

GUDLAVALLERU ENGINEERING COLLEGE

(An Autonomous Institute with Permanent Affiliation to JNTUK, Kakinada)

Seshadri Rao Knowledge Village, Gudlavalleru – 521 356.

Department of Computer Science and Engineering

III Year II Semester

2019 - 2020



DATA MINING-LAB FACULTY MANUAL

Prepared by

Dr. G.V.S.N.R.V.Prasad
Professor & Vice Prinicipal(Academics)

Dr.A.Jagadeswara Rao
Associate Professor

Mr. K. Bala Brahmeswara
Assistant Professor

Mrs. Y. Aditya
Assistant Professor

GUDLAVALLERU ENGINEERING COLLEGE

(An Autonomous Institution with Permanent Affiliation to JNTUK, Kakinada)

Seshadri Rao Knowledge Village, Gudlavalleru – 521356

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

INSTITUTE VISION & MISSION

Institute Vision:

To be a leading institution of engineering education and research, preparing students for

leadership in their fields in a caring and challenging learning environment.

Institute Mission:

- To produce quality engineers by providing state-of-the-art engineering education.
- To attract and retain knowledgeable, creative, motivated and highly skilled individuals whose leadership and contributions uphold the college tenets of education, creativity, research and responsible public service.
- To develop faculty and resources to impart and disseminate knowledge and information to students and also to society that will enhance educational level, which in turn, will contribute to social and economic betterment of society.
- To provide an environment that values and encourages knowledge acquisition and academic freedom, making this a preferred institution for knowledge seekers.
- To provide quality assurance.
- To partner and collaborate with industry, government, and R&D institutes to develop new knowledge and sustainable technologies and serve as an engine for facilitating the nation's economic development.
- To impart personality development skills to students that will help them to succeed and lead.
- To instil in students the attitude, values and vision that will prepare them to lead lives of personal integrity and civic responsibility.
- To promote a campus environment that welcomes and makes students of all races, cultures and civilizations feel at home.
- Putting students face to face with industrial, governmental and societal challenges.

DEPARTMENT VISION & MISSION

Vision

To be a Centre of Excellence in Computer Science and Engineering education and training to meet the challenging needs of the industry and society.

Mission

- To impart quality education through well-designed curriculum in tune with the growing software needs of the industry.
- To serve our students by inculcating in them problem solving, leadership, teamwork skills and the value of commitment to quality, ethical behavior & respect for others.
- To foster industry-academia relationship for mutual benefit and growth.

PROGRAMME EDUCATIONAL OBJECTIVES (PEOs)

PEO 1: Identify, analyze, formulate and solve computing problems both independently and in a team environment using appropriate modern tools.

PEO 2: Develop software systems with significant technical, legal, ethical, social, environmental and economic considerations.

PEO 3: Exhibit commitment in lifelong learning, professional development and leadership and communicate effectively with professional clients and the public.

PROGRAM SPECIFIC OUTCOMES (PSOs)

1. Design, develop, test and maintain reliable software systems and intelligent systems.
2. Design and develop websites, web apps and mobile apps.

Program Outcomes (Pos)

(B.Tech in Computer Science and Engineering)

Engineering students will be able to:

1. **Engineering knowledge:** Apply the knowledge of mathematics, science, engineering fundamentals, and an engineering specialization to the solution of complex engineering problems.
2. **Problem analysis:** Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural sciences, and engineering sciences.
3. **Design/development of solutions:** Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal, and environmental considerations.
4. **Conduct investigations of complex problems:** Use research-based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions.
5. **Modern tool usage:** Create, select, and apply appropriate techniques, resources, and modern engineering and IT tools including prediction and modelling to complex engineering activities with an understanding of the limitations.
6. **The engineer and society:** Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice.
7. **Environment and sustainability:** Understand the impact of the professional engineering solutions in societal and environmental

- contexts, and demonstrate the knowledge of, and need for sustainable development.
8. **Ethics:** Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice.
 9. **Individual and team work:** Function effectively as an individual, and as a member or leader in diverse teams, and in multidisciplinary settings.
 10. **Communication:** Communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions.
 11. **Project management and finance:** Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments.
 12. **Life-long learning:** Recognize the need for, and have the preparation and ability to engage in independent and life-long learning in the broadest context of technological change.

Department of Computer Science and Engineering

GENERAL LABORATORY INSTRUCTIONS

1. Students are advised to come to the laboratory at least 5 minutes before (to the starting time), those who come after 5 minutes will not be allowed into the lab.
2. Plan your task properly much before to the commencement, come prepared to the lab with the synopsis / program / experiment details.
3. Student should enter into the laboratory with:
 - a. Laboratory observation notes with all the details (Problem statement, Aim, Algorithm, Procedure, Program, Expected Output, etc.,) filled in for the lab session.
 - b. Laboratory Record updated up to the last session experiments and other utensils (if any) needed in the lab.
 - c. Proper Dress code and Identity card.
4. Sign in the laboratory login register, write the TIME-IN, and occupy the computer system allotted to you by the faculty.
5. Execute your task in the laboratory, and record the results / output in the lab observation note book, and get certified by the concerned faculty.
6. All the students should be polite and cooperative with the laboratory staff, must maintain the discipline and decency in the laboratory.
7. Computer labs are established with sophisticated and high end branded systems, which should be utilized properly.
8. Students / Faculty must keep their mobile phones in SWITCHED OFF mode during the lab sessions. Misuse of the equipment, misbehaviors with the staff and systems etc., will attract severe punishment.
9. Students must take the permission of the faculty in case of any urgency to go out ; if anybody found loitering outside the lab / class without permission during working hours will be treated seriously and punished appropriately.

10. Students should LOG OFF/ SHUT DOWN the computer system before he/she leaves the lab after completing the task (experiment) in all aspects. He/she must ensure the system / seat is kept properly.

COURSE NAME: DATA MINING LAB

COURSE CODE: CS2514

COURSE OBJECTIVES:

To exercise the data mining techniques such as classification, clustering, pattern mining etc with different datasets and dynamic parameters using WEKA tool.

COURSE OUTCOMES:

Upon successful completion of the course, the students will be able to

- Learn to execute data mining tasks using a data mining toolkit (such as WEKA) and visualize the results.
- Demonstrate the working of algorithms for data mining tasks such as association rule mining, classification, clustering and regression.

MAPPING OF COURSE OUTCOMES WITH PROGRAM OUTCOMES:

COURSE OUTCOMES	PO1	PO2	PO3	PO4	PO5	PO6	PO7	PO8	PO9	P10	P11	P12
CO1	√	√		√	√							√
CO2	√	√			√							√

DATAMINING LAB INDEX

SNO	Experiment Name	Page No
1	Explore WEKA Data Mining/Machine Learning Toolkit	11
2	Perform data preprocessing tasks on i. Add attribute ii. Add expression iii. Copy attribute iv. Remove attribute	31
3	Demonstrate performing classification on data sets	49
4	Demonstrate performing association rule mining on data sets	62
5	Demonstrate performing regression on data sets	78
6	Demonstrate performing SVM classification on data sets	88
7	Demonstrate performing clustering on data sets	96
8	Demonstrate performing knowledge flow on WEKA	113

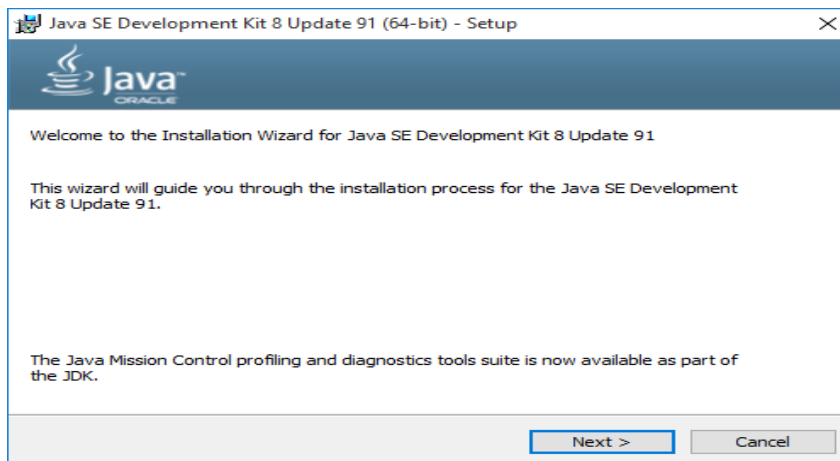
ADDITIONAL -EXPERIMENT LIST

SNO	Experiment Name	Page No
1	FILE FORMATS FOR WEKA	154
2	Perform Data Preprocessing tasks on Data Cleaning & Noisy Data.	156
3	Demonstrate performing Hierarchical clustering on data sets	161

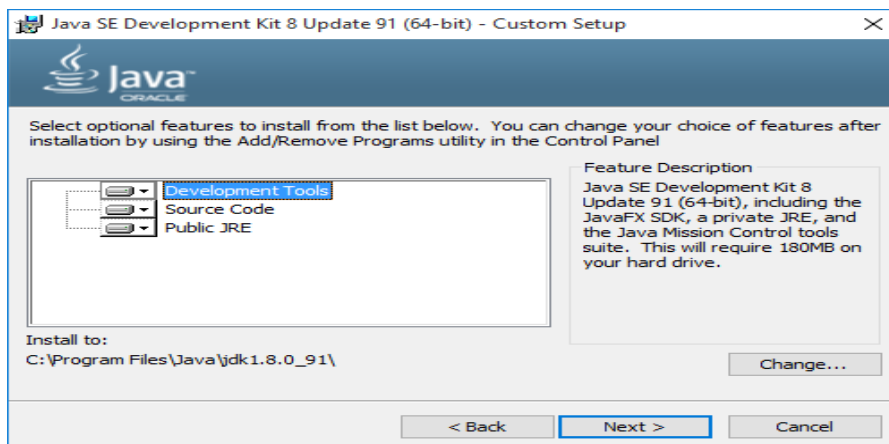
AIM:**i) Downloading and/or installation of WEKA data mining toolkit**

Note: if weka installation package don't have java included in it then compatible version of the java has to be installed initially as given below.

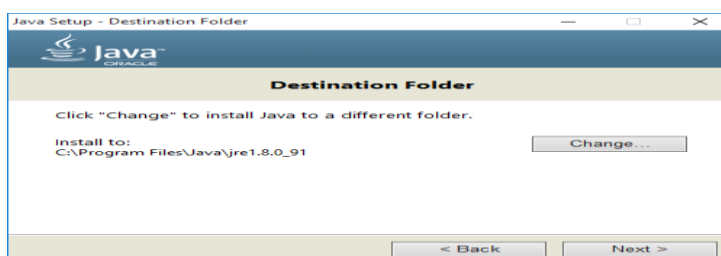
Step 1: Select appropriate JAVA version (Here we are using:jdk1.8). Operate the executable file of welcome window to the installation wizard of Development kit ,click next button



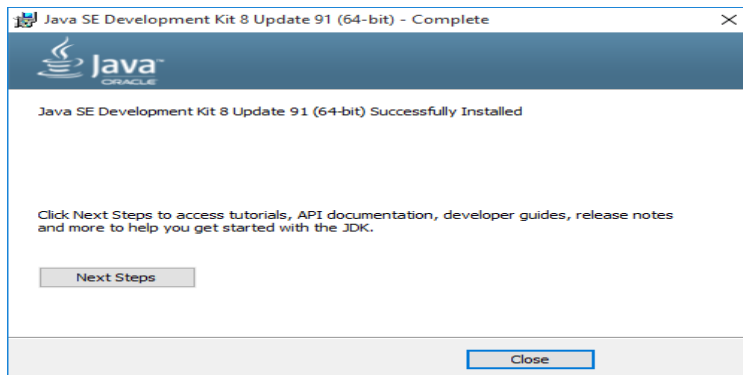
Step 2: Custom setup window to select path and click next.



Step 3: In the destination folder window we check java file stored location and click next.



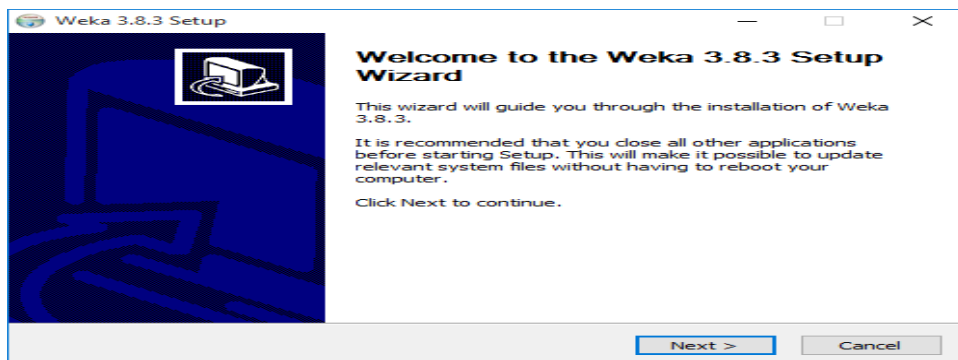
Step 4: Installation process completion window is displayed and click close.



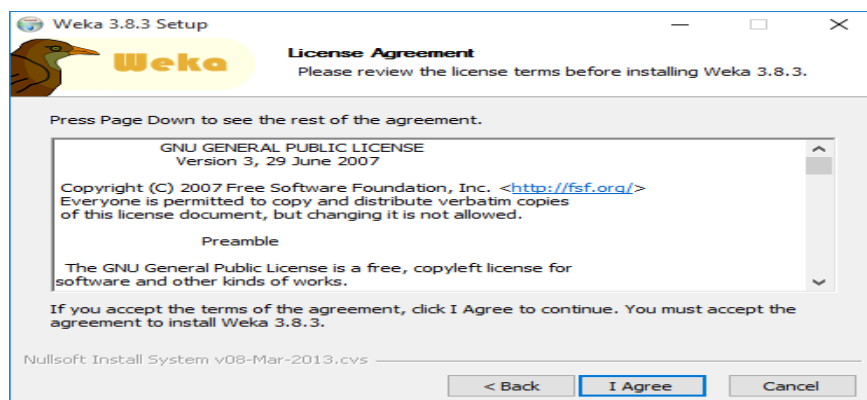
Weka –Installation:

Go to the Weka website, <http://www.cs.waikato.ac.nz/ml/weka/>, and download the software. On the left hand side, click on the link that says download. Select the appropriate link corresponding to the version of the software based on your operating system and whether or not you already have Java VM running on your machine. Save the self extracting executable to the disk and double-click on it to install weka.

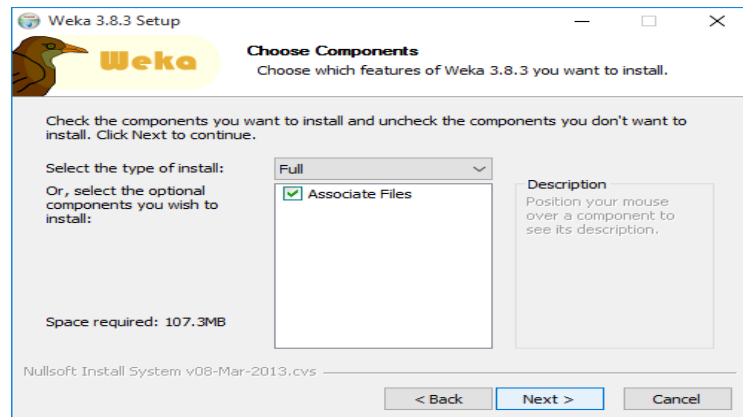
Step 1: Welcome setup wizard window is viewed where next button has to be clicked



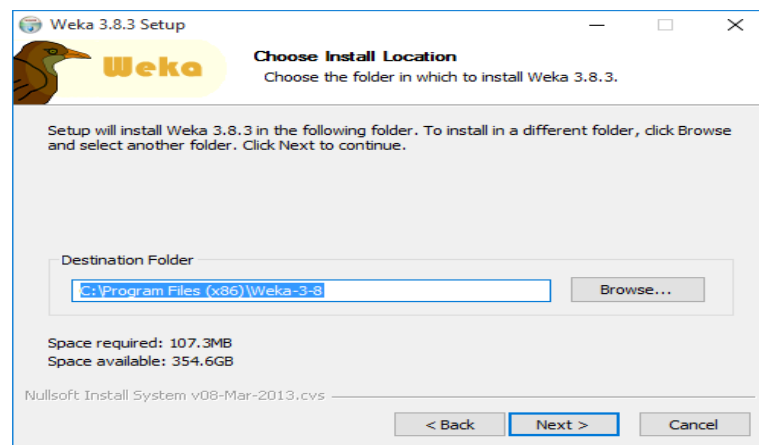
Step2: In the license agreement window accept the agreement with “I Agree” button.



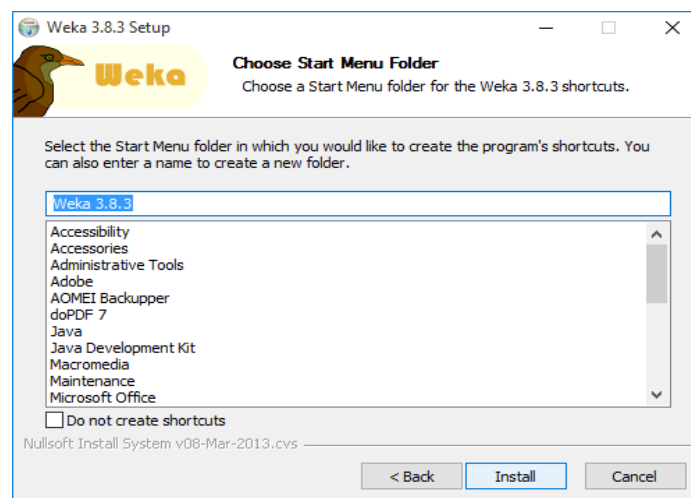
Step 3: The component window is displayed where the type of installation must be “Full” and click on next.



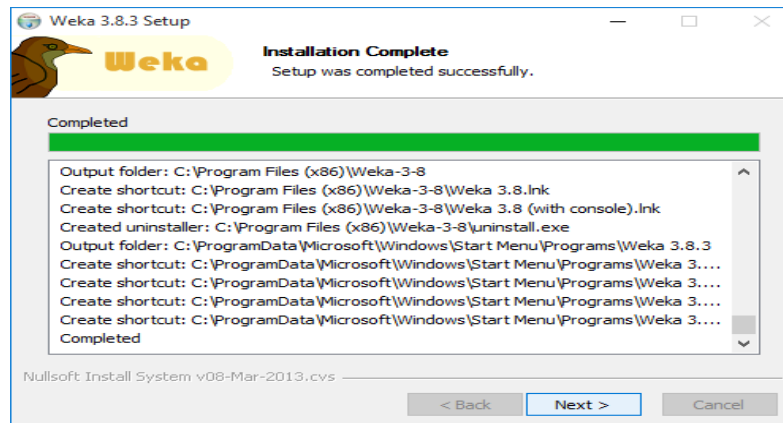
Step 4: In the installation window to select the path where the software loads can be selected through Browse button to the destination folder and click next button.



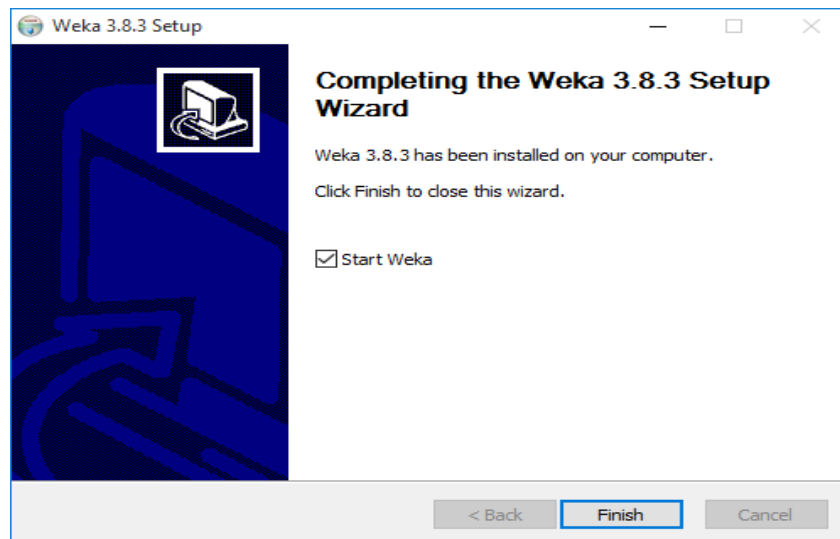
Step 5: To add to the list of programs from the start menu check the create shortcut and press the Install button



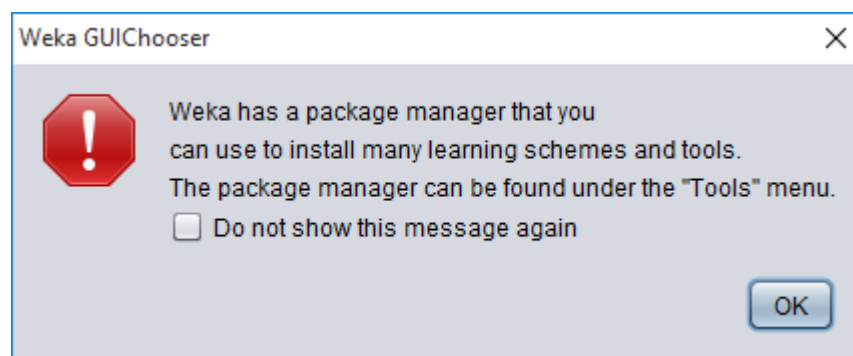
Step 6: completion of installation is indicated in the window where the next button is selected.



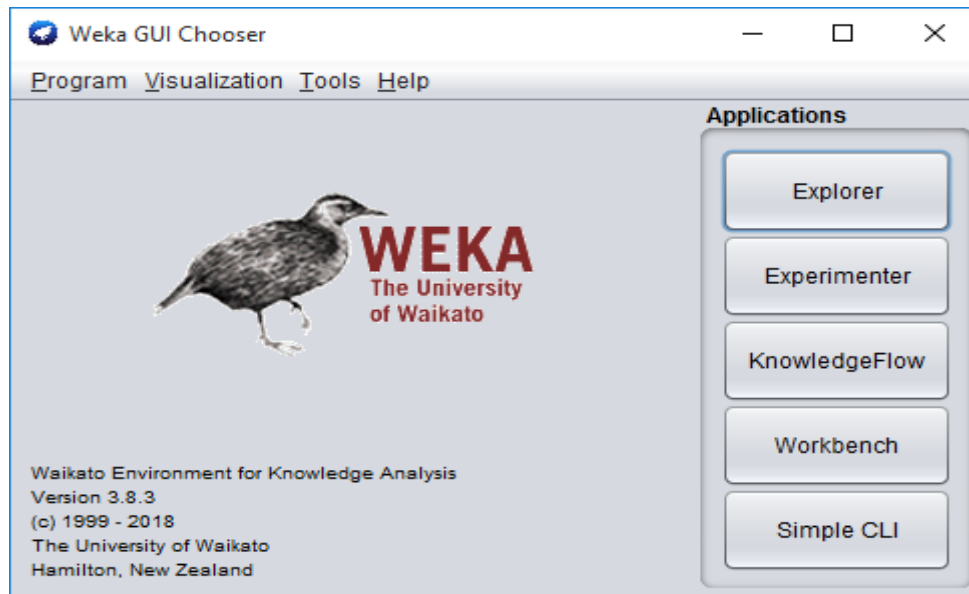
Step 7: The setup wizard with the start up of Weka is chosen and finish the process.



Step 8: An warning window tells about the feature of package manager, click ok.



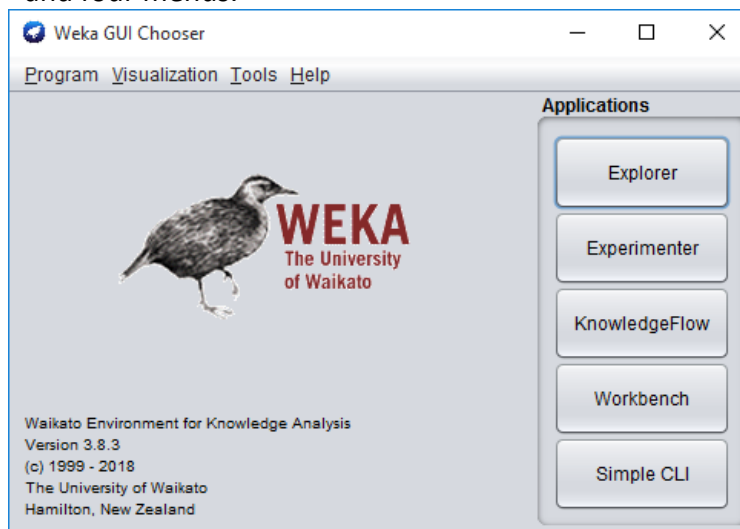
Step 9: Weka home page is viewed.



Step 10: The installation process is completed fully.

2. Understand the feature of Weka toolkit

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.
- **KnowledgeFlow** 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.
- **SimpleCLI** 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:

1. Program

Program	
LogWindow	Ctrl+L
Memory usage	Ctrl+M
Exit	Ctrl+E

- **LogWindow** Opens a log window that captures all that is printed to stdout or stderr. Useful for environments like MS Windows, where WEKA is normally not started from a terminal.
- **Exit** Closes WEKA.

2. Visualization

Visualization	
Plot	Ctrl+P
ROC	Ctrl+R
TreeVisualizer	Ctrl+T
GraphVisualizer	Ctrl+G
BoundaryVisualizer	Ctrl+B

- **Plot** For plotting a 2D plot of a dataset.
- **ROC** Displays a previously saved ROC curve.
- **TreeVisualizer** For displaying directed graphs, e.g., a decision tree.
- **GraphVisualizer** Visualizes XML BIF or DOT format graphs, e.g., for Bayesian networks.
- **BoundaryVisualizer** Allows the visualization of classifier decision boundaries in two dimensions.

3. Tools

Tools	
ArffViewer	Ctrl+A
SqlViewer	Ctrl+S
Bayes net editor	Ctrl+N

- **ArffViewer** An MDI application for viewing ARFF files in spread- sheet format.
- **SqlViewer** Represents an SQL worksheet, for querying databases via JDBC.
- **Bayes net editor** An application for editing, visualizing and learn- ing Bayes nets.

4. Help

Help	
Weka homepage	Ctrl+H
HOWTOs, code snippets, etc.	Ctrl+W
Weka on Sourceforge	Ctrl+F
SystemInfo	Ctrl+I

- **Weka homepage** Opens a browser window with WEKA's home-page.
- **HOWTOs, code snippets, etc.** The general WekaWiki [2], con-taining lots of examples and HOWTOs around the development and use of WEKA.
- **Weka on Sourceforge** WEKA's project homepage on Sourceforge.net.
- **SystemInfo** Lists some internals about the Java/WEKA environment, e.g., the CLASSPATH.

3. Navigate the options available in the WEKA (ex. Select attributes panel, preprocess panel, classify panel, cluster panel, associate panel and visualize panel).

Explorer

1. The user interface

1.1 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 greyed 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.

1.2 Status Box

The status box appears at the very bottom of the window. It displays messages that keep you informed about what's going on. For example, if the Explorer is busy loading a file, the status box will say that.

- **Memory information.** Display in the log box the amount of memory available to WEKA.
- **Run garbage collector.** Force the Java garbage collector to search for memory that is no longer needed and free it up, allowing more memory for new tasks. Note that the garbage collector is constantly running as a background task anyway.

1.3 Log Button

Clicking on this button brings up a separate window containing a scrollable text field. Each line of text is stamped with the time it was entered into the log. As you perform actions in WEKA, the log keeps a record of what has happened.

1.4 WEKA Status Icon

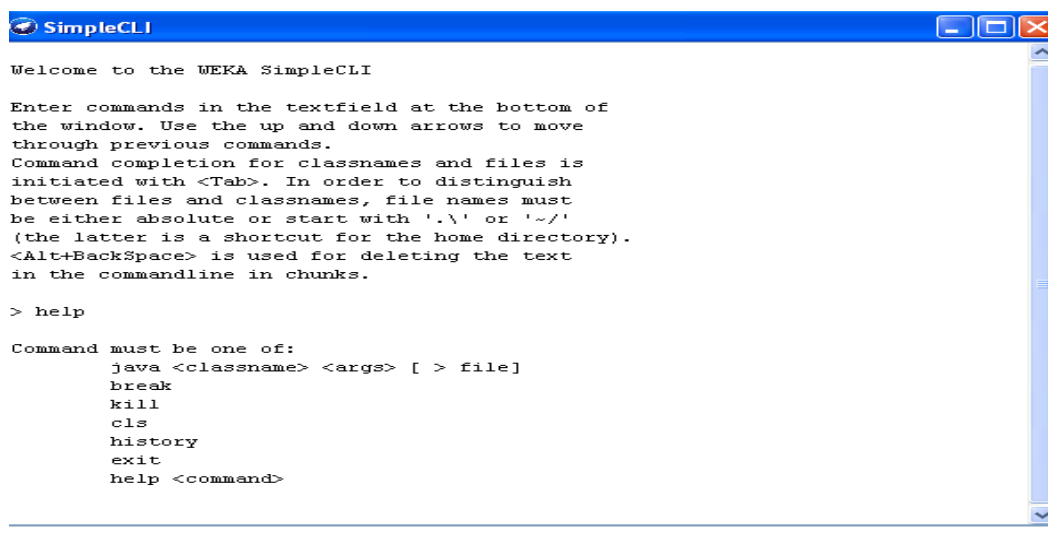
To the right of the status box is the WEKA status icon. When no processes are running, the bird sits down and takes a nap. The number beside the × symbol gives the number of concurrent processes running.

1.5 Graphical output

Most graphical displays in WEKA, e.g., the Graph Visualizer or the Tree Visualizer, support saving the output to a file. A dialog for saving the output can be brought up with Alt+Shift+left-click. Supported formats are currently Windows Bitmap, JPEG, PNG and EPS (encapsulated Postscript). The dialog also allows you to specify the dimensions of the generated image.

Simple CLI

The Simple CLI provides full access to all Weka classes, i.e., classifiers, filters, clusterers, etc., but without the hassle of the CLASSPATH (it facilitates the one, with which Weka was started). It offers a simple Weka shell with separated command line and output.



```

SimpleCLI

Welcome to the WEKA SimpleCLI

Enter commands in the textfield at the bottom of
the window. Use the up and down arrows to move
through previous commands.
Command completion for classnames and files is
initiated with <Tab>. In order to distinguish
between files and classnames, file names must
be either absolute or start with './' or '~/ '
(the latter is a shortcut for the home directory).
<Alt+BackSpace> is used for deleting the text
in the commandline in chunks.

> help

Command must be one of:
  java <classname> <args> [ > file]
  break
  kill
  cls
  history
  exit
  help <command>

```

Commands

The following commands are available in the Simple CLI:

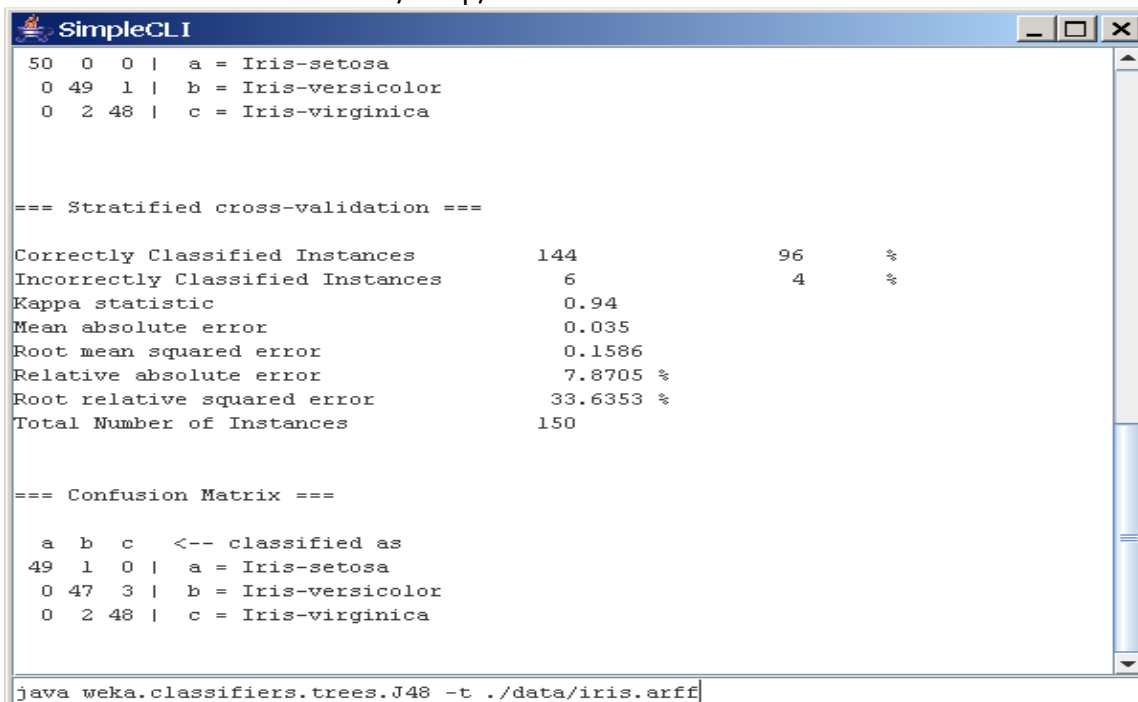
- **java <classname> [<args>]**
invokes a java class with the given arguments (if any)

- **break** stops the current thread, e.g., a running classifier, in a friendly manner.
- **Kill** stops the current thread in an unfriendly fashion
- **cls** clears the output area
- **exit** exits the Simple CLI
- **help** [<command>] provides an overview of the available commands if without a command name as argument, otherwise more help on the specified command.

Invocation

In order to invoke a Weka class, one has only to prefix the class with "java". This command tells the Simple CLI to load a class and execute it with any given parameters. E.g., the J48 classifier can be invoked on the iris dataset with the following command:

```
java weka.classifiers.trees.J48 -t c:/temp/iris.arff
```



```

SimpleCLI
50 0 0 | a = Iris-setosa
0 49 1 | b = Iris-versicolor
0 2 48 | c = Iris-virginica

=== Stratified cross-validation ===

Correctly Classified Instances      144           96   %
Incorrectly Classified Instances     6            4   %
Kappa statistic                    0.94
Mean absolute error                 0.035
Root mean squared error             0.1586
Relative absolute error             7.8705 %
Root relative squared error        33.6353 %
Total Number of Instances          150

=== Confusion Matrix ===

 a b c  <-- classified as
49 1 0 | a = Iris-setosa
0 47 3 | b = Iris-versicolor
0 2 48 | c = Iris-virginica

java weka.classifiers.trees.J48 -t ./data/iris.arff

```

Command redirection

Starting with this version of Weka one can perform a basic redirection:

```
java weka.classifiers.trees.J48 -t test.arff > j48.txt
```

Note: the > must be preceded and followed by a space, otherwise it is not recognized as redirection, but part of another parameter.

WEEK 2:

- iv) Study the arff file format
- v) Explore the available data sets in WEKA.
- vi) Load a data set (ex. Weather dataset, Iris dataset, etc.)
- vii) Load each dataset and observe the following:

- a) List the attribute names and their types
 - b) Number of records in each dataset
 - c) Identify the class attribute (if any)
 - d) Plot histogram
 - e) Determine the number of records for each class.
 - f) Visualize the data in various dimensions
- iv) Study the arff file format

Attribute-Relation File Format (ARFF)

An ARFF (Attribute-Relation File Format) file is an ASCII text file that describes a list of instances sharing a set of attributes. ARFF files were developed by the Machine Learning Project at the Department of Computer Science of The University of Waikato for use with the Weka machine learning software. This document describes the version of ARFF used with Weka versions 3.2 to 3.3; this is an extension of the ARFF format as described in the data mining book written by Ian H. Witten and Eibe Frank (the new additions are string attributes, date attributes, and sparse instances).

This explanation was cobbled together by Gordon Paynter (gordon.paynter at ucr.edu) from the Weka 2.1 ARFF description, email from Len Trigg (lenbok at myrealbox.com) and Eibe Frank (eibe at cs.waikato.ac.nz), and some datasets. It has been edited by Richard Kirkby (rkirkby at cs.waikato.ac.nz). Contact Len if you're interested in seeing the ARFF 3 proposal.

Overview

ARFF files have two distinct sections. The first section is the **Header** information, which is followed by the **Data** information.

The **Header** of the ARFF file contains the name of the relation, a list of the attributes (the columns in the data), and their types. An example header on the standard IRIS dataset looks like this:

```
% 1. Title: Iris Plants Database
%
% 2. Sources:
%   (a) Creator: R.A. Fisher
%   (b) Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)
%   (c) Date: July, 1988
%
@RELATION iris
```

```

@ATTRIBUTE sepallength NUMERIC
@ATTRIBUTE sepalwidth NUMERIC
@ATTRIBUTE petallength NUMERIC
@ATTRIBUTE petalwidth NUMERIC
@ATTRIBUTE class {Iris-setosa,Iris-versicolor,Iris-virginica}

```

The **Data** of the ARFF file looks like the following:

```

@DATA
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
5.0,3.6,1.4,0.2,Iris-setosa
5.4,3.9,1.7,0.4,Iris-setosa
4.6,3.4,1.4,0.3,Iris-setosa
5.0,3.4,1.5,0.2,Iris-setosa
4.4,2.9,1.4,0.2,Iris-setosa
4.9,3.1,1.5,0.1,Iris-setosa

```

Lines that begin with a % are comments. The **@RELATION**, **@ATTRIBUTE** and **@DATA** declarations are case insensitive.

Examples

Several well-known machine learning datasets are distributed with Weka in the \$WEKAHOME/data directory as ARFF files.

The ARFF Header Section

The ARFF Header section of the file contains the relation declaration and attribute declarations.

The @relation Declaration

The relation name is defined as the first line in the ARFF file. The format is:

```
@relation <relation-name>
```

where <relation-name> is a string. The string must be quoted if the name includes spaces.

The @attribute Declarations

Attribute declarations take the form of an ordered sequence of **@attribute** statements. Each attribute in the data set has its own **@attribute** statement which uniquely defines the name of that attribute and its data type. The order the attributes are declared indicates the column position in the data section of the file. For example, if an attribute is the third one declared then Weka expects that all that attribute's values will be found in the third comma delimited column.

The format for the **@attribute** statement is:

```
@attribute <attribute-name> <datatype>
```

where the *<attribute-name>* must start with an alphabetic character. If spaces are to be included in the name then the entire name must be quoted.

The *<datatype>* can be any of the four types currently (version 3.2.1) supported by Weka:

- numeric
- <nominal-specification>
- string
- date [<date-format>]

where *<nominal-specification>* and *<date-format>* are defined below. The keywords **numeric**, **string** and **date** are case insensitive.

Numeric attributes

Numeric attributes can be real or integer numbers.

Nominal attributes

Nominal values are defined by providing an *<nominal-specification>* listing the possible values: {<nominal-name1>, <nominal-name2>, <nominal-name3>, ... }

For example, the class value of the Iris dataset can be defined as follows:

```
@ATTRIBUTE class {Iris-setosa,Iris-versicolor,Iris-virginica}
```

Values that contain spaces must be quoted.

String attributes

String attributes allow us to create attributes containing arbitrary textual values. This is very useful in text-mining applications, as we can create datasets with string attributes, then write

Weka Filters to manipulate strings (like StringToWordVectorFilter). String attributes are declared as follows:

```
@ATTRIBUTE LCC string
```

Date attributes

Date attribute declarations take the form:

```
@attribute <name> date [<date-format>]
```

where <name> is the name for the attribute and <date-format> is an optional string specifying how date values should be parsed and printed (this is the same format used by SimpleDateFormat). The default format string accepts the ISO-8601 combined date and time format: "yyyy-MM-dd'T'HH:mm:ss".

Dates must be specified in the data section as the corresponding string representations of the date/time (see example below).

ARFF Data Section

The ARFF Data section of the file contains the data declaration line and the actual instance lines.

The @data Declaration

The **@data** declaration is a single line denoting the start of the data segment in the file. The format is:

```
@data
```

The instance data

Each instance is represented on a single line, with carriage returns denoting the end of the instance.

Attribute values for each instance are delimited by commas. They must appear in the order that they were declared in the header section (i.e. the data corresponding to the nth **@attribute** declaration is always the nth field of the attribute).

Missing values are represented by a single question mark, as in:

```
@data
4.4,?,1.5,?,Iris-setosa
```

Values of string and nominal attributes are case sensitive, and any that contain space must be quoted, as follows:

```
@relation LCCvsLCSH
```

```
@attribute LCC string
@attribute LCSH string
```

```
@data
AG5, 'Encyclopedias and dictionaries.;Twentieth century.'
AS262, 'Science -- Soviet Union -- History.'
AE5, 'Encyclopedias and dictionaries.'
AS281, 'Astronomy, Assyro-Babylonian.;Moon -- Phases.'
AS281, 'Astronomy, Assyro-Babylonian.;Moon -- Tables.'
```

Dates must be specified in the data section using the string representation specified in the attribute declaration. For example:

```
@RELATION Timestamps

@ATTRIBUTE timestamp DATE "yyyy-MM-dd HH:mm:ss"

@DATA
"2001-04-03 12:12:12"
"2001-05-03 12:59:55"
```

Sparse ARFF files

Sparse ARFF files are very similar to ARFF files, but data with value 0 are not be explicitly represented.

Sparse ARFF files have the same header (i.e **@relation** and **@attribute** tags) but the data section is different. Instead of representing each value in order, like this:

```
@data
0, X, 0, Y, "class A"
```


0, 0, W, 0, "class B"

the non-zero attributes are explicitly identified by attribute number and their value stated, like this:

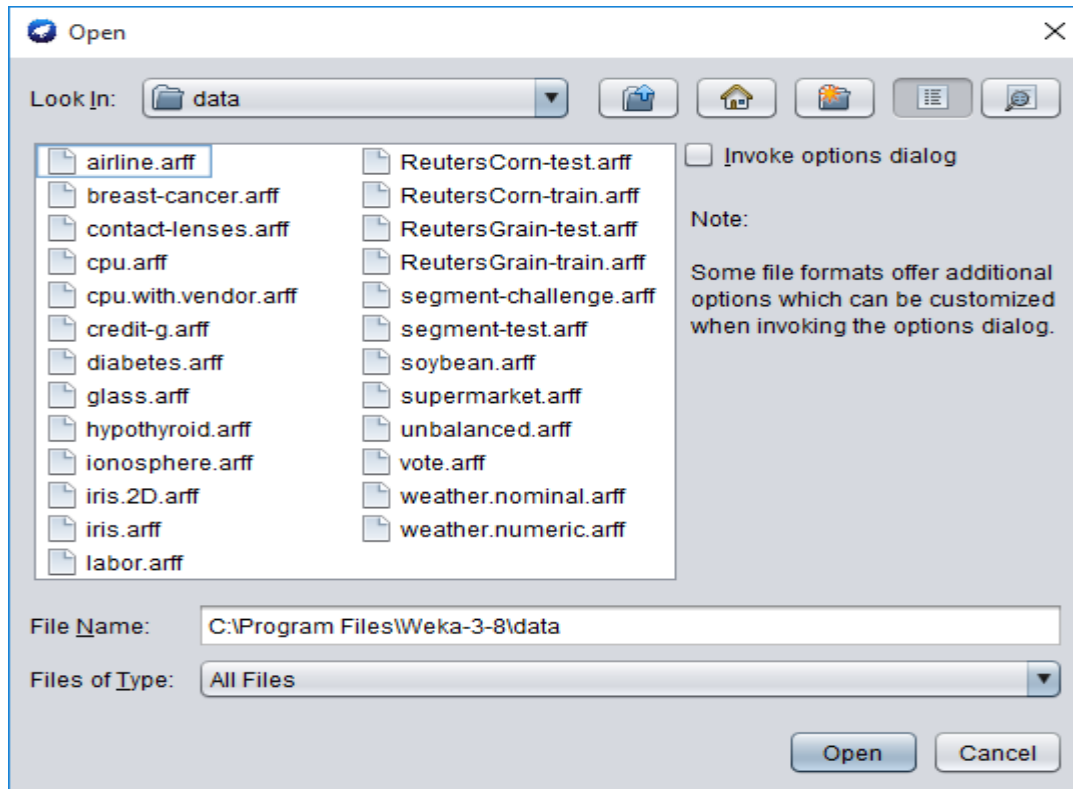
```
@data
{1 X, 3 Y, 4 "class A"}
{2 W, 4 "class B"}
```

Each instance is surrounded by curly braces, and the format for each entry is: <index> <space> <value> where index is the attribute index (starting from 0).

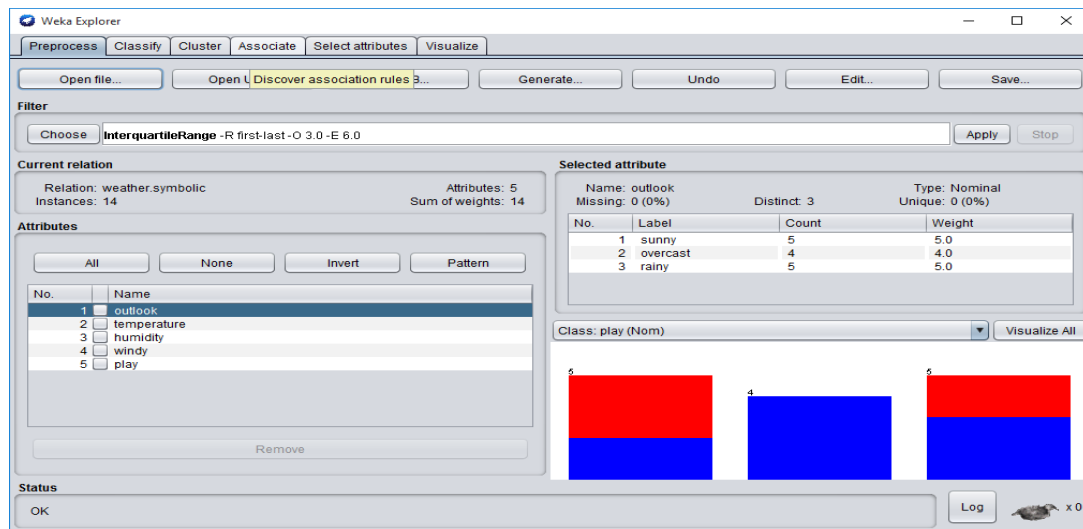
Note that the omitted values in a sparse instance are **0**, they are not "missing" values! If a value is unknown, you must explicitly represent it with a question mark (?).

Warning: There is a known problem saving SparseInstance objects from datasets that have string attributes. In Weka, string and nominal data values are stored as numbers; these numbers act as indexes into an array of possible attribute values (this is very efficient). However, the first string value is assigned index 0: this means that, internally, this value is stored as a 0. When a SparseInstance is written, string instances with internal value 0 are not output, so their string value is lost (and when the arff file is read again, the default value 0 is the index of a different string value, so the attribute value appears to change). To get around this problem, add a dummy string value at index 0 that is never used whenever you declare string attributes that are likely to be used in SparseInstance objects and saved as Sparse ARFF files.

v) Explore the available data sets in WEKA.



vi) Load a data set (ex. Weather dataset, Iris dataset, etc.)



Viewer

Relation: weather.symbolic

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

Add instance Undo OK Cancel

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

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

Filter: Choose Apply Stop

Current relation: Relation: Iris Instances: 150 Attributes: 5 Sum of weights: 150

Selected attribute: Name: sepalength Missing: 0 (0%) Distinct: 35 Type: Numeric Unique: 9 (6%)

Statistic	Value
Minimum	4.3
Maximum	7.9
Mean	5.843
StdDev	0.828

Class: class (Nom) Visualize All

Status: OK Log x0

Viewer

Relation: Iris

No.	1: sepalength	2: sepalwidth	3: petalength	4: petalwidth	5: class
	Numeric	Numeric	Numeric	Numeric	Nominal
1	5.1	3.5	1.4	0.1	Iris-S...
2	4.9	3.0	1.4	0.1	Iris-S...
3	4.7	3.2	1.3	0.1	Iris-S...
4	4.6	3.1	1.5	0.2	Iris-S...
5	5.0	3.6	1.6	0.2	Iris-S...
6	4.4	3.0	1.4	0.1	Iris-S...
7	4.6	3.4	1.4	0.1	Iris-S...
8	5.4	4.4	1.7	0.3	Iris-S...
9	4.9	3.7	1.4	0.2	Iris-S...
10	5.4	4.1	1.5	0.2	Iris-S...
11	4.9	3.0	1.4	0.1	Iris-S...
12	4.8	3.0	1.4	0.1	Iris-S...
13	4.8	3.0	1.4	0.1	Iris-S...
14	4.3	3.0	1.3	0.1	Iris-S...
15	5.8	4.0	1.6	0.3	Iris-S...
16	5.1	3.5	1.4	0.2	Iris-S...
17	4.4	3.0	1.4	0.1	Iris-S...
18	5.7	4.4	1.7	0.3	Iris-S...
19	5.1	3.5	1.4	0.2	Iris-S...
20	4.1	3.0	1.3	0.1	Iris-S...
21	4.4	3.0	1.4	0.1	Iris-S...
22	5.0	3.6	1.6	0.2	Iris-S...
23	4.4	3.0	1.4	0.1	Iris-S...
24	5.0	3.6	1.6	0.2	Iris-S...
25	4.4	3.0	1.4	0.1	Iris-S...
26	5.0	3.6	1.6	0.2	Iris-S...
27	4.9	3.5	1.5	0.2	Iris-S...
28	5.0	3.6	1.6	0.2	Iris-S...
29	5.4	4.4	1.7	0.3	Iris-S...
30	4.1	3.0	1.3	0.1	Iris-S...
31	5.1	3.5	1.4	0.2	Iris-S...
32	4.9	3.0	1.4	0.1	Iris-S...
33	5.4	4.4	1.7	0.3	Iris-S...
34	4.4	3.0	1.4	0.1	Iris-S...
35	5.0	3.6	1.6	0.2	Iris-S...
36	4.9	3.5	1.5	0.2	Iris-S...
37	5.0	3.6	1.6	0.2	Iris-S...

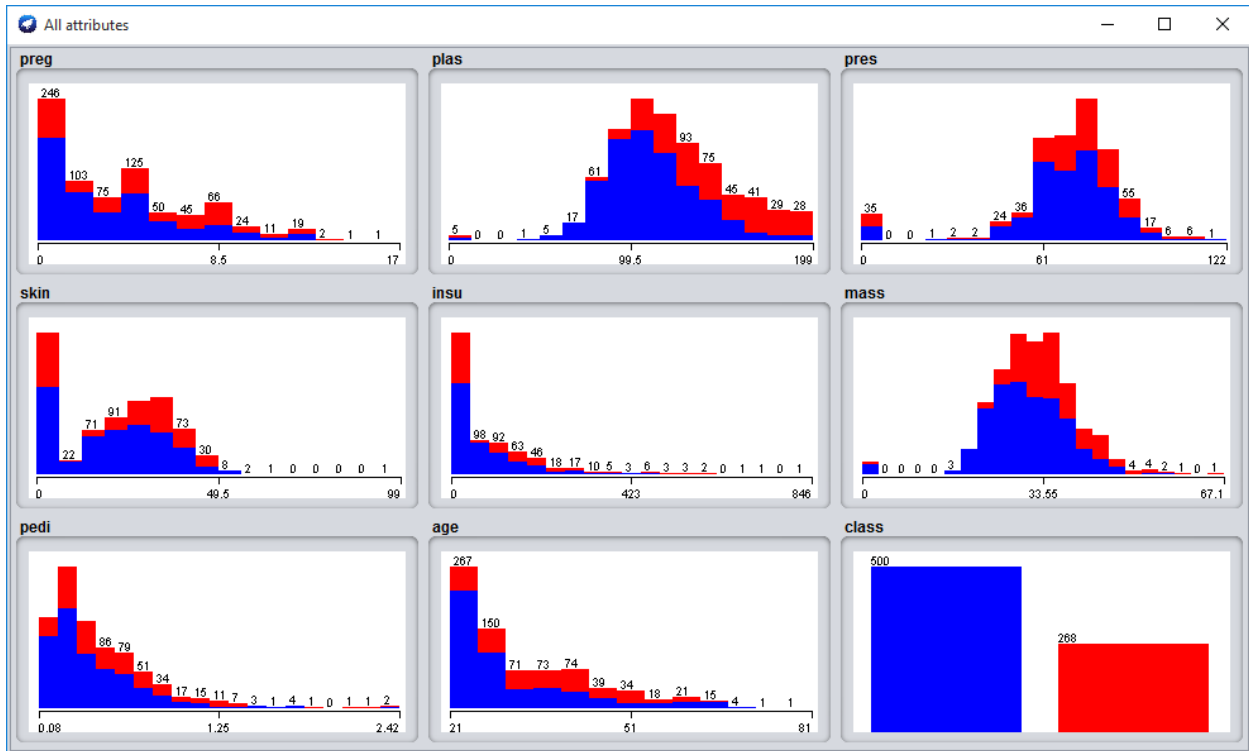
Add instance Undo OK Cancel

- vii) Load each dataset and observe the following:
- List the attribute names and they types
 - Number of records in each dataset
 - Identify the class attribute (if any)
 - Plot histogram
 - Determine the number of records for each class.
 - Visualize the data in various dimensions

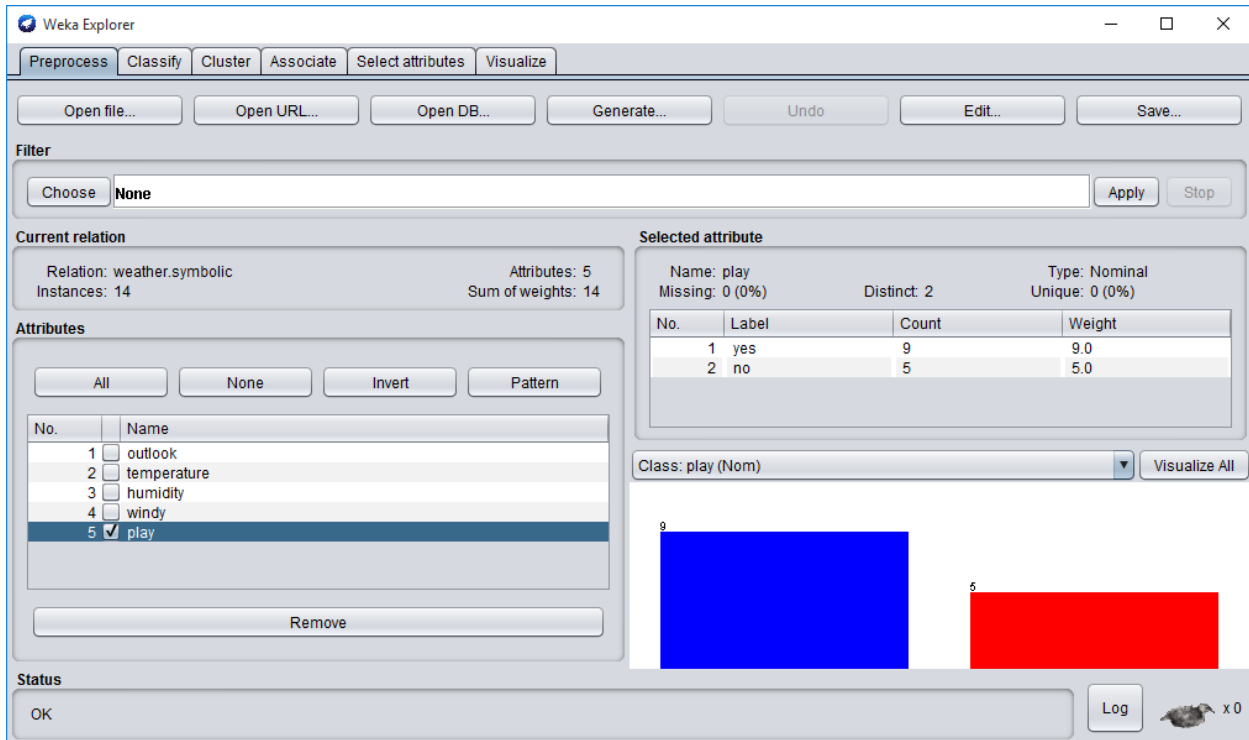
The screenshot shows the Weka Explorer interface with the 'Visualize' tab selected. The 'Current relation' is 'weather.symbolic' with 14 instances and 5 attributes. The 'Selected attribute' is 'outlook', which is a nominal attribute with 3 distinct values and 0 missing values. The 'Attributes' list includes outlook, temperature, humidity, windy, and play. The 'Class' is set to 'play (Nom)'. The visualization shows a stacked bar chart for the 'outlook' attribute, with the 'play' class represented by blue and the 'no play' class by red. The counts for each outlook category are: sunny (5 sunny, 5 play), overcast (4 overcast, 4 play), and rainy (5 sunny, 5 play).

No.	Label	Count	Weight
1	sunny	5	5.0
2	overcast	4	4.0
3	rainy	5	5.0

d) Plot histogram:



e) Determine the number of records for each class.



f) Visualize the data in various dimensions

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **None** [Apply] [Stop]

Current relation
Relation: pima_diabetes
Instances: 768
Attributes: 9
Sum of weights: 768

Selected attribute
Name: insu
Missing: 0 (0%)
Distinct: 186
Type: Numeric
Unique: 93 (12%)

Statistic	Value
Minimum	0
Maximum	846
Mean	79.799
StdDev	115.244

Class: insu (Num) [Visualize All]

Attributes
All | None | Invert | Pattern

No.	Name
1	preg
2	plas
3	pres
4	skin
5	<input checked="" type="checkbox"/> insu
6	mass
7	pedi
8	age
9	class

[Remove]

Status
OK [Log] x 0

VIVA-QUESTIONS

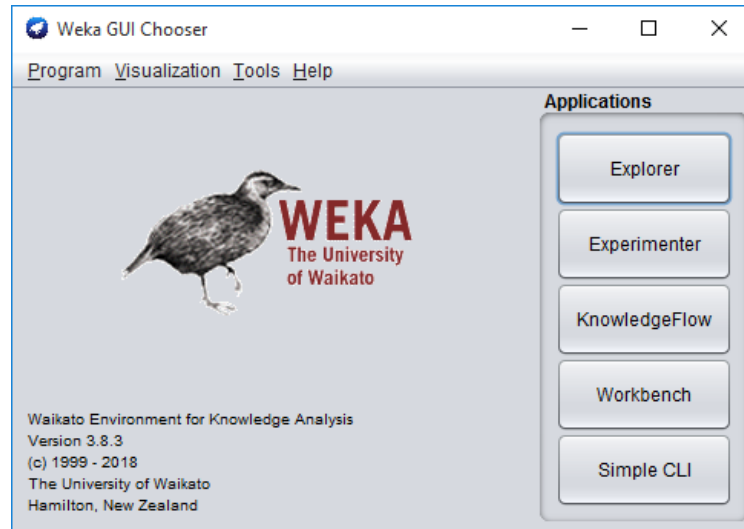
- WEKA Stands for
- ARFF Stands for
- CSV Stands for
- Present Version of WEKA Software?
- Which of the following is not a DATA Mining Software ?
A. WEKA B. SPSS C. Orange D. Cogno's
- Which of the following is not a Data Type in Weka ?
A. Numeric B. String C. Date D. Real
- Which of the following is not a support files in Weka ?
A. ARFF B. CSV C..cpp
- The full form of KDD is
- Data Mining helps in
- Extreme values that occur infrequently are called as

EXP2: Perform data preprocessing tasks on

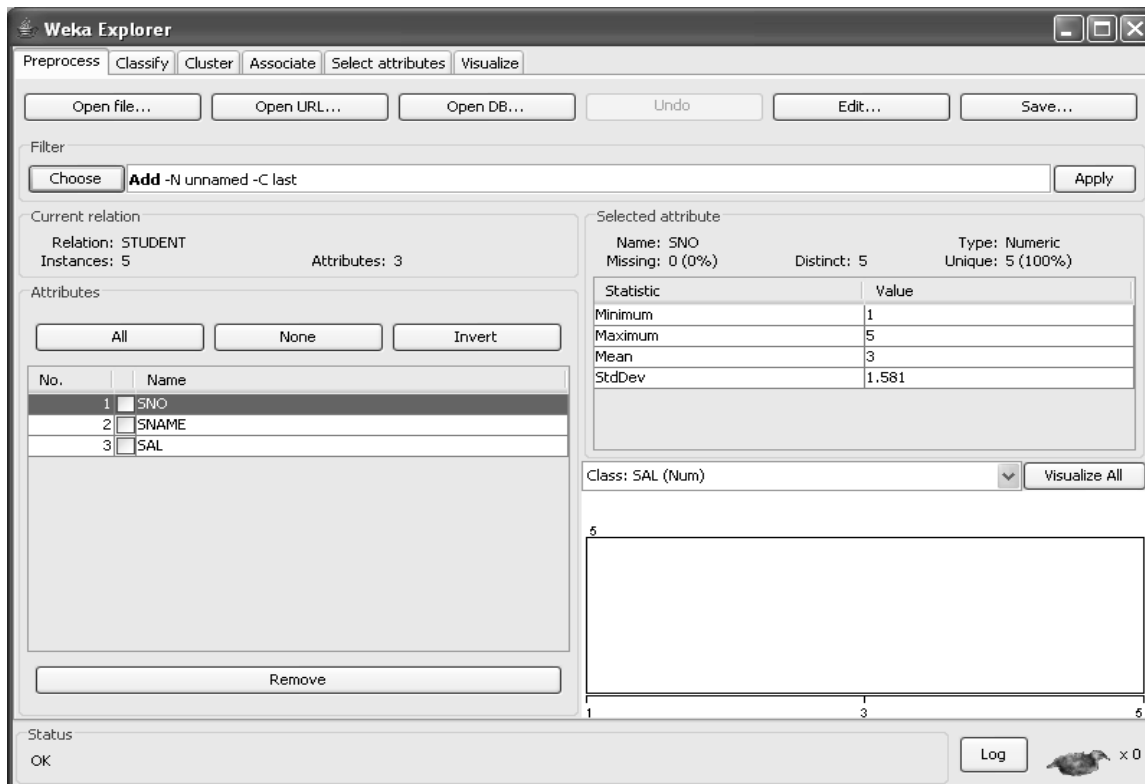
i. Add attribute ii. Add expression iii. Copy attribute iv. Remove attribute

Filters

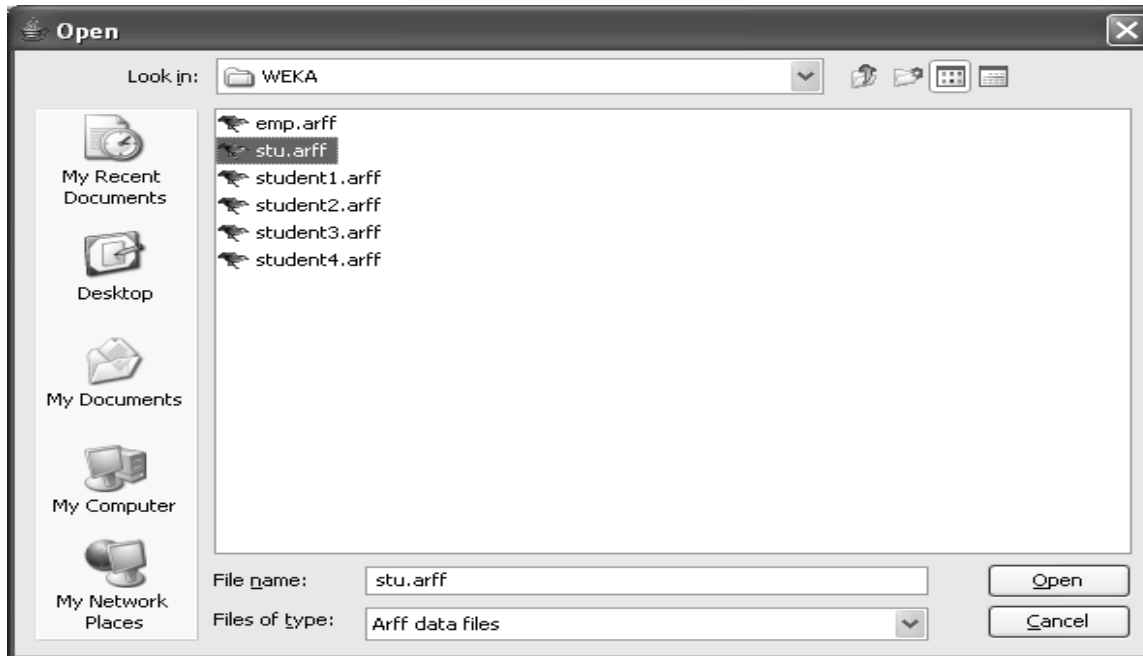
Step-1: - Go to start button then select All Programs and then select weka 3.8.3



Click on Explorer

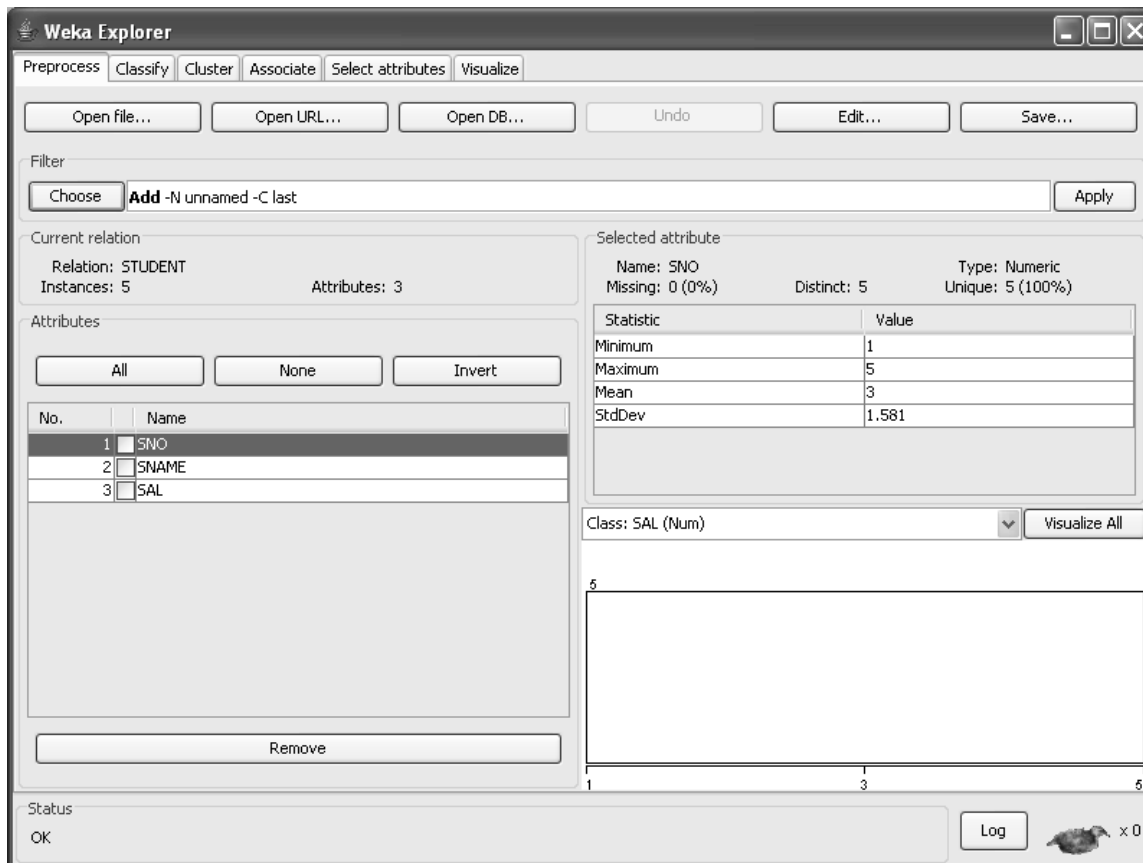


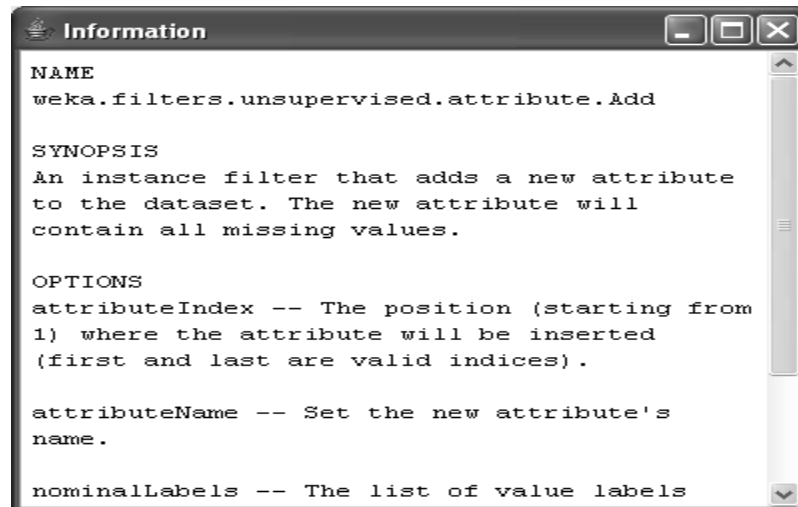
To open the required file



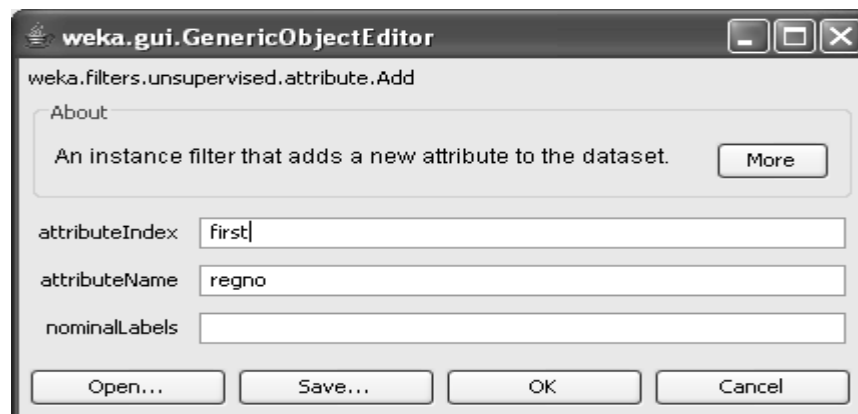
ADD ATTRIBUTE

Step-2: - Go to preprocess menu and choose then select attribute of Add option.

