

**AMERICAN INTERNATIONAL
UNIVERSITY-
BANGLADESH**



Faculty of Science and Technology

Assignment Cover Page

Assignment Title:	DATA WAREHOUSING AND DATA MINING FINAL PROJECT			
Assignment No:	1	Date of Submission:		20 April, 2022
Course Title:	DATA WAREHOUSING AND DATA MINING			
Course Code:	Click here to enter text.	Section:	A	
Semester:	Spring	2019-20	Course Teacher:	AKINUL ISLAM JONY

Declaration and Statement of Authorship:

1. I/we hold a copy of this Assignment/Case-Study, which can be produced if the original is lost/damaged.
2. This Assignment/Case-Study is my/our original work and no part of it has been copied from any other student's work or from any other source except where due acknowledgement is made.
3. No part of this Assignment/Case-Study has been written for me/us by any other person except where such collaboration has been authorized by the concerned teacher and is clearly acknowledged in the assignment.
4. I/we have not previously submitted or currently submitting this work for any other course/unit.
5. This work may be reproduced, communicated, compared and archived for the purpose of detecting plagiarism.
6. I/we give permission for a copy of my/our marked work to be retained by the Faculty for review and comparison, including review by external examiners.
7. I/we understand that Plagiarism is the presentation of the work, idea or creation of another person as though it is your own. It is a form of cheating and is a very serious academic offence that may lead to expulsion from the University. Plagiarized material can be drawn from, and presented in, written, graphic and visual form, including electronic data, and oral presentations. Plagiarism occurs when the origin of them arterial used is not appropriately cited.
8. I/we also understand that enabling plagiarism is the act of assisting or allowing another person to plagiarize or to copy my/our work.

* Student(s) must complete all details except the faculty use part.

** Please submit all assignments to your course teacher or the office of the concerned teacher.

Group Name/No. : 1

No	Name	ID	Program	Signature
----	------	----	---------	-----------

1 MD SAKIB HOSSAIN SHOVON 19-39526-1 BSc [CSE] SHOVON

Faculty use only

FACULTY COMMENTS	Marks Obtained	
	Total Marks	

Assignment/Case-Study Cover; © AIUB-2020

Introduction:

The main objective of this project is to analyzing data and to find out hidden pattern of the data, discovering new findings by applying Data Mining techniques. Here we have used several Data Mining techniques such as K-Nearest Neighbours, Decision Trees, Naïve Bayes, Random Forest, Support Vector Machine, and Neural Network. We have collected the “Drug200” dataset from kaggle. We are using weka tool to analyse the data.

Dataset:

The dataset we have taken is Drug200. We have 6 attributes in this dataset. These are Age, Sex, BP, Cholesterol, Na_to_K and Drug. Drug is the target attributes. There are 200 instances. There are no missing values in our dataset. All the values are in numeric form.

Age attributes: Age is numerical. Age of the patient is 15 to 74.

Sex attributes: Gender of the patients is male and female. Male percentage is 52% and Female percentage is 48%.

BP attributes: 39% patient blood pressure level is high, 32% patient blood pressure level is low and 30% patient blood pressure level is normal.

Cholesterol attributes: 52% patient cholesterol is high and 49% patient cholesterol is normal.

Na_to_K: Sodium to potassium Ration in Blood varies from 6.27 to 38.2

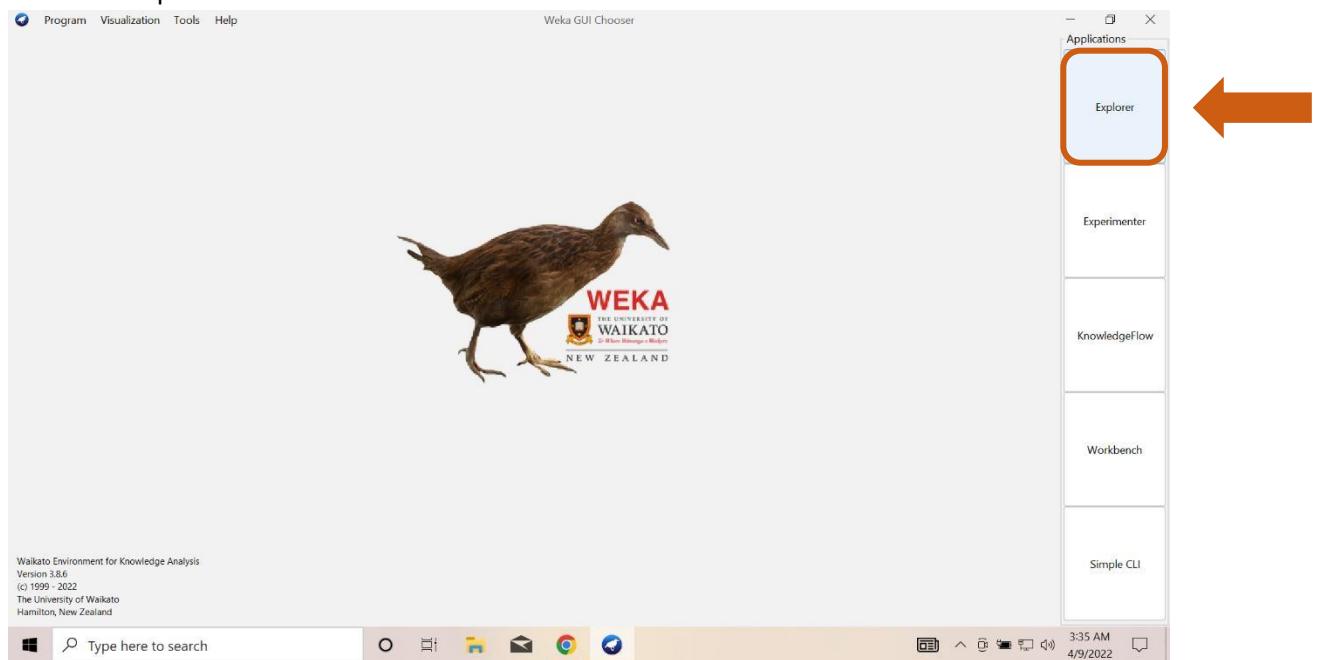
Drug: It's the target attribute that we want to classify. Here drag X is 27%, drag Y is 46% and other minor drags are 28% all together.

URL:

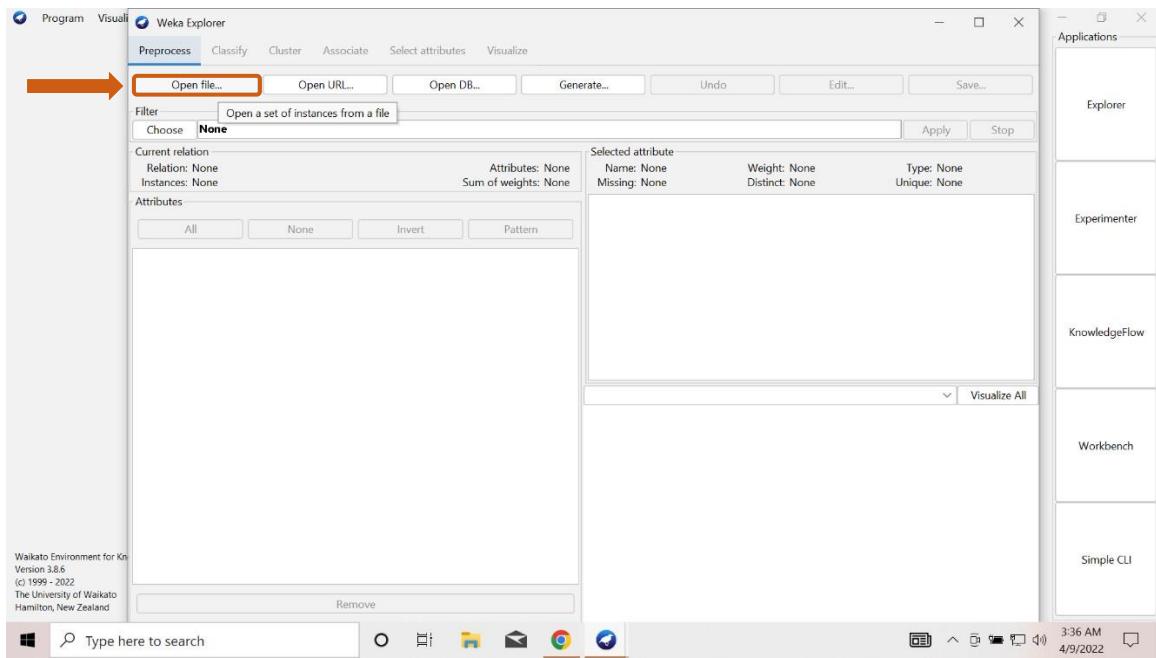
<https://www.kaggle.com/code/ashwinshetgaonkar/drug-classification-random-forest-100-acc/data>

Model Development:

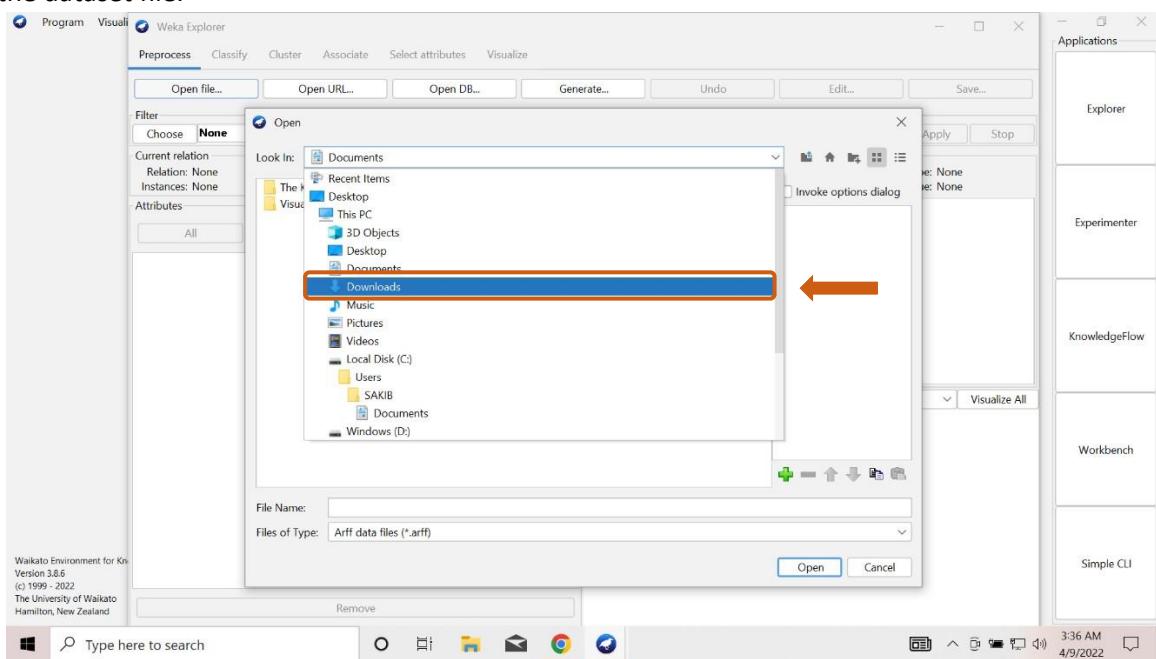
Weka was used in this project. After opening the software, the interface will look like this. We have to click on "Explorer".



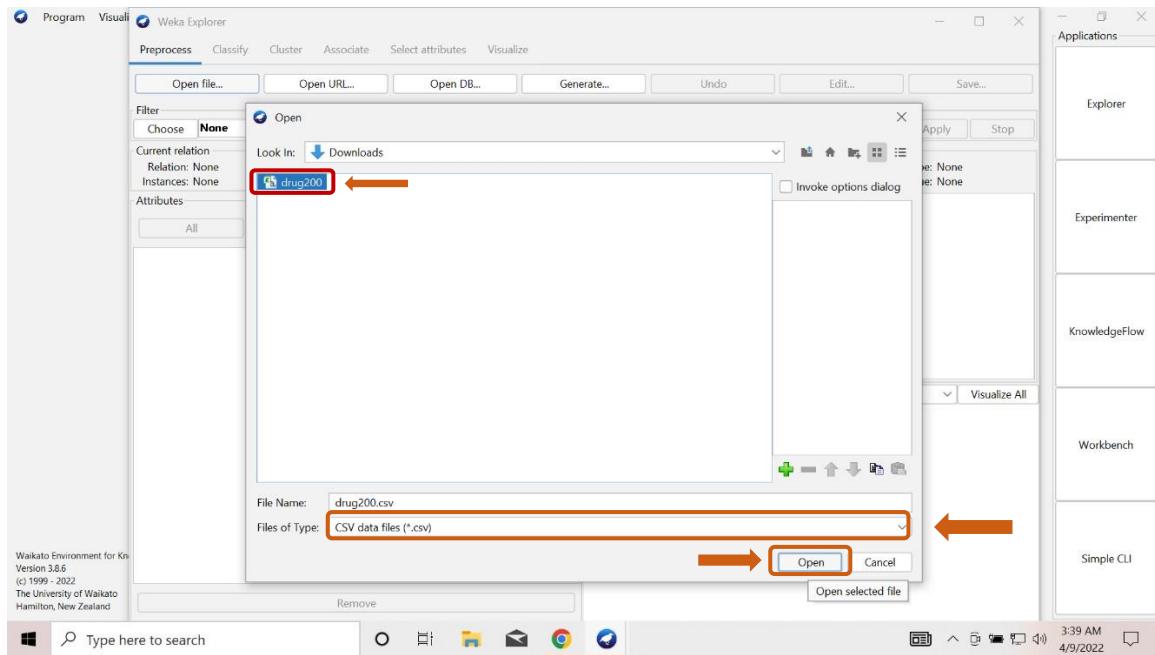
Then this interface will appear and we have to choose "Open file" option.



