

Supervised Learning on Bakary Data Using WEKA

Shariq Syed & Sudhin Domala



Outline

- Classification Tool: WEKA
Waikato Environment for Knowledge Analysis
- Can be downloaded from
<http://www.cs.waikato.ac.nz/~ml/weka/index.html>



Data Preparation: Removing Attributes

- Removing Attributes: We ended up removing attributes which had roughly 33% or more of its data missing but ideally
- On WEKA, we looked at the percentage of data which was missing and removed those from the dataset. These attributes were:
 - Mo: 88% missing
 - Co: 84% missing
 - As: 72% missing
 - Cd: 70% missing
 - Pb: 55% missing
 - Sc: 50% missing
 - Ni: 40% missing
 - Li: 39% missing



Data Preparation: Missing Data

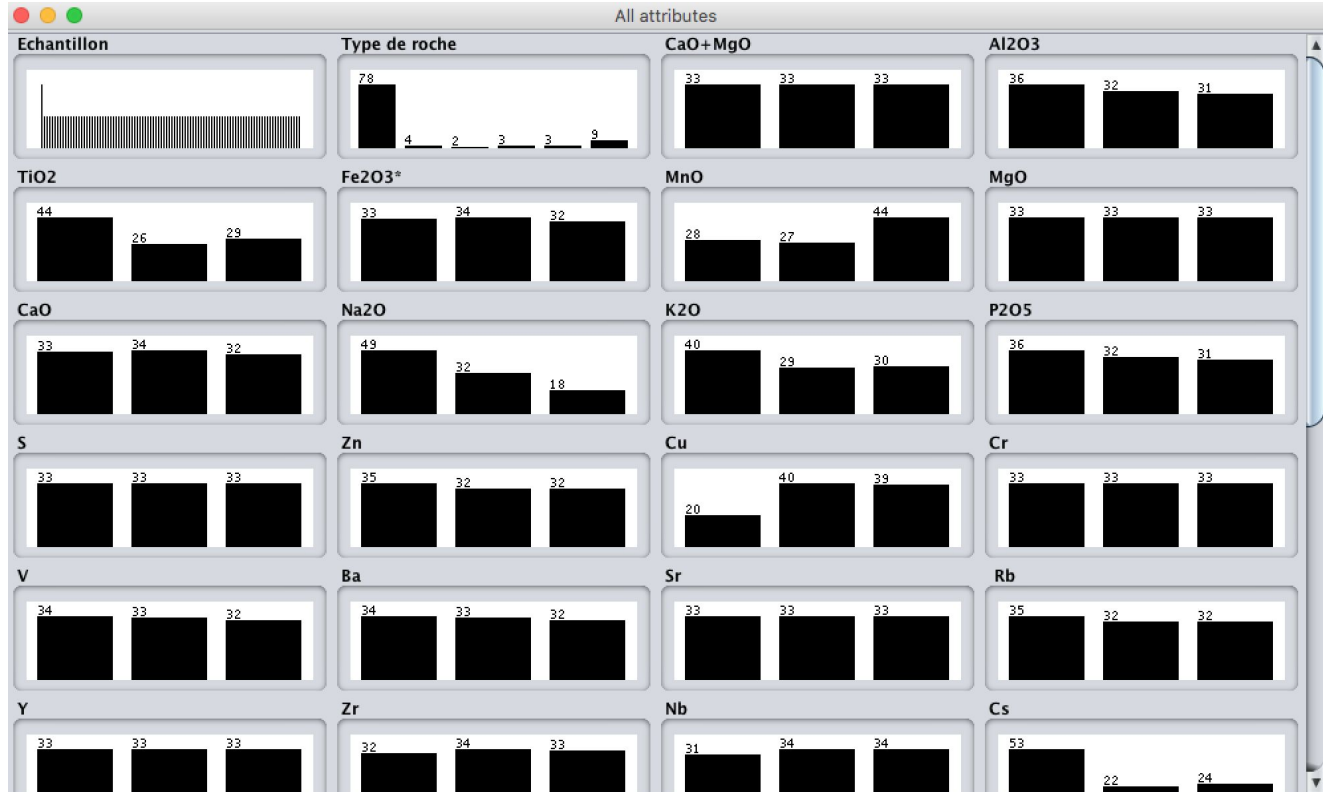
- Many of the attributes contained missing data which can throw off the classifier and so it is imperative that these values be filled in with either 0's or averages given the other classes. To achieve this, we used WEKA's ReplaceMissingValues filter which replaced the missing values with the mean of the numeric values.



Data Preprocessing: Discretization

- In order to discretize the data, we applied another filter on WEKA which is called Discretize. The only modification I made was change the number of bins from 10 to 3.
- Data #1
 - Binning Method(Discrete Bin Count = 3)
 - useEqualFrequency = True
 - Equal Width Bins
- Data #2
 - Binning Method(Discrete Bin Count = 3,)
 - useEqualFrequency = False
 - Equal Depth Bins