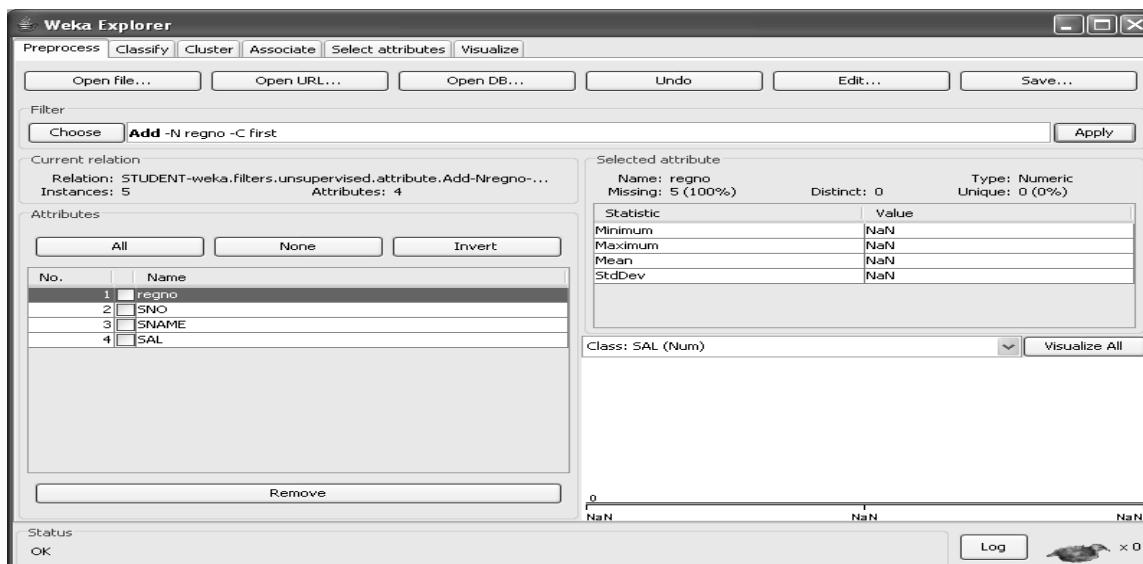




Step-3: -To enter the required fields is added to table.



Click on ok button



Click on Apply button.

No.	regno	SNO	SNAME	SAL
	Numeric	Numeric	String	Numeric
1	1.0	DEEPTHI	1000.0	
2	2.0	SIRI	2000.0	
3	3.0	RAJI	4000.0	
4	4.0	LAKS...	3000.0	
5	5.0	PRAS...	5000.0	

Step-4: - if the index position is last to enter the fields

weka.gui.GenericObjectEditor

weka.filters.unsupervised.attribute.Add

About

An instance filter that adds a new attribute to the dataset. More

attributeIndex: last

attributeName: avg

nominalLabels: Set the new attribute's name

Open... Save... OK Cancel

Click on Ok button.

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **Add -N avg -C last** Apply

Current relation

Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Nregno-...
Instances: 5 Attributes: 5

Attributes

All None Invert

No.	Name
1	<input type="checkbox"/> regno
2	<input type="checkbox"/> SNO
3	<input type="checkbox"/> SNAME
4	<input type="checkbox"/> SAL
5	<input type="checkbox"/> avg

Remove

Status: OK

Selected attribute

Name: regno Type: Numeric
Missing: 5 (100%) Distinct: 0 Unique: 0 (0%)

Statistic	Value
Minimum	NaN
Maximum	NaN
Mean	NaN
StdDev	NaN

Class: avg (Num) Visualize All

Log x 0

Click on Apply button.

Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Nregno-Cfirst-weka.filters.uns...

No.	regno	SNO	SNAME	SAL	avg
	Numeric	Numeric	String	Numeric	Numeric
1		1.0	DEEPTHI	1000.0	
2		2.0	SIRI	2000.0	
3		3.0	RAJI	4000.0	
4		4.0	LAKS...	3000.0	
5		5.0	PRAS...	5000.0	

Buttons: Undo, OK, Cancel

Click on ok button.

Step-5: -if the index position is middle

weka.gui.GenericObjectEditor

weka.filters.unsupervised.attribute.Add

About

An instance filter that adds a new attribute to the dataset. More

attributeIndex:

attributeName:

nominalLabels:

Buttons: Open..., Save..., OK, Cancel

Click on ok button.

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **Add -N "max sal" -C 4** Apply

Current relation: Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navig-Cla... Instances: 5 Attributes: 6

Attributes: All None Invert

No.	Name
1	regno
2	SNO
3	SNAME
4	'max sal'
5	SAL
6	avg

Inverts the current attribute selection

Remove

Selected attribute

Name: regno
Missing: 5 (100%)
Distinct: 0
Type: Numeric
Unique: 0 (0%)

Statistic	Value
Minimum	NaN
Maximum	NaN
Mean	NaN
	NaN

Class: avg (Num) Visualize All

Status: OK Log x 0

Click on Apply button.

Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Clast-weka.filters.unsup...

No.	regno Numeric	SNO Numeric	SNAME String	'max sal' Numeric	SAL Numeric	avg Numeric
1		1.0	DEEPTHI		1000.0	
2		2.0	SIRI		2000.0	
3		3.0	RAJI		4000.0	
4		4.0	LAKS...		3000.0	
5		5.0	PRAS...		5000.0	

Buttons: Undo, OK, Cancel

Click on ok button.

Add Expression

Step-6: - click on choose button and then select AddExpression option.

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **AddExpression** -E a1^2 -N expression Apply

Current relation
Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Cla...
Instances: 5 | Attributes: 6

Selected attribute
Name: regno | Type: Numeric
Missing: 0 (0%) | Distinct: 5 | Unique: 5 (100%)

Statistic	Value
Minimum	12
Maximum	16
Mean	14
StdDev	1.581

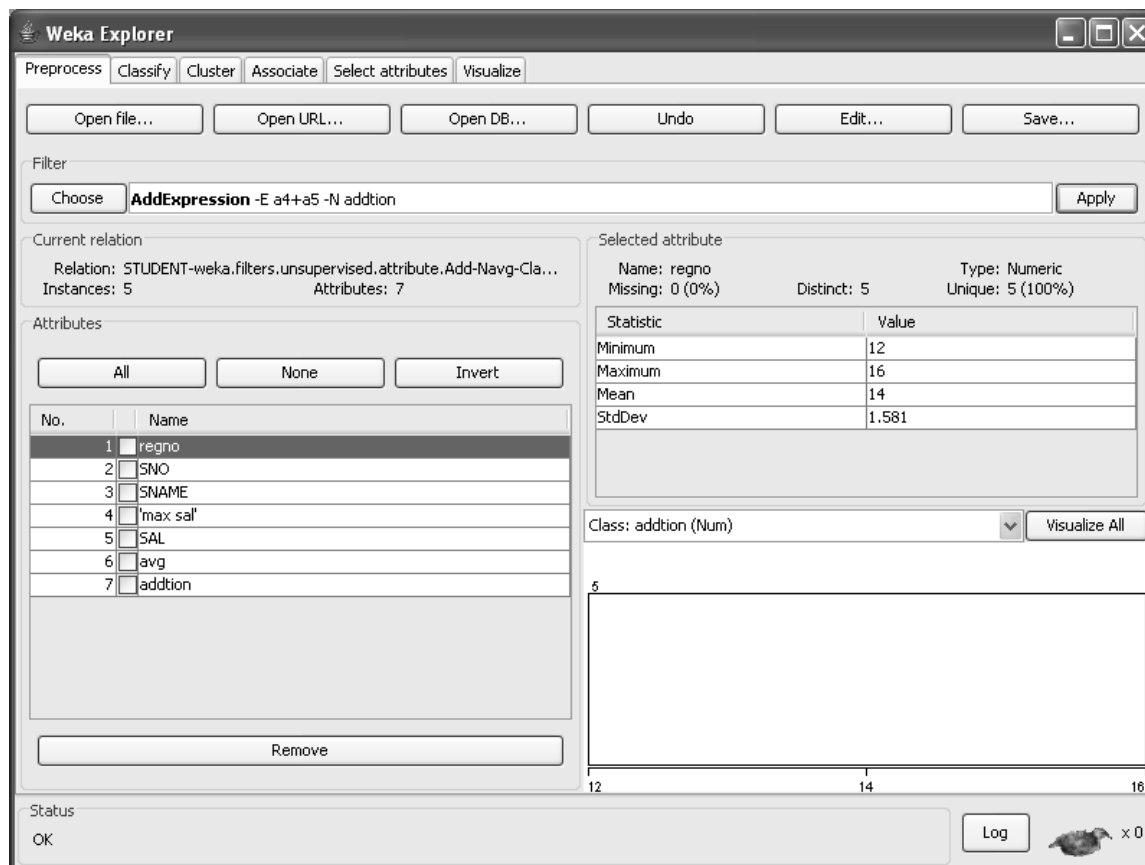
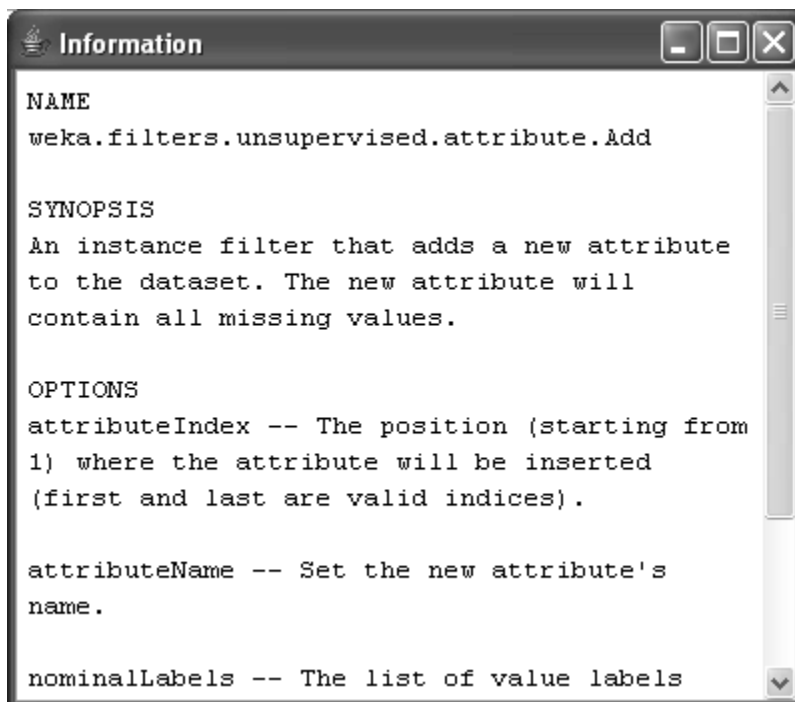
Class: avg (Num) Visualize All

Attributes: All | None | Invert

No.	Name
1	<input checked="" type="checkbox"/> regno
2	<input type="checkbox"/> SNO
3	<input type="checkbox"/> SNAME
4	<input type="checkbox"/> 'max sal'
5	<input type="checkbox"/> SAL
6	<input type="checkbox"/> avg

Remove

Status: OK | Log | x 0



Click on Apply option.

Viewer [X]

Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Clast-weka.filters.unsup...

No.	regno Numeric	SNO Numeric	SNAME String	'max sal' Numeric	SAL Numeric	avg Numeric	addition Numeric
1	12.0	1.0	DEEPTHI	2000.0	1000.0	10.2	3000.0
2	13.0	2.0	SIRI	4000.0	2000.0	12.2	6000.0
3	14.0	3.0	RAJI	6000.0	4000.0	13.3	10000.0
4	15.0	4.0	LAKS...	6000.0	3000.0	14.2	9000.0
5	16.0	5.0	PRAS...	10000.0	5000.0		15000.0

Undo OK Cancel

Click on ok button.

Weka Explorer [] [] [X]

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **AddExpression -E a4-a5 -N subtraction** Apply

Current relation
Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Cla...
Instances: 5 Attributes: 8

Selected attribute
Name: regno Type: Numeric
Missing: 0 (0%) Distinct: 5 Unique: 5 (100%)

Statistic	Value
Minimum	12
Maximum	16
Mean	14
StdDev	1.581

Class: subtraction (Num) Visualize All

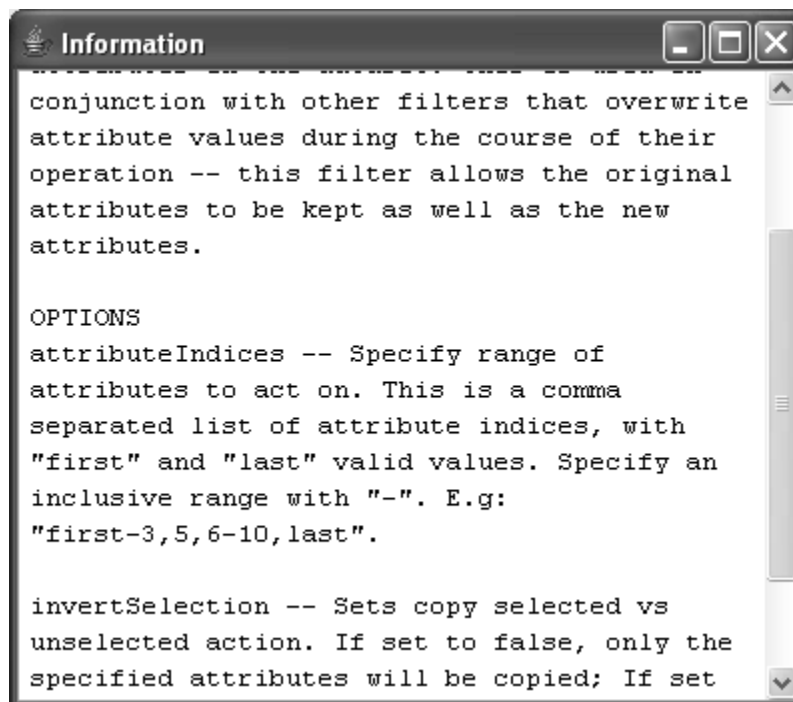
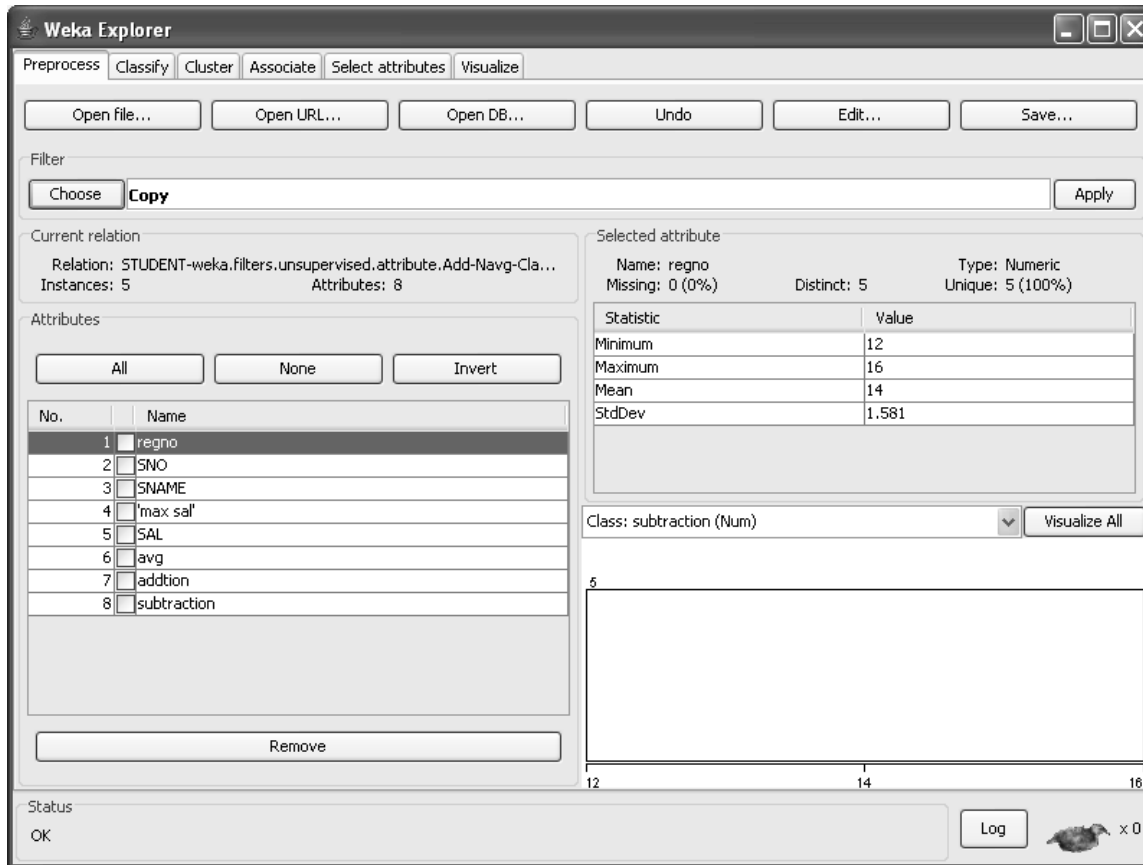
5

12 14 16

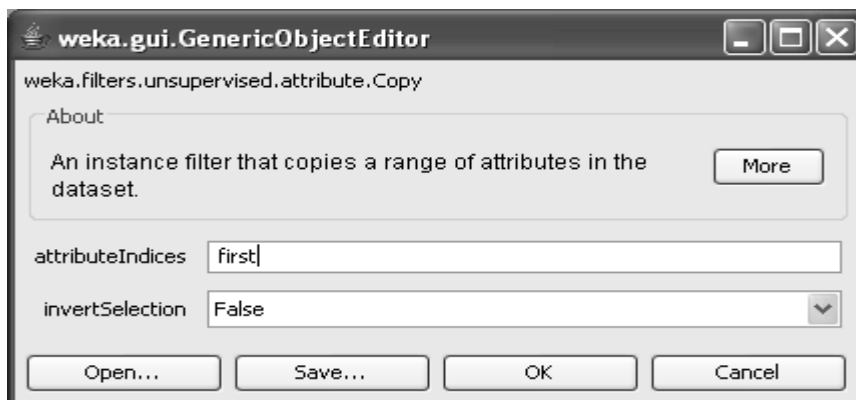
Status: OK Log [] x 0

COPY

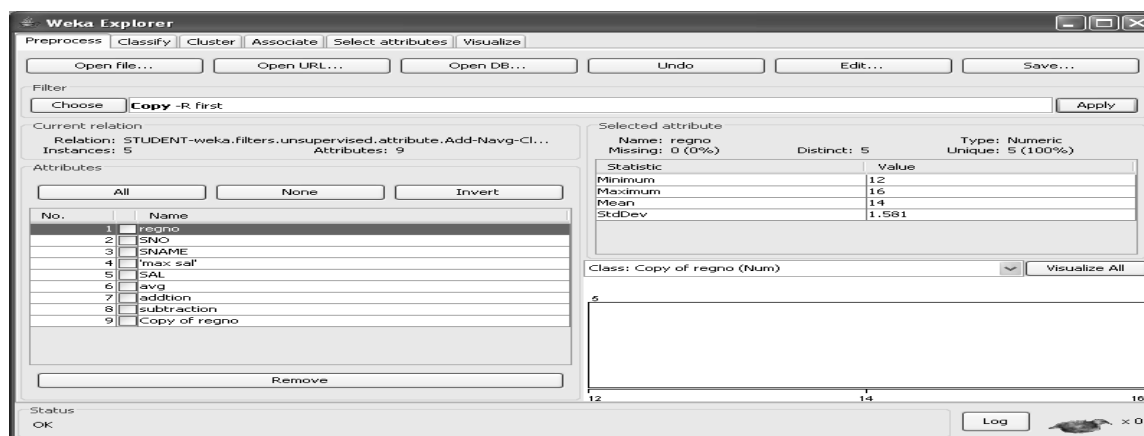
Step-7: - Click on choose button and then select copy option.



If the index place is any



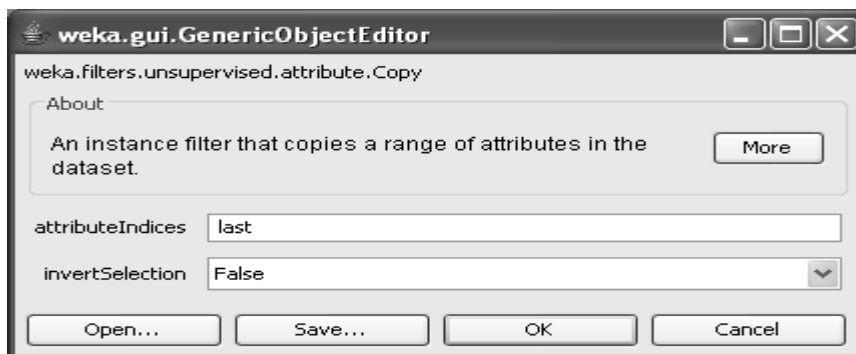
Click on ok button.



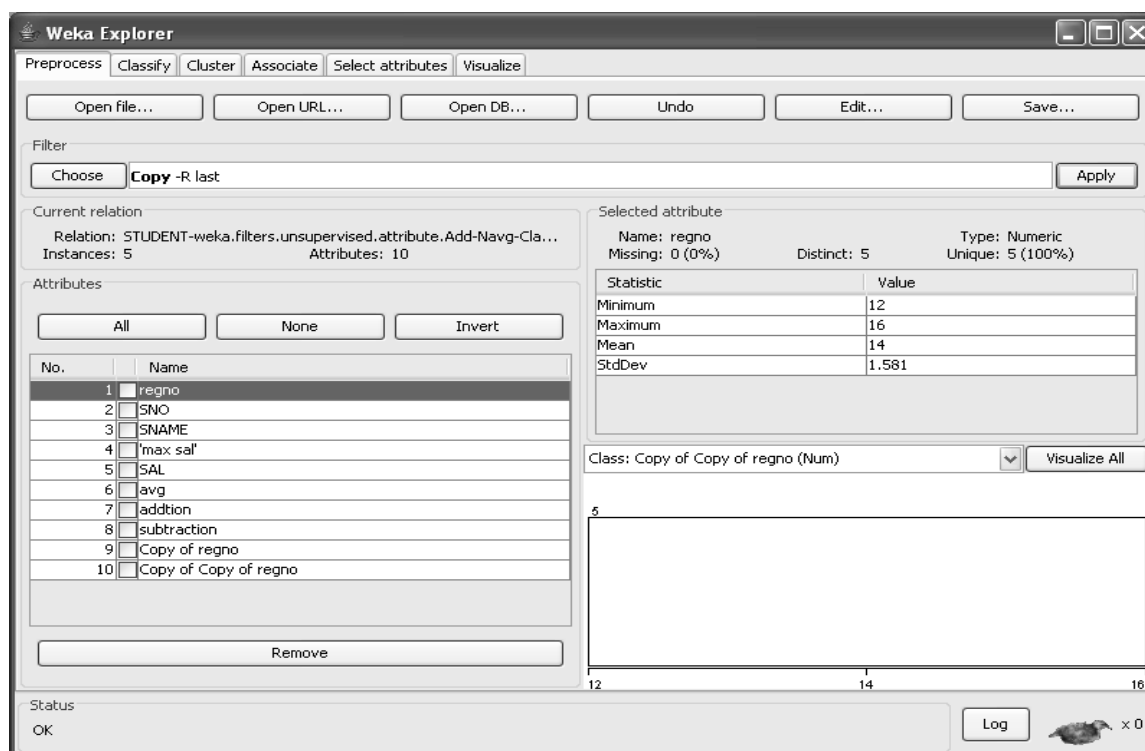
Click On ok button.

Click on ok button.

If the last position.



Click on ok button.

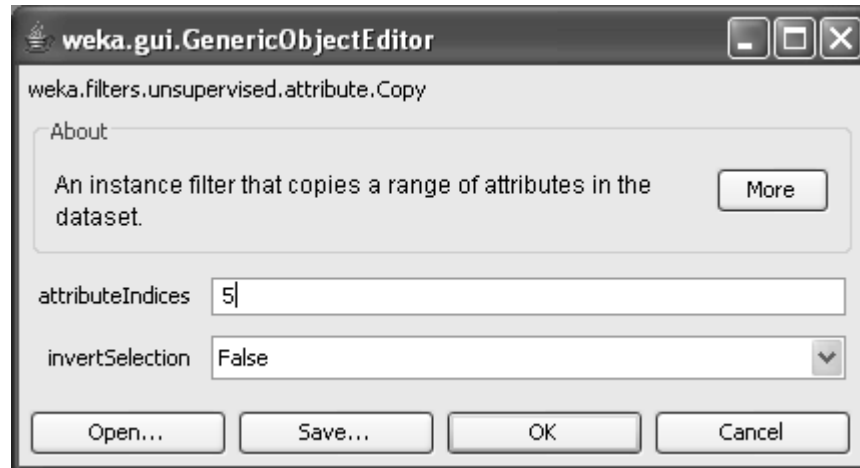


Click on Apply option.

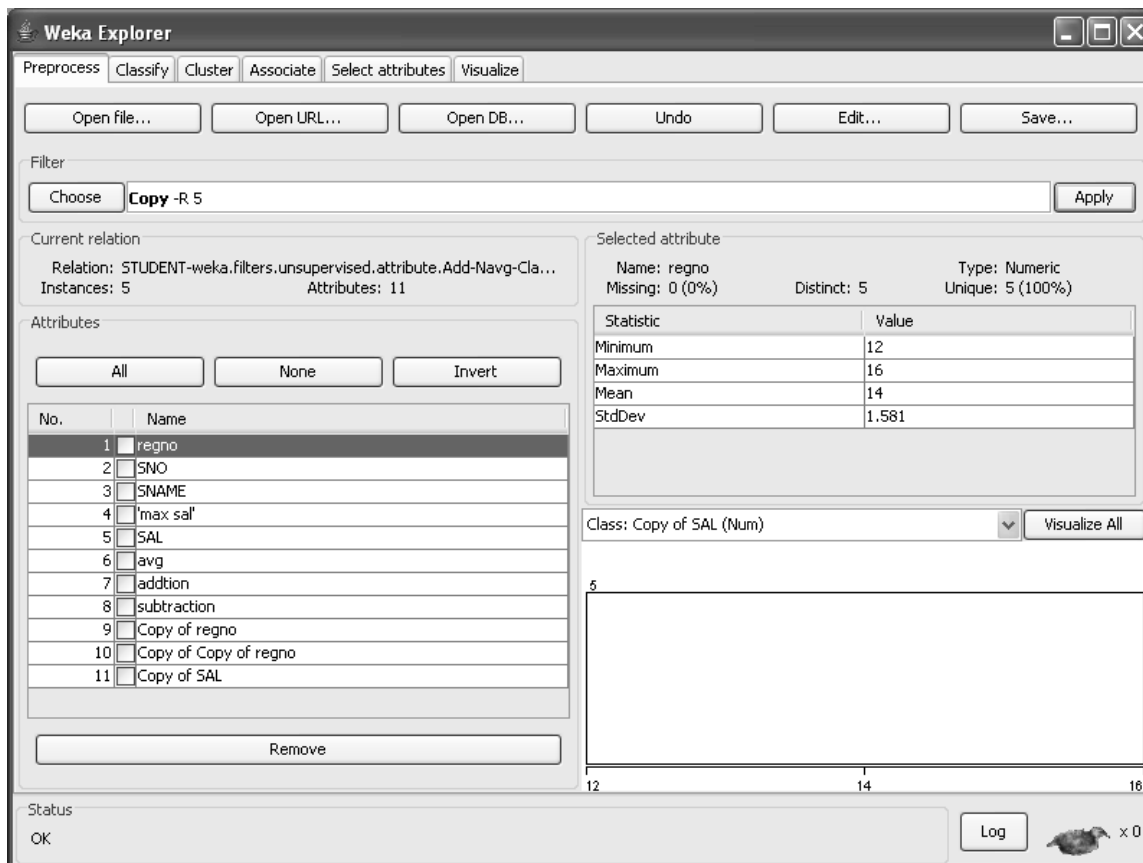
SNO	SNAME	max sal	SAL	avg	addition	subtraction	Copy of regno	Copy
1.0	DEEPTHI	2000.0	1000.0	10.2	3000.0	1000.0	12.0	12.0
2.0	SIRI	4000.0	2000.0	12.2	6000.0	2000.0	13.0	13.0
3.0	RAJI	6000.0	4000.0	13.3	10000.0	2000.0	14.0	14.0
4.0	LAKS...	6000.0	3000.0	14.2	9000.0	3000.0	15.0	15.0
5.0	PRAS...	10000.0	5000.0	15000.0	5000.0	16.0	16.0	16.0

Click on ok button.

If the index position is specified.



Click on ok button.



Click on Apply option.

Viewer

Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Clast-weka.filters.unsupervised.attribute.Add-Nregno-Cfirst-weka.Fil...

No.	regno Numeric	SNO Numeric	SNAME String	'max sal' Numeric	SAL Numeric	avg Numeric	addition Numeric	subtraction Numeric	Copy of regno Numeric	Copy of Copy of regno	Copy of SAL Numeric
1	12.0	1.0	DEEPTHI	2000.0	1000.0	10.2	3000.0	1000.0	12.0	12.0	1000.0
2	13.0	2.0	SIRI	4000.0	2000.0	12.2	6000.0	2000.0	13.0	13.0	2000.0
3	14.0	3.0	RAJI	6000.0	4000.0	13.3	10000.0	2000.0	14.0	14.0	4000.0
4	15.0	4.0	LAKS...	6000.0	3000.0	14.2	9000.0	3000.0	15.0	15.0	3000.0
5	16.0	5.0	PRAS...	10000.0	5000.0		15000.0	5000.0	16.0	16.0	5000.0

Undo OK Cancel

Click on ok button.

If the index position range.

weka.gui.GenericObjectEditor

weka.filters.unsupervised.attribute.Copy

About

An instance filter that copies a range of attributes in the dataset. More

attributeIndices

invertSelection Specify range of attributes t

Open... Save... OK Cancel

Click on ok button.

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter Apply

Current relation
Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navg-Cl...
Instances: 5 Attributes: 15

Attributes

No.	Name
<input type="checkbox"/>	'max sal'
<input type="checkbox"/>	SAL
<input type="checkbox"/>	avg
<input type="checkbox"/>	addition
<input type="checkbox"/>	subtraction
<input type="checkbox"/>	Copy of regno
<input type="checkbox"/>	Copy of Copy of regno
<input type="checkbox"/>	Copy of SAL
<input type="checkbox"/>	Copy of SAL
<input type="checkbox"/>	Copy of avg
<input type="checkbox"/>	Copy of addition
<input type="checkbox"/>	Copy of subtraction

Selected attribute

Name: regno	Type: Numeric
Missing: 0 (0%)	Distinct: 5
	Unique: 5 (100%)

Statistic	Value
Minimum	12
Maximum	16
Mean	14
StdDev	1.581

Class: Copy of subtraction (Num) Visualize All

5

12 14 16

Status
OK Log x 0

Click on Apply option.

	'max sal'	SAL	avg	addition	subtraction	Copy of regno	Copy of Copy of regno
HI	2000.0	1000.0	10.2	3000.0	1000.0	12.0	
	4000.0	2000.0	12.2	6000.0	2000.0	13.0	
	6000.0	4000.0	13.3	10000.0	2000.0	14.0	
..	6000.0	3000.0	14.2	9000.0	3000.0	15.0	
..	10000.0	5000.0		15000.0	5000.0	16.0	

Click on ok button.

REMOVE

Step-8:- Click on choose button and then select Remove option

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

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

Filter: Choose **Remove** Apply

Current relation: Relation: STUDENT-weka.filters.unsupervised.attribute.Add-Navig-Clas...
Instances: 5 | Attributes: 15

Attributes: All | None | Invert

No.	Name
4	'max sal'
5	SAL
6	avg
7	addition
8	subtraction
9	Copy of regno
10	Copy of Copy of regno
11	Copy of SAL
12	Copy of SAL
13	Copy of avg
14	Copy of addition
15	Copy of subtraction

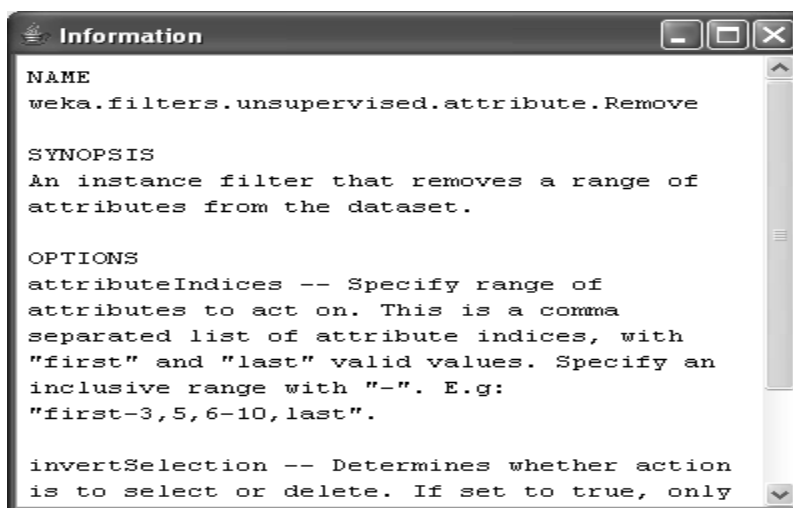
Remove

Selected attribute: Name: regno | Missing: 0 (0%) | Distinct: 5 | Type: Numeric | Unique: 5 (100%)

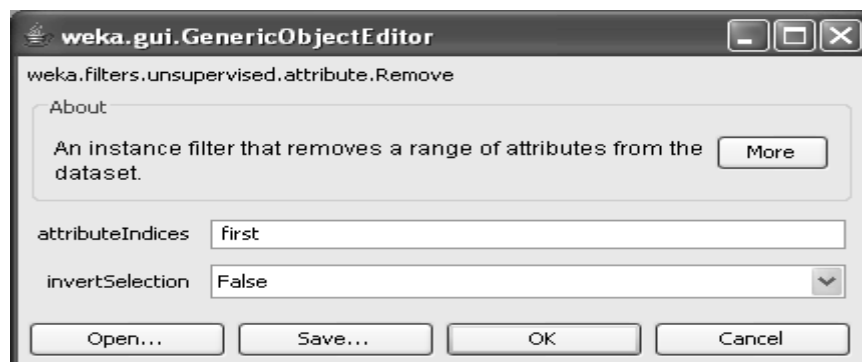
Statistic	Value
Minimum	12
Maximum	16
Mean	14
StdDev	1.581

Class: Copy of subtraction (Num) Visualize All

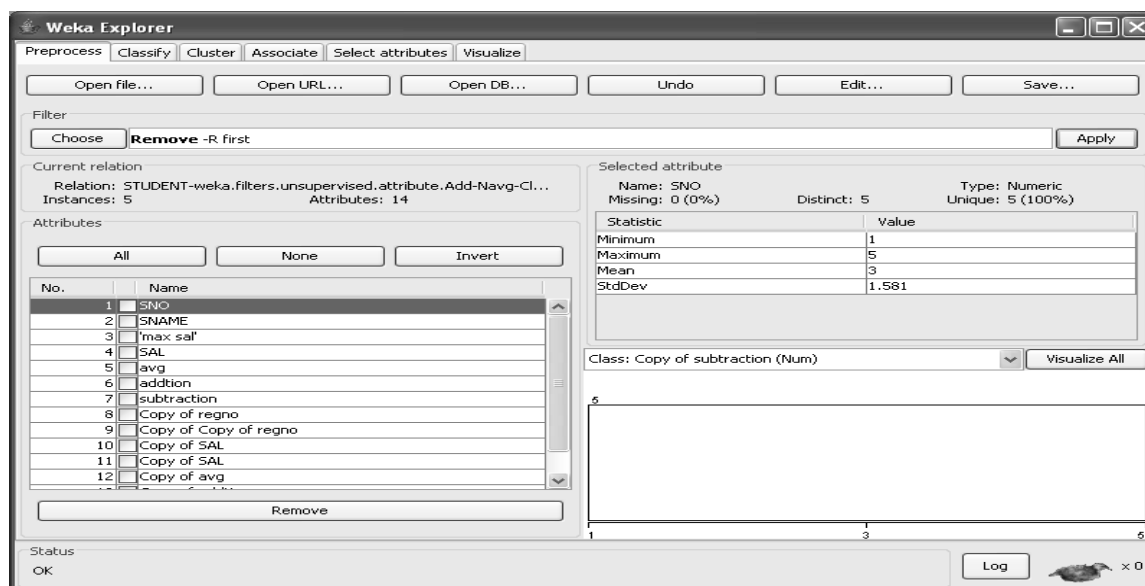
Status: OK | Log | x 0



If the index position is first.



Click on ok button.

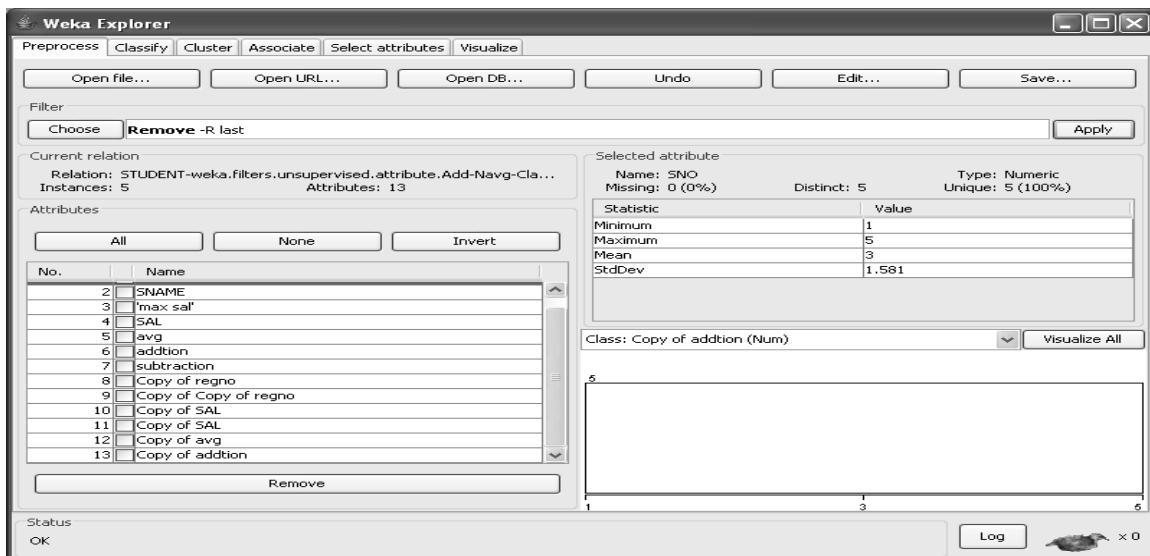


Click on Apply option.

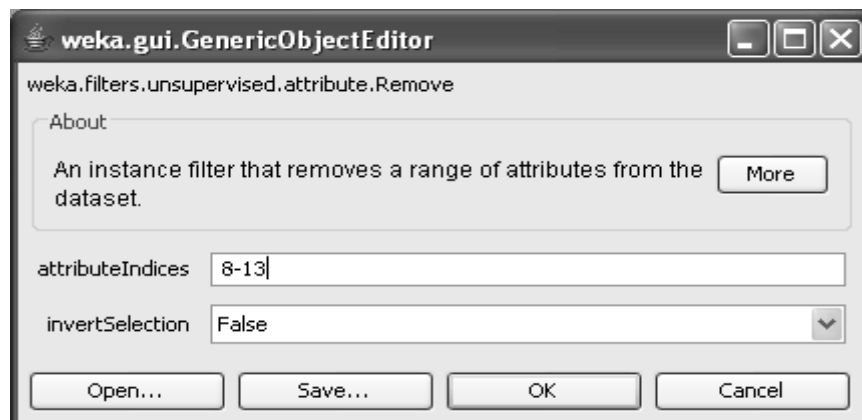
If the index position is last.



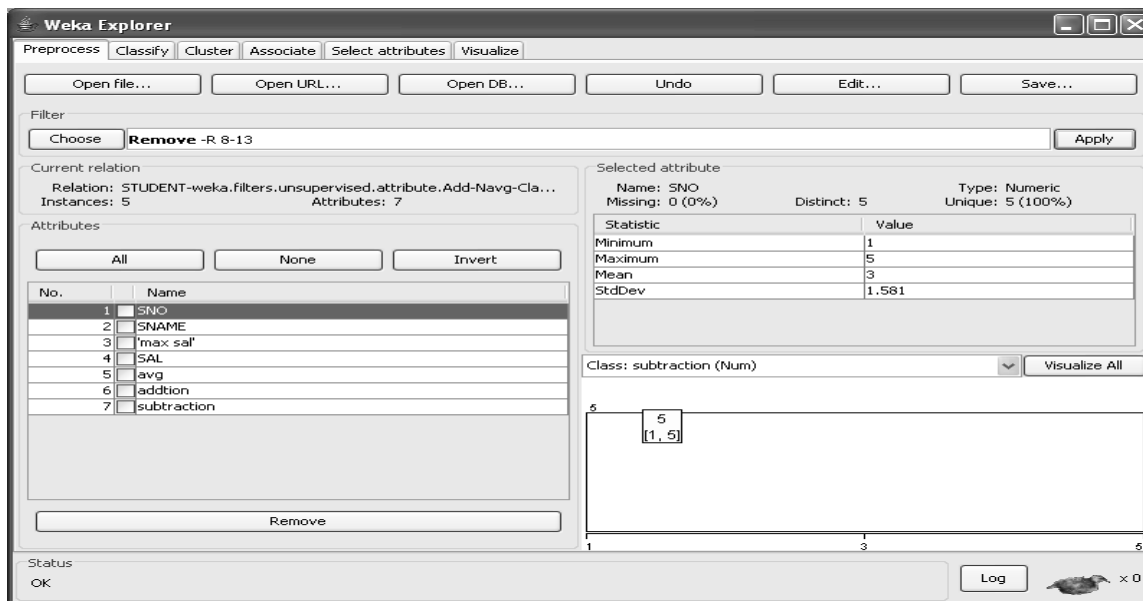
Click on ok button.



If the range .



Click on ok button.



Click on Apply option.

Viewer interface showing a table of data. The table has 8 columns: No., SNO, SNAME, 'max sal', SAL, avg, addition, and subtraction. The data rows are:

No.	SNO	SNAME	'max sal'	SAL	avg	addition	subtraction
1	1.0	DEEPTHI	2000.0	1000.0	10.2	3000.0	1000.0
2	2.0	SIRI	4000.0	2000.0	12.2	6000.0	2000.0
3	3.0	RAJI	6000.0	4000.0	13.3	10000.0	2000.0
4	4.0	LAKS...	6000.0	3000.0	14.2	9000.0	3000.0
5	5.0	PRAS...	10000.0	5000.0		15000.0	5000.0

Click on ok button.

VIVA – QUESTIONS

1. When to apply the data preprocessing techniques for mining the data
2. Which of the following is not a data preprocessing methods:
 - A. Data Visualization
 - B. Data Discretization
 - C. Data Cleaning
 - D. Data Reduction
3. Use the attribute mean to fill the missing value of data
1,2,3,4,5,6,__,7,8,9,10.
4. Data for Attendance : 50,55,60,65,70,75,80,85,90,95
Partition the above attendance data into equidepth bins of depth 5.
5. Data for Attendance : 4,8,15 Smoot by bin boundaries
6. Incorrect or invalid data is known as_____
7. PCA stands for
8. The Minimum and maximum values for the attribute income are \$12,000and \$98,000 , respectively. We would like to map income to the range[0.0,1.0].By min-max normalization, a value of \$73,600 for income is transformed to
9. The mean and standard deviation of the values for the attribute income are \$54,000 and \$16,000 respectively. With Z-Score normalization a value of \$73,600 for income is transformed to
10. ____ reduces the data set size by removing irrelevant or redundant attributes(dimensions)

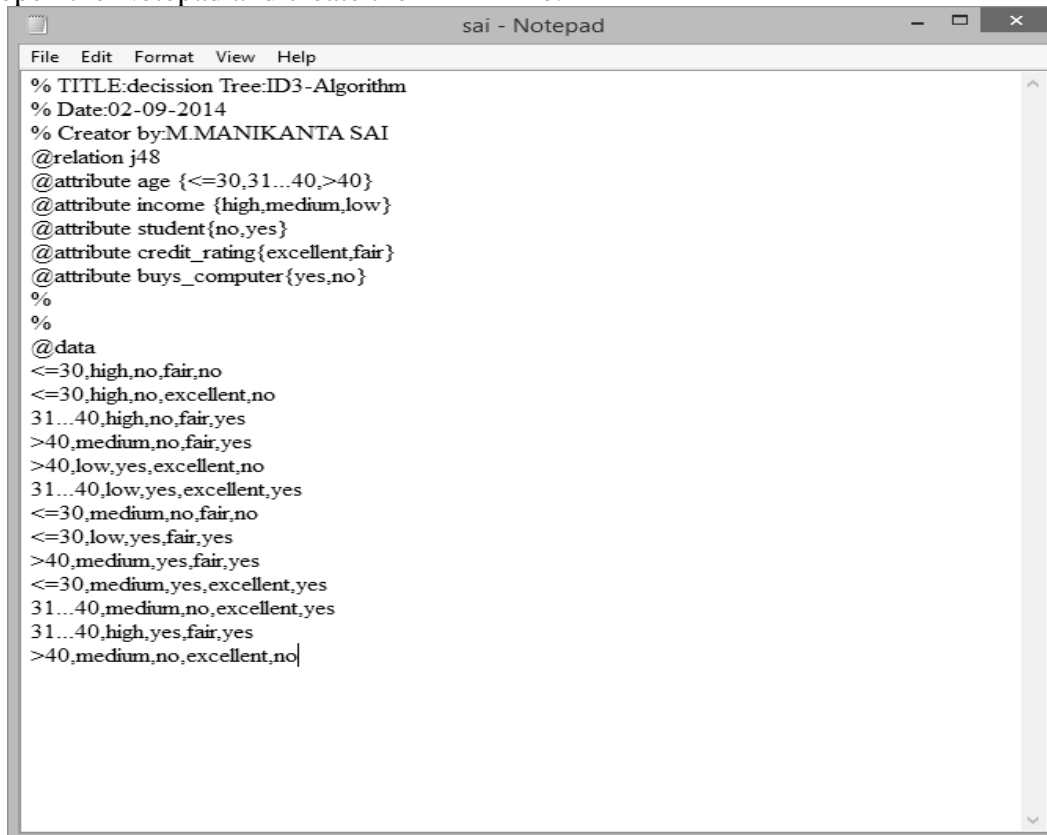
Experiment No. 3: Demonstrate performing classification on data sets

A.

- i. Load weather dataset into WEKA and run Id3, J48 classification algorithm. Study the classifier output. Compute entropy values, kappa statistic.
- ii. Extract if-then rules from the decision tree generated by the classifier, observe the confusion matrix and derive Accuracy, F-measure, TPrate, FPrate, precision and recall values. Apply cross-validation strategy with various fold levels and compare the accuracy results

ID3 Algorithm

Step1: open the Notepad and create the ARFF File.



```

sai - Notepad
File Edit Format View Help
% TITLE:decision Tree:ID3-Algorithm
% Date:02-09-2014
% Creator by:M.MANIKANTA SAI
@relation j48
@attribute age {<=30,31...40,>40}
@attribute income {high,medium,low}
@attribute student {no,yes}
@attribute credit_rating {excellent,fair}
@attribute buys_computer {yes,no}
%
%
@data
<=30,high,no,fair,no
<=30,high,no,excellent,no
31...40,high,no,fair,yes
>40,medium,no,fair,yes
>40,low,yes,excellent,no
31...40,low,yes,excellent,yes
<=30,medium,no,fair,no
<=30,low,yes,fair,yes
>40,medium,yes,fair,yes
<=30,medium,yes,excellent,yes
31...40,medium,no,excellent,yes
31...40,high,yes,fair,yes
>40,medium,no,excellent,no
  
```

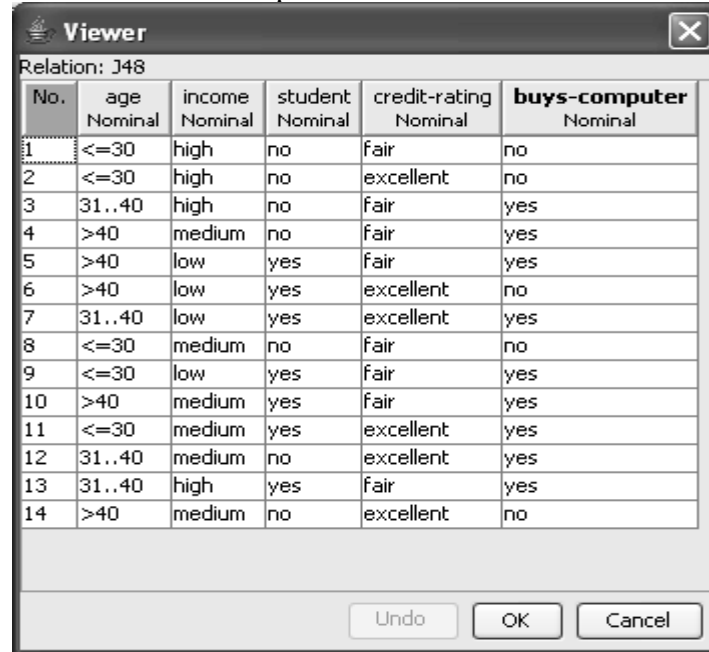
Step 2: To open All Programms → weka3.8.3 → weka3.8(with console)

STEP 3: Click on **EXPLORER** .



Step 4: open the ARFF file from the folder and click open.

Step 5: In the Preprocessor menu select open file and select Edit

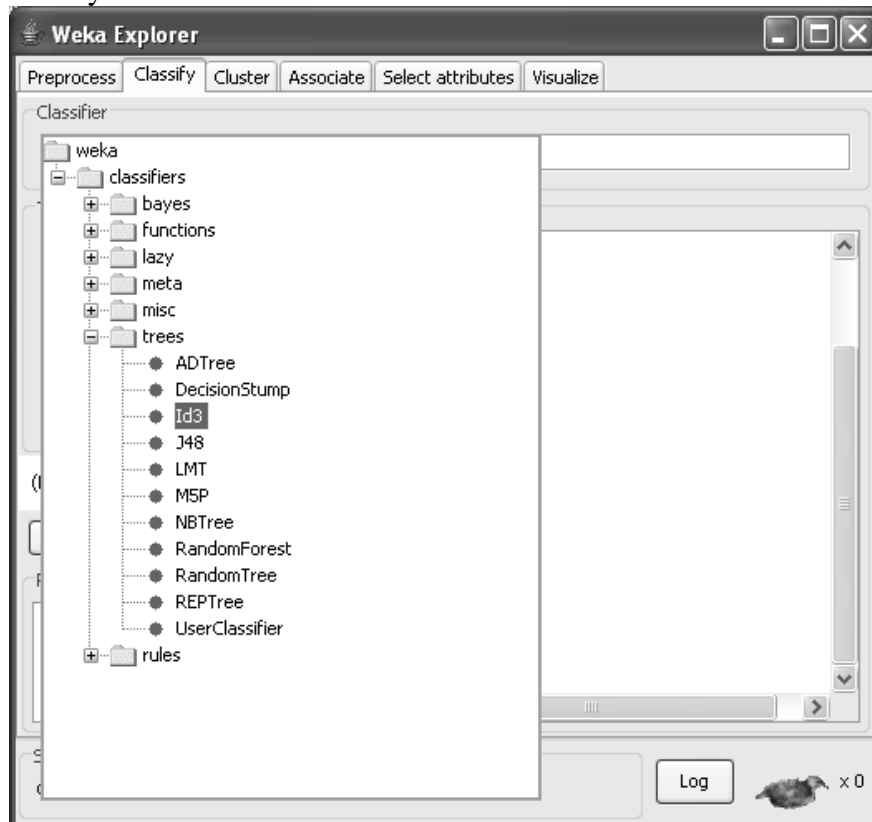


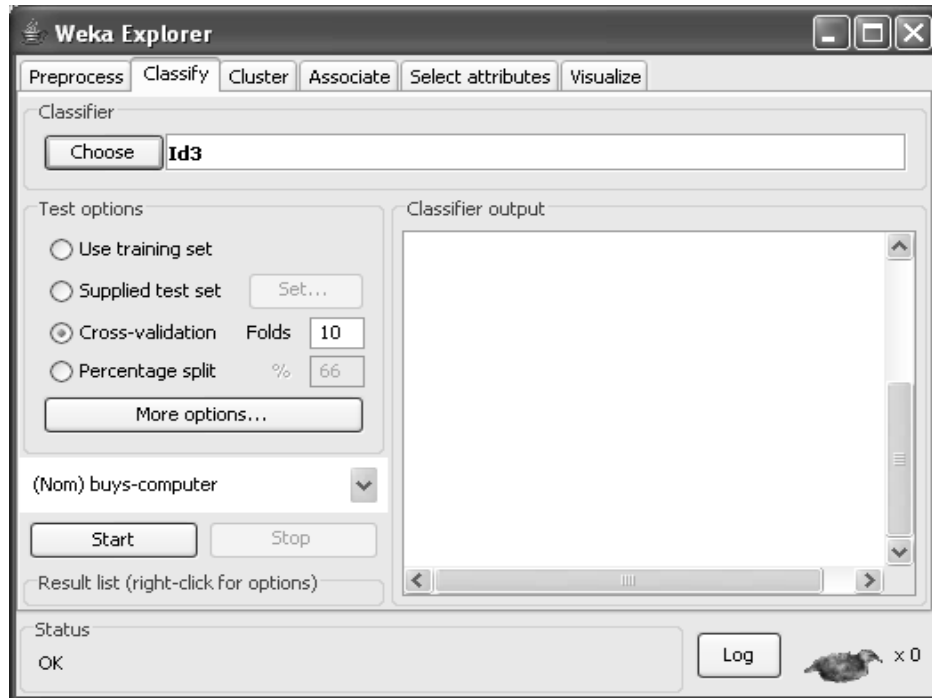
The screenshot shows the 'Viewer' window for the 'Relation: J48'. It displays a table with 14 rows and 6 columns. The columns are: No., age (Nominal), income (Nominal), student (Nominal), credit-rating (Nominal), and buys-computer (Nominal). The data is as follows:

No.	age Nominal	income Nominal	student Nominal	credit-rating Nominal	buys-computer Nominal
1	<=30	high	no	fair	no
2	<=30	high	no	excellent	no
3	31..40	high	no	fair	yes
4	>40	medium	no	fair	yes
5	>40	low	yes	fair	yes
6	>40	low	yes	excellent	no
7	31..40	low	yes	excellent	yes
8	<=30	medium	no	fair	no
9	<=30	low	yes	fair	yes
10	>40	medium	yes	fair	yes
11	<=30	medium	yes	excellent	yes
12	31..40	medium	no	excellent	yes
13	31..40	high	yes	fair	yes
14	>40	medium	no	excellent	no

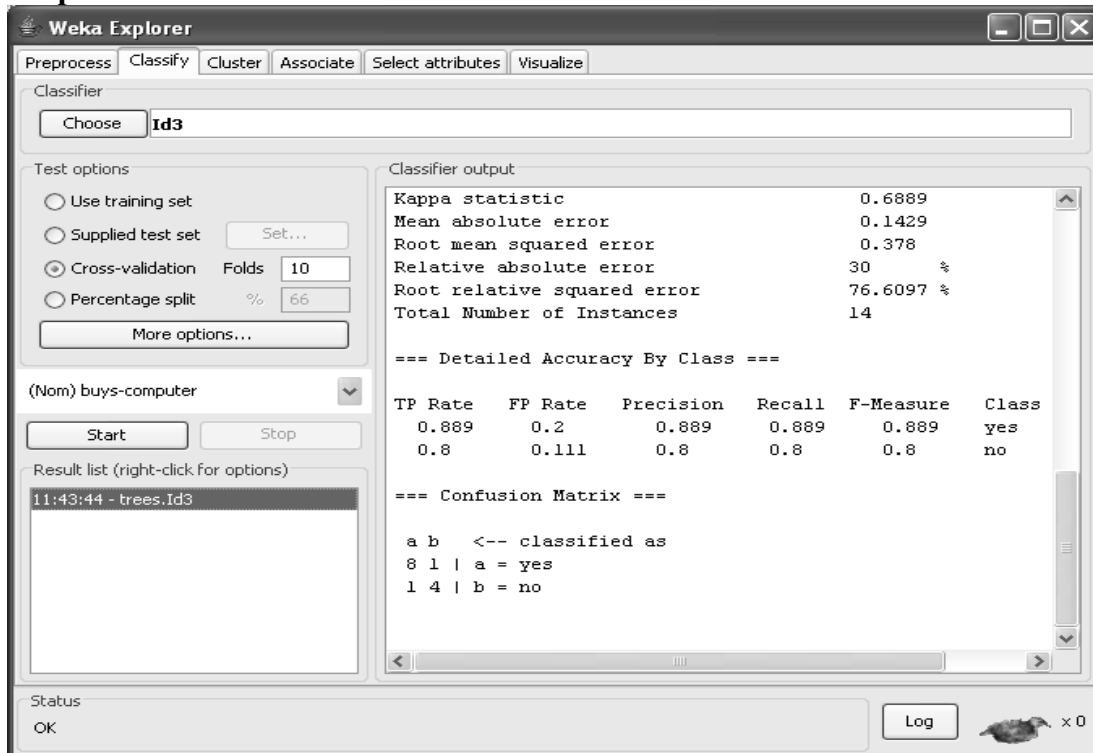
Buttons at the bottom: Undo, OK, Cancel.

Step 6: Go to classify menu click choose on that select trees Click ID3



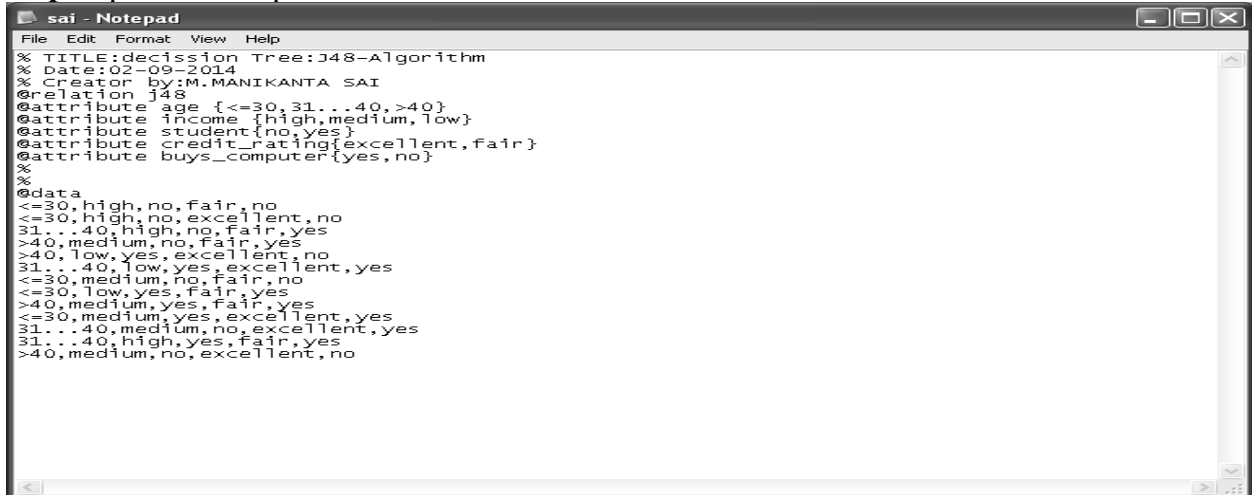


Step 7: And then click start button.



J48

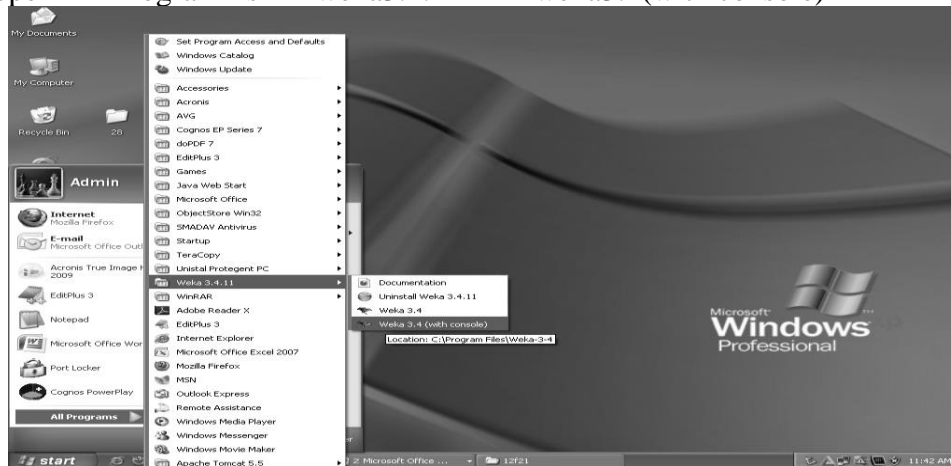
Step1: open the Notepad and create the ARFF File.



