Otto Friedrich University Bamberg

# Data Streams and Complex Event Processing

# Assignment - 05

*Supervisor:*
**Prof. Dr. Daniela Nicklas**
Mobile Software Systems

*Submitted By :*
**Manasa Reddy Bandari - 1942692**
**Akshay Sharma - 1943984**
**Kummari Daniel Raj - 1973448**

Feb 4, 2019

# Contents

# 1    Exercise 1:

To investigate the influence of streamed continuous data on the efficiency of the learning results, we use two different datasets: one is from cars sales data and the second one is from the health sector (both data sets can be found on the VC under the folder "Datasets"). The goal of the task is to observe how a classification algorithm creates classification models based on varying sizes of data. By using different sizes for the training and testing, we try to imitate the behavior of streamed data. The class to be observed in the cars dataset is the price and in the second dataset, the final attribute is the class. (increased binding protein, decreased binding protein and negative)

## 1.1    Have a look at both datasets and give a brief explanation on both datasets and what do they describe?

Car Dataset includes characteristics like number of people, number of doors, overall cost, safety values and much more. There are dependencies in between many characteristics like capacity and boot space. It shows an entire domain to check and decide the best option for the desired input.

Health Dataset includes information about the person's age and gender and will determine on the number of factors whether they have thyroid or not.

## 1.2   To learn the effect of decreasing training set on the data, we use the first dataset (cars) with the classification algorithm J48 and record the precision results in three different settings:

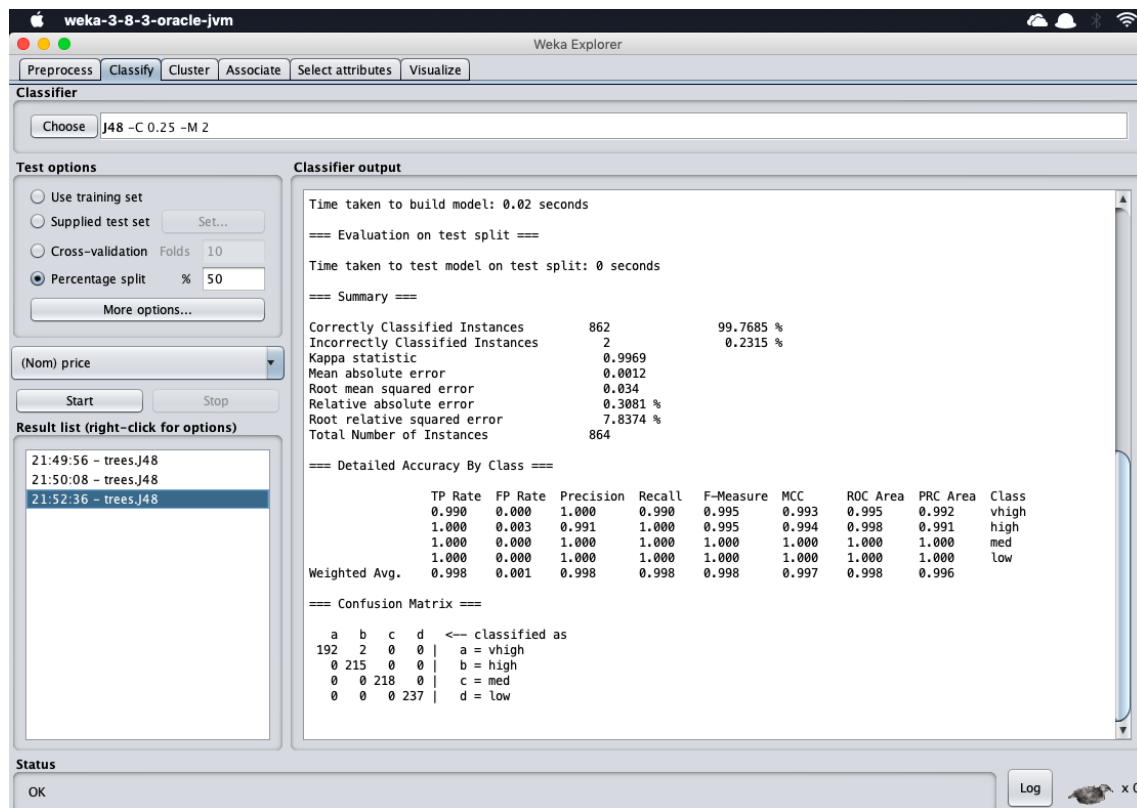### 1.2.1   Use 50% of the dataset for training and the other 50% for testing