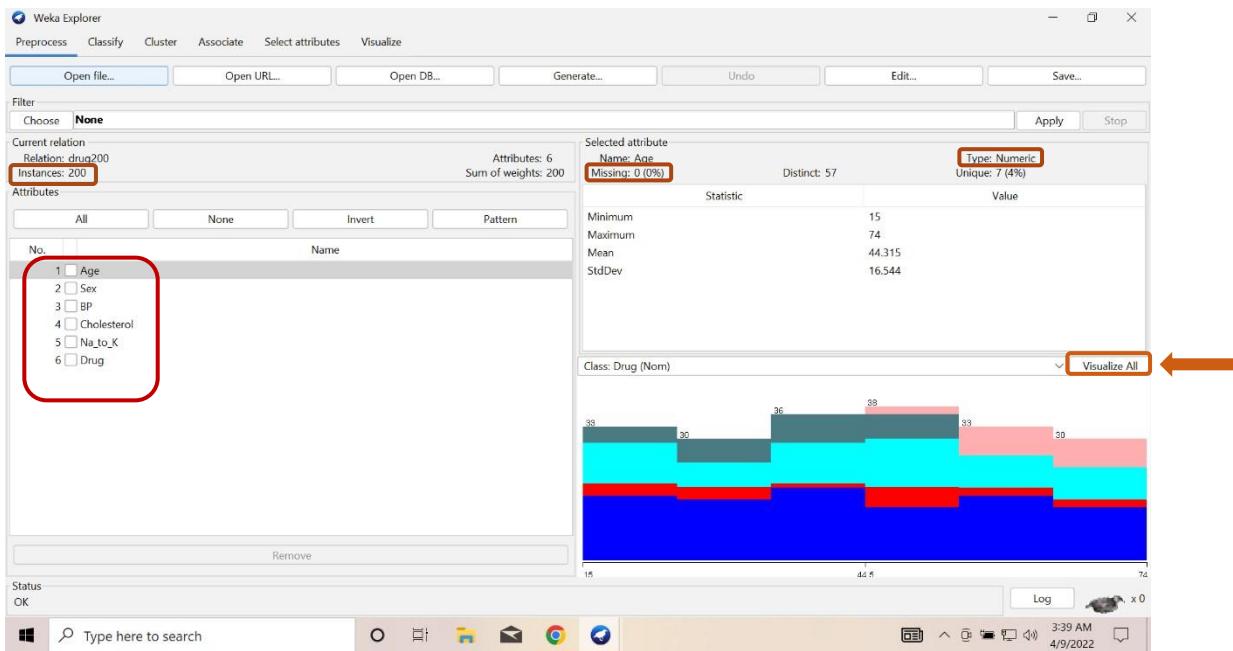
After clicking on “Open file” this interface will appear and we will be selecting the folder that have the dataset file.



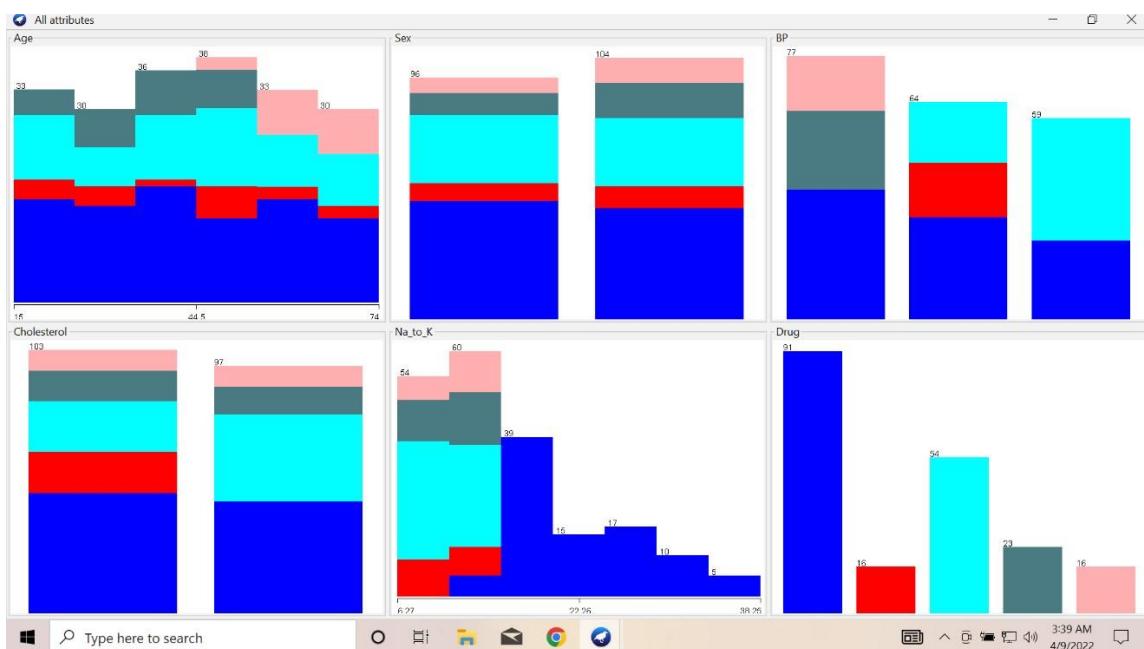
We have to make sure that we choose the correct files type. Our dataset was in .csv format so we have selected “CSV data files(.csv)”. Then our required dataset will be visible and we have to select it and by clicking on “Open” it will be open a new interface.



Here we can see that all the information is perfectly visible. All the six attributes are showing here. We can see the number of instances is 200. There is no missing value in this dataset and the type is numeric.



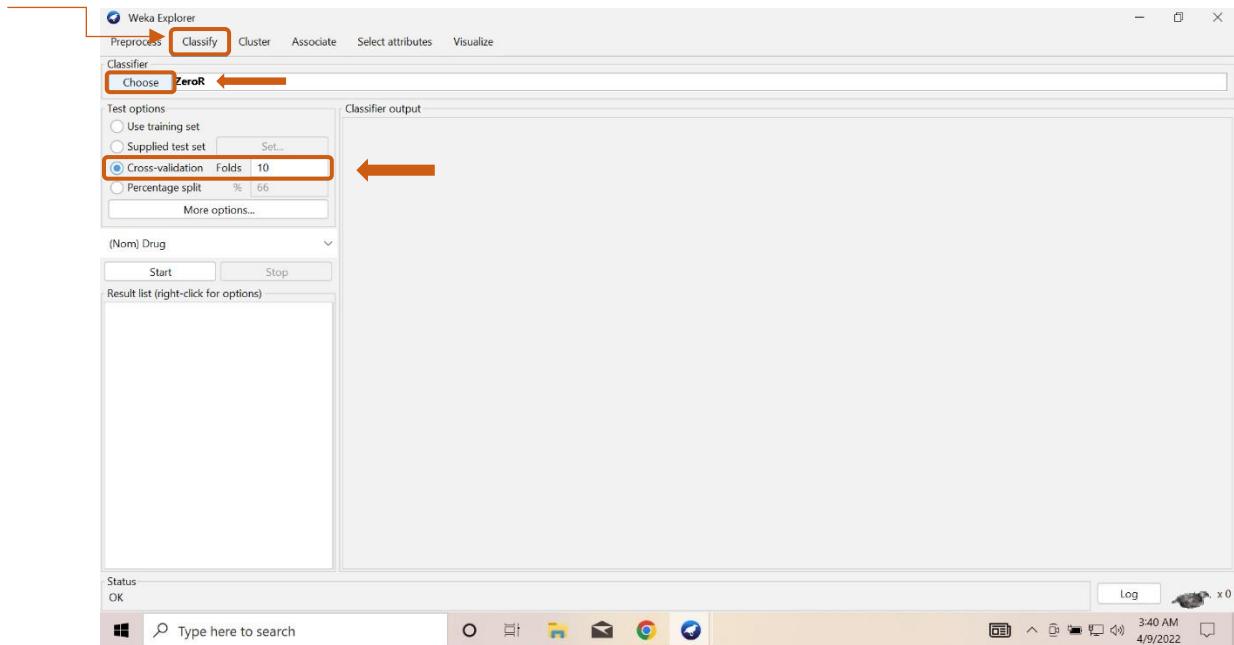
By clicking on “Visualize All” a new interface will appear. We can observe six types of graph for six types of attribute.



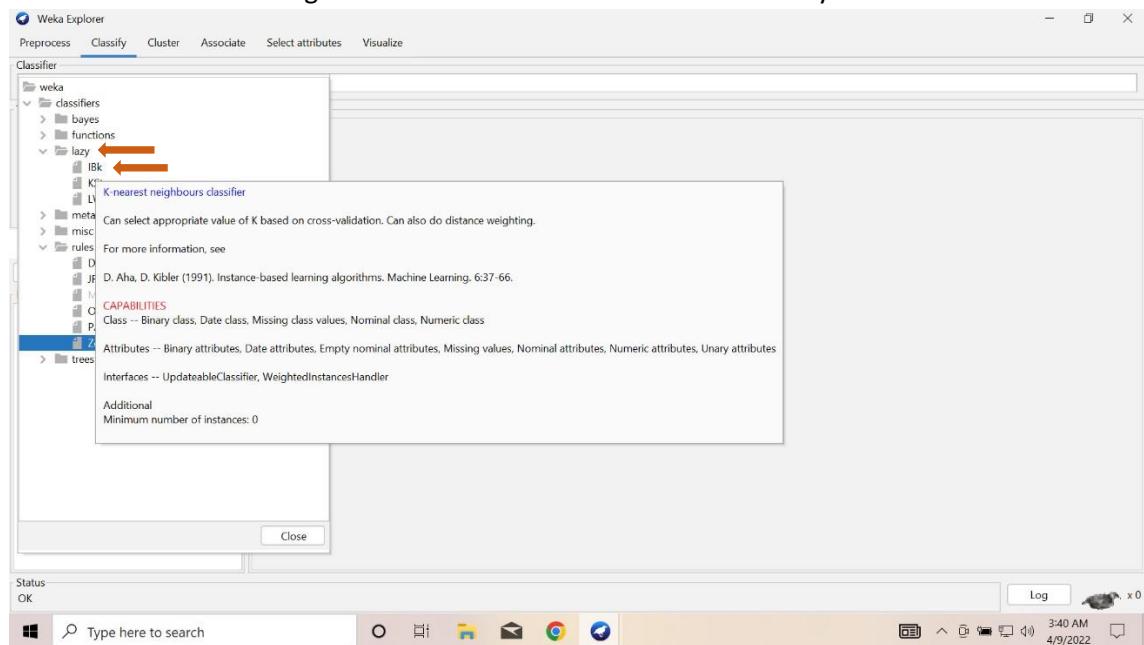
Till now we haven't trained our model yet. For all kind of classification these steps will be same. Now are going to apply algorithm.

K-Nearest Neighbours (KNN):

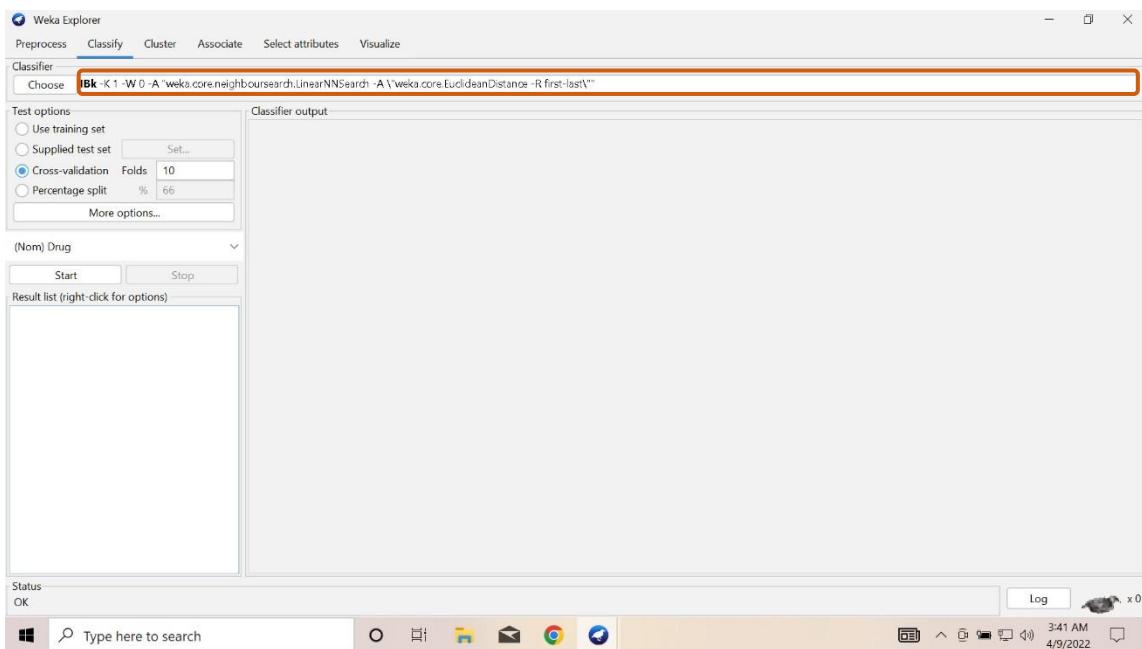
Now if we want to apply KNN the below steps will be followed. Firstly, choosing the option “Classify” and then choosing “Cross-validation, Fold 10”. Then we will click the option “Choose”.



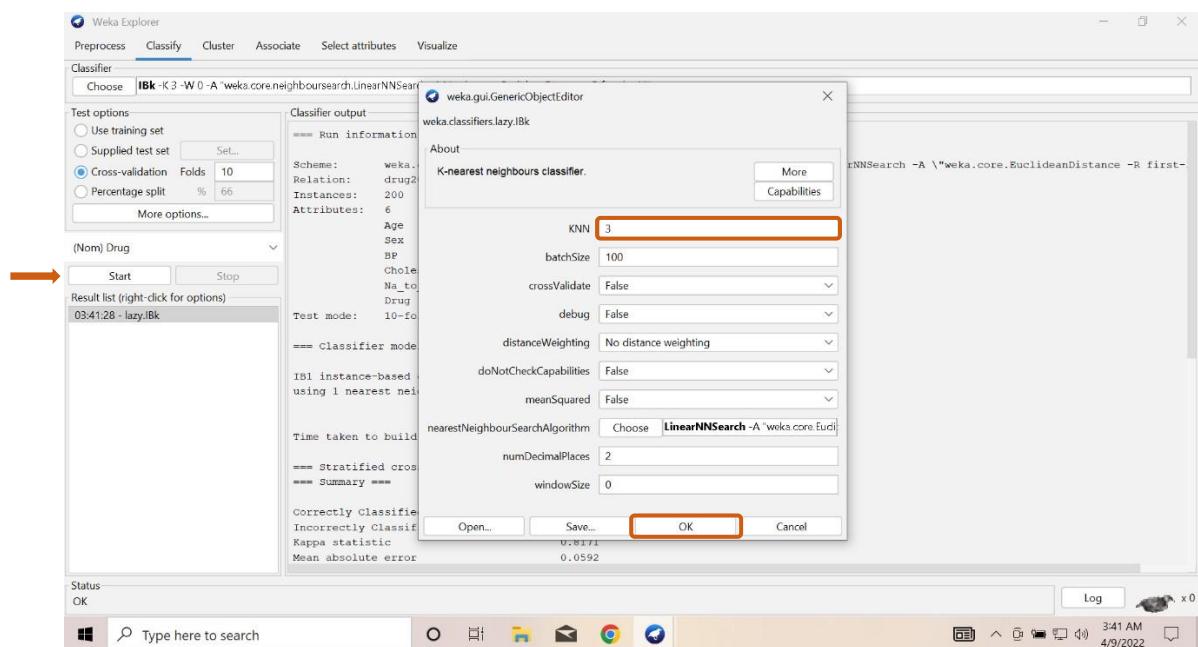
To work with K-Nearest Neighbour Classification we have chosen “lazy” then “IBk”



Then “IBk” will be loaded.

