

Aim : To familiarize

- Download, Installation of WEKA data mining tool kit
- Understand the features of WEKA toolkit such as Explorer, Knowledge Flow interface, Experimenter, command line interface
- Navigate the options available in WEKA

Theory

WEKA (Waikato Environment for Knowledge Analysis) is an open source software that provides tools for data preprocessing, implementation of several ML algorithms and visualization tools. It can be used to develop ML techniques and apply them to real world data mining problems. This is fully developed in Java language and provides access to SQL databases using Java Database Connectivity (JDBC).

⇒ Weka functionalities

The various functionalities include many stages in dealing with big data

- i) Raw data collection
- ii) Data Preprocessing to clean the raw data
- iii) Application of ML algorithms that is specific to the application, parameters
- iv) Output Visualization to inspect the data
- v) Selection of best ML model specific to the application

⇒ Installation of WEKA

Visit WEKA's official website and download the installation file and run it to complete the installation

⇒ Launching WEKA applications

From the screen displayed (GUI application) choose how to run the application. The five different types are

- i) Explorer
- ii) Experimenter
- iii) Knowledge Flow
- iv) Workbench
- v) Simple GUI

→ WEKA explore

The different types of ML tabs that can be seen are

- i) pre process ii) classify iii) cluster iv) Associate v) Visualize
- vi) Select attributes

Under these tabs there are pre implemented ML algorithms:

- i) Pre process Tab: Select data file, process it and make it fit for applying the various ML algorithms
- ii) Classify Tab: Several supervised and unsupervised algorithms like Linear Regression, Support Vector Machine, Decision Trees, Random Forest, Naive Bayes etc may be applied.
- iii) Cluster Tab: Some clustering algorithms are k-means, fuzzy, hierarchical and so on.
- iv) Associate Tab: This includes Apriori, Fuzzy Association, PGrowth.
- v) Select attributes Tab: Allows feature selection based on algorithms such as Classification, Principal components etc.
- vi) Visualize Tab: Allows to visualize the processed data for analysis.

Aim: To understand and perform the following operations

- Study the csv file format
- Explore the available data sets in WEKA
- Load the dataset (breast cancer) and observe the following:
 - a. List the attribute names and types
 - b. Number of records in each dataset
 - c. Identify the class attribute (if any)
 - d. Perform preprocessing (min 3)
 - e. Plot Histogram

Theory

→ Loading Data

This can be done from the following sources -

i) Local File System

Under the ML tab, click on 'Open file' button. A directory navigator window opens up through which we can navigate to the desired folder. WEKA installation comes with some sample databases (C:\Program Files\WEKA-3-8-6\data). The contents of the file would be loaded in the WEKA environment.

ii) Web

On clicking 'Open URL...' button which opens up a pop up box. Type any URL where data is stored. The Explorer will load the data from the remote site.

iii) DB

On clicking 'Open DB...' button, a window opens up where we can set the connection string to the input database, set up the query for data selection, process the query and load selected records in WEKA.

>> File formats

- The different types of files supported by WEKA includes arff, arff.gz, bci, csv, dat, data, geom, geom.gz, libsvm, m, namu, ruff, ruff.gz. The default type is arff.

- Arff format (Attribute-Relation file format)

→ This contains 2 sections - header and data where the header describes the attribute types and the data section contains a comma separated list of data.

>> Exploring datasets in WEKA

- Open the breast cancer dataset using 'open file...' option.

- Current relation sub window

→ It shows the name of the loaded dataset

→ There are 14 instances (no. of rows)

→ The table contains 5 attributes (5 fields)

- Attributes sub window

→ This appears on the left side and displays various fields in database.

→ The weather database contains 5 fields - outlook, temperature, humidity, windy and play.

- Selected attributes sub window

On selecting an attribute from list, further details can be displayed on the right side. For example, in the temperature attribute the following can be observed

→ The name and type of attribute

→ The type here is Nominal

→ The number of missing values is 0

→ There are 3 distinct values with no unique value

→ The table underneath shows the nominal values - hot, mild etc

→ It shows count, weight or frequency of 1 for each nominal value

The following can be offered

- The @relation tag defines the name of DB
- The @attribute tag defines the attributes
- The @data tag starts the list of data rows separated by comma.
- The attribute can take nominal values
@attribute outlook (sunny, overcast, rainy)
- The attribute can take real values
@attribute temperature real
- A target, class variable can also be a set
@attribute play (yes, no)

• Visualization of attributes

This can be seen at the bottom of the window, on clicking the 'Visualize All' button.

Data preprocessing is a data mining technique to transform raw data into an efficient format. The processed data is then fed to different algorithms for analysis. The various steps involved in data preprocessing are

i) Data cleaning: It involves handling of missing data, noisy data etc

a. Missing Data: It can be handled in various ways

→ Ignore the tuples (suitable for large data)

→ Fill missing values using mean or most probable value

b. Noisy data: These are the meaningless data that can't be interpreted by machines and is generated due to faulty data collection, data entry error etc.

→ Binning method (works on sorted data).

→ Regression (Linear or multiple)

→ clustering which helps recognize outliers.

ii) Data Transformation: This involves

a. Normalization (scale data values in a specified range)

b. Attribute selection (new attributes are constructed)

c. Discretization (replace raw values by interval or conceptual levels)

d. Concept hierarchy generation (connects from lower to higher hierarchy).

iii) Data reduction: This aims to increase the storage efficiency and reduce data storage and analysis costs.

a. Data cube Aggregation (for construction of data cube)

b. Attribute subset selection (discards attributes other than highly relevant data)

- c. Numerosity reduction (enables storing of data models rather than data)
- e. Dimensionality reduction (By encoding mechanisms like wavelet and PCA)

>> Data cleaning using WEKA.

- i) Replace Missing Values - with modes, means from training data
- ii) RemoveWithValues - filters instances according to the value of attributes
- iii) InterQuartileRange - Detects outliers and extreme values
- iv) Discretize - converts a range of numeric attributes to nominal attributes
- v) Normalize - Normalizes all numeric values for the given dataset
- vi) NominalToBinary - converts all nominal attributes to binary numeric attributes
- vii) NumericToNominal - converts numeric to nominal attributes
- viii) Remove - Removes a range of attributes
- ix) RemoveBy Name - Removes based on a regular expression matched but will not remove the class attribute
- x) RenameAttribute - Used for renaming attributes

Conclusion