

DATA MINING

Project Report

Group Members:

- | | |
|-----------------------|-----------|
| 1) Dhivahar Perumal | 111465042 |
| 2) Soham Mehta | 111496015 |
| 3) Manu Mathew | 111492994 |
| 4) Arjun Mathew Dan | 111492985 |
| 5) Shilpa Mary George | 111492833 |

1. INTRODUCTION

The aim of the project is to learn about the various classification tools and machine learning techniques. We build two types of classifiers i.e. a descriptive and a non-descriptive classifier. We are using the Weka(Waikato Environment for Knowledge Analysis) tool which provides a suite of machine learning algorithms for data mining tasks. We use it for data preparation, pre-processing, classification and visualization.

2. ABOUT DATA

The data provided to us is a real-life classification data with TYPE DE ROCHE (Rock Type) as a CLASS attribute. There are total 98 records with 48 attributes and 6 classes.

The Classes for which we have train our models are:

- C1: R. Carbonatees AND R. Carbonatees impures
- C2: Pyrate
- C3: Charcopyrite
- C4: Galene
- C5: Spahlerite
- C6: Sediments terrigenes

Most important attributes are: S, Zn, Pb, Cu, CaO+MgO, CaO, MgO, Fe₂O₃ as they have the maximum values provided for training our model.

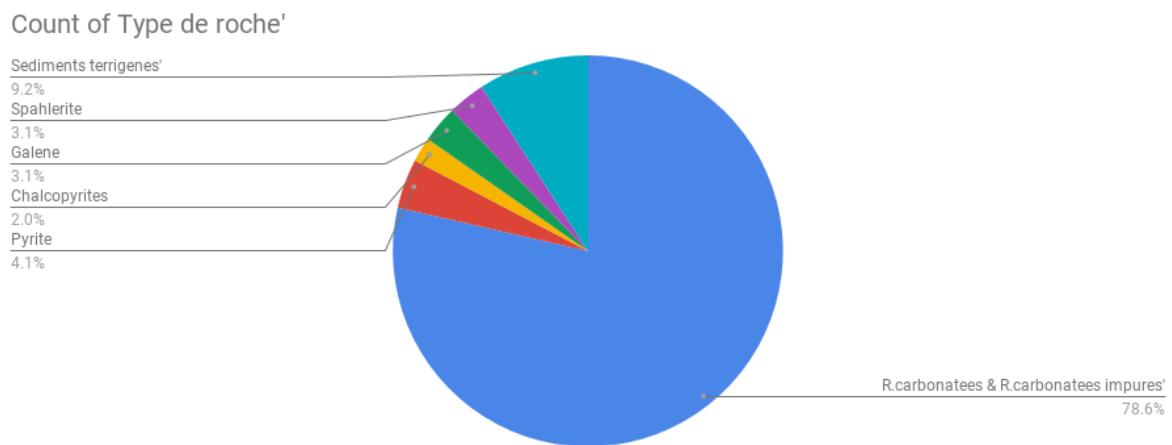
Since this is a real life experimental data, it also contains a lot of missing data and duplicate data. There are also many instances where the data is not properly formatted or in the correct format. These all corrections are made in the data preparation step.

3. OBSERVATIONS

The data provided has following characteristics. The insights can be obtained about the same by plotting these relevant values. The main observation from the given data is that there is a high variability in the types of rocks in the given data set.

Type de roche'	Number of Rocks
Chalcopyrites	2
Galene	3
Pyrite	4
R.carbonatees & R.carbonatees impures'	77
Sediments terrigenes'	9
Spahlerite	3
Grand Total	98

The data provided has such high variance that 78% of the rocks belong only a single category. The plot for the given data is shown below:

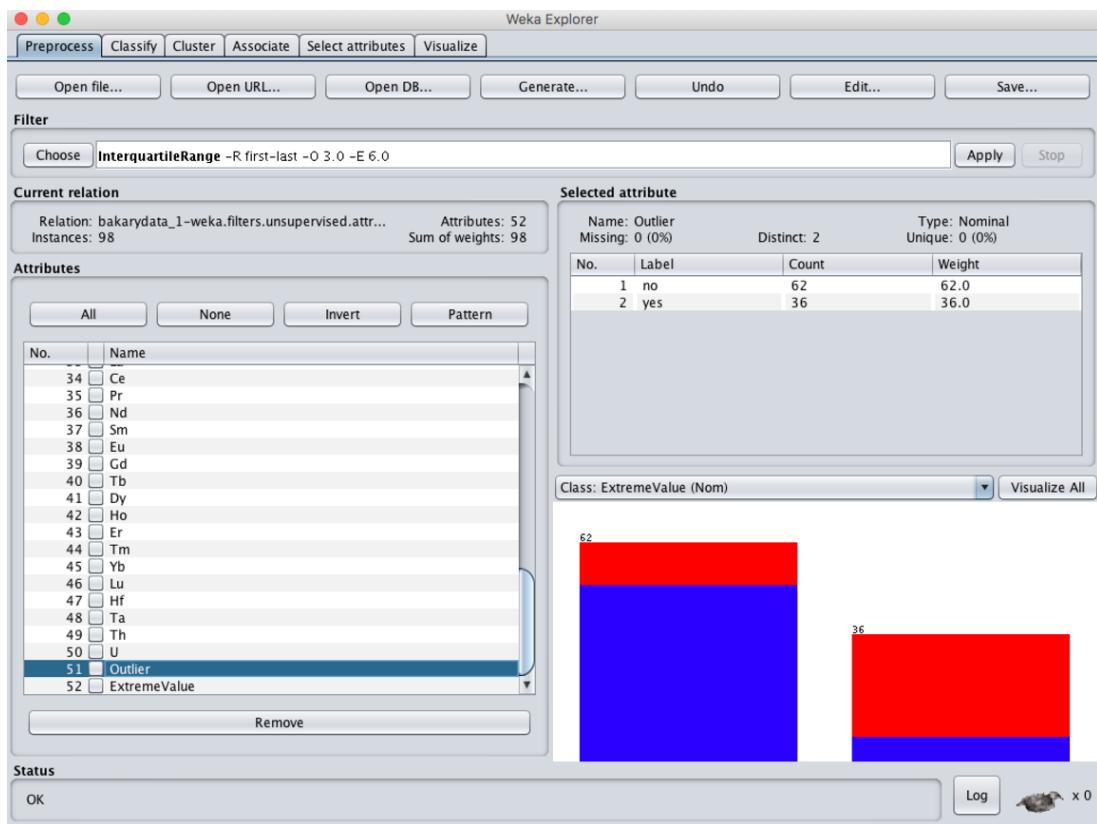


4. DATA PREPARATION

The data given initially had a lot of missing values, duplicate values and values with incorrect formatting. All these inaccuracies are removed while preparing the data.

4.1 Finding Outliers

We first identified the values which were acting as outliers in the given dataset. This can be performed using Weka method “InterQuartile Range”[1] to find the Outlier and Extreme values in a given dataset.



4.2 Cleaning Misformatted Data Values

Some values are not in a proper format as required. Weka can train the model only based on values which are either numerical or nominal. *E.g. Li- Sédiments terrigènes is given a value <0.3 which is replaced with 0.25.*