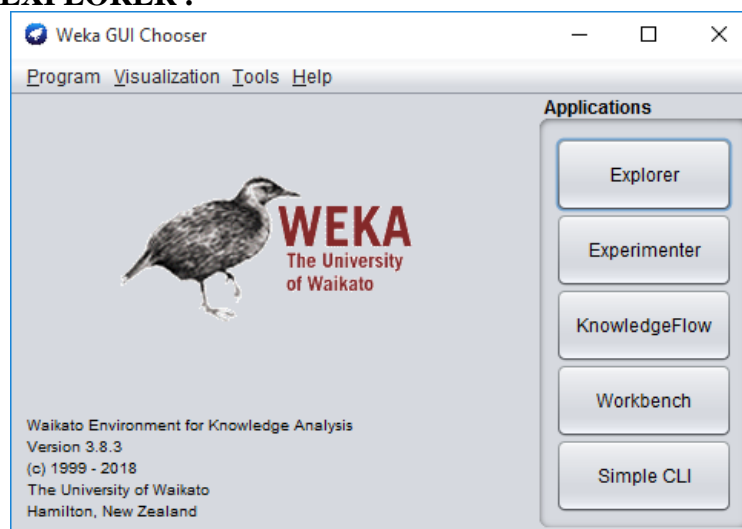
```

sai - Notepad
File Edit Format View Help
% TITLE:decision Tree:J48-Algorithm
% Date:02-09-2014
% Creator by:M.MANIKANTA SAI
@relation j48
@attribute age {<=30,31,..,40,>40}
@attribute income {high,medium,low}
@attribute student{no,yes}
@attribute credit_rating{excellent,fair}
@attribute buys_computer{yes,no}
%
%
@data
<=30,high,no,fair,no
<=30,high,no,excellent,no
31,..,40,high,no,fair,yes
>40,medium,no,fair,yes
>40,low,yes,excellent,no
31,..,40,low,yes,excellent,yes
<=30,medium,no,fair,no
<=30,low,yes,fair,yes
>40,medium,yes,fair,yes
<=30,medium,yes,excellent,yes
31,..,40,medium,no,excellent,yes
31,..,40,high,yes,fair,yes
>40,medium,no,excellent,no
  
```

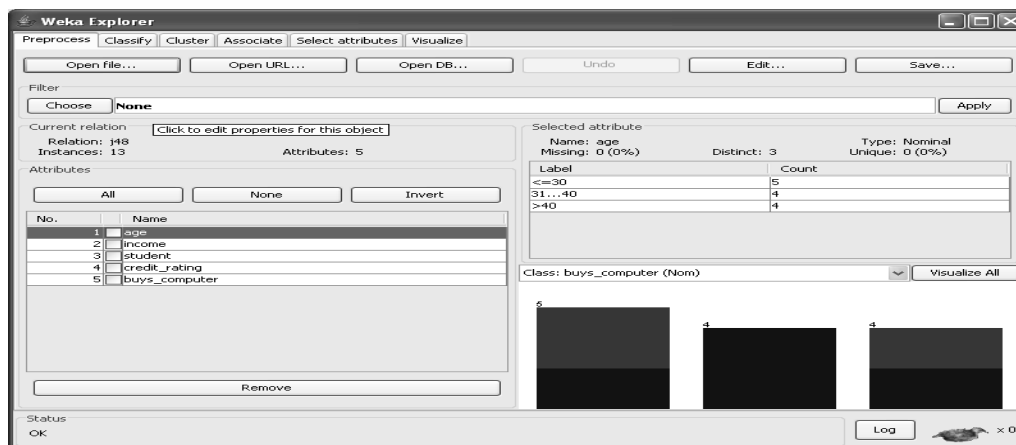
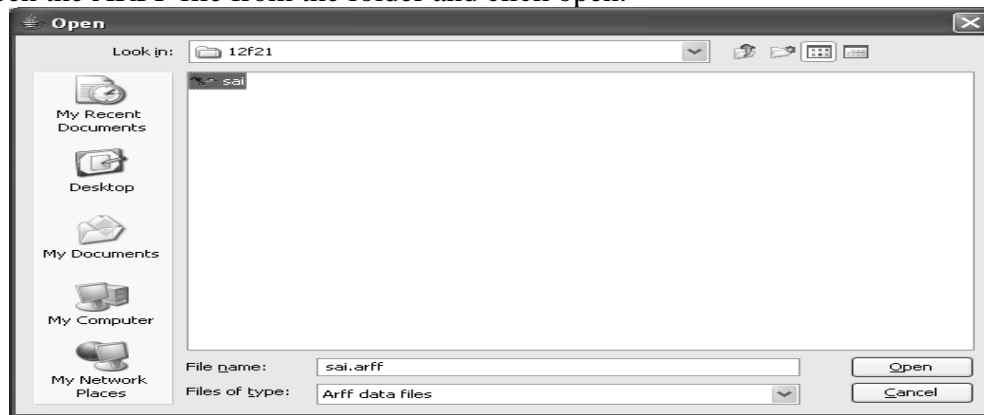
Step2: To open All Programms → weka3.4.11 → weka3.4(with console)



STEP 3: Click on **EXPLORER** .

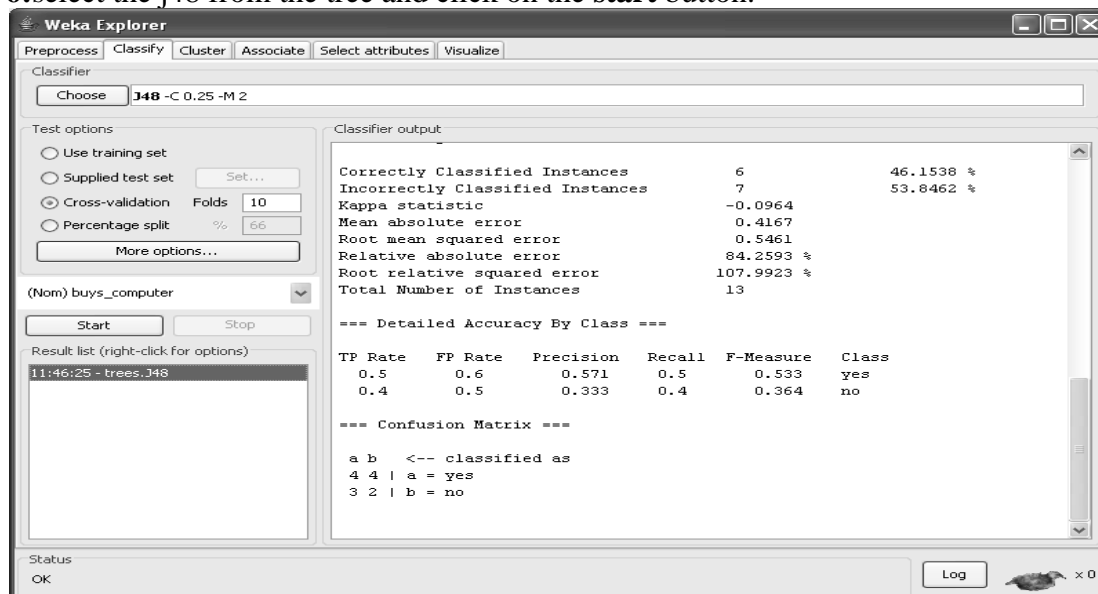


Step 4: open the ARFF file from the folder and click open.

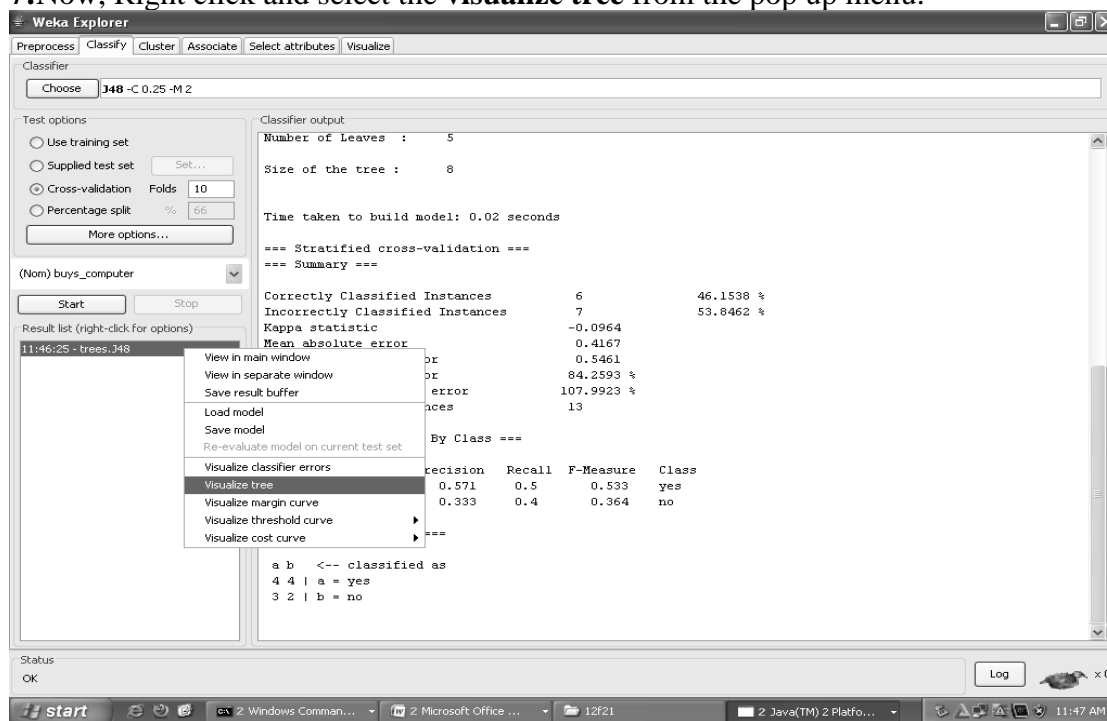


Step 5: click on the **classify** from the menu and click on the **choose**.

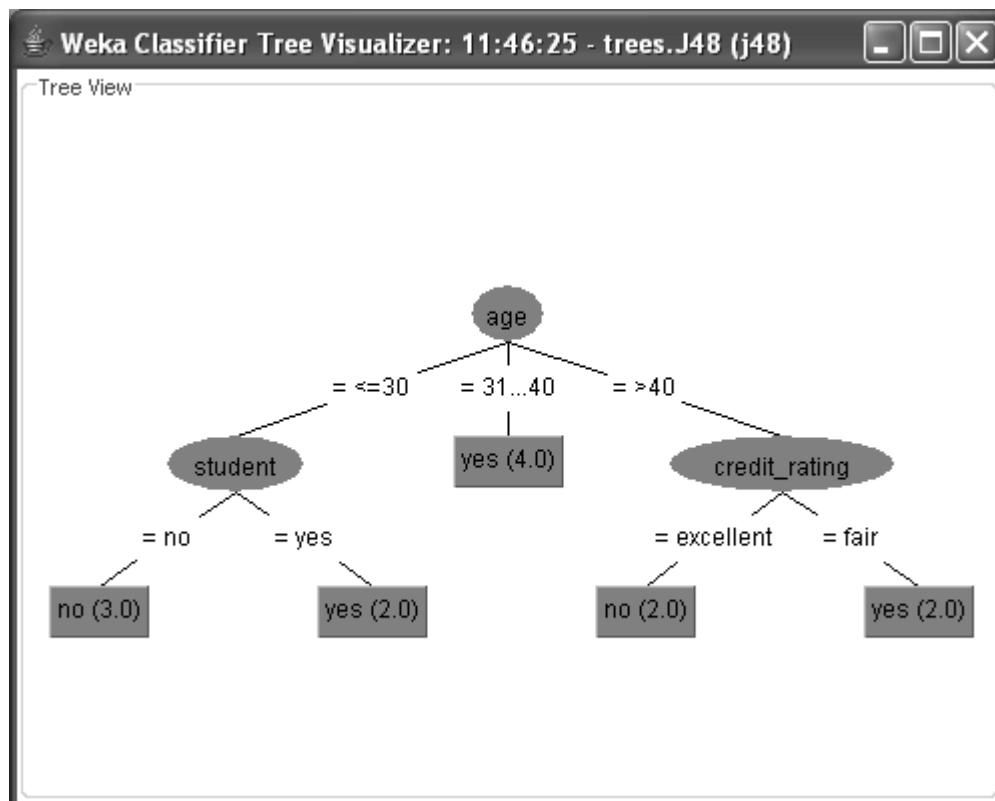
Step 6: select the **j48** from the tree and click on the **start** button.



Step 7: Now, Right click and select the **visualize tree** from the pop up menu.



Output:



B.

- i. Load weather dataset into WEKA and perform Naive-Bayes classification
- ii. K-nearest neighbour classification. Interpret the results obtained.

Naïve Bayesian Classification

The naïve Bayesian classifier, or simple Bayesian classifier, works as follows:

1. Let D be a training set of tuples and their associated class labels. As usual, each tuple is represented by an n -dimensional attribute vector, $X = (x_1, x_2, \dots, x_n)$, depicting n measurements made on the tuple from n attributes, respectively, A_1, A_2, \dots, A_n .
2. Suppose that there are m classes, C_1, C_2, \dots, C_m . Given a tuple, X , the classifier will predict that X belongs to the class having the highest posterior probability, conditioned on X . That is, the naïve Bayesian classifier predicts that tuple X belongs to the class C_i if and only if

$$P(C_i|X) > P(C_j|X) \quad \text{for } 1 \leq j \leq m, j \neq i.$$

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)}.$$

3. As $P(X)$ is constant for all classes, only $P(X|C_i)P(C_i)$ need be maximized. If the class prior probabilities are not known, then it is commonly assumed that the classes are

$$P(X|C_i) = \prod_{k=1}^n P(x_k|C_i)$$

$$= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i).$$

4. Given data sets with many attributes, it would be extremely computationally expensive to compute $P(X|C_i)$. In order to reduce computation in evaluating $P(X|C_i)$, the naive assumption of class conditional independence is made. This presumes that the values of the attributes are conditionally independent of one another, given the class label of the tuple (i.e., that there are no dependence relationships among the attributes). Thus,

We wish to predict the class

label of a tuple using naïve Bayesian classification, given the same training data as in Example weather data set for decision tree induction. The training data are in weather data set. The data tuples are described by the attributes *age*, *income*, *student*, and *credit rating*. The class label attribute, *buys computer*, has two distinct values (namely, *yes*, *no*). Let C_1 correspond to the

class *buys computer = yes* and *C2* correspond to *buys computer = no*. The tuple we wish to classify is

$X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit rating} = \text{fair})$

$P(\text{buys computer} = \text{yes}) = 9/14 = 0.643$

$P(\text{buys computer} = \text{no}) = 5/14 = 0.357$

To compute $P(X_j|C_i)$, for $i = 1, 2$, we compute the following conditional probabilities:

$P(\text{age} = \text{youth} | \text{buys computer} = \text{yes}) = 2/9 = 0.222$

$P(\text{age} = \text{youth} | \text{buys computer} = \text{no}) = 3/5 = 0.600$

$P(\text{income} = \text{medium} | \text{buys computer} = \text{yes}) = 4/9 = 0.444$

$P(\text{income} = \text{medium} | \text{buys computer} = \text{no}) = 2/5 = 0.400$

$P(\text{student} = \text{yes} | \text{buys computer} = \text{yes}) = 6/9 = 0.667$

$P(\text{student} = \text{yes} | \text{buys computer} = \text{no}) = 1/5 = 0.200$

$P(\text{credit rating} = \text{fair} | \text{buys computer} = \text{yes}) = 6/9 = 0.667$

$P(\text{credit rating} = \text{fair} | \text{buys computer} = \text{no}) = 2/5 = 0.400$

We need to maximize $P(X_j|C_i)P(C_i)$, for $i = 1, 2$. $P(C_i)$, the prior probability of each class, can be computed based on the training tuples:

Using the above probabilities, we obtain

$P(X | \text{buys computer} = \text{yes}) = P(\text{age} = \text{youth} | \text{buys computer} = \text{yes}) \cdot$

$P(\text{income} = \text{medium} | \text{buys computer} = \text{yes}) \cdot$

$P(\text{student} = \text{yes} | \text{buys computer} = \text{yes}) \cdot$

$P(\text{credit rating} = \text{fair} | \text{buys computer} = \text{yes})$

$= 0.222 \cdot 0.444 \cdot 0.667 \cdot 0.667 = 0.044.$

Similarly,

$P(X | \text{buys computer} = \text{no}) = 0.600 \cdot 0.400 \cdot 0.200 \cdot 0.400 = 0.019.$

To find the class, C_i , that maximizes $P(X_j|C_i)P(C_i)$, we compute

$P(X | \text{buys computer} = \text{yes})P(\text{buys computer} = \text{yes}) = 0.044 \cdot 0.643 = 0.028$

$P(X | \text{buys computer} = \text{no})P(\text{buys computer} = \text{no}) = 0.019 \cdot 0.357 = 0.007$

Therefore, the naïve Bayesian classifier predicts *buys computer = yes* for tuple X .

Step1: Open Notepad and write the program and save with **.arff**

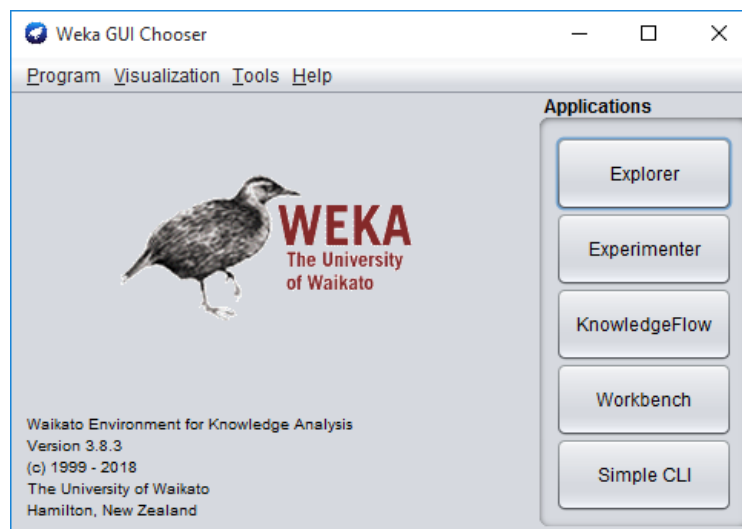


```

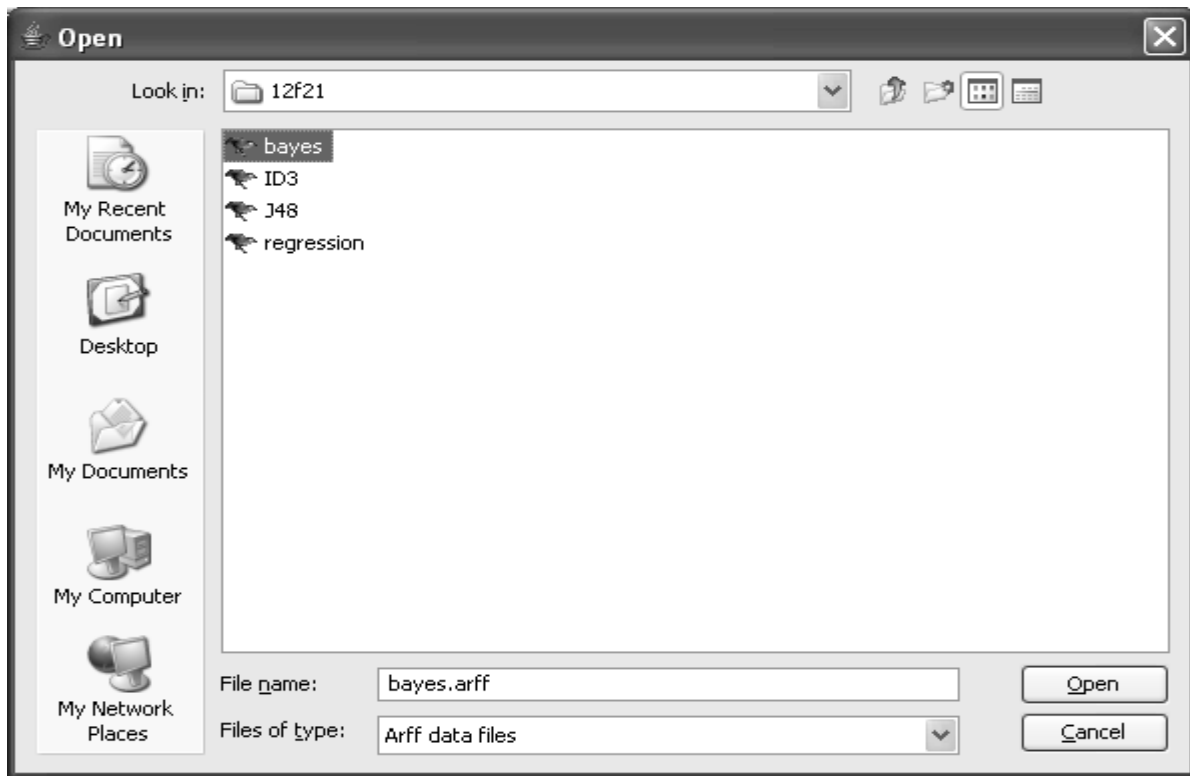
bayer - Notepad
File Edit Format View Help
% TITLE:decision Tree:Bayes
% Date:09-09-2014
% Creator by:M.MANIKANTA SAI
@relation Bayes
@attribute age {<=30,31...40,>40}
@attribute income {high,medium,low}
@attribute student{no,yes}
@attribute credit_rating{excellent,fair}
@attribute buys_computer{yes,no}
%
%
@data
<=30,high,no,fair,no
<=30,high,no,excellent,no
31...40,high,no,fair,yes
>40,medium,no,fair,yes
>40,low,yes,excellent,no
31...40,low,yes,excellent,yes
<=30,medium,no,fair,no
<=30,low,yes,fair,yes
>40,medium,yes,fair,yes
<=30,medium,yes,excellent,yes
31...40,medium,no,excellent,yes
31...40,high,yes,fair,yes
>40,medium,no,excellent,no

```

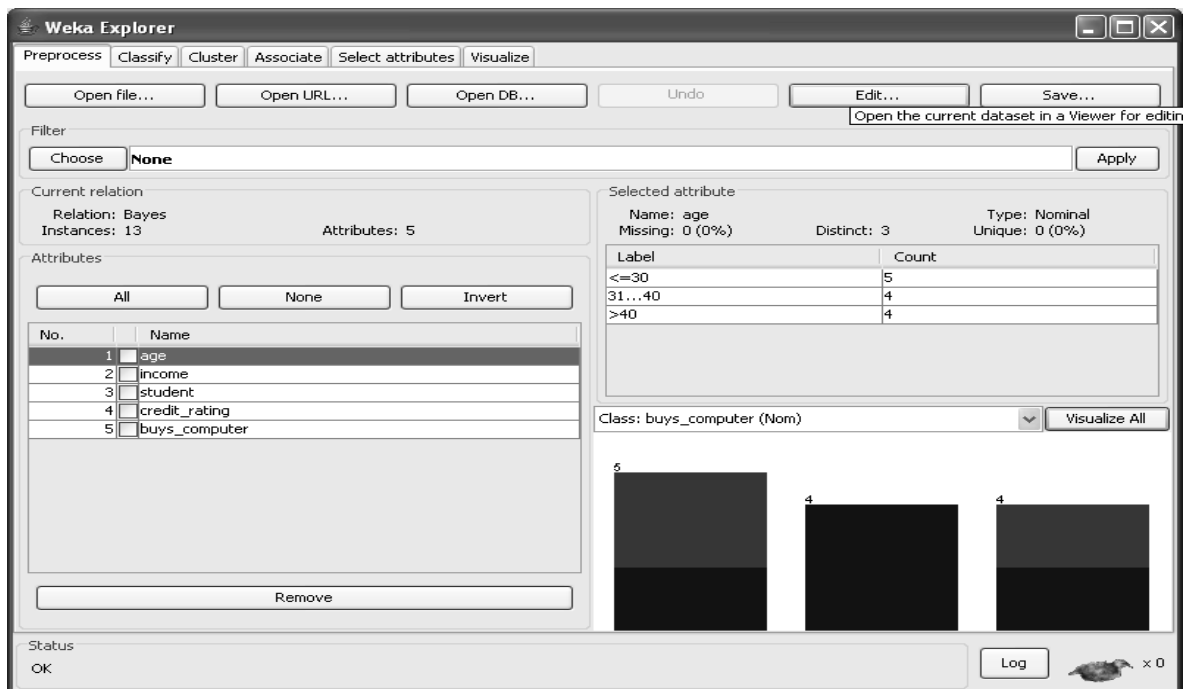
Step2: Go to Start menu and select Weka 3.8.3, on that select Weka 3.8(with console).



Step3:Go to Preprocessor and select the open file.



Step4: In the Preprocessor menu select Edit.



Viewer

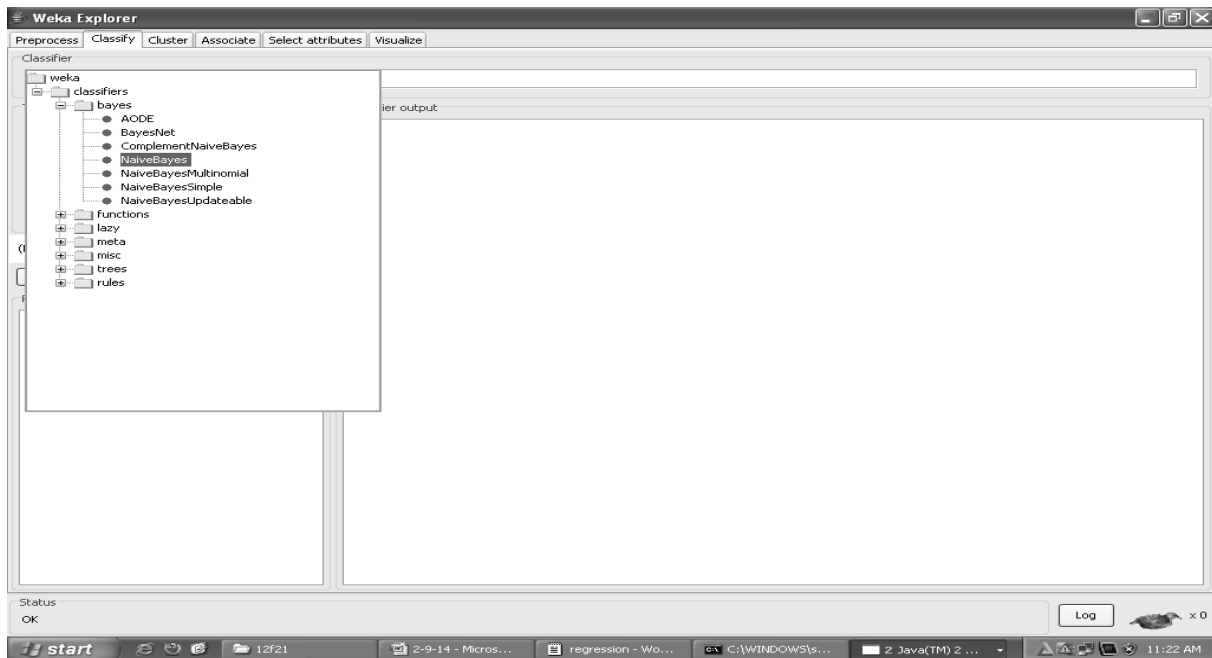
Relation: Bayes

No.	age Nominal	income Nominal	student Nominal	credit-rating Nominal	buys-computer Nominal
1	<=30	high	no	fair	no
2	<=30	high	no	excellent	no
3	31..40	high	no	fair	yes
4	>40	medium	no	fair	yes
5	>40	low	yes	fair	yes
6	>40	low	yes	excellent	no
7	31..40	low	yes	excellent	yes
8	<=30	medium	no	fair	no
9	<=30	low	yes	fair	yes
10	>40	medium	yes	fair	yes
11	<=30	medium	yes	excellent	yes
12	31..40	medium	no	excellent	yes
13	31..40	high	yes	fair	yes
14	>40	medium	no	excellent	no

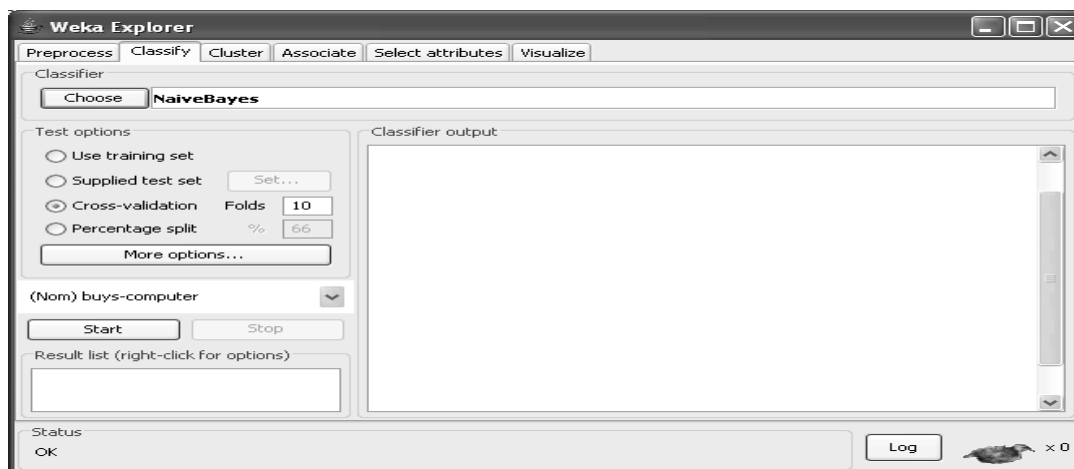
Right click (or left+alt)

Undo OK Cancel

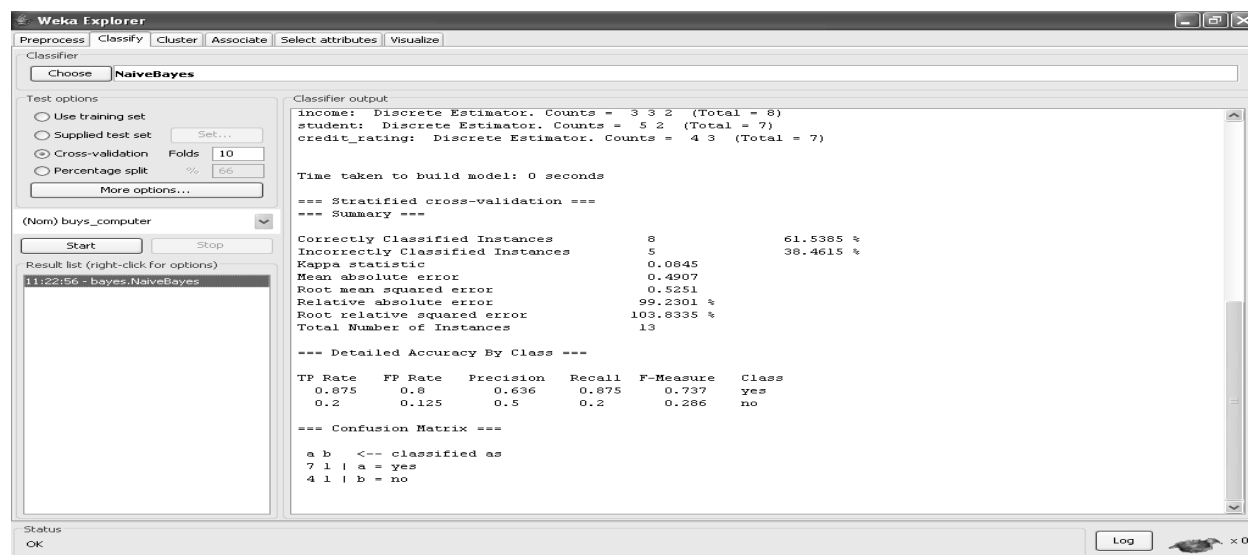
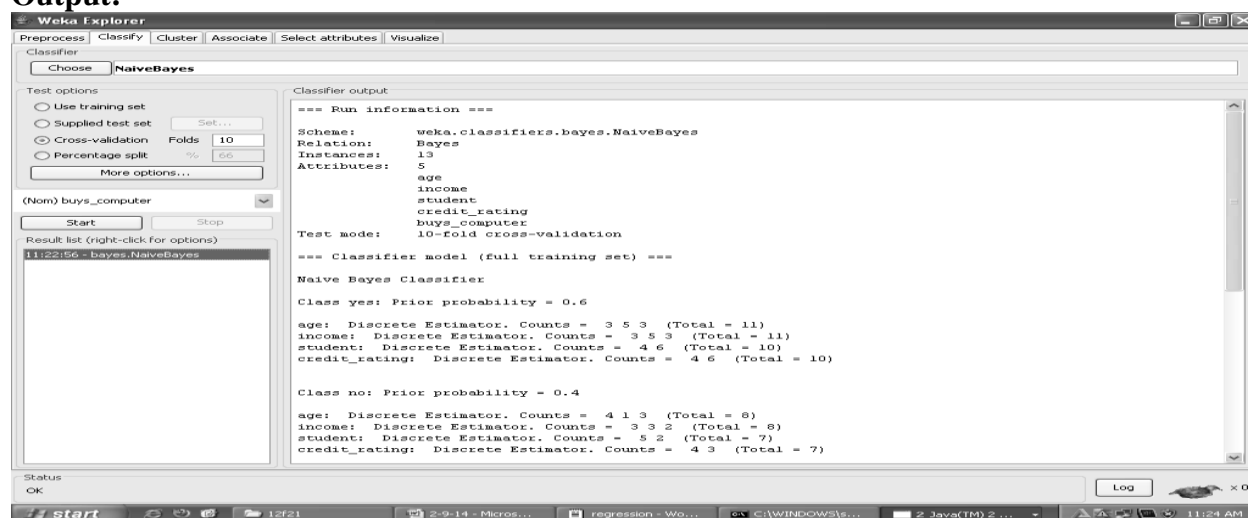
Step5: Go to **classify** tab menu ,click no **choose** button and then select **Trees** from the list Click on the **NaiveBayes** option.



Step 6:click on the **start** button to view the result.



Output:



VIVA-QUESTIONS

1. Formulae for Information gain=_____
2. Formulae for GainRatio=_____
3. Compute entropy value of(9,5) = _____
4. Classification is supervised learning.
5. _____ measure is used to select the test attribute at each node in the decision tree.
6. Posterior probability can be calculated by _____ theorem.
7. Decision tree is a type of _____ algorithm.
8. KNN stand for
9. How does KNN calculate distance?
10. ID3 stands for

Experiment No. 4: Demonstrate performing association rule mining on data sets

- i. Load transaction dataset into WEKA and run Apriori algorithm with different support and confidence values.
- ii. The Apriori algorithm uses a generate and count strategy for deriving frequent items sets and generate association rules.

The Apriori Algorithm

Apriori is a seminal algorithm proposed by R. Agrawal and R. Srikant in 1994 for mining frequent itemsets for Boolean association rules. The name of the algorithm is based on the fact that the algorithm uses *prior knowledge* of frequent itemset properties, as we shall see following. Apriori employs an iterative approach known as a *level-wise* search, where k -itemsets are used to explore $(k+1)$ -itemsets. First, the set of frequent 1-itemsets is found by scanning the database to accumulate the count for each item, and collecting those items that satisfy minimum support. The resulting set is denoted L_1 . Next, L_1 is used to find L_2 , the set of frequent 2-itemsets, which is used to find L_3 , and so on, until no more frequent k -itemsets can be found. The finding of each L_k requires one full scan of the database.

Apriori property: *All nonempty subsets of a frequent itemset must also be frequent.* The Apriori property is based on the following observation. By definition, if an itemset I does not satisfy the minimum support threshold, $min\ sup$, then I is not frequent; that is, $P(I) < min\ sup$. If an item A is added to the itemset I , then the resulting itemset (i.e., $I \cup \{A\}$) cannot occur more frequently than I . Therefore, $I \cup \{A\}$ is not frequent either; that is, $P(I \cup \{A\}) < min\ sup$.

The join step: To find L_k , a set of candidate k -itemsets is generated by joining L_{k-1} with itself. This set of candidates is denoted C_k . Let l_1 and l_2 be itemsets in L_{k-1} . The notation $li[j]$ refers to the j th item in li (e.g., $l_1[k-2]$ refers to the second to the last item in l_1). For the $(k-1)$ -itemset, li , this means that the items are sorted such that $li[1] < li[2] < \dots < li[k-1]$. The join, L_{k-1} on L_{k-1} , is performed, where members of L_{k-1} are joinable if their first $(k-2)$ items are in common. That is, members l_1 and l_2 of L_{k-1} are joined if $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2]) \wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$. The condition $l_1[k-1] < l_2[k-1]$ simply ensures that no duplicates are generated. The resulting itemset formed by joining l_1 and l_2 is $l_1[1], l_1[2], \dots, l_1[k-2], l_1[k-1], l_2[k-1]$.

The prune step: C_k is a superset of L_k , that is, its members may or may not be frequent, but all of the frequent k -itemsets are included in C_k . A scan of the database to determine the count of each

candidate in C_k would result in the determination of L_k (i.e., all candidates having a count no less than the minimum support count are frequent by definition, and therefore belong to L_k).

Transactional data for an *Allelectronics* branch.

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

In the first iteration of the algorithm, each item is a member of the set of candidate 1-itemsets, C_1 . The algorithm simply scans all of the transactions in order to count the number of occurrences of each item.

2. Suppose that the minimum support count required is 2, that is, $min\ sup = 2$. (Here, we are referring to *absolute* support because we are using a support count. The corresponding relative support is $2/9 = 22\%$). The set of frequent 1-itemsets, L_1 , can then be determined. It consists of the candidate 1-itemsets satisfying minimum support. In our example, all of the candidates in C_1 satisfy minimum support.

3. To discover the set of frequent 2-itemsets, L_2 , the algorithm uses the join L_1 on L_1 to generate a candidate set of 2-itemsets, C_2 . C_2 consists of $\square_j L_1 \square_j$ 2-itemsets. Note that no candidates are removed from C_2 during the prune step because each subset of the candidates is also frequent.

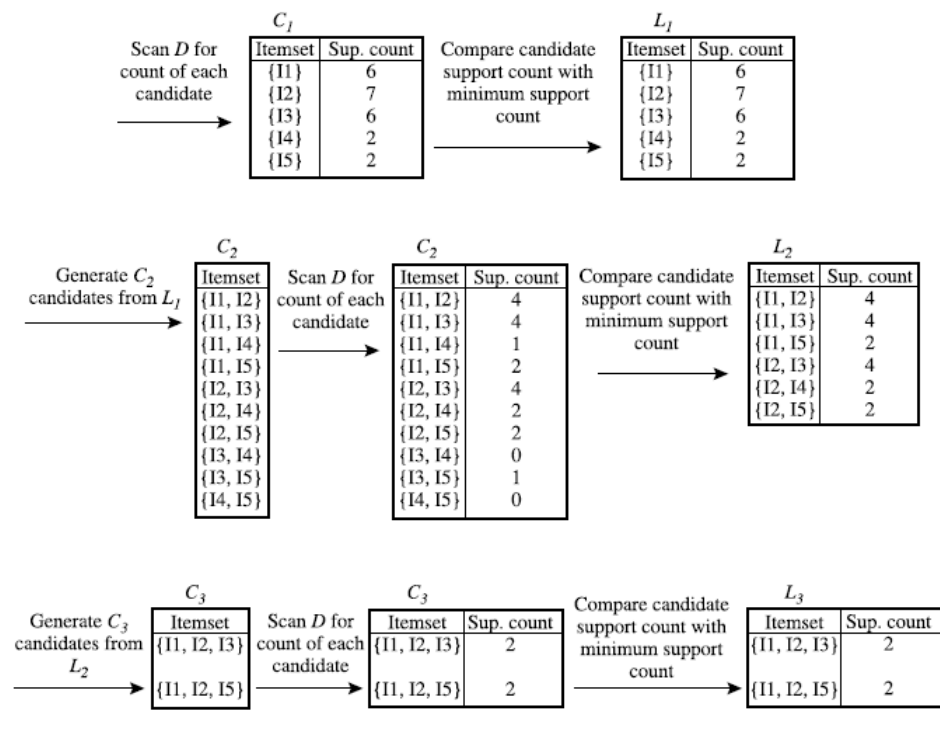
4. Next, the transactions in D are scanned and the support count of each candidate itemset in C_2 is accumulated.

5. The set of frequent 2-itemsets, L_2 , is then determined, consisting of those candidate 2-itemsets in C_2 having minimum support.

6. The generation of the set of candidate 3-itemsets, C_3 . From the join step, we first get $C_3 = L_2 \text{ on } L_2 = \{I_1, I_2, I_3\}, \{I_1, I_2, I_5\}, \{I_1, I_3, I_5\}, \{I_2, I_3, I_4\}, \{I_2, I_3, I_5\}, \{I_2, I_4, I_5\}$. Based on the Apriori property that all subsets of a frequent itemset must also be frequent, we can determine that the four latter candidates cannot possibly be frequent. We therefore remove them from C_3 , thereby saving the effort of unnecessarily obtaining their counts during the subsequent scan of D to determine L_3 . Note that when given a candidate k -itemset, we only need to check if its $(k-1)$ -subsets are frequent since the Apriori algorithm uses a level-wise search strategy.

7. The transactions in D are scanned in order to determine L_3 , consisting of those candidate 3-itemsets in C_3 .

8. The algorithm uses L_3 on L_3 to generate a candidate set of 4-itemsets, C_4 . Although the join results in $\{I_1, I_2, I_3, I_5\}$, this itemset is pruned because its subset $\{I_2, I_3, I_5\}$ is not frequent. Thus, $C_4 = \emptyset$, and the algorithm terminates, having found all of the frequent itemsets.



Algorithm: Apriori. Find frequent itemsets using an iterative level-wise approach based on candidate generation.

Input:

- D , a database of transactions;
- min_sup , the minimum support count threshold.

Output: L , frequent itemsets in D .

Method:

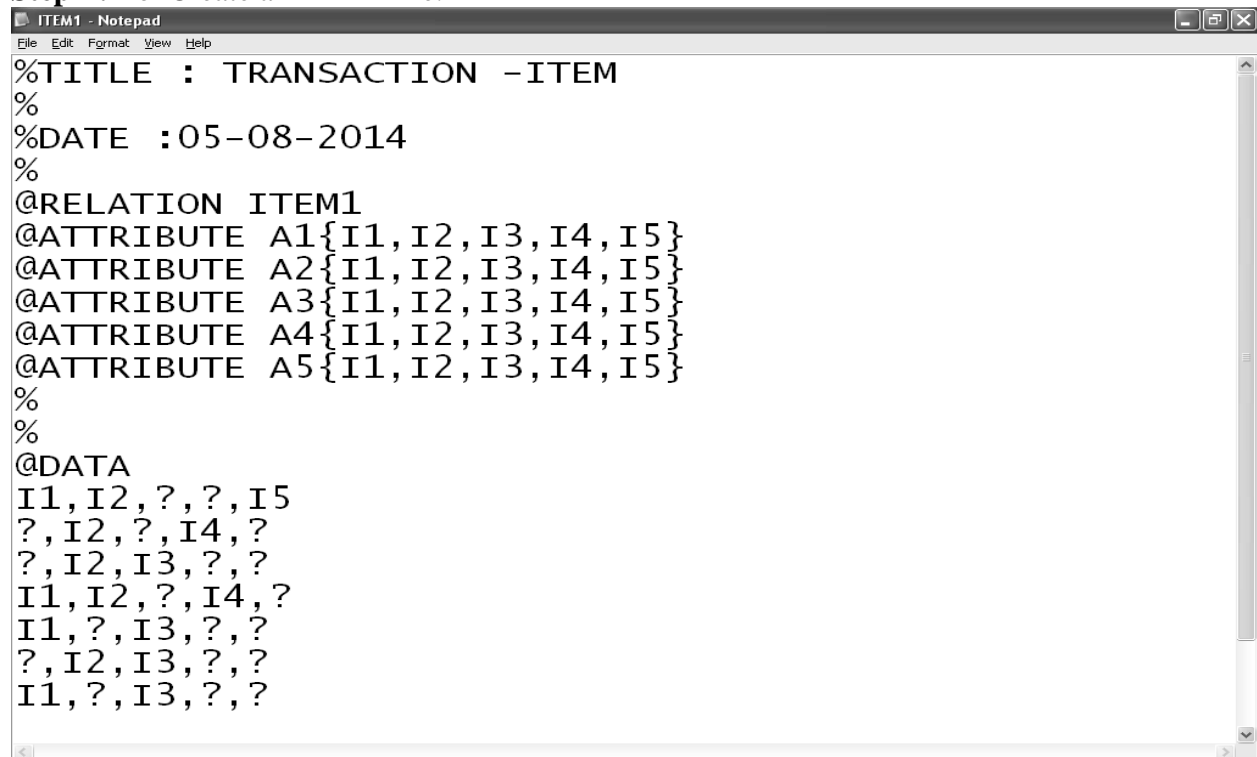
```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;
(2) for ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) {
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates
(6)     for each candidate  $c \in C_t$ 
(7)        $c.\text{count}++$ ;
(8)   }
(9)    $L_k = \{c \in C_k | c.\text{count} \geq min\_sup\}$ 
(10) }
(11) return  $L = \cup_k L_k$ ;

procedure apriori_gen( $L_{k-1}$ :frequent ( $k-1$ )-itemsets)
(1) for each itemset  $t_1 \in L_{k-1}$ 
(2)   for each itemset  $t_2 \in L_{k-1}$ 
(3)     if ( $(t_1[1] = t_2[1]) \wedge (t_1[2] = t_2[2]) \wedge \dots \wedge (t_1[k-2] = t_2[k-2]) \wedge (t_1[k-1] < t_2[k-1])$ ) then {
(4)        $c = t_1 \bowtie t_2$ ; // join step: generate candidates
(5)       if has_infrequent_subset( $c, L_{k-1}$ ) then
(6)         delete  $c$ ; // prune step: remove unfruitful candidate
(7)       else add  $c$  to  $C_k$ ;
(8)     }
(9) return  $C_k$ ;

procedure has_infrequent_subset( $c$ : candidate  $k$ -itemset;
                                 $L_{k-1}$ : frequent ( $k-1$ )-itemsets); // use prior knowledge
(1) for each ( $k-1$ )-subset  $s$  of  $c$ 
(2)   if  $s \notin L_{k-1}$  then
(3)     return TRUE;
(4) return FALSE;
```

Association rule mining using the Apriori-Algorithm

Step-1: To Create an ARFF File.

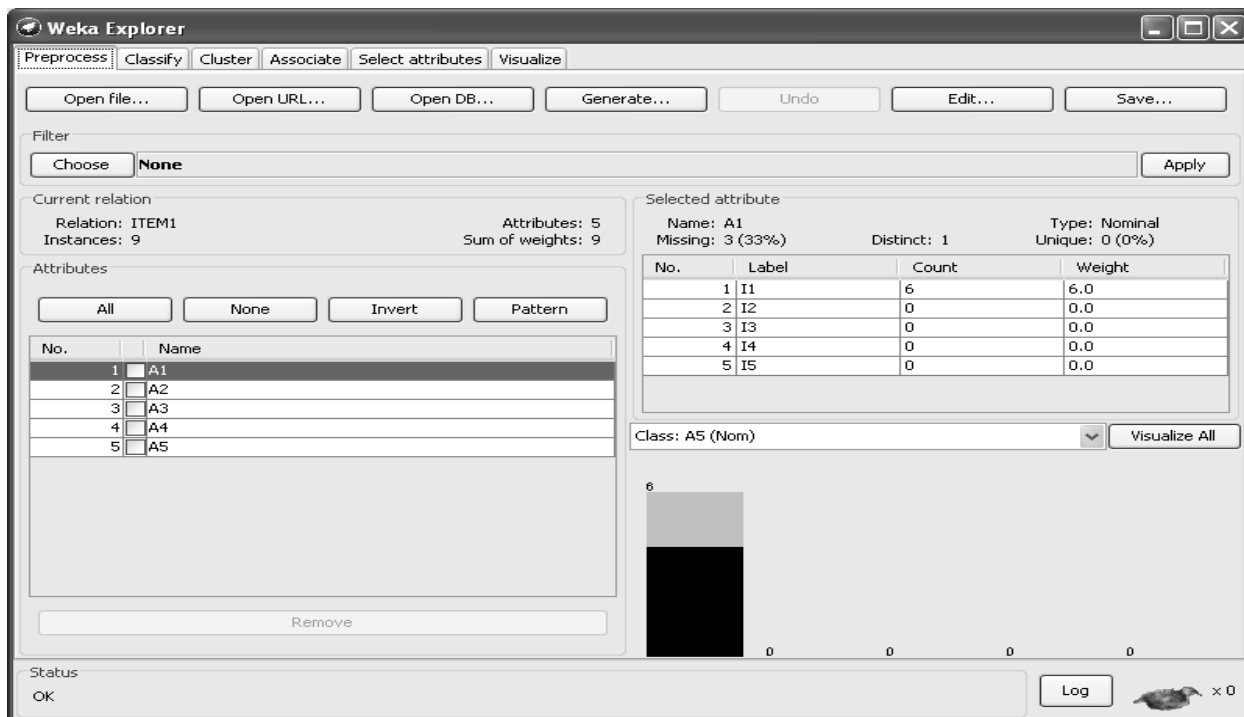


```
ITEM1 - Notepad
File Edit Format View Help
%TITLE : TRANSACTION -ITEM
%
%DATE :05-08-2014
%
@RELATION ITEM1
@ATTRIBUTE A1{I1,I2,I3,I4,I5}
@ATTRIBUTE A2{I1,I2,I3,I4,I5}
@ATTRIBUTE A3{I1,I2,I3,I4,I5}
@ATTRIBUTE A4{I1,I2,I3,I4,I5}
@ATTRIBUTE A5{I1,I2,I3,I4,I5}
%
%
@DATA
I1,I2,?,?,I5
?,I2,?,I4,?
?,I2,I3,?,?
I1,I2,?,I4,?
I1,?,I3,?,?
?,I2,I3,?,?
I1,?,I3,?,?
```

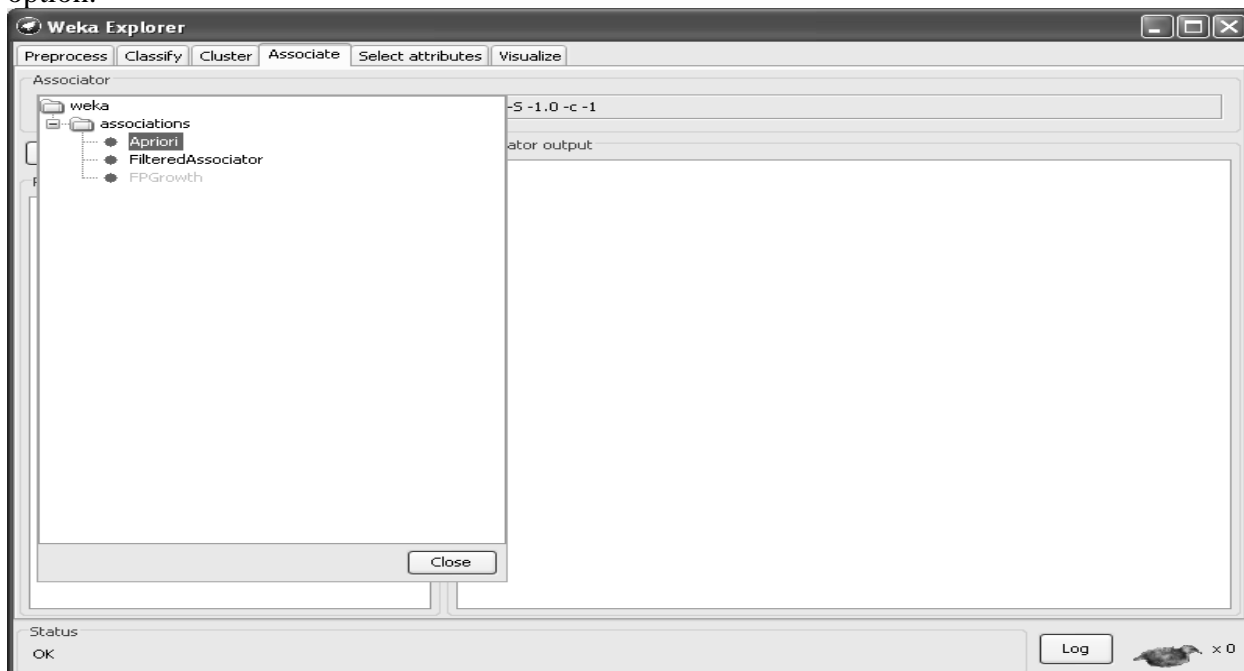
Step-2: -Save as the file .ARFF file format

Step-3: - Go to All Programs in the start menu bar and then choose weka 3.8.3 in the weka 3.8(with console).

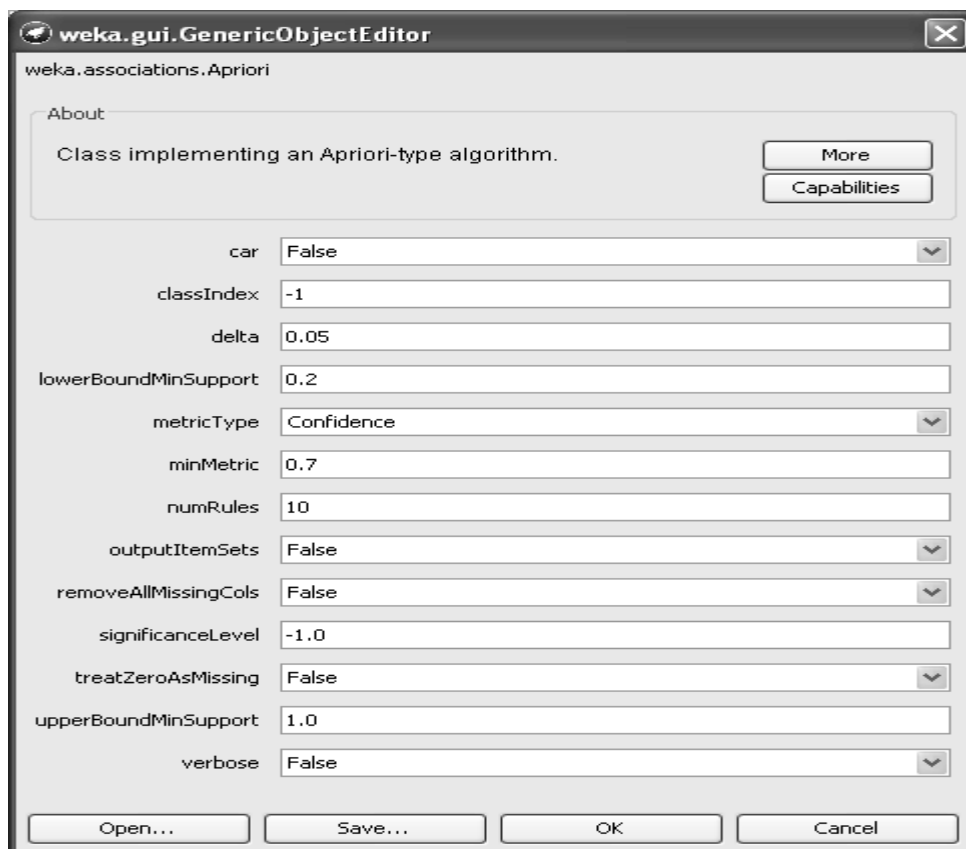
Step-4: In the weka window to choose the Explorer and then select the Preprocess in the weka explorer window.



Step-5: Go to Association option in the menu and then select choose button to select the Apriori option.

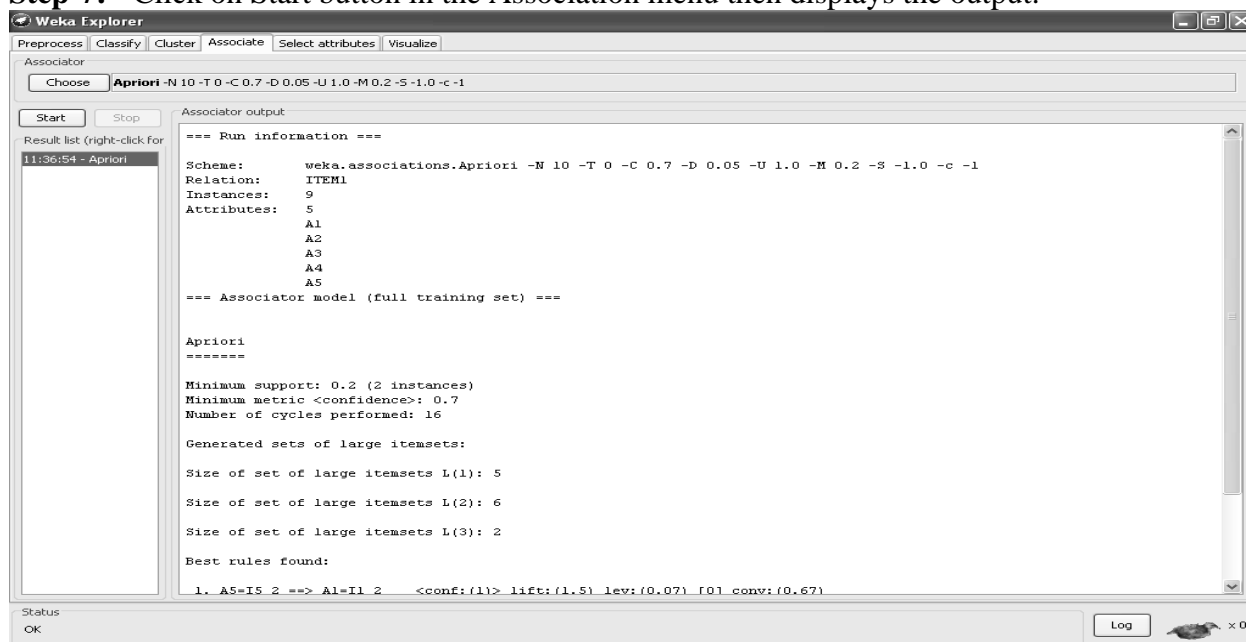


Step-6: - Click on Apriori field to configure the various options.



Click on OK button.

Step-7: - Click on Start button in the Association menu then displays the output.



Step-8: - Various options are calculated for minimum support value like this.

The image shows two windows from the Weka software. The top window is the 'weka.gui.GenericObjectEditor' for the 'weka.associations.Apriori' class. It displays various configuration options for the Apriori algorithm:

- car: False
- classIndex: -1
- delta: 0.05
- lowerBoundMinSupport: 0.3
- metricType: Confidence
- minMetric: 0.7
- numRules: 10
- outputItemSets: False
- removeAllMissingCols: False
- significanceLevel: -1.0
- treatZeroAsMissing: False
- upperBoundMinSupport: 1.0
- verbose: False

The bottom window is 'Weka Explorer' showing the 'Associator' tab. The 'Choose' button is set to 'Apriori -N 10 -T 0 -C 0.7 -D 0.05 -U 1.0 -M 0.3 -S -1.0 -c -1'. The 'Associator output' pane shows the following text:

```

A4
A5
=== Associator model (full training set) ===

Apriori
=====

Minimum support: 0.3 (3 instances)
Minimum metric <confidence>: 0.7
Number of cycles performed: 14

Generated sets of large itemsets:

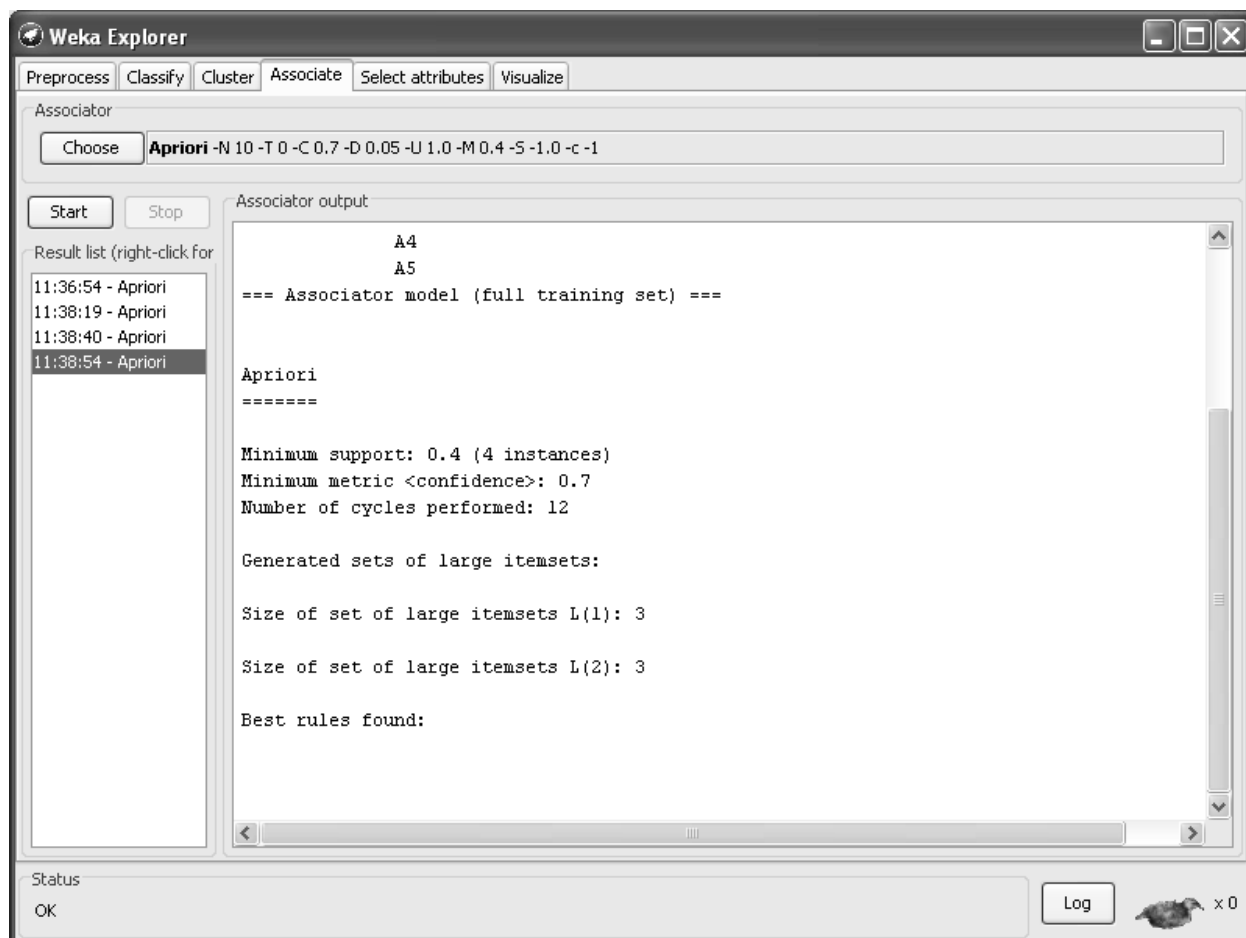
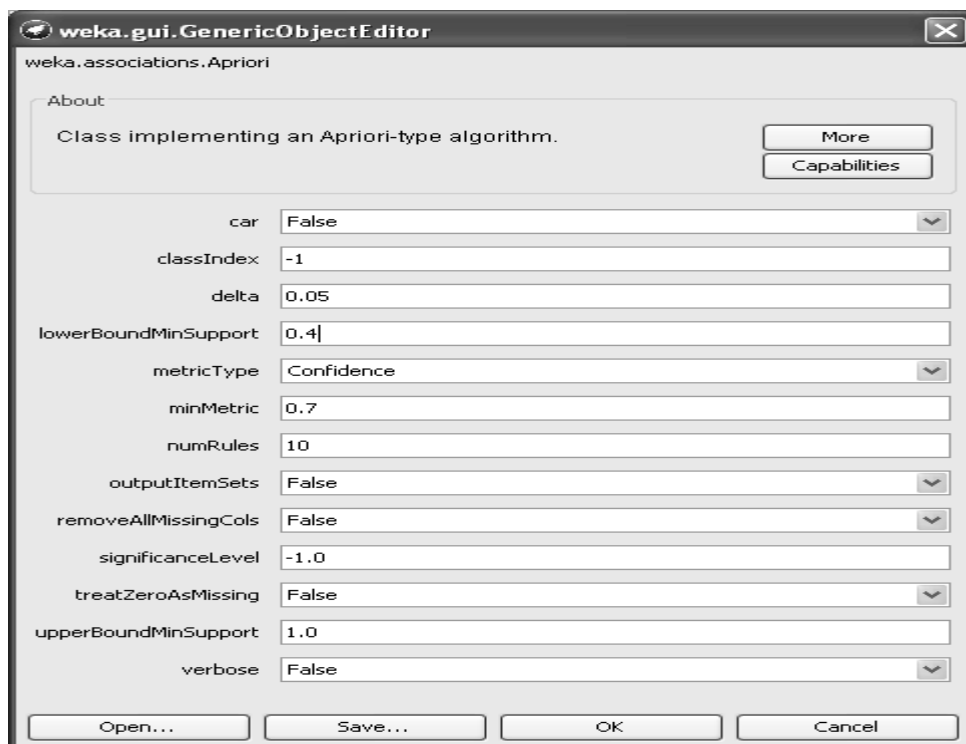
Size of set of large itemsets L(1): 3

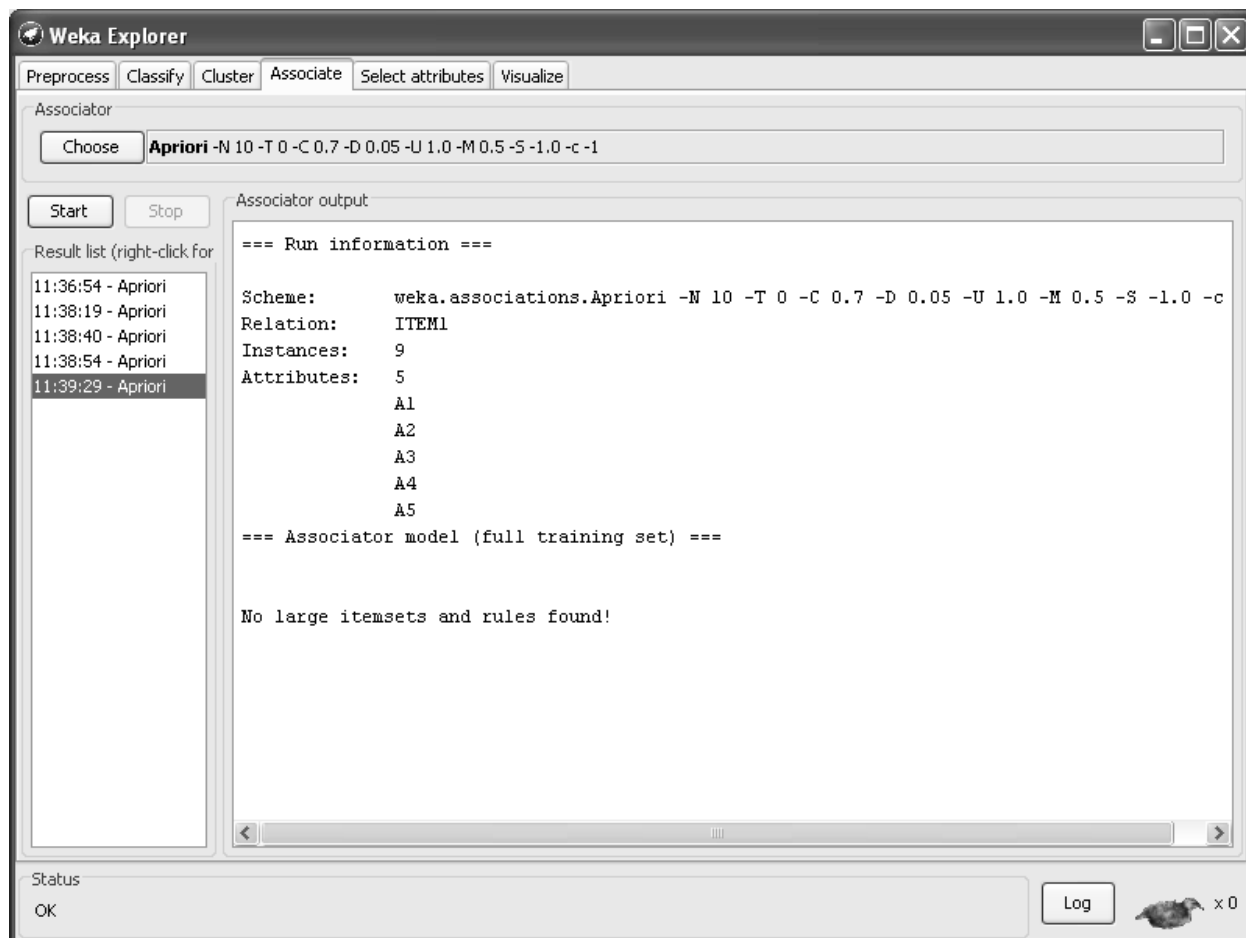
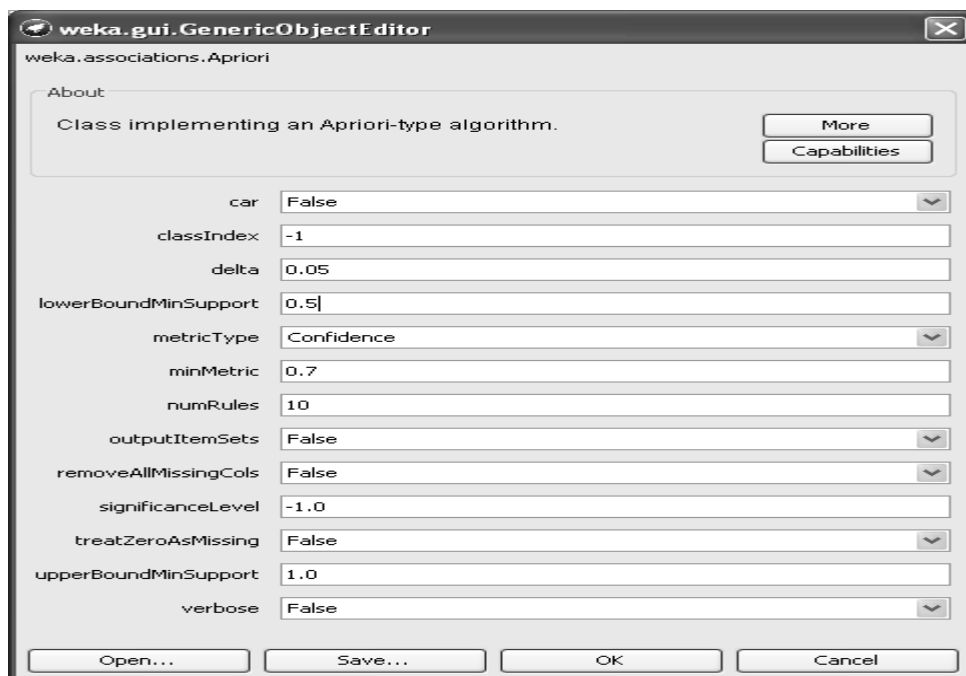
Size of set of large itemsets L(2): 3

Best rules found:

```

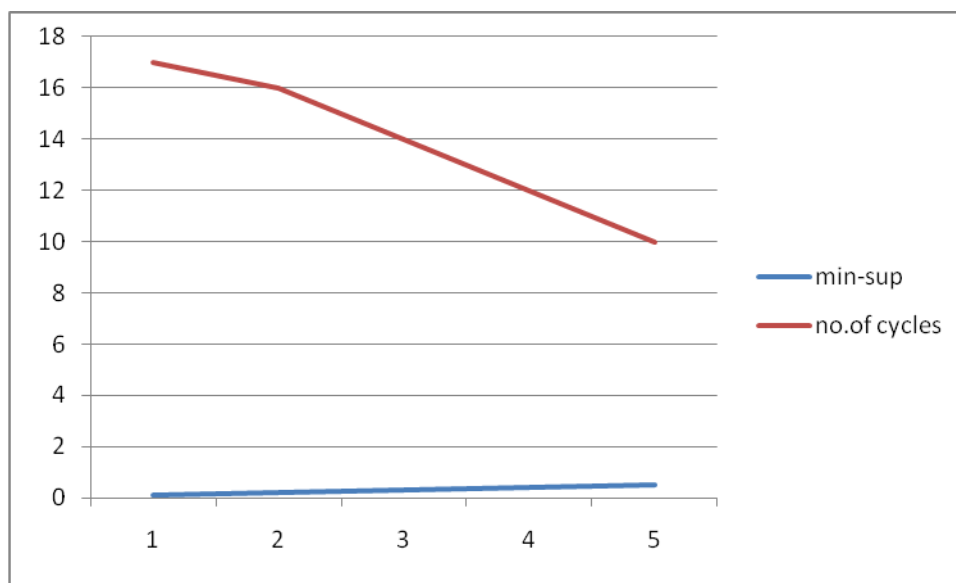
The 'Result list' on the left shows two entries for 'Apriori' with timestamps 11:36:54 and 11:38:19. The status bar at the bottom indicates 'OK' and a 'Log' button.





