

Experiment 1

Aim

To perform data pre-processing task using Weka data mining tool.

Theory

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

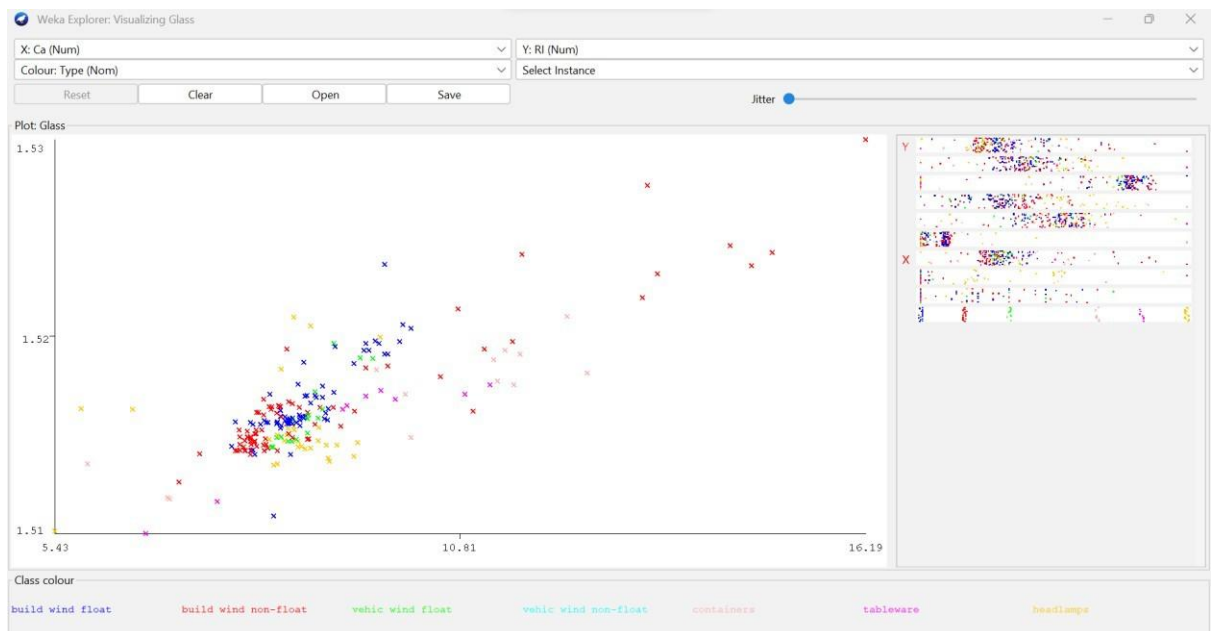
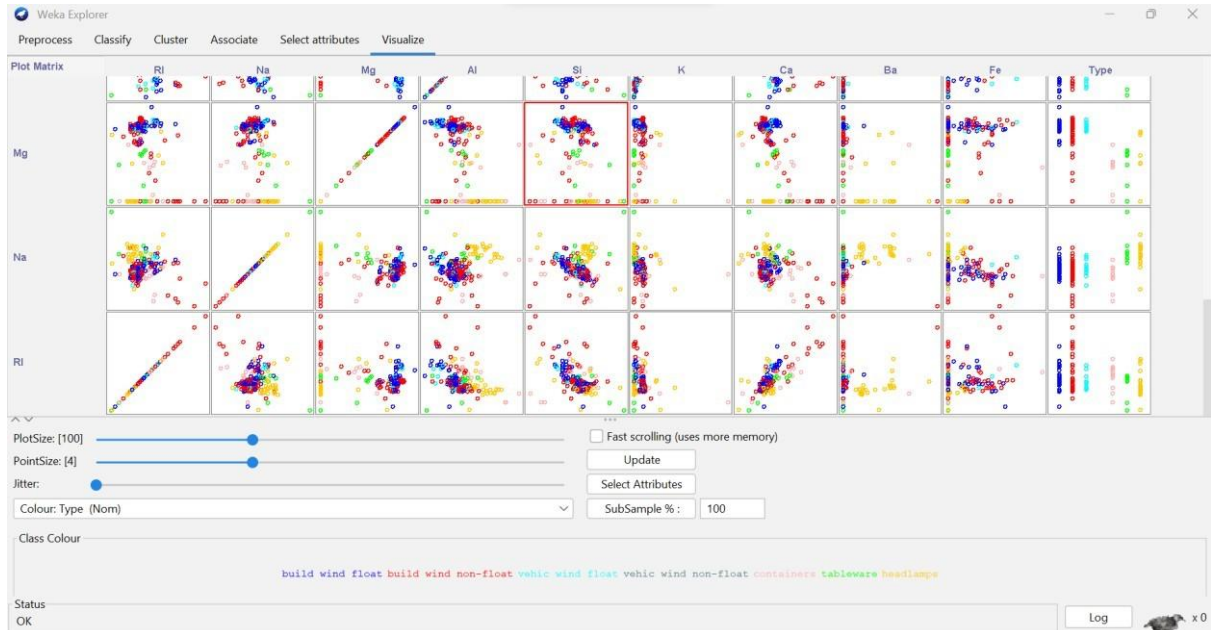
Tasks performed through Weka:

1. **Pre-processing:** Data pre-processing is a data mining technique which is used to transform the raw data in a useful and efficient format. The steps involved in this include Data Cleaning, Data Transformation and Data Reduction.
2. **Classification:** Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data.
3. **Clustering:** The process of making a group of abstract objects into classes of similar objects is known as clustering. In the process of cluster analysis, the first step is to partition the set of data into groups with the help of data similarity, and then groups are assigned to their respective labels.
4. **Association Rule:** Association rule mining finds interesting associations and relationships among large sets of data items. This rule shows how frequently an itemset occurs in a transaction. A typical example is a Market Based Analysis.
5. **Select Attributes:** The attribute selection task essentially consists in selecting a subset of originally available attributes to be subsequently used for model creation. A search strategy is needed for any attribute selection technique that is based on evaluating attribute subsets rather than single attributes.
6. **Visualization:** Data visualization is the graphical representation of information and data. By using visual elements like charts, graphs, and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data.