Data Preprocessing: Discretization Set



Data Preprocessing: Discretization Set





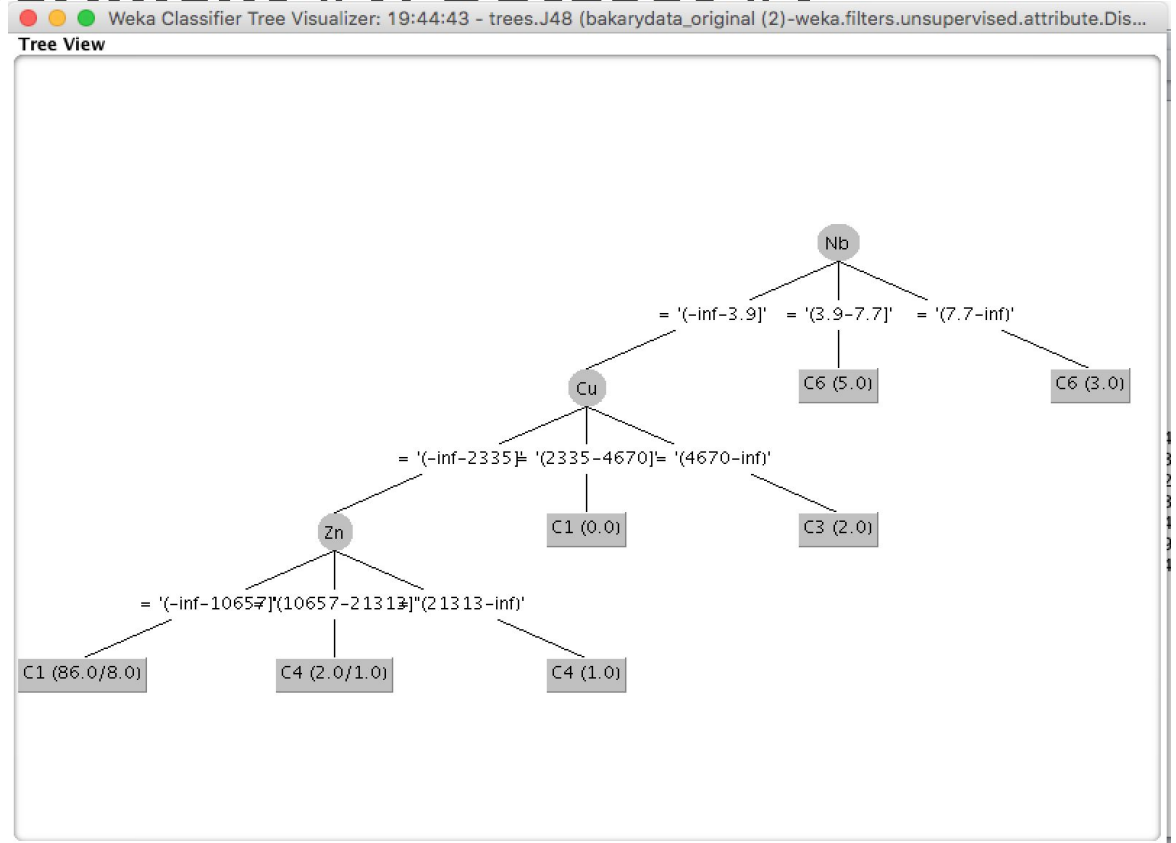
Experiments

- Experiment 1 -> Full Learning
 - Use all attributes to classify all classes (C1-C6)
 - Done by Discretization and usage of J48 Algorithm
- Experiment 2 -> Contrast Learning
 - Using all attributes to compare class C1 with the rest of the classes
 - Done by Discretization and Renaming Nominal Values on Type de roche, and J48
- Experiment 3 -> Limited Learning
 - Construction of decision tree using only the major attributes
 - Removal of non-major attributes in preprocessing and prior to experiment 1&2 implementations



Experiment #1. Dataset #1

- Decision tree using J48 algorithm
- Used k-fold(k=10)
- Predictive Accuracy: 84.8485%





Discriminant Rules

- IF Nb = “(-inf-3.9)” AND Cu = “(-inf-2335] – AND Zn = “(-inf-10657)”

THEN class = “C1”

- IF Nb = “(-inf-3.9)” AND Cu = “(-inf-2335] – AND Zn = “(10657-21313)”

THEN class = “C4”

- IF Nb = “(-inf-3.9)” AND Cu = “(-inf-2335] – AND Zn = “(21313-inf)”

THEN class = “C4”



Discriminant Rules(cont.)

- IF Nb = “(-inf-3.9]” AND Cu = “(2335-4670)”

THEN class = “C1”

- IF Nb = “(-inf-3.9]” AND Cu = “(4670-inf)”

THEN class = “C3”

- IF Nb = “(3.9-7.7]”

THEN class = “C6”