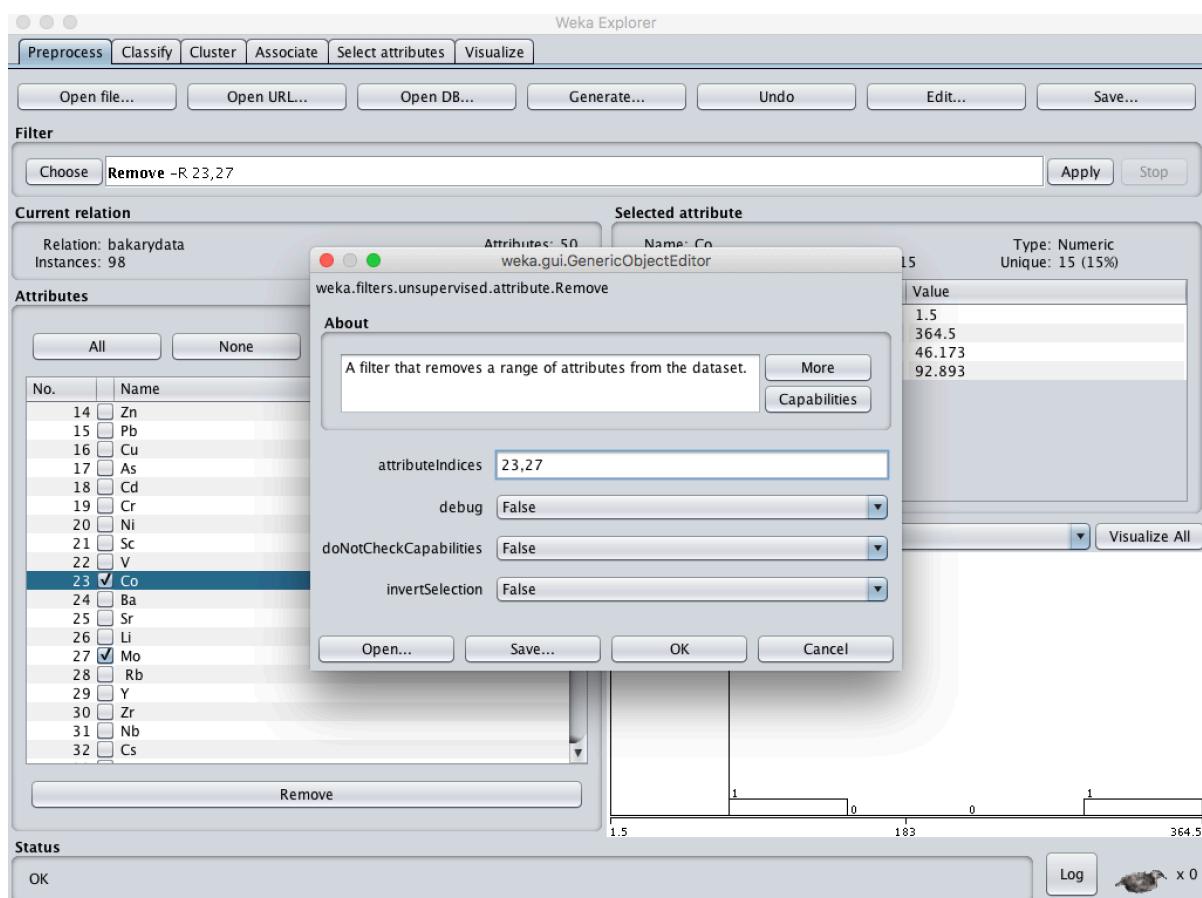
4.3 Solving Class Label Discrepancies

There were a few values in which the class labels were misspelled. This caused the model to consider this misspelled value as another class label. These values were corrected to obtain a uniform class label *E.g. pyrrite and pyrite*

There were a few values which belong to the same class. These were combined to obtain a uniform label for both. *E.g. carbonates and impures belong to a single group.*

4.4 Dropping Irrelevant Attributes

There were a few attributes which had a lot of missing values. Imputing these values would have caused an incorrect prediction. Its best to drop these attributes since they provide none or sometimes negative impact on training[2]. *E.g. MO and CO have less than 20 values*

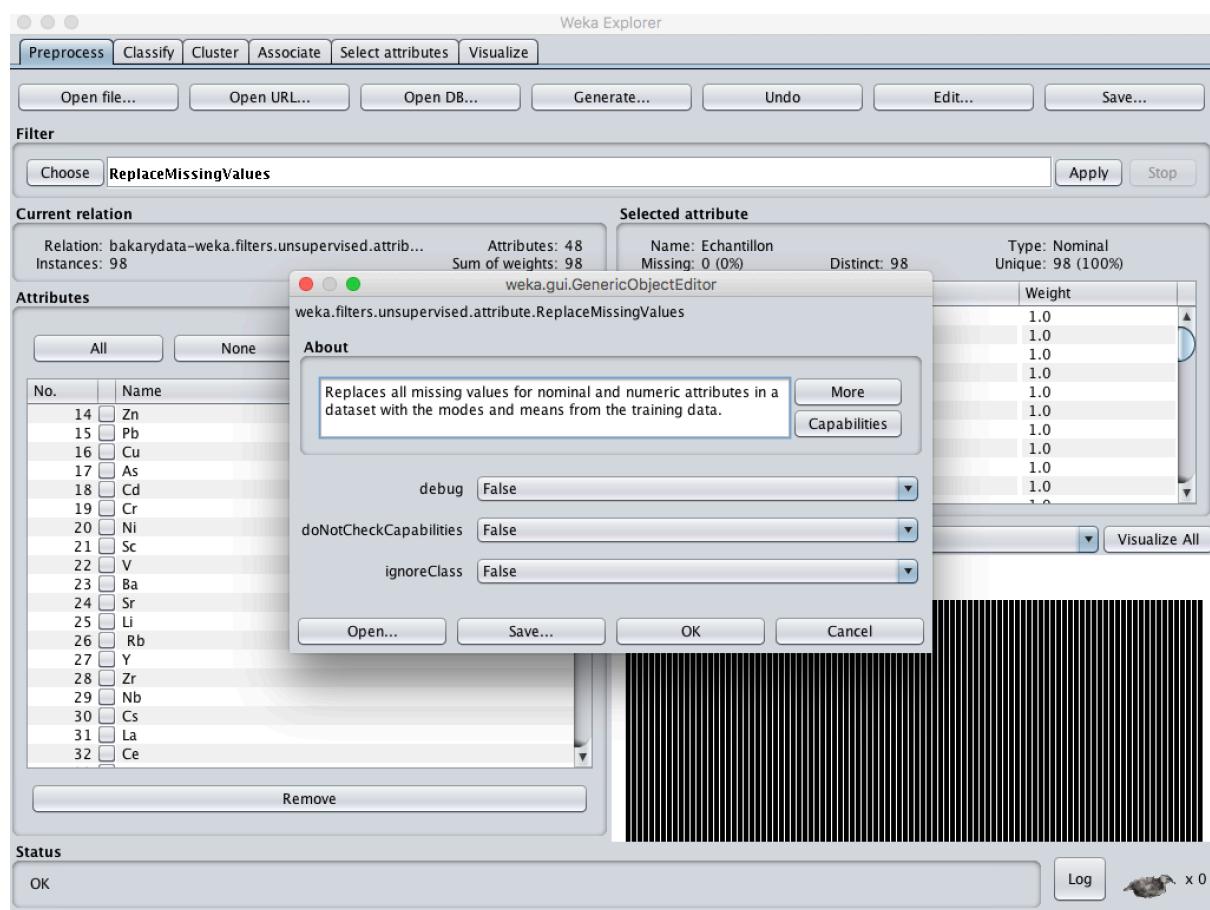


4.5 Dropping Duplicate and Irrelevant Attributes

There are multiple rows in the given dataset which are repeating or they have no values at all. These all rows are just dropped so as to create a rigid structure where all rows and attributes are filled. E.g. There is a complete empty row and a duplicate header.

4.6 Filling Missing Values

Finally, after all the cleaning is done, the missing values need to be filled so that the model can be trained. There are multiple ways in which these values can be filled. One method is to fill 0 in the missing value. However, most of the times it provides inaccurate results as 0 is itself a value used by some of the other cells. Hence it is best approach to fill it with either mean, median or mode. We have filled these missing values with the mean of that attribute.

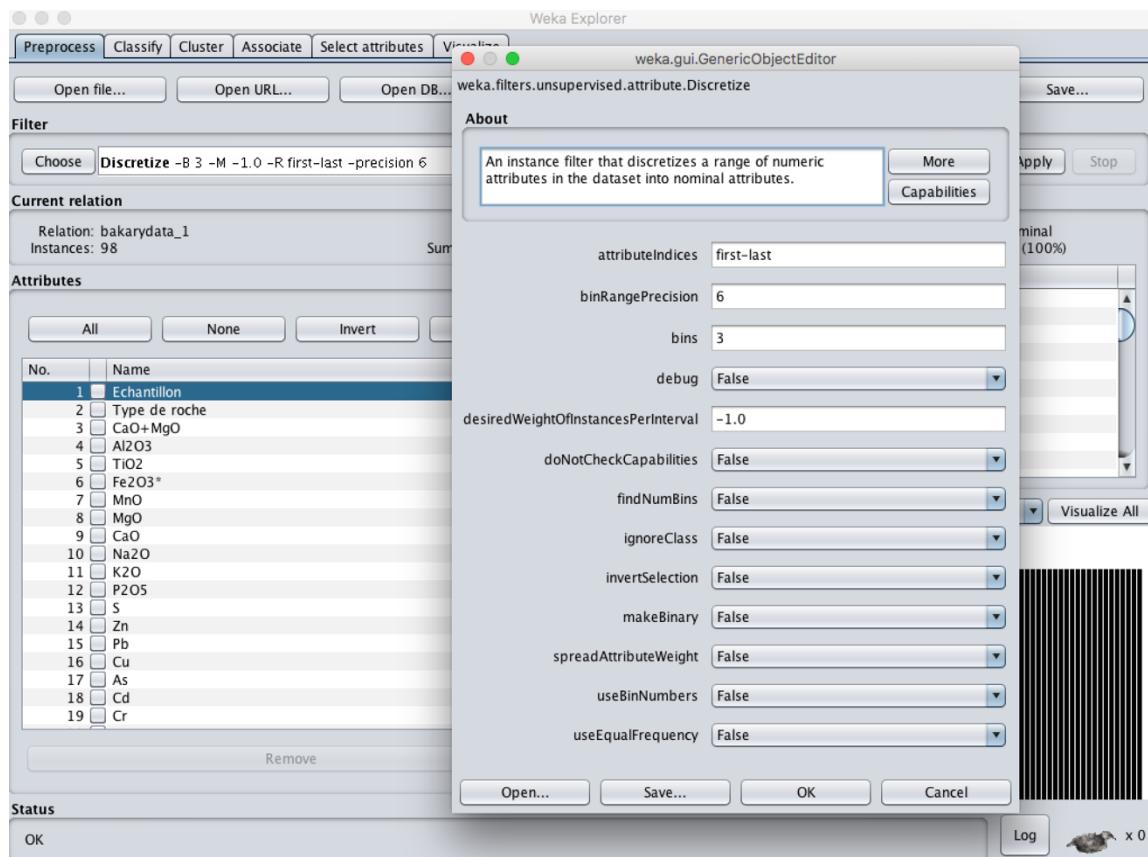


5. DATA PREPROCESSING

Here, we use our prepared data and process it using data discretization. The dataset after cleaning is divided into two data sets namely PD1 and PD2. The discretization is performed using the method of binning with equal width and equal frequency respectively. Also the cleaned project data is normalized to get the dataset PD.

