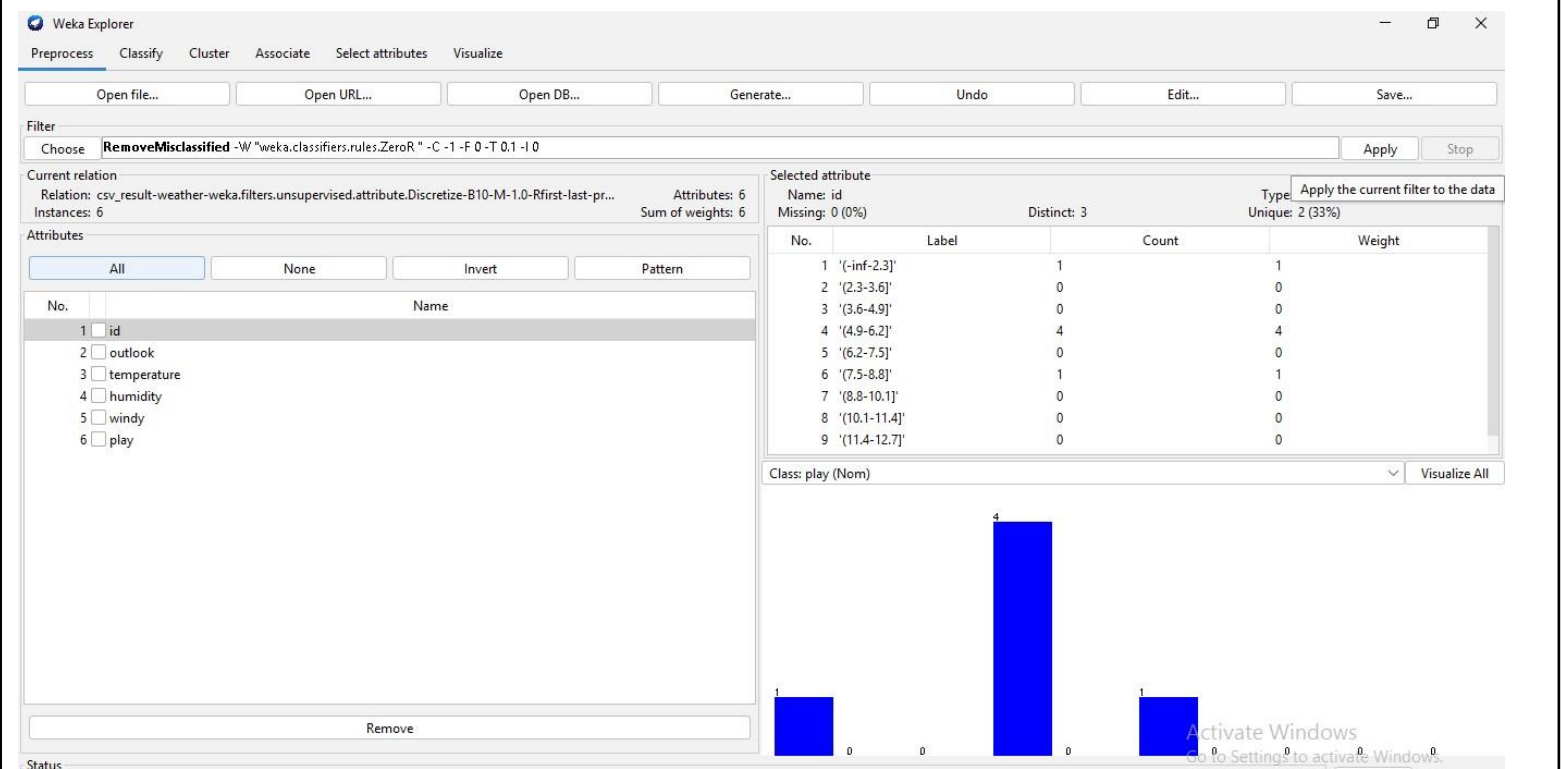
5) REMOVE DUPLICATES:
Removes all duplicate instances from the first batch of data it receives



6) RESAMPLE:
Produces a random subsample of a dataset using either sampling with replacement orwithout replacement



7) REMOVE MISCLASSIFIED:
A filter that removes instances which are incorrectly classified.