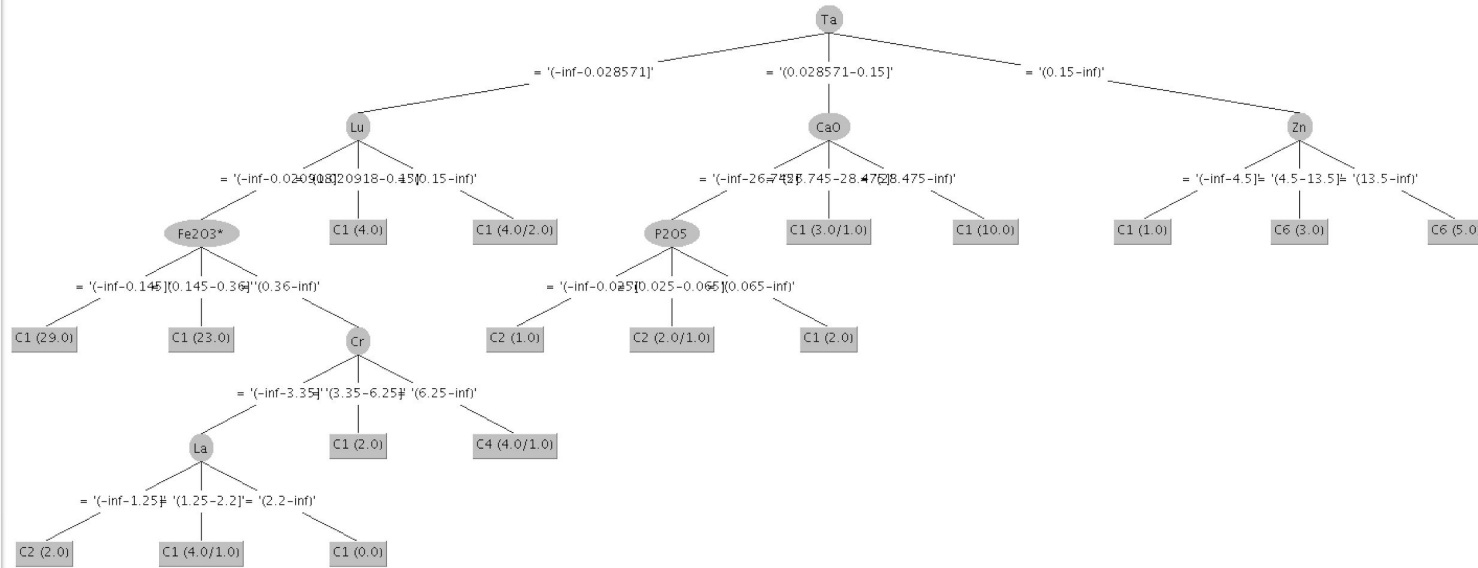
- IF Nb = “(7.7-inf)”

THEN class = “C6”

Experiment #1. Dataset #2

Weka Classifier Tree Visualizer: 21:42:58 - trees.J48 (bakarydata_original (1)-weka.filters.unsupervised.attribute.Discretize-F-B3-M-1.0-Rfirst-last-precision6-unset-class-temporarily)

Tree View



- Used Training Set method for PA calculation
- Predictive Accuracy: 75.7576%



Discriminant Rules

- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe2O3="(-inf-0.145]"
THEN Class="C1"
- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe2O3="(0.145-0.36]"
THEN Class="C1"
- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe2O3="(0.145-0.36]" AND
Cr="(-inf-3.35]" AND La="(-inf-1.25]"
THEN Class="C2"



Discriminant Rules(cont.)

- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe₂O₃="(0.145-0.36]" AND Cr="(-inf-3.35]" AND La="(1.25-2.2]"

THEN Class="C1"

- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe₂O₃="(0.145-0.36]" AND Cr="(-inf-3.35]" AND La="(2.2-inf]"

THEN Class="C1"

- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe₂O₃="(0.145-0.36]" AND Cr="(3.35,6.25]"

THEN Class="C1"



Discriminant Rules(cont.)

- IF Ta="(-inf-0.028571]" AND Lu="(-inf-0.020918]" AND Fe₂O₃="(0.145-0.36]" AND Cr="(6.25-inf]"

THEN Class="C4"

- IF Ta="(-inf-0.028571]" AND Lu="(0.020918-0.15]"

THEN Class="C1"

- IF Ta="(-inf-0.028571]" AND Lu="(0.15-inf)"

THEN Class="C1"



Discriminant Rules(cont.)

- IF $Ta = "(0.028571-0.15]"$ AND $CaO = "(-inf-26.745]"$ AND $P2O5 = "(-inf-0.025]"$
THEN Class="C2"
- IF $Ta = "(0.028571-0.15]"$ AND $CaO = "(-inf-26.745]"$ AND $P2O5 = "(0.025-0.065]"$
THEN Class="C2"
- IF $Ta = "(0.028571-0.15]"$ AND $CaO = "(-inf-26.745]"$ AND $P2O5 = "(0.065-inf]"$
THEN Class="C1"
- IF $Ta = "(0.028571-0.15]"$ AND $CaO = "(26.745-28.475]"$
THEN Class="C3"



Discriminant Rules(cont.)

- IF Ta="(0.028571-0.15]" AND CaO="(28.475-inf]"

THEN Class="C1"

- IF Ta="(0.15-inf]" AND Zn="(inf-4.5]"

THEN Class="C1"