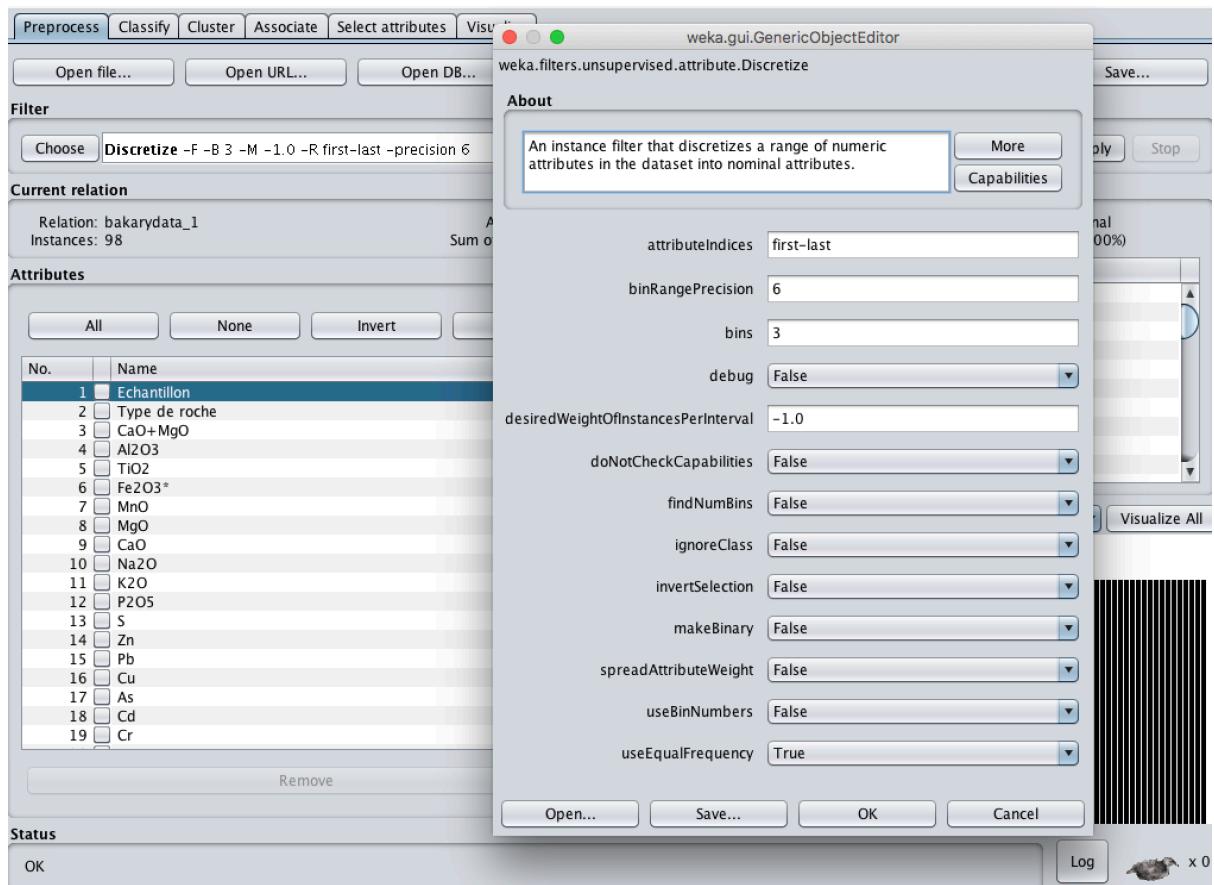
5.1 Discretization by Equal Width Binning (PD1)

Equal width binning divides the range into N intervals/ bins of the same size. We use this dataset for the Decision Trees Descriptive Classifier. The following setting were used in Weka for Discretization by Equal Width Binning.



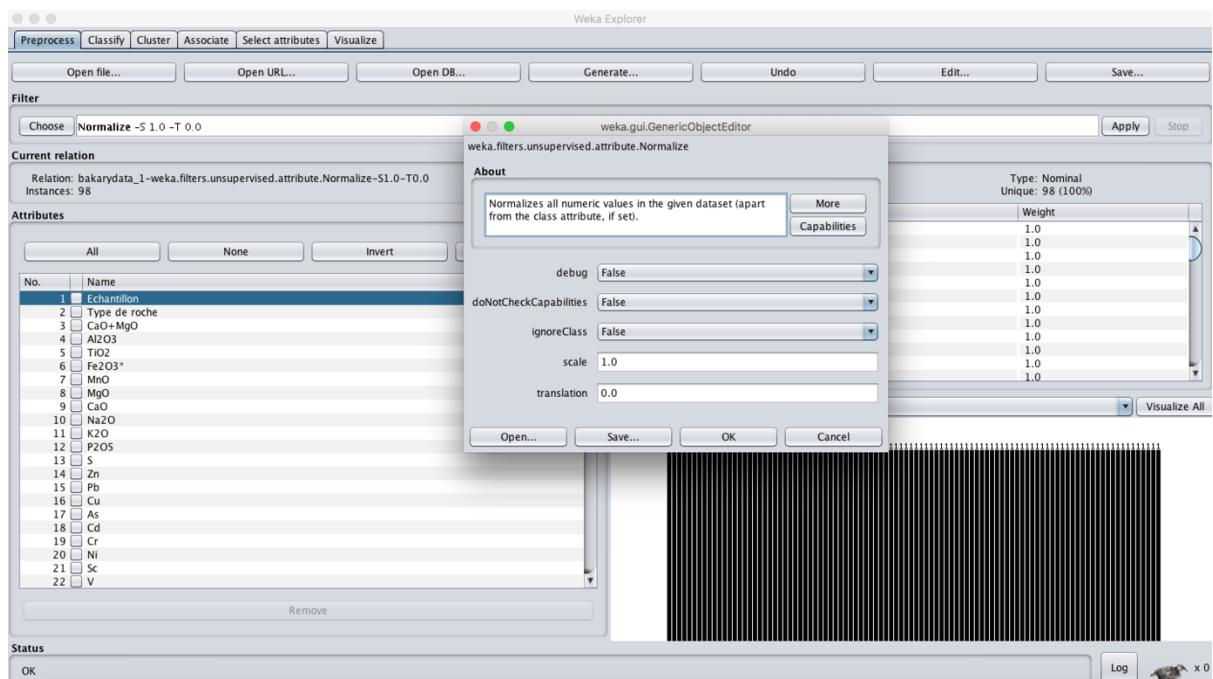
5.2 Discretization by Equal Frequency Binning (PD2)

This is also known as Equal Depth binning. It divides the range into N intervals/ bins, each containing nearly the same number of samples. We use this dataset for the Decision Trees Descriptive Classifier. The following setting were used in Weka for Discretization by Equal Frequency Binning.



5.3 Normalization (PD)

Here we adjust the values measured on different scales to a notionally common scale. We use this dataset for the Neural Network Non-Descriptive Classifier. The following setting were used in Weka for Normalizing.