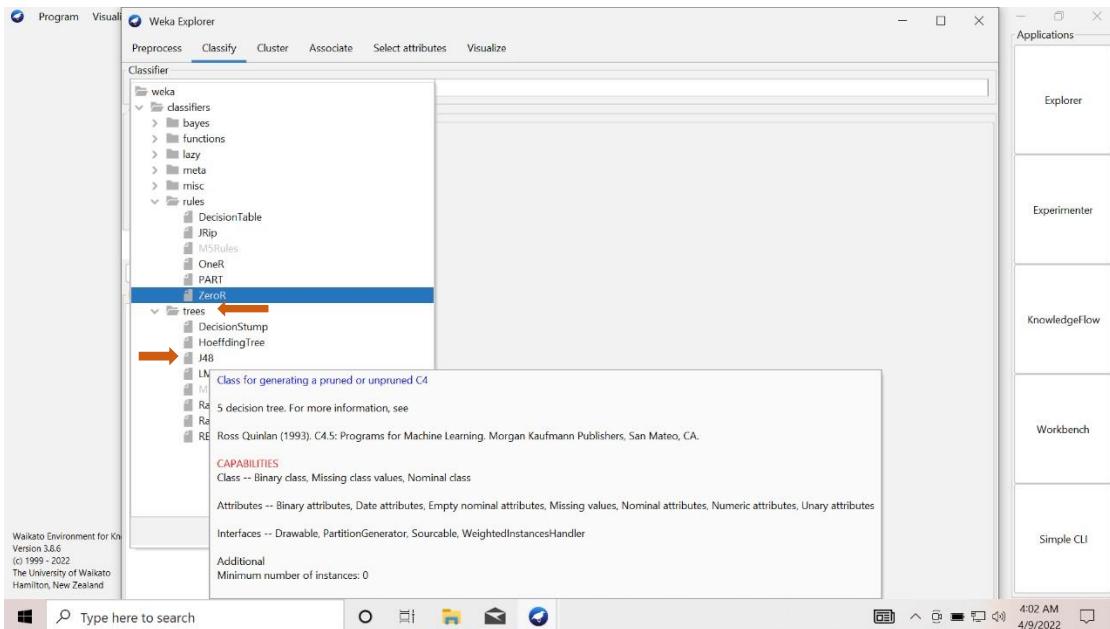


After left click in “IBk” this window will appear and we are choosing KNN value 3 then “OK” then “Start”.

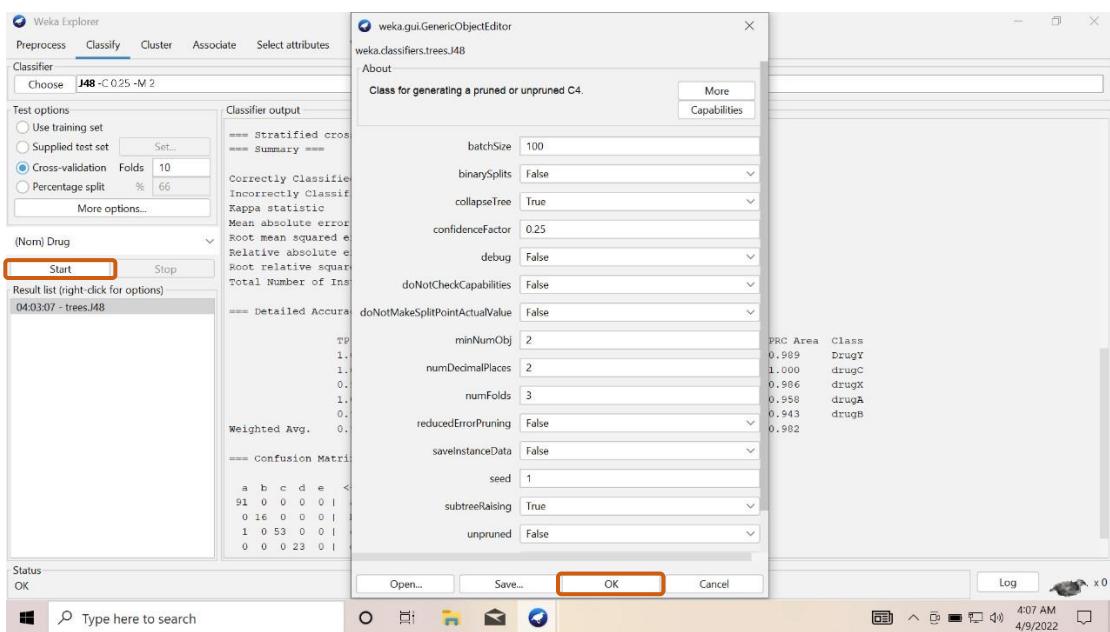


Decision Tree Classification Algorithm:

Now if we want to apply Decision Tree Classification Algorithm the below steps will be followed. Firstly, choosing the option “Classify” and then choosing “trees” then “J48”.

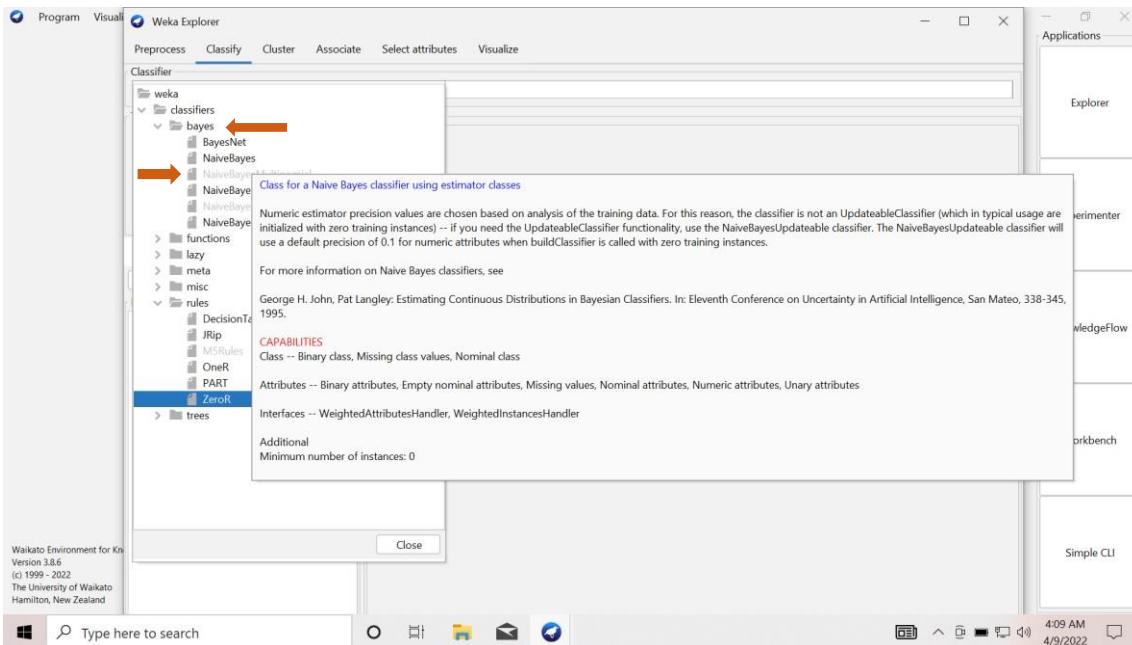


After left click in “J48” this window will appear and we are choosing “OK” then “Start”.

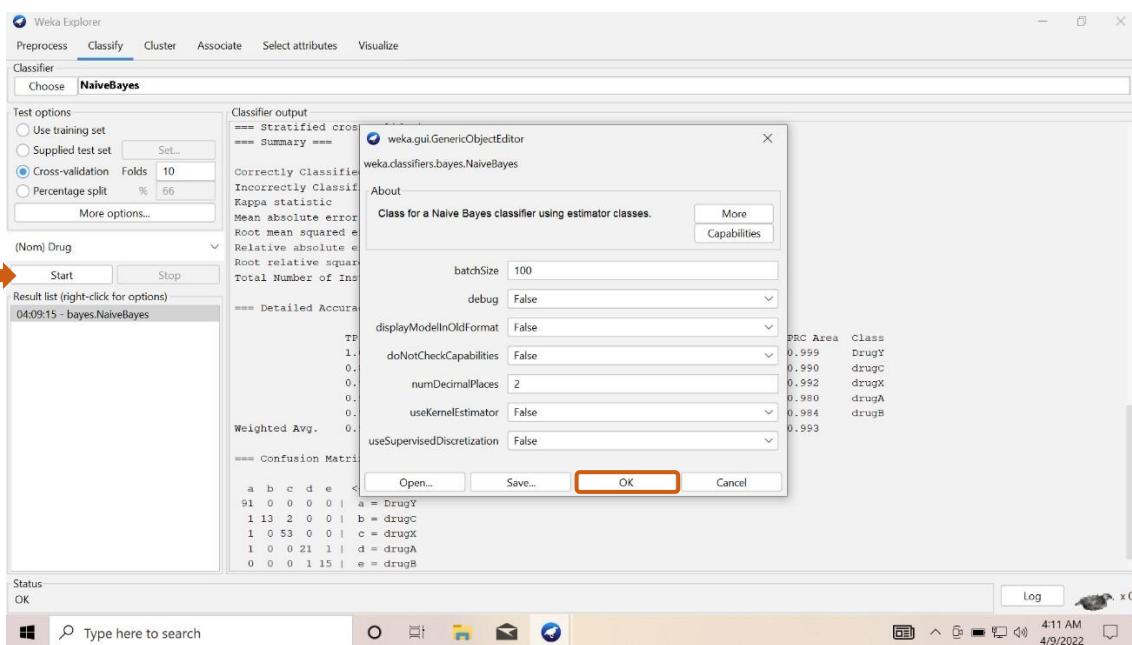


Naïve-Bayes Classification Algorithm:

Now if we want to apply Naïve-Bayes Classification Algorithm the below steps will be followed. Firstly, choosing the option “Classify” and then choosing “bayes”.

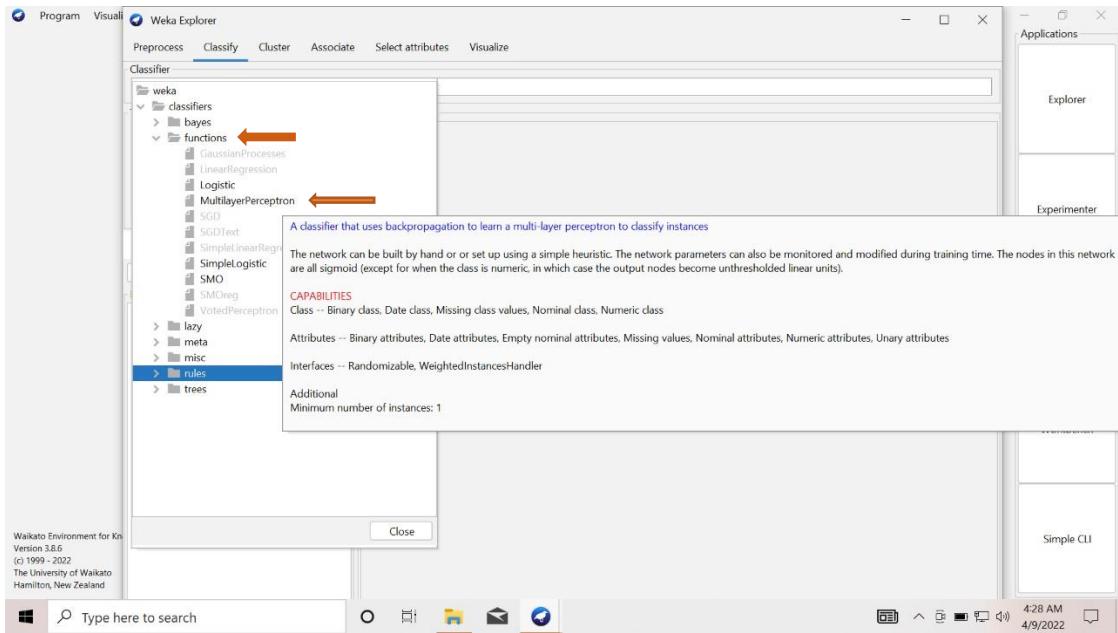


After left click in “NaiveBayes” this window will appear and we are choosing “OK” then “Start”.

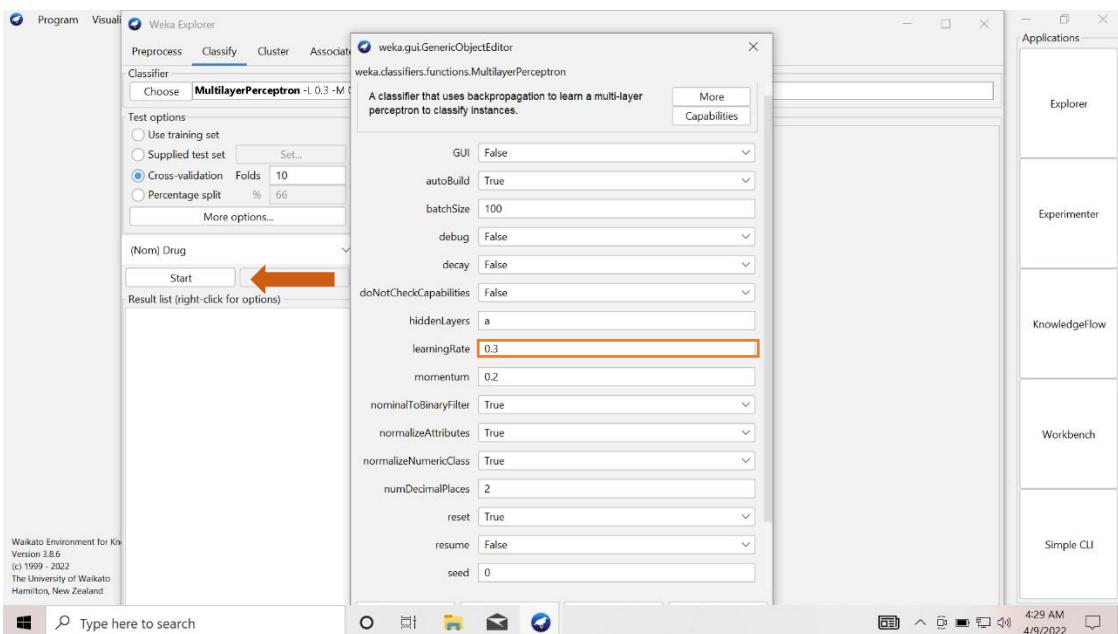


Neural Network Classification Algorithm:

Now if we want to apply Neural Network Classification Algorithm the below steps will be followed.
Firstly, choosing the option “function” & then MultilayerPerceptron.