PRACTICAL - 5

AIM: To Perform Frequent Pattern Mining using Apriori Algorithm over anydataset

Apriori algorithm refers to the algorithm which is used to calculate the association rules betweenobjects. It means how two or more objects are related to one another. In other words, we can say that the apriori algorithm is an association rule learning that analyzes that people who bought product A also bought product B.

The given three components comprise the apriori algorithm.

- 1. Support
- 2. Confidence
- 3. Lift

Support

Support (Item 1) = (Transactions relating Item 1) / (Total transactions)

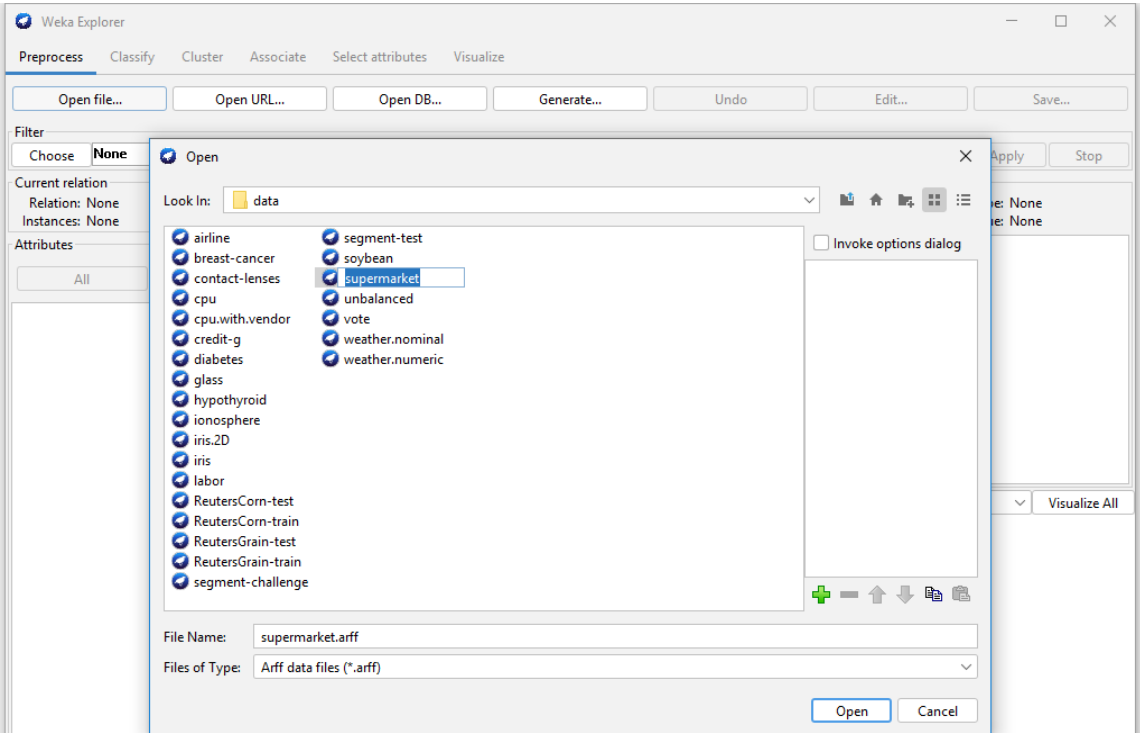
Confidence

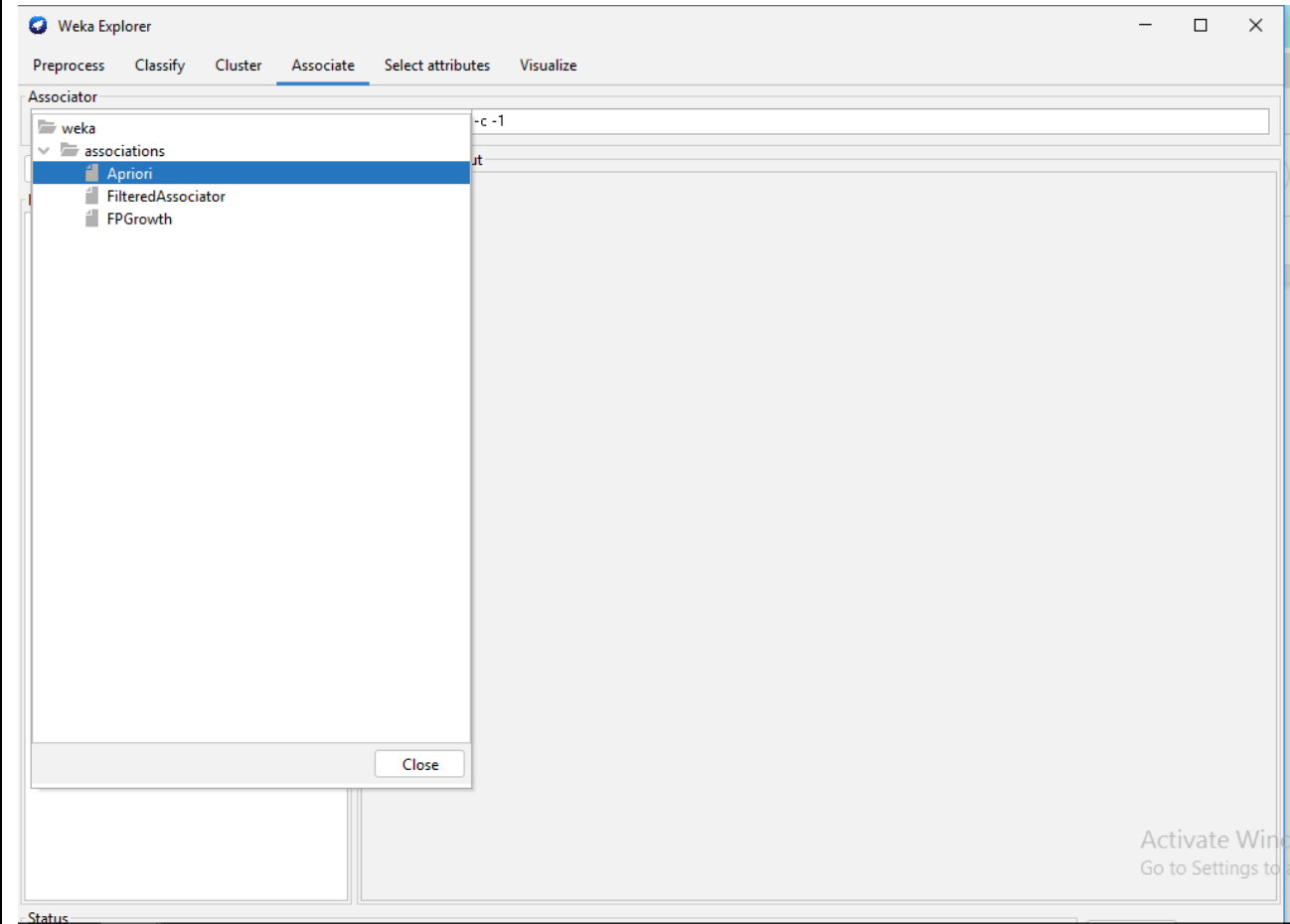
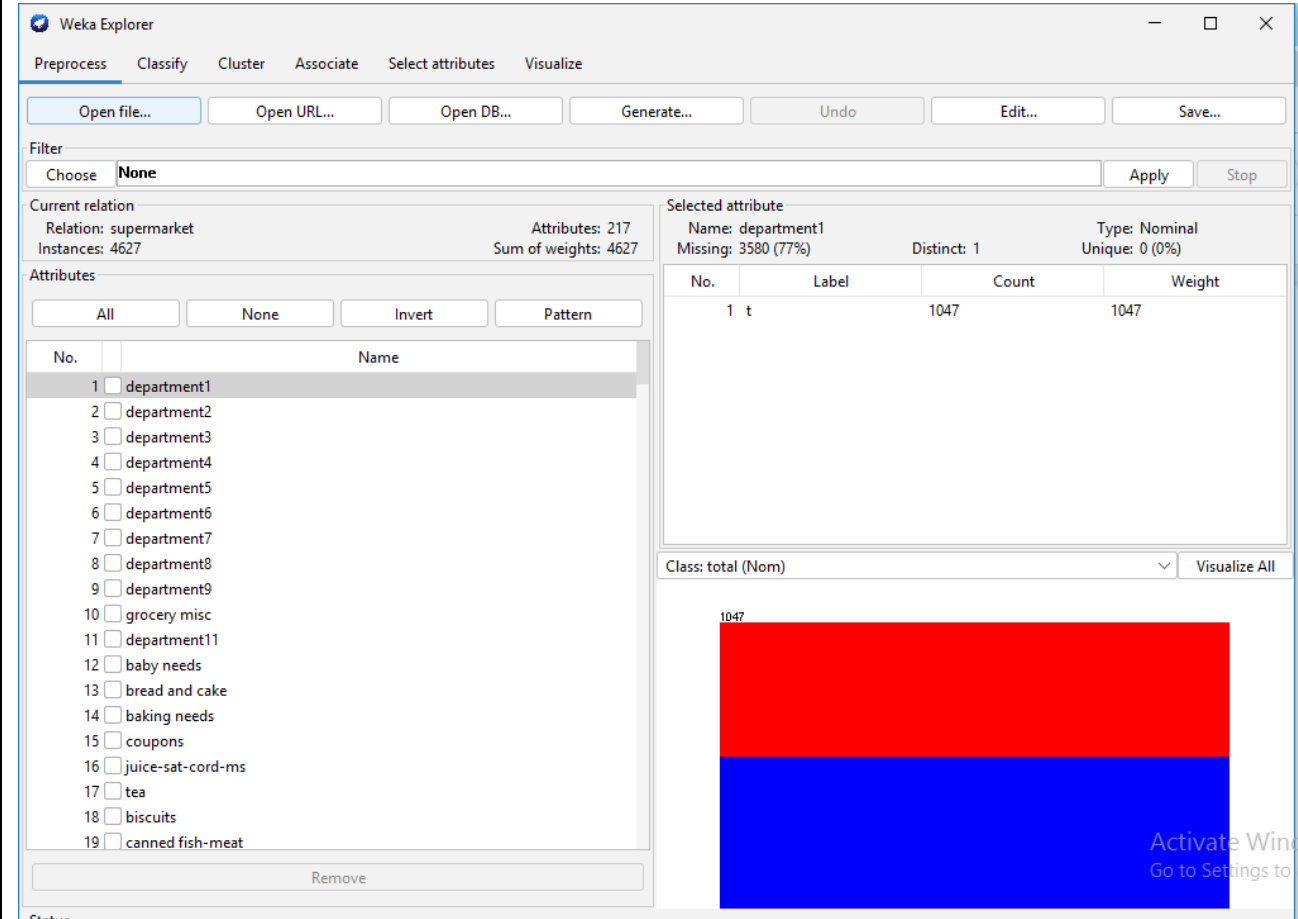
Confidence = (Transactions relating both Item 1 and Item 2) / (Total transactions involving Item1)