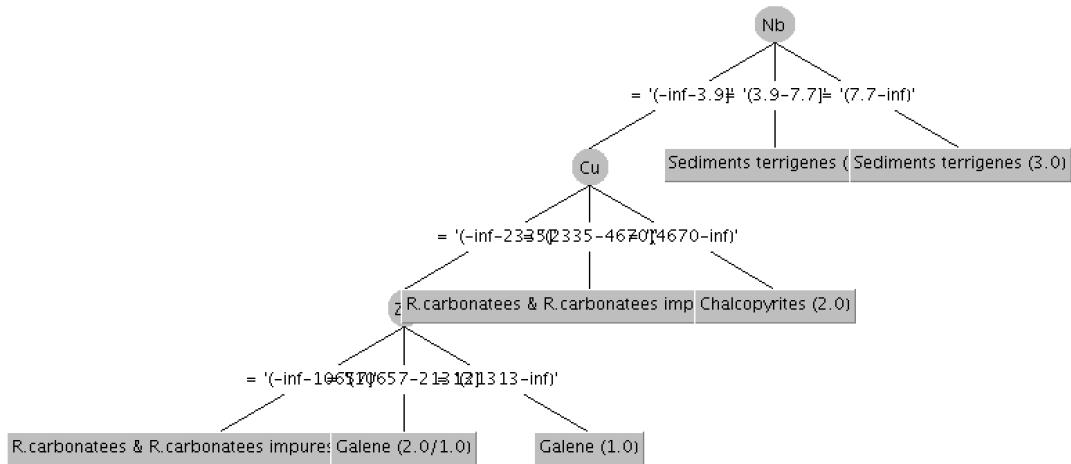
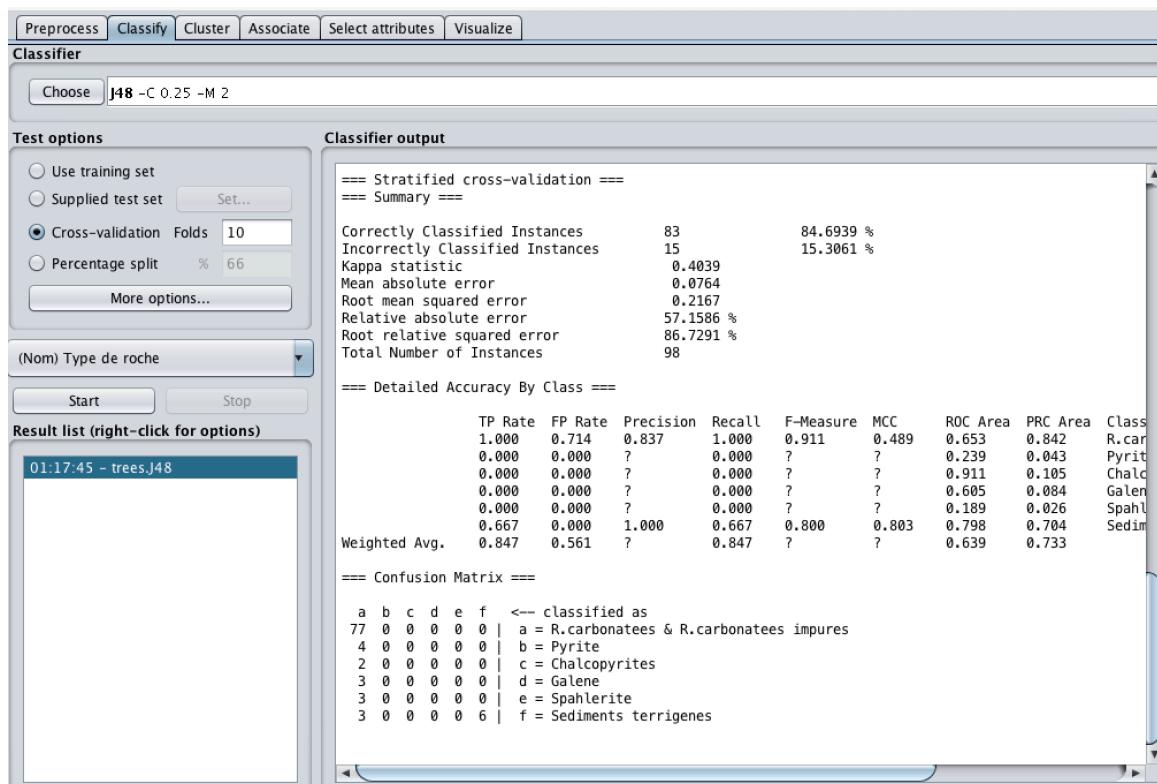
6. EXPERIMENTS

The following 3 experiments were conducted as part of this project.

6.1. EXPERIMENT 1 (Non-Contrast Learning):

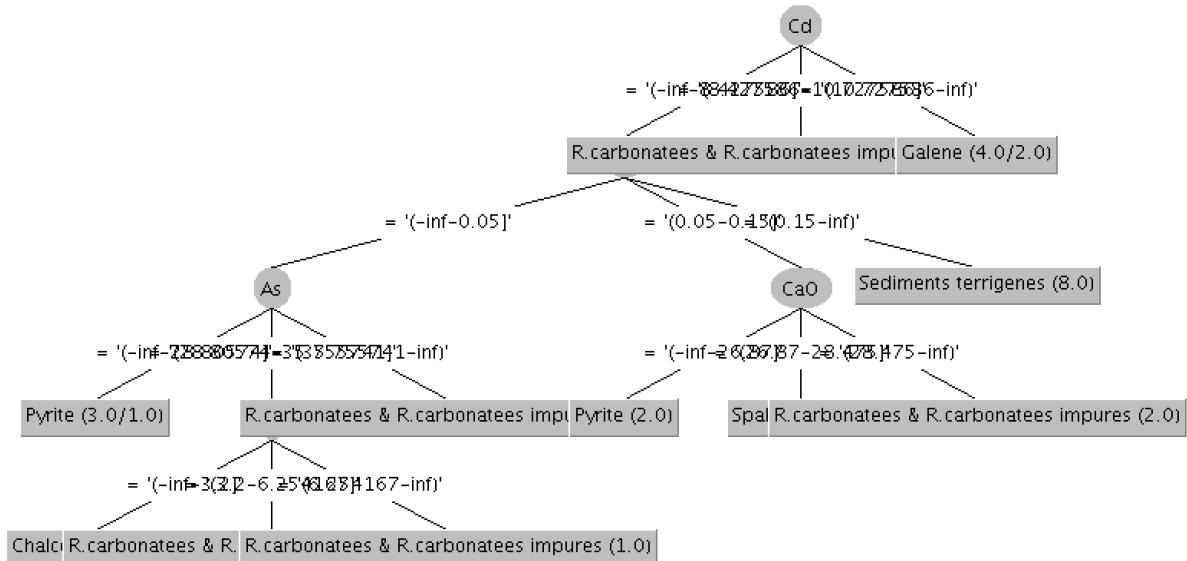
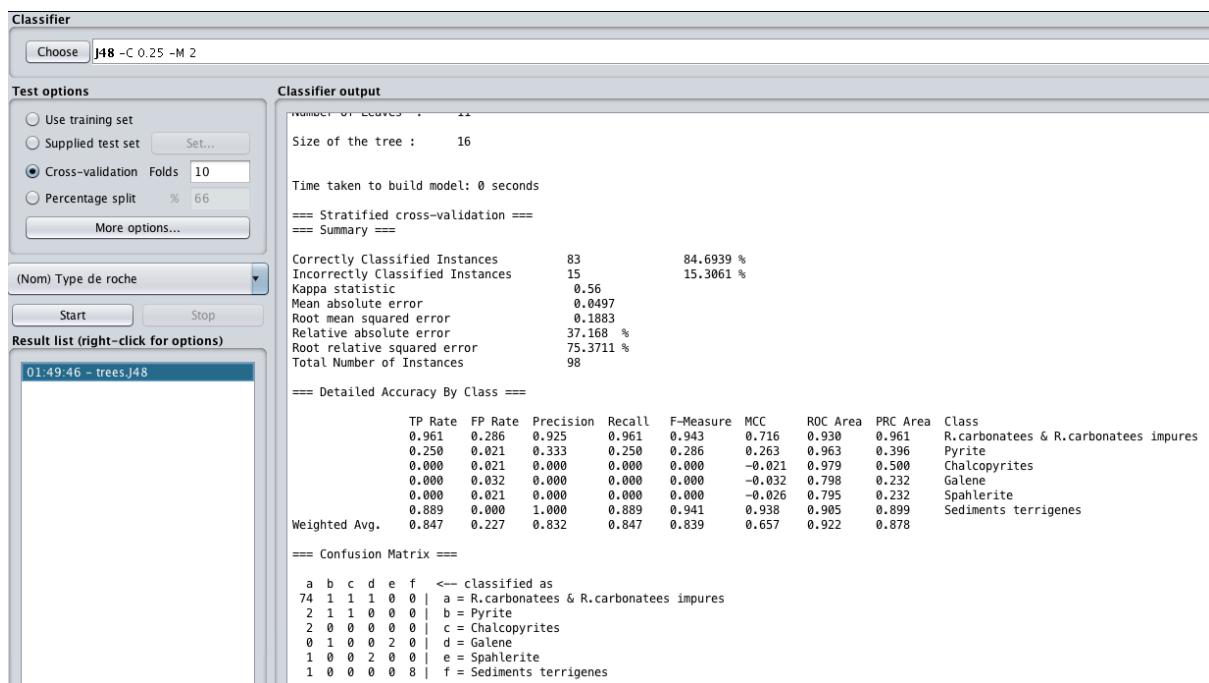
All the records were used to perform the full classification (learning), i.e. built a classifier for all classes C1- C6 simultaneously. This experiment was repeated for the 3 datasets PD1, PD2 and PD.

a) With PD1- Equal Width Binning:



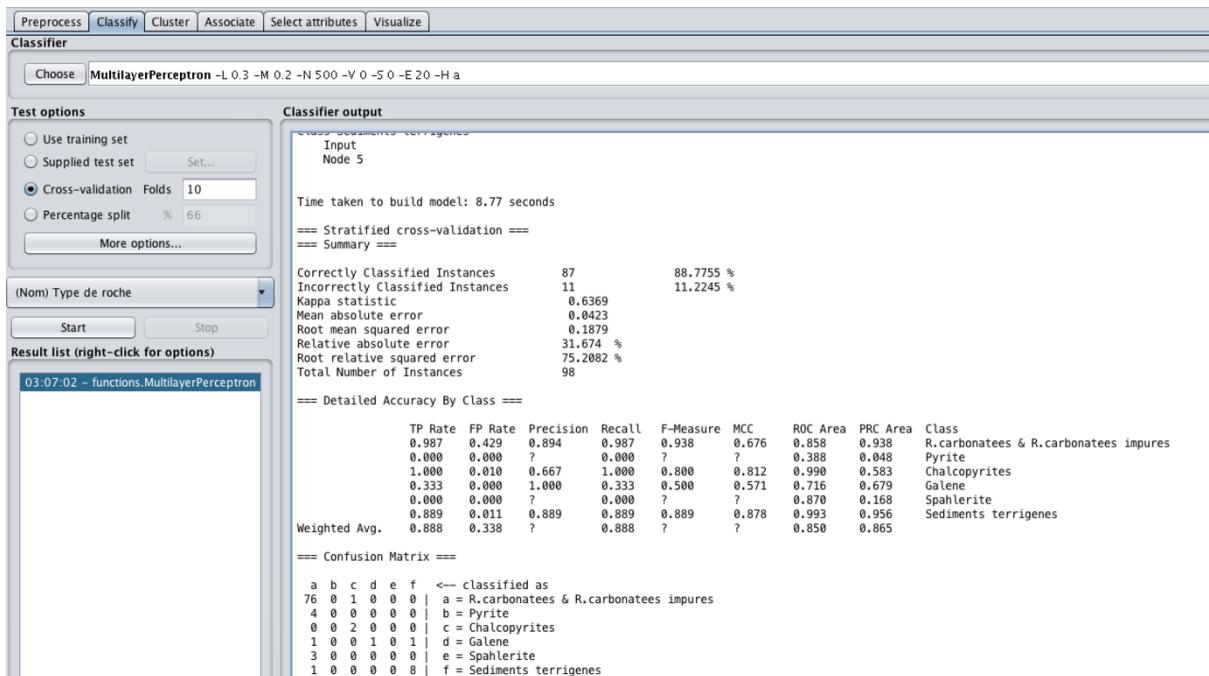
ACCURACY: 84.69%

b) With PD2 - Equal Frequency Binning:



ACCURACY: 84.69%

c) With PD - Normalized:



ACCURACY: 88.77%

6.2. EXPERIMENT 2 (Contrast Learning)

All the records were used to perform the contrast classification i.e. contrasting class **C1**(R. carbonatees and R. carbonatees impures) with a class **not C1** that contains other classes. This experiment was repeated for the 3 datasets PD1, PD2 and PD.

a) With PD1- Equal Width Binning:

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 -C 0.25 -M 2

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

(Nom) Type de roche

Start Stop

Result list (right-click for options)

```
01:29:11 - treesJ48
Correctly Classified Instances      85          86.7347 %
Incorrectly Classified Instances   13          13.2653 %
Kappa statistic                   0.4916
Mean absolute error               0.1978
Root mean squared error           0.3396
Relative absolute error            58.0778 %
Root relative squared error       82.693 %
Total Number of Instances         98
```

Time taken to build model: 0 seconds

== Stratified cross-validation ==

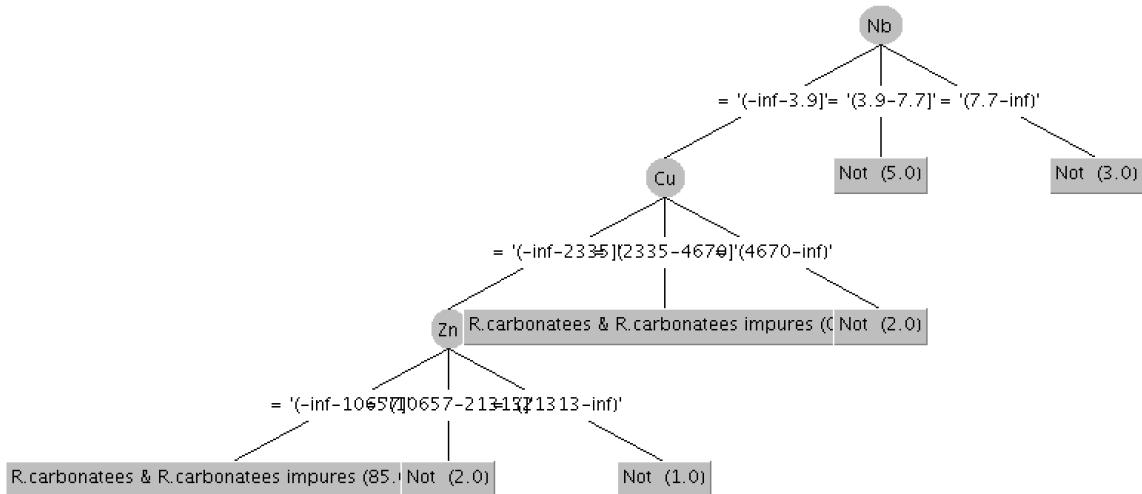
== Summary ==

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
1.000	0.619	0.856	1.000	0.922	0.571	0.724	0.860	R.carbonatees & R.carbonatees impures	
0.381	0.000	1.000	0.381	0.552	0.571	0.724	0.610	Not	
Weighted Avg.	0.867	0.486	0.887	0.867	0.843	0.571	0.724	0.806	

== Detailed Accuracy By Class ==

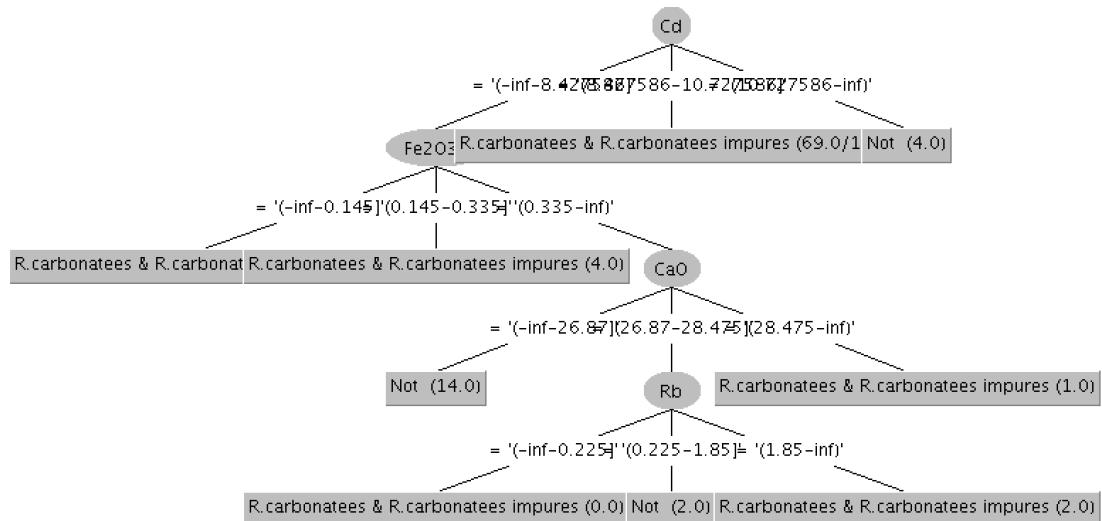
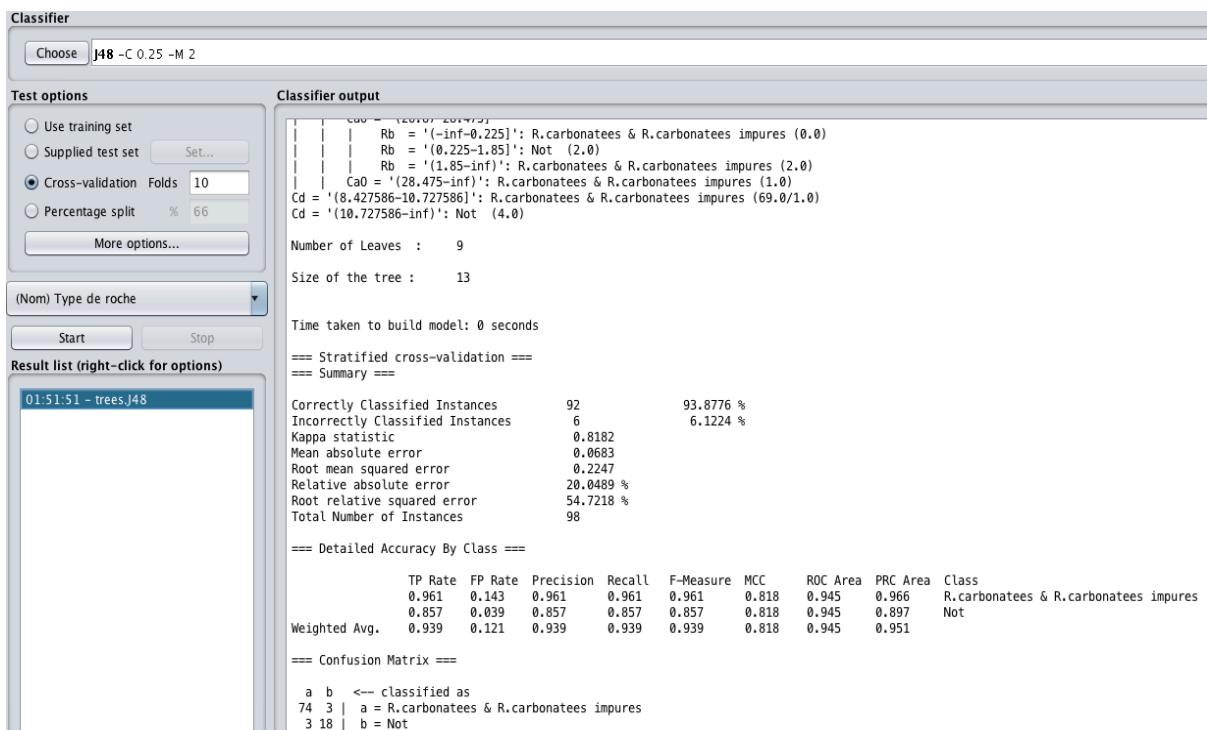
	a	b	<-- classified as
77	0	1	a = R.carbonatees & R.carbonatees impures
13	8	1	b = Not