Figure 1: 50% Training Dataset

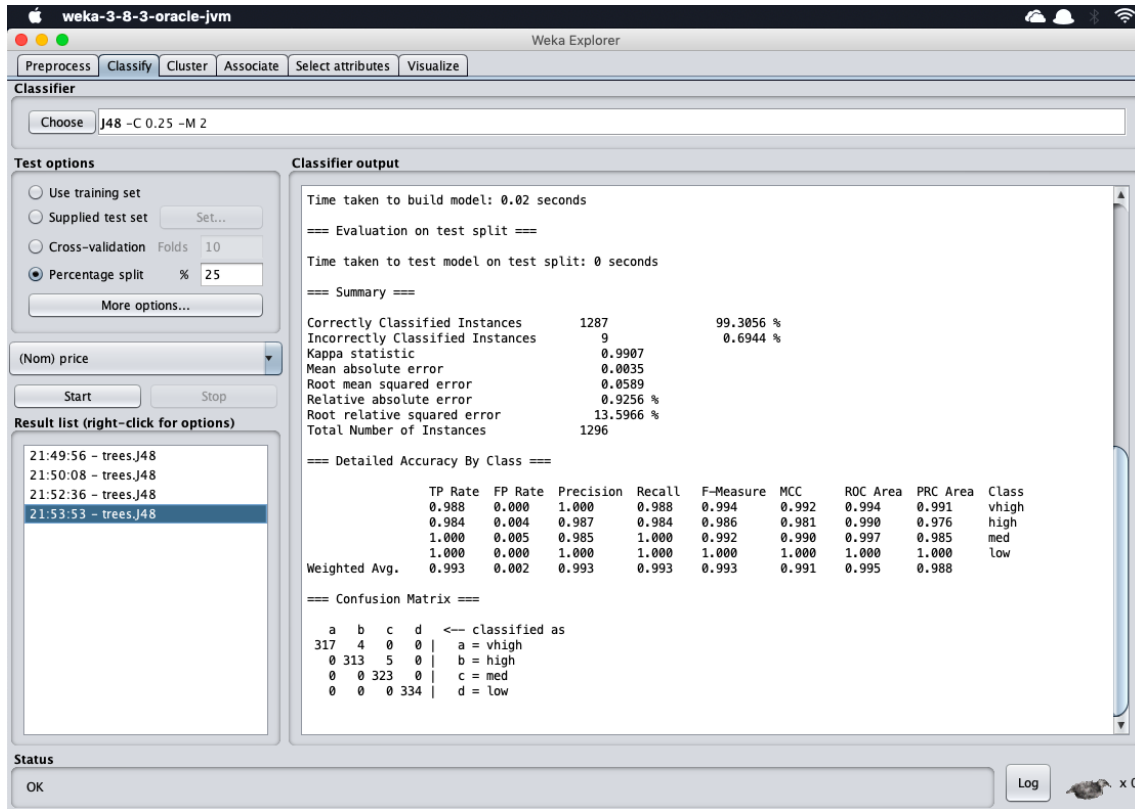### 1.2.2 Use 25% of the dataset for training and the other 75% for testing