TASKS

Pre-processing activities to be observed in Weka (performed on the 'glass' dataset) :

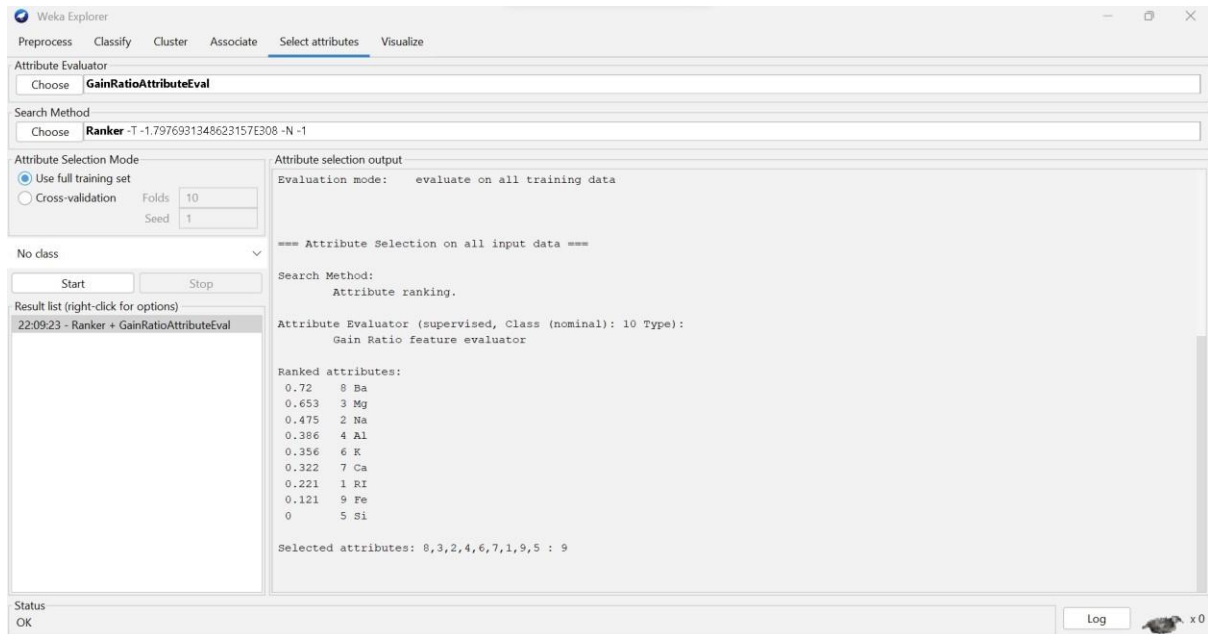
1. Visualization : Visualize scatter plot for all the attributes from dataset selected from Weka. Determine correlation if any using these plots for different datasets.



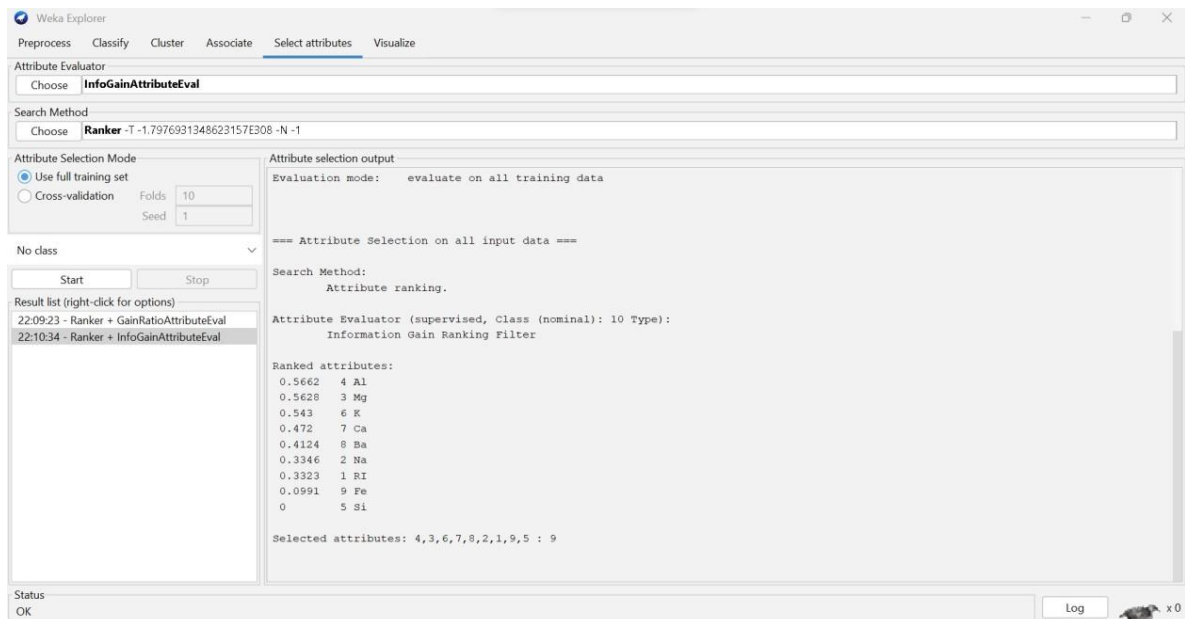
Observation – There is positive correlation in this graph which means that the relationship between two variables(types) tend to move in the same direction.