Graph for Apriori Algorithm: -

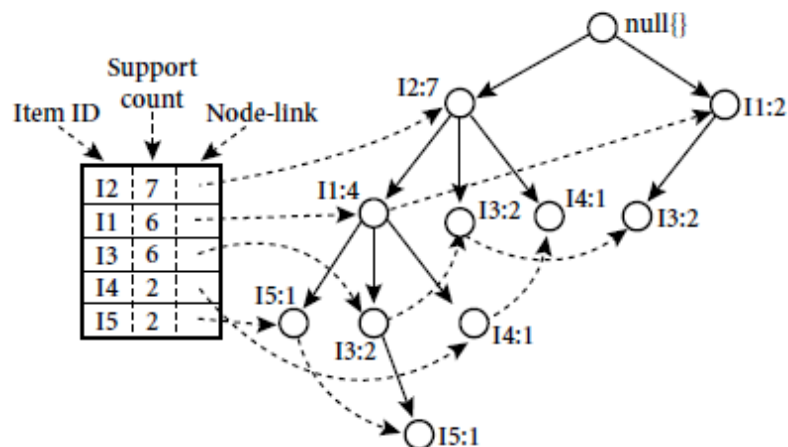
min-sup	no.of cycles
0.1	17
0.2	16
0.3	14
0.4	12
0.5	10

**FP-growth.**

The first scan of the database is the same as Apriori, which derives the set of frequent items (1-itemsets) and their support counts (frequencies). Let the minimum support count be 2. The set of frequent items is sorted in the order of descending support count. This resulting set or *list* is denoted L . Thus, we have $L = I2:7, I1: 6, I3: 6, I4: 2, I5: 2$.

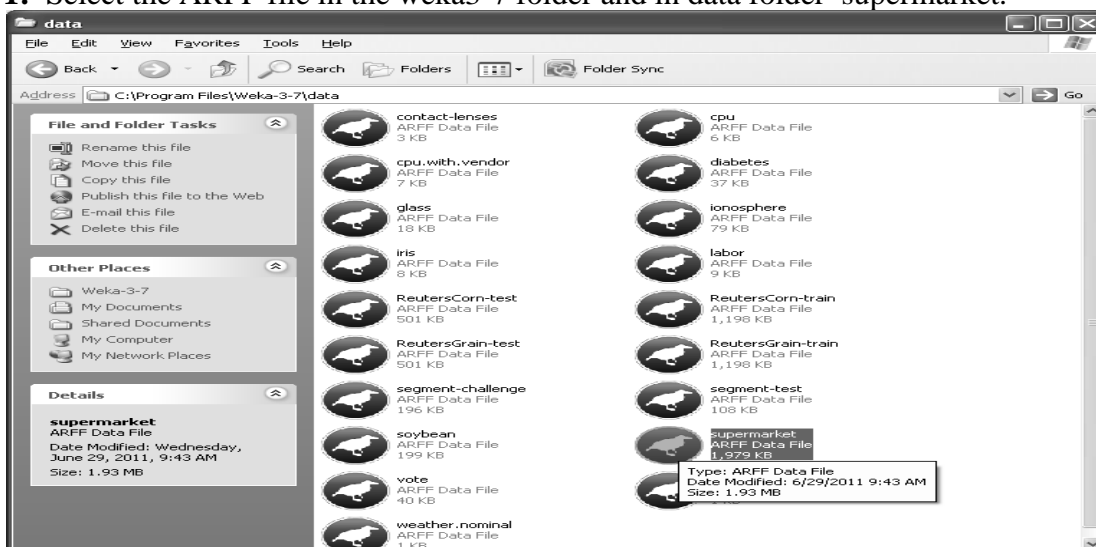
An FP-tree is then constructed as follows. First, create the root of the tree, labeled with “null.” Scan database D a second time. The items in each transaction are processed in L order (i.e., sorted according to descending support count), and a branch is created for each transaction. For example, the scan of the first transaction, “T100: I1, I2, I5,” which contains three items (I2, I1, I5 in L order), leads to the construction of the first branch of the tree with three nodes, I2: 1, I1:1, and I5: 1, where I2 is linked as a child of the root, I1 is linked to I2, and I5 is linked to I1. The second transaction, T200, contains the items I2 and I4 in L order, which would result in a branch where I2 is linked to the root and I4 is linked to I2. However, this branch would share a common prefix, I2, with the existing path for T100. Therefore, we instead increment the count of the I2

node by 1, and create a newnode, I4: 1,which is linked as a child of I2: 2. In general,when considering the branch to be added for a transaction, the count of each node along a common prefix is incremented by 1, and nodes for the items following the prefix are created and linked accordingly.

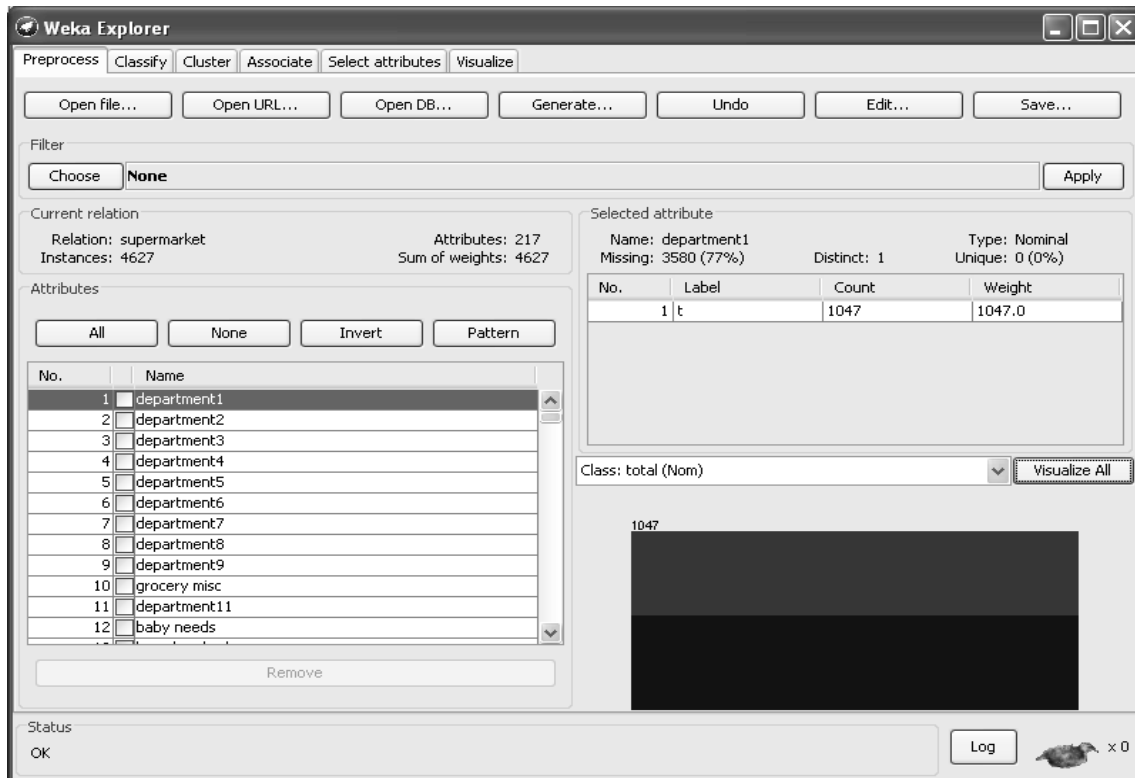


Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	{{I2, I1: 1}, {I2, I1, I3: 1}}	$\langle I2: 2, I1: 2 \rangle$	{I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2}
I4	{{I2, I1: 1}, {I2: 1}}	$\langle I2: 2 \rangle$	{I2, I4: 2}
I3	{{I2, I1: 2}, {I2: 2}, {I1: 2}}	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	{I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2}
I1	{{I2: 4}}	$\langle I2: 4 \rangle$	{I2, I1: 4}

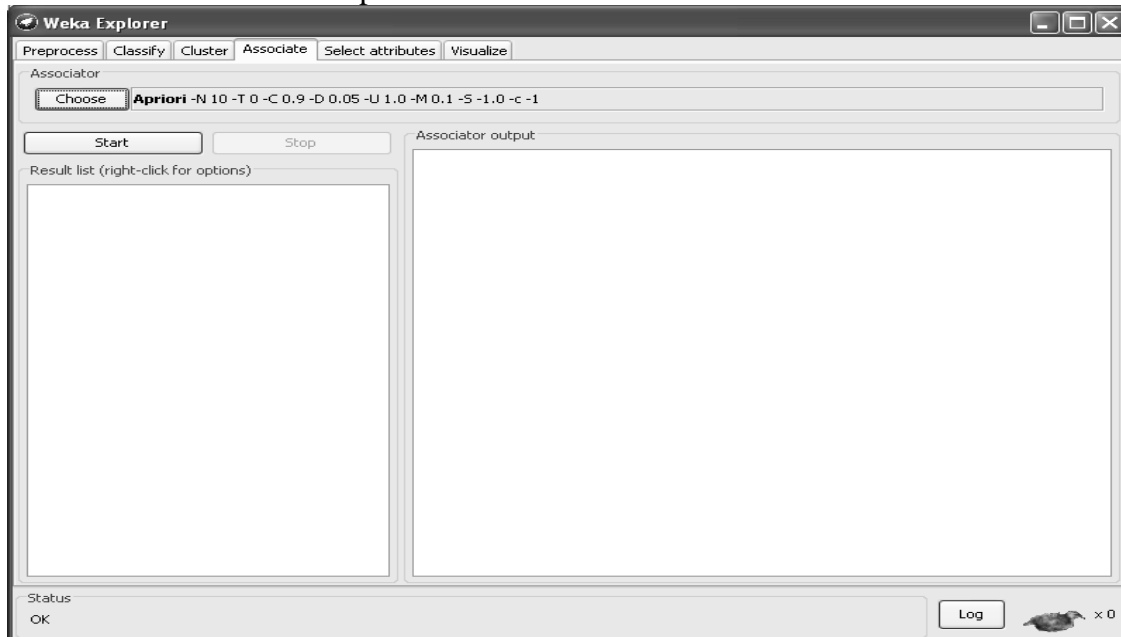
Step 1: Select the ARFF file in the weka3-7 folder and in data folder supermarket.



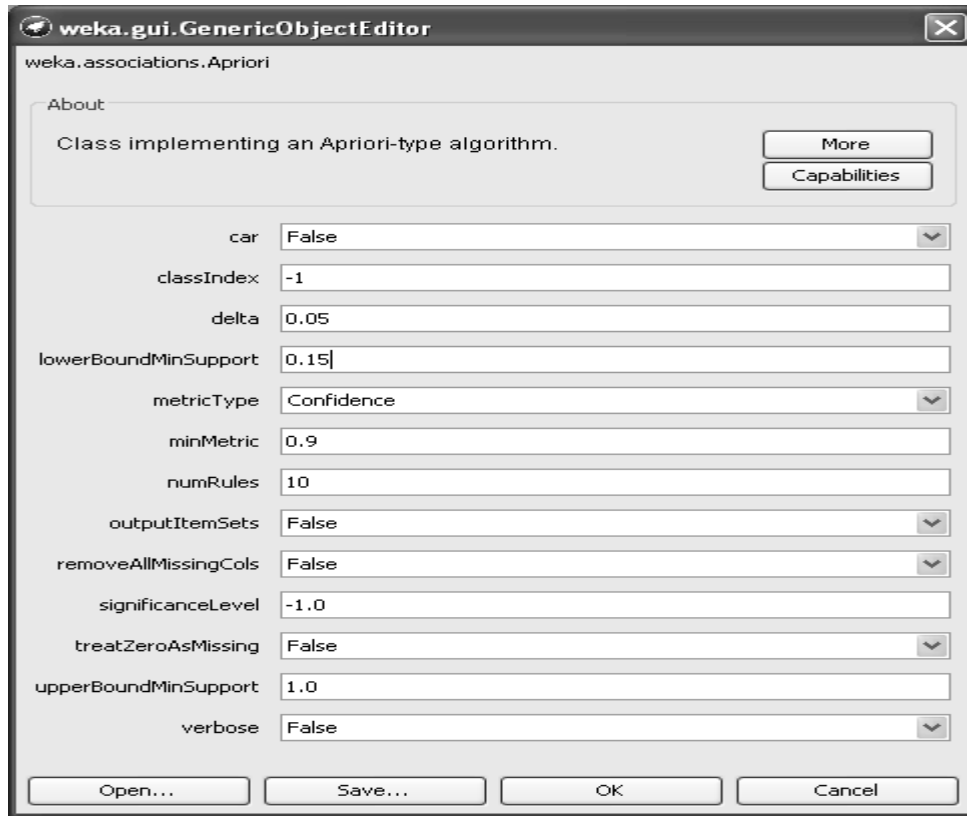
Double Click on the supermarket file and then open.



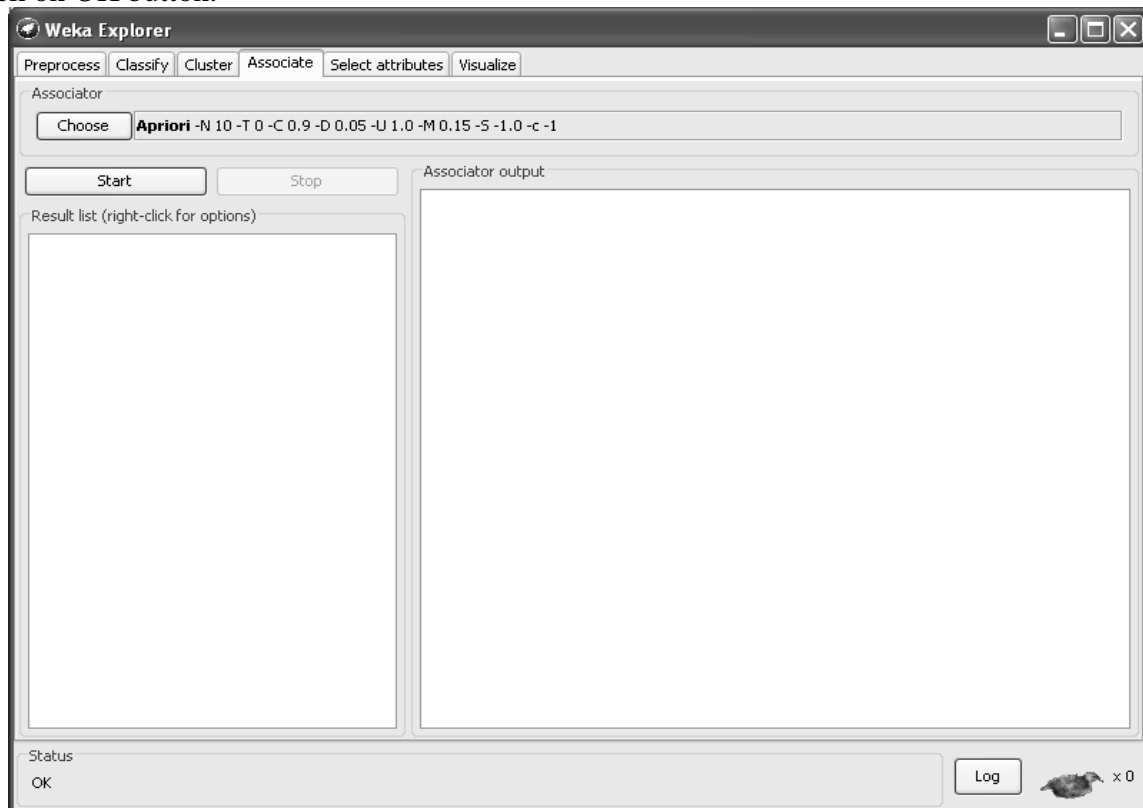
Step 2: Go to the Association option in the menu bar.

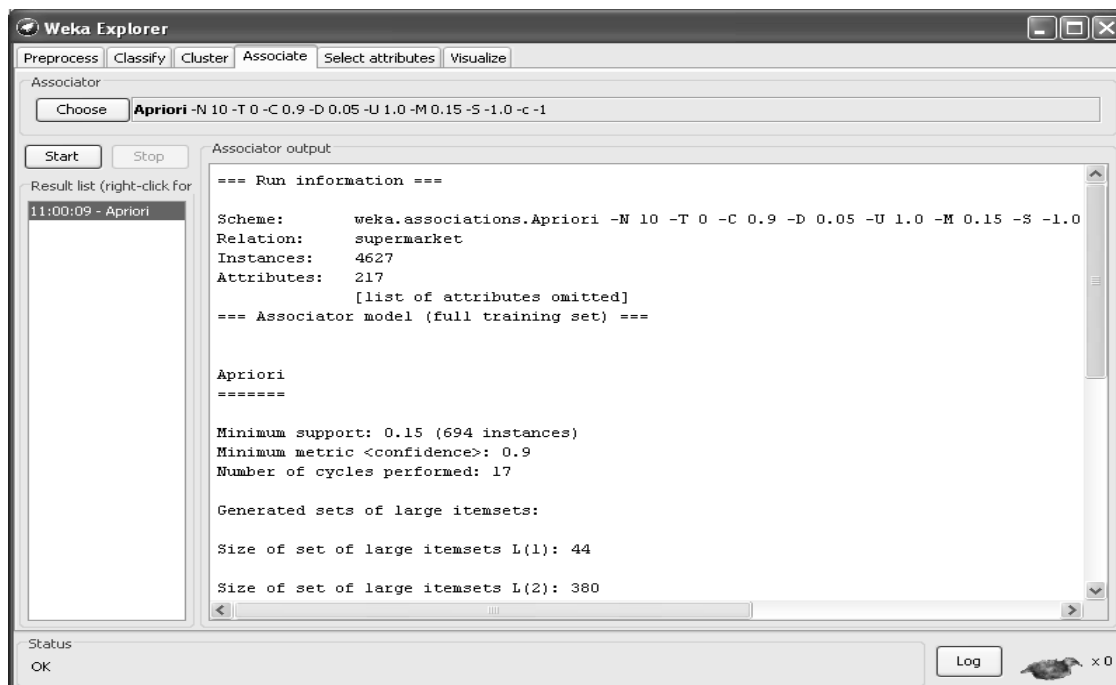


Step 3: Click on choose button and configure the various options.

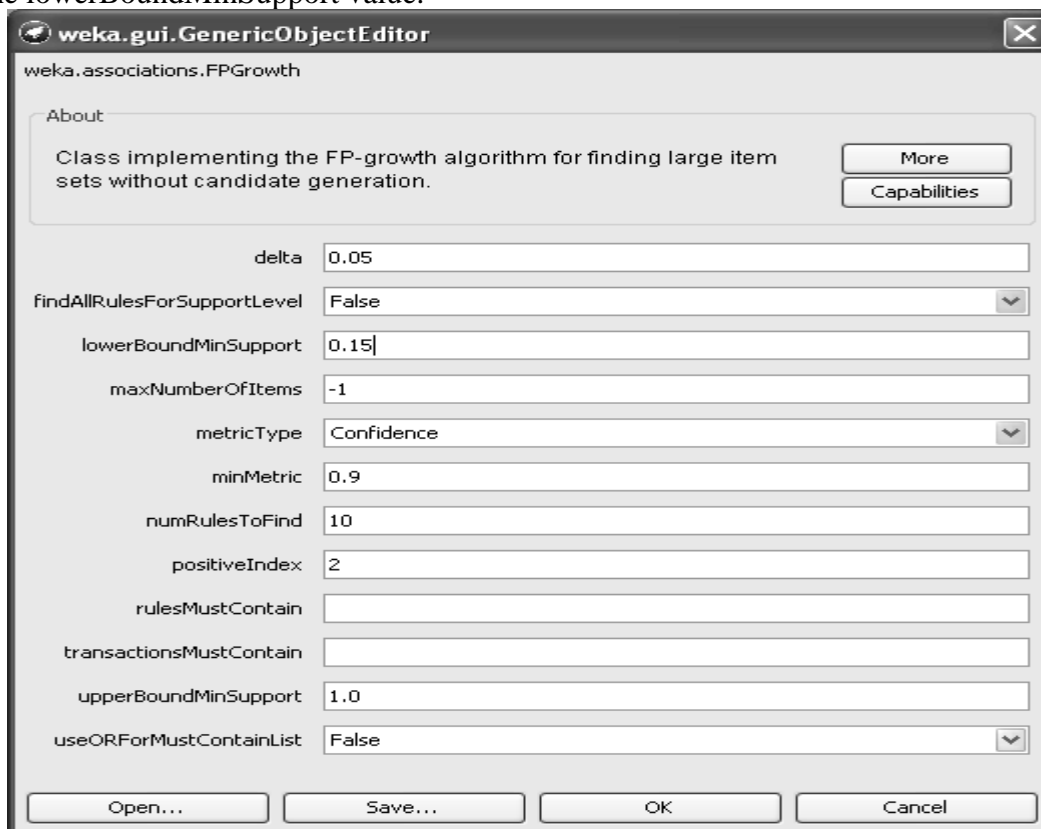


Click on OK button.

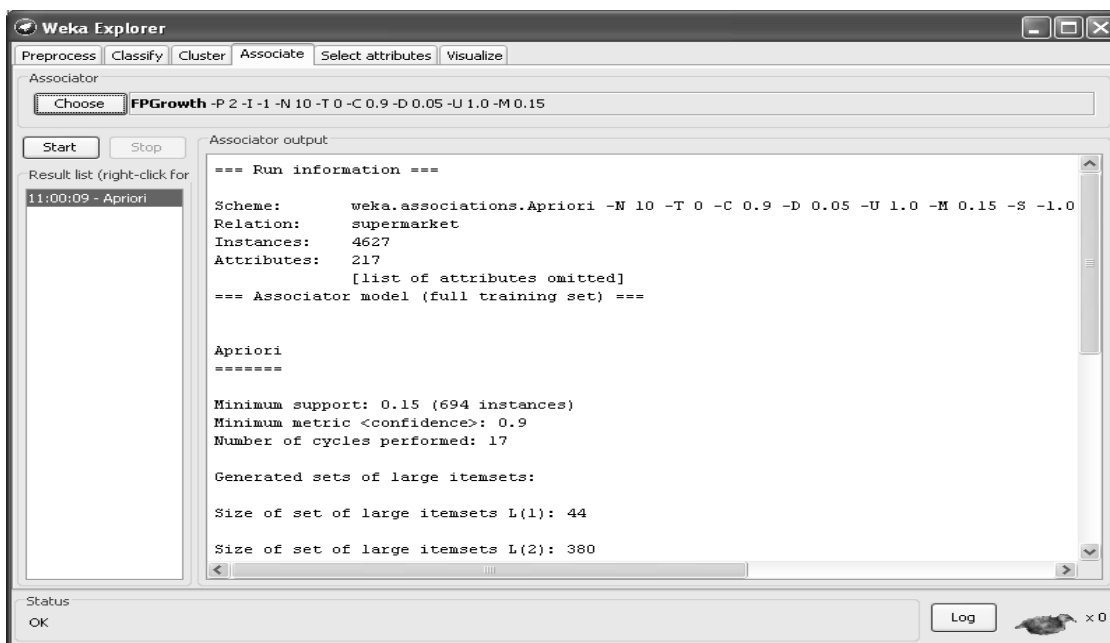




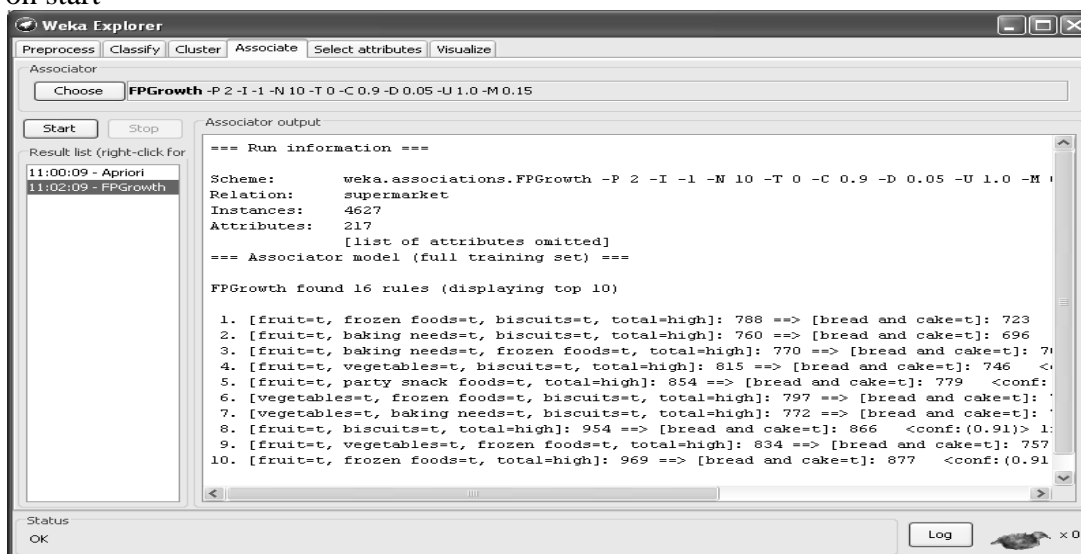
Give the lowerBoundMinSupport value.



Step 4: Click on choose button and then select the FP-growth Algorithm and configure to various minimum support values.



Click on start



min-sup	apriori	fp-growth
0.15	0.44	0.2
0.2	0.3	0.1
0.3	0.2	0.1
0.4	0.15	0.1
0.5	0.11	0.1
0.6	0.8	0.1
0.7	0.6	0.1
0.8	0.5	0.1
0.9	0.3	0.1

VIVA-QUESTIONS

1. Market-basket problem was formulated by_____.
2. The Market Basket Analysis is a typical example of_____
3. An itemset satisfying the support criterion is known as:
4. In Apriori algorithm, for generating k-item sets, we use:
5. Mathematically, $\text{Support}(\{X\} \rightarrow \{Y\})$ is
6. All nonempty subset of a frequent itemset must also be frequent is _____ property
7. FP-Growth algorithm is based on _____.
8. The formula for $\text{Support}(A \Rightarrow B)=$
9. The formula for $\text{Confidence}(A \Rightarrow B)=$
10. Apriori method mines the frequent itemsets without candidate generation

Experiment No. 5: Demonstrate performing regression on data sets

- i. Load numeric dataset into WEKA and build linear regression model.
- ii. Explore simple linear regression technique that only looks at one variable

Linear Regression

Straight-line regression analysis involves a response variable, y , and a single predictor variable, x . It is the simplest form of regression, and models y as a linear function of x . That is,

$$y = b + wx;$$

where the variance of y is assumed to be constant, and b and w are regression coefficients specifying the Y-intercept and slope of the line, respectively. The regression coefficients, w and b , can also be thought of as weights, so that we can equivalently write,

$$y = w_0 + w_1x$$

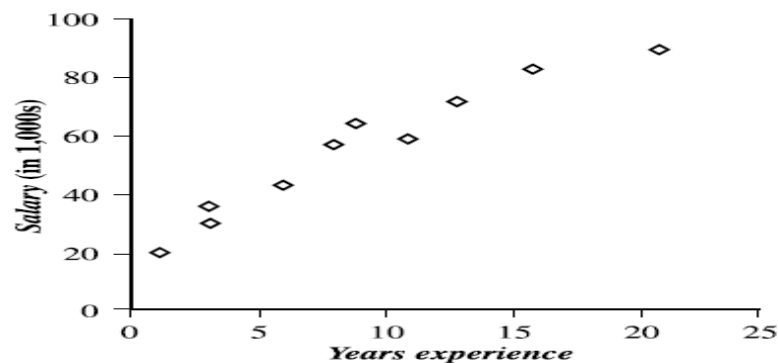
These coefficients can be solved for by the method of least squares, which estimates the best-fitting straight line as the one that minimizes the error between the actual data and the estimate of the line.

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

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

where \bar{x} is the mean value of x_1, x_2, \dots, x_j , and \bar{y} is the mean value of y_1, y_2, \dots, y_j . The coefficients w_0 and w_1 often provide good approximations to otherwise complicated regression equations.

x years experience	y salary (in \$1000s)
3	30
8	57
9	64
13	72
3	36
6	43
11	59
21	90
1	20
16	83



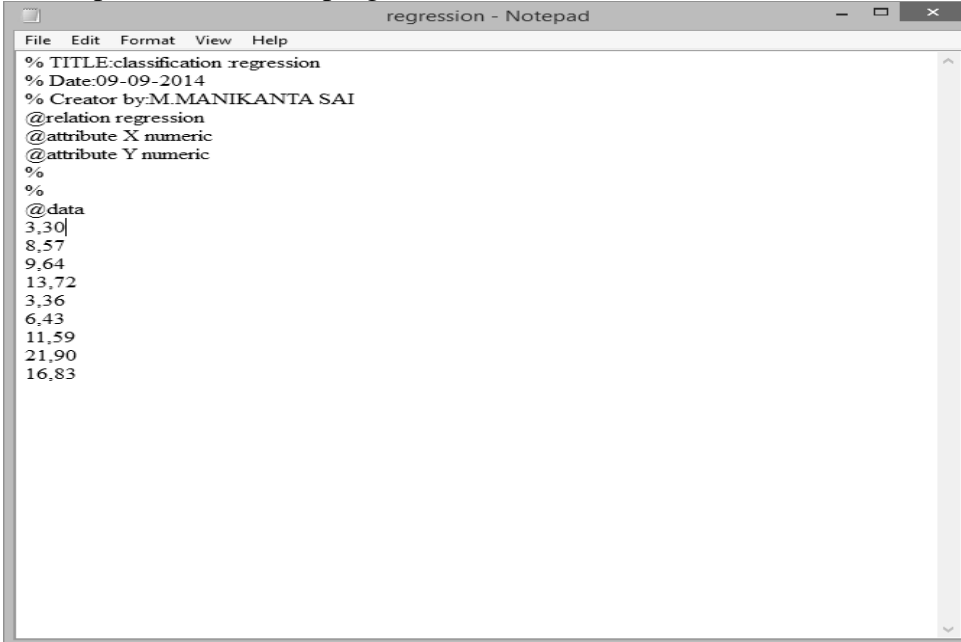
Given the above data, we compute $x = 9:1$ and $y = 55:4$. Substituting these values into Equations (6.50) and (6.51), we get

$$w_1 = 3:5$$

$$w_0 = 23:6$$

Thus, the equation of the least squares line is estimated by $y = 23:6 + 3:5x$. Using this equation, we can predict that the salary of a college graduate with, say, 10 years of experience is \$58,600.

Step1: Open Notepad and write the program and save with **.arff**

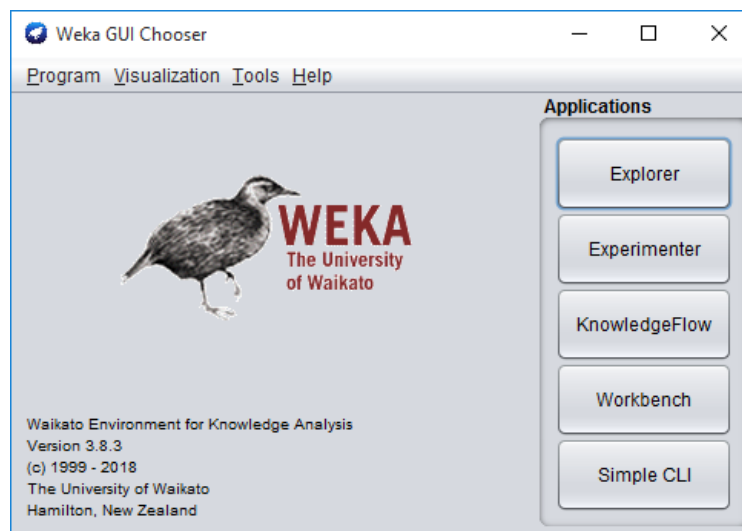


```

regression - Notepad
File Edit Format View Help
% TITLE:classification regression
% Date:09-09-2014
% Creator by:M.MANIKANTA SAI
@relation regression
@attribute X numeric
@attribute Y numeric
%
%
@data
3,30|
8,57
9,64
13,72
3,36
6,43
11,59
21,90
16,83

```

Step2: Go to Start menu and select Weka 3.8.3, on that select Weka 3.8(with console).select Knowledge Flow.



Step 3: In Weka Select Explorer window select **preprocessor** tab and select the **open file** for selection of desired file.



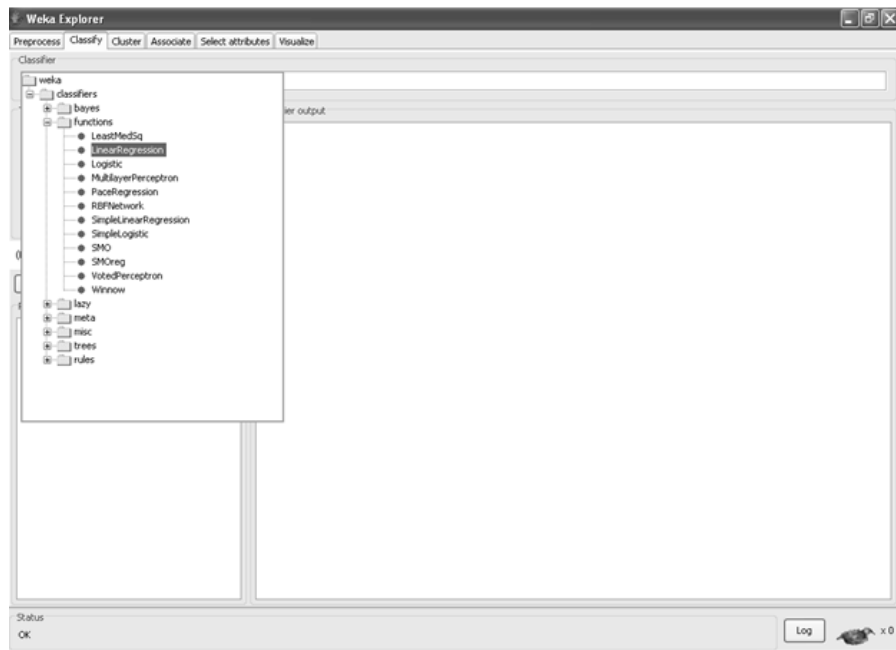
Step 4: In the Preprocessor menu select open file and select Edit

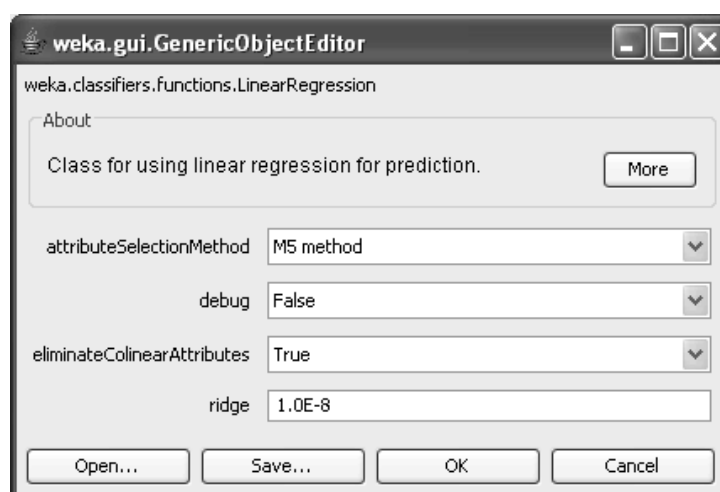
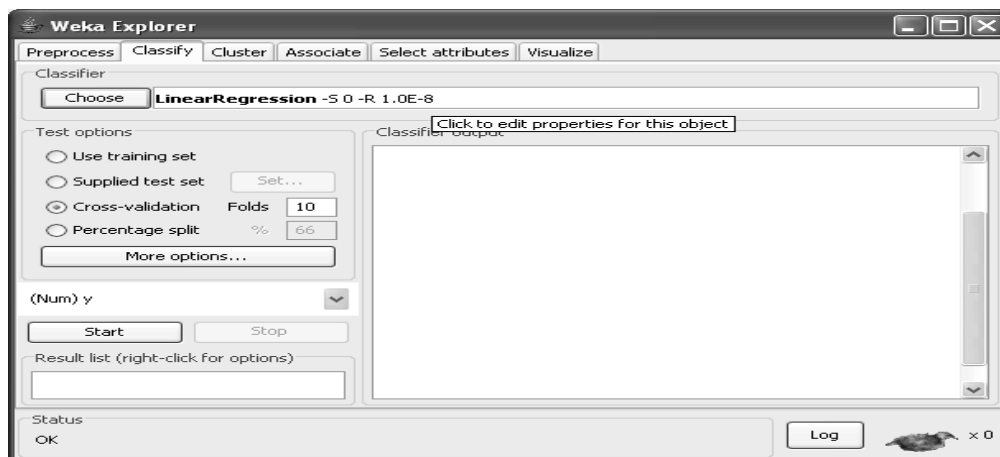
Relation: regression

No.	X Numeric	Y Numeric
1	3.0	30.0
2	8.0	57.0
3	9.0	64.0
4	13.0	72.0
5	3.0	36.0
6	6.0	43.0
7	11.0	59.0
8	21.0	90.0
9	16.0	83.0

Buttons: Undo, OK, Cancel

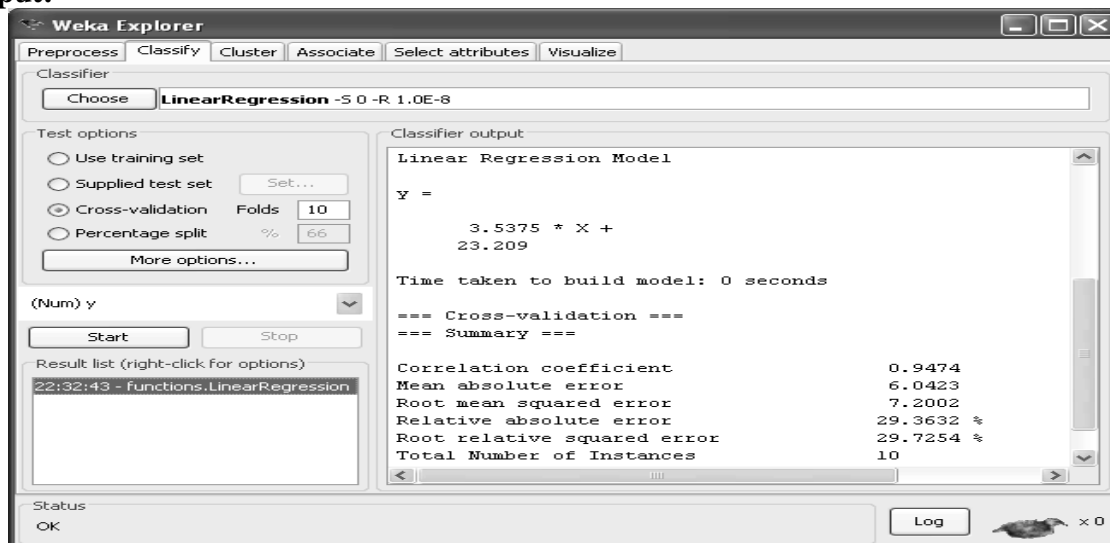
Step 5: Go to classify menu click choose classifiers on that functions and that Linear Regression.






When we click on the Linear Regression option on the above the dialog box appears then click OK, and click on start button.

Output:



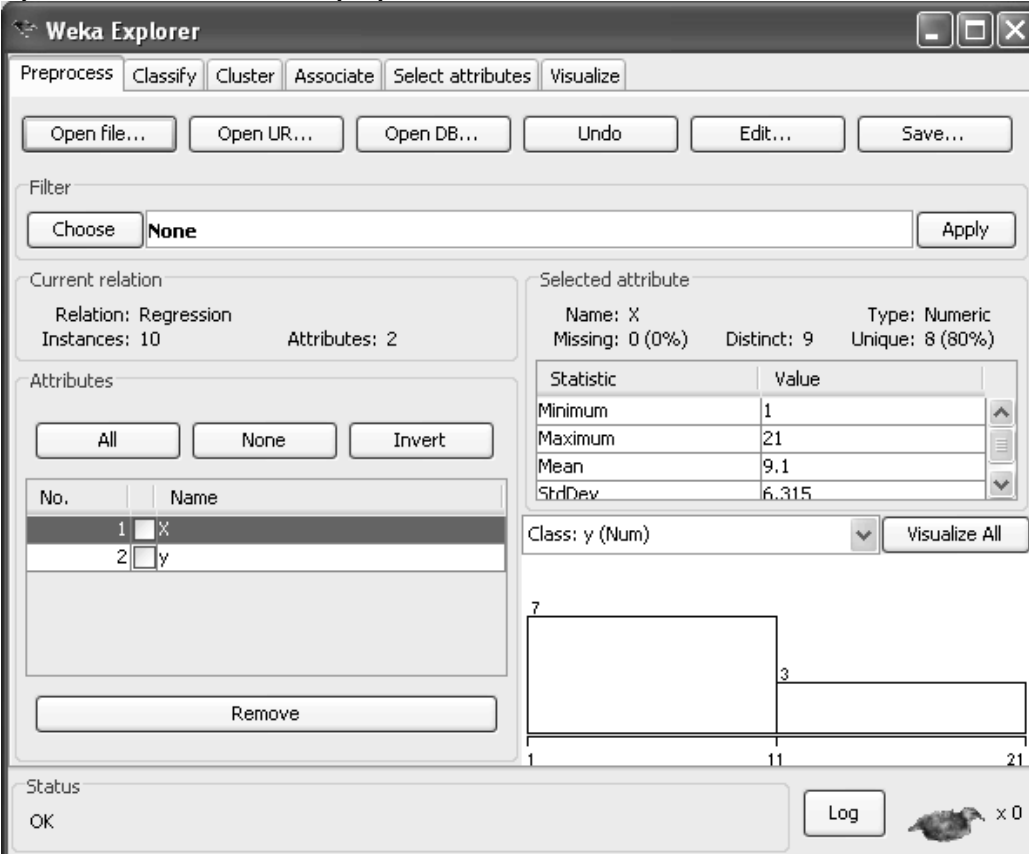
Step1: Open Notepad and write the program and save with **.arff**



```

Regression - Notepad
File Edit Format View Help
%Title:Classification:Regression
%Date:09-09-2014
%Created by:Eswar
@relation Regression
@attribute X numeric
@attribute y numeric
%
%
@data
72,84
50,63
81,77
74,78
94,90
86,75
59,49
83,79
69,77
33,52
88,74
81 90
  
```

Step2: Go to Start menu and select Weka 3.4.11, on that select Weka 3.4(with console) In Weka Select Explorer and on that select preprocessor and select the desired file



The screenshot shows the Weka Explorer interface with the Preprocess tab selected. The 'Current relation' section displays 'Relation: Regression' and 'Instances: 10'. The 'Attributes' section shows a list with 'X' and 'y' selected. The 'Selected attribute' section shows 'Name: X' with statistics: Minimum: 1, Maximum: 21, Mean: 9.1, StdDev: 6.315. A visualization area shows a bar chart with values 7 and 3. The status bar at the bottom shows 'OK' and a 'Log' button.

In the Preprocess open file is present on that select the desired file



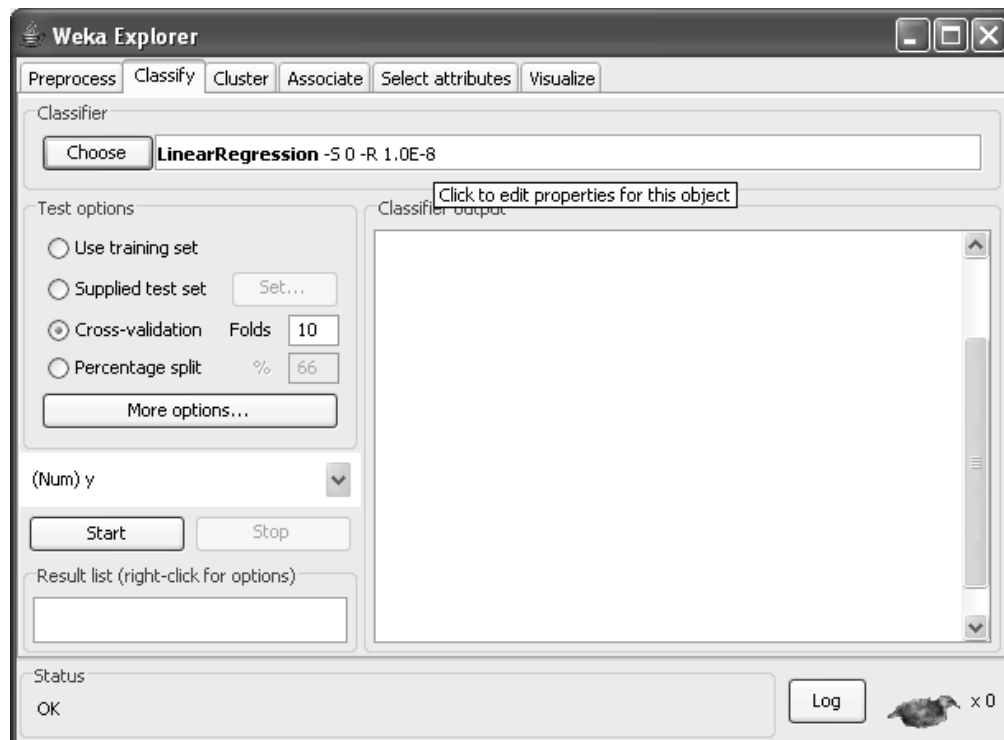
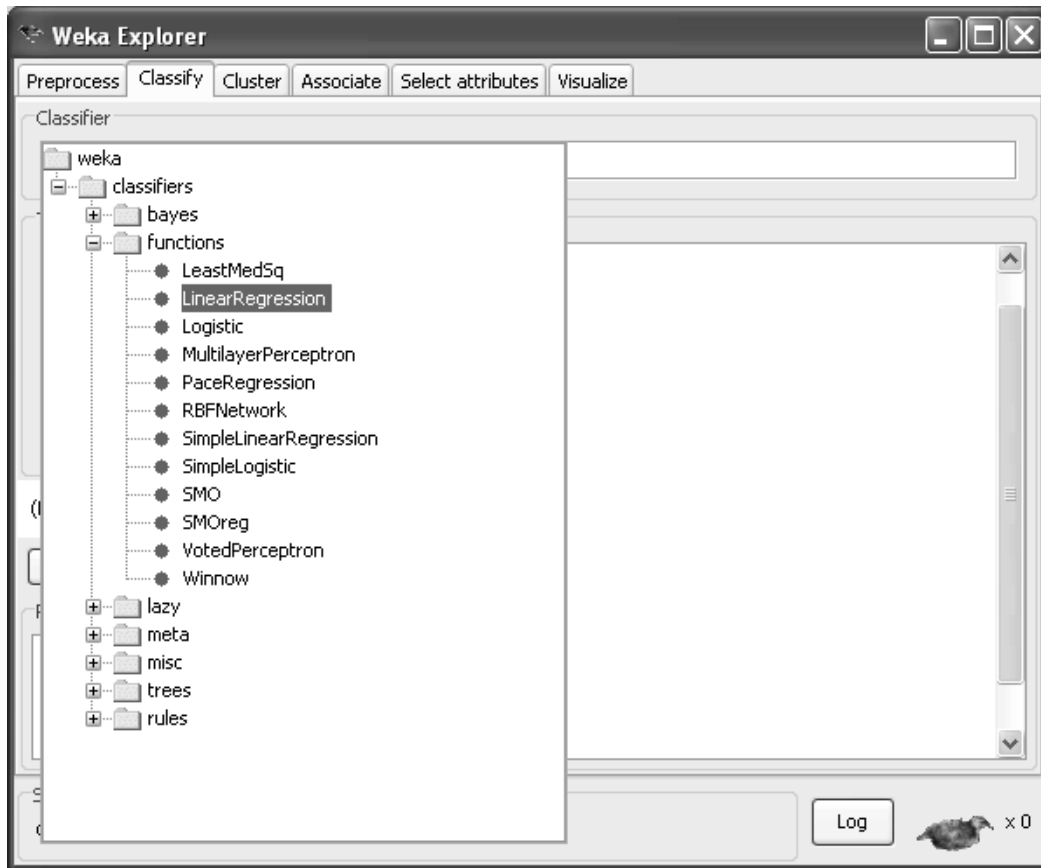
Step3: In the Preprocessor menu select open file and select Edit.

Relation: Regression

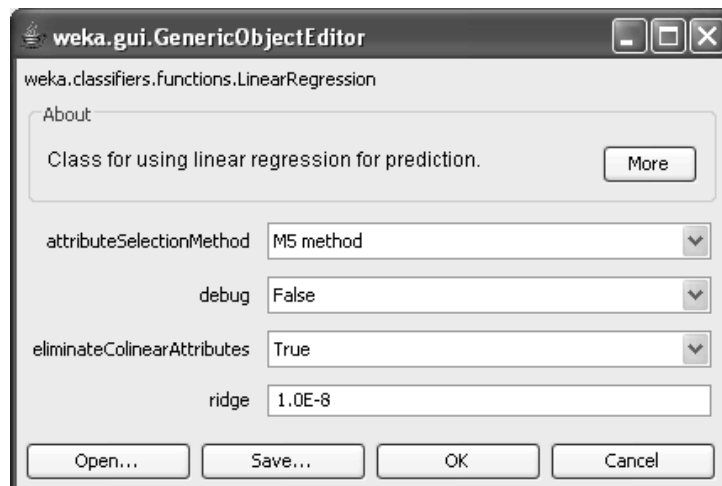
No.	X Numeric	Y Numeric
1	72.0	84.0
2	50.0	63.0
3	81.0	77.0
4	74.0	78.0
5	94.0	90.0
6	86.0	75.0
7	59.0	49.0
8	83.0	79.0
9	69.0	77.0
10	33.0	52.0
11	88.0	74.0
12	81.0	90.0

Undo OK Cancel

Step4: Go to classify menu click choose classifiers on that functions and that Linear Regression

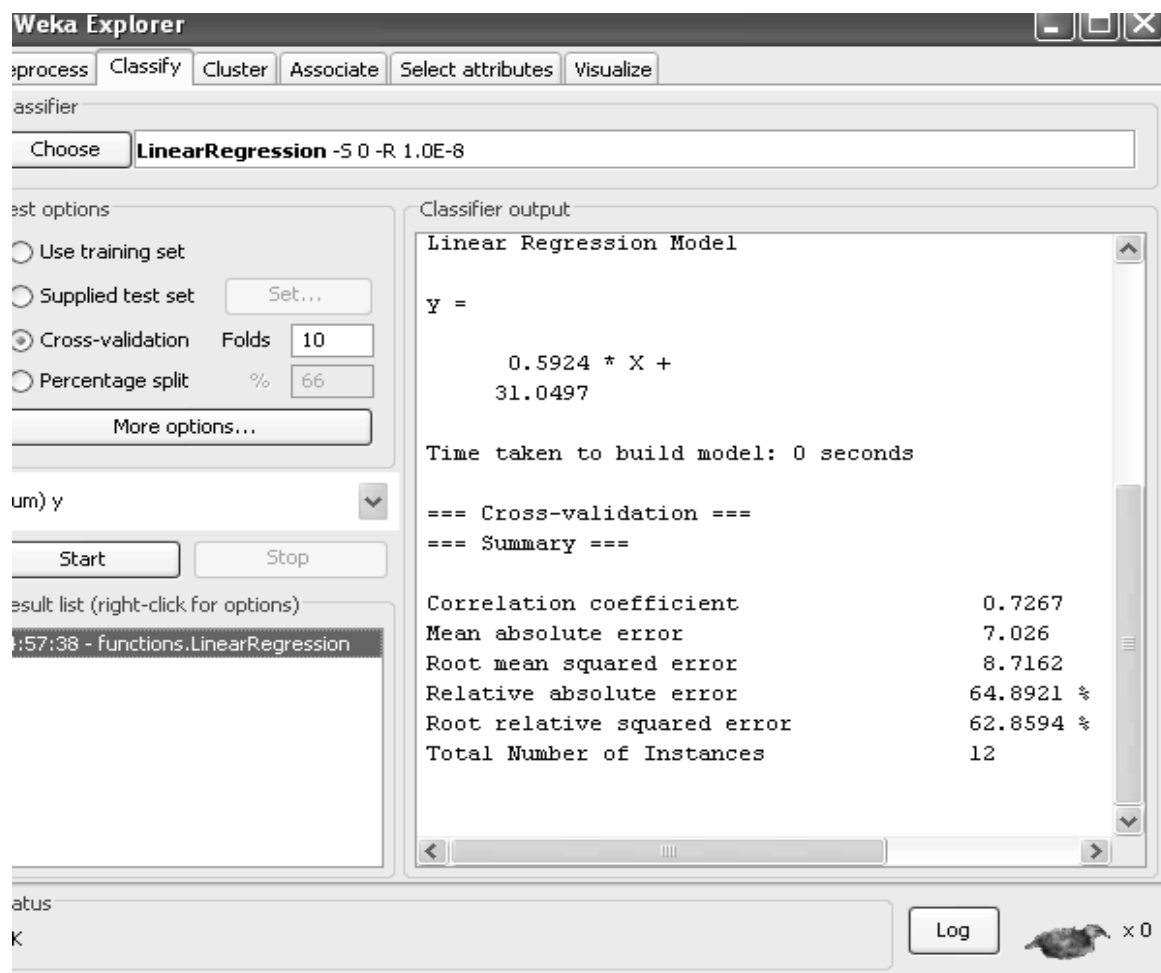


Step 5: When we click on the Linear Regression option on the above the dialog box appears then click O.K



Step 6: And then click OK. button, and then click start in the classify menu present in Weka Explorer window.

Output:



VIVA-QUESTIONS

1. The process of constructing a mathematical model or function that can be used to predict or determine one variable by another variable is called.
2. In the regression equation $Y = 12 - 21X$, the slope is
3. In the regression equation $Y = 55.65 + 0.50X$, the intercept is
4. Which of the following methods do we use to find the best fit line for data in Linear Regression
5. In the mathematical Equation of Linear Regression $Y = \beta_1 + \beta_2X + \epsilon$, (β_1, β_2) refers to
6. A regression between foot length (response variable in cm) and height (explanatory variable in inches) for 33 students resulted in the following regression equation.
 $\hat{y} = 10.9 + 0.23x$
One student in the sample was 100 inches tall with a foot length of 29 cm. What is the predicted foot length for this student?
7. Consider x_1, x_2 to be the independent variables and y the dependent variable, which of the following represents a linear regression model?
8. The sum of the squared residuals is called as:
9. The value of the coefficient of determination ranges between
10. How many coefficients do you need to estimate in a simple linear regression model (One independent variable)?

Experiment No. 6: Demonstrate performing SVM classification on data sets

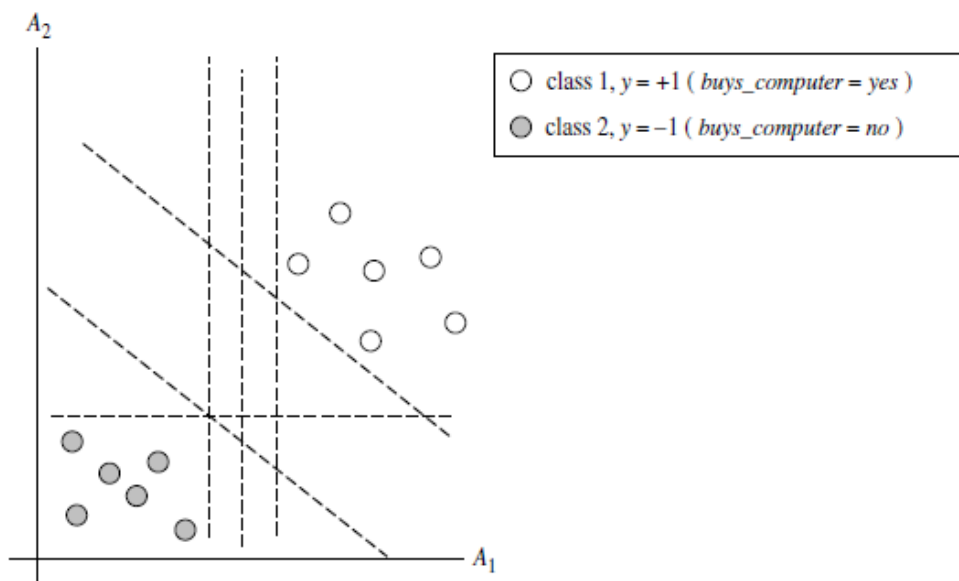
- i. Load each dataset into WEKA and run proximal SVM to find a classifier.
- ii. Load each dataset into WEKA and run linear classifier using linear Programming.

SVM

Support Vector Machines, a promising new method for the classification of both linear and nonlinear data. In a nutshell, a support vector machine (or SVM) is an algorithm that works as follows. It uses a nonlinear mapping to transform the original training data into a higher dimension. Within this new dimension, it searches for the linear optimal separating hyperplane (that is, a “decision boundary” separating the tuples of one class from another). With an appropriate nonlinear mapping to a sufficiently high dimension, data from two classes can always be separated by a hyperplane.

The Case When the Data Are Linearly Separable

Let the data set D be given as $(X_1, y_1), (X_2, y_2), \dots, (X_j, y_j), \dots, (X_n, y_n)$, where X_i is the set of training tuples with associated class labels, y_i . Each y_i can take one of two values, either $+1$ or -1 (i.e., $y_i \in \{+1, -1\}$), corresponding to the classes *buys computer = yes* and *buys computer = no*, respectively.

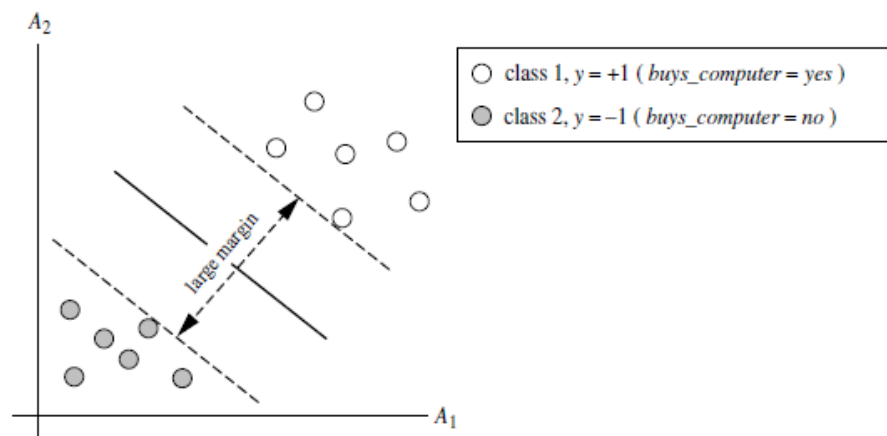
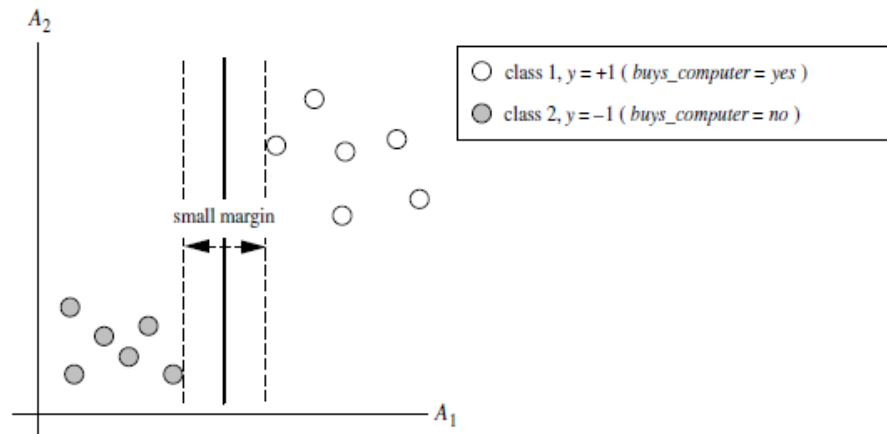


Before we get into the definition of margins, let's take an intuitive look at this figure. Both hyperplanes can correctly classify all of the given data tuples. Intuitively, however, we expect the hyperplane with the larger margin to be more accurate at classifying future data tuples than

the hyperplane with the smaller margin. This is why (during the learning or training phase), the SVM searches for the hyperplane with the largest margin, that is, the *maximum marginal hyperplane* (MMH). The associated margin gives the largest separation between classes. Getting to an informal definition of margin, we can say that the shortest distance from a hyperplane to one side of its margin is equal to the shortest distance from the hyperplane to the other side of its margin, where the “sides” of the margin are parallel to the hyperplane. When dealing with the MMH, this distance is, in fact, the shortest distance from the MMH to the closest training tuple of either class.

$$W \cdot X + b = 0,$$

$$w_0 + w_1x_1 + w_2x_2 = 0.$$

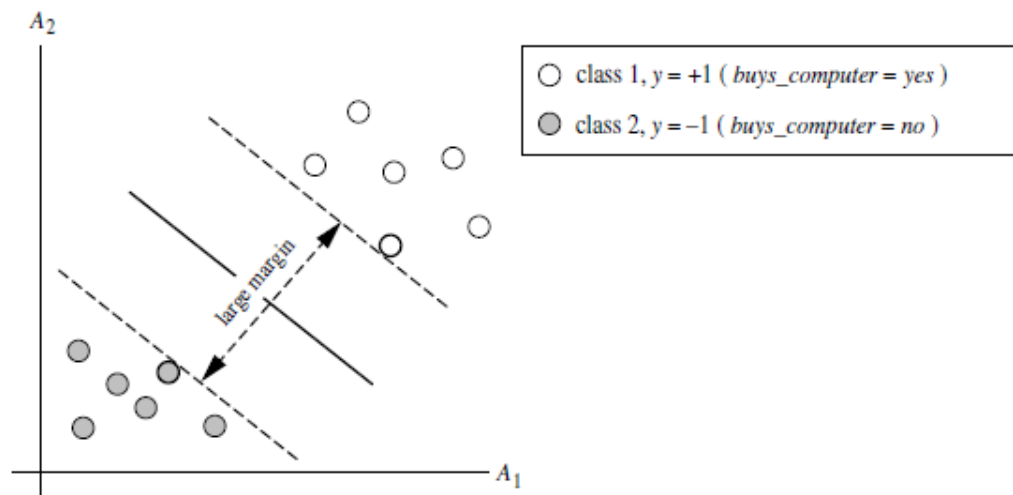


Thus, any point that lies above the separating hyperplane satisfies

$$w_0 + w_1x_1 + w_2x_2 > 0.$$

Similarly, any point that lies below the separating hyperplane satisfies

$$w_0 + w_1x_1 + w_2x_2 < 0.$$



The weights can be adjusted so that the hyperplanes defining the “sides” of the margin can be written as

$$H_1 : w_0 + w_1x_1 + w_2x_2 \geq 1 \text{ for } y_i = +1, \text{ and}$$

$$H_2 : w_0 + w_1x_1 + w_2x_2 \leq -1 \text{ for } y_i = -1.$$

$$^9\text{If } \mathbf{W} = \{w_1, w_2, \dots, w_n\} \text{ then } \sqrt{\mathbf{W} \cdot \mathbf{W}} = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}.$$

Algorithm

- Define an optimal hyper plane: maximize margin
- Extend the above definition for non-linearly separable problems: have a penalty term for misclassifications.
- Map data to high dimensional space where it is easier to classify with linear decision surfaces: reformulate problem so that data is mapped implicitly to this space.