Lift

Lift = Confidence/ Support

Uploading Data set





Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Associator

ChooseApriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

StartStop

Result list (right-click for ...)

09:52:16 - Apriori

Associator output

Minimum support: 0.15 (694 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 17

Generated sets of large itemsets:

Size of set of large itemsets L(1): 44
Size of set of large itemsets L(2): 380
Size of set of large itemsets L(3): 910
Size of set of large itemsets L(4): 633
Size of set of large itemsets L(5): 105
Size of set of large itemsets L(6): 1

Best rules found:

1. biscuits=t frozen foods=t fruit=t total=high 788 ==> bread and cake=t 723 <conf:(0.92)> lift:(1.27) lev:(0.03)
2. baking needs=t biscuits=t fruit=t total=high 760 ==> bread and cake=t 696 <conf:(0.92)> lift:(1.27) lev:(0.03)
3. baking needs=t frozen foods=t fruit=t total=high 770 ==> bread and cake=t 705 <conf:(0.92)> lift:(1.27) lev:(0.03)
4. biscuits=t fruit=t vegetables=t total=high 815 ==> bread and cake=t 746 <conf:(0.92)> lift:(1.27) lev:(0.03) [179]
5. party snack foods=t fruit=t total=high 854 ==> bread and cake=t 779 <conf:(0.91)> lift:(1.27) lev:(0.04) [164]
6. biscuits=t frozen foods=t vegetables=t total=high 797 ==> bread and cake=t 725 <conf:(0.91)> lift:(1.26) lev:(0.04) [179]
7. baking needs=t biscuits=t vegetables=t total=high 772 ==> bread and cake=t 701 <conf:(0.91)> lift:(1.26) lev:(0.04) [179]
8. biscuits=t fruit=t total=high 954 ==> bread and cake=t 866 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
9. frozen foods=t fruit=t vegetables=t total=high 834 ==> bread and cake=t 757 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
10. frozen foods=t fruit=t total=high 969 ==> bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)

Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Associator

ChooseApriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

StartStop

Result list (right-click for ...)

09:52:16 - Apriori

Associator output

Minimum support: 0.15 (694 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 17

Generated sets of large itemsets:

Size of set of large itemsets L(1): 44
Size of set of large itemsets L(2): 380
Size of set of large itemsets L(3): 910
Size of set of large itemsets L(4): 633
Size of set of large itemsets L(5): 105
Size of set of large itemsets L(6): 1

Best rules found:

1. biscuits=t frozen foods=t fruit=t total=high 788 ==> bread and cake=t 723 <conf:(0.92)> lift:(1.27) lev:(0.03)
2. baking needs=t biscuits=t fruit=t total=high 760 ==> bread and cake=t 696 <conf:(0.92)> lift:(1.27) lev:(0.03)
3. baking needs=t frozen foods=t fruit=t total=high 770 ==> bread and cake=t 705 <conf:(0.92)> lift:(1.27) lev:(0.03)
4. biscuits=t fruit=t vegetables=t total=high 815 ==> bread and cake=t 746 <conf:(0.92)> lift:(1.27) lev:(0.03) [179]
5. party snack foods=t fruit=t total=high 854 ==> bread and cake=t 779 <conf:(0.91)> lift:(1.27) lev:(0.04) [164]
6. biscuits=t frozen foods=t vegetables=t total=high 797 ==> bread and cake=t 725 <conf:(0.91)> lift:(1.26) lev:(0.04) [179]
7. baking needs=t biscuits=t vegetables=t total=high 772 ==> bread and cake=t 701 <conf:(0.91)> lift:(1.26) lev:(0.04) [179]
8. biscuits=t fruit=t total=high 954 ==> bread and cake=t 866 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
9. frozen foods=t fruit=t vegetables=t total=high 834 ==> bread and cake=t 757 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
10. frozen foods=t fruit=t total=high 969 ==> bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)

weka.gui.GenericObjectEditor

weka.associations.Apriori

About

Class implementing an Apriori-type algorithm.

More

Capabilities

carFalse

classIndex-1

delta0.05

doNotCheckCapabilitiesFalse

lowerBoundMinSupport0.1

metricTypeConfidence

minMetric0.9

numRules10

outputItemSetsFalse

removeAllMissingColsFalse

significanceLevel-1.0

treatZeroAsMissingFalse

upperBoundMinSupport1.0

verboseFalse

Open...Save...OKCancelSettings to

PRACTICAL - 6

AIM: Implementation of any one classifier using JAVA (Bayesianclassifier) and verify results with WEKA.

Naive Bayes classifiers are a collection of classification algorithms based on Bayes’ Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other.

Bayes' Theorem:

- The formula for Bayes' theorem is given as:

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