2. Select Attributes : Apply suitable feature selection filter like GainRatio etc to choose relevant attributes from the list of attributes. Observe the ranks / priority provided by the filter.

GainRatioAttributeEval –

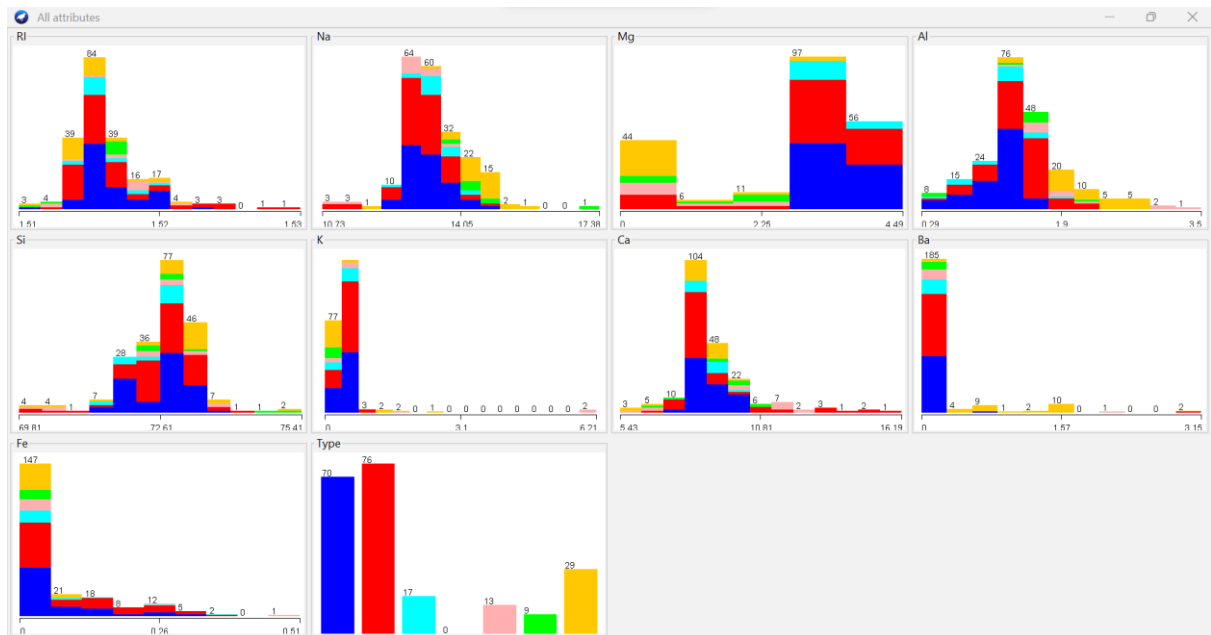


InfoGainAttributeEval –



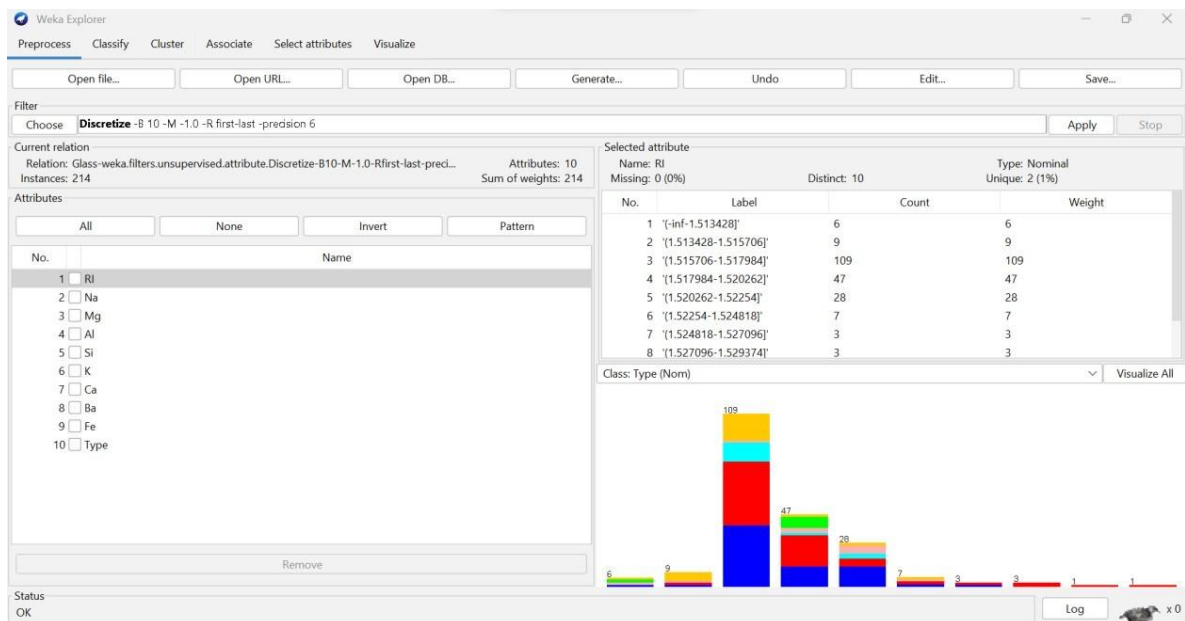
3. Pre-processing :