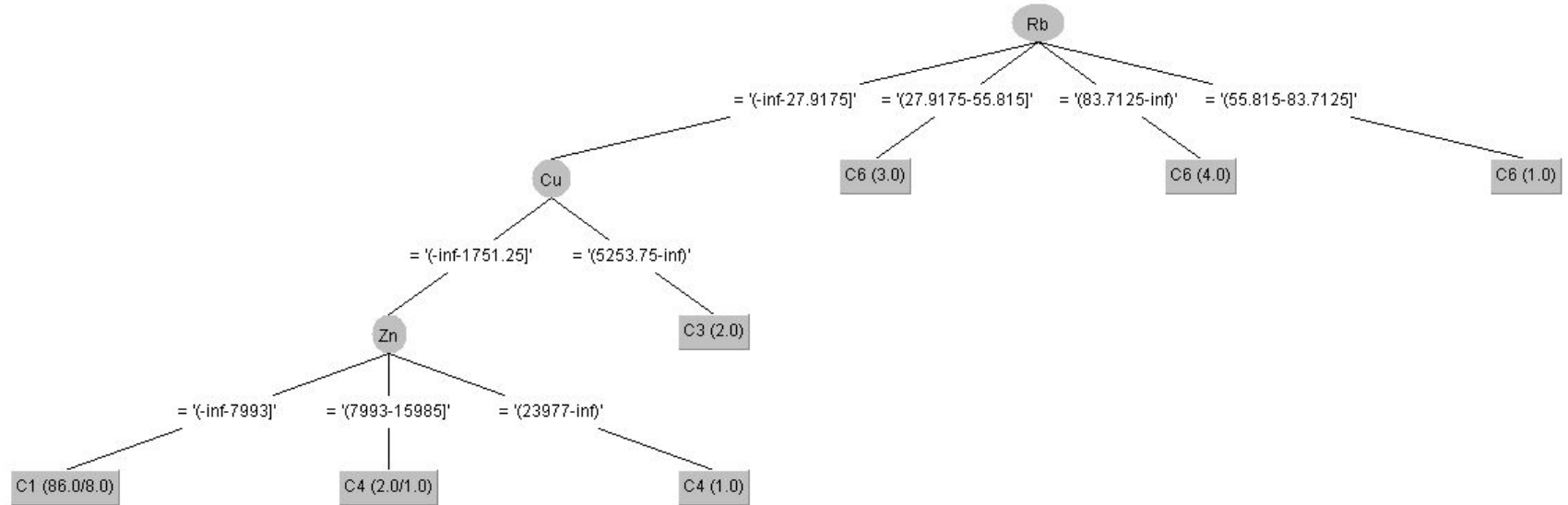
- IF Ta="(0.15-inf]" AND Zn="(4.5-13.5]"

THEN Class="C6"

- IF Ta="(0.15-inf]" AND Zn="(13.5-inf]"

THEN Class="C6"

Experiment #2, Data Set #1



- Predictive Accuracy: 83.8384%



Discriminant Rules

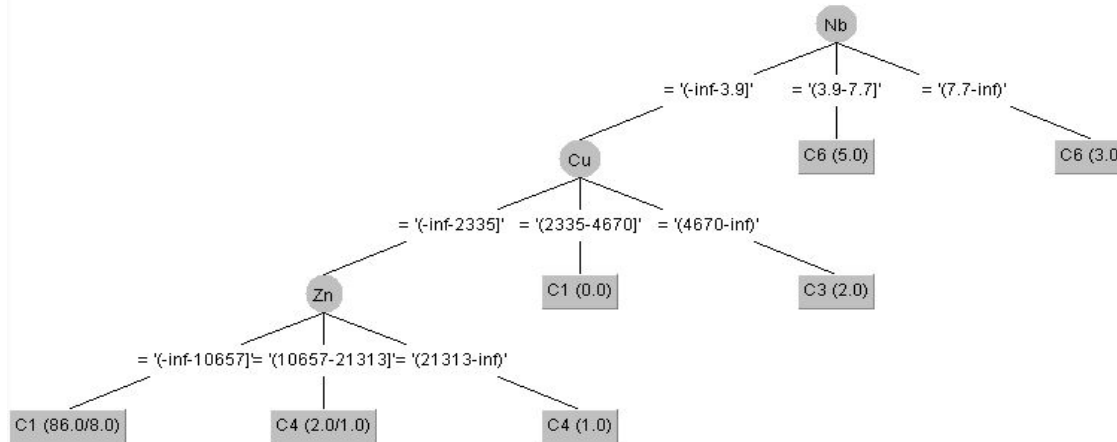
- IF Rb= $(-\infty-27.9175]$ AND Cu= $(-\infty-1751.25]$ AND Zn= $(-\infty-7993]$
THEN Class="C1"
- IF Rb= $(-\infty-27.9175]$ AND Cu= $(-\infty-1751.25]$ AND Zn= $(7993-15985]$
THEN Class="C4"
- IF Rb= $(-\infty-27.9175]$ AND Cu= $(-\infty-1751.25]$ AND Zn= $(23977-\infty]$
THEN Class="C4"
- IF Rb= $(-\infty-27.9175]$ AND Cu= $(5253.75-\infty]$
THEN Class="C3"



Discriminant Rules(cont.)

- IF Rb="(27.9175-55.815]"
THEN Class="C6"
- IF Rb="(55.815-83.7125]"
THEN Class="C6"
- IF Rb="(83.7125-inf]"
THEN Class="C6"

