# BIRZEIT UNIVERSITY

Faculty of Engineering & Technology – Computer Systems Engineering Department

Second Semester 2023-2024

## Artificial Intelligent ENCS3340

---

Project #2:
Machine Learning for Classification

---

**Prepared by**:

Rasha Daoud - 1210382

Nadia Thaer - 1210021

**Instructor:** Aziz Qaroush

**Section:** 3

**Date:** 18st June 2024

# Table of Contents

# List of Figures:

# Introduction

## The Dataset used:

**Id #1: 1210382**

**Id #2: 1210021 ➡ 1%3=1**

As requested, based on the last digit of the lowest student ID number in the team, we select dataset #1 which is about "Early stage diabetes risk prediction Dataset".

## The Attributes:



Figure 1: Dataset on WEKA.

Number of attributes= 17, Number of instance= 520

The set of attributes are listed in the second box.

## Dataset Description

The "Early Stage Diabetes Risk Prediction Dataset" consists of health-related attributes aimed at predicting the risk of diabetes. The dataset includes the following attributes:

1. **Age**: Age of the patient (years).
2. **Gender**: Male or Female.
3. **Polyuria**: Frequent urination (Yes/No).
4. **Polydipsia**: Excessive thirst (Yes/No).
5. **Sudden weight loss**: Sudden weight loss (Yes/No).
6. **Weakness**: General body weakness (Yes/No).
7. **Polyphagia**: Excessive hunger (Yes/No).
8. **Genital thrush**: Yeast infection (Yes/No).
9. **Visual blurring**: Blurred vision (Yes/No).
10. **Itching**: Itching (Yes/No).
11. **Irritability**: Irritability (Yes/No).
12. **Delayed healing**: Delay in wound healing (Yes/No).
13. **Partial paresis**: Muscle weakness (Yes/No).
14. **Muscle stiffness**: Muscle stiffness (Yes/No).
15. **Alopecia**: Hair loss (Yes/No).
16. **Obesity**: Obesity (Yes/No).
17. **Class**: Diabetes risk (Positive/Negative).

The goal is to predict the "Class" attribute based on the other features.

# Discretization

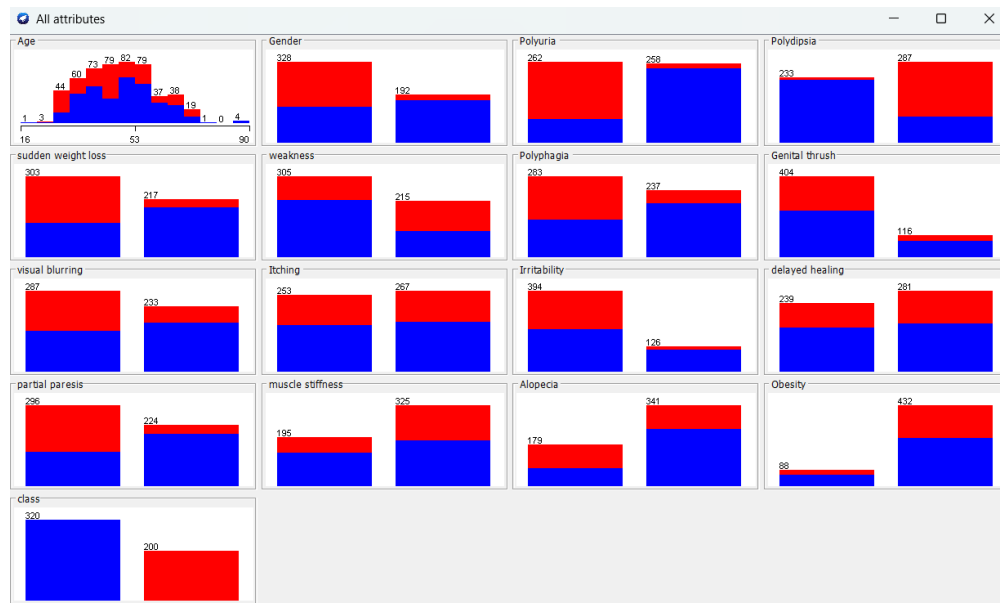In all three branches, the process of discretization in the first part is the same:



Figure 2: Distribution for all attributes before Discretization.