Figure 2: 25% Training Dataset

### 1.2.3   Use 10% of the dataset for training and the other 90% for testing
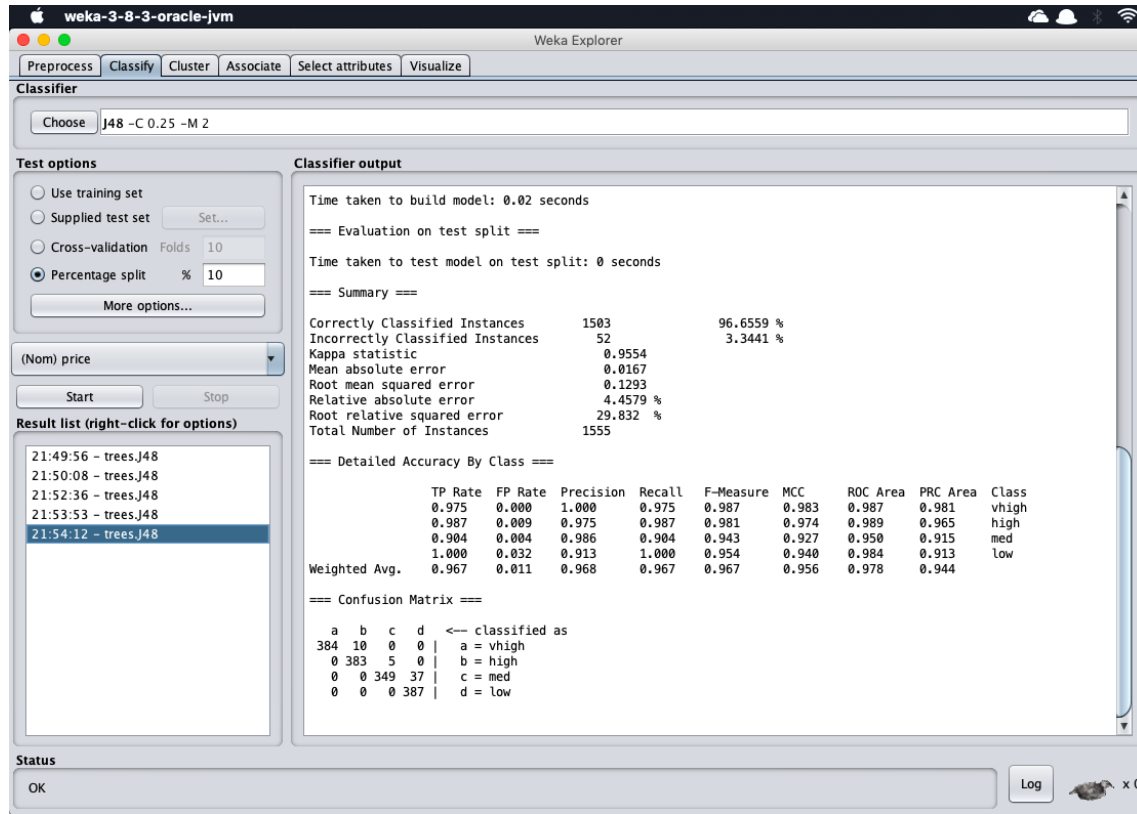


Figure 3: 10% Training Dataset

## 1.3 Use the second dataset (thyroid) with the same algorithm and record the precision results in three different settings:

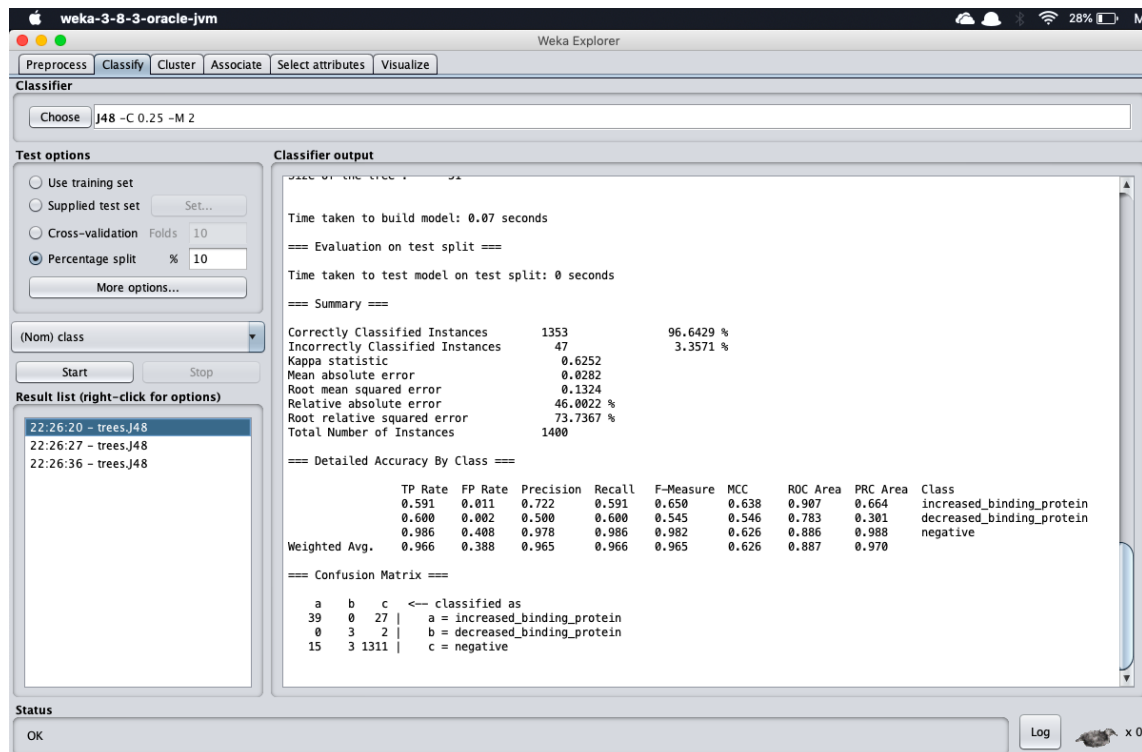### 1.3.1 Use 50% of the dataset for training and the other 50% for testing