Experiment #2, Data Set #2



- Predictive Accuracy:
84.8485%



Discriminant Rules

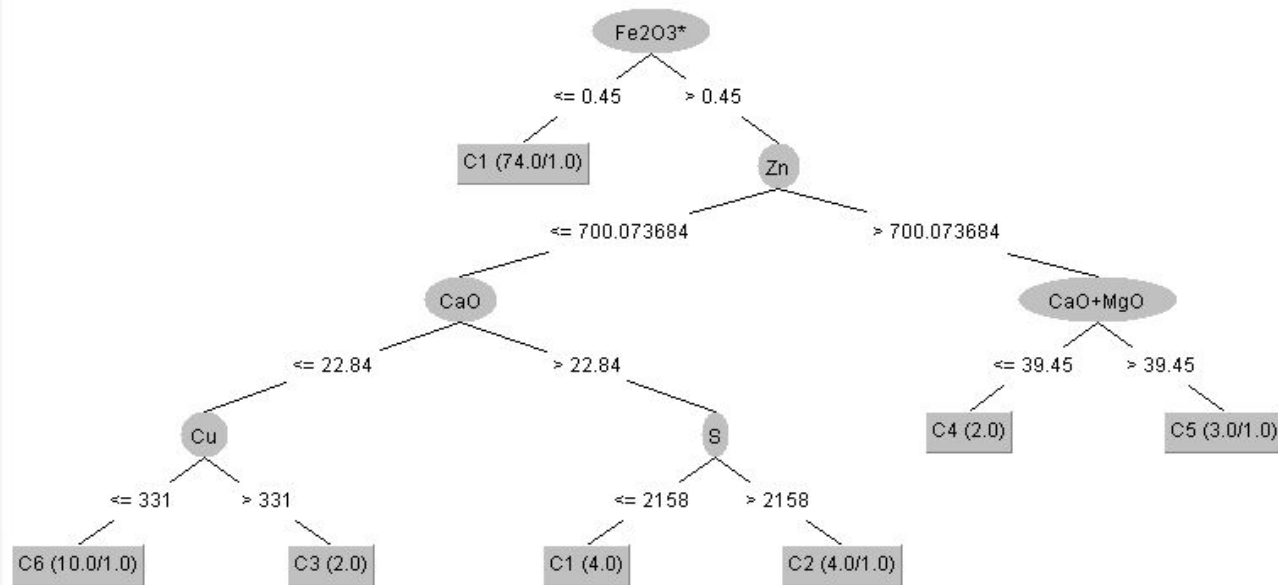
- IF Nb="(-inf-3.9]" AND Cu="(-inf-2335] AND Zn="(-inf-10657]"
THEN Class="C1"
- IF Nb="(-inf-3.9]" AND Cu="(-inf-2335] AND Zn="(10657-21313]"
THEN Class="C4"
- IF Nb="(-inf-3.9]" AND Cu="(-inf-2335] AND Zn="(21313-inf]"
THEN Class="C4"
- IF Nb="(-inf-3.9]" AND Cu="(2335-4670]"
THEN Class="C1"



Discriminant Rules(cont.)

- IF Nb="(-inf-3.9]" AND Cu="(4670-inf]
THEN Class="C3"
- IF Nb="(3.9-7.7]"
THEN Class="C6"
- IF Nb="(7.7-inf]"
THEN Class="C3"

Experiment #3, Data Set #1



- Predictive Accuracy: 81.8182%



Discriminant Rules

- IF $\text{Fe}_2\text{O}_3 = "(\leq .45]"$
THEN Class="C1"
- IF $\text{Fe}_2\text{O}_3 = "(> .45]"$ AND $\text{Zn} = "<= 700.073684"$ AND $\text{CaO} = "<= 22.84"$ AND $\text{Cu} = "<= 331"$
THEN Class="C6"
- IF $\text{Fe}_2\text{O}_3 = "(> .45]"$ AND $\text{Zn} = "<= 700.073684"$ AND $\text{CaO} = "<= 22.84"$ AND $\text{Cu} = "> 331"$
THEN Class="C3"



Discriminant Rules(cont.)

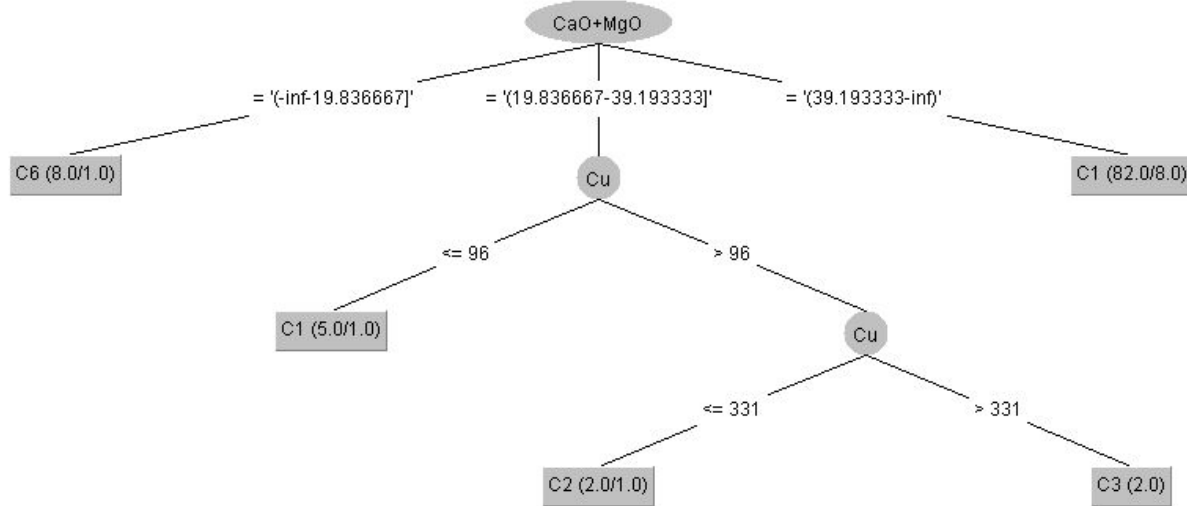
- IF $\text{Fe}_2\text{O}_3 = (>.45]$ AND $\text{Zn} = \leq 700.073684$ AND $\text{CaO} = > 22.84$ AND $\text{S} = \leq 2158$
THEN Class="C1"
- IF $\text{Fe}_2\text{O}_3 = (>.45]$ AND $\text{Zn} = \leq 700.073684$ AND $\text{CaO} = > 22.84$ AND $\text{S} = > 2158$
THEN Class="C2"
- IF $\text{Fe}_2\text{O}_3 = (>.45]$ AND $\text{Zn} = > 700.073684$ AND $\text{CaO} + \text{MgO} = \leq 39.45$
THEN Class="C4"



Discriminant Rules(cont.)