For the first attribute "Age" it is the only continuous attribute, so we can discretize it:
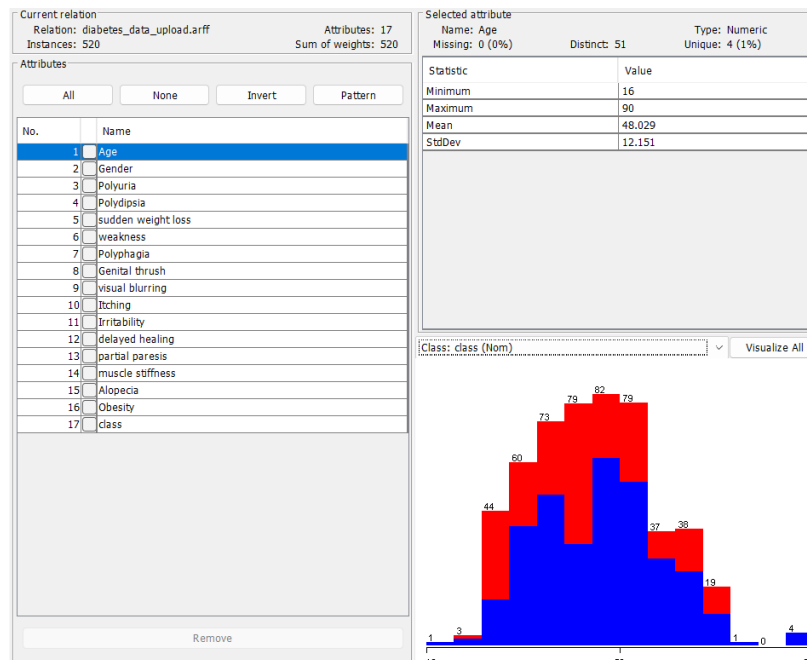


Figure 3: Age attribute before discretization.

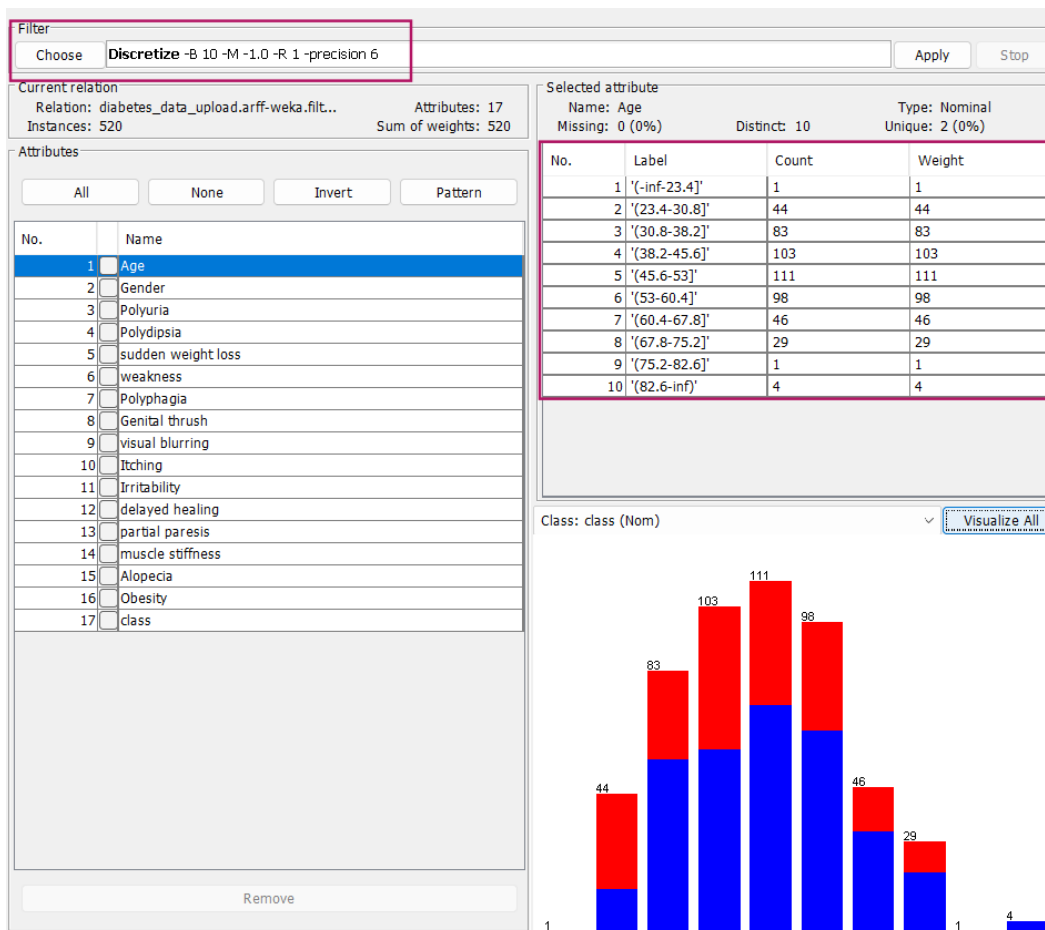Applying Discretize filter on "Age" attribute using 10 bins, give us the following result:

Figure 4: Age attribute after Discretization with 10 bins.