Figure 4: 50% Training Dataset

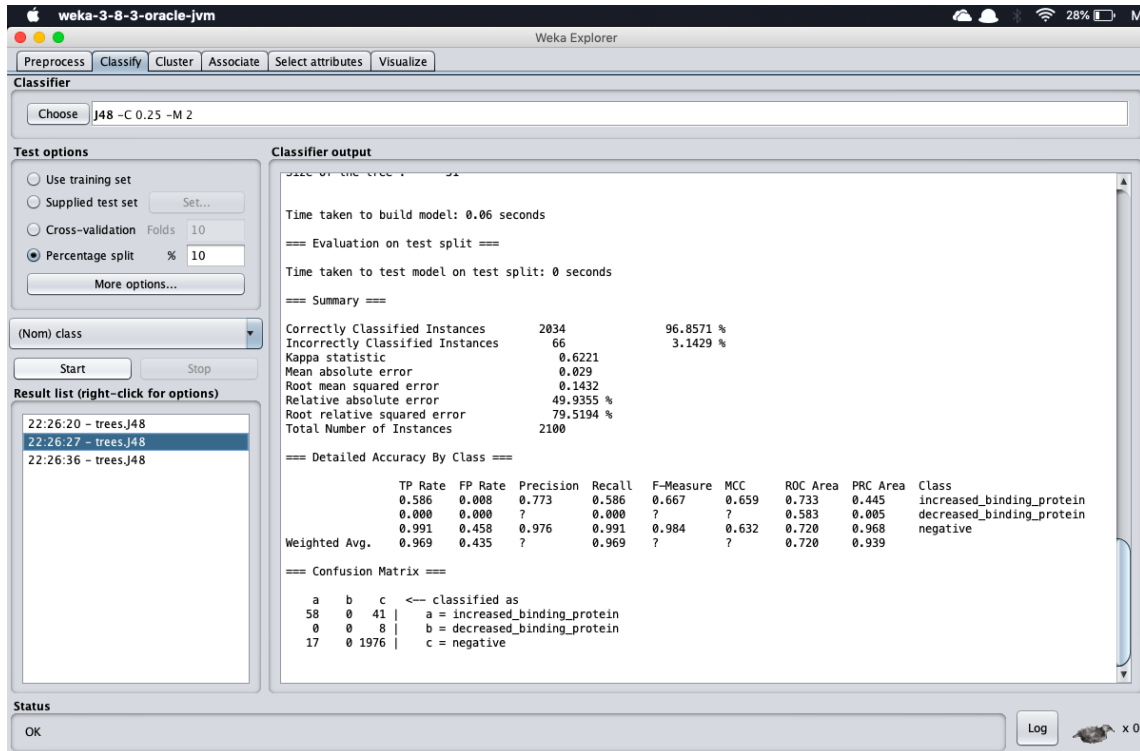### 1.3.2   Use 25% of the dataset for training and the other 75% for testing



