| |
|---|
| Experiment No.8 |
| Pre-processing using WEKA tool |
| Date of Performance: |
| Date of Submission: |

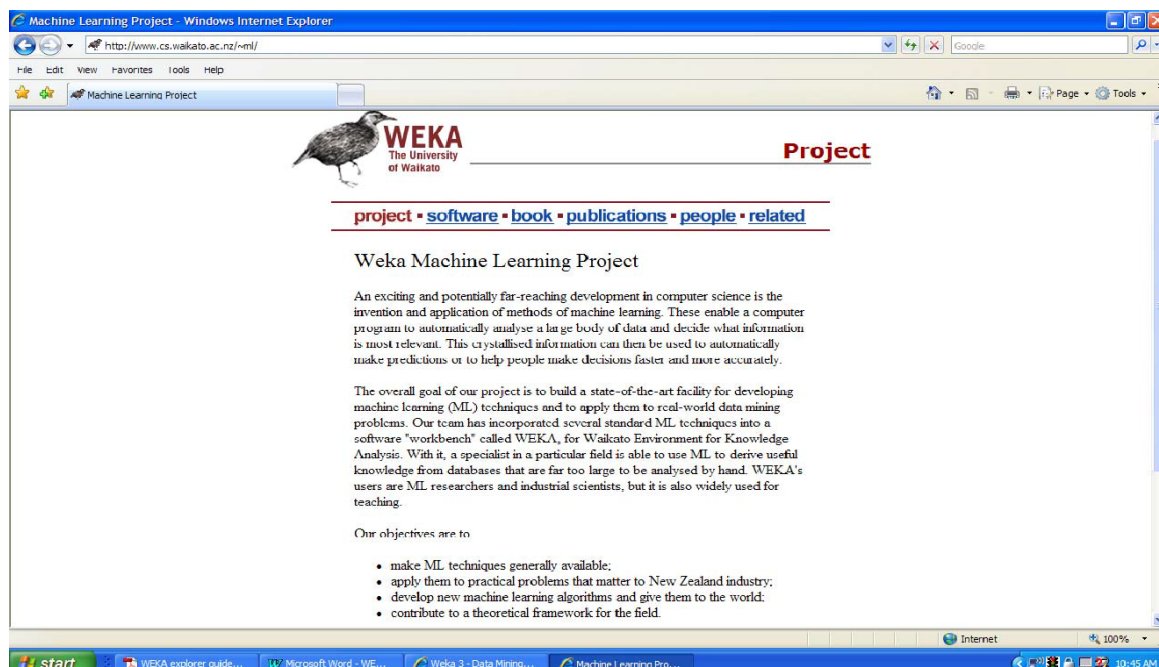**Aim:** To implement Pre-processing using WEKA tool

**Objective:** Pre-processing using WEKA tool

**Theory:**

WEKA, formally called Waikato Environment for Knowledge Learning, is a computer program that was developed at the University of Waikato in New Zealand for the purpose of identifying information from raw data gathered from agricultural domains. WEKA supports many different standard data mining tasks such as data preprocessing, classification, clustering, regression, visualization and feature selection. The basic premise of the application is to utilize a computer application that can be trained to perform machine learning capabilities and derive useful information in the form of trends and patterns. WEKA is an open source application that is freely available under the GNU general public license agreement. Originally written in C the WEKA application has been completely rewritten in Java and is compatible with almost every computing platform. It is user friendly with a graphical interface that allows for quick set up and operation. WEKA operates on the predication that the user data is available as a flat file or relation, this means that each data object is described by a fixed number of attributes that usually are of a specific type, normal alpha-numeric or numeric values. The WEKA application allows novice users a tool to identify hidden information from database and file systems with simple to use options and visual interfaces.

**Installation**

The program information can be found by conducting a search on the Web for WEKA Data Mining . The site has a very large amount of useful information on the program's benefits and background. New users might find some benefit from investigating the user manual for the program. The main WEKA site has links to this information as well as past experiments for new users to refine the potential uses that might be of particular interest to them. When prepared to download the software it is best to select the latest application from the selection offered on the site. The format for downloading the application is offered in a self installation package and is a simple procedure that provides the complete program on the end users machine that is ready to use when extracted.

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 its 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.

**Step 1: Data Pre Processing or Cleaning**

1. Launch Weka-> click on the tab Explorer
2. Load a dataset. (Click on "Open File" & locate the datafile)
3. Click on PreProcess tab & then look at your lower R.H.S. bottom window click on drop down arrow and choose "No Class"
4. Click on "Edit" tab, a new window opens up that will show you the loaded datafile. By looking at your dataset you can also find out if there are missing values in it or not. Also please note the attribute types on the column header. It would either be 'nominal' or 'numeric'.
a. If your data has missing values then its best to clean it first before you apply any forms of mining algorithm to it. In the below figure , you will see the highlighted fields are blank that means the data at hand is dirty and it first needs to be cleaned.

| | | | | | | |
|---|---|---|---|---|---|---|
| | | Translational Kin... KT* | | KT7A | 1.0 | |
| Vector | RESULT | ...lational Dy... DT* | | DT4A | 1.0 | 3.0 CORRECT |
| Equation | RESULT | | | E2A | 1.0 | 2.0 CORRECT |
| Explain-Further | HINT_MSG | missing values | | FLUIDS14 | 1.0 | 2.0 HINT |
| Equation | RESULT | Vectors | VEC* | VEC5B | 1.0 | 2.0 CORRECT |
| Equation | RESULT | Rotational Kine... KR* | | KR3A | 1.0 | 2.0 INCORRECT |
| Equation | RESULT | Translational Kin... KT* | | KT10A | 1.0 | 1.0 CORRECT |
| Equation | RESULT | Vectors | VEC* | VEC1D | 1.0 | 3.0 CORRECT |
| Vector | RESULT | Translational Kin... KT* | | KT9A | 1.0 | 12.0 CORRECT |
| Explain-Further | HINT_MSG | Circular Motion | ROTS* | ROTS3A | 1.0 | 2.0 HINT |

b.   Data Cleaning: To clean the data, you apply "Filters" to it. Generally the data will be missing with values, so the filter to apply is "ReplaceMissingWithUserConstant" (the filter choice may vary according to your need, for more information on it please consult the resources).Click on Choose button below Filters-> Unsupervised->attribute————————> ReplaceMissingWithUserConstant



**Code and output**:

**Conclusion**: In conclusion, WEKA's usefulness lies in its ability to simplify various aspects of the machine learning workflow, from data preprocessing to model evaluation and deployment. The outputs it generates, including performance metrics, visualizations, and model files, empower users to make data-driven decisions and effectively communicate their results, making it a valuable tool for data scientists, researchers, and analysts.