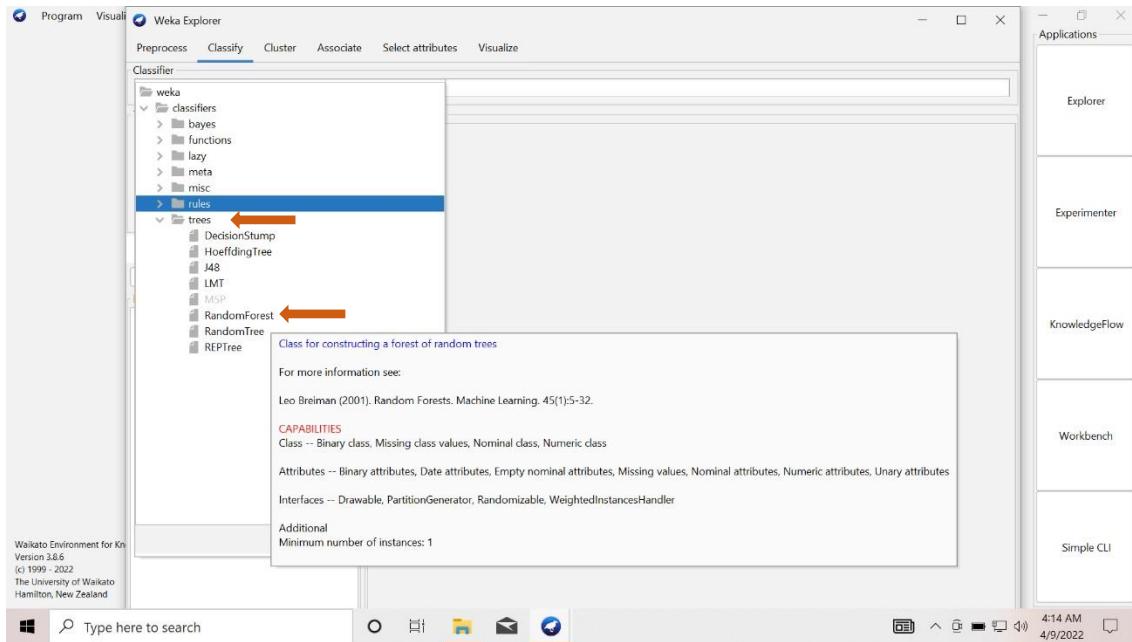


After left click in “MultilayerPerceptron” this window will appear and we are choosing “OK” then “Start”. Here learning rate is 0.3

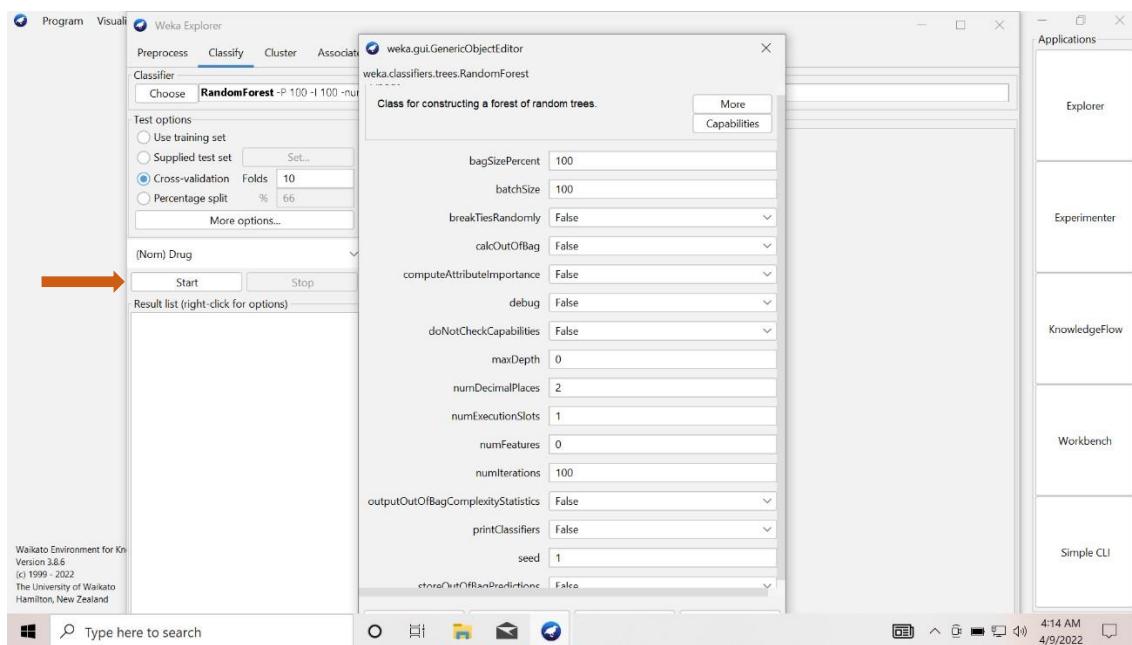


Random-Forest Classification Algorithm:

Now if we want to apply RandomForest Classification Algorithm the below steps will be followed. Firstly, choosing the option “rules” and then choosing “trees” then choosing “RandomForest”.

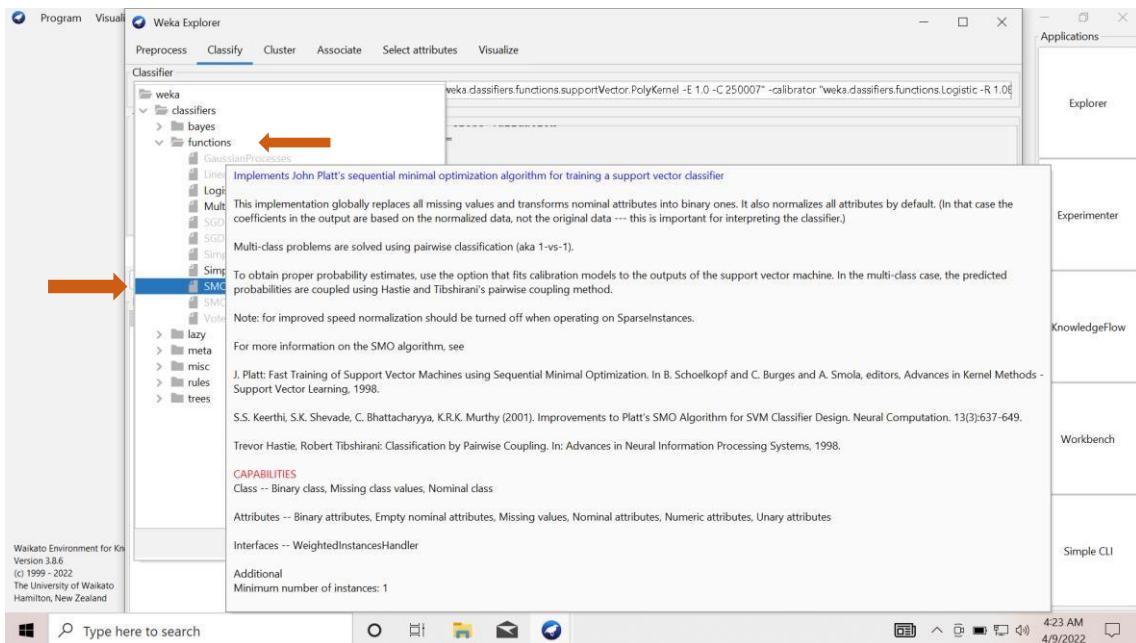


After left click in “RandomForest” this window will appear and we are choosing “OK” then “Start”.

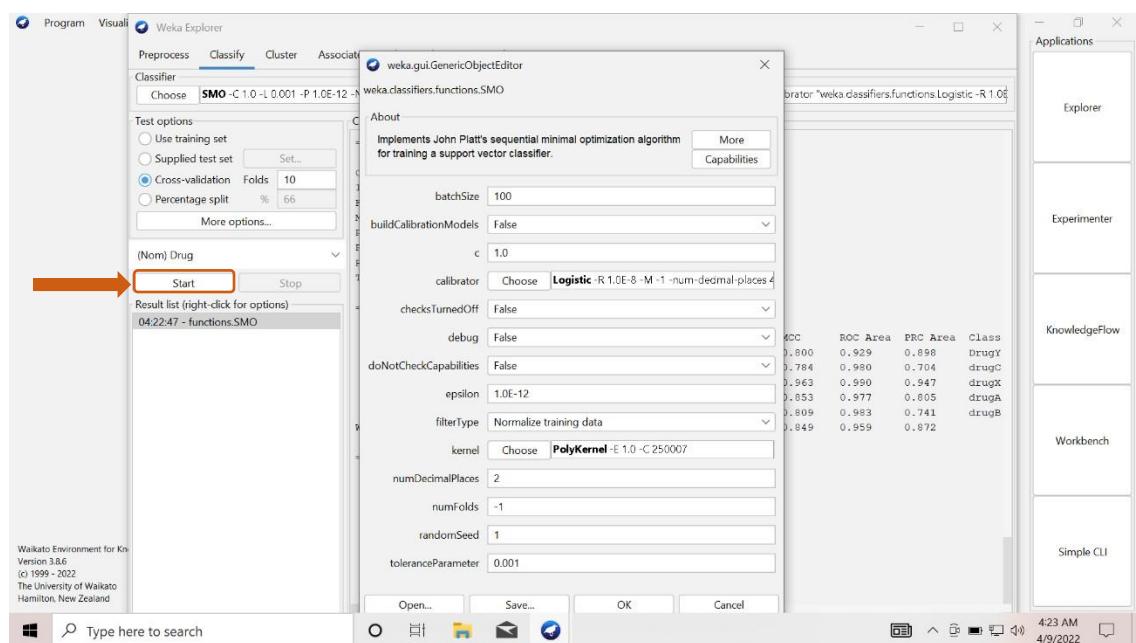


Support Vector Machine Classification Algorithm:

Now if we want to apply Support Vector Machine Classification Algorithm the below steps will be followed. Firstly, choosing the option “functions” then “SMC”.

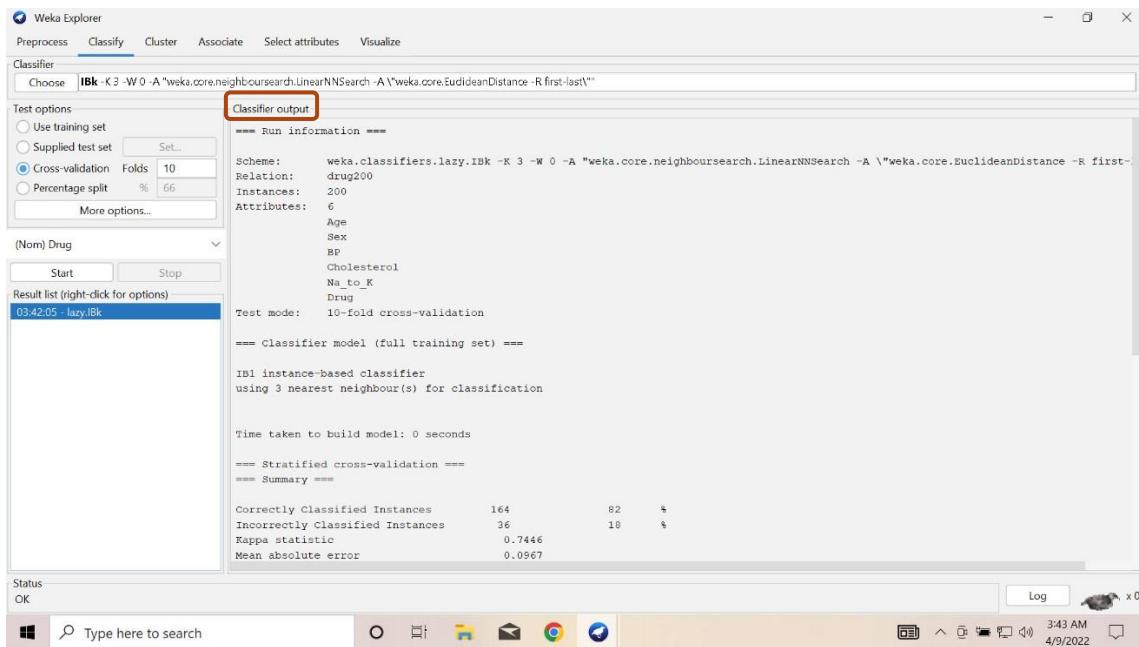


After left click in “SMO” this window will appear and we are choosing “OK” then “Start”.