Figure 5: 25% Training Dataset

### 1.3.3   Use 10% of the dataset for training and the other 90% for testing
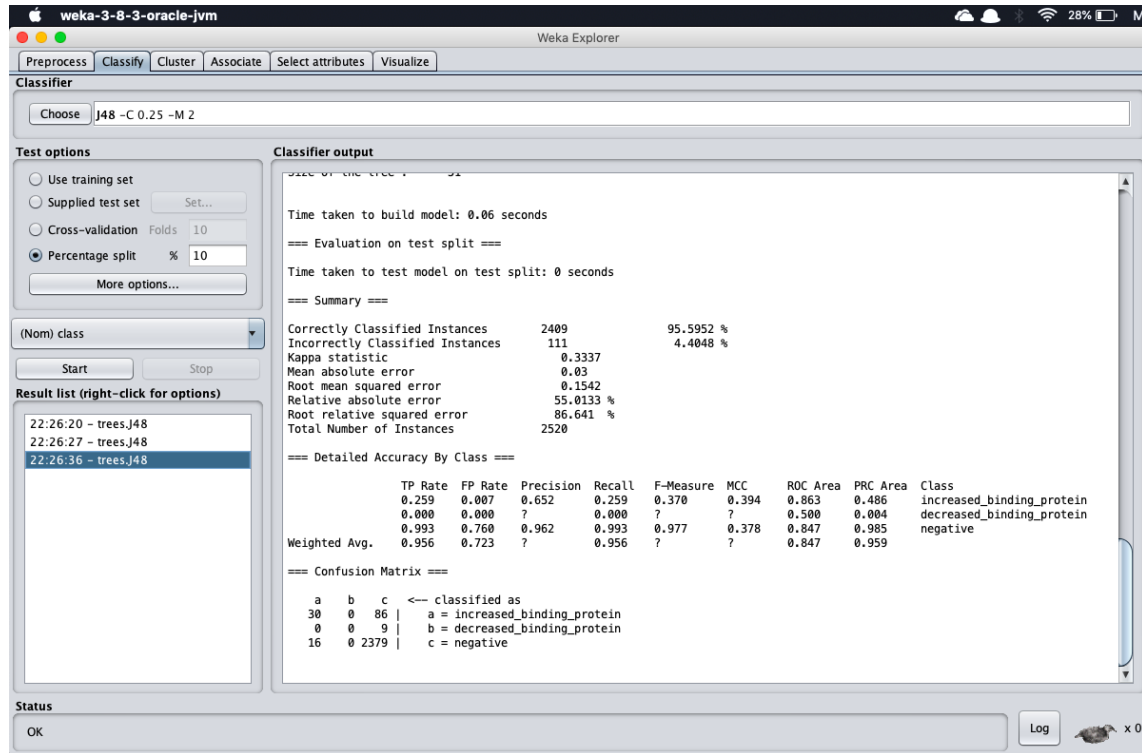


Figure 6: 10% Training Dataset

## 1.4   What are the effects of the different sizes of training and testing datasets on the learned models and accuracy?

Data is necessary for training to predict the accurate results so we reduce the training data percentage then we saw a gradual decline of the precision. In order to produce near accurate results, more training data is required.

Also, we noticed that the correctly classified instances were getting increased as the training percentage was getting increased.

## 1.5 Compare the different precision results between the two datasets and explain the change of accuracy.

With change in training dataset, there will always be change in precision and accuracy altogether. We see the decline of precision and accuracy when the training dataset was reduced.

| Dataset | Precision | | |
|---------|-----------|--------|--------|
| | 50 % Training | 25 % Training | 10 % Training |
| Car | 0.998 | 0.993 | 0.968 |
| Health | 0.965 | - | - |