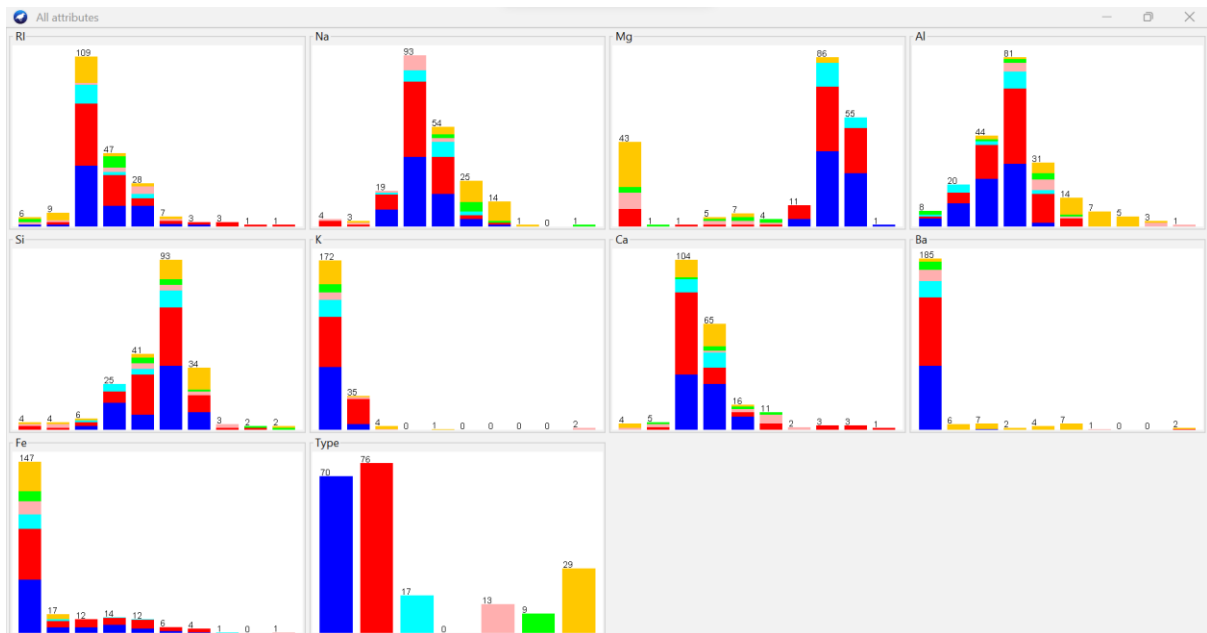
a) Visualize All: Select this button to visualize histograms of all attributes.



b) Filter: Choose Discretization under Unsupervised and Supervised methods. Observe the discretization and the outliers.

Discretization under Unsupervised method –





Discretization under Supervised method –

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

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

Filter Choose **Discretize** -R first-last -precision 6 Apply Stop

Current relation
Relation: Glass-weka.filters.supervised.attribute.Discretize-Rfirst-last-precision6
Instances: 214
Attributes: 10
Sum of weights: 214

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> RI
2	<input type="checkbox"/> Na
3	<input type="checkbox"/> Mg
4	<input type="checkbox"/> Al
5	<input type="checkbox"/> Si
6	<input type="checkbox"/> K
7	<input type="checkbox"/> Ca
8	<input type="checkbox"/> Ba
9	<input type="checkbox"/> Fe
10	<input type="checkbox"/> Type

Remove

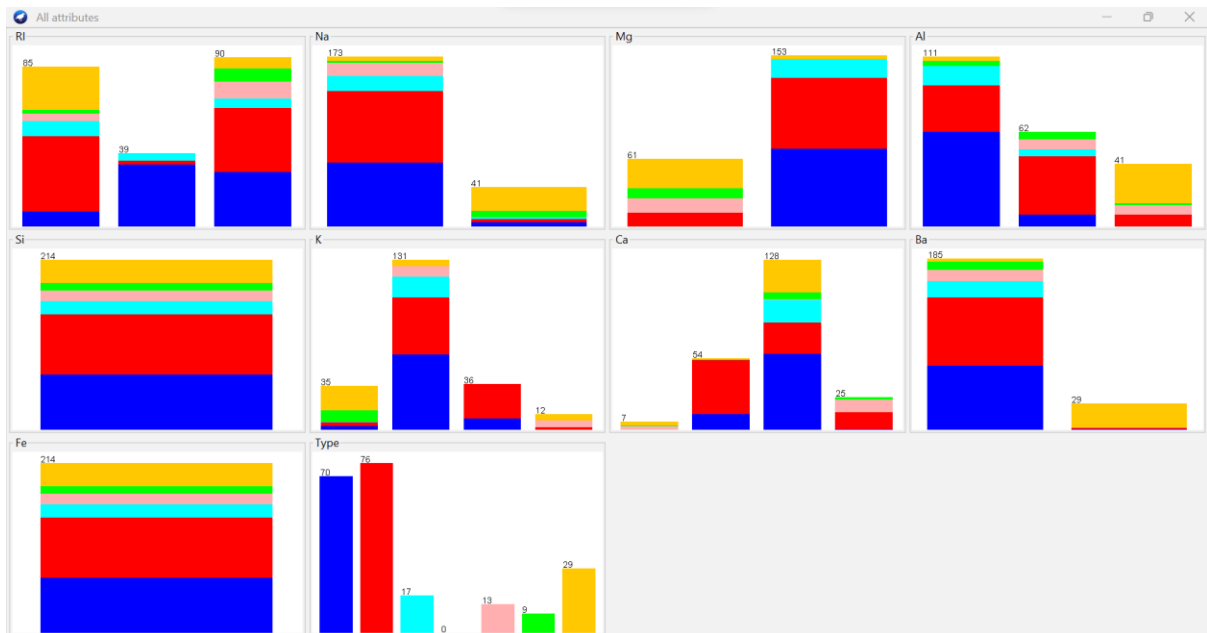
Status OK

Selected attribute
Name: RI
Missing: 0 (0%)
Distinct: 3
Type: Nominal
Unique: 0 (0%)