- IF $\text{Fe}_2\text{O}_3 = (>.45]$ AND $\text{Zn} = >700.073684$ AND $\text{CaO} + \text{MgO} = >39.45$
THEN Class="C3"

Experiment #3, Data Set #2



- Predictive Accuracy: 84.8485%



Discriminant Rules

- IF $\text{CaO} + \text{MgO} = "(-\infty - 19.836667]"$
THEN Class="C6"
- IF $\text{CaO} + \text{MgO} = "(19.836667 - 39.193333]"$ AND $\text{Cu} = "<=96"$
THEN Class="C1"
- IF $\text{CaO} + \text{MgO} = "(19.836667 - 39.193333]"$ AND $\text{Cu} = ">96"$ AND $\text{Cu} = "<=331"$
THEN Class="C2"
- IF $\text{CaO} + \text{MgO} = "(19.836667 - 39.193333]"$ AND $\text{Cu} = ">96"$ AND $\text{Cu} = ">331"$
THEN Class="C3"



Discriminant Rules(cont.)

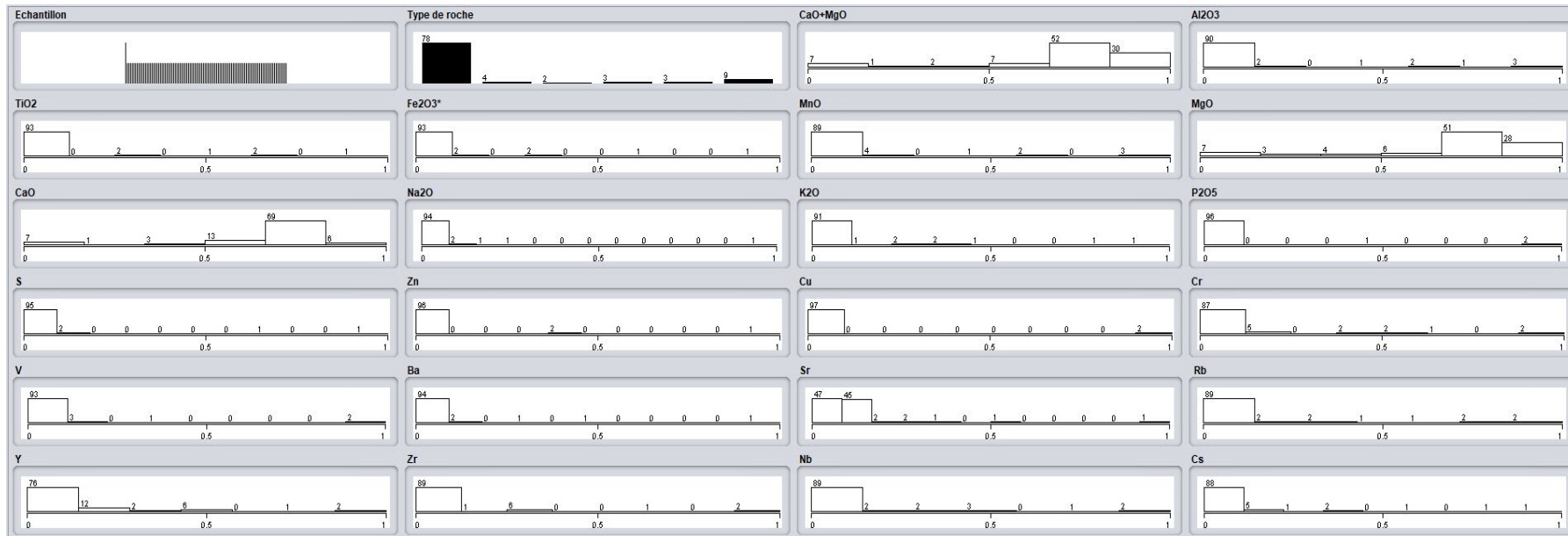
- IF $\text{CaO} + \text{MgO} = (39.193333 - \text{inf}]$
THEN Class="C1"



NN Data Preprocessing: Normalization

- In order to normalize the data, we applied another filter on WEKA which is called Normalize.
 - No modifications and no selections have been made.
 - Scale: 1.0
 - Translation: 0.0
- For experiment 2 (and thus experiment 3), we used “RenameNomial Values” to distinguish the classes on “Type de roche”
- Neuron networks display the input nodes in green, connections in grey lines, red nodes for hidden layers, and yellow boxes for classes.