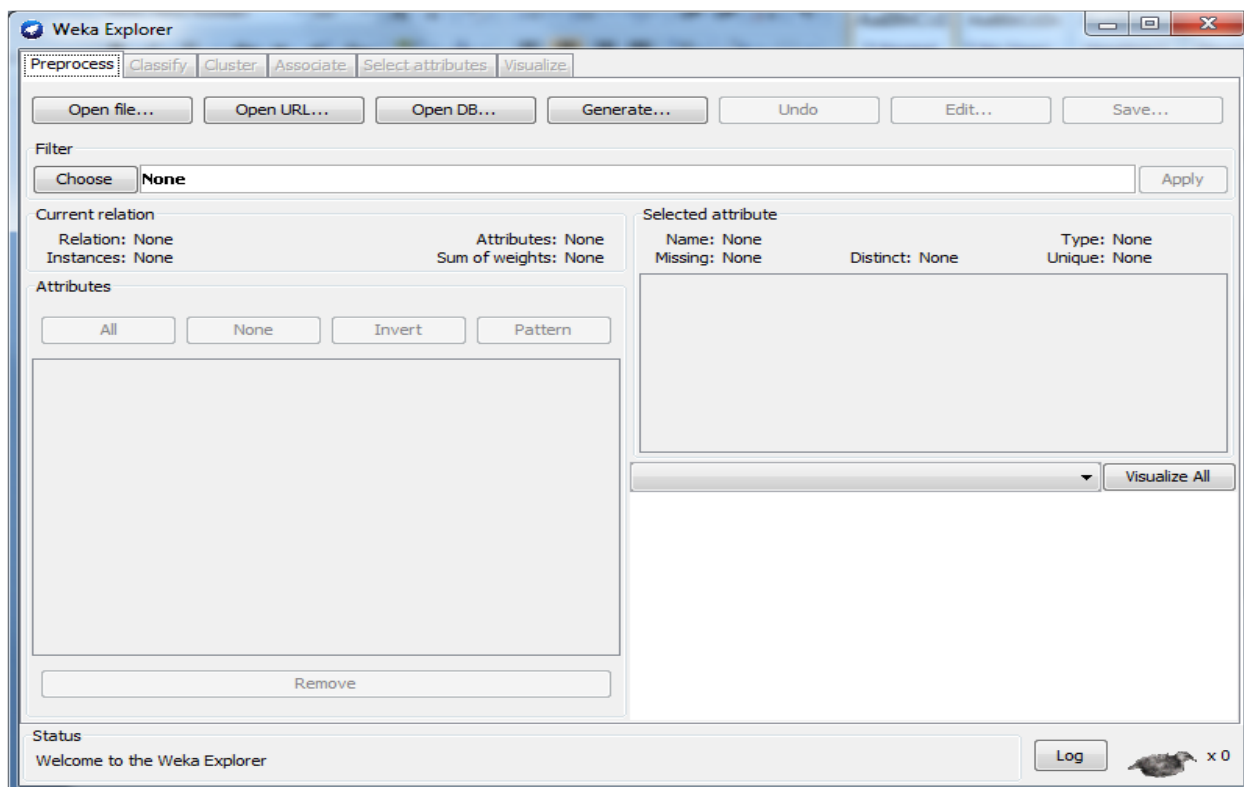
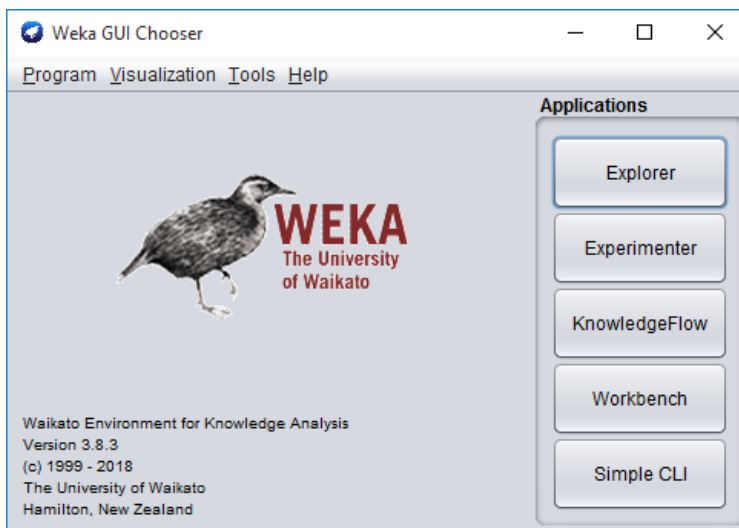
PROCEDURE:

1. We begin the experiment by loading the data (vote.arff) into weka.
2. Next we select the classify tab and click choose function button to select the Support vector machine (SMO).
3. Now we specify the various parameters. These can be specified by clicking in the text box to the right of the chose button.
4. 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.

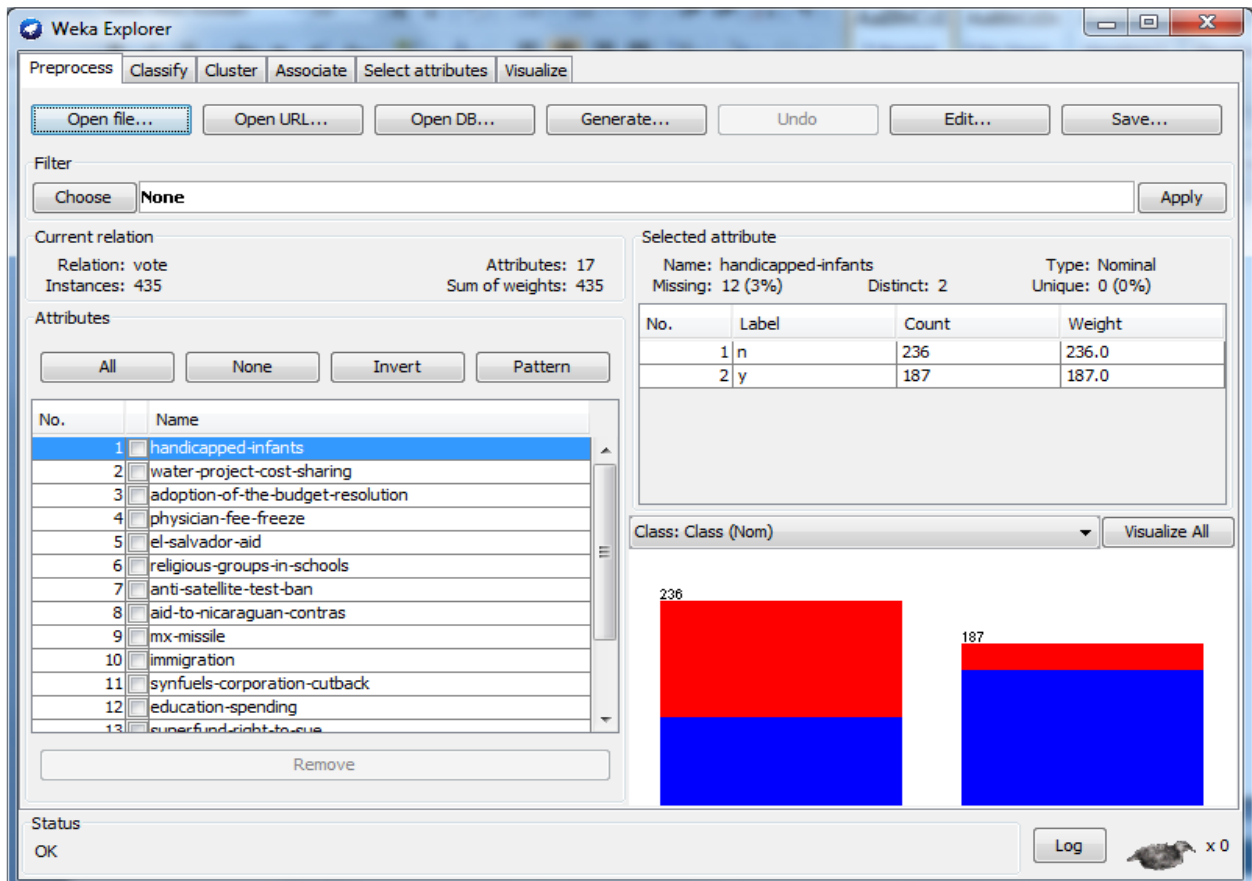
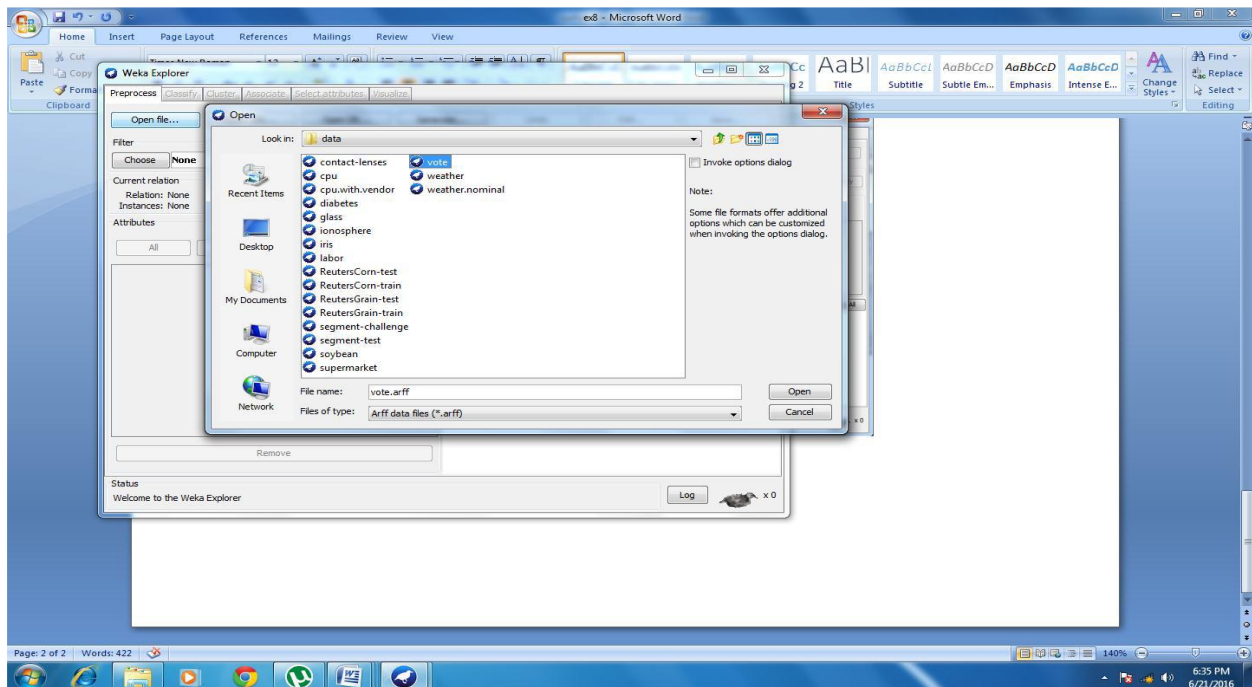
5. We now click "start" to generate the model .the ASCII version of the tree as well as evaluation statistic will appear in the right panel when the model construction is complete.
6. Note that the classification accuracy of model is about 69%.this indicates that we may find more work. (Either in preprocessing or in selecting current parameters for the classification)
7. The run information of the support vector classifier will be displayed with the correctly and incorrectly classified instances.

STEPS:

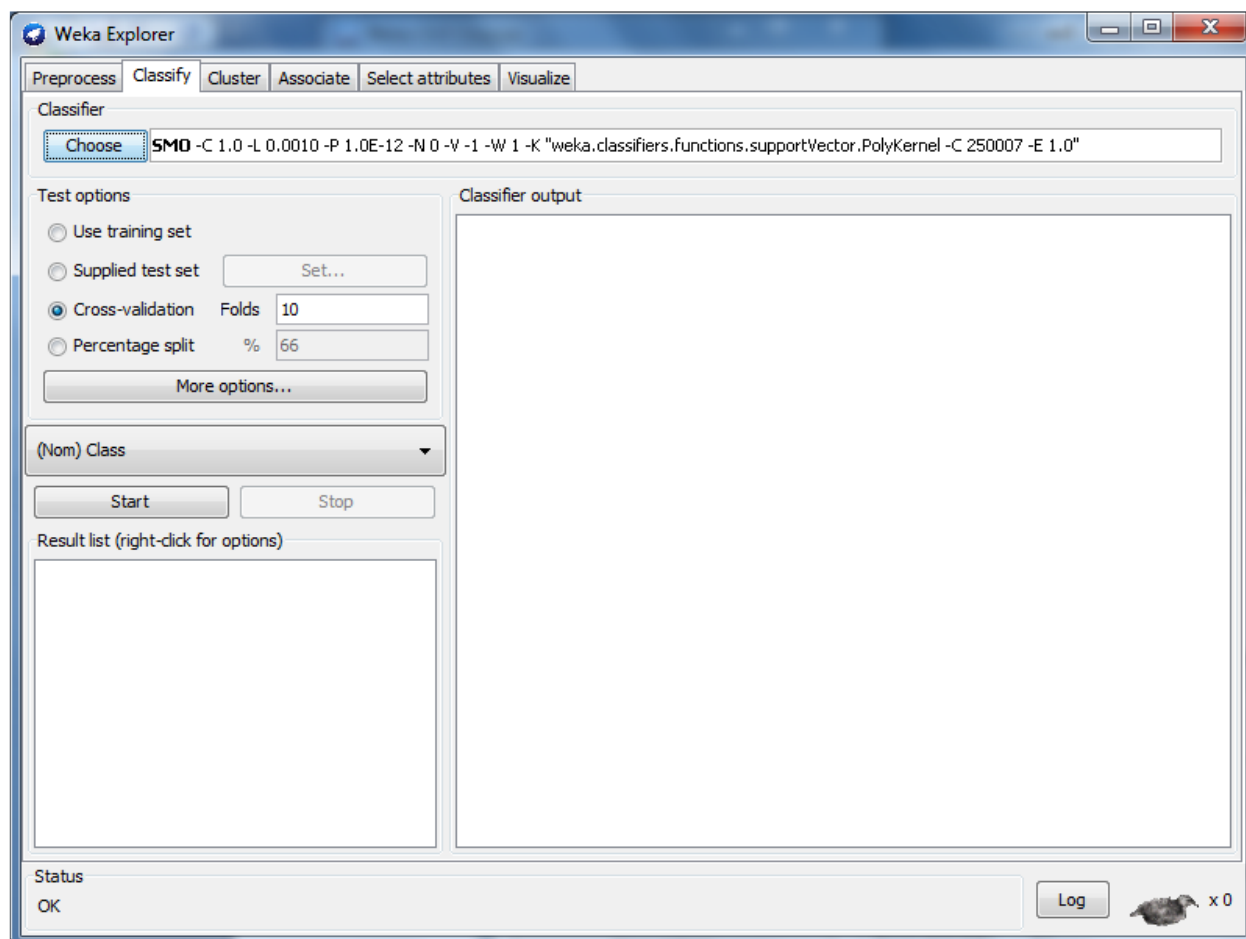
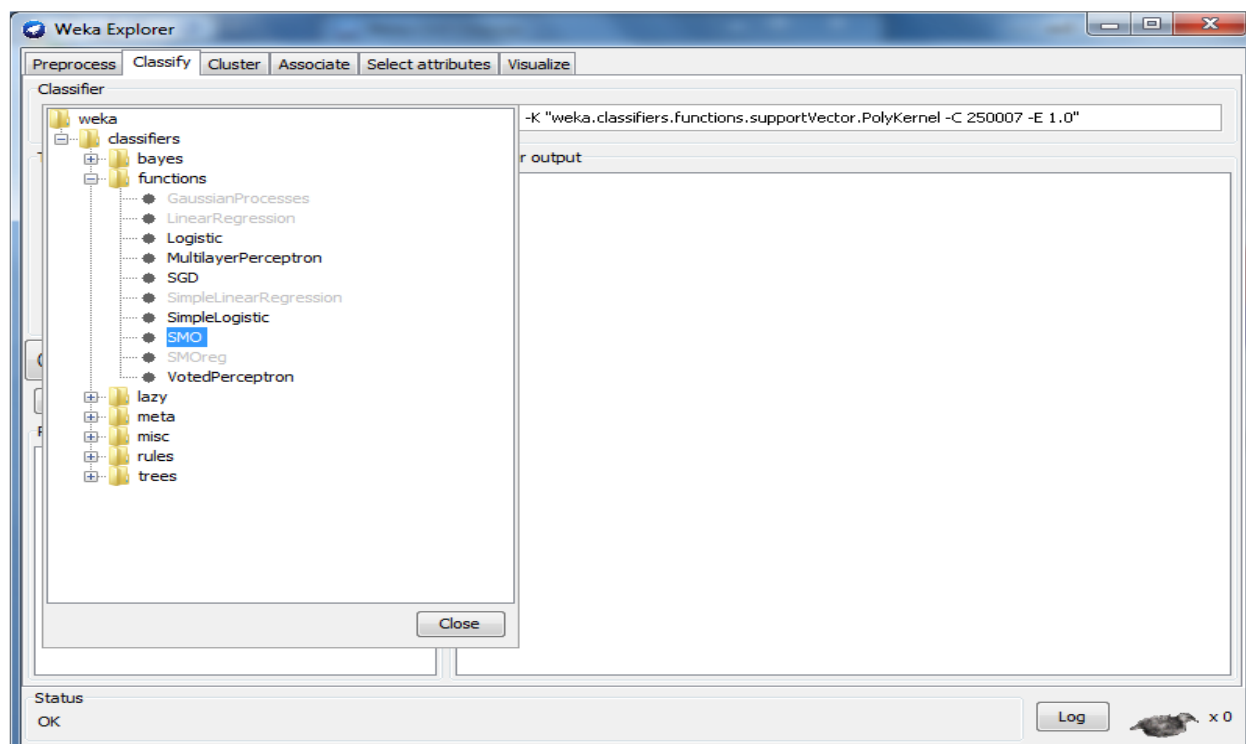
Open **Weka** tool – click **Explorer** in **Weka GUI Chooser**.



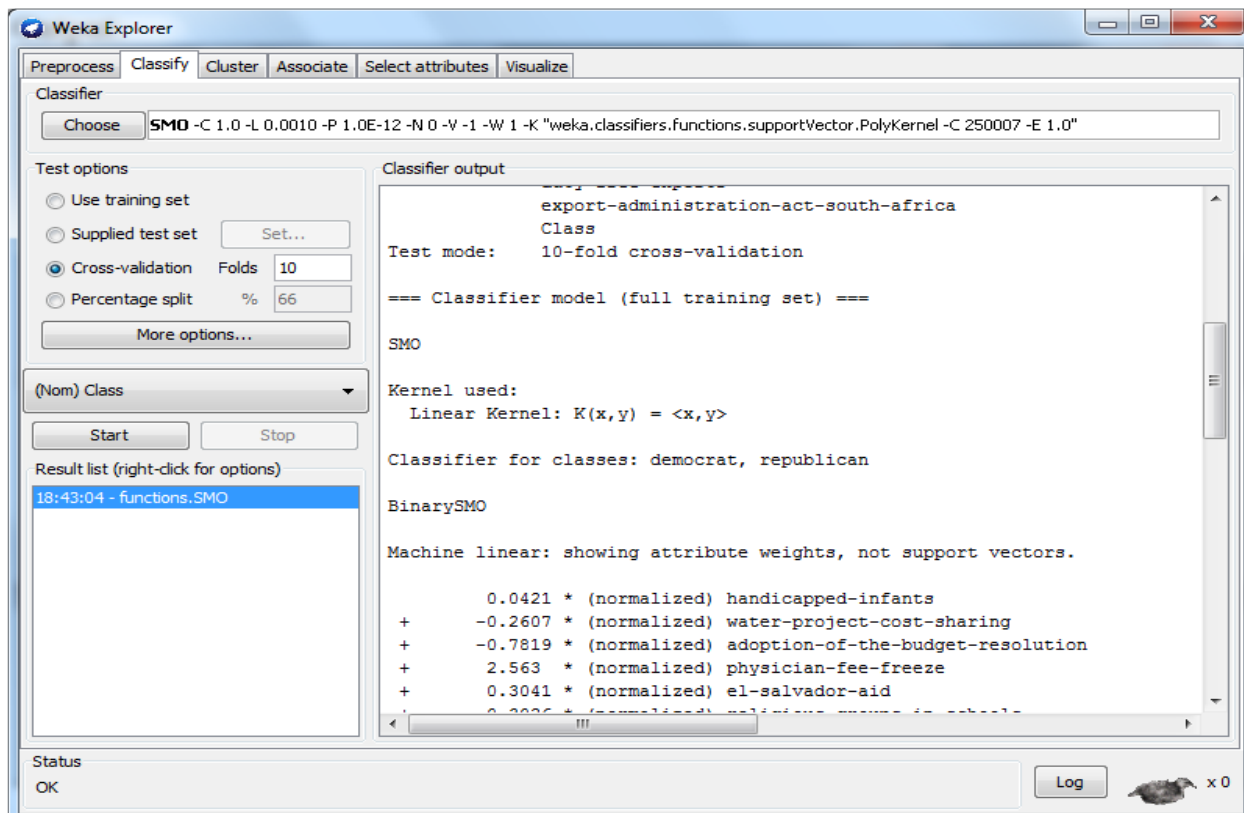
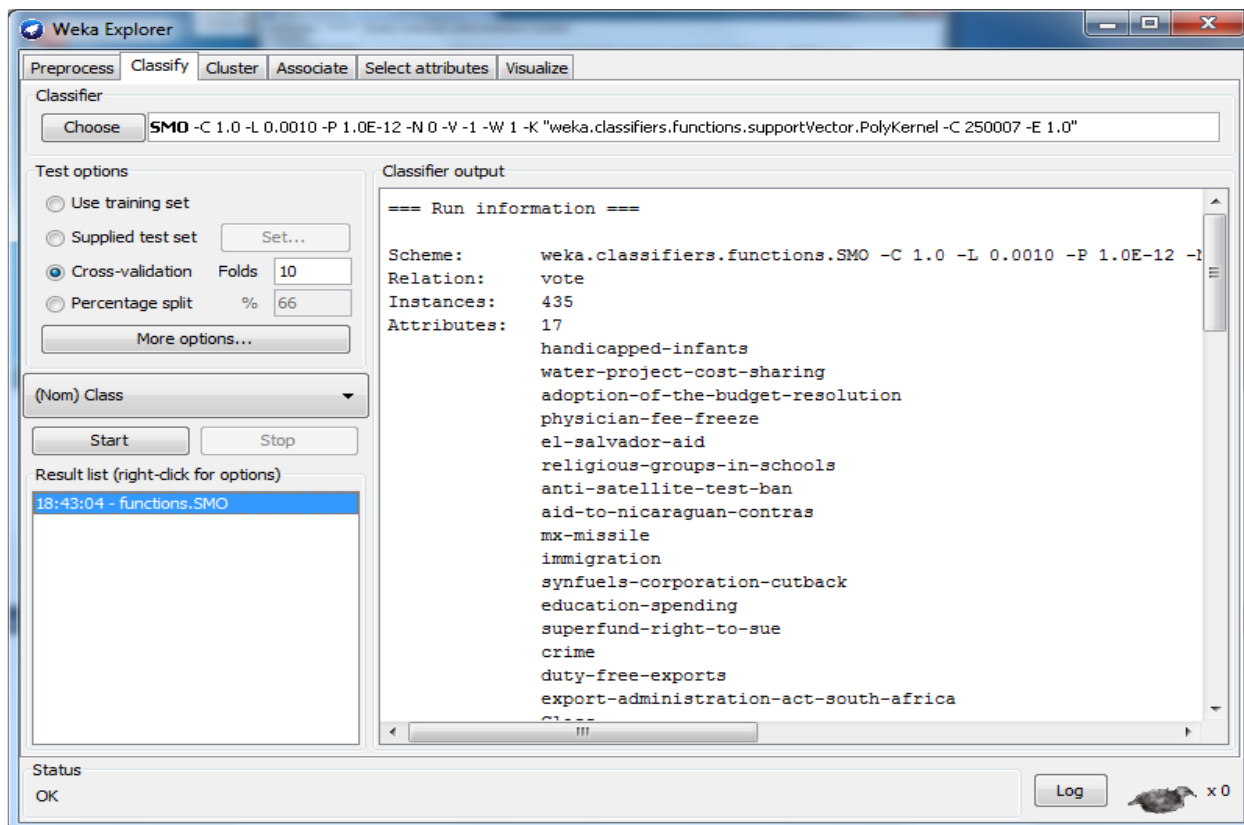
In **preprocess** tab click **Open file..** – choose **file name** (vote.arff).



Goto **Classify** tab - click **Choose** button – select **SMO** option.



Click **Start** button.



Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose SMO -C 1.0 -L 0.0010 -P 1.0E-12 -N 0 -V -1 -W 1 -K "weka.classifiers.functions.supportVector.PolyKernel -C 250007 -E 1.0"

Test options

Use training set

Supplied test set Set...

Cross-validation Folds 10

Percentage split % 66

More options...

(Nom) Class

Start Stop

Result list (right-click for options)

18:43:04 - functions.SMO

Classifier output

Time taken to build model: 0.05 seconds

=== Stratified cross-validation ===

=== Summary ===

Correctly Classified Instances	418	96.092 %
Incorrectly Classified Instances	17	3.908 %
Kappa statistic	0.9178	
Mean absolute error	0.0391	
Root mean squared error	0.1977	
Relative absolute error	8.2405 %	
Root relative squared error	40.6018 %	
Coverage of cases (0.95 level)	96.092 %	
Mean rel. region size (0.95 level)	50 %	
Total Number of Instances	435	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area
	0.963	0.042	0.973	0.963	0.968	0.96
	0.958	0.037	0.942	0.958	0.95	0.96
Weighted Avg.	0.961	0.04	0.961	0.961	0.961	0.96

Status OK

Log x 0

VIVA-QUESTIONS

1. Support vector machine may be termed as:
2. What do you mean by generalization error in terms of the SVM?
3. Margin of a hyperplane is defined as:
4. In a hard margin support vector machine
5. The effectiveness of an SVM depends upon:
6. If the hyper plane $WTX+b=0$ coorectly classifies all the training points (X_i, Y_i) , where $Y_i = \{+1, -1\}$, then
7. Support vectors are the data points that lie closest to the decision surface
8. Consider a two class problem, whose training points are distributed in the figure below. One possible separating hyperplane.
9. For a two-class classification problem, we use an SVM classifier and obtain the following separating hyperplane. We have marked 4 instances of the training data. Identify the point which will have the most impact on the shape of the boundary on it's removal.

Experiment No. 7: Demonstrate performing clustering on data sets

- i. Load each dataset into WEKA and run simple k-means clustering algorithm with different values of k (number of desired clusters).
- ii. Explore visualization features of WEKA to visualize the clusters. Derive interesting insights and explain.

CLUSTERING

Cluster analysis or clustering is the assignment of a set of observations into subsets (called clusters) so that observations in the same cluster are similar in some sense. Clustering is a method of unsupervised learning, and a common technique for statistical data analysis used in many fields, including machine learning, data mining, pattern recognition, image analysis and bioinformatics.

Types of clustering

Data clustering algorithms can be hierarchical. Hierarchical algorithms find successive clusters using previously established clusters. These algorithms can be either agglomerative ("bottom-up") or divisive ("top-down"). Agglomerative algorithms begin with each element as a separate cluster and merge them into successively larger clusters. Divisive algorithms begin with the whole set and proceed to divide it into successively smaller clusters.

Partitional algorithms typically determine all clusters at once, but can also be used as divisive algorithms in the hierarchical clustering.

Density-based clustering algorithms are devised to discover arbitrary-shaped clusters. In this approach, a cluster is regarded as a region in which the density of data objects exceeds a threshold. DBSCAN and OPTICS are two typical algorithms of this kind.

Two-way clustering, co-clustering or bi-clustering are clustering methods where not only the objects are clustered but also the features of the objects, i.e., if the data is represented in a data matrix, the rows and columns are clustered simultaneously.

Many clustering algorithms require specification of the number of clusters to produce in the input data set, prior to execution of the algorithm. Barring knowledge of the proper value beforehand, the appropriate value must be determined, a problem for which a number of techniques have been developed.

***k*-means clustering:**

In statistics and machine learning, k -means clustering is a method of cluster analysis which aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean.

ALGORITHM:

Regarding computational complexity, the k -means clustering problem is:

- NP-hard in general Euclidean space d even for 2 clusters
- NP-hard for a general number of clusters k even in the plane

- If k and d are fixed, the problem can be exactly solved in time $O(ndk+1 \log n)$, where n is the number of entities to be cluster

Example

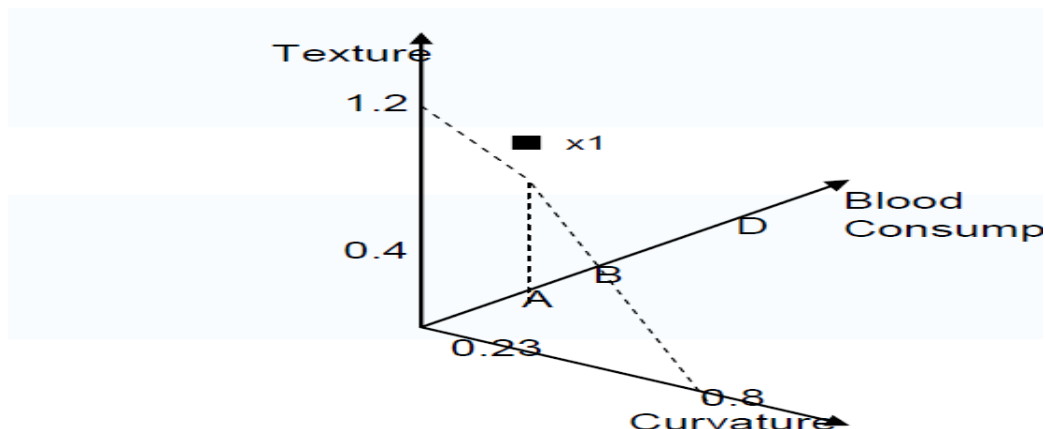
for example, if our class (decision) attribute is tumor Type and its values are: malignant, benign, etc. - these will be the classes. They will be represented by cluster1, cluster2, etc. However, the class information is never provided to the algorithm. The class information can be used later on, to evaluate how accurately the algorithm classified the objects.

Curvature	Texture	Blood Consump	Tumor Type
0.8	1.2	A	Benign
0.75	1.4	B	Benign
0.23	0.4	D	Malignant
0.23	0.5	D	Malignant

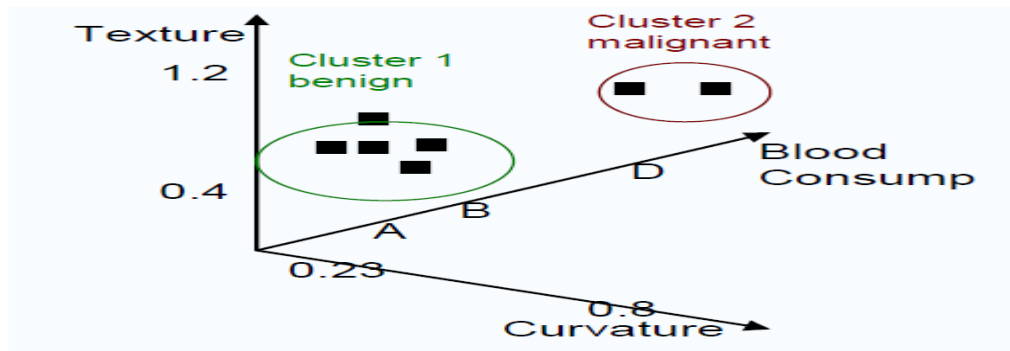
→

	Curvature	Texture	Blood Consump	Tumor Type
x1	0.8	1.2	A	Benign
x2	0.75	1.4	B	Benign
x3	0.23	0.4	D	Malignant
x4	0.23	0.5	D	Malignant
.				
.				

The way we do that, is by plotting the objects from the database into space. Each attribute is one dimension



After all the objects are plotted, we will calculate the distance between them, and the ones that are close to each other – we will group them together, i.e. place them in the same cluster.



Problem: Cluster the following eight points (with (x, y) representing locations) into three clusters A1(2, 10) A2(2, 5) A3(8, 4) A4(5, 8) A5(7, 5) A6(6, 4) A7(1, 2) A8(4, 9). Initial cluster centers are: A1(2, 10), A4(5, 8) and A7(1, 2).

The distance function between two points $a=(x1, y1)$ and $b=(x2, y2)$ is defined as:

$$\rho(a, b) = |x2 - x1| + |y2 - y1|.$$

Use k-means algorithm to find the three cluster centers after the second iteration.

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean 1	Dist Mean 2	Dist Mean 3	Cluster
A1	(2, 10)				
A2	(2, 5)				
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

First we list all points in the first column of the table above. The initial cluster centers – means, are (2, 10), (5, 8) and (1, 2) - chosen randomly. Next, we will calculate the distance from the first point (2, 10) to each of the three means, by using the distance function:

point	mean1
$x1, y1$	$x2, y2$
(2, 10)	(2, 10)

$$\begin{aligned} \rho(a, b) &= |x2 - x1| + |y2 - y1| \\ \rho(\text{point}, \text{mean1}) &= |x2 - x1| + |y2 - y1| \\ &= |2 - 2| + |10 - 10| \\ &= 0 + 0 \\ &= 0 \end{aligned}$$

point	mean2
$x1, y1$	$x2, y2$
(2, 10)	(5, 8)

$$\begin{aligned} \rho(a, b) &= |x2 - x1| + |y2 - y1| \\ \rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |5 - 2| + |8 - 10| \\ &= 3 + 2 \\ &= 5 \end{aligned}$$

point	mean3
$x1, y1$	$x2, y2$
(2, 10)	(1, 2)

$$\begin{aligned} \rho(a, b) &= |x2 - x1| + |y2 - y1| \\ \rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |1 - 2| + |2 - 10| \\ &= 1 + 8 \\ &= 9 \end{aligned}$$

So, we fill in these values in the table:

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean	Dist Mean	Dist Mean	Cluster
		1	2	3	
A1	(2, 10)	0	5	9	1
A2	(2, 5)				
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

So, which cluster should the point (2, 10) be placed in? The one, where the point has the shortest distance to the mean – that is mean 1 (cluster 1), since the distance is 0.

Cluster 1 Cluster 2 Cluster 3
(2, 10)

So, we go to the second point (2, 5) and we will calculate the distance to each of the three means, by using the distance function:

point mean1
 $x1, y1$ $x2, y2$
(2, 5) (2, 10)

$$\begin{aligned} \rho(a, b) &= |x2 - x1| + |y2 - y1| \\ \rho(\text{point}, \text{mean1}) &= |x2 - x1| + |y2 - y1| \\ &= |2 - 2| + |10 - 5| \\ &= 0 + 5 \\ &= 5 \end{aligned}$$

point mean2
 $x1, y1$ $x2, y2$
(2, 5) (5, 8)

$$\begin{aligned} \rho(a, b) &= |x2 - x1| + |y2 - y1| \\ \rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |5 - 2| + |8 - 5| \\ &= 3 + 3 \\ &= 6 \end{aligned}$$

point mean3
 $x1, y1$ $x2, y2$
 (2, 5) (1, 2)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\rho(\text{point}, \text{mean2}) = |x2 - x1| + |y2 - y1|$$

$$= |1 - 2| + |2 - 5|$$

$$= 1 + 3$$

$$= 4$$

So, we fill in these values in the table:

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean	Dist Mean	Dist Mean	Cluster
		1	2	3	
A1	(2, 10)	0	5	9	1
A2	(2, 5)	5	6	4	3
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

So, which cluster should the point (2, 5) be placed in? The one, where the point has the shortest distance to the mean – that is mean 3 (cluster 3), since the distance is 0.

Cluster 1 Cluster 2 Cluster 3
 (2, 10) (2, 5)

Analogically, we fill in the rest of the table, and place each point in one of the clusters:

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean	Dist Mean	Dist Mean	Cluster
		1	2	3	
A1	(2, 10)	0	5	9	1
A2	(2, 5)	5	6	4	3
A3	(8, 4)	12	7	9	2
A4	(5, 8)	5	0	10	2
A5	(7, 5)	10	5	9	2
A6	(6, 4)	10	5	7	2
A7	(1, 2)	9	10	0	3
A8	(4, 9)	3	2	10	2

Cluster 1	Cluster 2	Cluster 3
(2, 10)	(8, 4)	(2, 5)
	(5, 8)	(1, 2)
	(7, 5)	
	(6, 4)	
	(4, 9)	

Next, we need to re-compute the new cluster centers (means). We do so, by taking the mean of all points in each cluster.

For Cluster 1, we only have one point A1(2, 10), which was the old mean, so the cluster center remains the same.

For Cluster 2, we have $((8+5+7+6+4)/5, (4+8+5+4+9)/5) = (6, 6)$

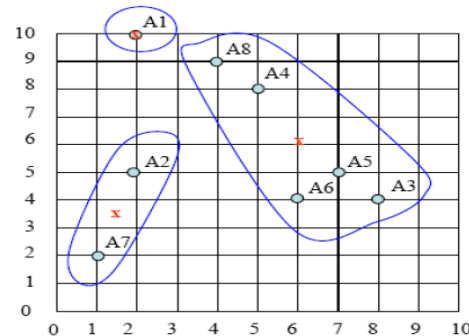
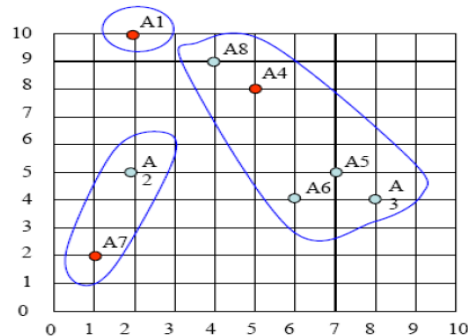
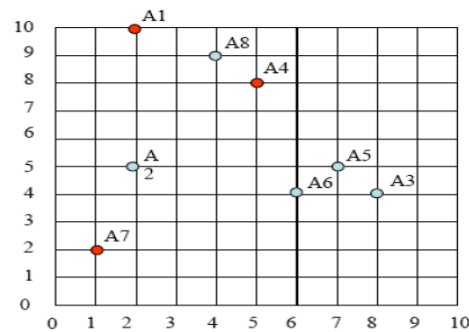
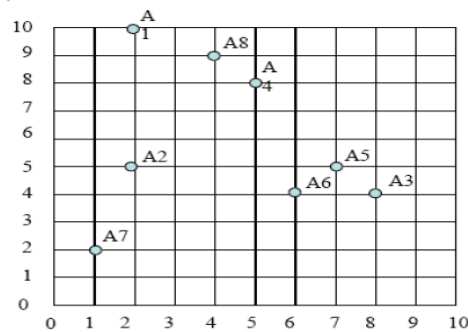
For Cluster 3, we have $((2+1)/2, (5+2)/2) = (1.5, 3.5)$

new clusters: 1: {A1}, 2: {A3, A4, A5, A6, A8}, 3: {A2, A7}

b) centers of the new clusters:

C1= (2, 10), C2= $((8+5+7+6+4)/5, (4+8+5+4+9)/5) = (6, 6)$, C3= $((2+1)/2, (5+2)/2) = (1.5, 3.5)$

c)



The initial cluster centers are shown in red dot. The new cluster centers are shown in red x.

That was Iteration1 (epoch1). Next, we go to Iteration2 (epoch2), Iteration3, and so on until the means do not change anymore.

In Iteration2, we basically repeat the process from Iteration1 this time using the new means we computed.

d)

We would need two more epochs. After the 2nd epoch the results would be:

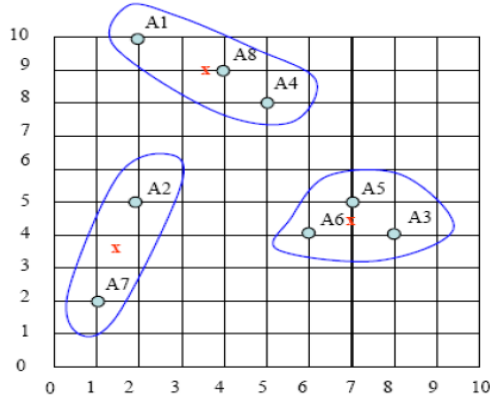
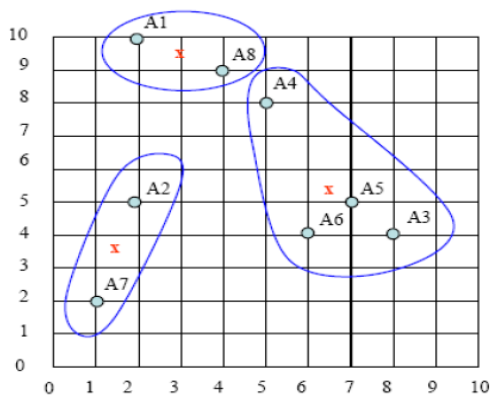
1: {A1, A8}, 2: {A3, A4, A5, A6}, 3: {A2, A7}

with centers $C1=(3, 9.5)$, $C2=(6.5, 5.25)$ and $C3=(1.5, 3.5)$.

After the 3rd epoch, the results would be:

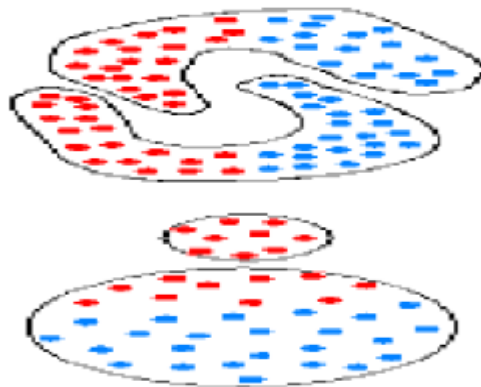
1: {A1, A4, A8}, 2: {A3, A5, A6}, 3: {A2, A7}

with centers $C1=(3.66, 9)$, $C2=(7, 4.33)$ and $C3=(1.5, 3.5)$.



K-means: Problems and Limitations

Based on minimizing within cluster error -a criterion that is not appropriate for many situations.-
Unsuitable when clusters have widely different sizes or have convex shapes.



K-mediod method:

The k -medoids algorithm is a clustering algorithm related to the k -means algorithm and the medoidshift algorithm. Both the k -means and k -medoids algorithms are partitional (breaking the dataset up into groups) and both attempt to minimize squared error, the distance between points labeled to be in a cluster and a point designated as the center of that cluster. In contrast to the k -means algorithm k medoids chooses datapoints as centers (medoids or exemplars).

ALGORITHM:

The most common realization of k -medoid clustering is the Partitioning Around Medoids (PAM) algorithm and is as follows:

1. Initialize: randomly select k of the n data points as the mediods

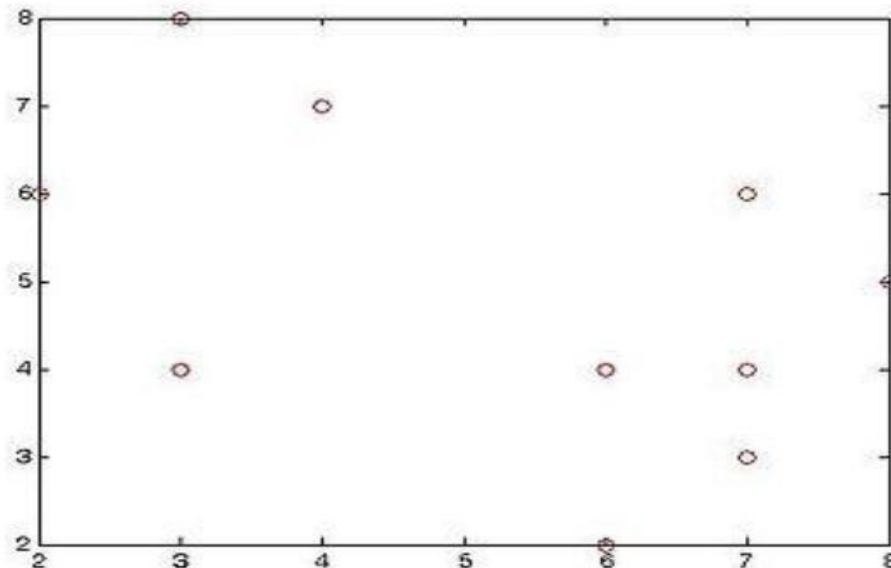
2. Associate each data point to the closest medoid. ("closest" here is defined using any valid distance metric, most commonly Euclidean distance, Manhattan distance or Minkowski distance)
3. For each medoid m
 1. For each non-medoid data point o
 1. Swap m and o and compute the total cost of the configuration
4. Select the configuration with the lowest cost.
5. Repeat steps 2 to 5 until there is no change in the medoid.

EXAMPLE:

Cluster the following data set of ten objects into two clusters i.e $k = 2$.

Consider a data set of ten objects as follows:

X1	2	6
X2	3	4
X3	3	8
X4	4	7
X5	6	2
X6	6	4
X7	7	3
X8	7	4
X9	8	5
X10	7	6

**Step 1**

Initialise k centre

Let us assume $c_1 = (3,4)$ and $c_2 = (7,4)$

So here c_1 and c_2 are selected as medoid.

Calculating distance so as to associate each data object to its nearest medoid. Cost is calculated using Minkowski distance metric with $r = 1$.

C1		Data Objects (X _i)		Cost (distance)
3	4	2	6	3
3	4	3	8	4
3	4	4	7	4
3	4	6	2	5
3	4	6	4	3
3	4	7	3	5
3	4	8	5	6
3	4	7	6	6

C2		Data Objects (X _i)		Cost (distance)
7	4	2	6	7
7	4	3	8	8
7	4	4	7	6
7	4	6	2	3
7	4	6	4	1
7	4	7	3	1
7	4	8	5	2
7	4	7	6	2

Then so the clusters become:

Cluster1 = {(3,4)(2,6)(3,8)(4,7)}

Cluster2 = {(7,4)(6,2)(6,4)(7,3)(8,5)(7,6)}

Since the points (2,6) (3,8) and (4,7) are close to c1 hence they form one cluster whilst remaining points form another cluster.

So the total cost involved is 20.

Where cost between any two points is found using formula

$\text{cost}(x,c) = \sum |x-c|$ from $i=1$ to d

where x is any data object, c is the medoid, and d is the dimension of the object which in this case is 2.

Total cost is the summation of the cost of data object from its medoid in its cluster so here:

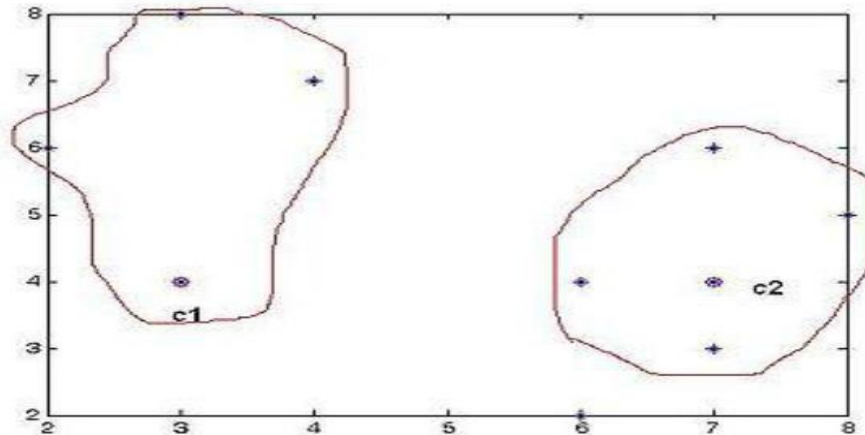
Total cost = { $\text{cost}((3,4),(2,6)) + \text{cost}((3,4),(3,8)) + \text{cost}((3,4),(4,7))$ }

+ $\text{cost}((7,4),(6,2)) + \text{cost}((7,4),(6,4)) + \text{cost}((7,4),(7,3))$

+ $\text{cost}((7,4),(8,5)) + \text{cost}((7,4),(7,6))$ }

= (3+4+4) + (3+1+1+2+2)

= 20



Step 2

Selection of nonmedoid O' randomly

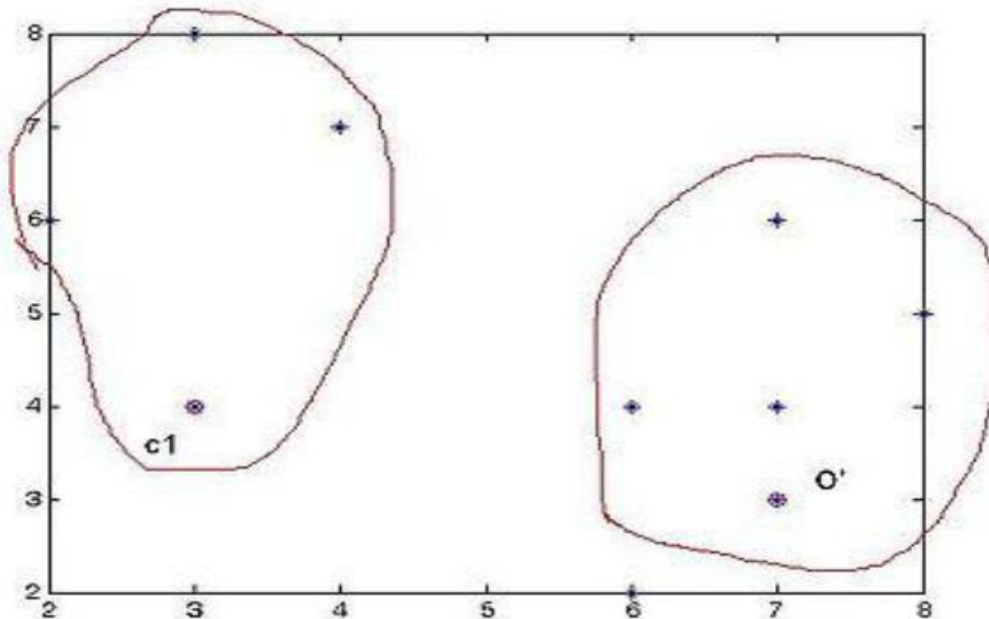
Let us assume $O' = (7,3)$

So now the medoids are $c1(3,4)$ and $O'(7,3)$

If $c1$ and O' are new medoids, calculate the total cost involved By using the formula in the step 1

C1		Data Objects (X_i)		Cost (distance)
3	4	2	6	3
3	4	3	8	4
3	4	4	7	4
3	4	6	2	5
3	4	6	4	3
3	4	7	4	4
3	4	8	5	6
3	4	7	6	6

O1		Data Objects (X_i)		Cost (distance)
7	3	2	6	8
7	3	3	8	9
7	3	4	7	7
7	3	6	2	2
7	3	6	4	2
7	3	7	4	1
7	3	8	5	3
7	3	7	6	3



total cost=3+4+4+2+2+1+3+3
=22

So cost of swapping medoid from c2 to O' is

S=current total cost-past total cost

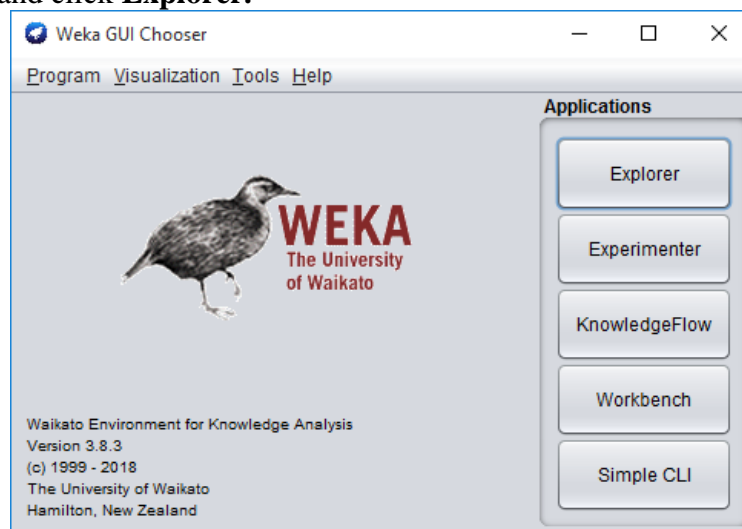
=22-20

=2

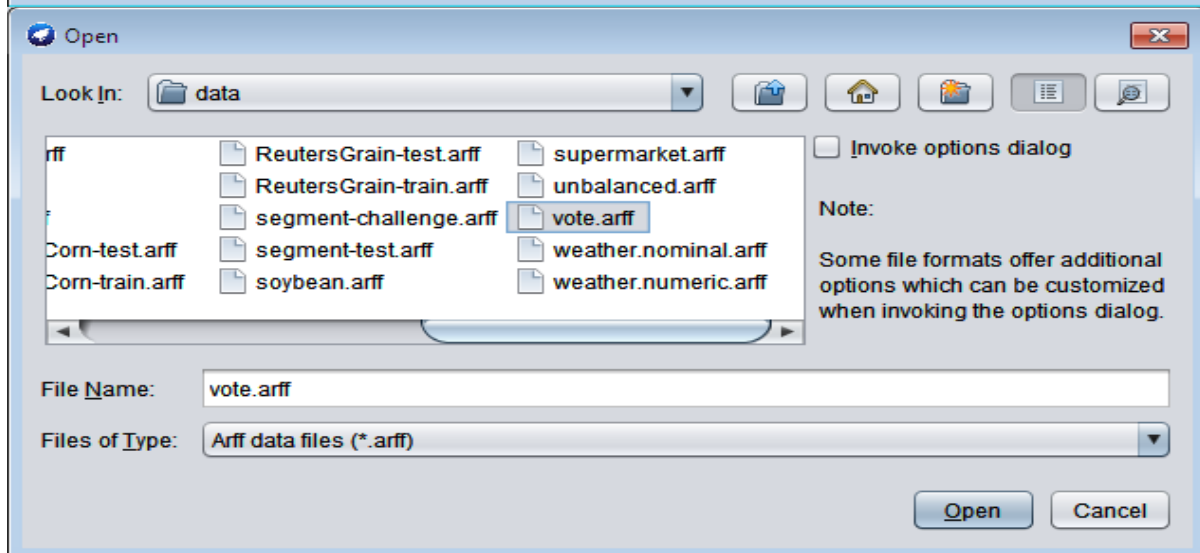
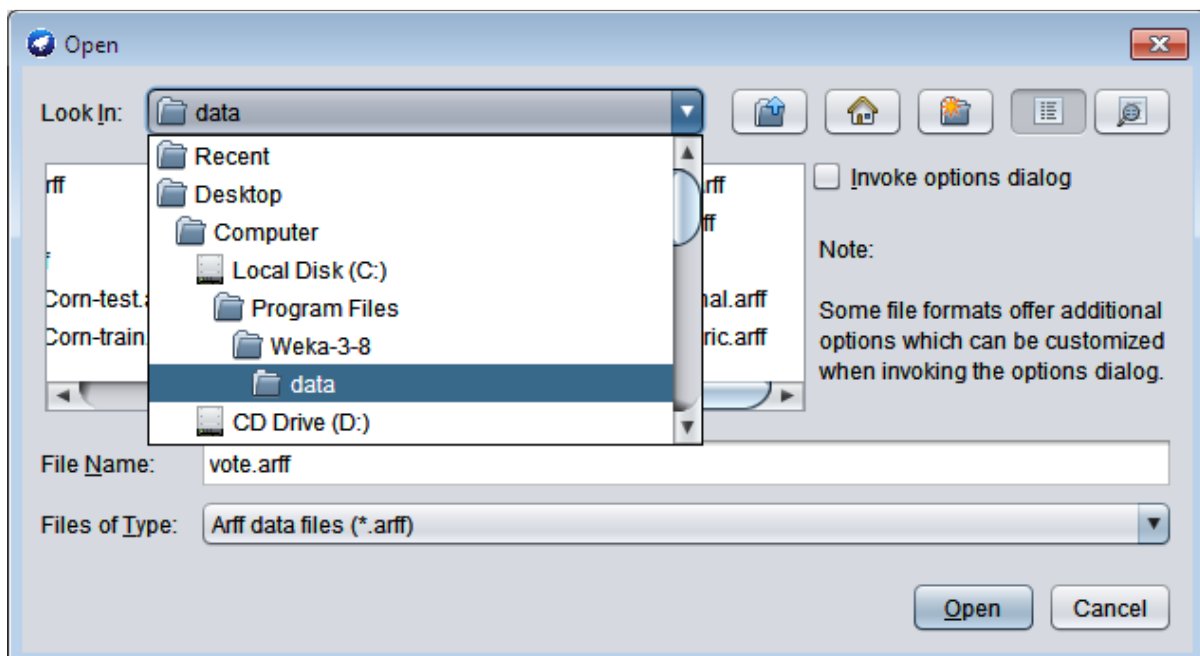
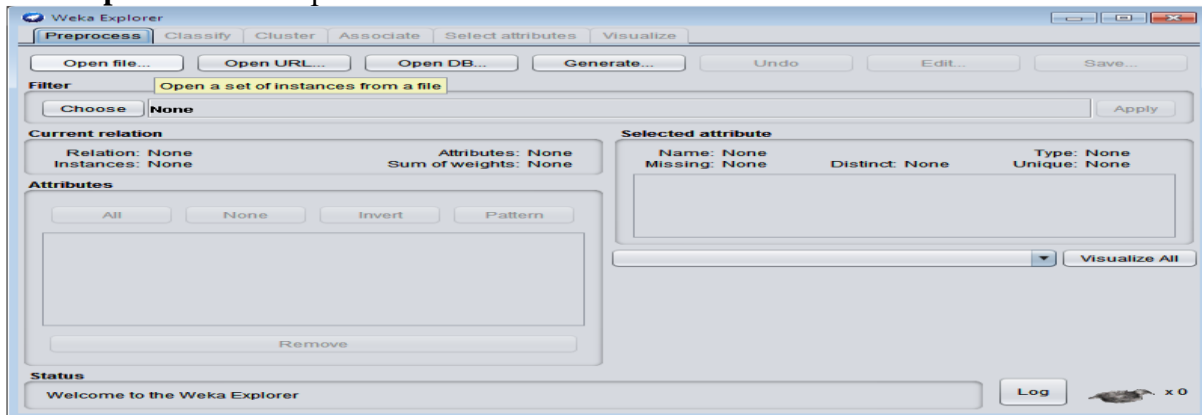
So moving to O' would be bad idea, so the previous choice was good and algorithm terminates here (i.e there is no change in the medoids).

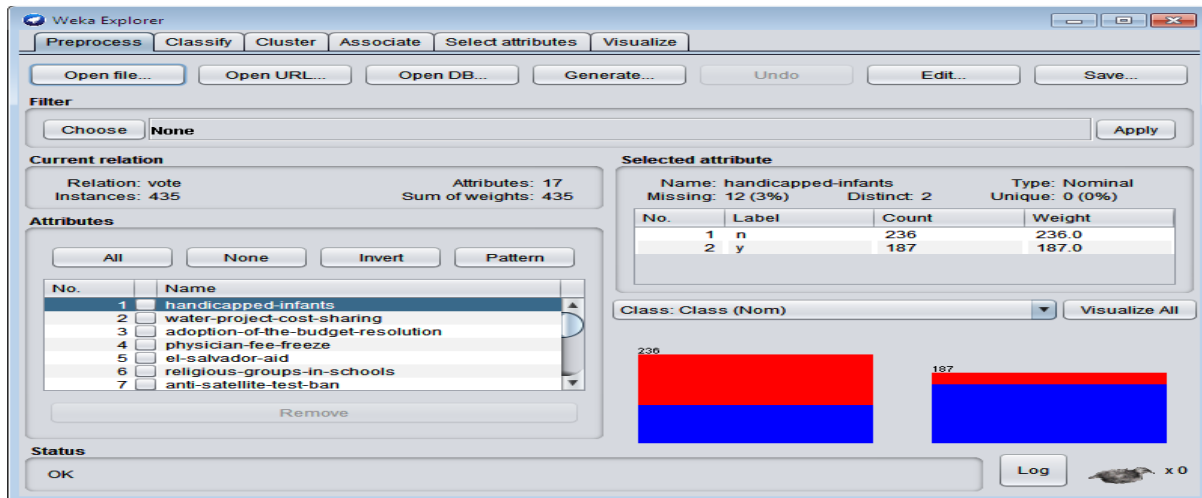
It may happen some data points may shift from one cluster to another cluster depending upon their closeness to medoid.

1. Open **weka** tool and click **Explorer**.

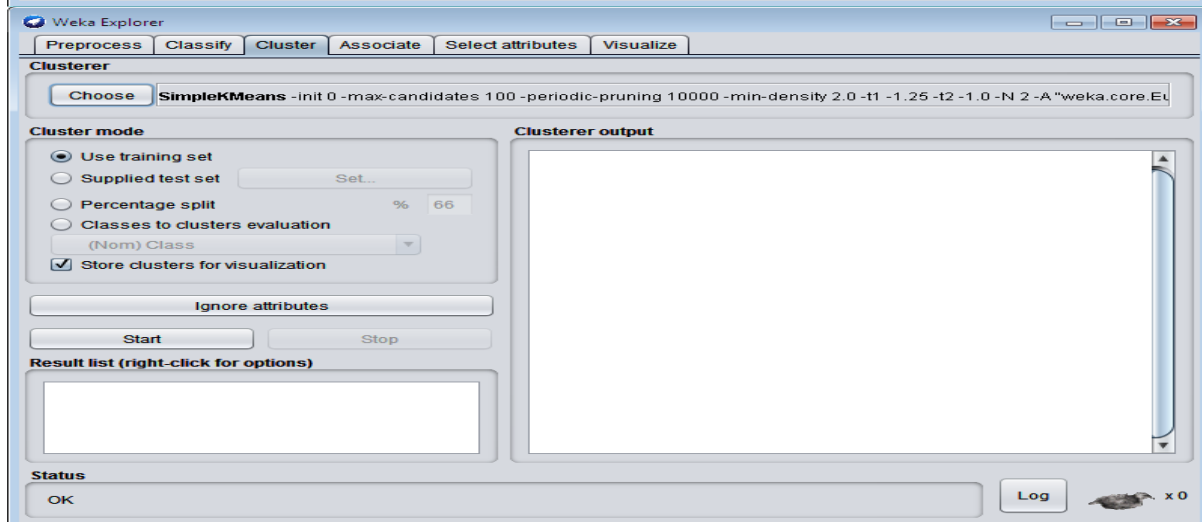
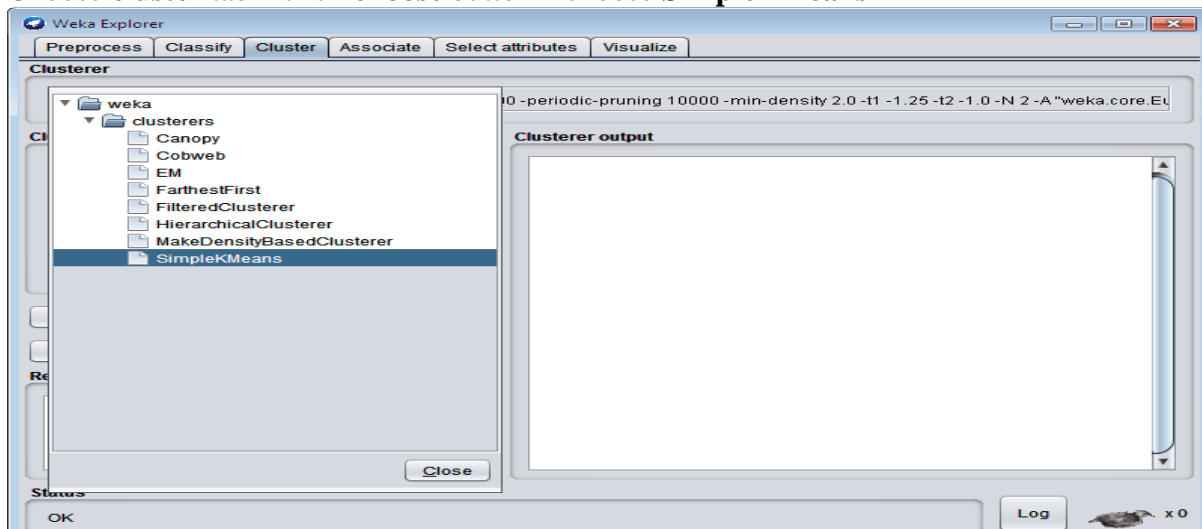


2. click **Open file...** in Preprocess tab.- choose **vote.arff**





Choose **cluster** tab – click **choose** button – choose **SimpleKmeans**



Click Start button

The screenshot shows the Weka Explorer interface with the 'Clusterer' tab selected. The 'Cluster mode' section has 'Use training set' selected. The 'Clusterer output' window displays the following information:

```

--- Run information ---
Scheme:      weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10
Relation:    vote
Instances:   495
Attributes:  17
handicapped-infants
water-project-cost-sharing
adoption-of-the-budget-resolution
physician-fee-freeze
el-salvador-aid
religious-groups-in-schools
anti-satellite-test-ban
aid-to-nicaraguan-contras
mx-missile
immigration
synfuels-corporation-cutback
education-spending
superfund-right-to-sue
crime
duty-free-exports
export-administration-act-south-africa
Class
Test mode:  evaluate on training data

--- Clustering model (full training set) ---

```

The 'Start' button is highlighted in the 'Cluster mode' section.