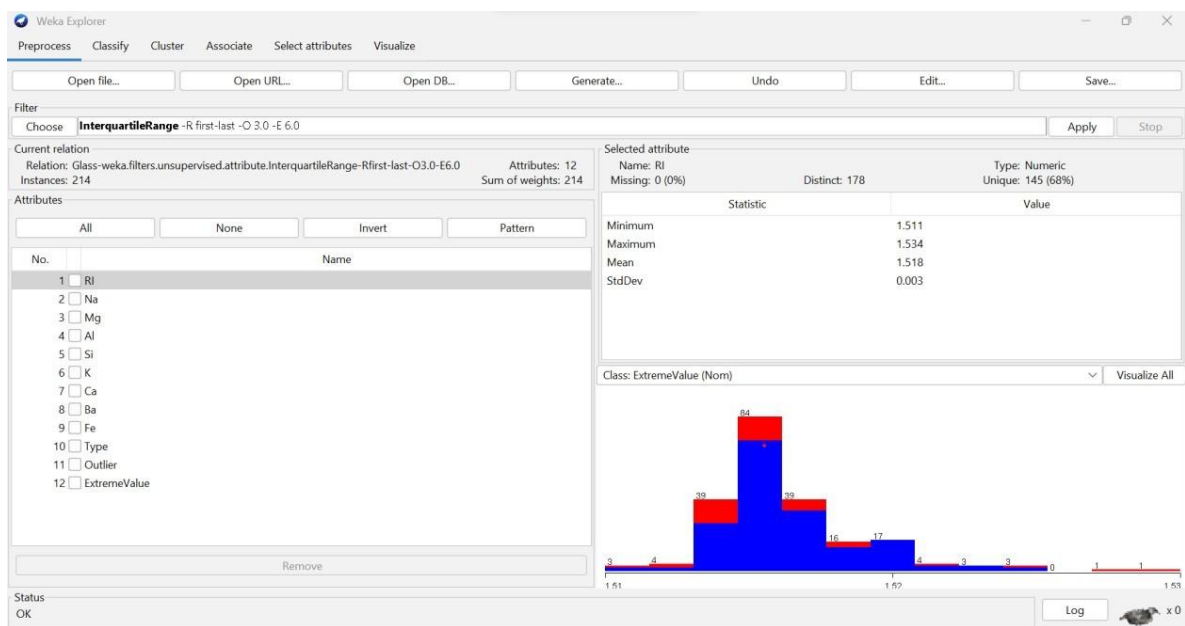
No.	Label	Count	Weight
1	'(-inf-1.517335]'	85	85
2	'(1.517335-1.517985]'	39	39
3	'(1.517985-inf)'	90	90

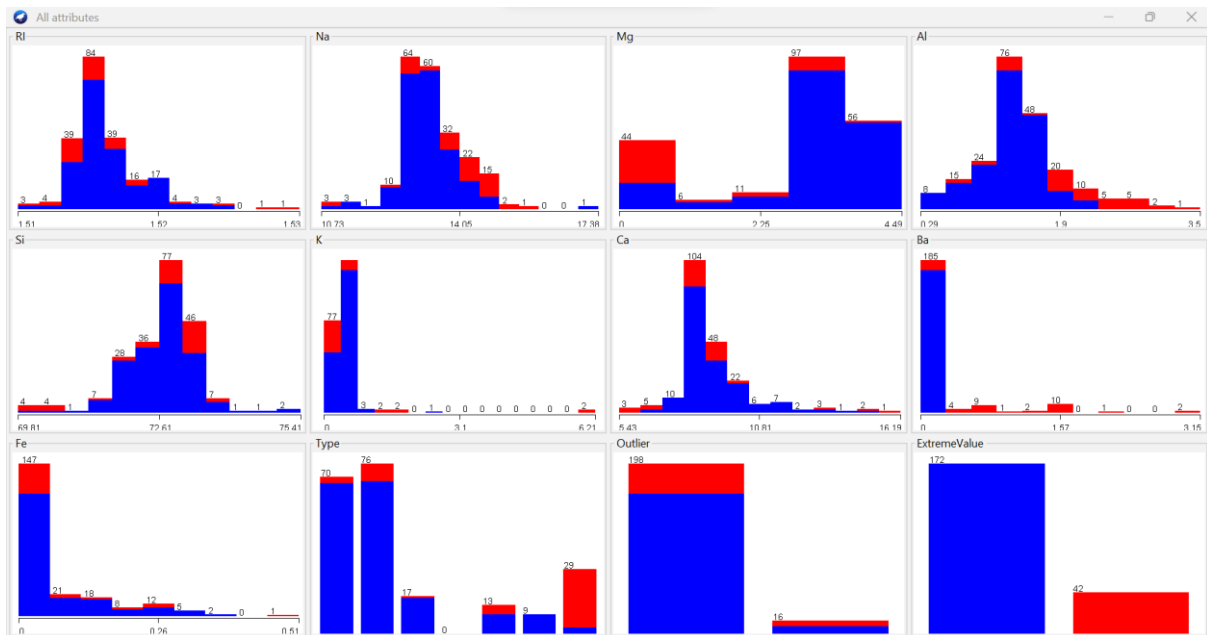
Class: Type (Nom) Visualize All

Log x 0



c) IQR: Observe the IQR values for a selected attribute. Observe the outlier and extreme values.