NN Data Preprocessing: Normalization

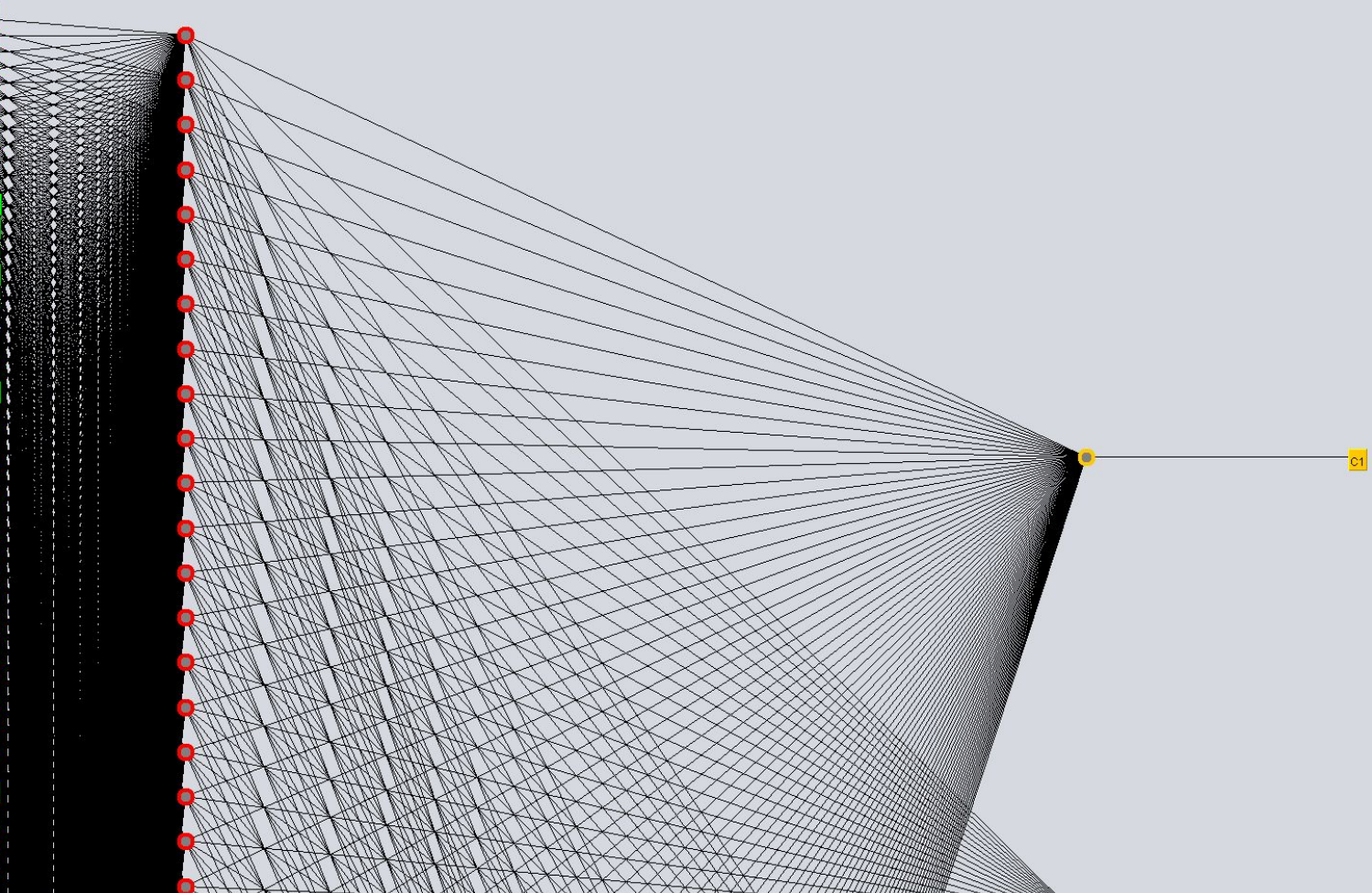




NN Set: Experiment #1

- Predictive Accuracy: 86.8687%
- To build neuron network we used: “MultilayerPerceptron” function in “Classify” section
- Had GUI set to true and Cross-Validation set of Folds of $k=10$

Echantillon=4-02BD304.1.3A1
Echantillon=17-02BD316.1.3A
Echantillon=37-02BD333.2.3B
Echantillon=39-02BD333.4.3D
Echantillon=48-02BD341.1.3B
Echantillon=47-02BD340.1.3B
Echantillon=49-02RM601.1.3A
Echantillon=85-02RM626.2.3A
Echantillon=109-02RM637.1.3C
Echantillon=114-02YH012.1.3C
Echantillon=128-02YH022.1.3A
Echantillon=135-02YH026.2.3C
Echantillon=138-02YH029.1.3A
Echantillon=3-02BD302.2.3A
Echantillon=19-02BD320.1.3B
Echantillon=21-02BD327.1.3A
Echantillon=28-02BD329.4.3D1
Echantillon=32-02BD331.1.3F
Echantillon=33-02BD331.2.3E
Echantillon=35-02BD332.1.3A
Echantillon=42-02BD337.1.3C
Echantillon=52-02RM601.4.3A
Echantillon=54-02RM602.2.3B
Echantillon=55-02RM602.3.3D
Echantillon=58-02RM602.6.3C
Echantillon=59-02RM602.7.3A
Echantillon=62-02RM605.2.5C
Echantillon=65-02RM605.5.3A
Echantillon=67-02RM605.7.3A
Echantillon=72-02RM612.1.3B
Echantillon=73-02RM613.1.3A
Echantillon=74-02RM615.1.3A
Echantillon=127-02YH020
Echantillon=132-02YH026.1.3A
Echantillon=16-02BD314.1.3A
Echantillon=18-02BD317.1.3A
Echantillon=22-02BD328.1.3C
Echantillon=125-02YH019



c1

Controls

Start

Epoch 0

Accept

Num Of Epochs 500

Error per Epoch = 0



NN Set: Experiment #2

- Predictive Accuracy: 86.8687%
- To build neuron network we used: “MultilayerPerceptron” with “Rename NominalValues” used prior to making NN.
- Had GUI set to true and Cross-Validation set of Folds of $k=10$
- Epoch set to 500 and implemented 10 fold (“start” X 10 and “accept” X 10)