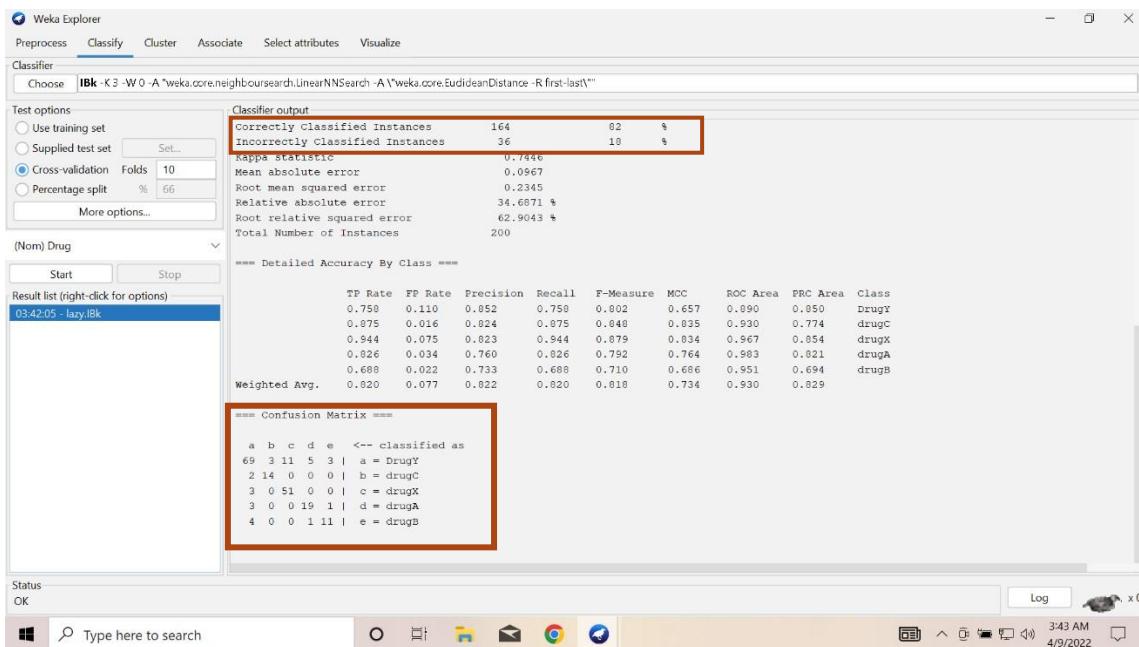


Discussion & Conclusion:

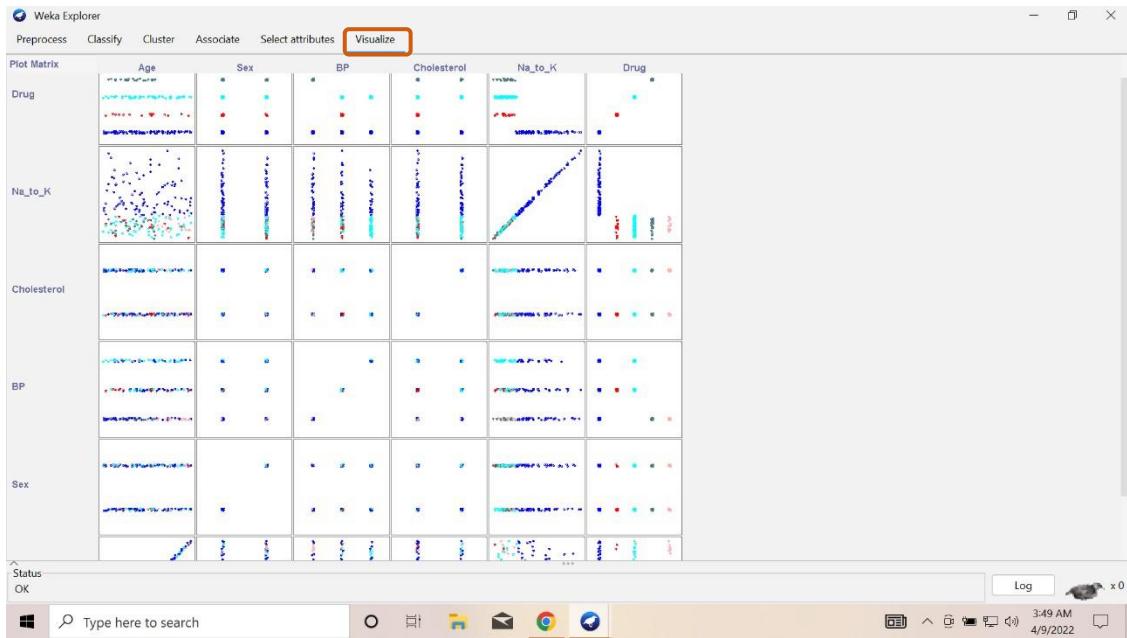
K-Nearest Neighbours:



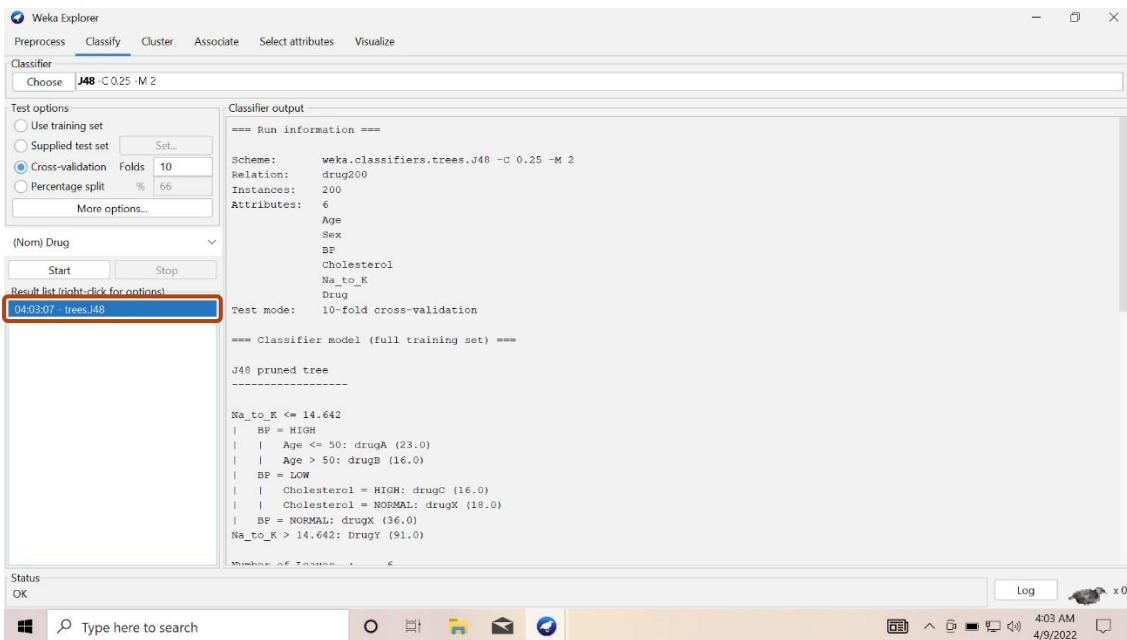
Correctly Classified Instances are 164 out of 200 instances. Incorrectly Classified Instances are 36. Accuracy correctly is 82% and incorrectly is 18%. We can also see the Confusion Matrix.



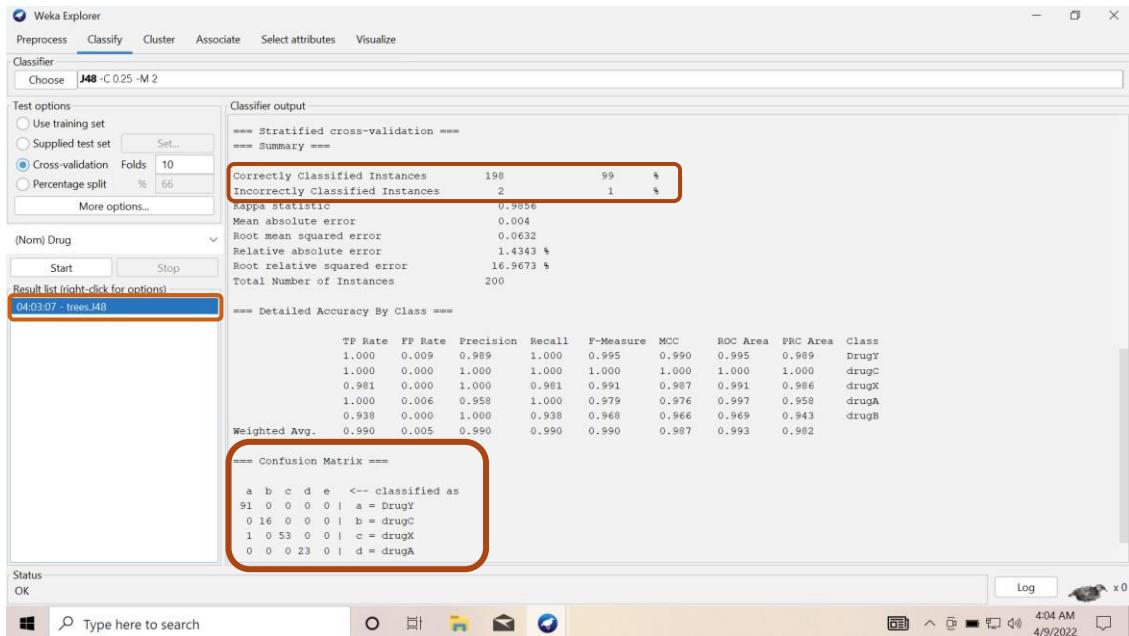
If we click on “Visualize” this interface will be shown.



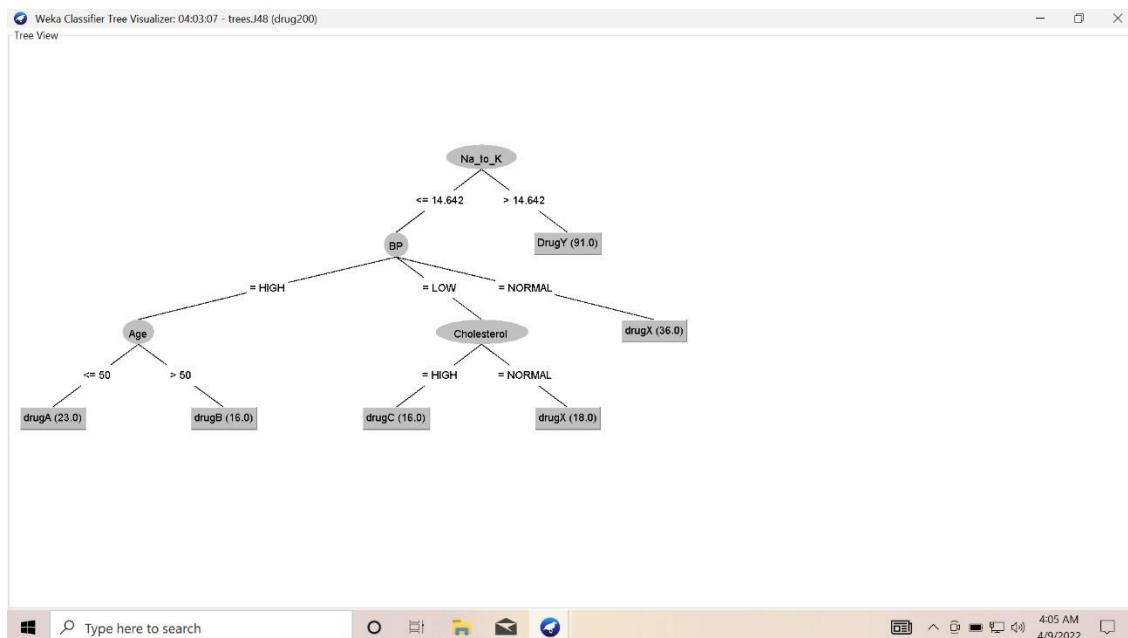
Decision Tree:



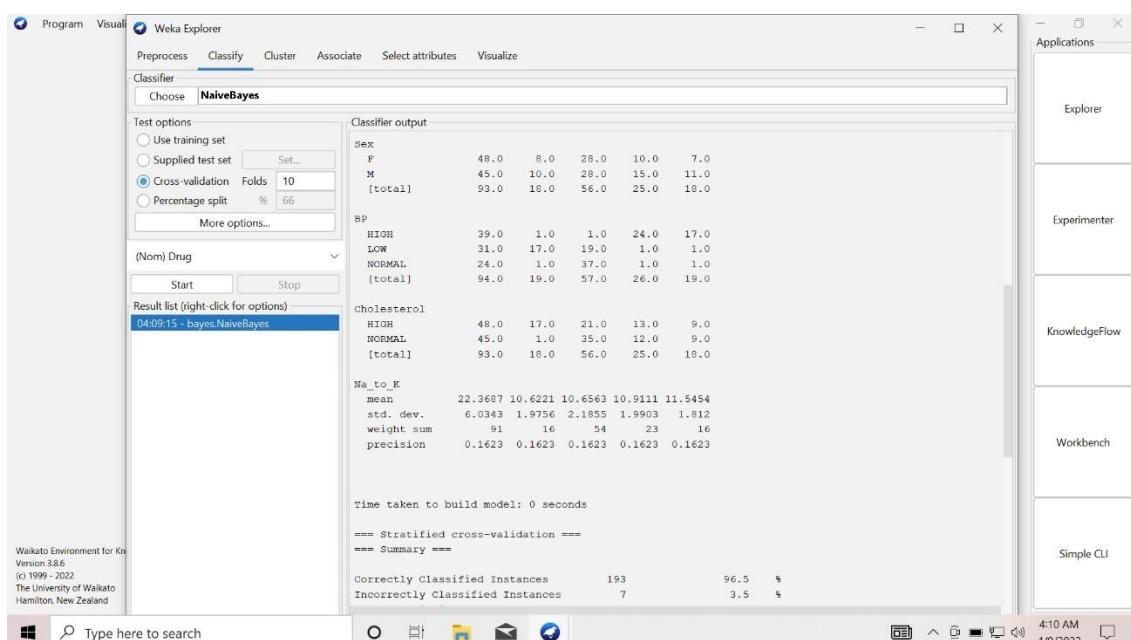
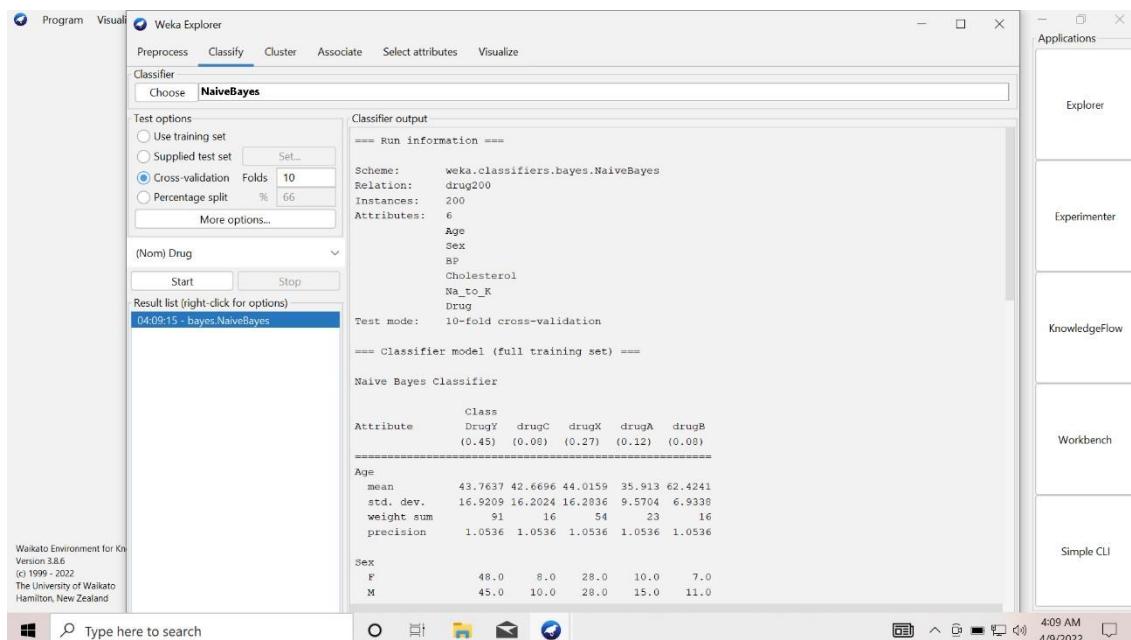
Here Correctly Classified Instances are 198 which is 99%.