== Confusion Matrix ==



ACCURACY: 86.73%

b) With PD2 - Equal Frequency Binning:



ACCURACY: 93.87%

c) With PD - Normalized:

The screenshot shows the Weka interface with the following configuration and output:

Classifier: Choose MultilayerPerceptron -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H a

Test options:

- (radio button selected) Cross-validation Folds 10
- Other options: Use training set, Supplied test set, Percentage split % 66, More options...

Result list (right-click for options):

03:12:27 – functions.MultilayerPerceptron

Classifier output:

```

Attrib Ld  0.002155423241200000
Attrib Hf  0.041674998118051365
Attrib Ta  0.01404410477128227
Attrib Th  0.001202385392645318
Attrib U   0.04777205691135958
Class R.carbonates & R.carbonates impures
Input
Node 0
Class Not
Input
Node 1

Time taken to build model: 8.06 seconds

== Stratified cross-validation ==
== Summary ==

Correctly Classified Instances      87          88.7755 %
Incorrectly Classified Instances    11          11.2245 %
Kappa statistic                     0.6051
Mean absolute error                 0.1185
Root mean squared error             0.3037
Relative absolute error              34.8022 %
Root relative squared error        73.9537 %
Total Number of Instances           98

== Detailed Accuracy By Class ==

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC    ROC Area  PRC Area  Class
0.987    0.476    0.884    0.987    0.933    0.639    0.920    0.975  R.carbonates & R.carbonates impures
0.524    0.013    0.917    0.524    0.667    0.639    0.920    0.827  Not
Weighted Avg.                      0.888    0.377    0.891    0.888    0.876    0.639    0.920    0.944

== Confusion Matrix ==

a  b  <-- classified as
76  1 |  a = R.carbonates & R.carbonates impures
10 11 |  b = Not
  
```

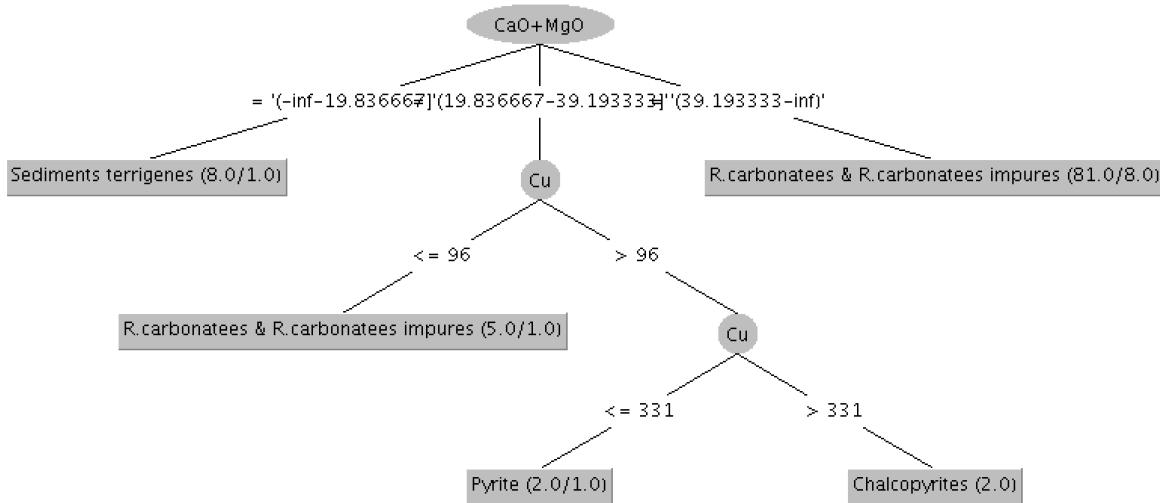
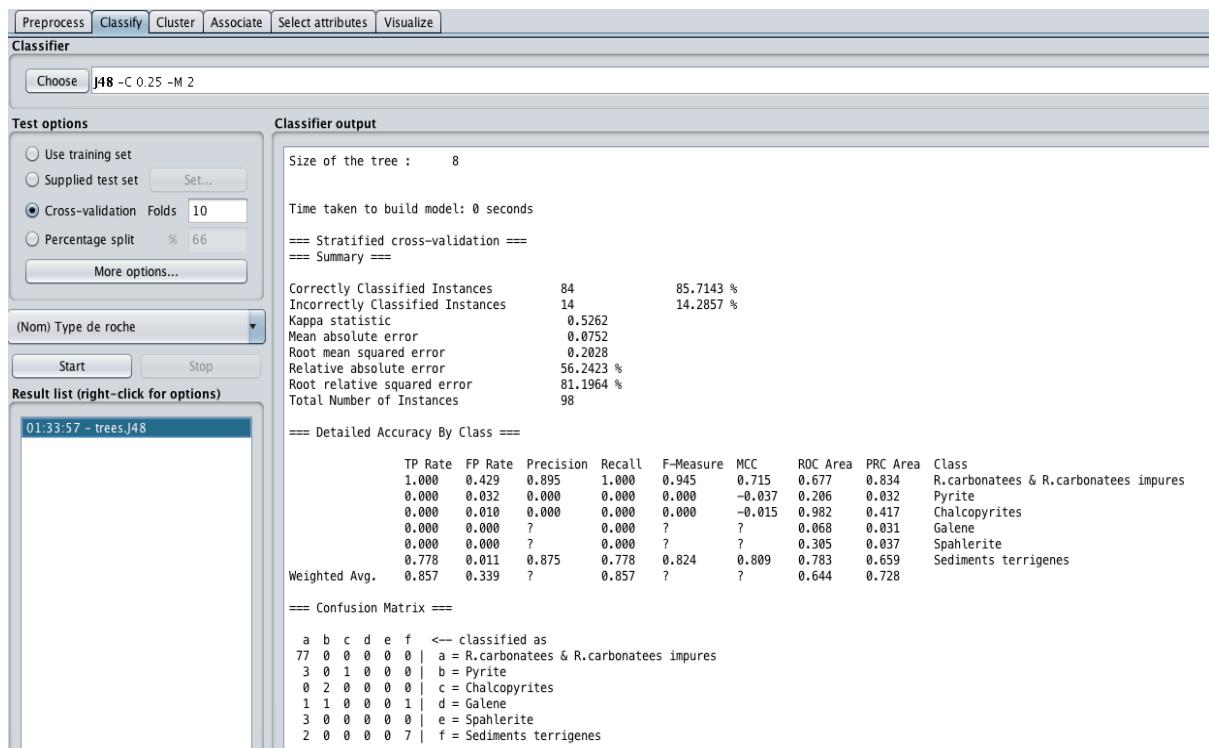
ACCURACY: 88.77%

6.3. EXPERIMENT 3

Experiments 1, 2 (both contrast and non-contrast learning) were repeated for all records with the most important attributes as defined by the expert (S, Zn, Pb, Cu, CaO+MgO, CaO, MgO, Fe₂O₃). This experiment was also repeated for the 3 datasets PD1, PD2 and PD.

a) With PD1- Equal Width Binning:

- Repeating experiment 1:



ACCURACY: 85.71%

- Repeating experiment 2:

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 -C 0.25 -M 2

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

Classifier output

```

CaO+MgO = '(-inf-19.836667]': Not (8.0)
CaO+MgO = '(19.836667-39.193333]' :
|   Cu <= 96: R.carbonates & R.carbonates impures (5.0/1.0)
|   Cu > 96: Not (4.0)
CaO+MgO = '(39.193333-inf)': R.carbonates & R.carbonates impures (81.0/8.0)

Number of Leaves : 4
Size of the tree : 6

Time taken to build model: 0 seconds

```

Result list (right-click for options)