Echantillon=4-02BD304 1.3A1
Echantillon=17-02BD316 1.3A
Echantillon=37-02BD333 2.3B
Echantillon=39-02BD333 4.3D
Echantillon=48-02BD341 1.3B
Echantillon=47-02BD340 1.3B
Echantillon=49-02RM601 1.3A
Echantillon=85-02RM625 2.3A
Echantillon=109-02RM637 1.3C
Echantillon=114-02YH012 1.3C
Echantillon=128-02YH022 1.3A
Echantillon=135-02YH026 2.3C
Echantillon=138-02YH029 1.3A
Echantillon=3-02BD302 2.3A
Echantillon=18-02BD320 1.3B
Echantillon=21-02BD327 1.3A
Echantillon=28-02BD329 4.3D1
Echantillon=32-02BD331 1.3F
Echantillon=33-02BD331 2.3E
Echantillon=35-02BD332 1.3A
Echantillon=42-02BD337 1.3C
Echantillon=52-02RM601 4.3A
Echantillon=54-02RM602 2.3B
Echantillon=55-02RM602 3.3D
Echantillon=58-02RM602 6.3C
Echantillon=59-02RM602 7.3A
Echantillon=62-02RM605 2.5C
Echantillon=65-02RM605 5.3A
Echantillon=67-02RM605 7.3A
Echantillon=72-02RM612 1.3B
Echantillon=73-02RM613 1.3A
Echantillon=74-02RM615 1.3A
Echantillon=127-02YH020
Echantillon=132-02YH025 1.3A
Echantillon=16-02BD314 1.3A
Echantillon=18-02BD317 1.3A
Echantillon=22-02BD328 1.3C
Echantillon=125-02YH018

Controls

Start

Epoch 500

Num Of Epochs 500

Accept

Error per Epoch = 0.0000287

C1



NN Set: Experiment #3

- Predictive Accuracy: 87.8788%
- Same parameters used like experiment 1 and 2 with just important elements
- Epoch set to 500 and implemented 10 fold (“start” X 10 and “accept” X 10)

Echantillon=4-02BD304.1.3A1
Echantillon=17-02BD316.1.3A
Echantillon=37-02BD333.2.3B
Echantillon=39-02BD333.4.3D
Echantillon=48-02BD341.1.3B
Echantillon=47-02BD340.1.3B
Echantillon=49-02RM601.1.3A
Echantillon=85-02RM625.2.3A
Echantillon=109-02RM637.1.3C
Echantillon=114-02YH012.1.3C
Echantillon=128-02YH022.1.3A
Echantillon=135-02YH026.2.3C
Echantillon=139-02YH029.1.3A
Echantillon=3-02BD302.2.3A
Echantillon=19-02BD320.1.3B
Echantillon=21-02BD327.1.3A
Echantillon=28-02BD329.4.3D1
Echantillon=32-02BD331.1.3F
Echantillon=33-02BD331.2.3E
Echantillon=35-02BD332.1.3A
Echantillon=42-02BD337.1.3C
Echantillon=52-02RM601.4.3A
Echantillon=54-02RM602.2.3B
Echantillon=55-02RM602.3.3D
Echantillon=58-02RM602.6.3C
Echantillon=59-02RM602.7.3A
Echantillon=62-02RM605.2.5C
Echantillon=65-02RM605.5.3A
Echantillon=67-02RM605.7.3A
Echantillon=72-02RM612.1.3B
Echantillon=73-02RM613.1.3A
Echantillon=74-02RM615.1.3A
Echantillon=127-02YH020
Echantillon=132-02YH025.1.3A
Echantillon=16-02BD314.1.3A
Echantillon=18-02BD317.1.3A
Echantillon=22-02BD328.1.3C
Echantillon=125-02YH018

Controls

Start

Epoch 0

Accept

Num Of Epochs 500

Error per Epoch = 0

C1

C2



Summary of Predictive Accuracy

	Data #1	Data #2
Experiment #1	84.8485%	75.7576%
Experiment #2	83.8384%	84.8485%
Experiment #3	81.8182%	84.8485%



Summary of Predictive Accuracy

	Data Set
Experiment #1	86.8687%
Experiment #2	86.8687%
Experiment #3	87.8788%



Analysis

- Dataset #1 was carried out with equal frequency bins
- Dataset #2 was carried out with equal width bins
- The higher accuracy can be misleading due to the high volume of data however in the decision tree, dataset #2 overall had a higher accuracy
- WEKA produces different rules depending on the methods used for data preparation.
- Decision Trees: Depending on the sets used, different experiments yielded higher accuracies. (Experiment 1 for Data Set #1 and Experiment 2&3 for Data Set #2)
- Neural Network: Experiment 3 regarding a select set of attributes yielded the highest number for accuracy