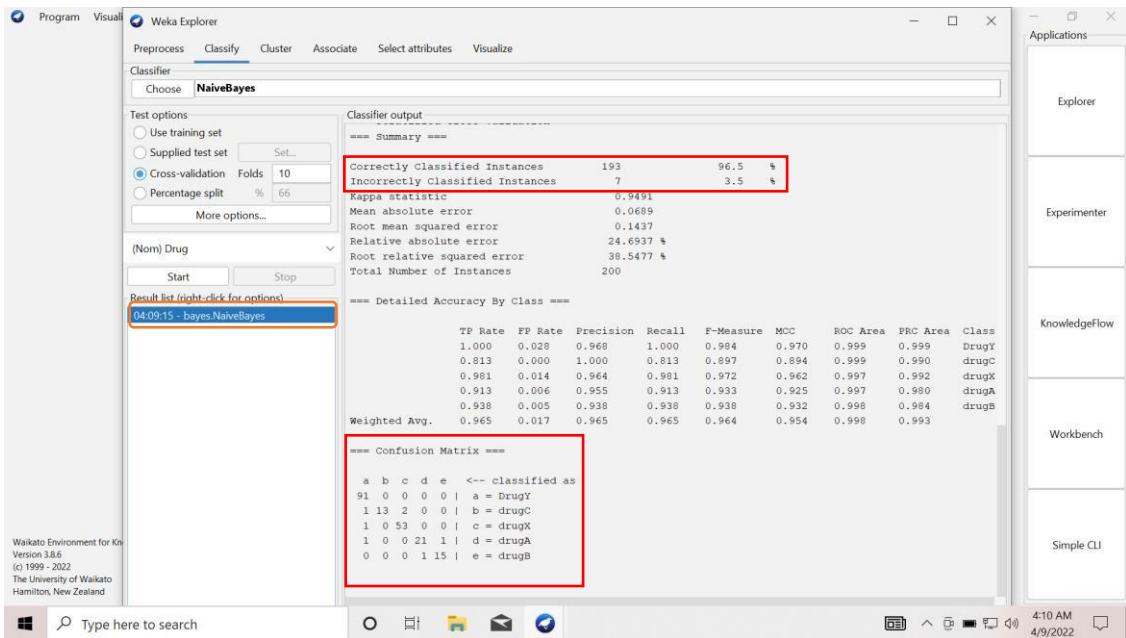
Here we can see the Tree.



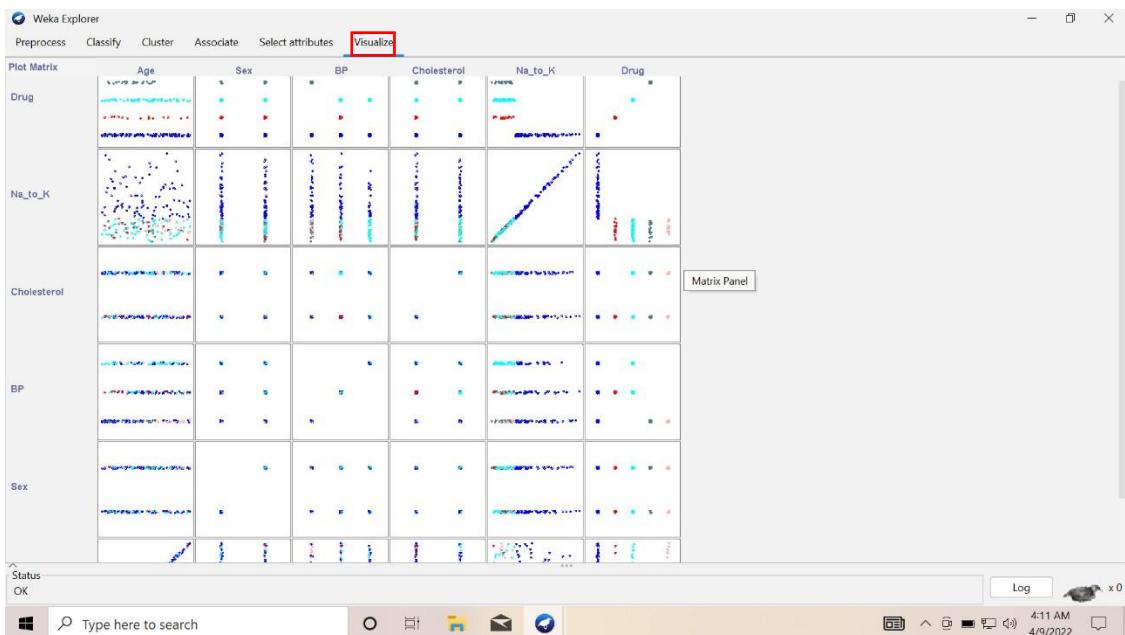
Naïve Bayes:



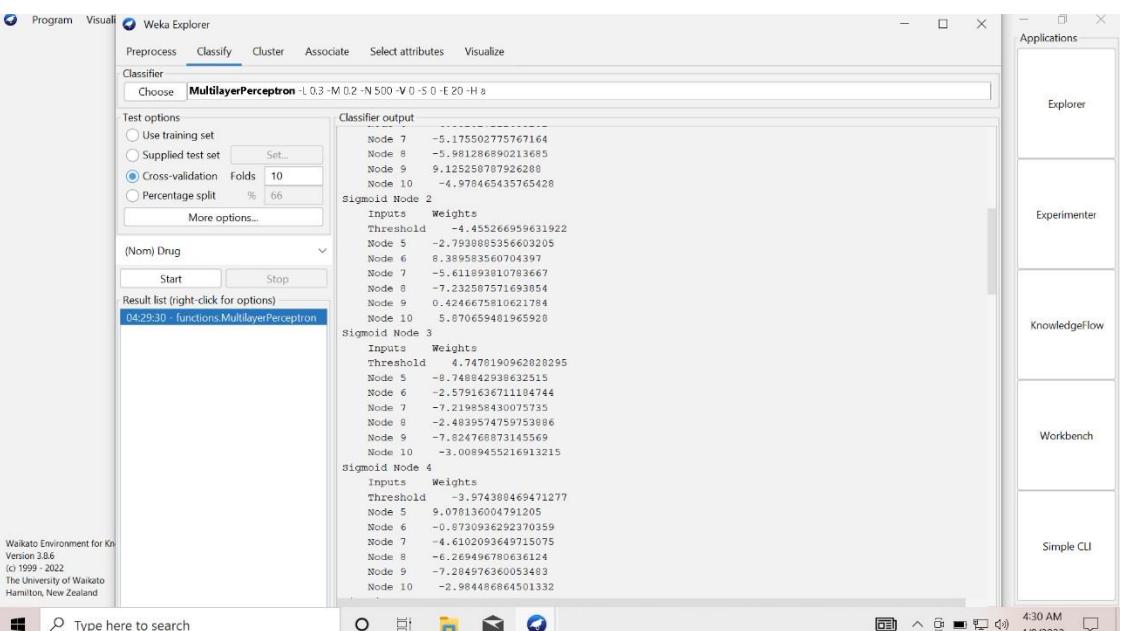
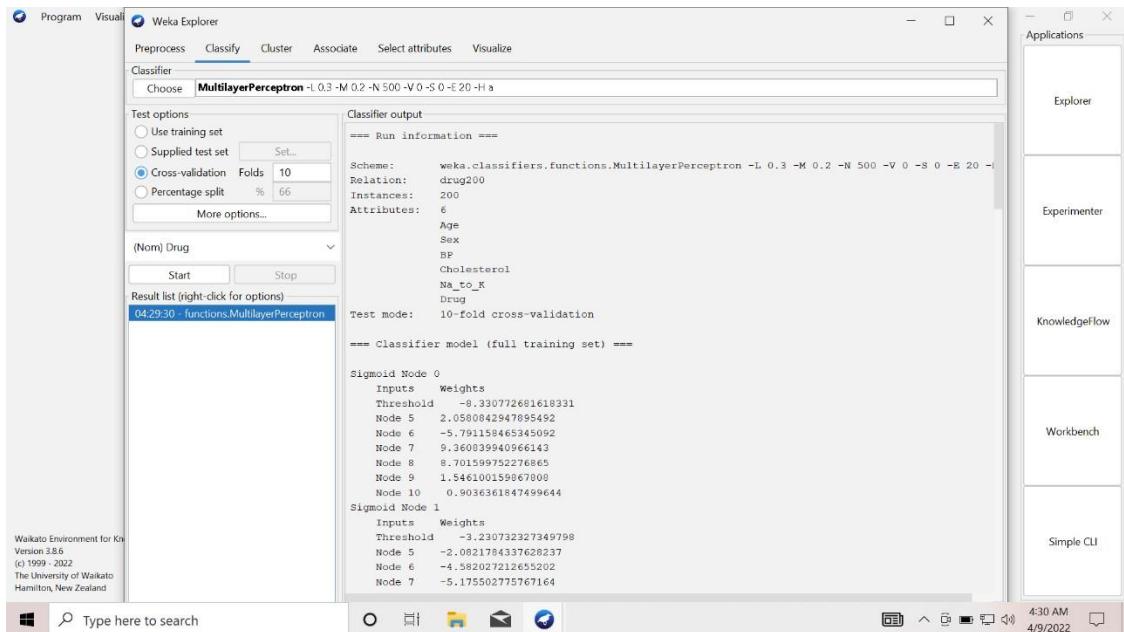
Correctly Classified Instances are 193 out of 200 instances. Incorrectly Classified Instances are 7. Accuracy is 96.5% and incorrectly is 3.5%. We can also see the Confusion Matrix.

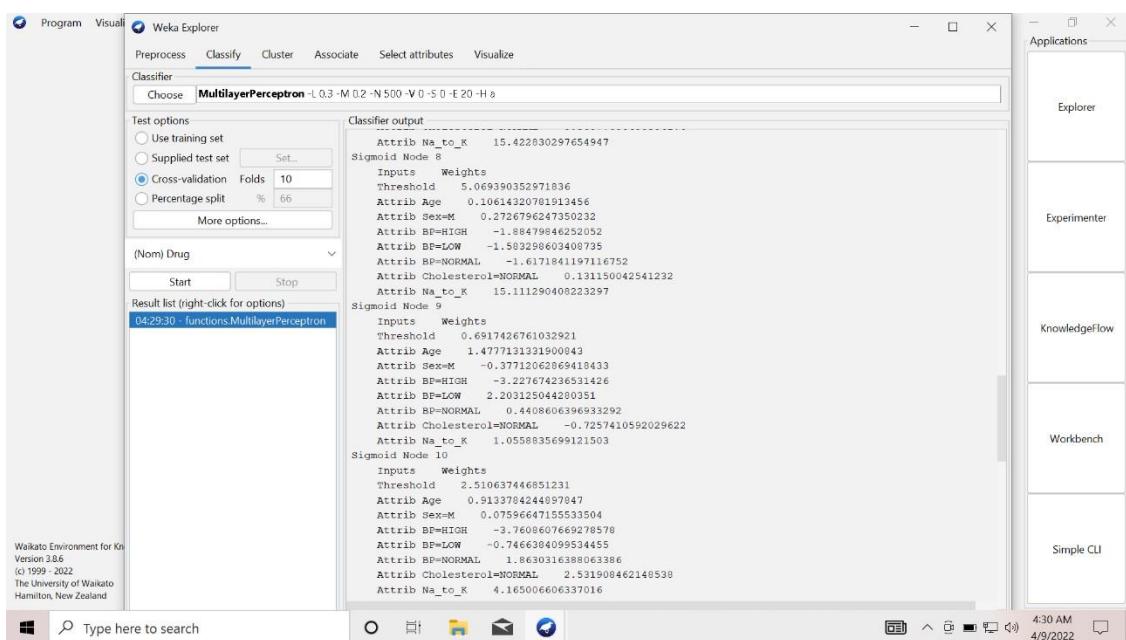
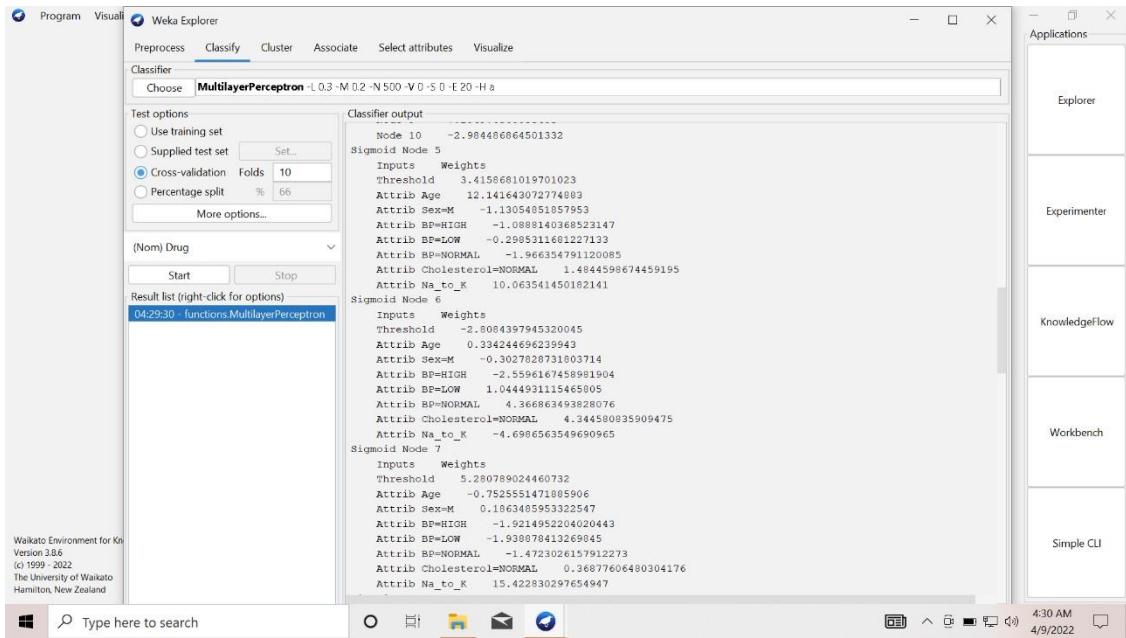


To Visualize click on the “Visualize” Option



Neural Network:





Correctly Classified Instances are 196 out of 200 instances. Incorrectly Classified Instances are 7. Accuracy is 98% and incorrectly is 2%. We can also see the Confusion Matrix.