The screenshot shows the Weka Explorer interface with the 'Clusterer' tab selected. The 'Clusterer output' window displays the following information:

```

KMeans
-----
Number of iterations: 3
Within cluster sum of squared errors: 1510.0
Initial starting points (random):
Cluster 0: n,n,y,y,y,y,n,n,y,n,n,n,y,y,y, democrat
Cluster 1: n,n,y,n,y,n,y,y,n,n,n,n,y,n,y, democrat
Missing values globally replaced with mean/mode

Final cluster centroids:
Attribute          Full Data          Cluster#
                   (435.0)            (214.0)            (221.0)
-----
handicapped-infants      n                   n                   y
water-project-cost-sharing y                   y                   n
adoption-of-the-budget-resolution y                   n                   y
physician-fee-freeze     n                   y                   n
el-salvador-aid          y                   y                   n
religious-groups-in-schools y                   y                   n
anti-satellite-test-ban  y                   n                   y
aid-to-nicaraguan-contras y                   n                   y
mx-missile                y                   y                   y
immigration               y                   y                   y
synfuels-corporation-cutback n                   n                   n
education-spending       n                   y                   n
superfund-right-to-sue   y                   y                   n
crime                     y                   y                   n
duty-free-exports         y                   n                   y
export-administration-act-south-africa y                   y                   y
Class                    democrat republican democrat

```

The 'Start' button is highlighted in the 'Cluster mode' section.

The screenshot shows the Weka Explorer interface with the 'Clusterer' tab selected. The 'Clusterer output' window displays the following information:

```

physician-fee-freeze      n                   y                   n
el-salvador-aid          y                   y                   n
religious-groups-in-schools y                   y                   n
anti-satellite-test-ban  y                   n                   y
aid-to-nicaraguan-contras y                   n                   y
mx-missile                y                   n                   y
immigration               y                   y                   y
synfuels-corporation-cutback n                   n                   n
education-spending       n                   y                   n
superfund-right-to-sue   y                   y                   n
crime                     y                   y                   n
duty-free-exports         y                   n                   y
export-administration-act-south-africa y                   y                   y
Class                    democrat republican democrat

Time taken to build model (full training data) : 0.02 seconds

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

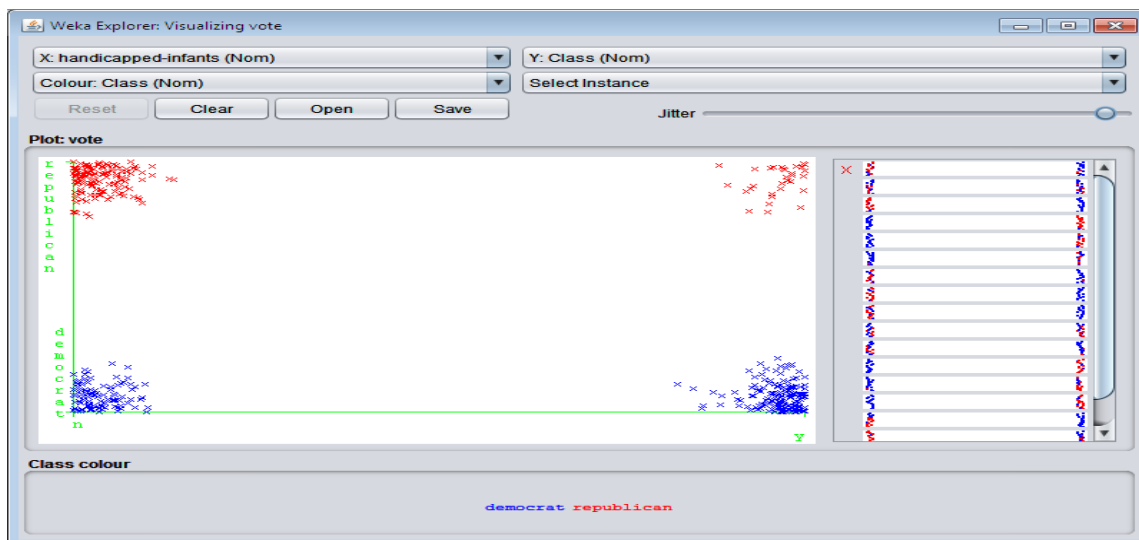
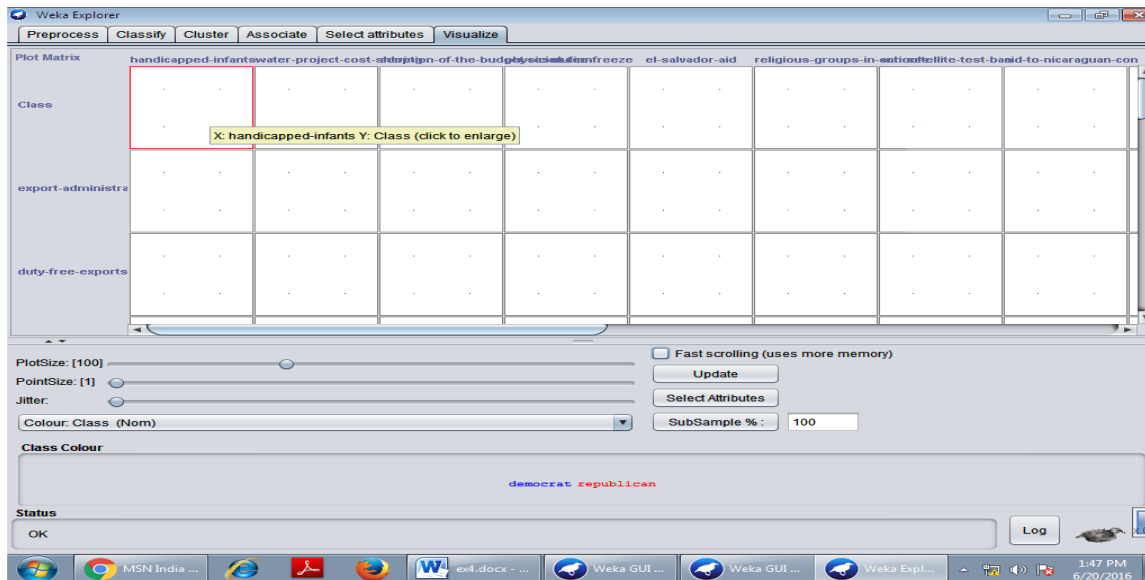
Clustered Instances

0      214 ( 49%)
1      221 ( 51%)

```

The 'Start' button is highlighted in the 'Cluster mode' section.

Goto - **Visualize** tab- click one box any visualize.



VIVA-QUESTIONS

1. A good clustering is one having
2. Which of the following is an exploratory data mining technique?
 - A. Classification
 - B. Clustering
 - C. Regression
3. Which of the following tasks can be best solved using Clustering.
 - A. Predicting the amount of rainfall based on various cues
 - B. Detecting fraudulent credit card transactions
 - C. Training a robot to solve a maze
4. Given two objects represented by the tuples (22, 1, 42, 10) and (20, 0, 36, 8);
5. Compute the Manhattan distance between the two objects

6. Consider a set of five 2-dimensional points $p_1=(0, 0)$, $p_2=(0, 1)$, $p_3=(5, 8)$, $p_4=(5, 7)$, and $p_5=(0, 0.5)$.

Euclidean distance is the distance function. The k-means algorithm is used to cluster the points into two clusters. The initial cluster centers are p_1 and p_4 . The clusters after two iterations of k-means are:

7. Clustering is a
8. Which of the following algorithm is most sensitive to outliers?

K-means clustering algorithm

K-medians clustering algorithm

K-modes clustering algorithm

K-medoids clustering algorithm

9. K means and K-medoids are example of which type of clustering method?

P1 P2 P3 P4

P1	0	1.000	1.414	1.000
P2	1.000	0	1.000	1.414
P3	1.414	1.000	0	1.000
P4	1.000	1.414	1.000	0

The Euclidean distance matrix between four 2-dimensional points, p_1 , p_2 , p_3 , and p_4 , is shown below. A possible set of co-ordinate values of these points are:

10. Consider a set of five 2-dimensional points $D_1=(2,0)$, $D_2=(1,3)$, $D_3=(3,5)$, $D_4=(2,2)$ and $D_5=(4,6)$.

Euclidean distance is the distance function. The k-means algorithm is used to cluster the points into two clusters. The initial cluster centers are D_2 and D_4 . The clusters after two iterations of k-means are:

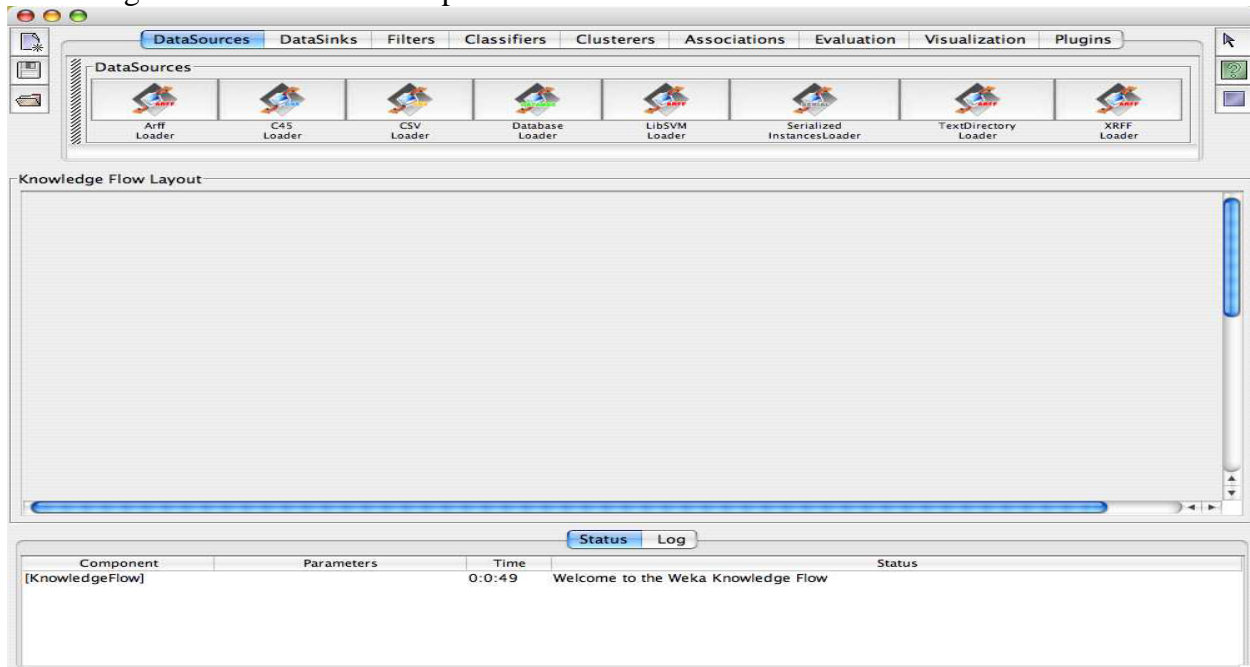
Experiment No. 8: Demonstrate performing knowledge flow on WEKA

- i. Perform pre-processing tasks.
- ii. Perform decision-tree
- iii. Perform clustering

KnowledgeFlow

Introduction

The KnowledgeFlow provides an alternative to the Explorer as a graphical front end to WEKA's core algorithms. The KnowledgeFlow is a work in progress so some of the functionality from the Explorer is not yet available. On the other hand, there are things that can be done in the KnowledgeFlow but not in the Explorer.



The Knowledge Flow presents a data-flow inspired interface to WEKA. The user can select WEKA components from a tool bar, place them on a layout canvas and connect them together in order to form a knowledge flow for processing and analyzing data. At present, all of WEKA's classifiers, filters, clusters, loaders and savers are available in the Knowledge Flow along with some extra tools.

WEKA there are ten classifiers that can handle data incrementally:

- AODE
- IB1
- IBk
- KStar
- NaiveBayesMultinomialUpdateable
- NaiveBayesUpdateable
- NNge
- Winnow

Features

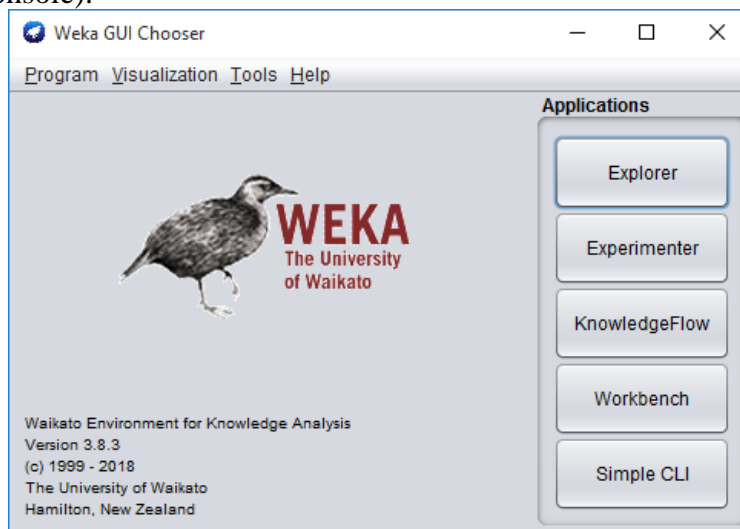
The Knowledge Flow offers the following features:

- intuitive data flow style layout

- process data in batches or incrementally
- process multiple batches or streams in parallel (each separate flow executes in its own thread)
- chain filters together
- view models produced by classifiers for each fold in a cross validation
- visualize performance of incremental classifiers during processing (scrolling plots of classification accuracy, RMS error, predictions etc.)
- plugin facility for allowing easy addition of new components to the Knowledge Flow

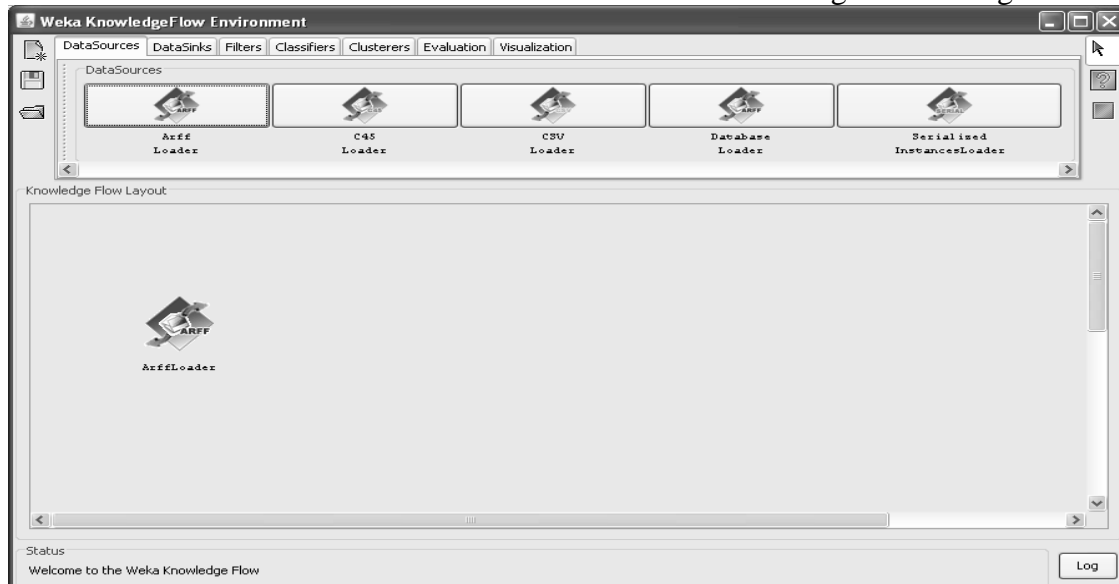
KNOWLEDGE FLOW FOR PRE-PROCESSING ADD ATTRIBUTE

Step-1: - Click on start button and then select All Programs and choose WEKA 3.8.3 in the WEKA 3.8(with console).

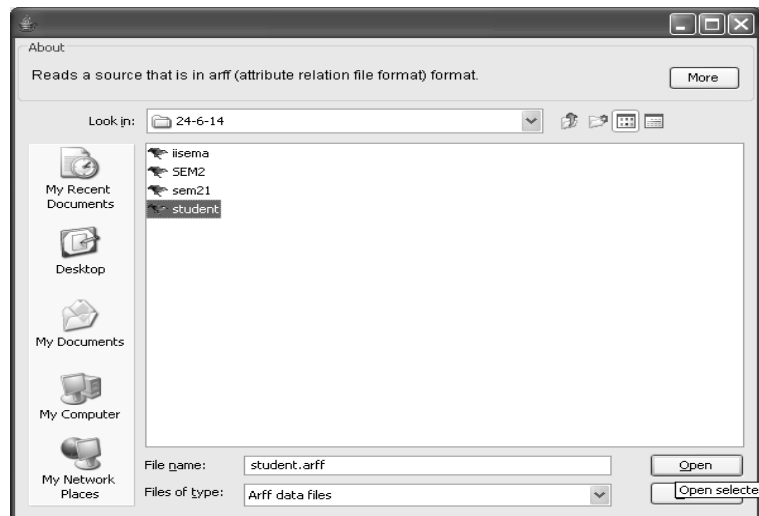


Step-2: - In the WEKA GUI to select the knowledge flow.

Step-3: - Select the data sources and then select ARFF loader and dragon knowledge flow layout.

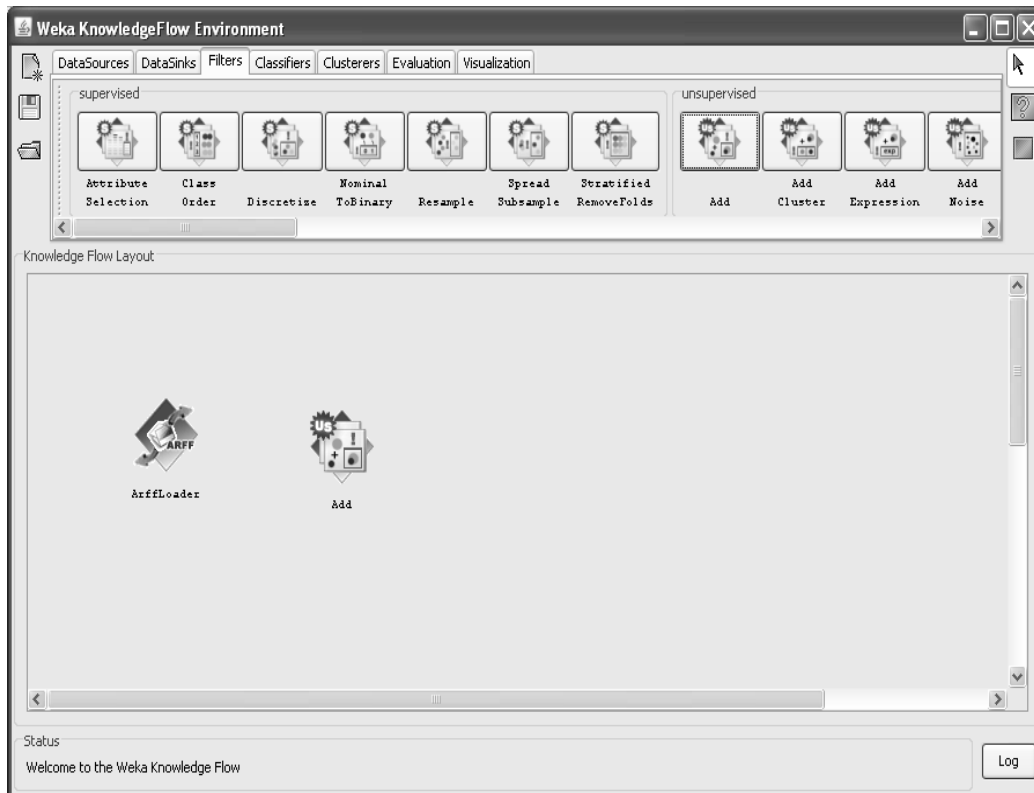


Step-4: - Select ARFF File right click in click configure to attach the file.

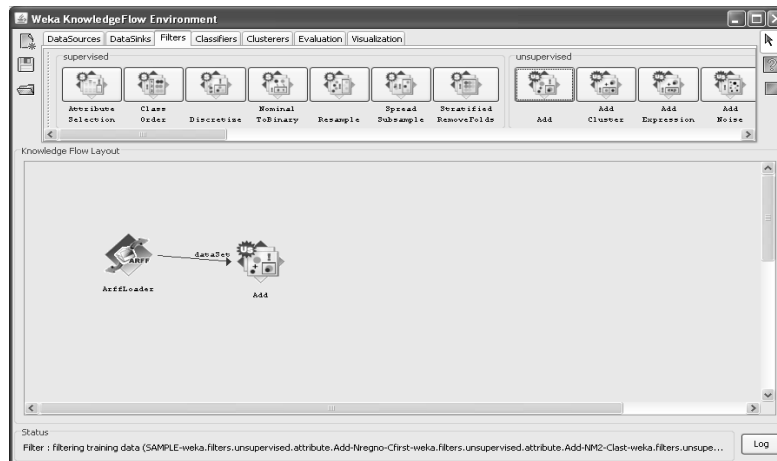


Click on open button.

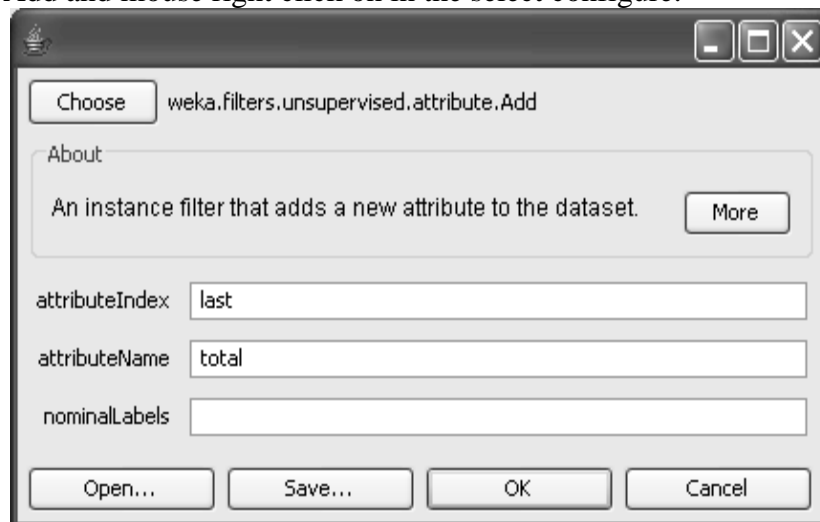
Step-5: - Go to Filters and select unsupervised in the Add Field and drag knowledge flow layout.



Step-6: - Select ARFF loader mouse right click in the connection to select data set and connection between ARFF loader to Add.

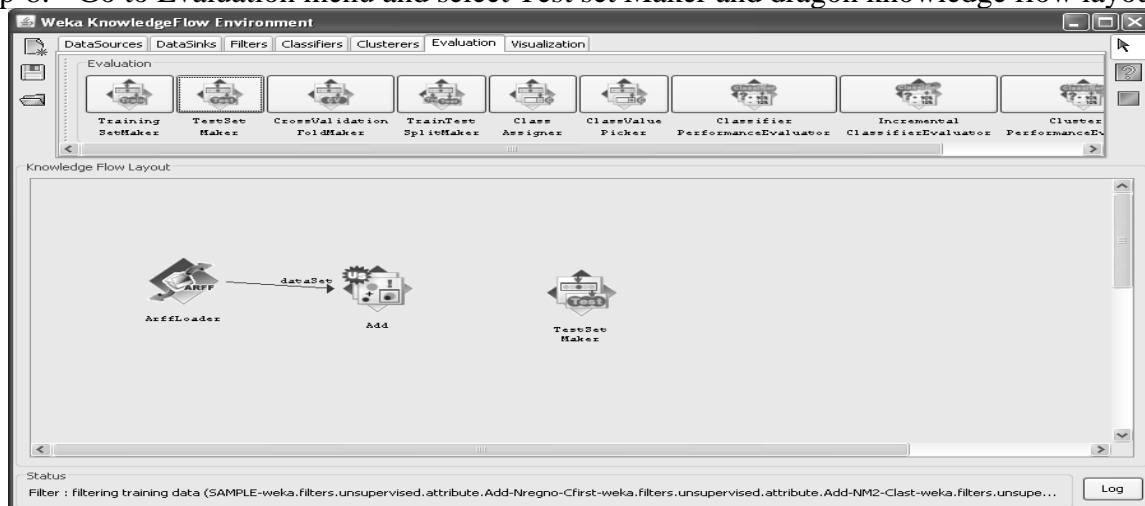


Step-7: - Select Add and mouse right click on in the select configure.

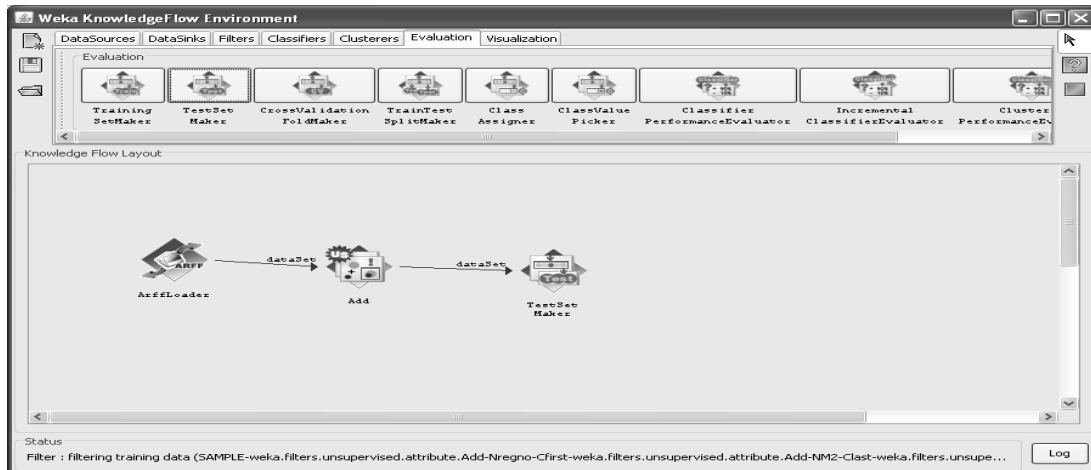


Click on ok button.

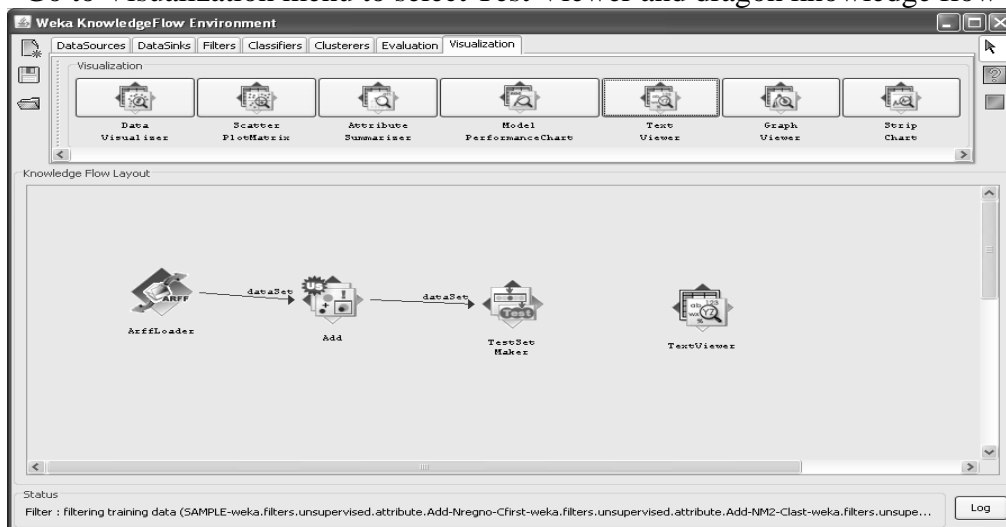
Step-8: - Go to Evaluation menu and select Test set Maker and dragon knowledge flow layout.



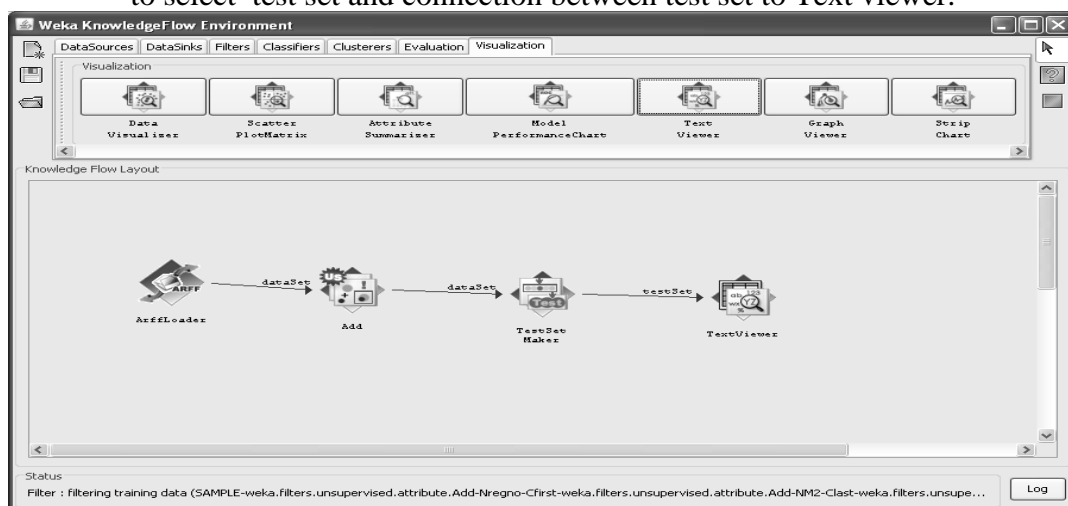
Step-9: - Select Add in the knowledge flow layout and mouse right click in the connection data set and connection between Add to Test Set Maker.



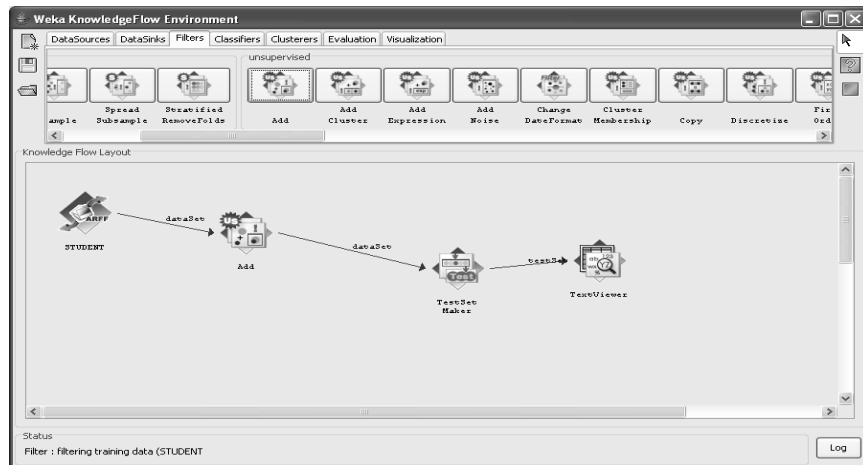
Step-10: - Go to Visualization menu to select Test Viewer and drag knowledge flow layout.



Step-11: - Select Test Set Maker in knowledge flow layout and mouse right click in connection to select test set and connection between test set to Text viewer.



Step-12: - Select ARFF Loader in the knowledge flow layout in the action select start loading.



Step-13: - Select Text Viewer in the knowledge flow layout mouse right click in the action to select show result.

The screenshot shows the 'Text Viewer' window. On the left, there is a 'Result list' section showing the time '04:38:27 - STUDENT-weka.filter'. On the right, there is a 'Text' section showing the output of the 'Text Viewer' operation. The output is a text representation of the data, starting with a schema definition and followed by a list of data rows. The schema defines attributes: SNO (string), SNAME (string), ENGLISH (numeric), CPDS (numeric), CO (numeric), DSGT (numeric), PeS (numeric), AcFM (numeric), CPDS-LAB (numeric), CO-LAB (numeric), and total (numeric). The data rows are numbered from 12481F0001 to 12481F0025, each followed by a list of values for the attributes, with some values being question marks.

```

@relation STUDENT-weka.filters.unsupervised.attribute.Add-Ntotal-Clas

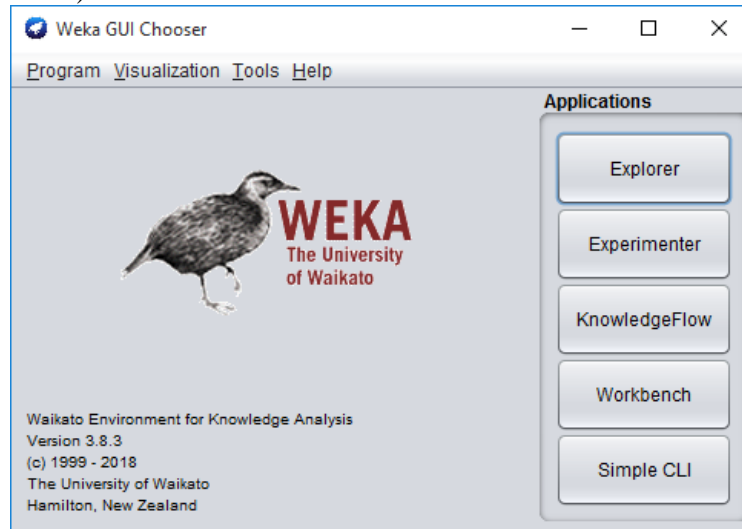
@attribute SNO string
@attribute SNAME string
@attribute ENGLISH numeric
@attribute CPDS numeric
@attribute CO numeric
@attribute DSGT numeric
@attribute PeS numeric
@attribute AcFM numeric
@attribute CPDS-LAB numeric
@attribute CO-LAB numeric
@attribute total numeric

@data
12481F0001,PLAKSHMI,74,75,79,80,85,78,74,80,?
12481F0002,BLAKSHMI,74,75,79,80,85,78,74,80,?
12481F0003,JEEVAN,84,85,89,80,86,77,74,80,?
12481F0004,PHANI,74,88,69,80,65,88,84,80,?
12481F0005,PAVANI,74,75,79,80,85,78,74,72,?
12481F0006,TRINADH,74,75,79,80,85,78,74,75,?
12481F0007,GOUTHAMI,74,75,79,80,85,78,74,80,?
12481F0008,LAKSHMI,74,75,79,80,85,78,74,60,?
12481F0009,GIRISH,74,75,79,80,85,78,74,70,?
12481F0010,GRACE,74,75,79,80,85,78,74,76,?
12481F0011,DLAKSHMI,74,75,79,80,85,78,74,75,?
12481F0012,SYAM,84,85,79,80,62,88,74,85,?
12481F0013,SPARJAN,74,75,79,80,85,78,74,70,?
12481F0014,ESWAR,74,75,79,80,85,78,74,70,?
12481F0015,SUBHADRA,74,75,69,80,85,78,74,70,?
12481F0016,PHANI,74,75,79,60,85,78,74,70,?
12481F0017,BRAMAHYA,74,75,89,80,85,78,74,70,?
12481F0018,DEEPTHI,74,75,70,80,85,78,74,70,?
12481F0019,SRAVYALATHA,74,75,79,80,85,78,74,70,?
12481F0020,NAGABABU,74,75,79,80,85,78,74,70,?
12481F0021,'HANIKANTA SAI',74,75,79,80,85,78,74,70,?
12481F0022,RAMAKRISHNA,74,75,89,80,85,78,74,70,?
12481F0023,NAGARJUNA,74,75,79,80,85,78,74,70,?
12481F0024,BHULAKSHMI,74,75,89,80,85,78,74,70,?
12481F0025,PLAKSHMI,74,75,79,70,85,78,74,70,?

```

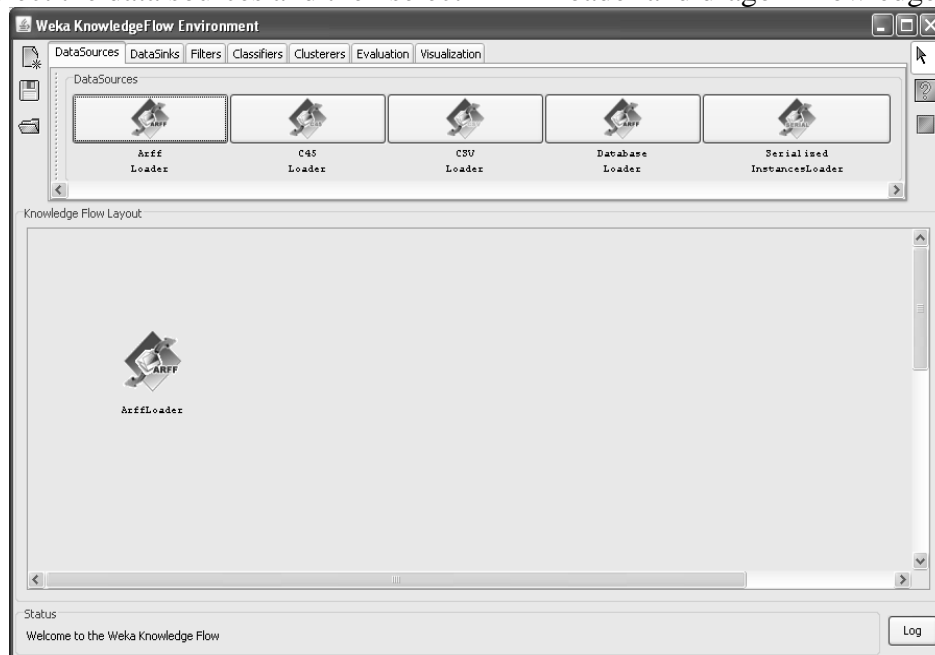
Add Expression.

Step-1: - Click on start button and then select All Programs and choose WEKA 3.8.3 in the WEKA 3.8(with console).

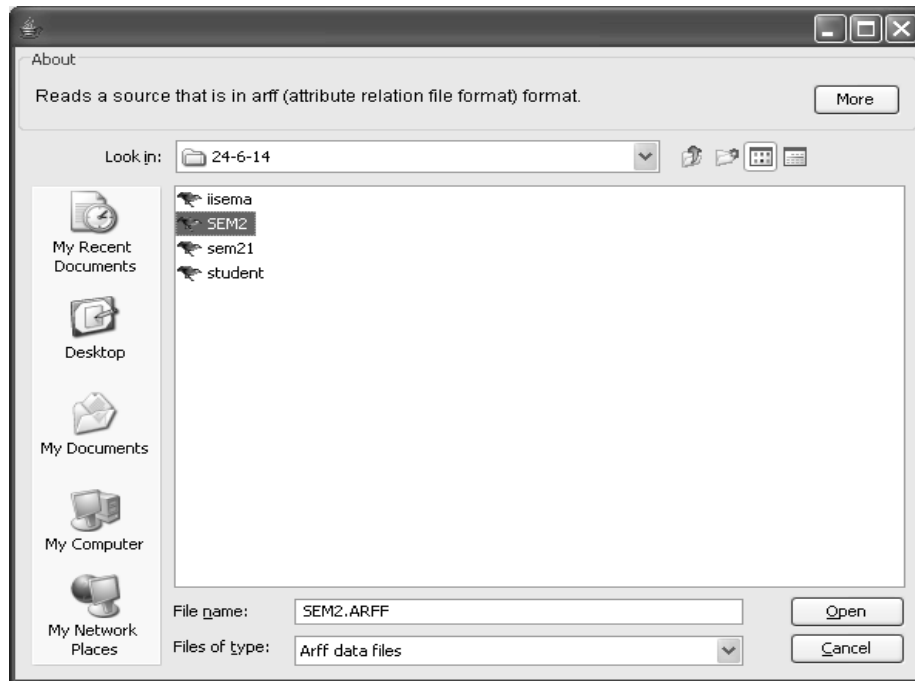


Step-2: - In the WEEKA GUI to select the knowledge flow.

Step-3: - Select the data sources and then select ARFF loader and dragon knowledge flow layout.

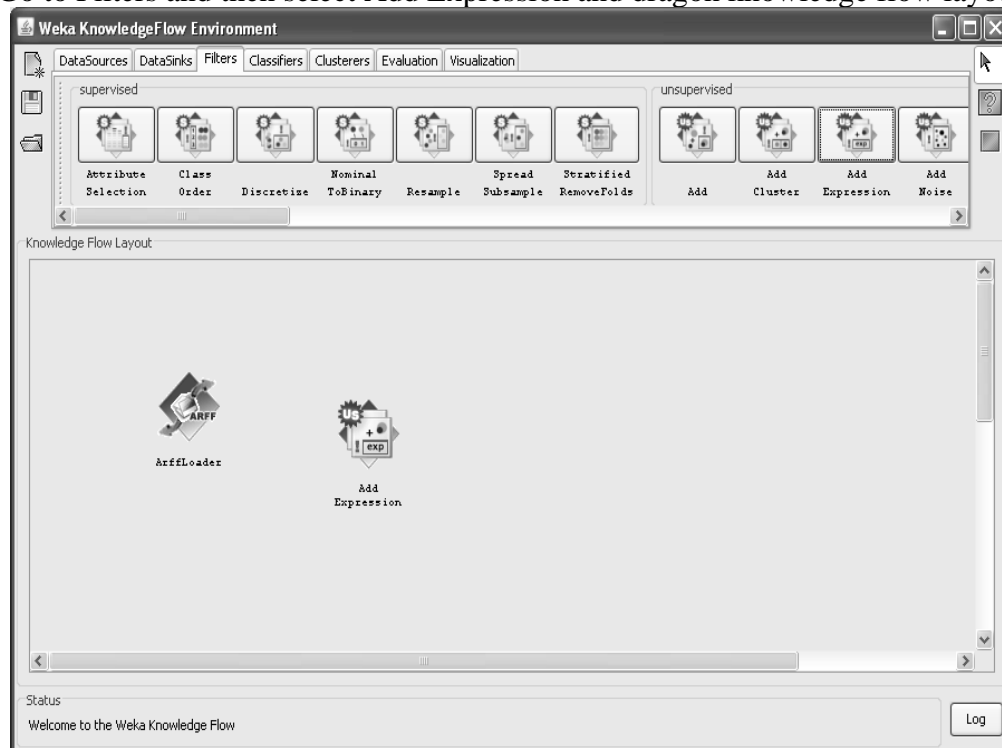


Step-4: - Select ARFF File right click in click configure to attach the file.

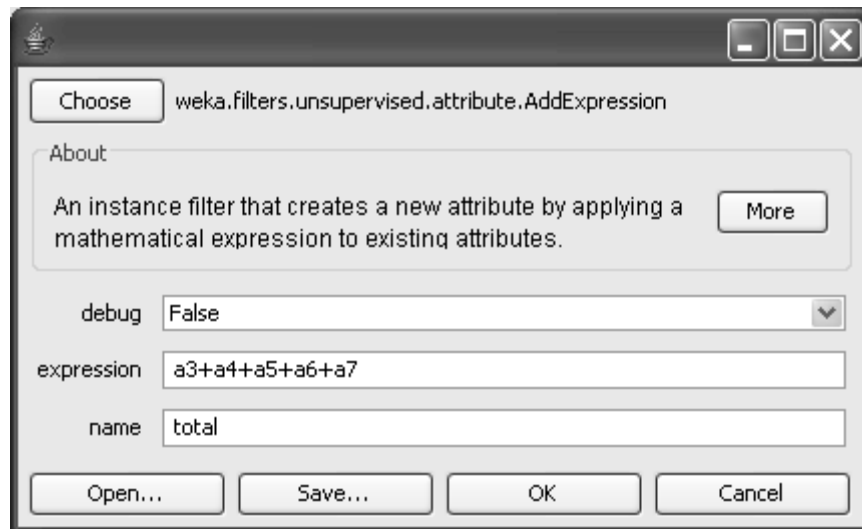


Click on open button.

Step-5: Go to Filters and then select Add Expression and dragon knowledge flow layout.

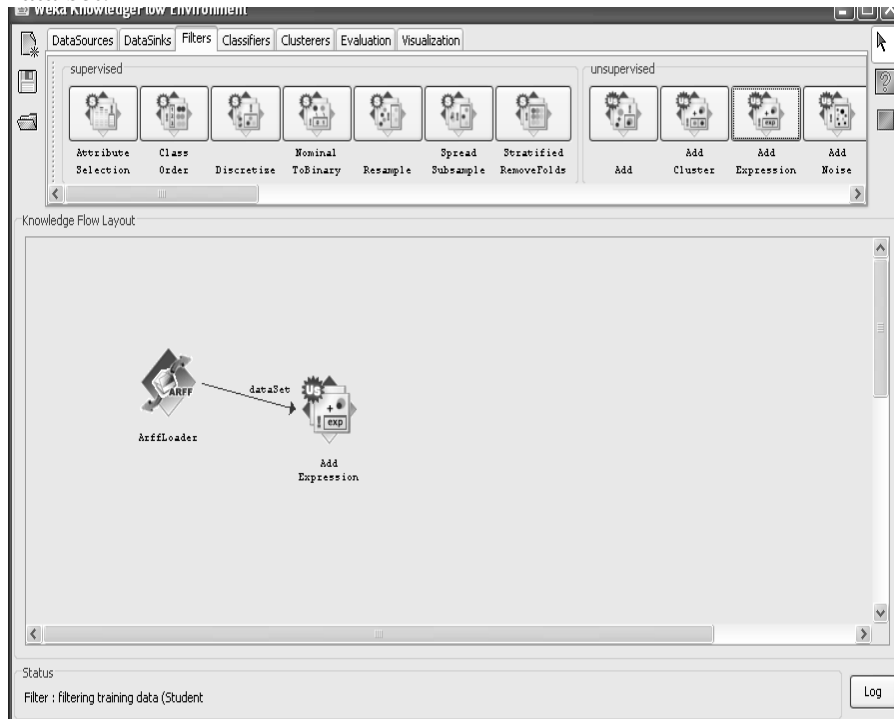


Step-6: - Select Add Expression and right click on select configure.

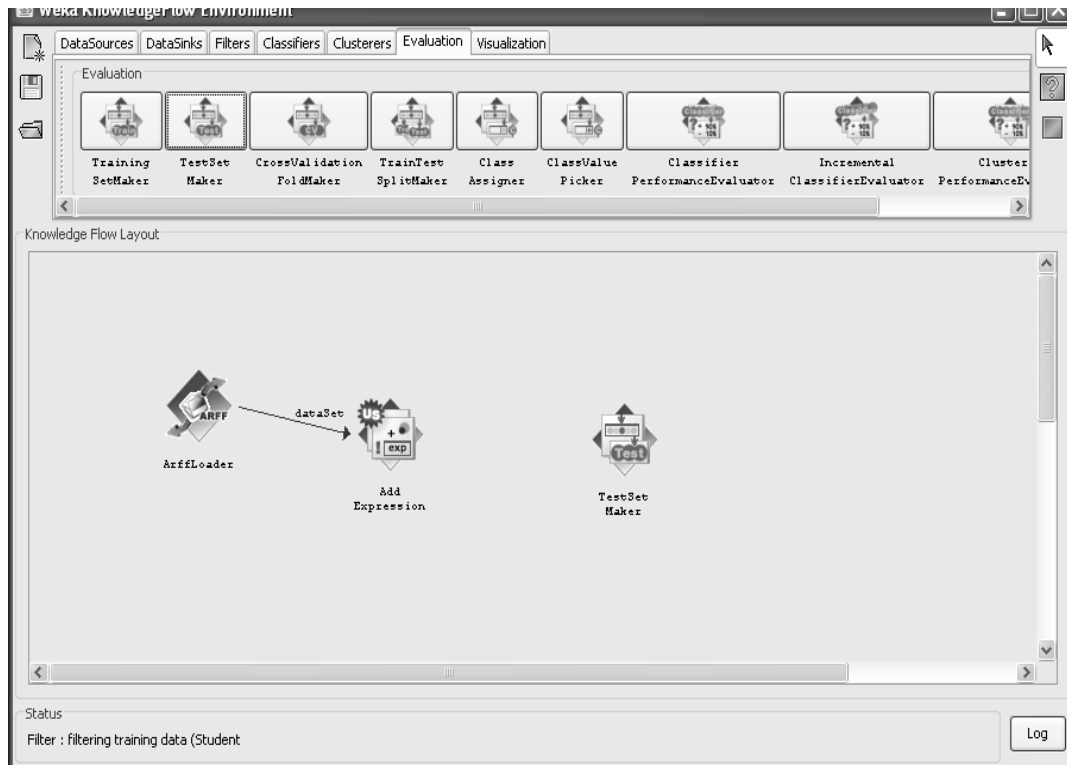


Click on ok button.

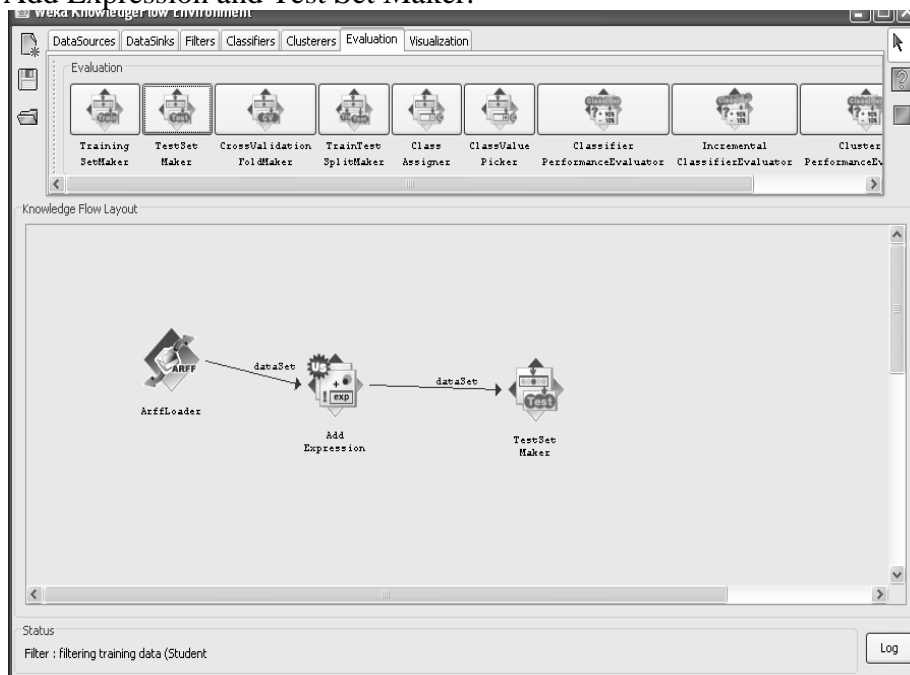
Step-7: - Select ARFF file to connection between the Add Expression to right click on ARFF file to select Data set.



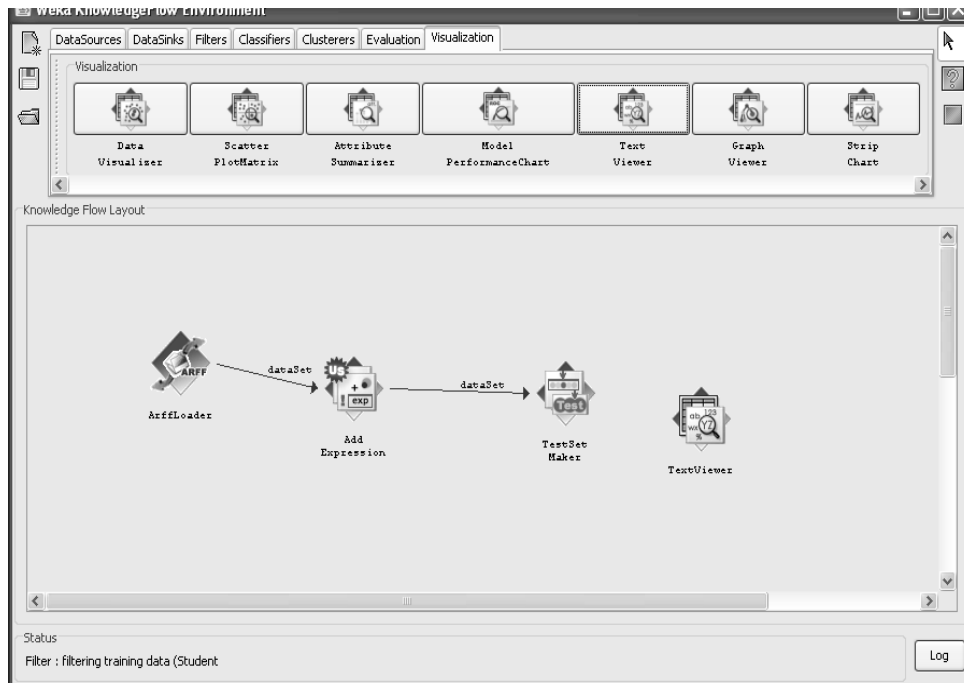
Step-8: - Go to Evaluation and select Test Set Maker drag on knowledge flow layout.



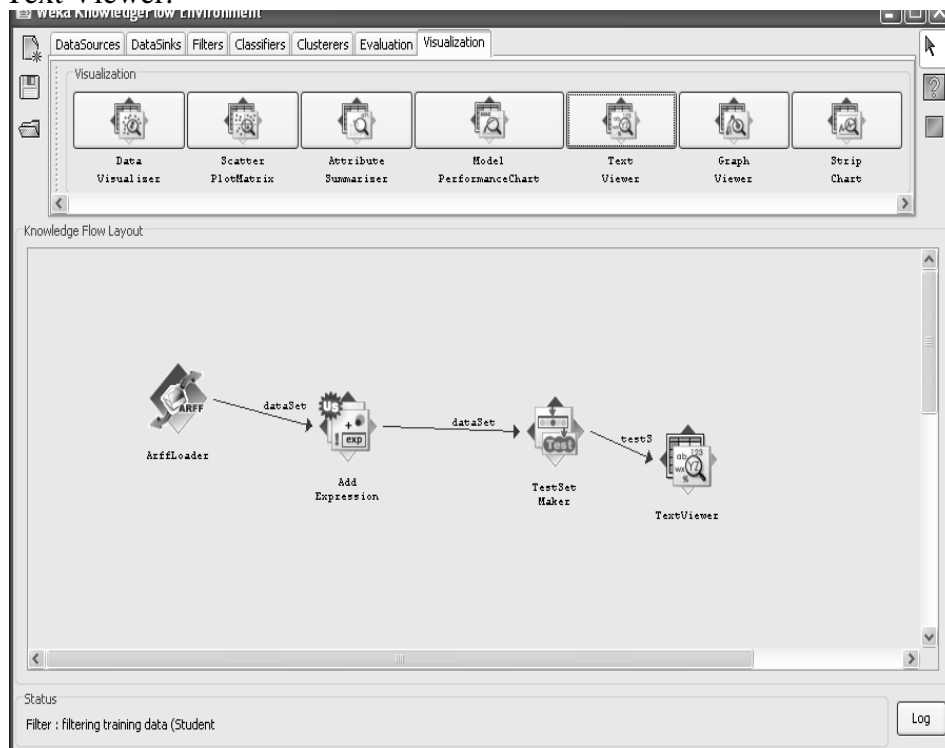
Step-9: - Select Add Expression and right click on to select Data set option to connection between the Add Expression and Test Set Maker.



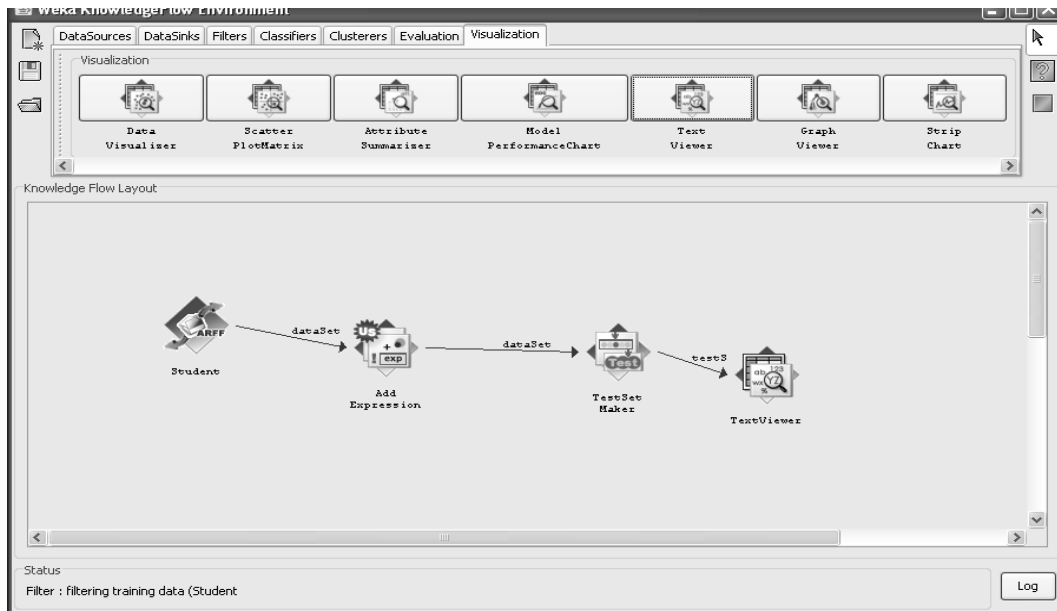
Step-10: - Go to visualization and select Text Viewer to drag on knowledge flow layout.



Step-11: - select Test Set maker to right click on to select test set to connction between Test Set Maker and Text Viewer.



Step-12: - select ARFF Loader to right click on start loading.



Step-13: - select Text Viewer to right Click on show result.

```

Result list
04:43:55 - STUDENT-weka.filter

Text
@relation STUDENT-weka.filters.unsupervised.attribute.AddExpression-Ea3+a4+a5+a6+a7-Nexpression

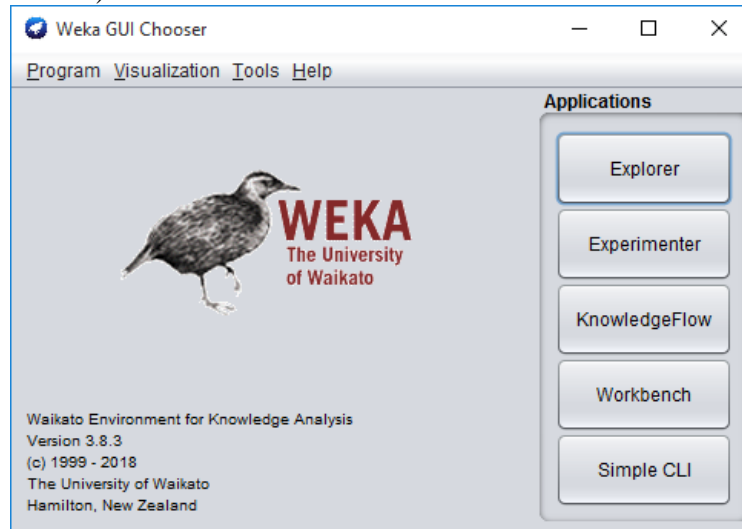
@attribute SNO string
@attribute SNAME string
@attribute JAVA numeric
@attribute COBOL numeric
@attribute OS numeric
@attribute OSHRM numeric
@attribute OR numeric
@attribute a3+a4+a5+a6+a7 numeric

@data
12481F0001,RLAKSHMI,85,86,87,79,89,426
12481F0002,BLAKSHMI,80,81,82,82,70,395
12481F0003,JEEVAN,79,75,85,86,90,415
12481F0004,PHANI,86,79,78,75,92,410
12481F0005,PAVANI,85,84,85,84,83,421
12481F0006,TRINADH,79,75,76,74,80,384
12481F0008,LAKSHMI,78,86,88,75,90,417
12481F0009,GIRISH,80,80,80,85,85,410
12481F0010,GRACE,85,86,85,80,85,421
12481F0011,LAKSHMI,85,81,90,79,80,415
12481F0012,SYAM,79,80,75,80,79,393
12481F0013,SPARJAN,95,76,80,75,79,405
12481F0014,ESWAR,90,79,79,76,80,404
12481F0015,SUBADHRA,85,80,76,80,70,391

```

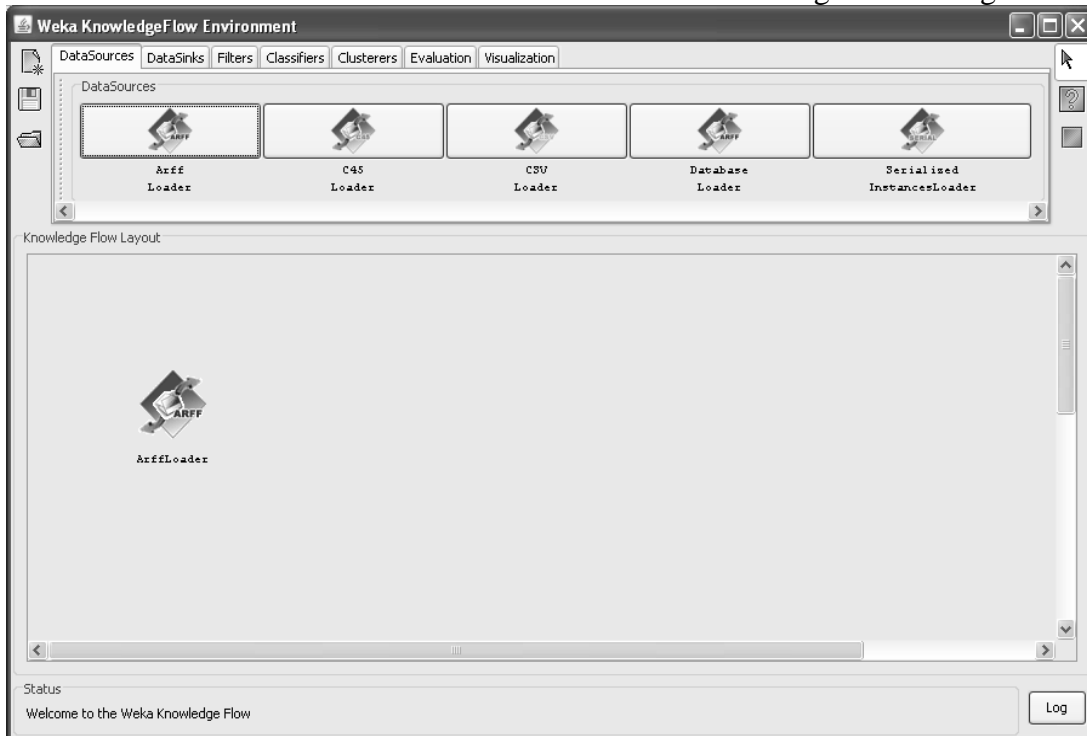
Copy Attribute.

Step-1: - Click on start button and then select All Programs and choose WEKA 3.8.3 in the WEKA 3.8(with console).

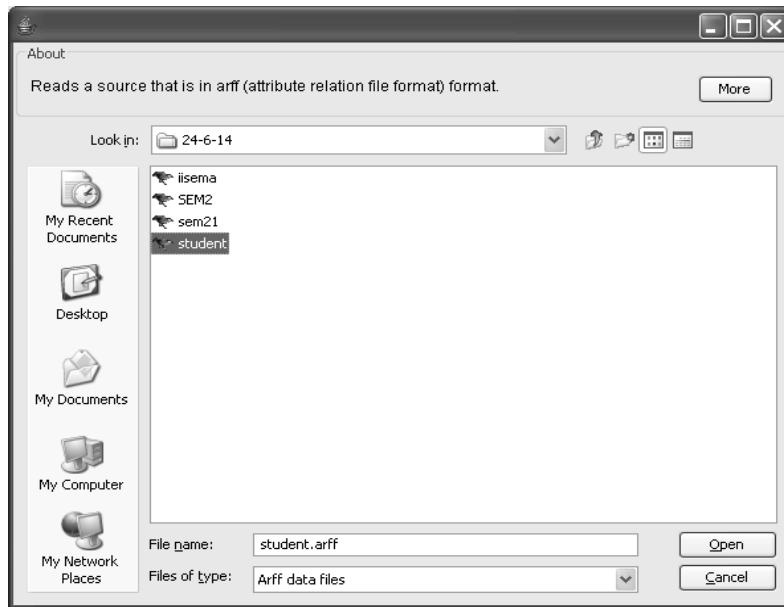


Step-2: - In the WEKA GUI to select the knowledge flow.

Step-3: - Select the data sources and then select ARFF loader and dragon knowledge flow layout.

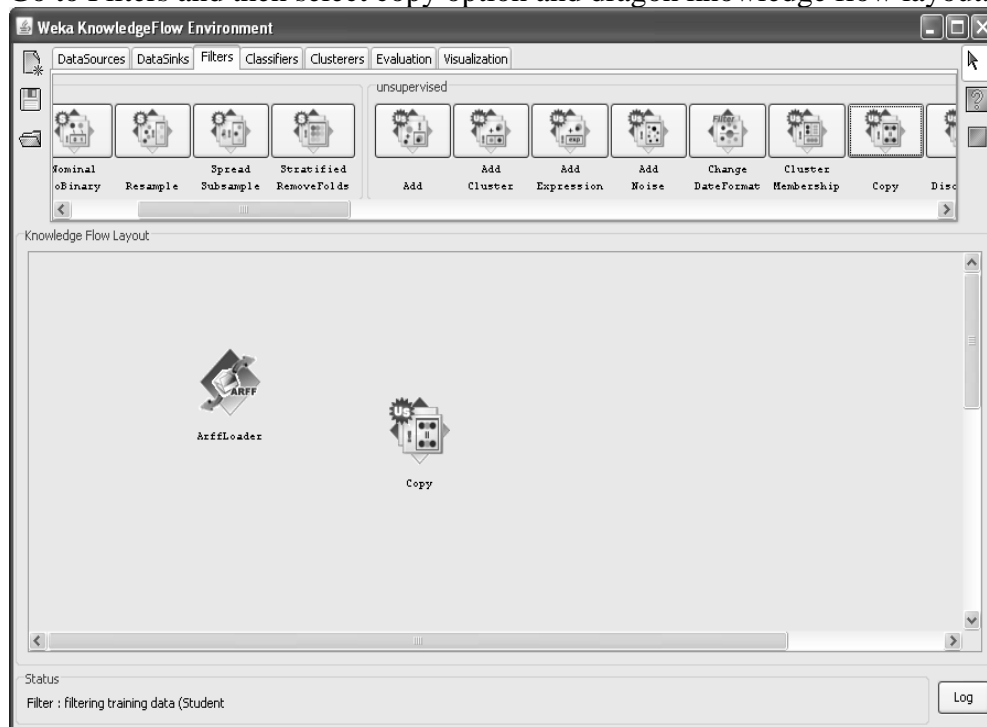


Step-4: - Select ARFF File right click in click configure to attach the file.