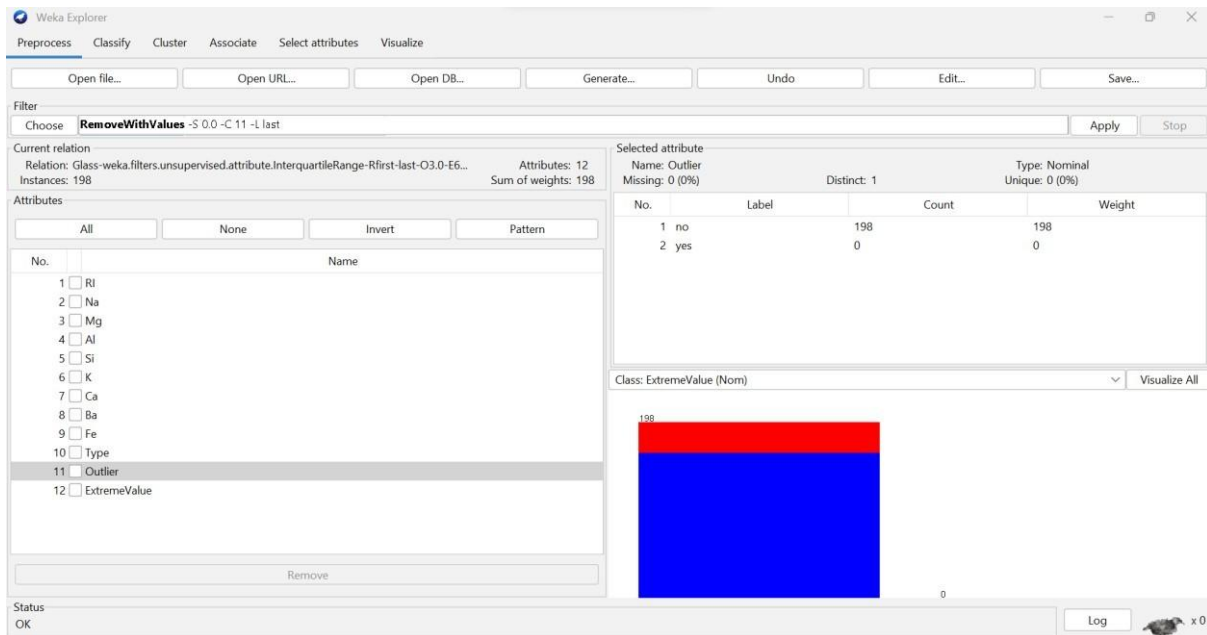




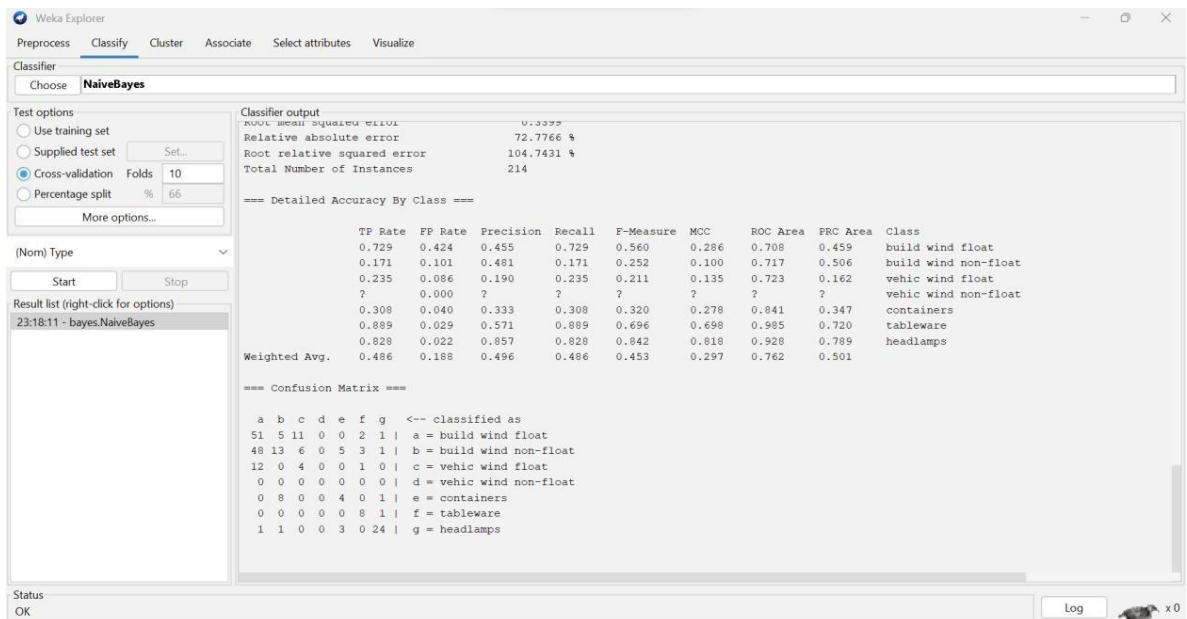
d) Remove the value: Remove instances with outlier values and show the screenshots of dataset before and after the removal.

The screenshot shows the Weka Explorer interface with the 'InterquartileRange' filter applied to the 'Outlier' attribute. The 'Filter' tab is active, and the 'InterquartileRange' filter is selected. The 'Current relation' shows 214 instances. The 'Attributes' list on the left includes RI, Na, Mg, Al, Si, K, Ca, Ba, Fe, Type, Outlier, and ExtremeValue. The 'Selected attribute' section shows the 'Outlier' attribute with a 'no' label and a count of 198, and a 'yes' label and a count of 16. The 'Class: ExtremeValue (Nom)' is selected, and a visualization of the data is shown.