Figure 7: Comparison

## 1.6 For the first dataset, use the safety attribute as a class, how accurate is the model with the same settings (50-50%, 25-75%. 10-90%)
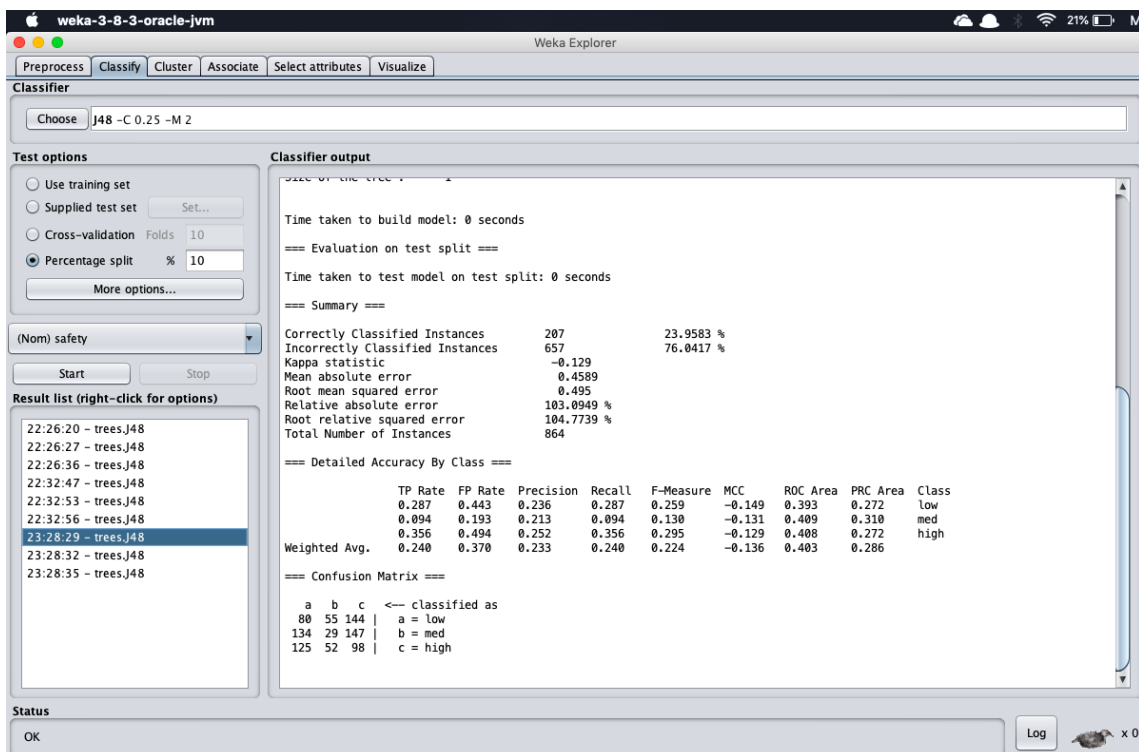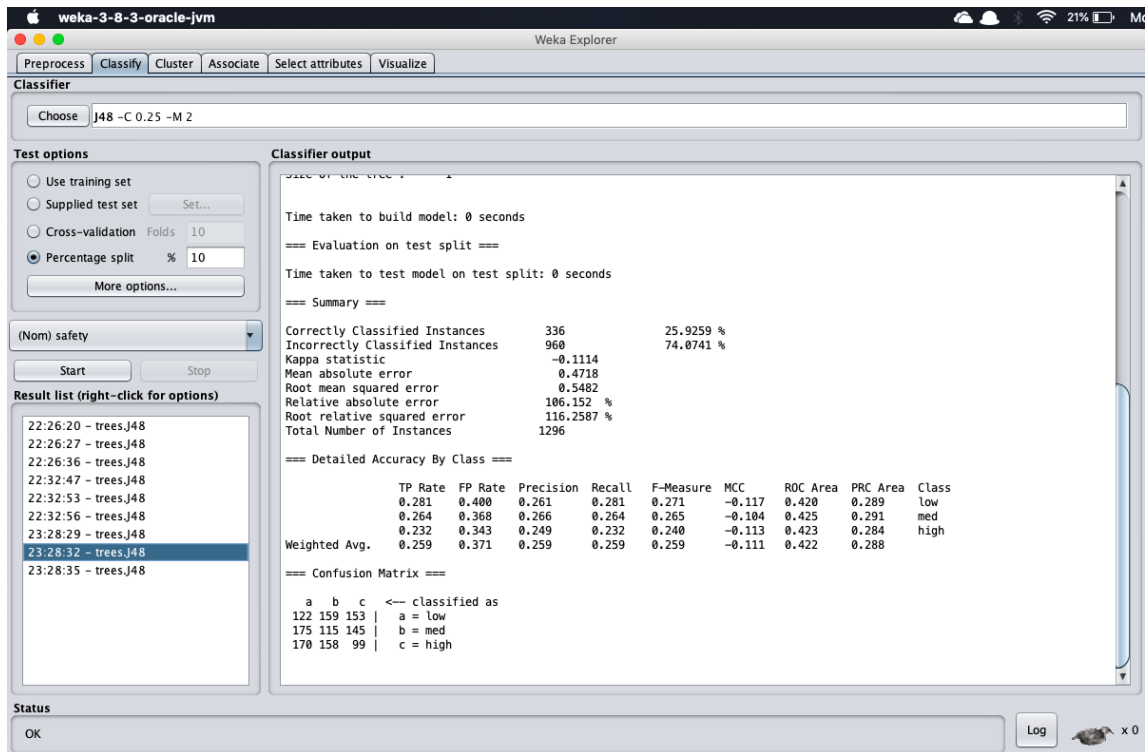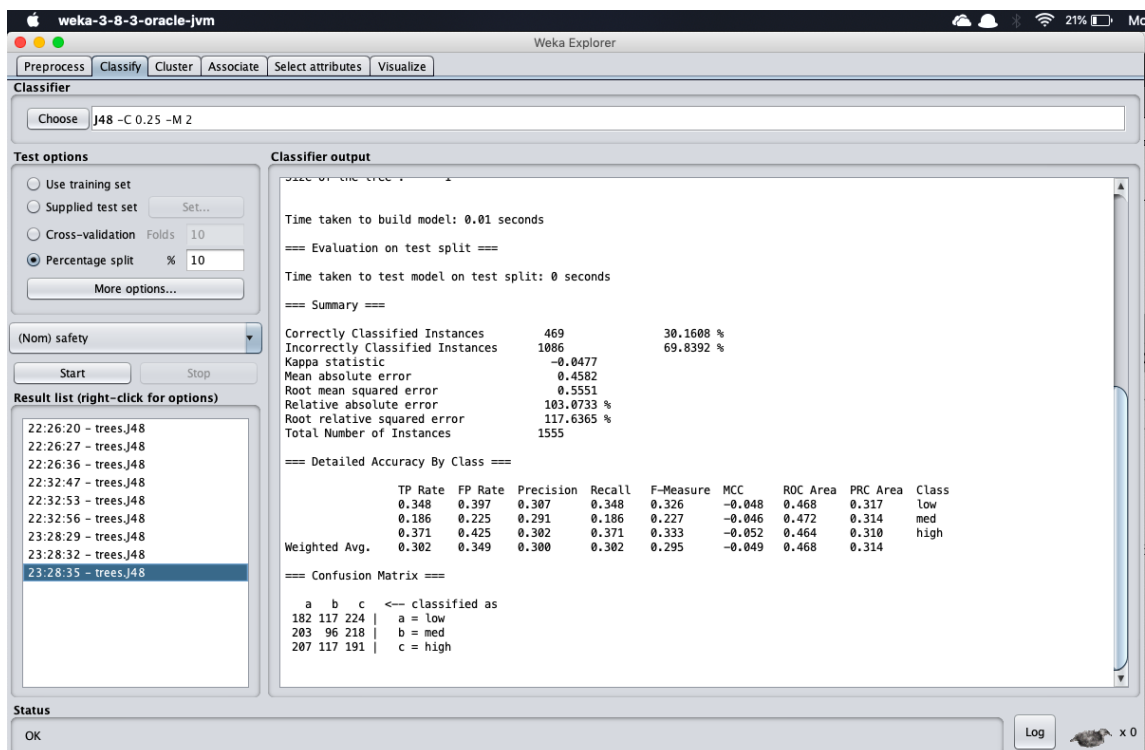


Figure 8: 50% Training Dataset

Figure 9: 25% Training Dataset



Figure 10: 10% Training Dataset