Click on open button.

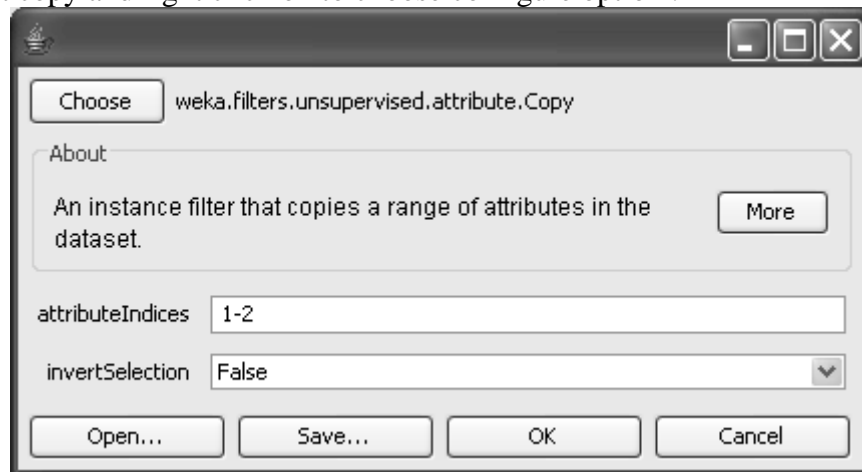
Step-5: - Go to Filters and then select copy option and dragon knowledge flow layout.



Step-6: - select the ARFF Loader and right click to choose Data Set to create connection between ARFF Loader and copy.

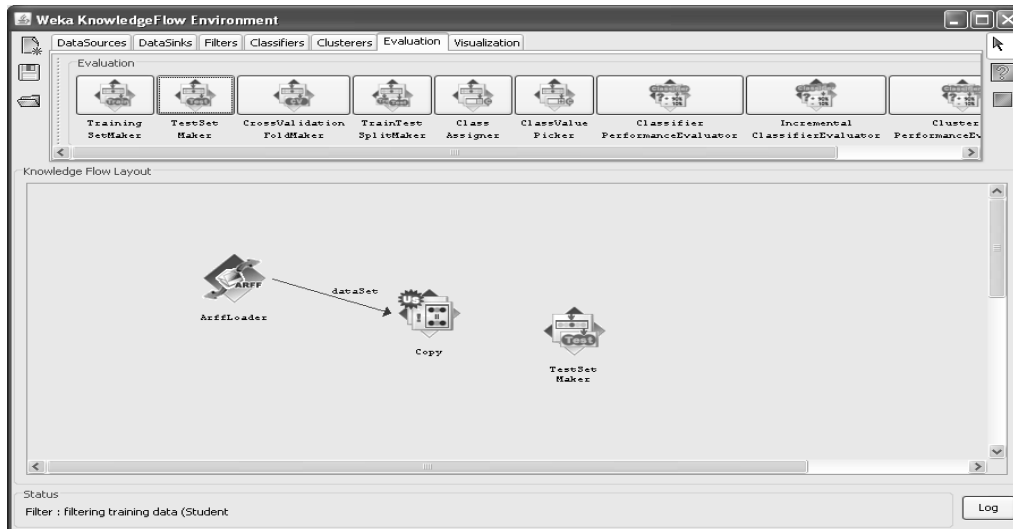


Step-7: - select copy and right click on to choose configure option .

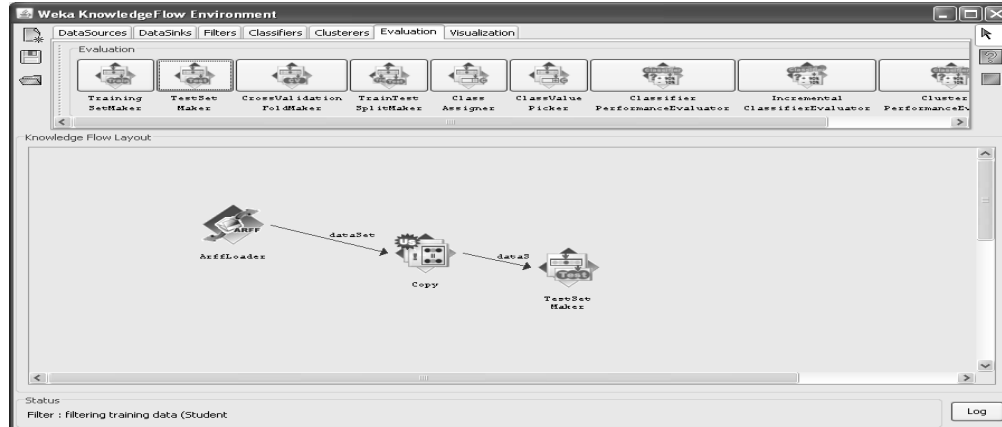


Click on ok button.

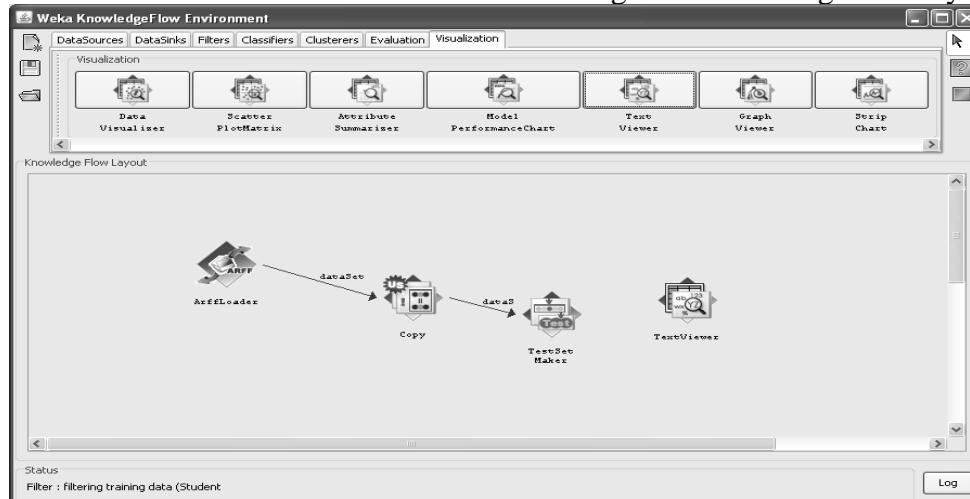
Step-8: - Go to Evaluation and select Test Set Maker drag on to knowledge flow layout.



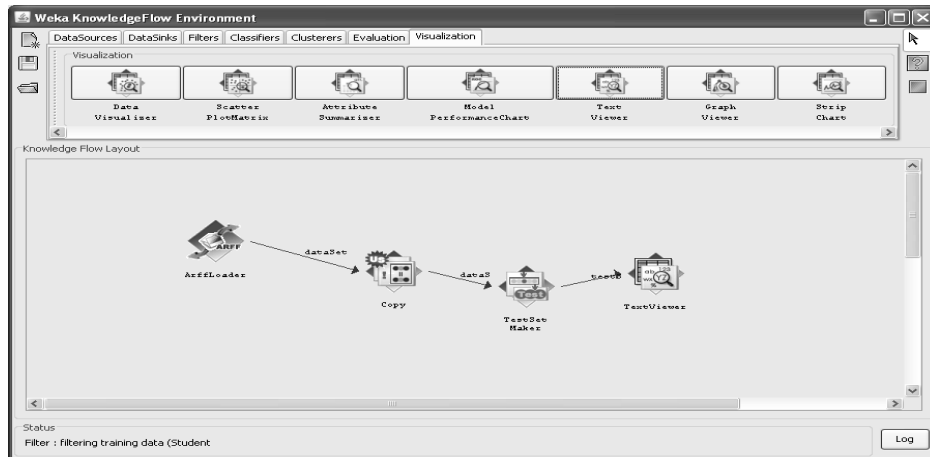
Step-9: - select copy and right click on to choose Data Set to connection between the copy and Test Set Maker.



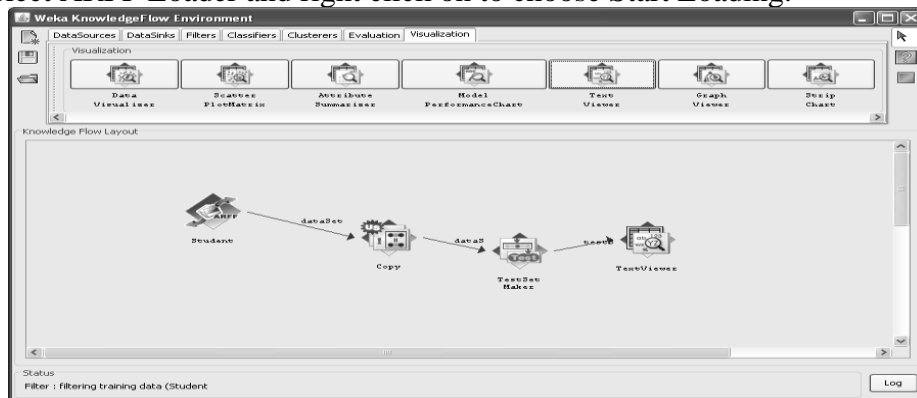
Step-10: - Go to Visualization and select Text Viewer drag on to knowledge flow layout.



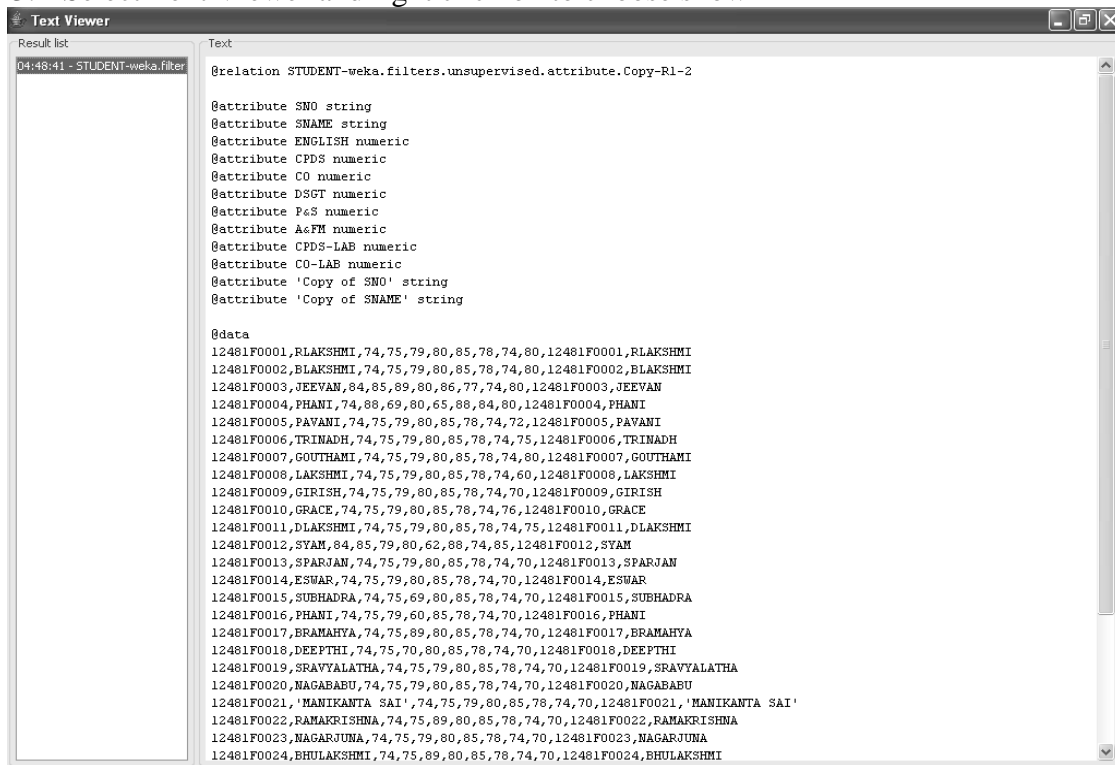
Step-11: - select Test Set Maker and right click on to choose Test Set to connection between Test Set Maker to Text Viewer.



Step-12: - select ARFF Loader and right click on to choose Start Loading.



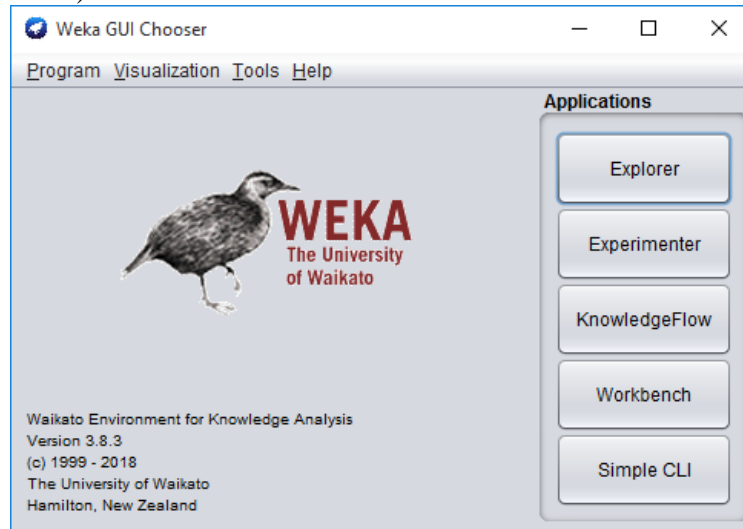
Step-13: - Select Text Viewer and right click on to choose show



result.

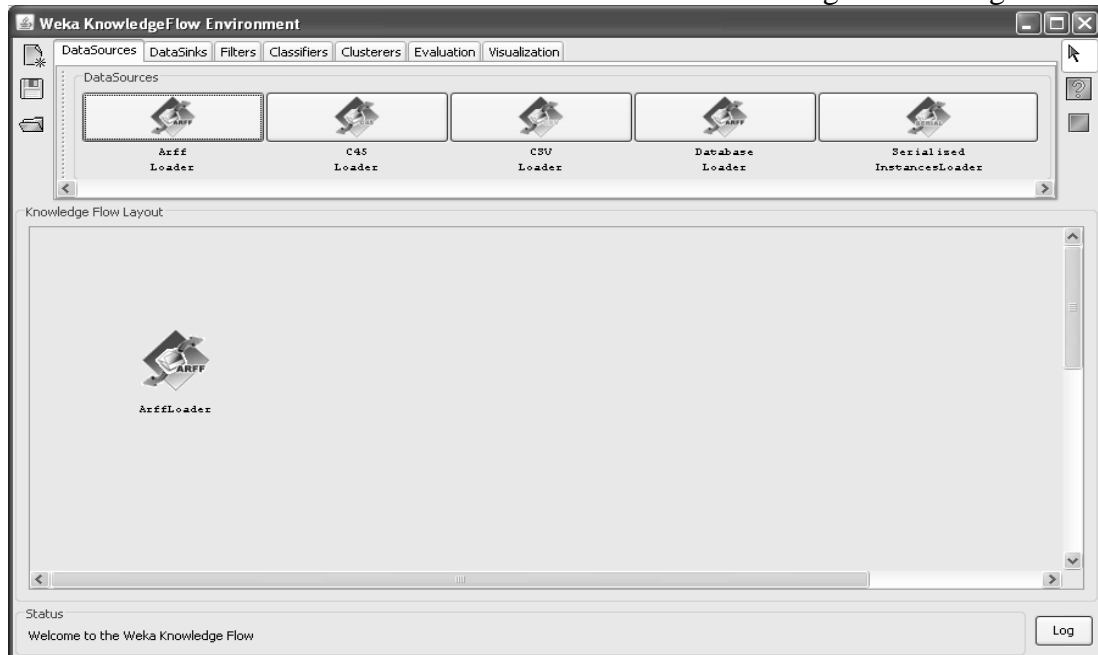
Remove Attribute

Step-1: - Click on start button and then select All Programs and choose WEKA 3.8.3 in the WEKA 3.8(with console).

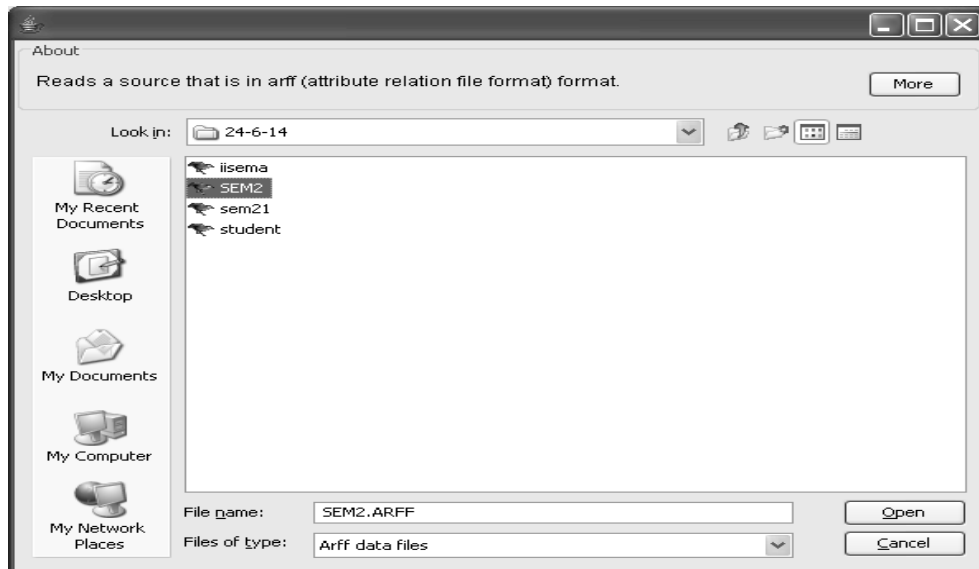


Step-2: - In the WEEKA GUI to select the knowledge flow.

Step-3: - Select the data sources and then select ARFF loader and dragon knowledge flow layout.



Step-4: - Select ARFF File right click in click configure to attach the file.

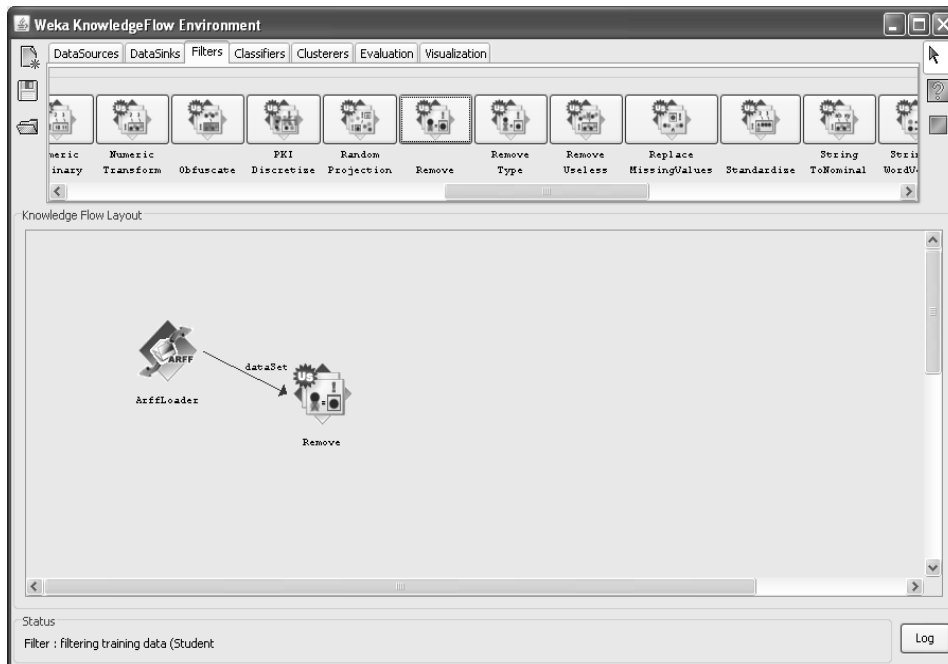


Click on open button.

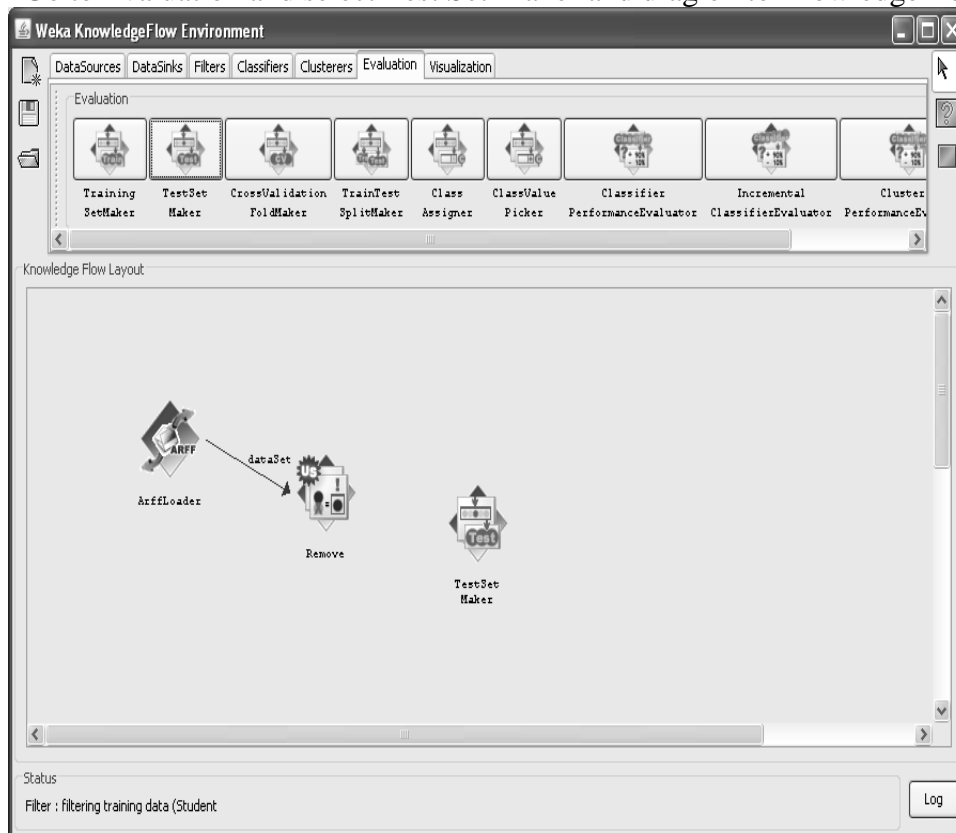
Step-5: - Go to Filters and then select Remove and dragon knowledge flow layout.



Step-6: - select ARFF Loader and right click on to choose Data Set to connection between the ARFF Loader and Remove.

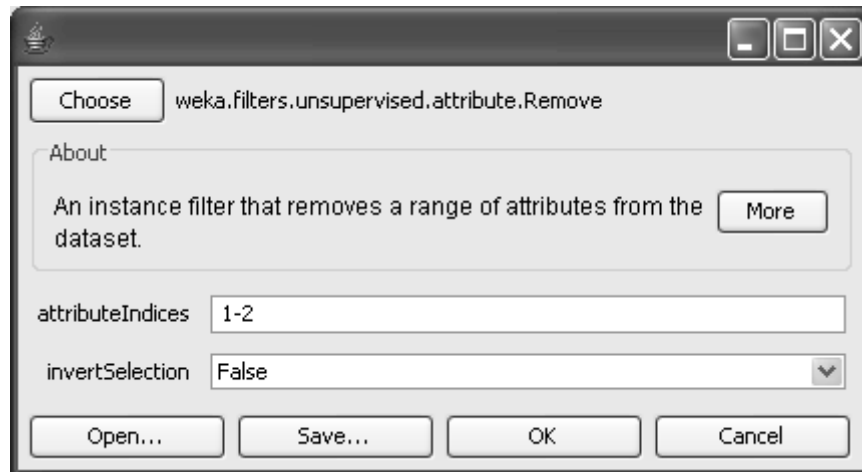


Step-7: - Go to Evaluation and select Test Set Maker and drag on to knowledge flow



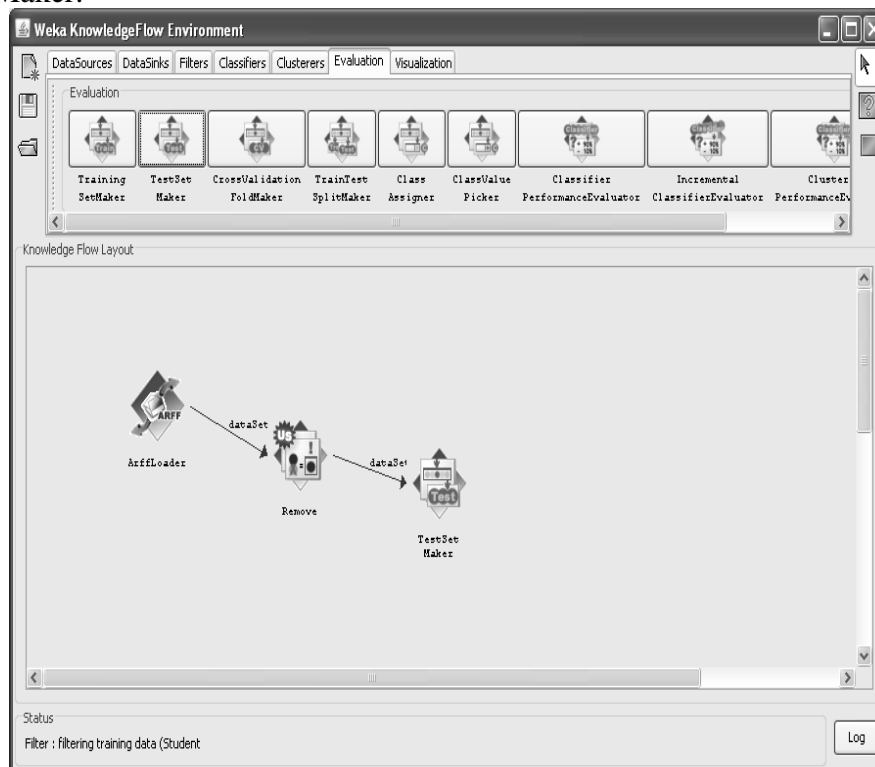
layout.

Step-8: -select Remove and right click on to choose configure .

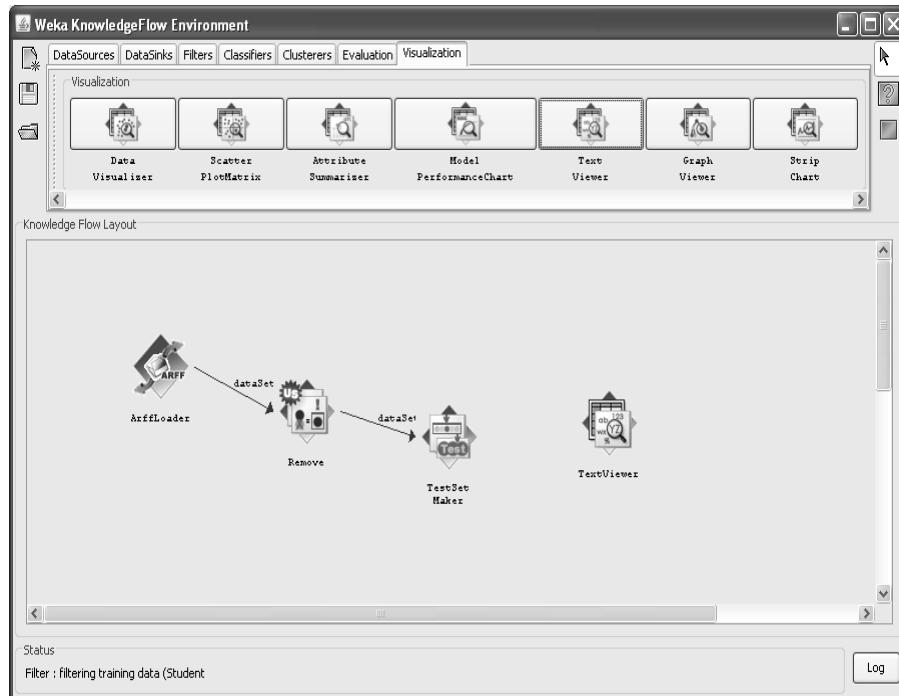


Click on ok button.

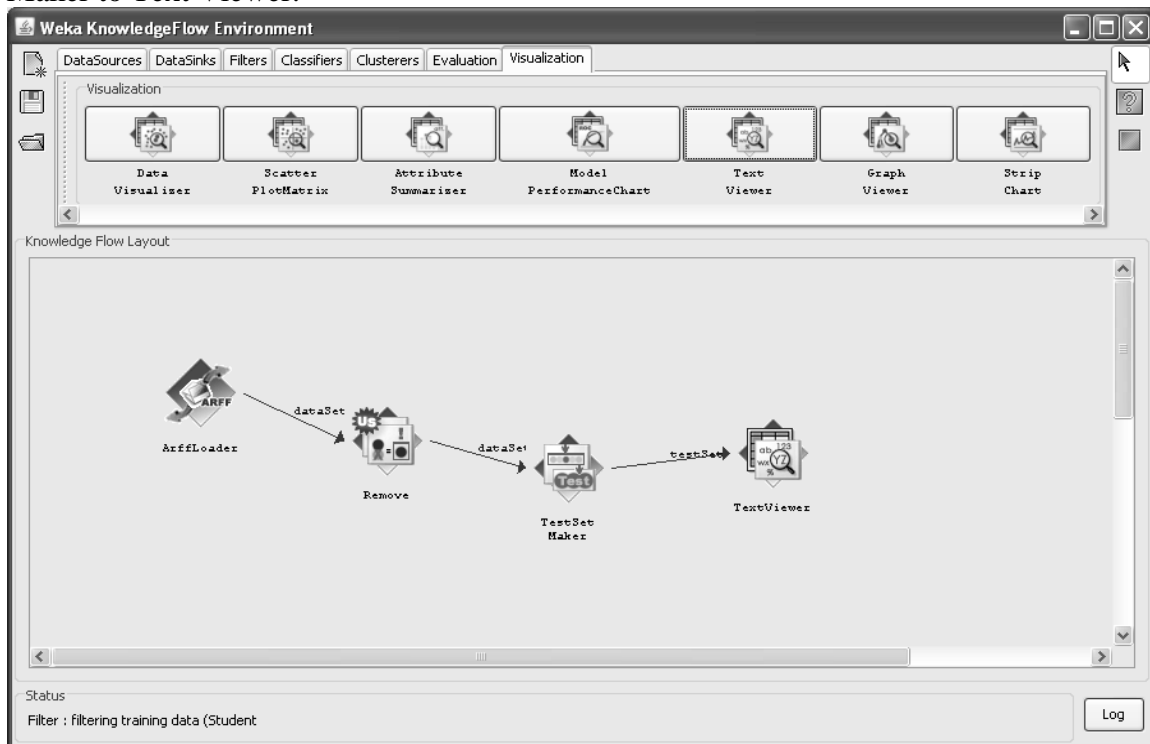
Step-9: - select Remove and right click on to choose Data Set to connection between to Remove and Test Set Maker.



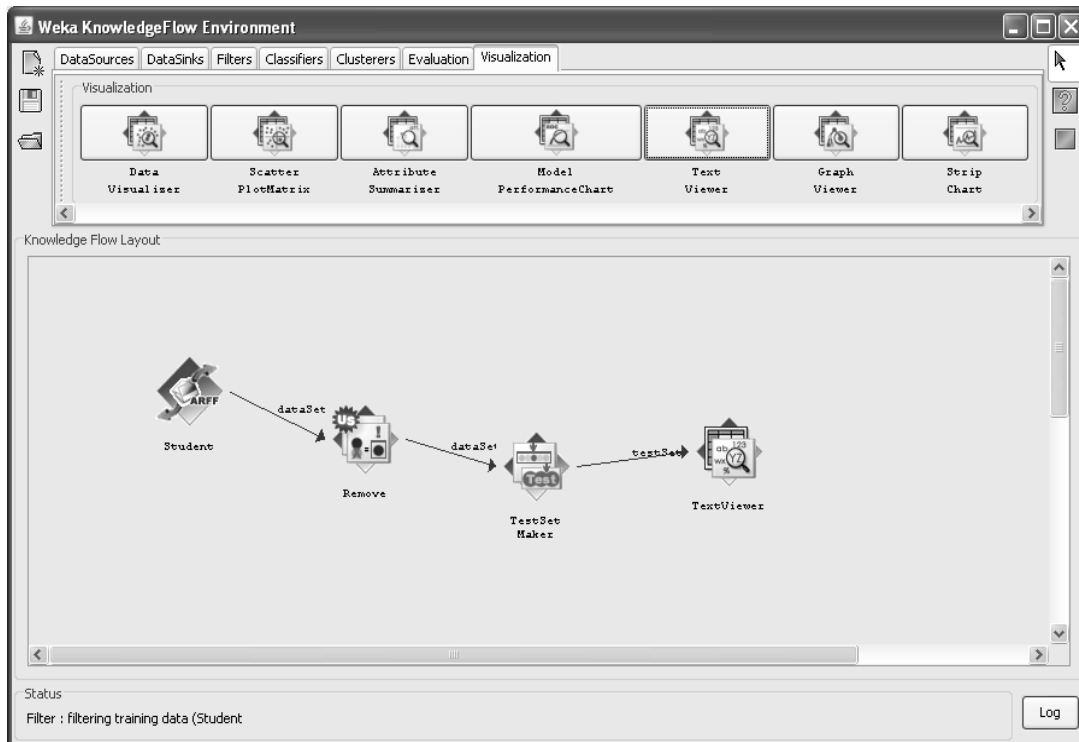
Step-10: -Go to Visualization and select Text Viewer drag on to knowledge flow layout.



Step-11: - select Test Set Maker and right click on choose Test Set to connection between Test Set Maker to Text Viewer.



Step-12: - Select ARFF Loader and right click on start loading.



Step-13: - select Text Viewer and right click on the show result.

```

@relation STUDENT-weka.filters.unsupervised.attribute.Remove-R1-2

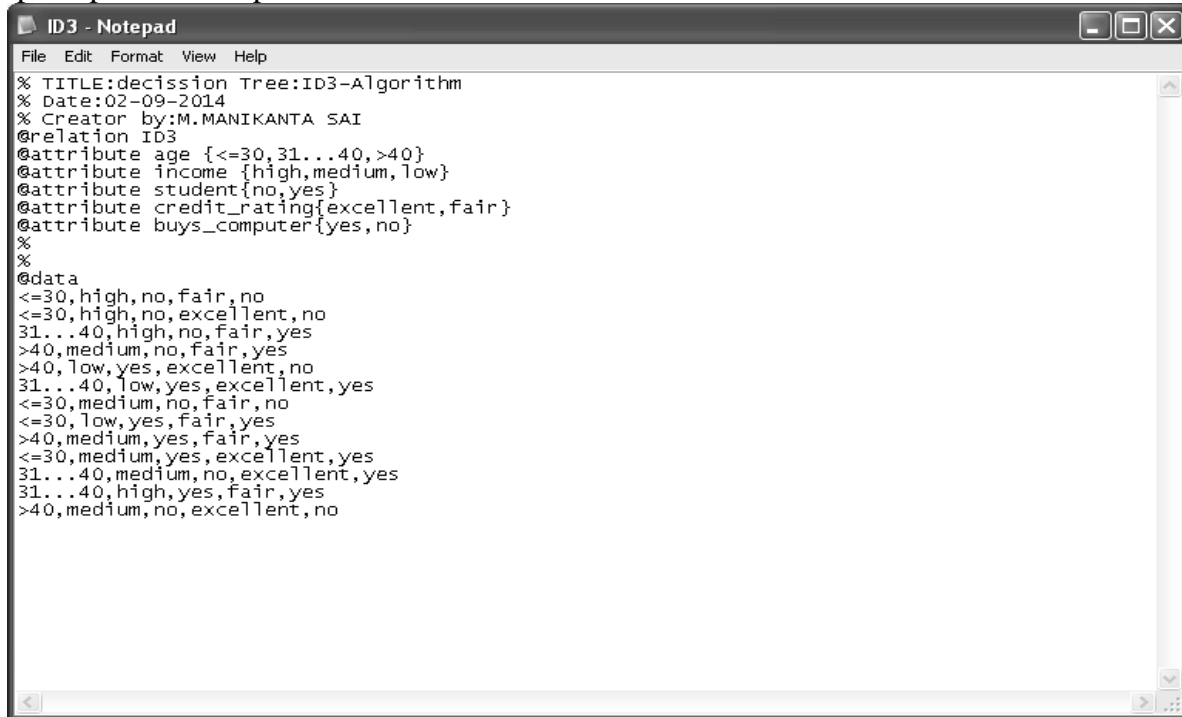
@attribute JAVA numeric
@attribute COBOL numeric
@attribute OS numeric
@attribute OSHEM numeric
@attribute OR numeric

@data
85,86,87,79,89
80,81,82,82,70
79,75,85,86,90
86,79,78,75,92
85,84,85,84,83
79,75,76,74,80
78,86,88,75,90
80,80,80,85,85
85,86,85,80,85
85,81,90,79,80
79,80,75,80,79
95,76,80,75,79
90,79,79,76,80
85,80,76,80,70
  
```

KNOWLEDGE FLOW FOR DECISION TREE

Knowledge flow for ID3 .

Step 1: open the Notepad and create the ARFF File.

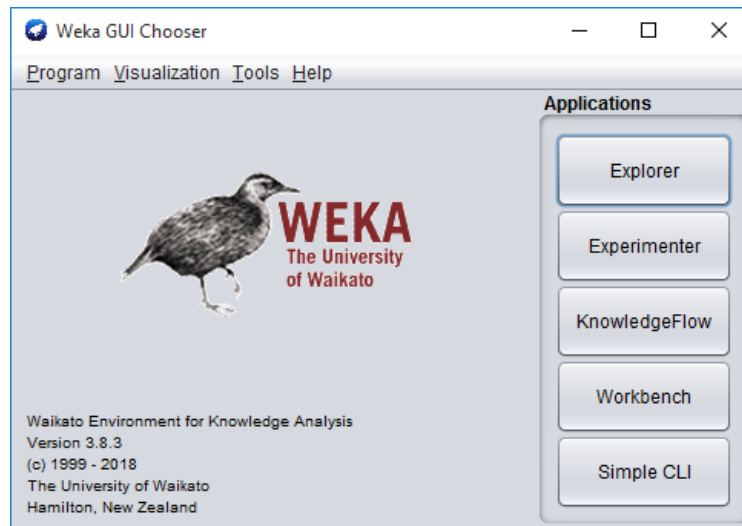


```

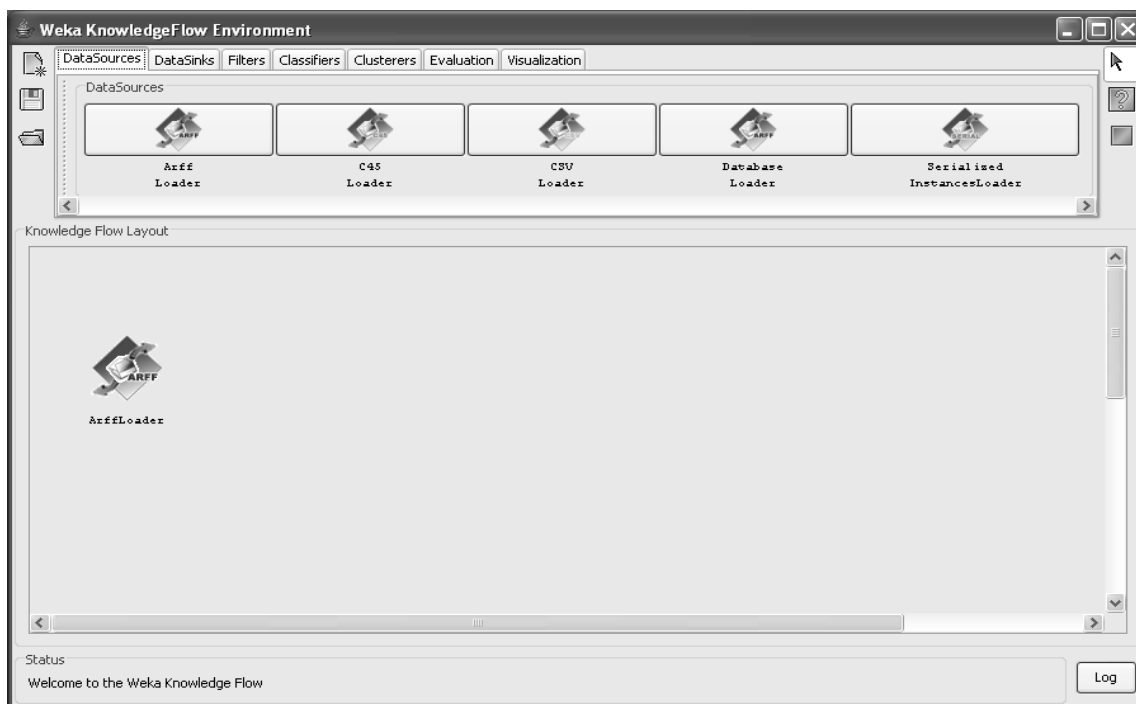
ID3 - Notepad
File Edit Format View Help
% TITLE:decision Tree:ID3-Algorithm
% Date:02-09-2014
% Creator by:M.MANIKANTA SAI
@relation ID3
@attribute age {<=30,31..40,>40}
@attribute income {high,medium,low}
@attribute student{no,yes}
@attribute credit_rating{excellent,fair}
@attribute buys_computer{yes,no}
%
%
@data
<=30,high,no,fair,no
<=30,high,no,excellent,no
31..40,high,no,fair,yes
>40,medium,no,fair,yes
>40,low,yes,excellent,no
31..40,low,yes,excellent,yes
<=30,medium,no,fair,no
<=30,low,yes,fair,yes
>40,medium,yes,fair,yes
<=30,medium,yes,excellent,yes
31..40,medium,no,excellent,yes
31..40,high,yes,fair,yes
>40,medium,no,excellent,no
  
```

Step2: To open All Programms → weka3.8.3 → weka3.8(with console) .

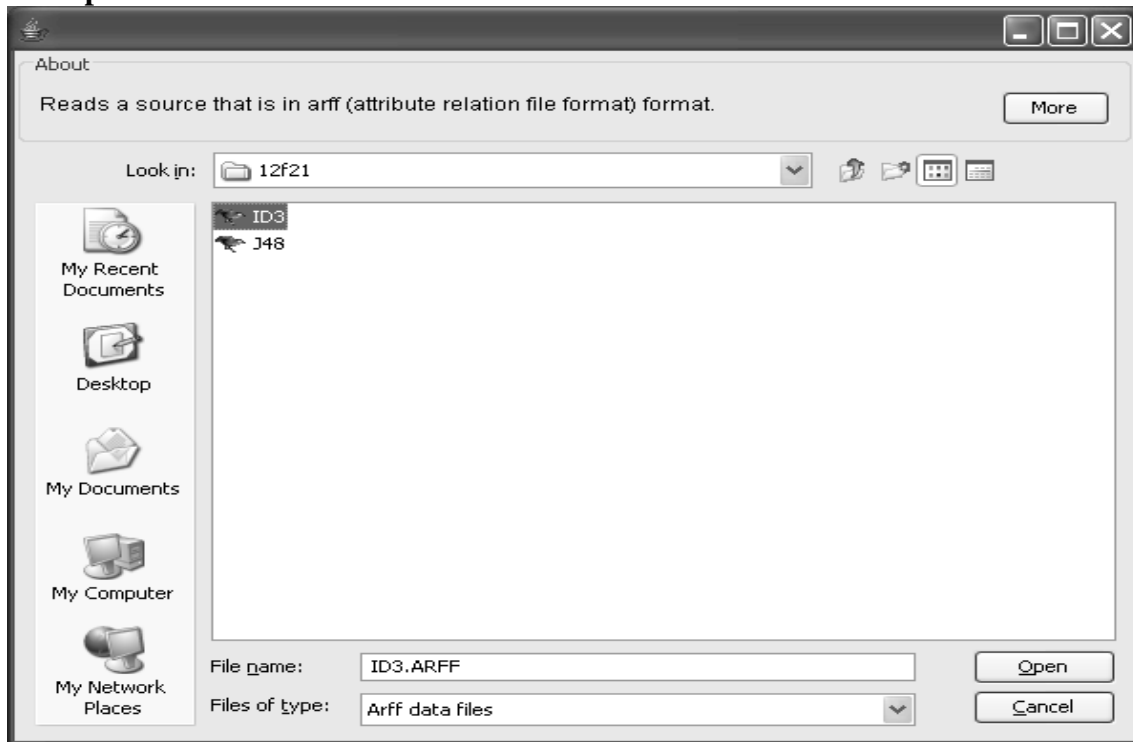
Step 3:click on the **Knowledge Flow**.



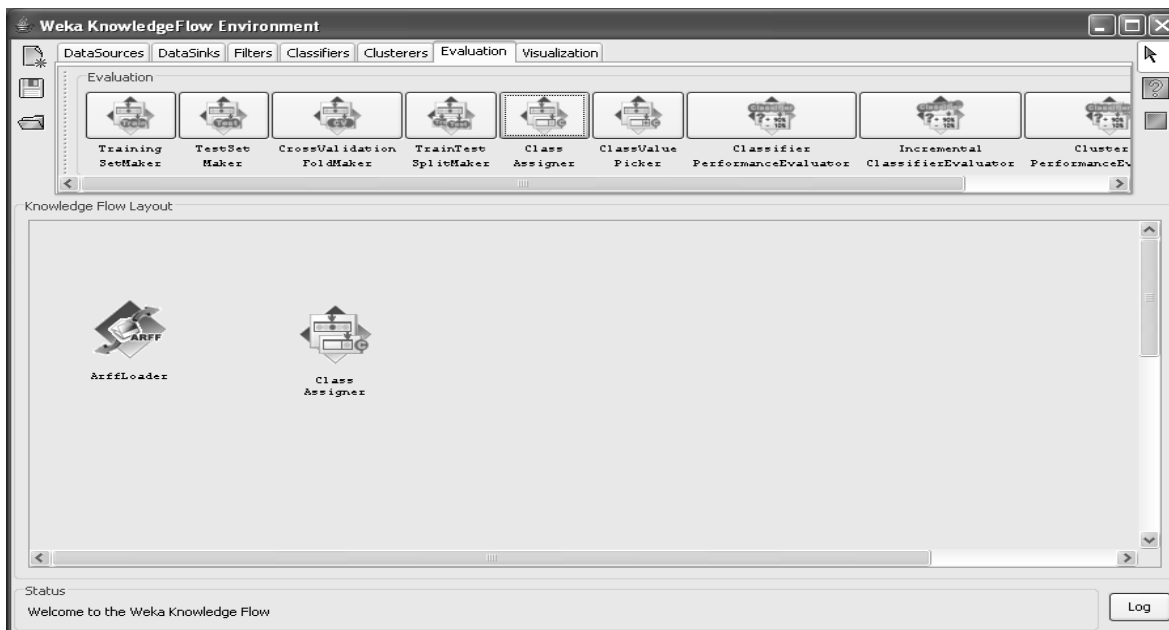
Step 4: click on the **Datasources** tab and drag the **Arff Loader** ,Right click over the Arff loader and select “configure” from the pop menu.



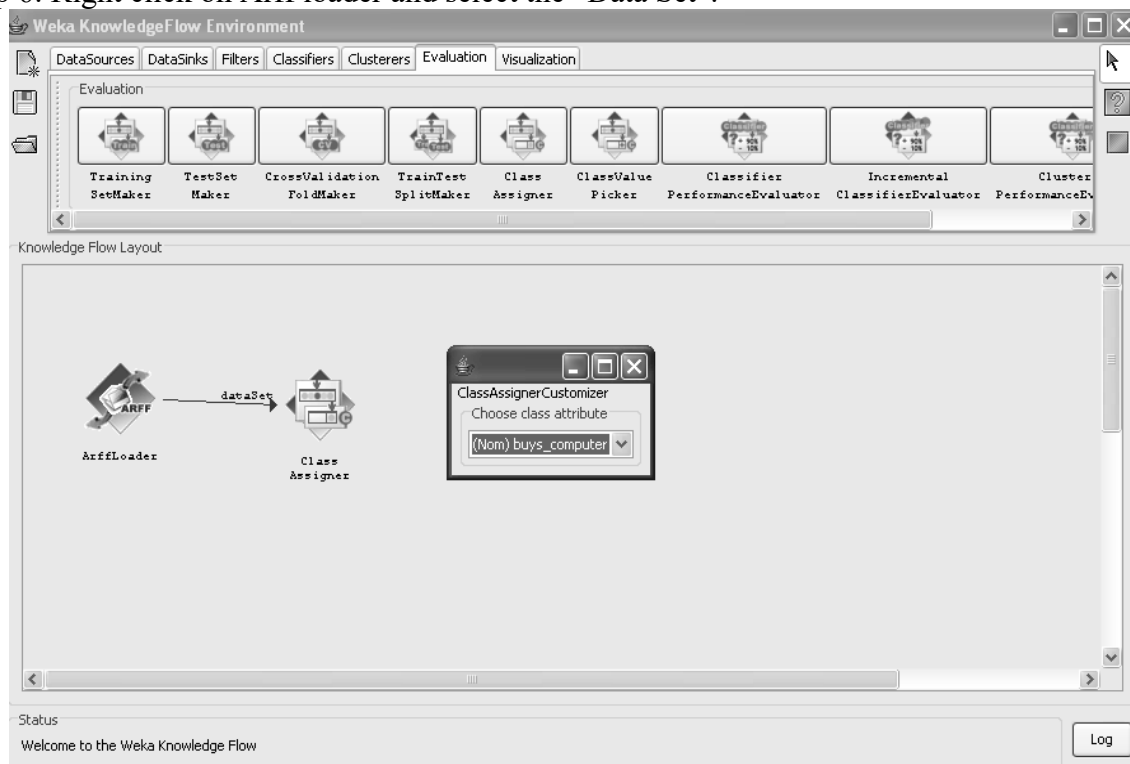
Click on **open**.



Step 5:click on the **Evaluation** tab at top of window,choose **class Assigner**,

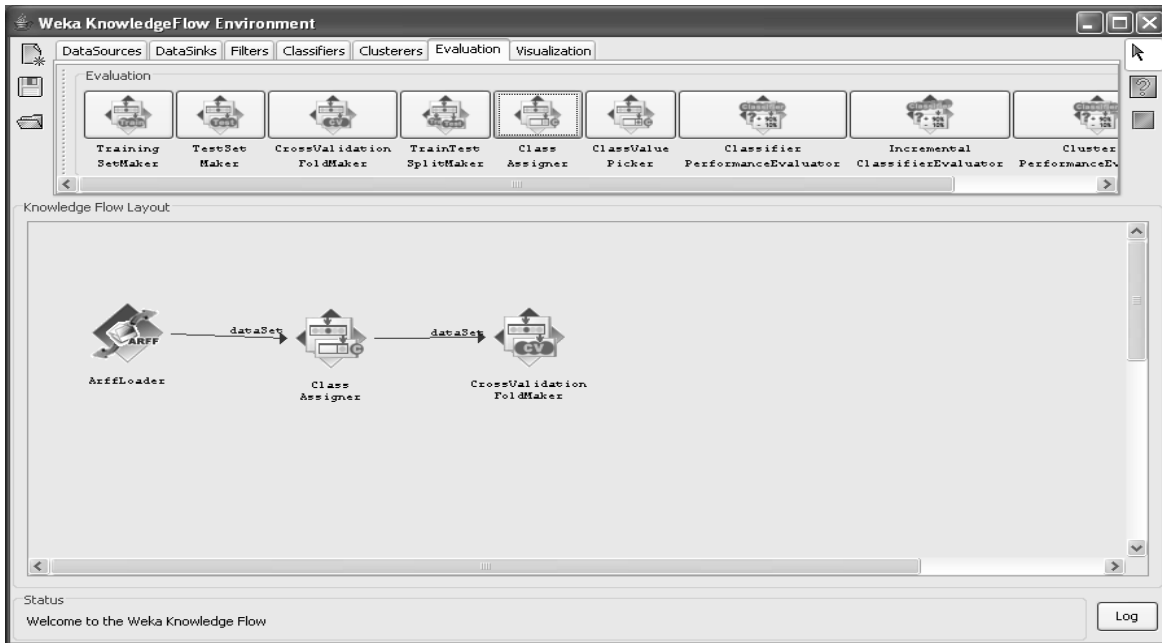


Step 6: Right click on Arff loader and select the “Data Set”.



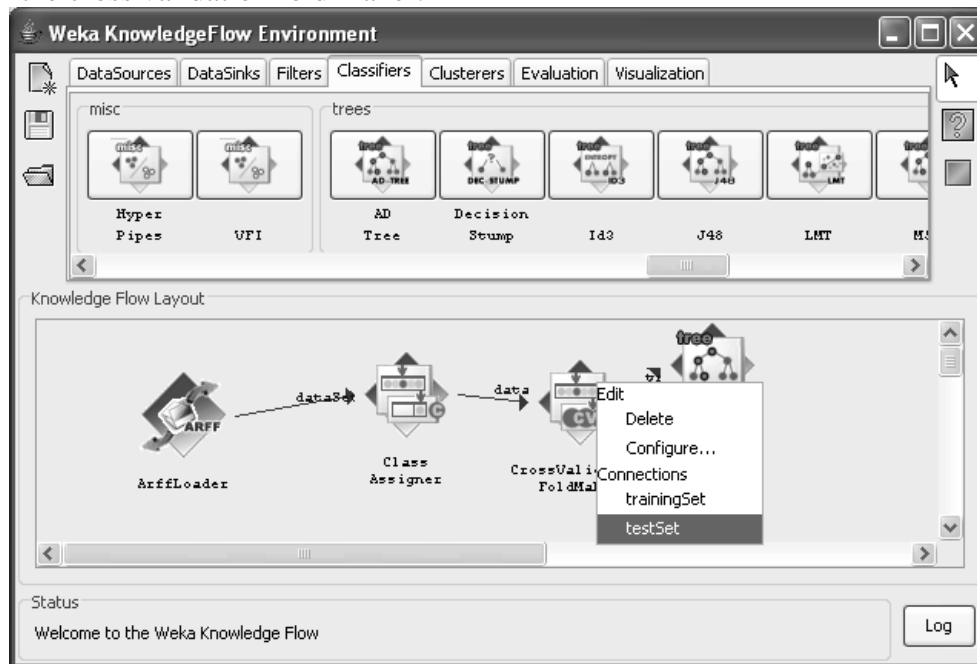
Step 7 : Choose cross validation fold maker, component from the evaluation toolbar

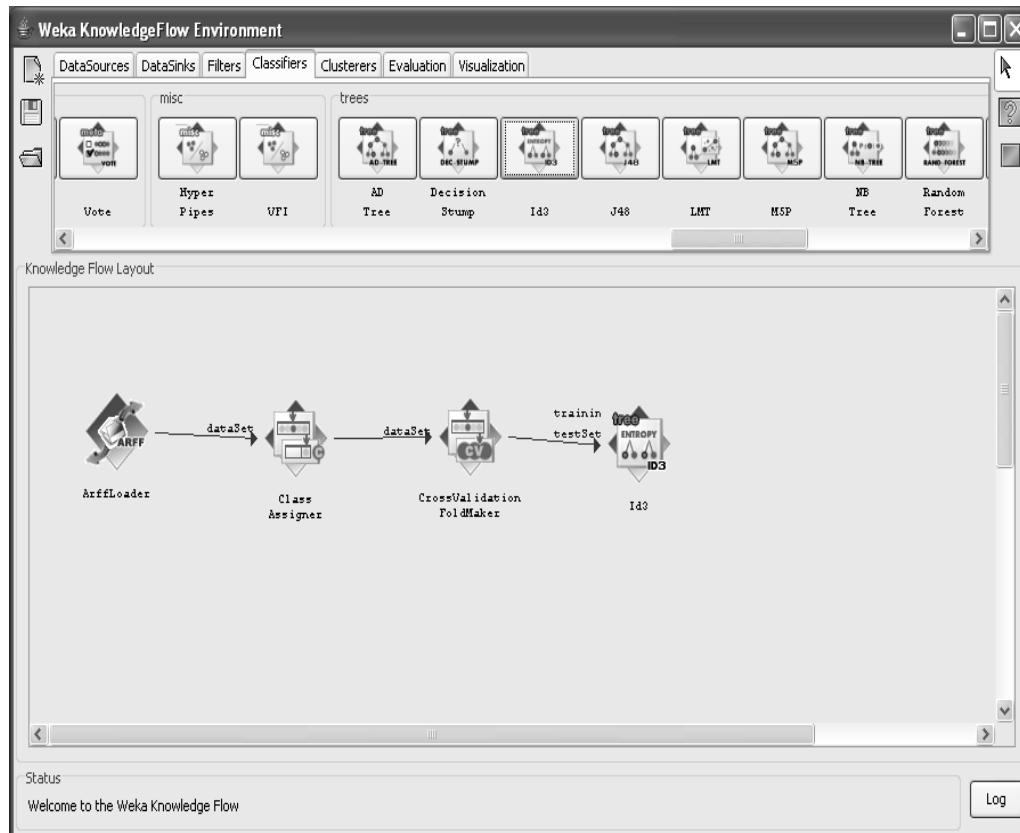
Step 8: connect class assigner to cross validation fold maker and right click over class assigner, select dataset under the connection menu.



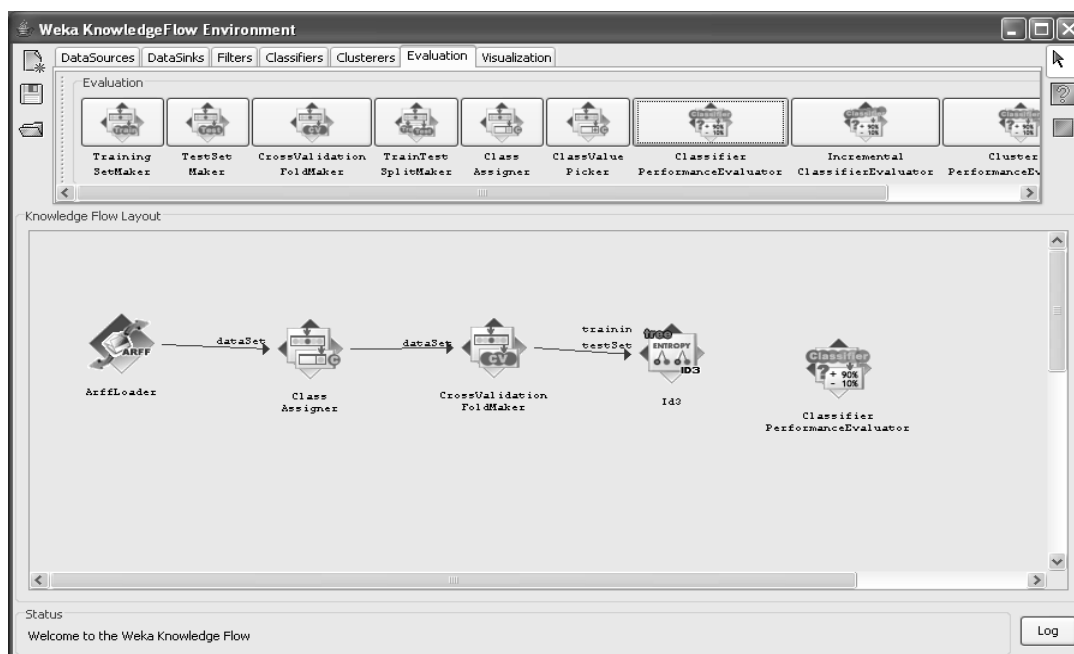
Step 9: click on classifiers tab on the window scroll the toolbar, specify ID3.

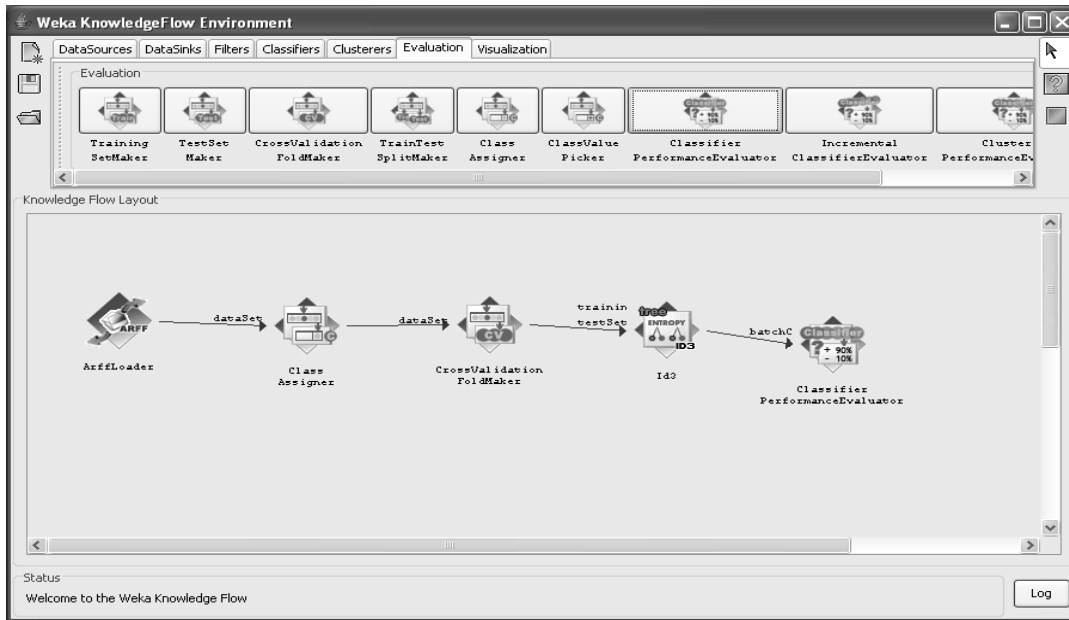
Step 10: connect the cross validation fold maker to ID3 by first choosing training set and then text set for the cross validation fold maker.





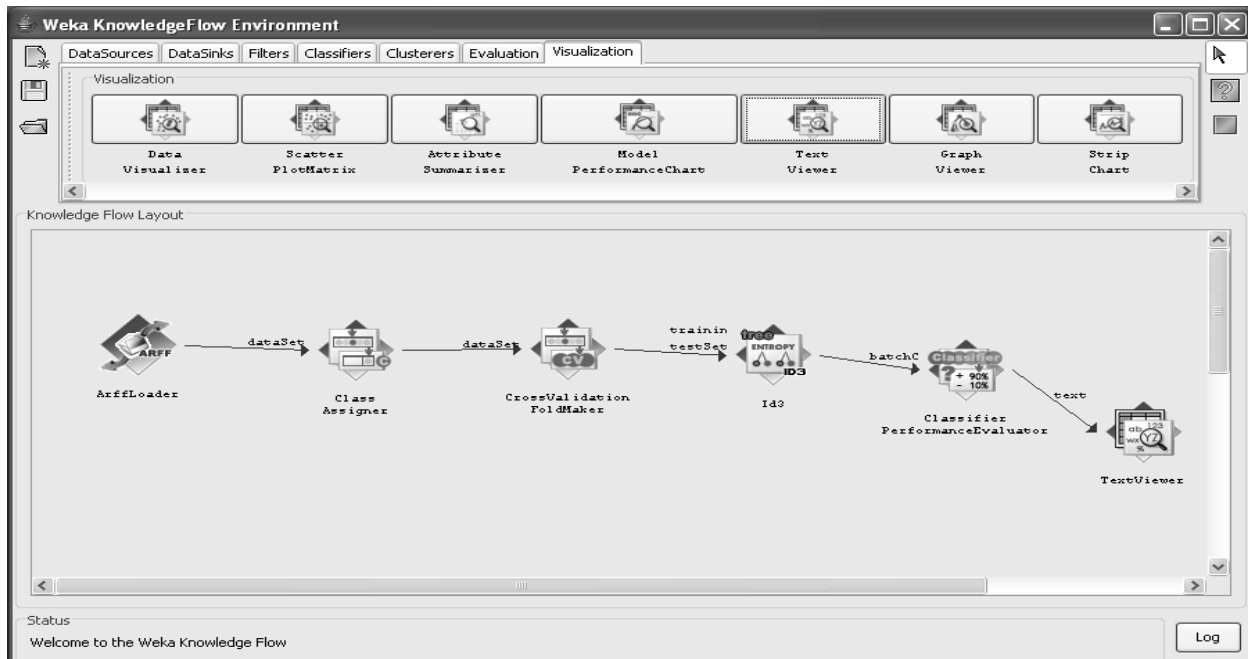
Step 11: Go to evaluation tab and place a classifier performance evaluator component on the layout Connect ID3 to this component by selecting the batch classifier entry from the popup menu for ID3.





Step13: Go to visualization toolbar and place text viewer component on the layout. connect the classifier performance evaluator to the text viewer by selecting the text entry form the popup for classifier performance evaluator.

step14: Now start the flow executing by selecting start loading from the popup menu for arffloader we also see some progress information in the status bar and log at bottom of the window.



step15: when finished, you can view the results by choosing show results from the popup menu for the test viewer component.