No.	Label	Count	Weight
1	no	198	198
2	yes	16	16



4. Classification : Perform NB, kNN and DT/rule-based classification. Naive Bayes –



kNN –

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'IBk -K 5 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A "weka.core.EuclideanDistance -R first-last"'. The test options are set to 'Cross-validation' with 'Folds' set to 10 and 'Percentage split' set to 66. The result list shows '23:20:41 - lazy.IBk'. The classifier output displays the following metrics:

Metric	Value
Root mean squared error	0.2303
Relative absolute error	51.243 %
Root relative squared error	78.9576 %
Total Number of Instances	214

Below the metrics is a 'Detailed Accuracy By Class' table:

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.843	0.229	0.641	0.843	0.728	0.582	0.867	0.713	build wind float
0.684	0.174	0.684	0.684	0.684	0.510	0.848	0.756	build wind non-float
0.000	0.010	0.000	0.000	0.000	-0.029	0.642	0.161	vehic wind float
?	0.000	?	?	?	?	?	?	vehic wind non-float
0.385	0.025	0.500	0.385	0.435	0.407	0.952	0.546	containers
0.667	0.010	0.750	0.667	0.706	0.695	0.909	0.565	tableware
0.793	0.016	0.885	0.793	0.836	0.814	0.890	0.843	headlamps
Weighted Avg.	0.678	0.142	0.635	0.678	0.651	0.533	0.853	

Below this is a 'Confusion Matrix' table:

a	b	c	d	e	f	g	<-- classified as
59	10	1	0	0	0	0	a = build wind float
20	52	1	0	3	0	0	b = build wind non-float
12	5	0	0	0	0	0	c = vehic wind float
0	0	0	0	0	0	0	d = vehic wind non-float
0	5	0	0	5	0	3	e = containers
0	2	0	0	1	6	0	f = tableware
1	2	0	0	1	2	23	g = headlamps

DecisionTable –

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The classifier chosen is 'DecisionTable -X 1 -S "weka.attributeSelection.BestFirst -D 1 -N 5"'. The test options are set to 'Cross-validation' with 'Folds' set to 10 and 'Percentage split' set to 66. The result list shows '23:22:40 - rules.DecisionTable'. The classifier output displays the following metrics:

Metric	Value
Root mean squared error	0.2768
Relative absolute error	81.4177 %
Root relative squared error	85.2945 %
Total Number of Instances	214

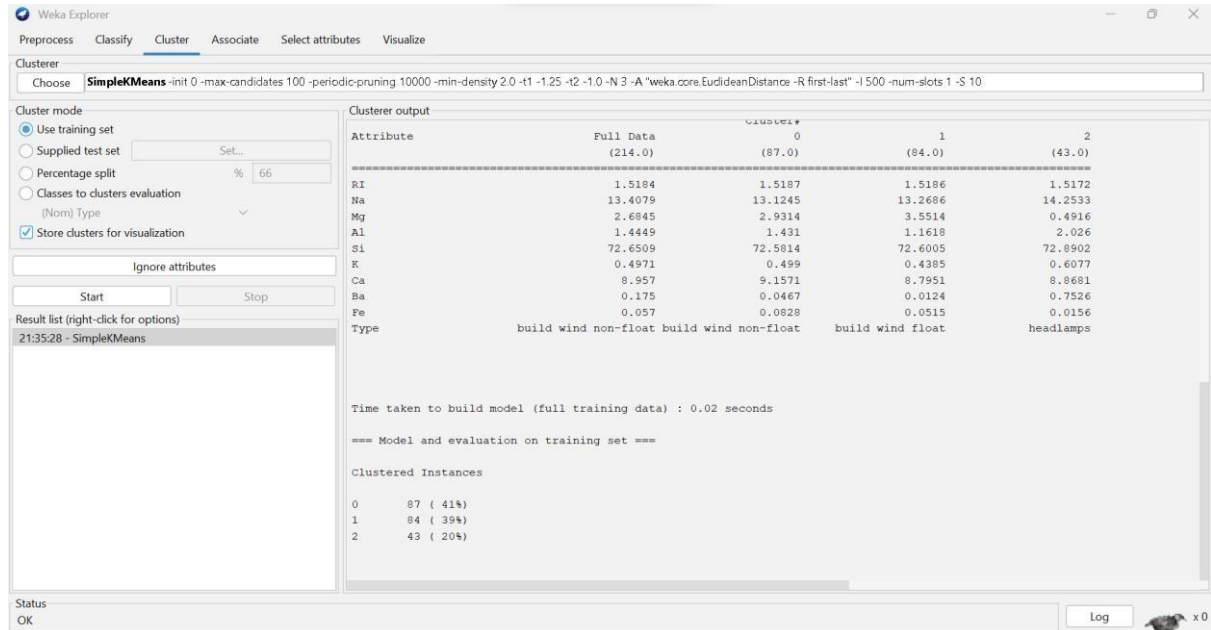
Below the metrics is a 'Detailed Accuracy By Class' table:

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.886	0.264	0.620	0.886	0.729	0.585	0.839	0.670	build wind float
0.671	0.167	0.689	0.671	0.680	0.507	0.792	0.720	build wind non-float
0.176	0.015	0.500	0.176	0.261	0.264	0.570	0.163	vehic wind float
?	0.000	?	?	?	?	?	?	vehic wind non-float
0.538	0.010	0.778	0.538	0.636	0.629	0.873	0.532	containers
0.667	0.010	0.750	0.667	0.706	0.695	0.853	0.514	tableware
0.586	0.000	1.000	0.586	0.739	0.742	0.926	0.813	headlamps
Weighted Avg.	0.682	0.148	0.702	0.682	0.669	0.560	0.815	

Below this is a 'Confusion Matrix' table:

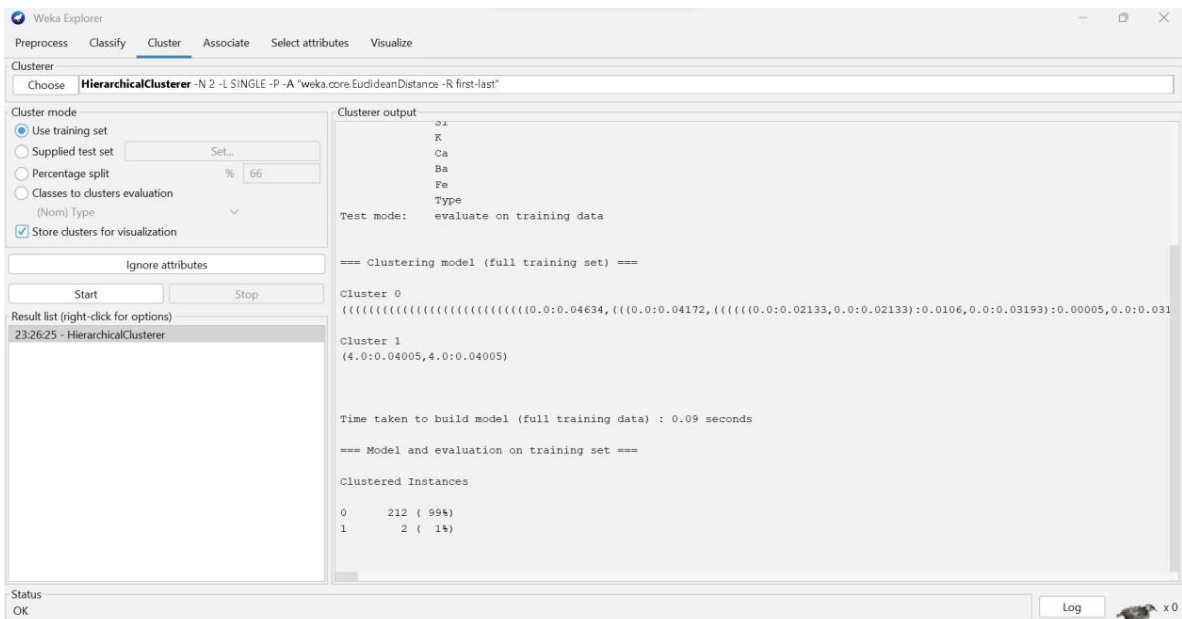
a	b	c	d	e	f	g	<-- classified as
62	7	1	0	0	0	0	a = build wind float
21	51	1	0	1	2	0	b = build wind non-float
11	3	3	0	0	0	0	c = vehic wind float
0	0	0	0	0	0	0	d = vehic wind non-float
1	5	0	0	7	0	0	e = containers
0	3	0	0	0	6	0	f = tableware
5	5	1	0	1	0	17	g = headlamps

5. Clustering : Perform kmeans, hierarchical clustering and explain the output. KMeans –



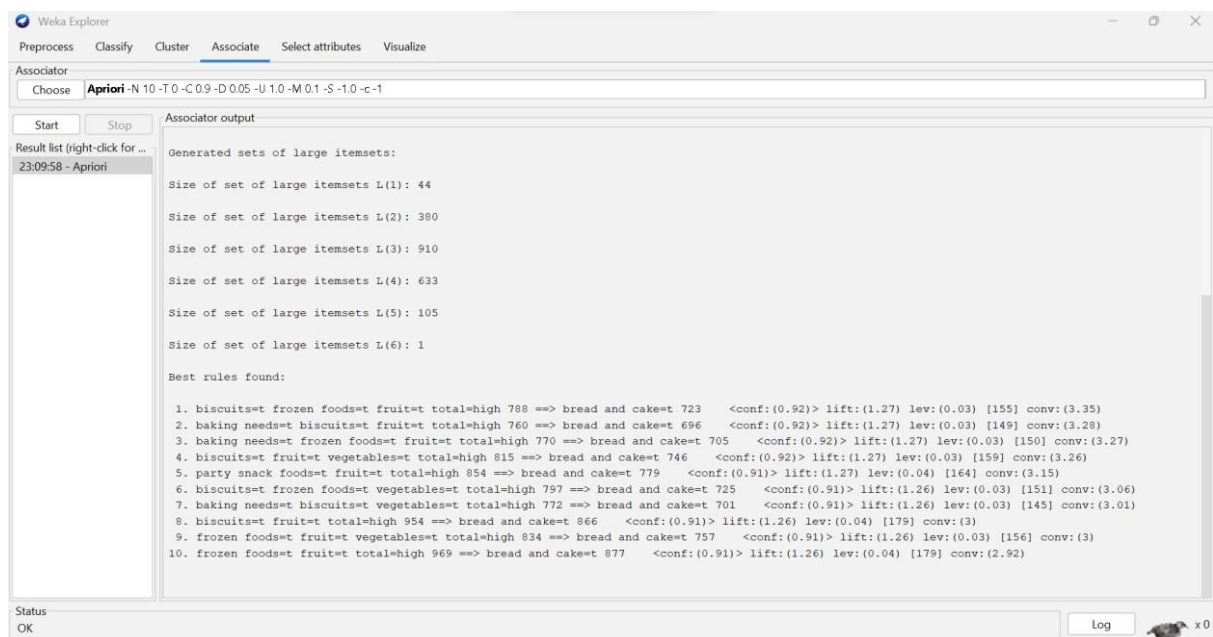
Observation – SimpleKMeans Clustering is an Unsupervised Learning Algorithm. It groups the unlabelled dataset into different clusters. We have to define the number of pre-defined clusters that need to be created in the process and then discovered the categories of groups in the unlabelled dataset on its own without the need for any training. Here, I defined three clusters and the output shows all the instances divided into the three different clusters.

HierarchicalCluster –



Observation - Hierarchical clustering is another unsupervised machine learning algorithm, which is used to group the unlabelled datasets into a cluster. In this algorithm, we develop the hierarchy of clusters in the form of a tree, and this tree-shaped structure is known as the dendrogram. In this algorithm, we do not need to have knowledge about the predefined number of clusters. Here, I set the number of clusters to be 2 and the output shows the instances are divided into the respective two clusters.

6. Association rule mining : Perform apriori algo and show the rules created. (Performed on the ‘supermarket’ dataset)



Conclusion :

Weka tool is useful and convenient tool to perform various preprocessing and data mining algorithms.