Weka Environment for KN

Version 3.8.6
(c) 1999 - 2022
The University of Waikato
Hamilton, New Zealand

4:30 AM
4/9/2022

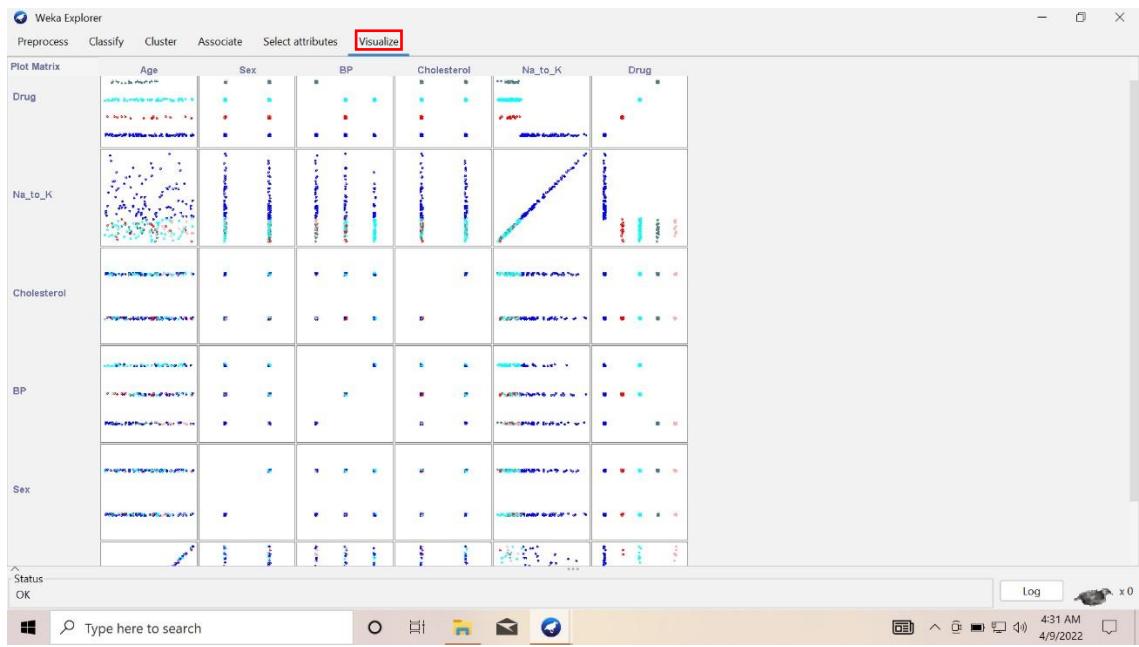
	a	b	c	d	e	<-- classified as	
a	90	1	0	0	0		a = DrugY
b	0	16	0	0	0		b = drugC
c	2	0	52	0	0		c = drugX
d	0	0	23	0	1		d = drugA
e	0	0	0	15	1		e = drugB

"/>

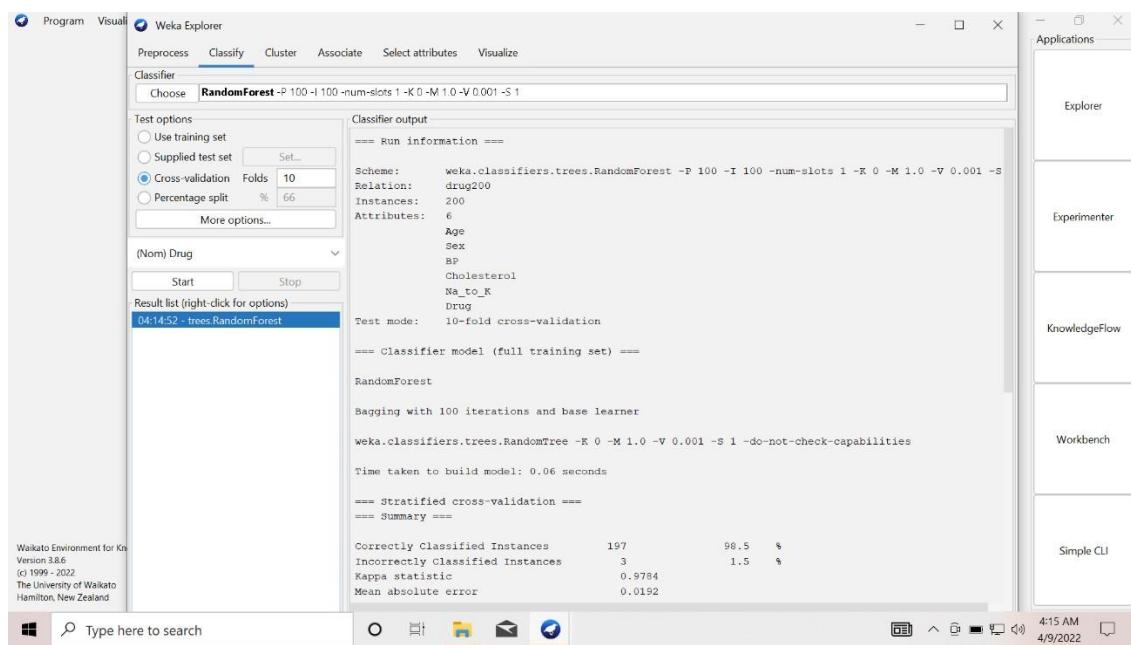
Weka Environment for KN

Version 3.8.6
(c) 1999 - 2022
The University of Waikato
Hamilton, New Zealand

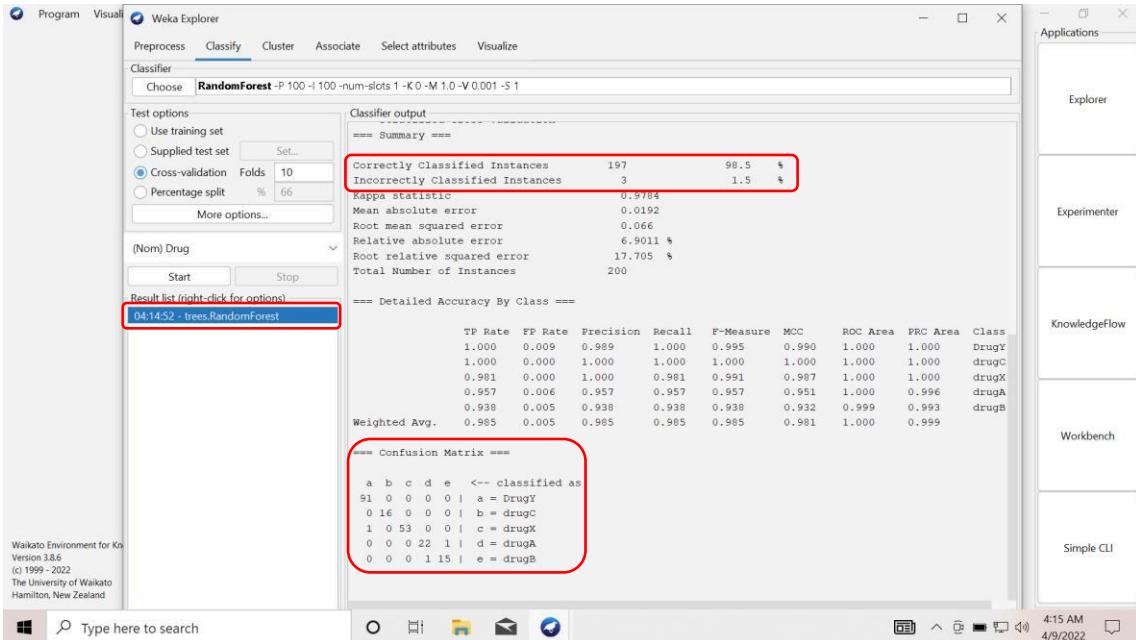
4:30 AM
4/9/2022



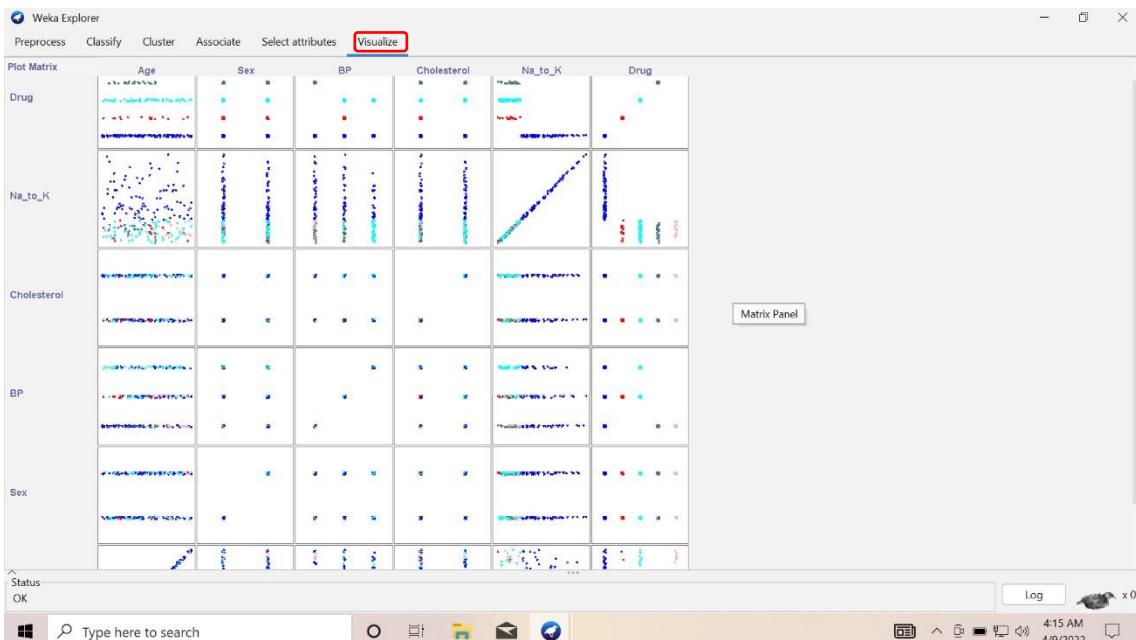
Random Forest:



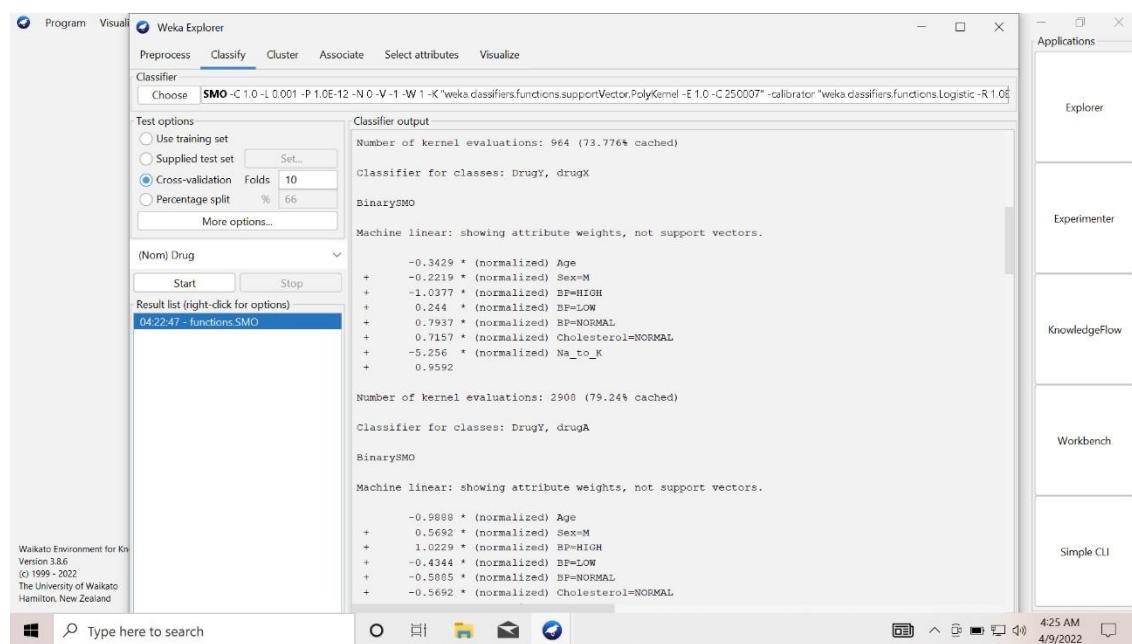
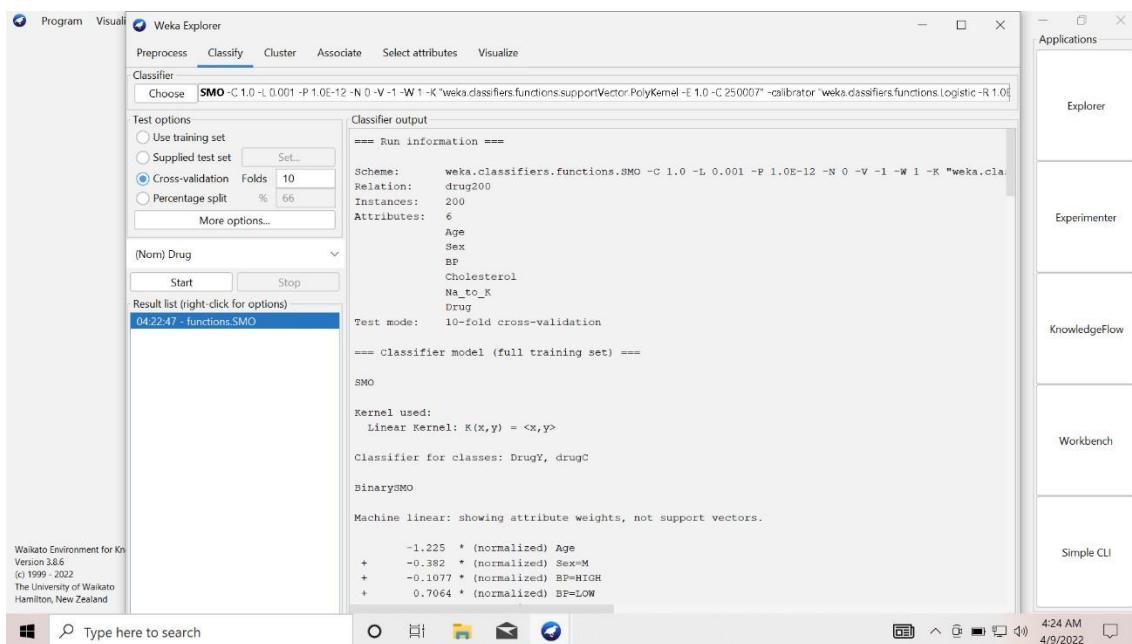
Correctly Classified Instances are 197 out of 200 instances. Incorrectly Classified Instances are 3. Accuracy is 98.5% and incorrectly is 1.5%. We can also see the Confusion Matrix