Output:

```

Text Viewer
Result list  Text
10:40:23 - Id3

=== Evaluation result ===

Scheme: Id3
Relation: ID3

Correctly Classified Instances      10          76.9231 %
Incorrectly Classified Instances    3           23.0769 %
Kappa statistic                    0.4935
Mean absolute error                 0.2308
Root mean squared error            0.4804
Relative absolute error             46.1538 %
Root relative squared error        96.0769 %
Total Number of Instances          13

=== Detailed Accuracy By Class ===

TP Rate  FP Rate  Precision  Recall  F-Measure  Class
0.875    0.4      0.778     0.875   0.824      yes
0.6      0.125    0.75      0.6     0.667      no

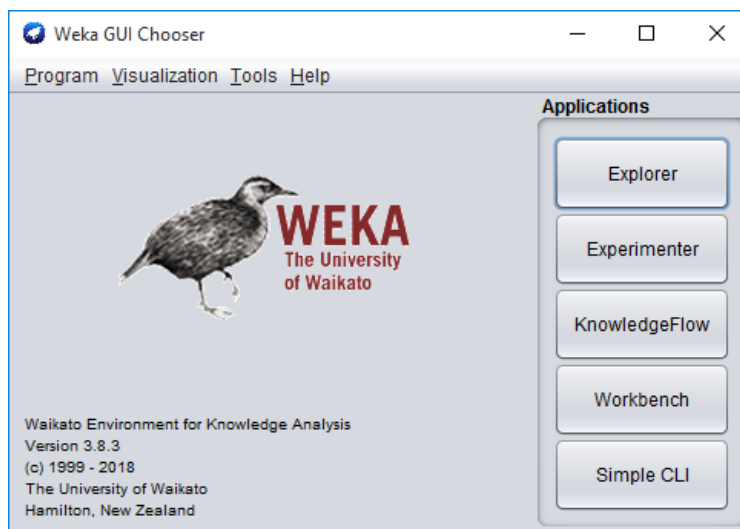
=== Confusion Matrix ===

a b  <-- classified as
7 1 | a = yes
2 3 | b = no

```

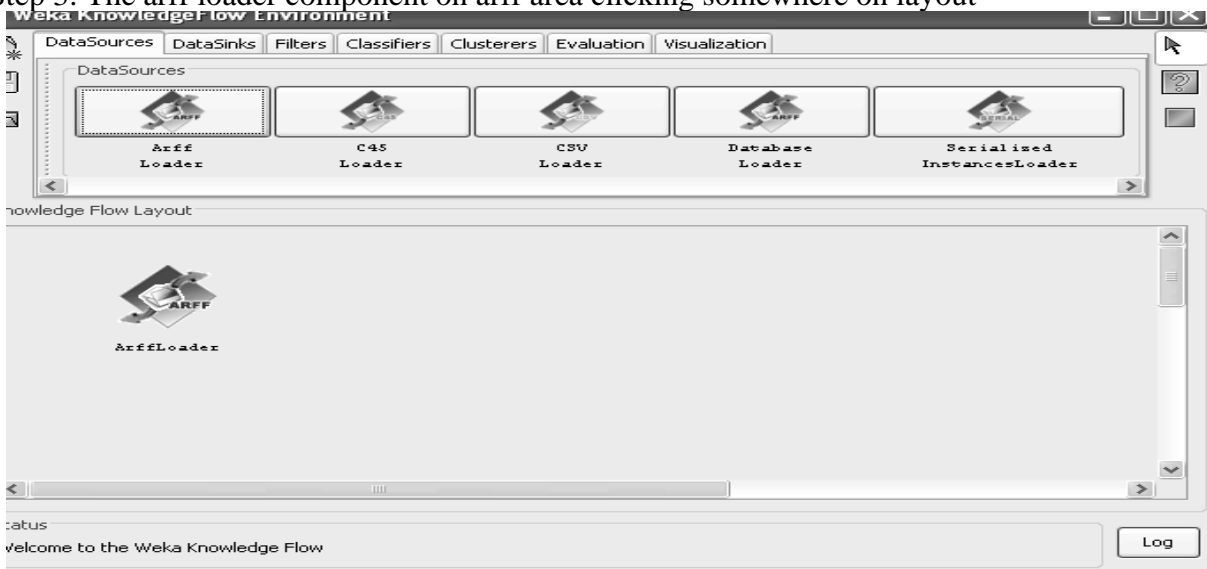
Knowledge flow for decision tree J48.

Step 1: Select knowledge flow in Weka GUI chooser and click on it.

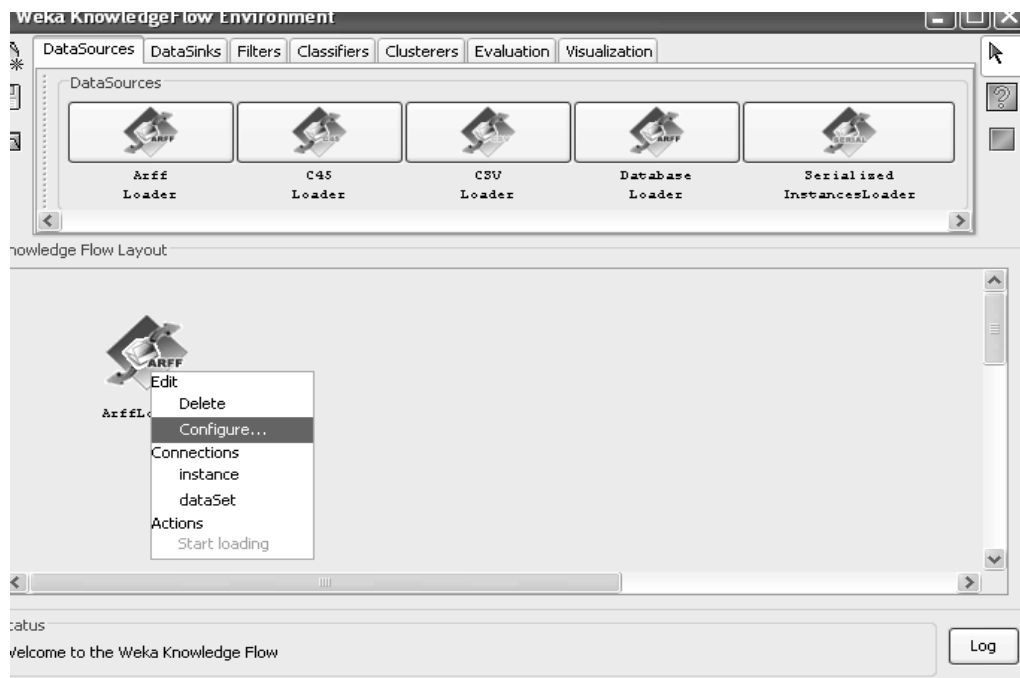


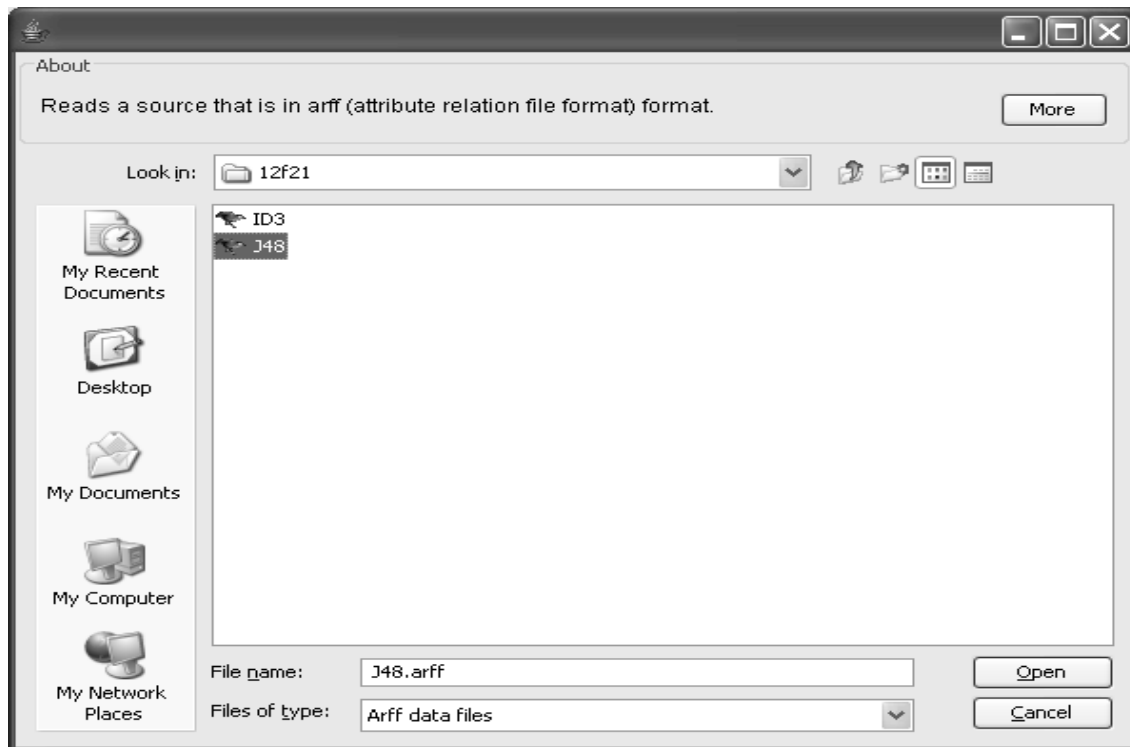
Step 2: click on data source tab and choose arff loader from the tool bar

Step 3: The arff loader component on arff area clicking somewhere on layout



Step 4: specify j48.arff file to load by right clicking mouse over the arffloader icon on the layout. A popup menu will appear select **configure**.

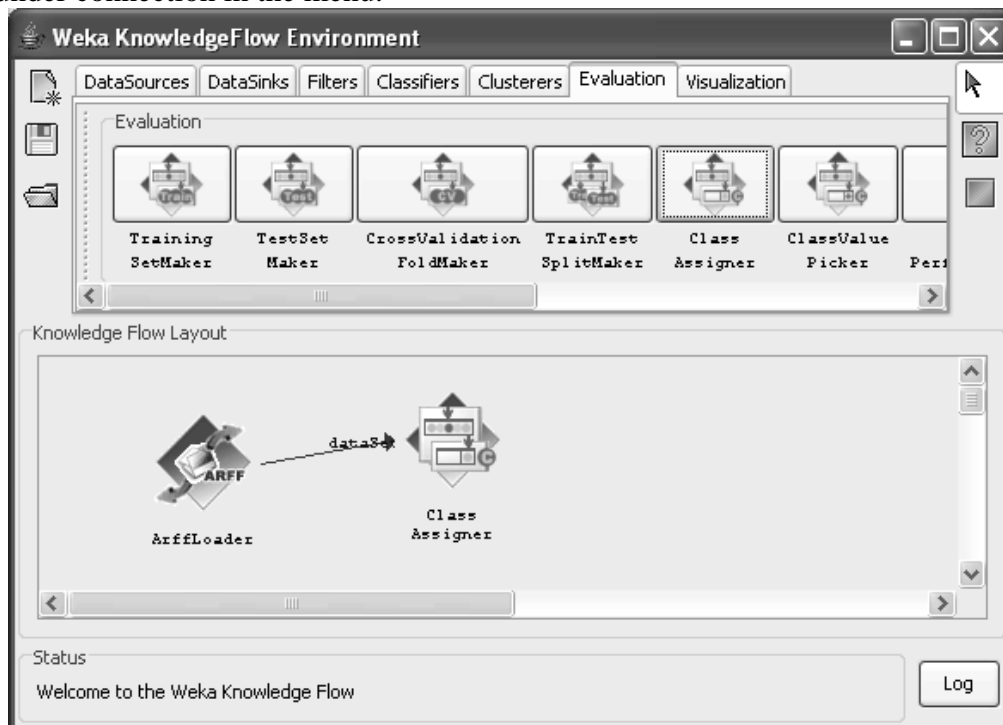




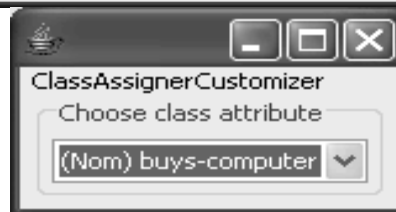
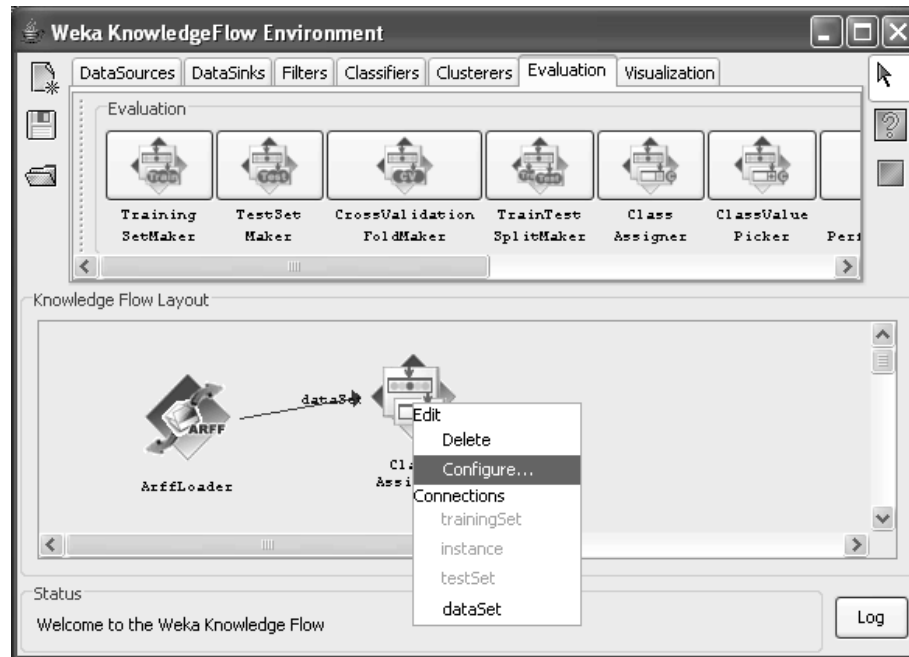
Click on open button.

Step 5: Click on evaluation tab, at the top of window, choose **Class Assigner**.

Step6: Now connect arff loader to Class Assigner , first right click over Arffloader and select **dataset** under connection in the menu.

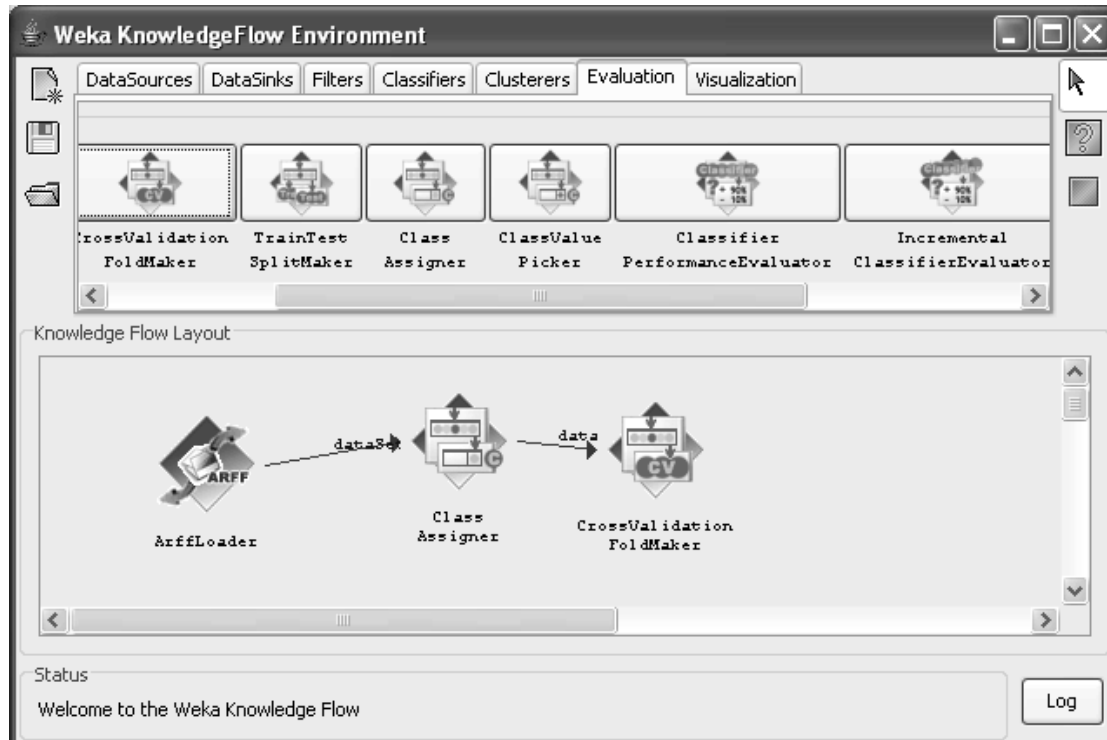


Step7:click on class assigner configure menu, choose specify columns is the class in our data given as(buys-computer).



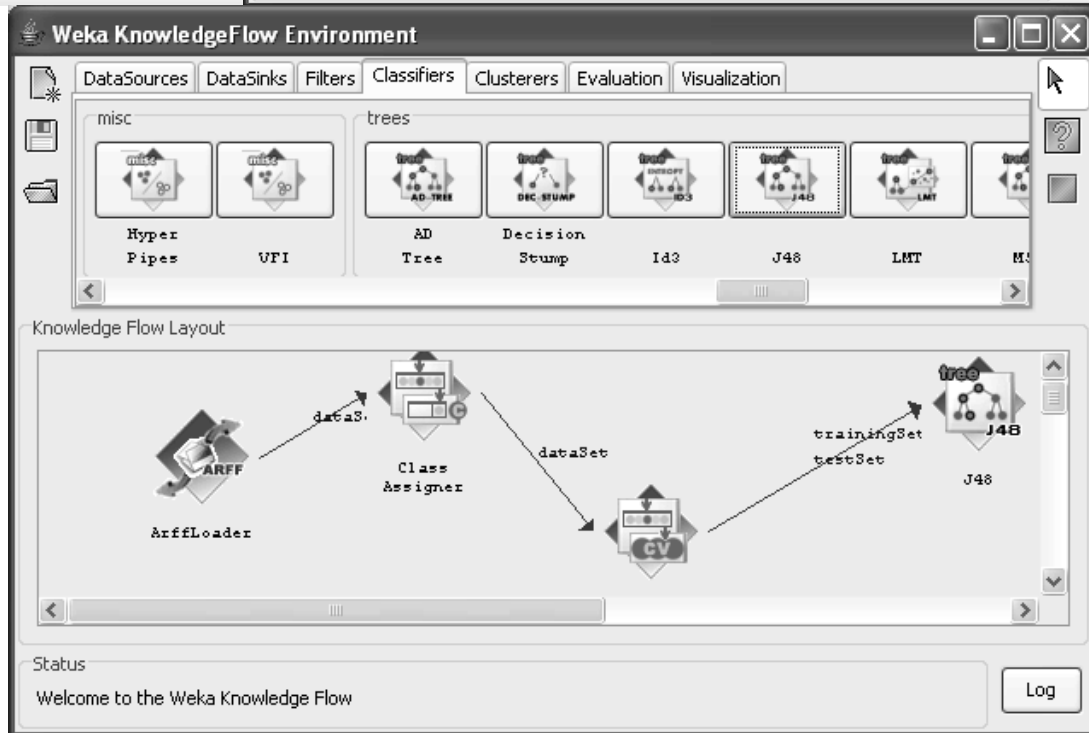
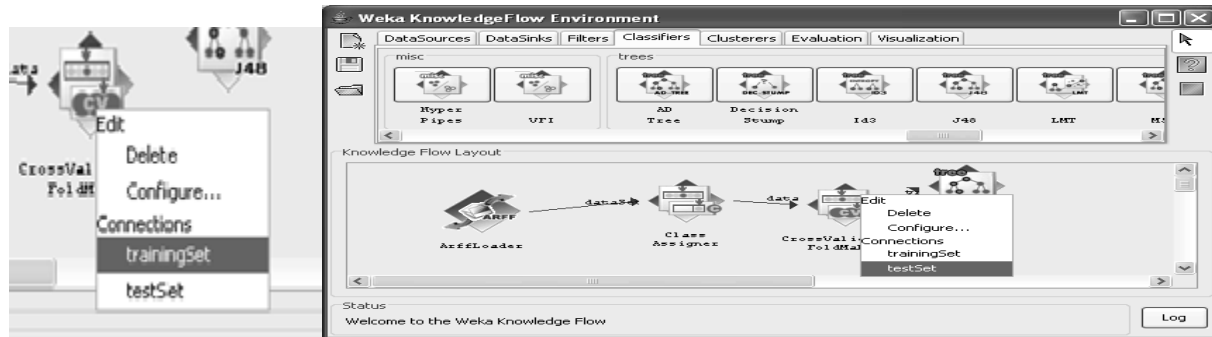
Step 8: Choose **cross validation fold maker**, component from the evaluation toolbar

Step 9: connect class assigner to cross validation fold maker and right click over class assigner, select dataset under the connection menu.

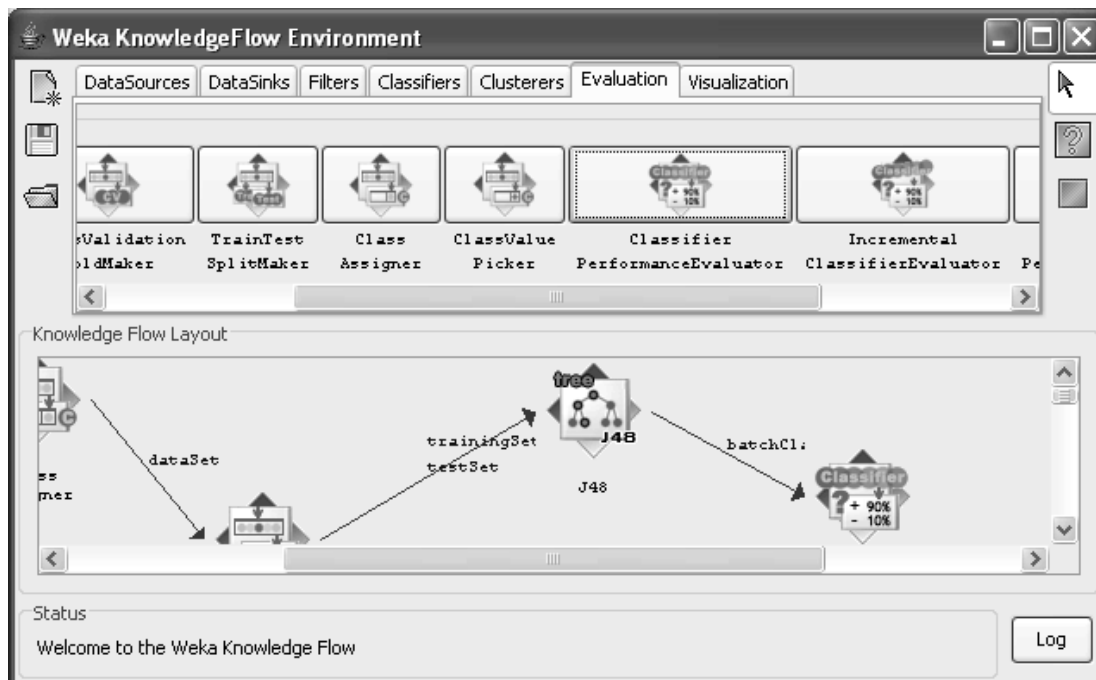
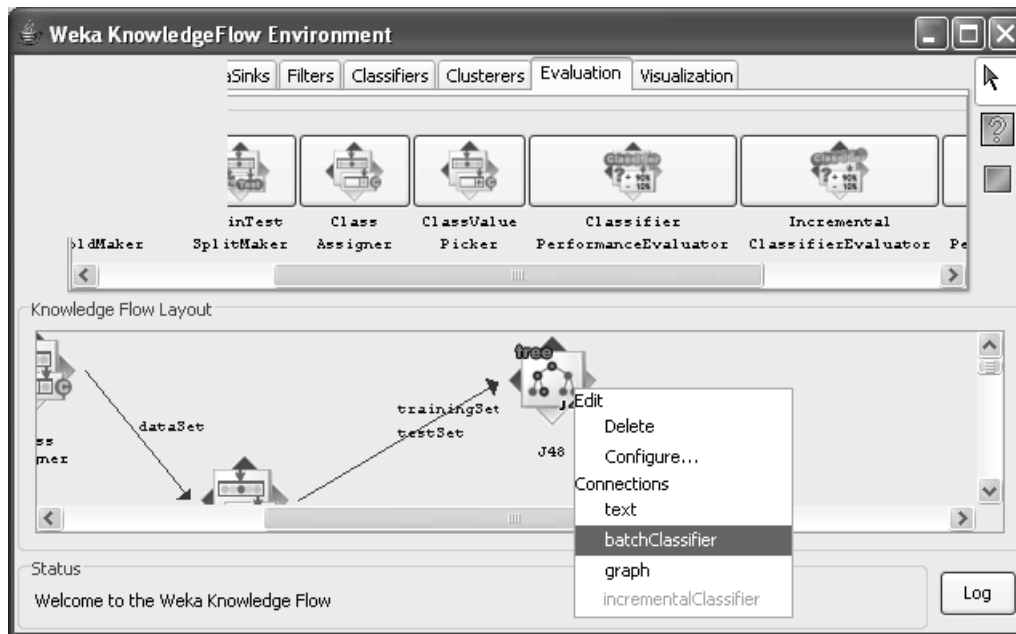


Step 10: click on classifiers tab on the window scroll to the toolbar, specify **j48**

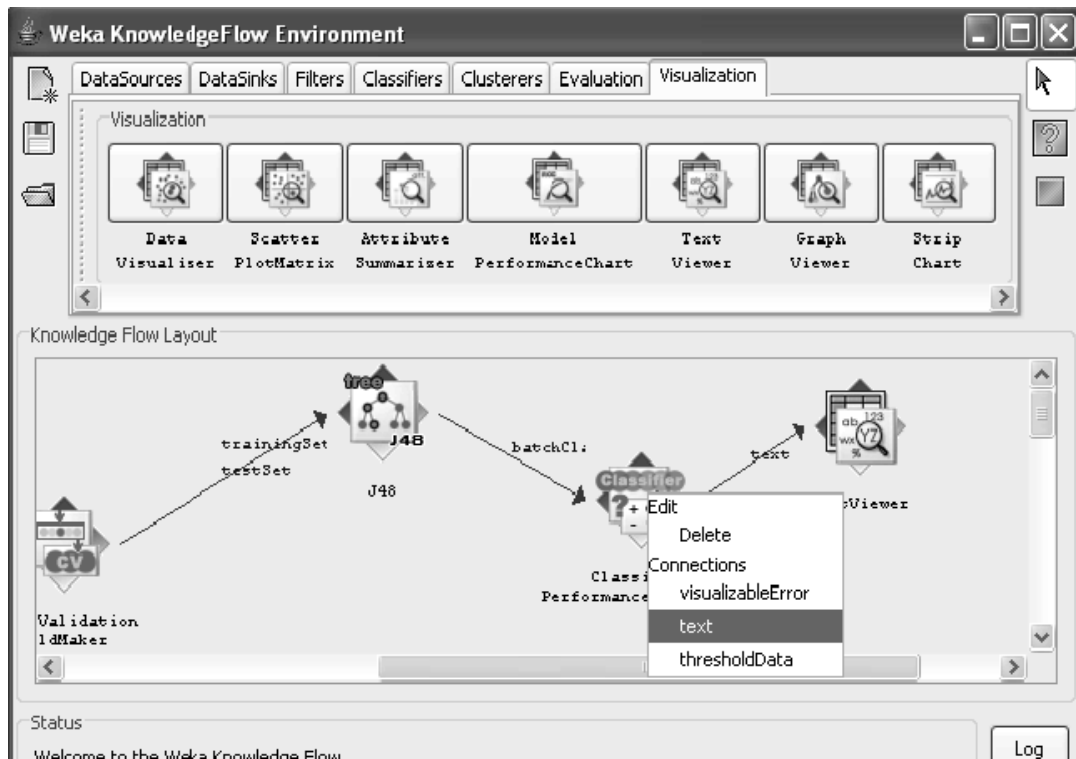
Step 11: connect the cross validation fold maker to j48 by first choosing **training set** and then **test set** for the cross validation fold maker



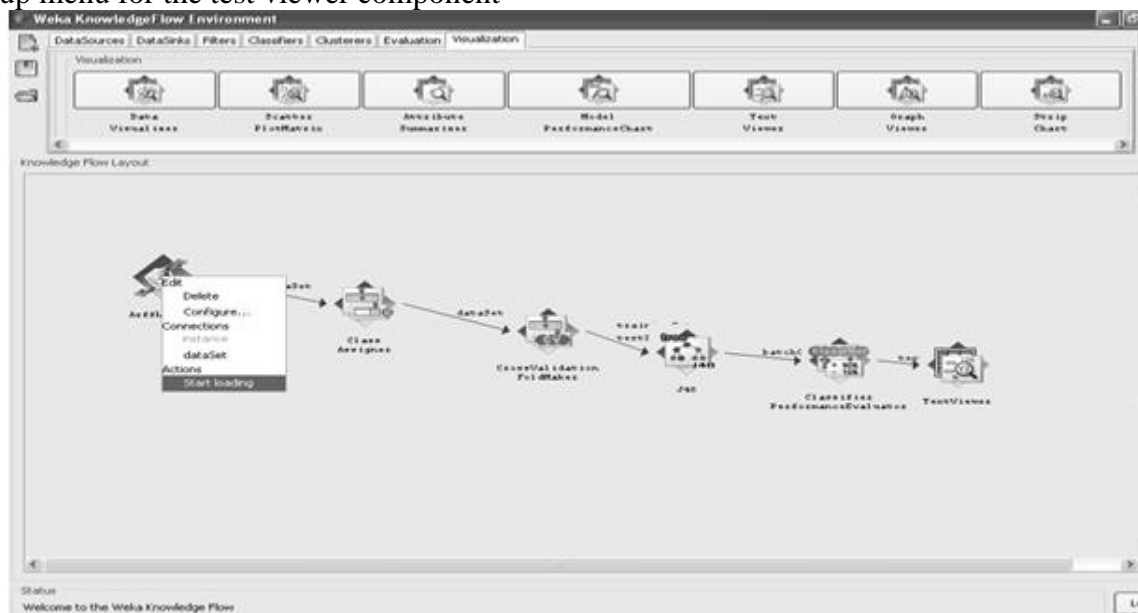
Step 12: Go to evaluation tab and place a classifier performance evaluator component on the layout Connect j48 to this component by selecting the **batch classifier** entry from the popup menu for j48



Step 13: Go to visualization toolbar and place **text viewer** component on the layout. connect the classifier performance evaluator to the text viewer by selecting the text entry form the popup for classifier performance evaluator



step14: Now start the flow executing by selecting **start loading** from the popup menu of arffloader. step15: when finished, you can view the results by choosing show results from the popup menu for the test viewer component



OUTPUT:

```

--- Evaluation result ---

Scheme: J48
Relation: j48

Correctly Classified Instances      6      46.1538 %
Incorrectly Classified Instances    7      53.8462 %
Kappa statistic                    -0.0964
Mean absolute error                 0.4167
Root mean squared error             0.5461
Relative absolute error             83.3333 %
Root relative squared error        109.2276 %
Total Number of Instances          13

=== Detailed Accuracy By Class ===

TP Rate  FP Rate  Precision  Recall  F-Measure  Class
0.5      0.6      0.571     0.5     0.533     yes
0.4      0.5      0.333     0.4     0.364     no

=== Confusion Matrix ===

a b  <-- classified as
4 4 | a = yes
3 2 | b = no

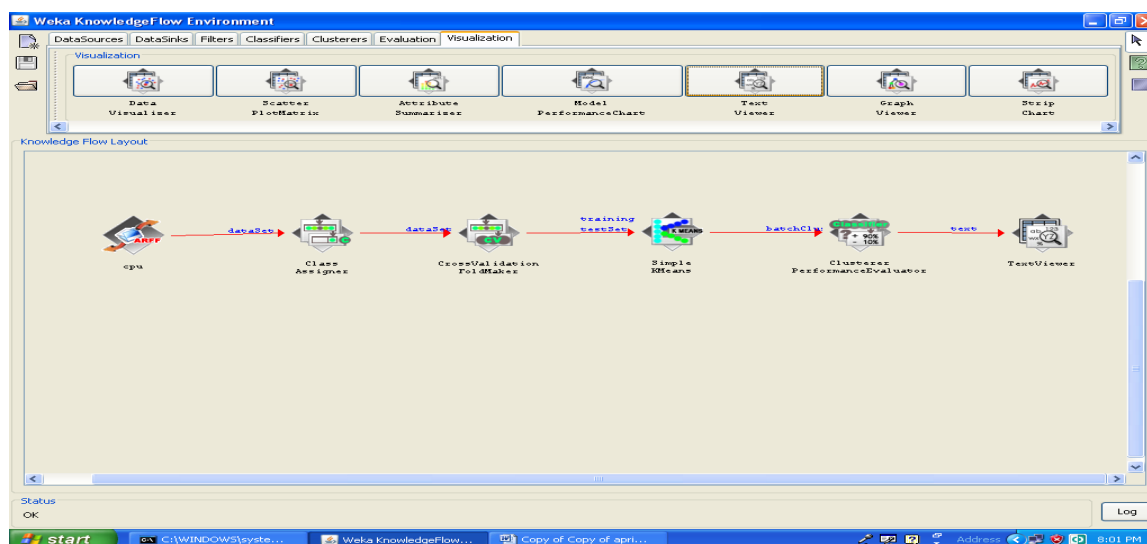
```

KNOWLEDGE FLOW FOR CLUSTERING

- First start the Knowledge Flow.
- Next click on the Data Sources tab and choose "ArffLoader" from the toolbar (the mouse pointer will change to a "cross hairs").
- Next place the ArffLoader component on the layout area by clicking somewhere on the layout (A copy of the ArffLoader icon will appear on the layout area).
- Next specify an arff file to load by first right clicking the mouse over the ArffLoader icon on the layout. A pop-up menu will appear. Select "Configure" under "Edit" in the list from this menu and browse to the location of your arff file.
- Next click the "Evaluation" tab at the top of the window and choose the "Class Assigner" (allows you to choose which column to be the class) component from the toolbar. Place this on the layout.
- Now connect the ArffLoader to the Class Assigner: first right click over the ArffLoader and select the "data Set" under "Connections" in the menu. A "rubber band" line will appear. Move the mouse over the Class Assigner component and left click - a red line labeled "data Set" will connect the two components.

- Next right click over the Class Assigner and choose "Configure" from the menu. This will pop up a window from which you can specify which column is the class in your data (last is the default).
- Next grab a "Cross Validation Fold Maker" component from the Evaluation toolbar and place it on the layout. Connect the Class Assigner to the Cross Validation Fold Maker by right clicking over "Class Assigner" and selecting "dataset" from under "Connections" in the menu.
- Next click the "Clusterer" tab the top of the window and choose the "SimpleKMeans" component from the toolbar. Place this on the layout.
- Connect the CrossValidationFoldMaker to SimpleKMeans by first choosing "training Set" and then "test Set" from the pop-up menu for the CrossValidationFoldMaker
- Next go back to the "Evaluation" tab and place a "ClustererPerformanceEvaluator" component on the layout. Connect SimpleKMeans to this component by selecting the "batch Clusterer" entry from the pop-up menu for simpleKMeans.
- Next go to the "Visualization" toolbar and place a "Text Viewer" component on the layout. Connect the ClustererPerformanceEvaluator to the Text Viewer by selecting the "text" entry from the pop-up menu for ClustererPerformanceEvaluator.
- Now start the flow executing by selecting "Start loading" from the pop-up menu for ArffLoader.
- When finished you can view the results by choosing show results from the pop-up menu for the TextViewer component.

Output:



```

Text
==== Evaluation result for training instances ====

Scheme: SimpleKMeans
Relation: cpu

KMeans
=====

Number of iterations: 17
Within cluster sum of squared errors: 17.681682670930194

Cluster centroids:

Cluster 0
  Mean/Mode: 256.7075 1549.2789 7129.6327 10.1088 2.3537 11.9184
  Std Devs: 285.9676 1349.555 4969.8984 16.307 2.2446 17.1098
Cluster 1
  Mean/Mode: 49.1707 7862.3415 28345.8537 80.5854 11.9024 40.2927
  Std Devs: 34.008 6109.2962 14976.0028 57.47 9.1373 38.7132

Clustered Instances

Unclustered instances : 188
KMeans
=====

Number of iterations: 4
Within cluster sum of squared errors: 19.50617730173854

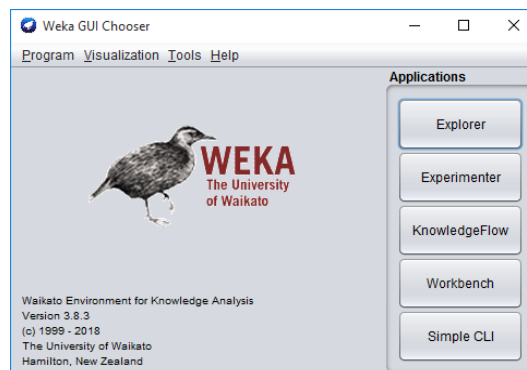
Cluster centroids:

Cluster 0
  Mean/Mode: 121.8471 3159.9294 12690.6118 27.4824 5 19.4765
  Std Devs: 95.3727 4131.1216 11539.21 41.1982 6.3925 24.085
Cluster 1
  Mean/Mode: 855.5556 653.3333 4584 2.8889 0.9444 2.1667
  Std Devs: 202.491 448.0924 3048.8873 7.7375 0.6391 1.7573

```

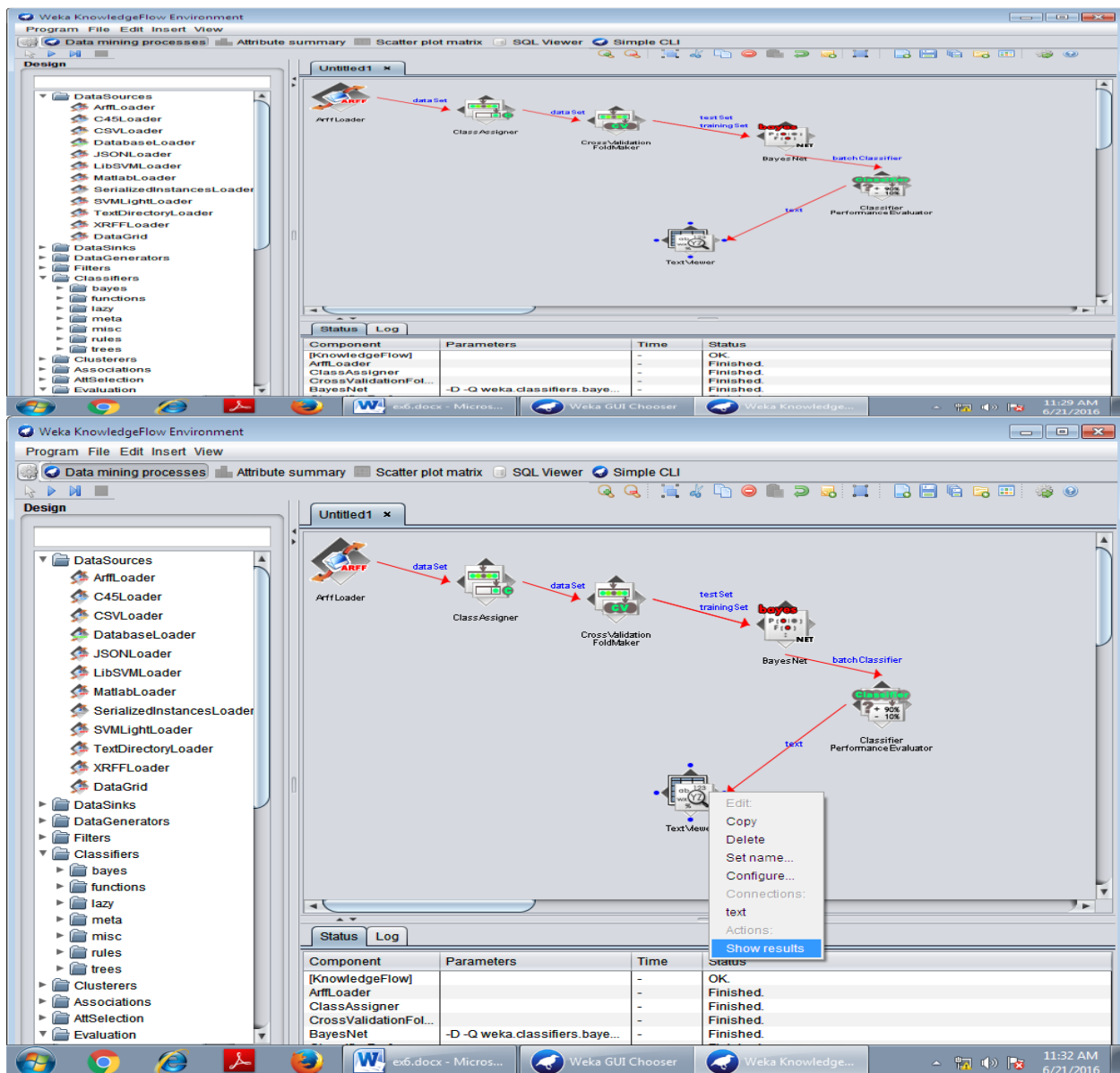
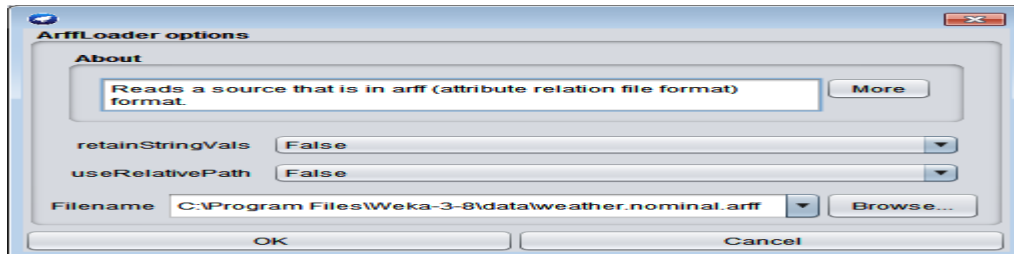
Naïve Bayes

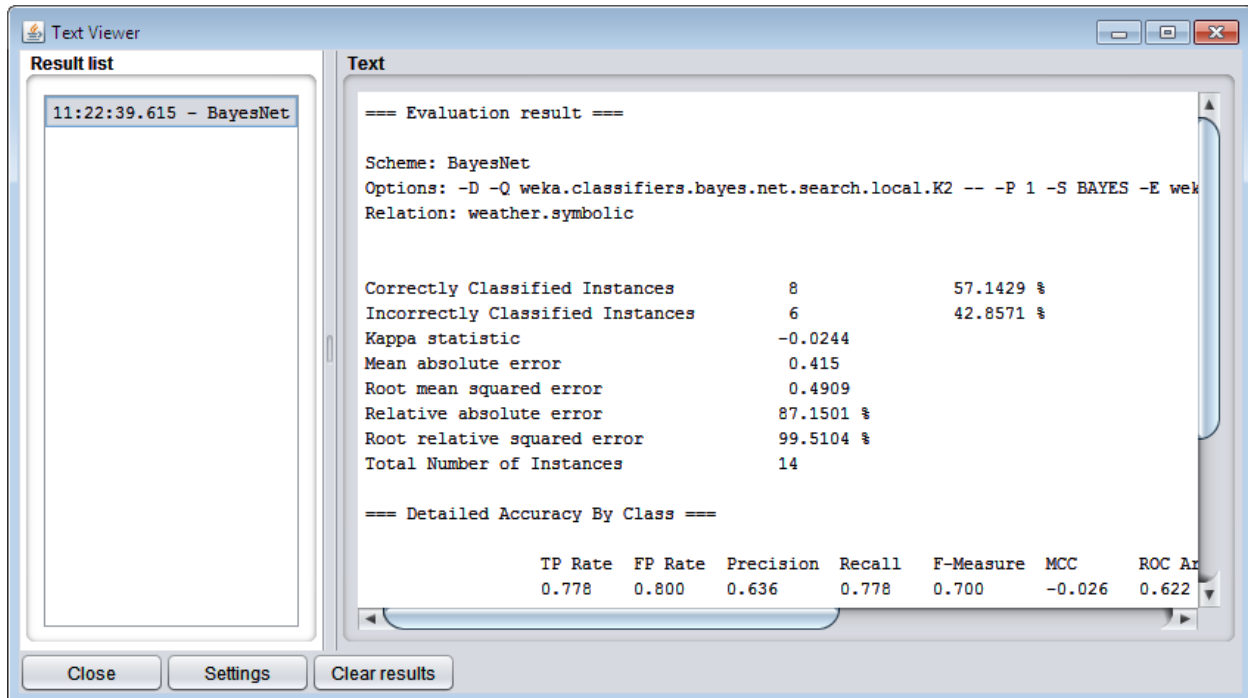
Click **Knowledge Flow** in **Weka GUI Chooser**



To create the following and click **Run** button and **Right click** the **TextViwer** and select **Show Result**.

Draw **ArrfLoader** and select the filename.





VIVA-QUESTIONS

- 1 ____ is the graphical representation of information and data.
- 2 Which types of chart shows the relationship between a numerical variable and categorical variable?
 - A. Line
 - B. Bar
 - C. Pie
 - D. x-y plot
- 3 Which of the following are Data sources of Knowledge Flow?
 - A. CSV
 - B. Notepad
 - C. Excel
 - D. Access
- 4 Data Visualization in mining cannot be done using
- 5 A bubble chart is a variation of
- 6 Which method shows hierarchical data in a nested format?
- 7 Data visualization is also an element of the broader _____.
- 8 Which of the following is a Trees in data visualization
 - A. B-Tree
 - B. B+-Tree
 - C. Heap-Tree
 - D. AD-Tree
- 9 which of the following not support data source
 - A. ARFF Loader
 - B. CSV Loader
 - C. Access Loader
 - D. XML Loader
- 10 what is use of class assigner in weka

ADDITIONAL EXPERIMENTS

FILE FORMATES FOR WEKA

1. Create CSV(Comma Separated Values) file.

Step1: Create an excel file and save with specified format as CSV(Comma Delimited).

Step2: Now open with notepad and check the values. Here, the fields of data in each row are delimited with a comma and individual rows are separated by new line.

2. Create arff(Attribute Relation File Format) file.

Step1: Open a notepad and type the data as instructed below:

ARFF files have two distinct sections. The first section is the **Header** information, which is followed the **Data** information.

The **Header** of the ARFF file contains the name of the relation, a list of the attributes (the columns in the data), and their types. An example header on the standard IRIS dataset looks like this:

```
% 1. Title: Iris Plants Database
%
% 2. Sources:
% (a) Creator: R.A. Fisher
% (b) Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)
% (c) Date: July, 1988
%
@RELATION iris
@ATTRIBUTE sepallength NUMERIC
@ATTRIBUTE sepalwidth NUMERIC
@ATTRIBUTE petallength NUMERIC
@ATTRIBUTE petalwidth NUMERIC
@ATTRIBUTE class {Iris-setosa,Iris-versicolor,Iris-virginica}
```

The **Data** of the ARFF file looks like the following:

```
@DATA
5.1,3.5,1.4,0.2,Iris-setosa
4.9,3.0,1.4,0.2,Iris-setosa
4.7,3.2,1.3,0.2,Iris-setosa
4.6,3.1,1.5,0.2,Iris-setosa
```

5.0,3.6,1.4,0.2,Iris-setosa

5.4,3.9,1.7,0.4,Iris-setosa

4.6,3.4,1.4,0.3,Iris-setosa

5.0,3.4,1.5,0.2,Iris-setosa

4.4,2.9,1.4,0.2,Iris-setosa

4.9,3.1,1.5,0.1,Iris-setosa

Lines that begin with a % are comments. The **@RELATION**, **@ATTRIBUTE** and **@DATA** declarations are case insensitive.

Step 2: Save the file as .arff.

Step 3: Open with Weka Explorer and check the file values.

3. Convert CSV to ARFF file format.

Step 1: Open CSV file with notepad

Step 2: To fill header and data section in CSV file.

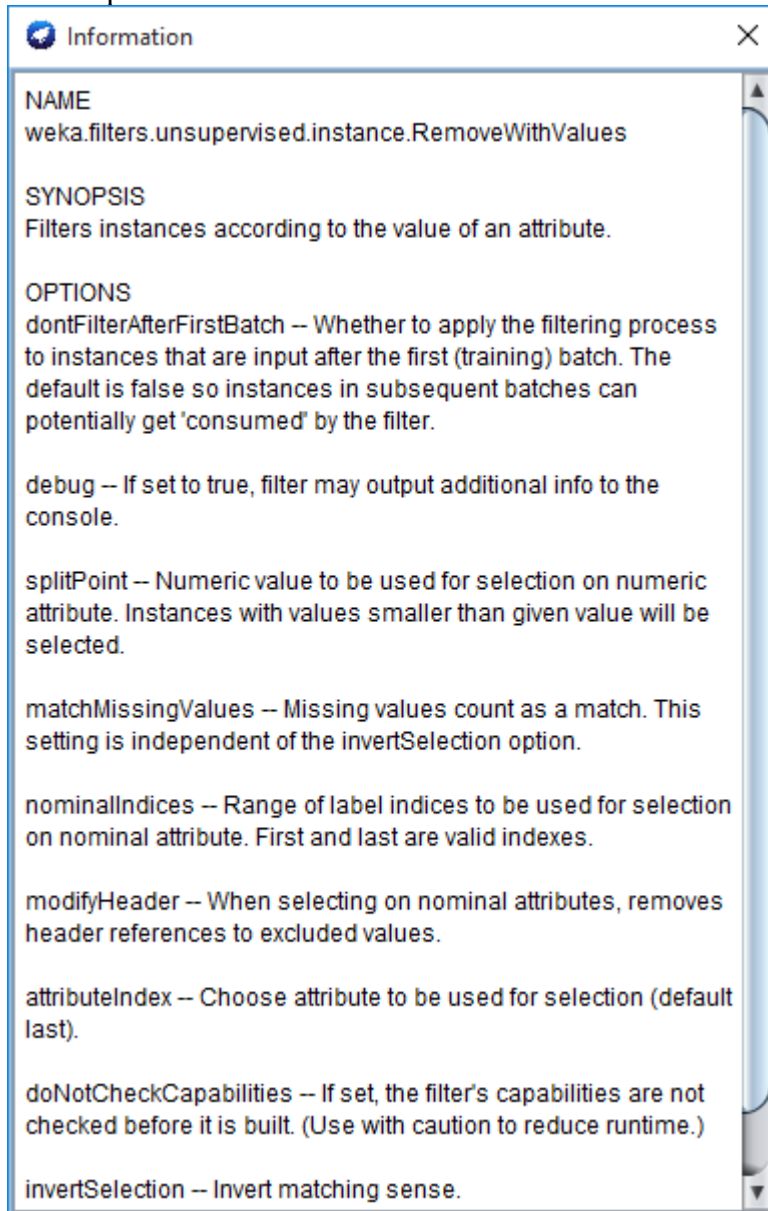
Step 3: save the file type as arff.

Step 4: Open with Weka Explorer.

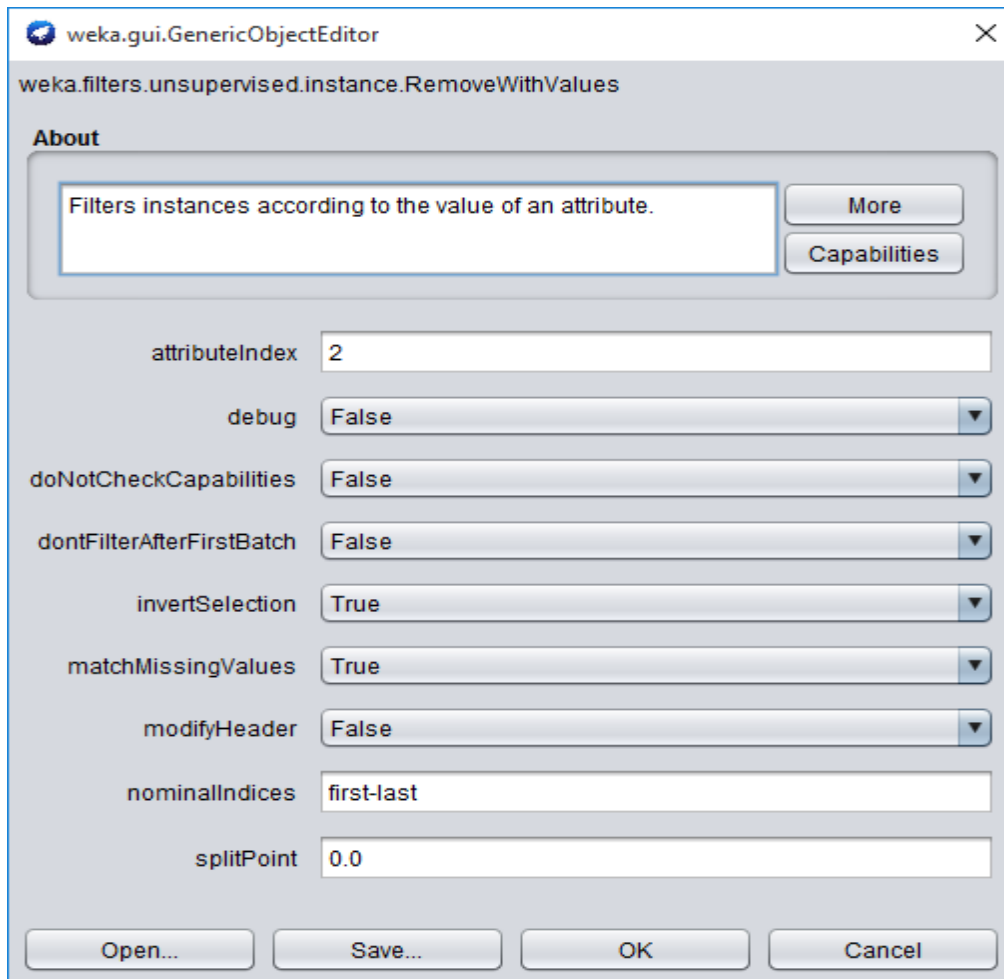
Missing Values

Ignore the tuple (Remove Missing Data)

1. Open the Weka Explorer.
2. Load the Student Data set
3. Click the “Choose” button for the filter and select “Remove with values” as under Unsupervised instance.Remove with Values.



4. Set Match Missing value to “True”



5. Click the “OK” button to use the configuration for the filter.

6. Click the “Apply” button to apply the filter click .

Viewer

Relation: report1-weka.filters.unsupervised.instance.RemoveWithValues-S0.0-C2-Lfirst-last-V-weka.filters.unsupervised.instance.RemoveWithValues-S0.0-C2-Lfirst-last-V-M

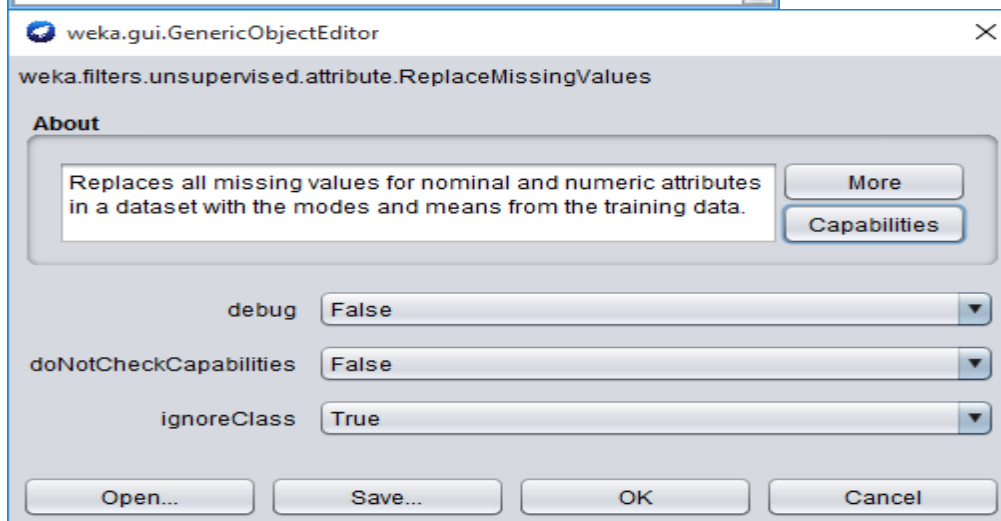
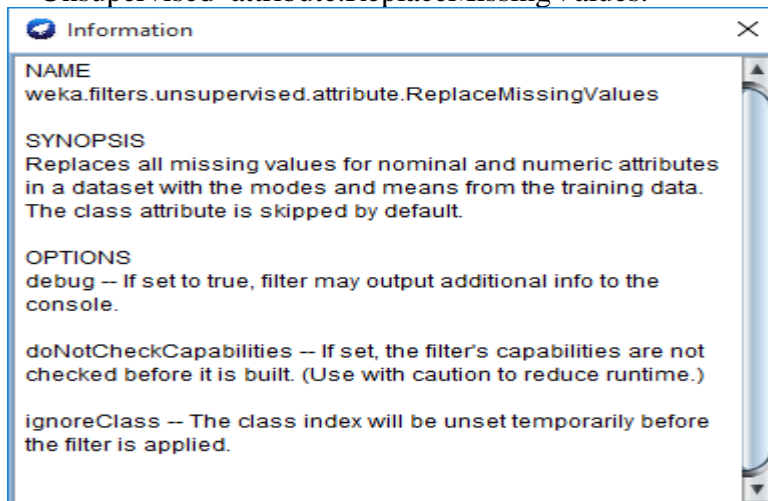
No.	1: SNO	2: SNAME	3: M1	4: M2	5: M3
	Numeric	Nominal	Numeric	Numeric	Numeric
1	1.0	balu	10.0	20.0	30.0
2	2.0	siya	100.0	40.0	50.0
3	3.0	raja	40.0	40.0	60.0

Add instance Undo OK Cancel

Impute Missing Values

How to Impute mean values for missing values

1. Open the Weka Explorer.
2. Load the Student Data set
3. Click the “Choose” button for the filter and select “ReplaceMissingValues” as under Unsupervised attribute.ReplaceMissingValues.



5. Click the “OK” button to use the configuration for the filter.
6. Click the “Apply” button to apply the filter click .

Viewer

Results: report-weka.filters.unsupervised.attribute.ReplaceMissingValues-unset-class-temporarily

No	1: SNO	2: SNAME	3: M1	4: M2	5: M3
	Nominal	Nominal	Numeric	Numeric	Numeric
1	1.0	DATA	10.0	20.0	30.0
2	2.0	DATA	100.0	40.0	50.0
3	3.0	DATA	40.0	40.0	50.0
4	7.0	DATA	40.0	38.5	46.25
5	4.0	DATA	34.0	38.5	46.25
6	6.0	DATA	48.0	54.0	46.25
7	5.0	DATA	45.0	38.5	46.0

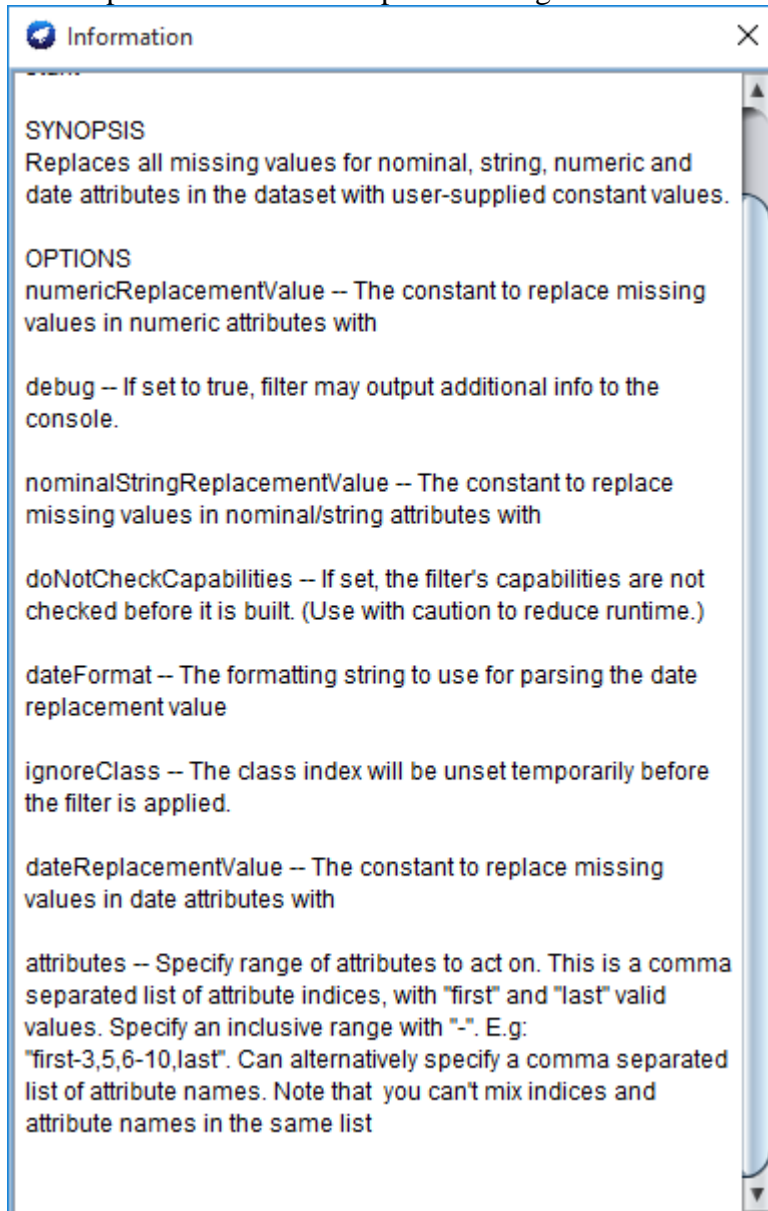
Add Instance Undo OK Cancel

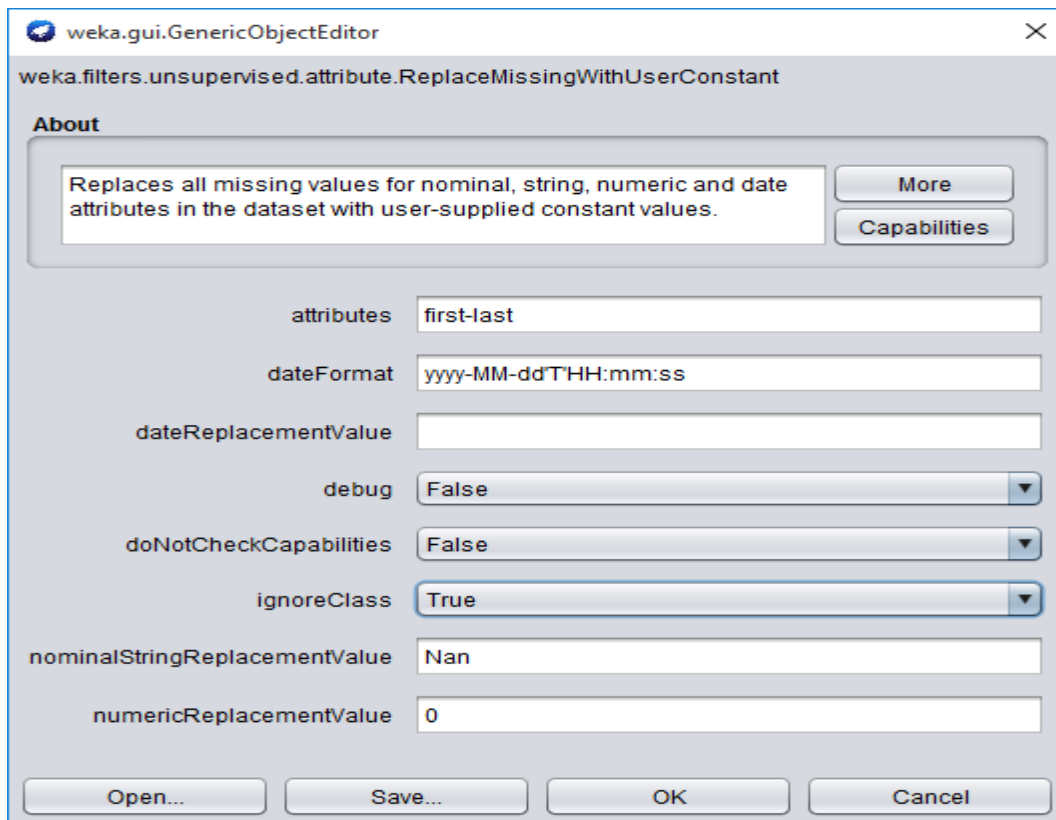
How to Impute constant values for missing values

1. Open the Weka Explorer.
2. Load the Student Data set

3. Click the “Choose” button for the filter and select “ReplaceMissingWithUserConstant” as under

Unsupervised attribute.ReplaceMissingWithUserConstant.





5. Click the “OK” button to use the configuration for the filter.
6. Click the “Apply” button to apply the filter click .

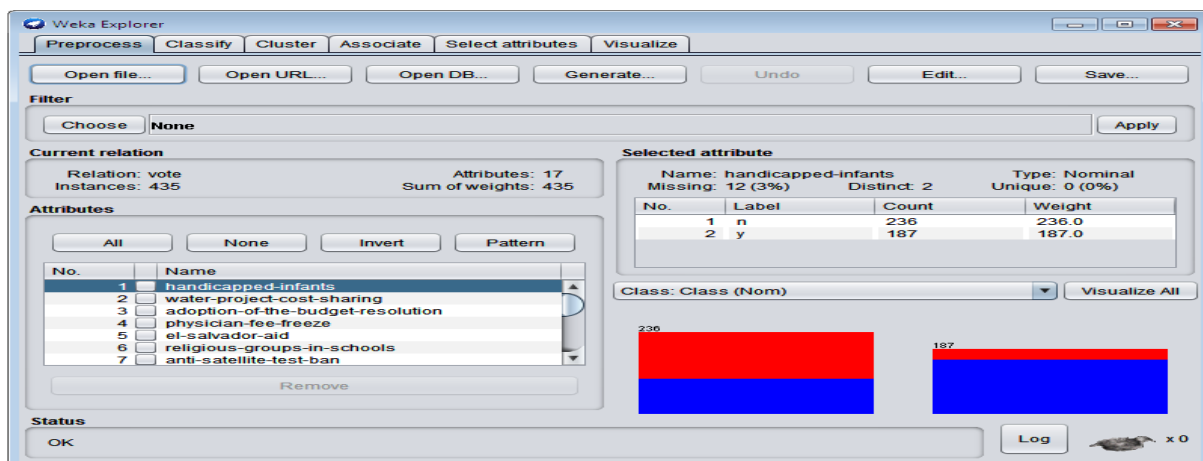
HIERARCHICAL CLUSTERING

Hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters. Strategies for hierarchical clustering generally fall into two types. Agglomerative is a "bottom up" approach: each observation starts in its own cluster, and pairs of clusters are merged as one moves up the hierarchy. Divisive is a "top down" approach: all observations start in one cluster, and splits are performed recursively as one moves down the hierarchy.

1. Open **Weka** tool and choose **Explorer**.



2. Click - **Open file...** in **preprocess** tab –choose **vote.arff**.



3. Goto **Cluster** tab – click **choose** button - select **HierarchicalClusterer**.