- 01:37:23 - treesJ48
- 01:38:01 - treesJ48
- 01:39:19 - treesJ48

Start Stop

```

== Stratified cross-validation ==
== Summary ==

```

	Correctly Classified Instances	89	90.8163 %
Incorrectly Classified Instances	9	9.1837 %	
Kappa statistic	0.6769		
Mean absolute error	0.1678		
Root mean squared error	0.2952		
Relative absolute error	49.2849 %		
Root relative squared error	71.8818 %		
Total Number of Instances	98		

```

== Detailed Accuracy By Class ==

```

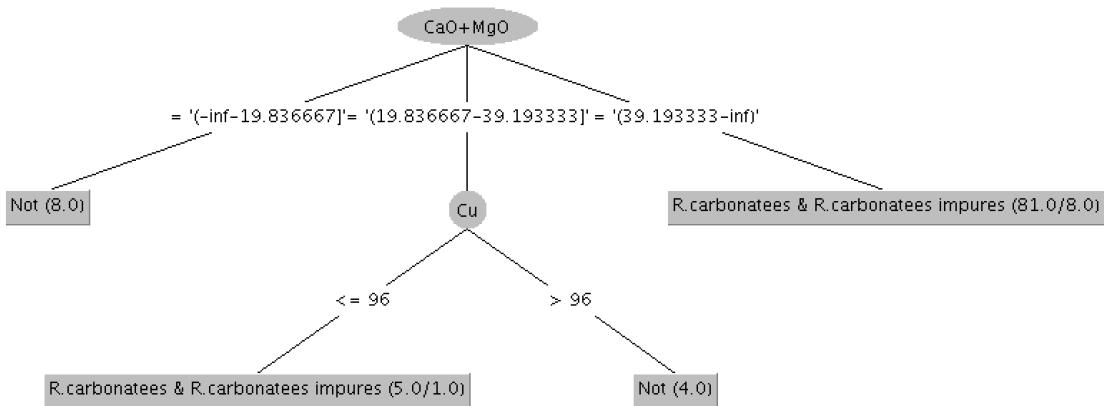
	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
1.000	0.429	0.895	1.000	0.945	0.715	0.684	0.843	0.661	R.carbonates & R.carbonates impures
0.571	0.000	1.000	0.571	0.727	0.715	0.684	0.661	0.661	Not
Weighted Avg.	0.908	0.337	0.918	0.908	0.898	0.715	0.684	0.804	

```

== Confusion Matrix ==

```

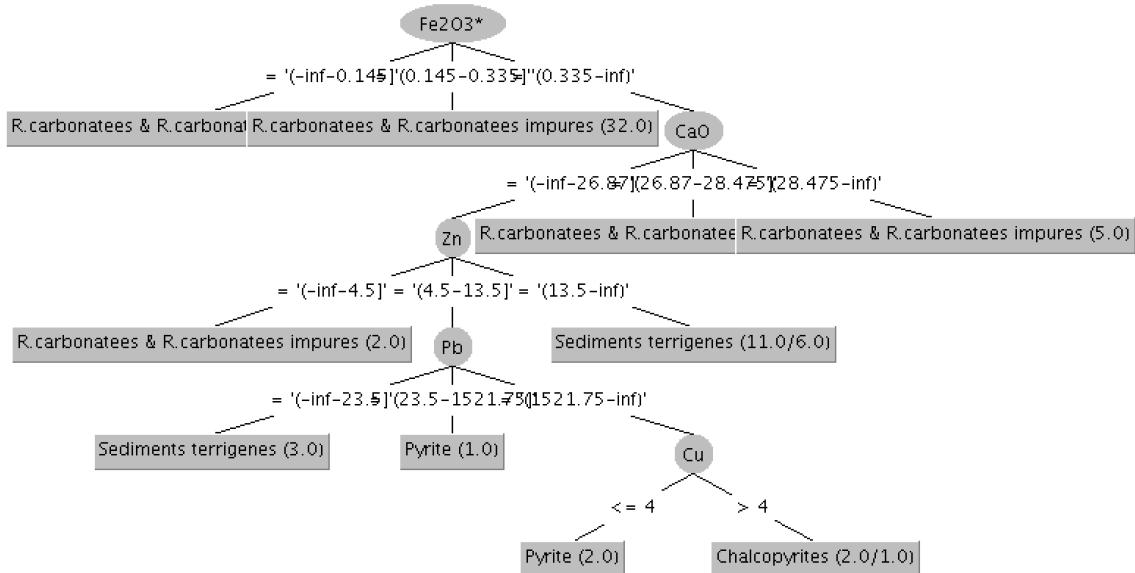
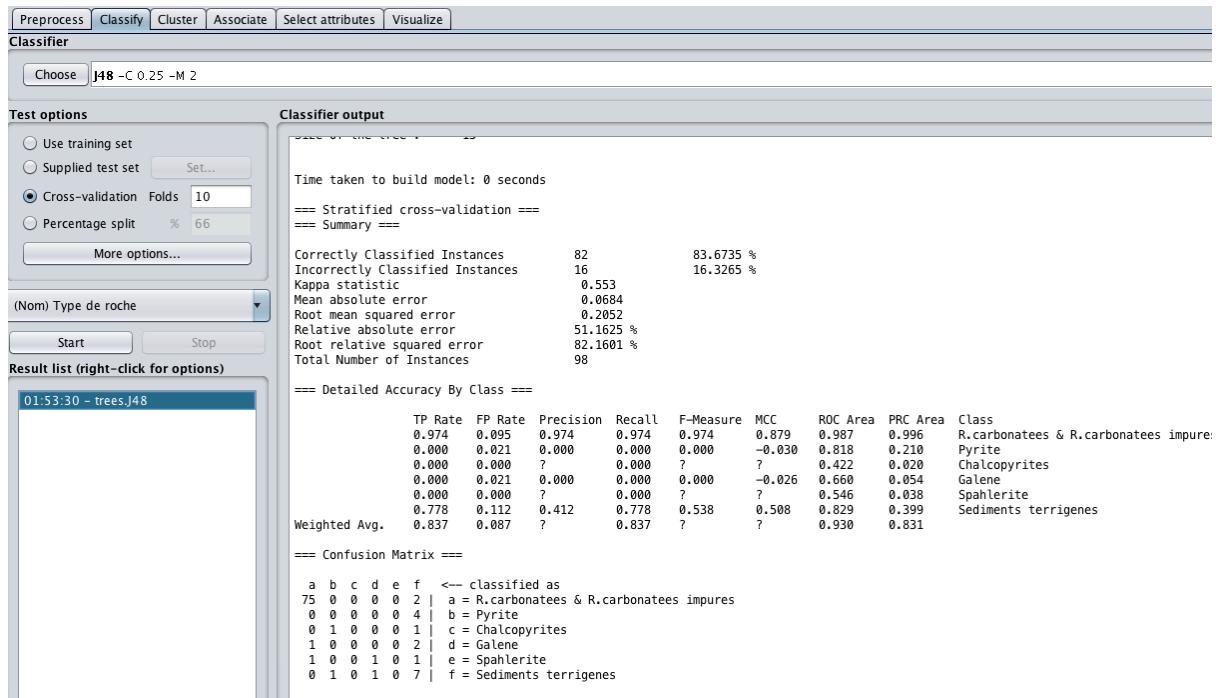
a	b	<-- classified as
77	0	a = R.carbonates & R.carbonates impures
9	12	b = Not



ACCURACY: 90.81%

b) With PD2 - Equal Frequency Binning:

- Repeating experiment 1:



ACCURACY: 83.67%

- Repeating experiment 2:

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier

Choose J48 -C 0.25 -M 2

Test options

- Use training set
- Supplied test set Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

Classifier output

```

Fe203* = '(-inf-0.145]'; R.carbonates & R.carbonates impures (33.0)
Fe203* = '(0.145-0.335]'; R.carbonates & R.carbonates impures (32.0)
Fe203* = '(0.335-inf)'
| CaO = '(-inf-26.87]': Not (21.0/2.0)
| CaO = '(26.87-28.475]': R.carbonates & R.carbonates impures (7.0/2.0)
| CaO = '(28.475-inf)': R.carbonates & R.carbonates impures (5.0)

Number of Leaves : 5
Size of the tree : 7

Time taken to build model: 0 seconds
==== Stratified cross-validation ====
==== Summary ====