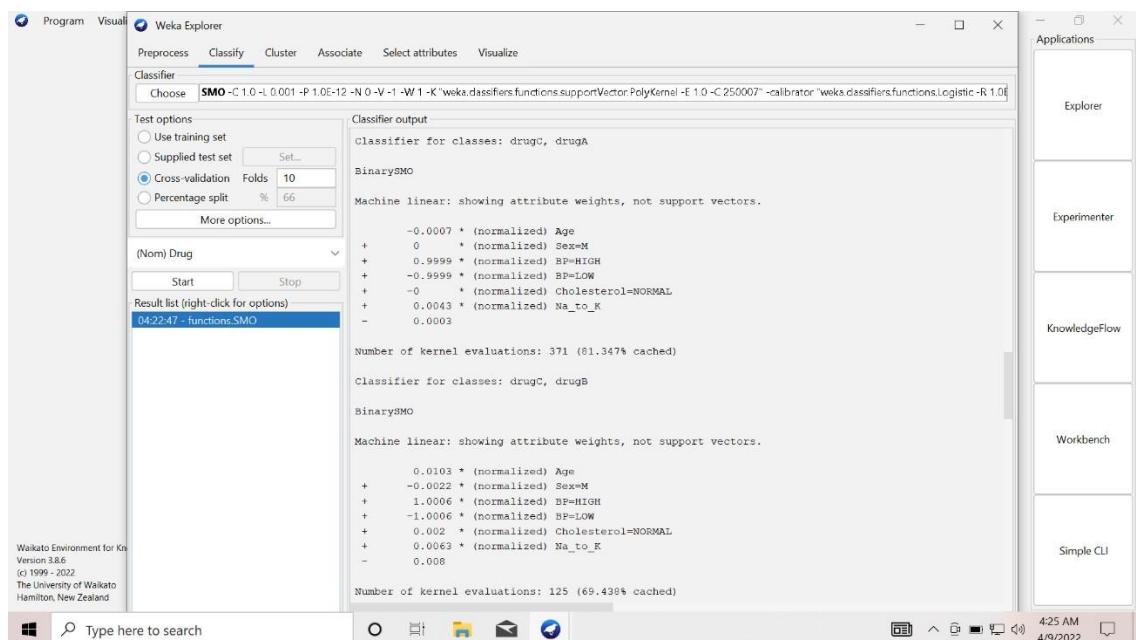
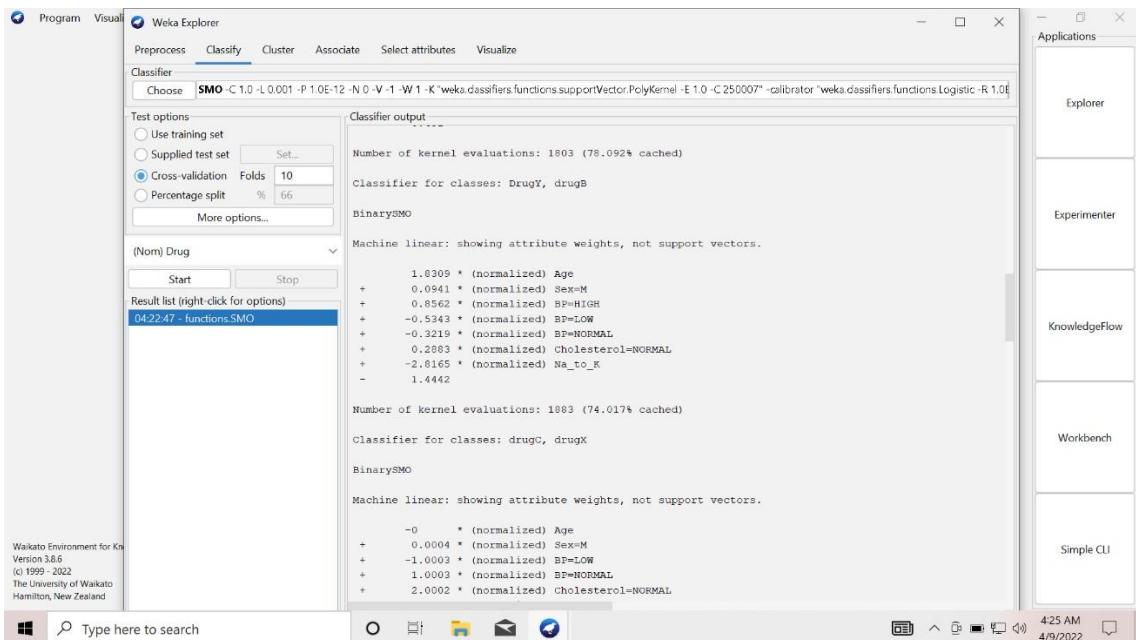


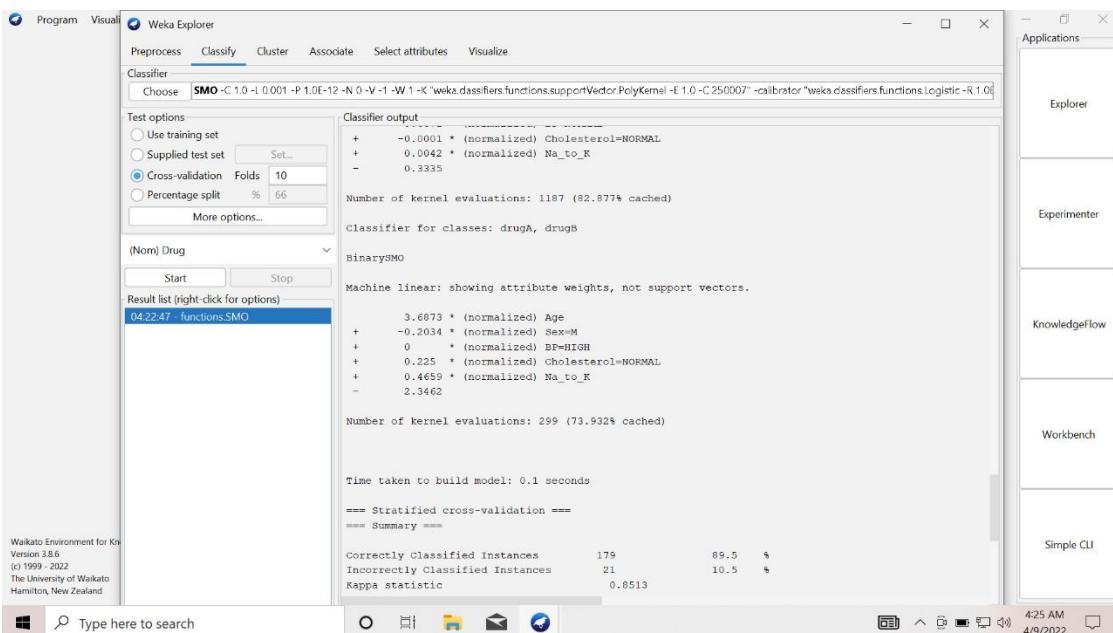
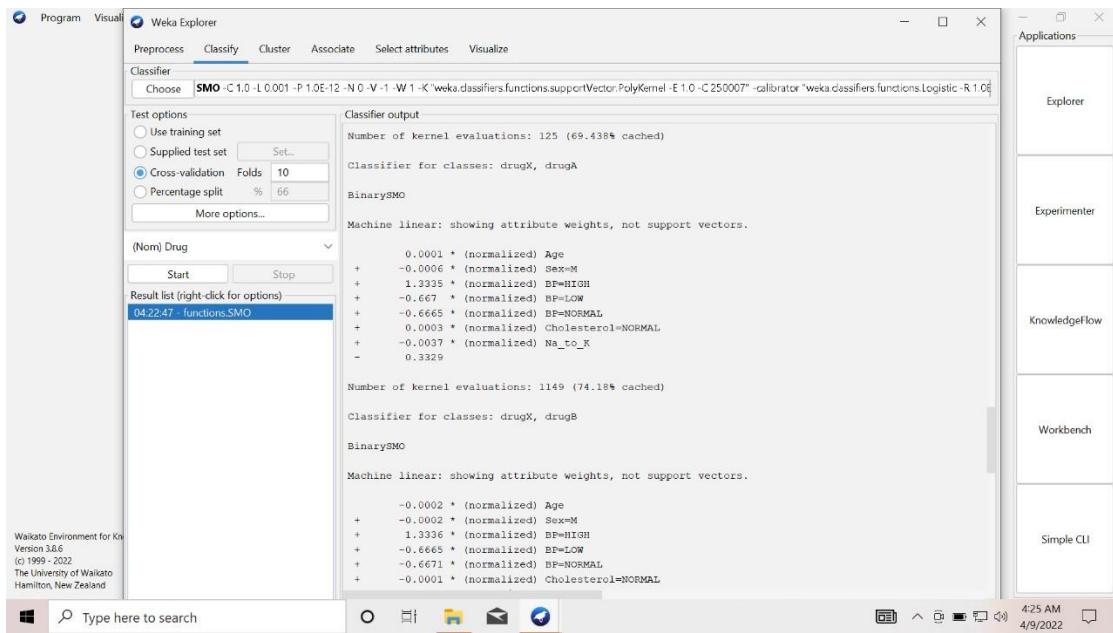
Clicking "Visualize" option to visualize data



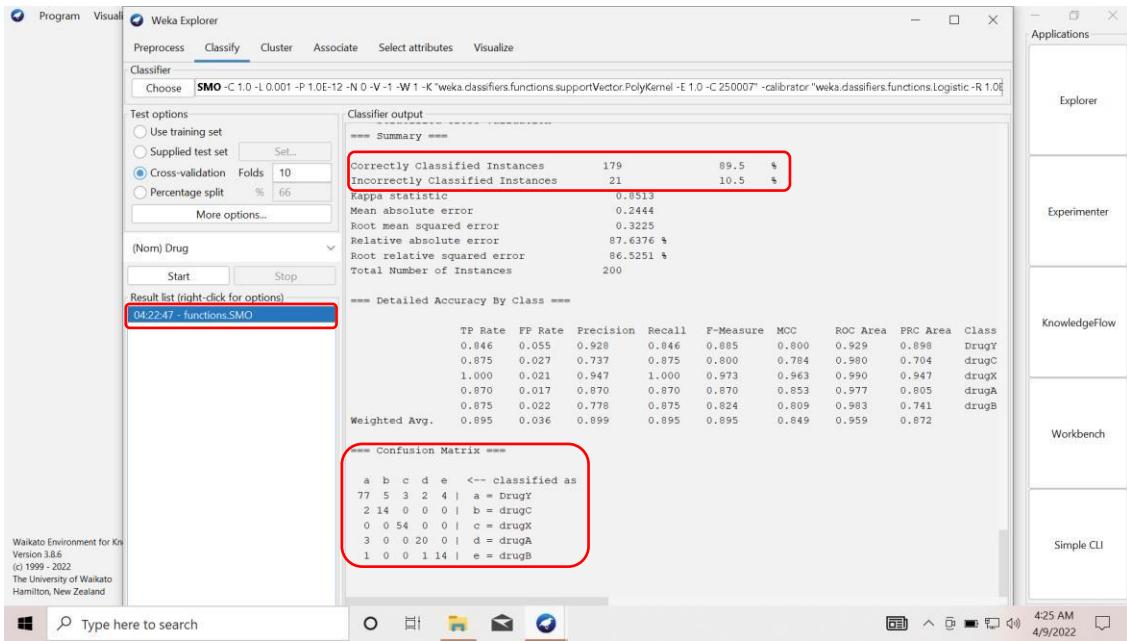
Support Vector Machine:



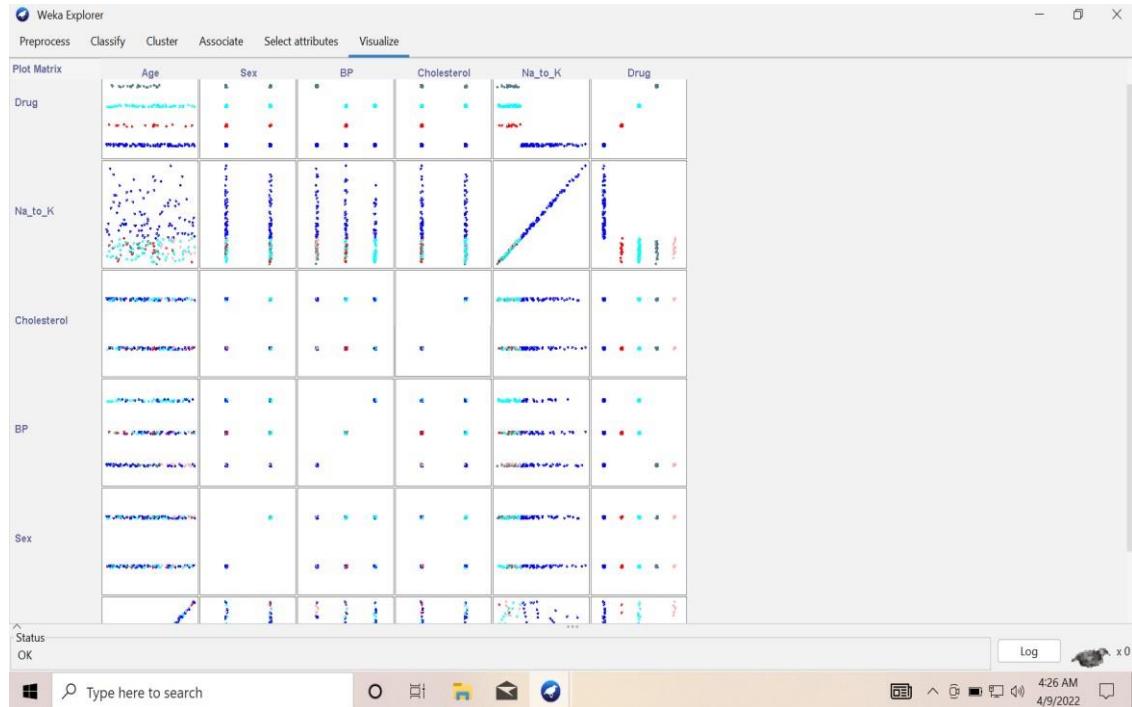




Correctly Classified Instances are 179 out of 200 instances. Incorrectly Classified Instances are 21. Accuracy is 89.5% and incorrectly is 10.5%. We can also see the Confusion Matrix.



To visualize the data we can select “Visualize” Option.



So, finally we have got the results according to the Accuracy are given below-

K-Nearest Neighbours: 82%

Decision Tree: 99%

Naïve Bayes: 96.5%

Neural Network: 98%

Random Forest: 98.5%

Support Vector Machine: 89.5%

From the given result we see Decision Tree gives highest accuracy to fit our model.