```

Correctly Classified Instances	94	95.9184 %
Incorrectly Classified Instances	4	4.0816 %
Kappa statistic	0.8788	
Mean absolute error	0.0556	
Root mean squared error	0.1919	
Relative absolute error	16.3127 %	
Root relative squared error	46.7229 %	
Total Number of Instances	98	

```
==== Detailed Accuracy By Class ====

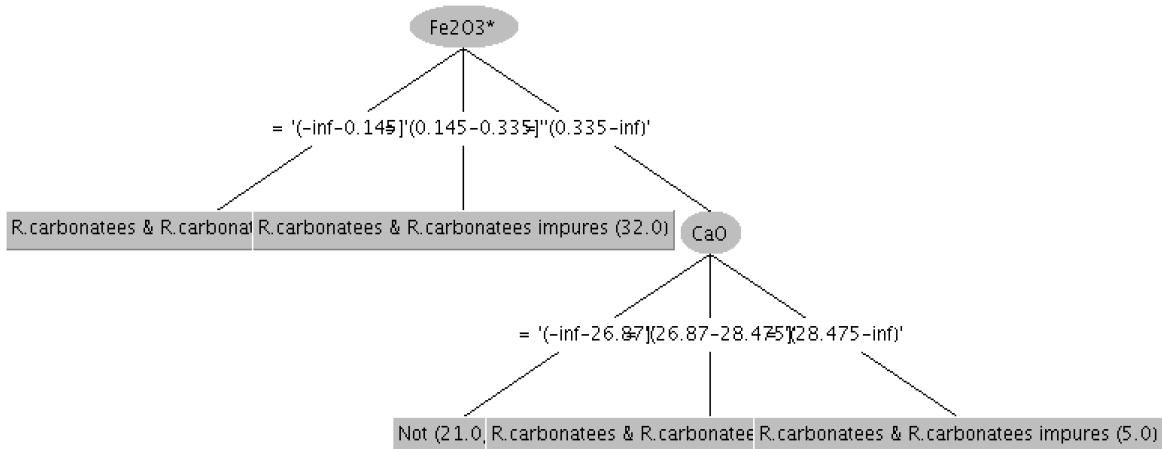
```

TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.974	0.095	0.974	0.974	0.974	0.879	0.988	0.996	R.carbonates & R.carbonates impures
0.905	0.026	0.905	0.905	0.905	0.879	0.988	0.947	Not
Weighted Avg.	0.959	0.080	0.959	0.959	0.879	0.988	0.985	

```
==== Confusion Matrix ====

```

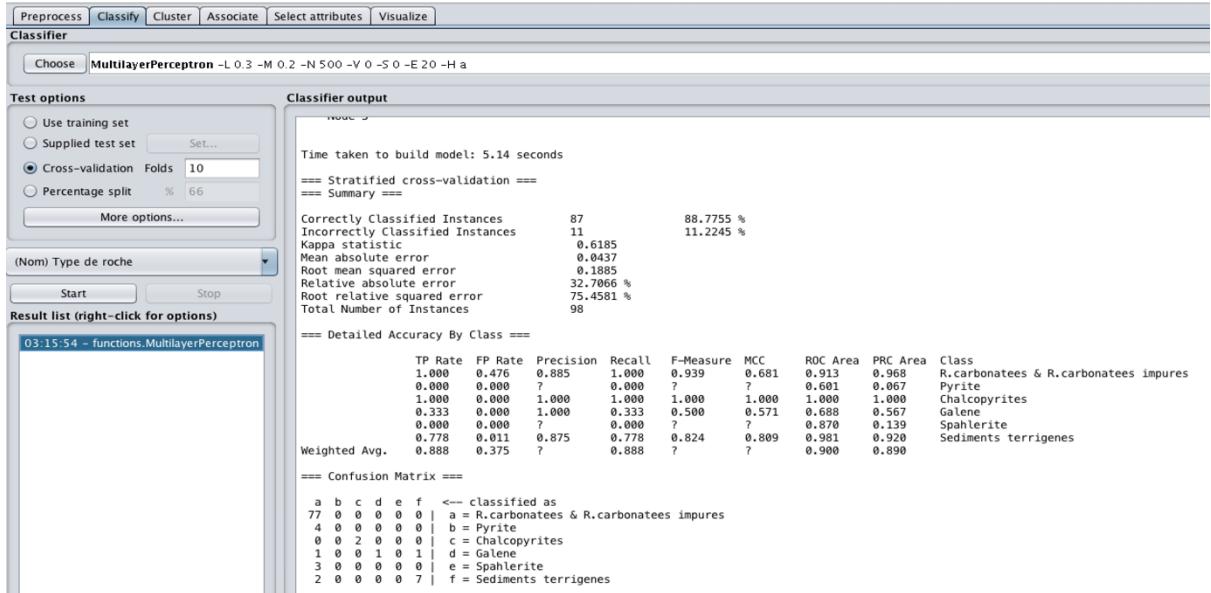
a	b	<-- classified as	
75	2	a = R.carbonates & R.carbonates impures	
2	19	b = Not	



ACCURACY: 95.91%

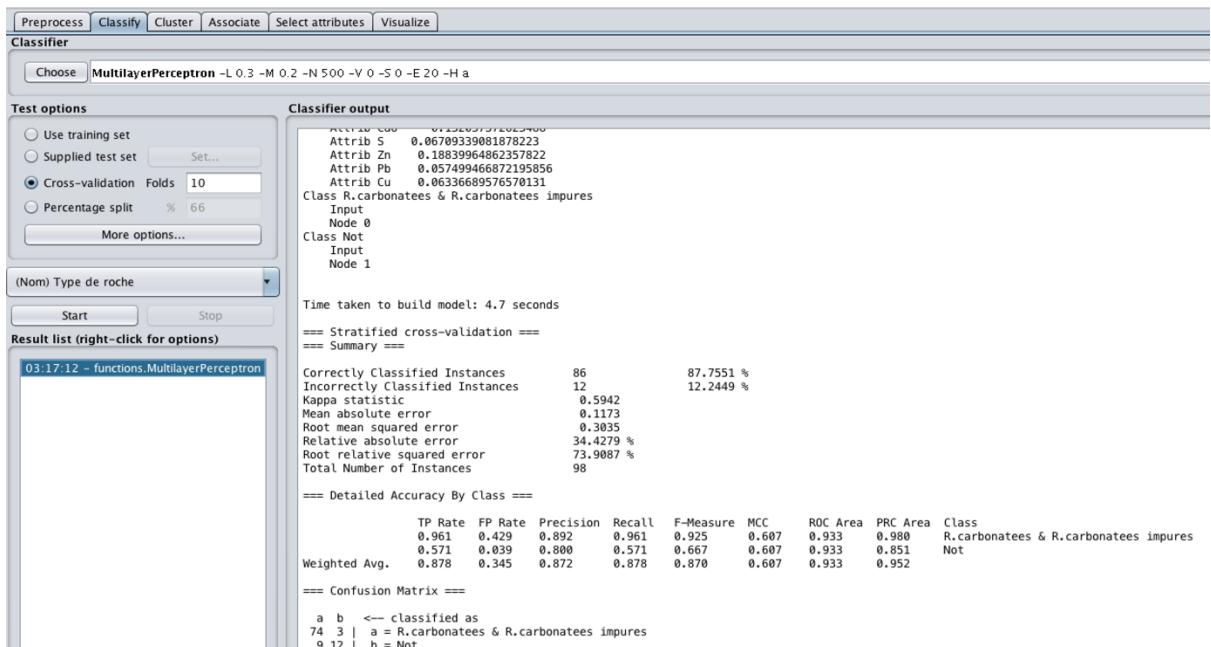
c) With PD - Normalized:

- Repeating experiment 1:



ACCURACY: 88.77%

- Repeating experiment 2:



ACCURACY: 87.75%

7. EVALUATING THE MODELS

A number of machine learning techniques were applied to the given dataset after cleaning.

Each of these methods were also applied on different versions of the same dataset:

PD data represents the original normalized data.

PD1 dataset represents Equal Width Binning Discretized data.

PD2 represents the Equal Frequency Binning Discretized data.

All the models that are implemented are evaluated on the basis of the accuracy as given in the table below.

Experiment →	1	2	3	
	Models ↓	Non Contrast Learning (All Attributes)	Contrast Learning (All Attributes)	Non Contrast Learning (Expert Selected Attributes)
Decision Tree (PD1)	84.69%	86.73%	85.71%	90.81%
Decision Tree (PD2)	84.69%	93.87%	83.67%	95.91%
Neural Network (PD)	88.77%	88.77%	88.77%	87.75%

8. SUMMARY

The aim of the project was to classify the given data on the basis of TYPE DE ROCHE (Rock Type). We used WEKA (Internet based Classification Tool) to build two types of classifiers: Descriptive and Non-Descriptive. For the descriptive classifier, we used Decision tree to generate sets of discriminant rules describing the content of the data. We used Neural Networks to build the Non-Descriptive Classifier.

The Data Preparation step included attributes selection, cleaning the data, filling the missing values, error correction etc. to build the Project Data. The Data Preprocessing step included 2 methods of Data Discretization Namely Equal Width Binning and Equal Frequency Binning to obtain the 2 datasets PD1 and PD2 respectively. Also, the project dataset was normalized to obtain the dataset PD. For each of these datasets 3 experiments were carried out in WEKA. The predictive accuracy was computed for the different classifiers built.

9. REFERENCES

- 1] Outlier Detection in Weka : <https://www.youtube.com/watch?v=WrjpO7CmUoQ>
- 2] Deleting an Attribute: <https://sourceforge.net/projects/weka/files/documentation/>
- 3] Data Cleaning using Weka : <https://machinelearningmastery.com/how-to-handle-missing-values-in-machine-learning-data-with-weka/>
- 4] Data Discretization : Data Discretization: <https://www.youtube.com/watch?v=P--UFzlNGeA&t=109s>
- 5] Equal Frequency Binning: <https://www.youtube.com/watch?v=aDMzPC5IO4c>