The discretize filter applied to the attribute shown above in the first box. And in the second one, the new discrete values and their count are shown in detail.

# Decision Tree

## 5-folds cross validation.

Using the 5-fold cross validation to train the model, and applying Naïve Bayes as a classifier, the result of cross matrix, accuracy, recall, precision, and f-source were calculated as shown below.
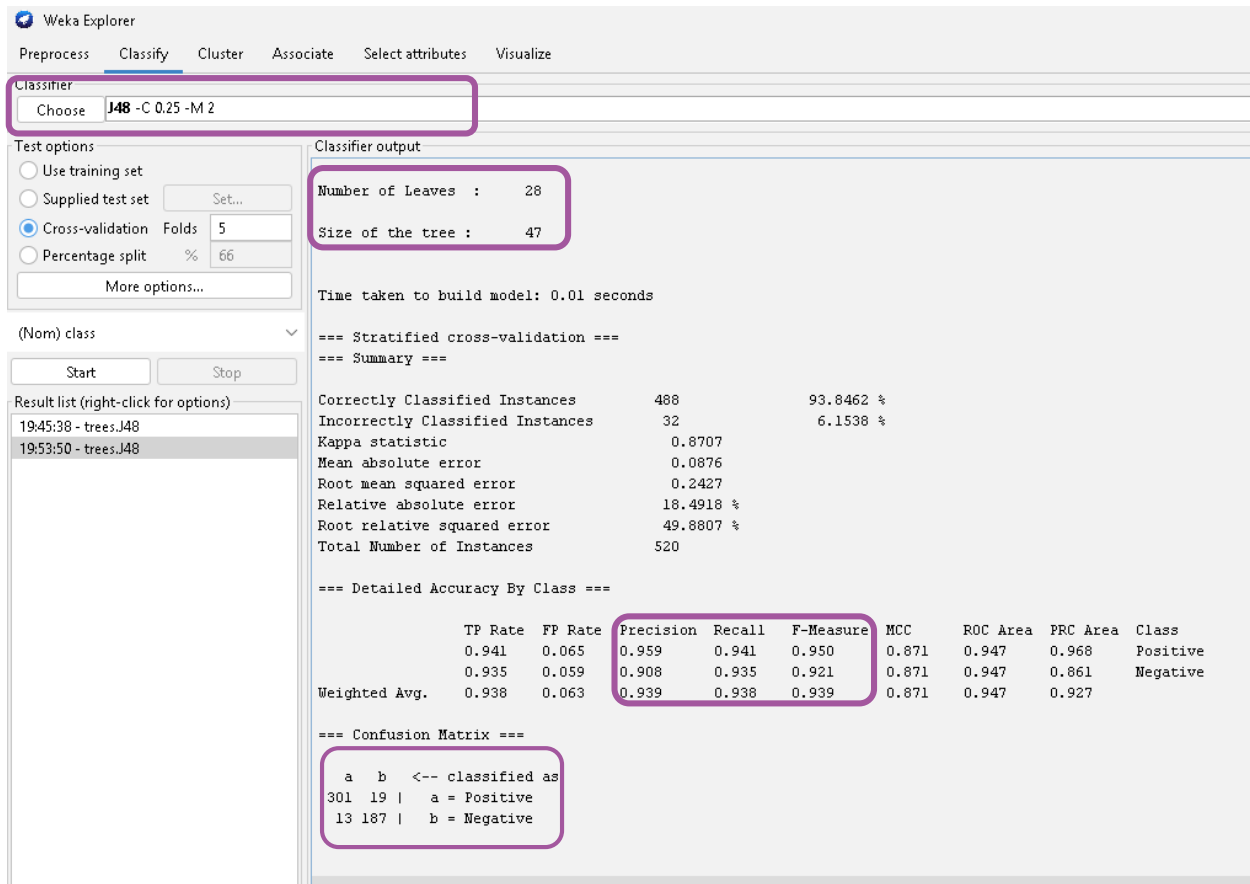


Figure 5: Test Decision Tree using 5-fold cross validation.

**From the Result of Classify**

- The Number of Leaves: 28, Size of the tree: 47
- The Classifier model use is → j48 (Type of Decision tree )
- From the figure, it is shown that the correctly classified instances are 93.8462 %
- The Result Found from the Confusion Matrix:

$$TP = 301, FN = 19, FP = 13 \text{ and } TN = 187.$$

As the precision is 0.959, recall 0.941, F-Measure = 0.950. These results could be acceptable.
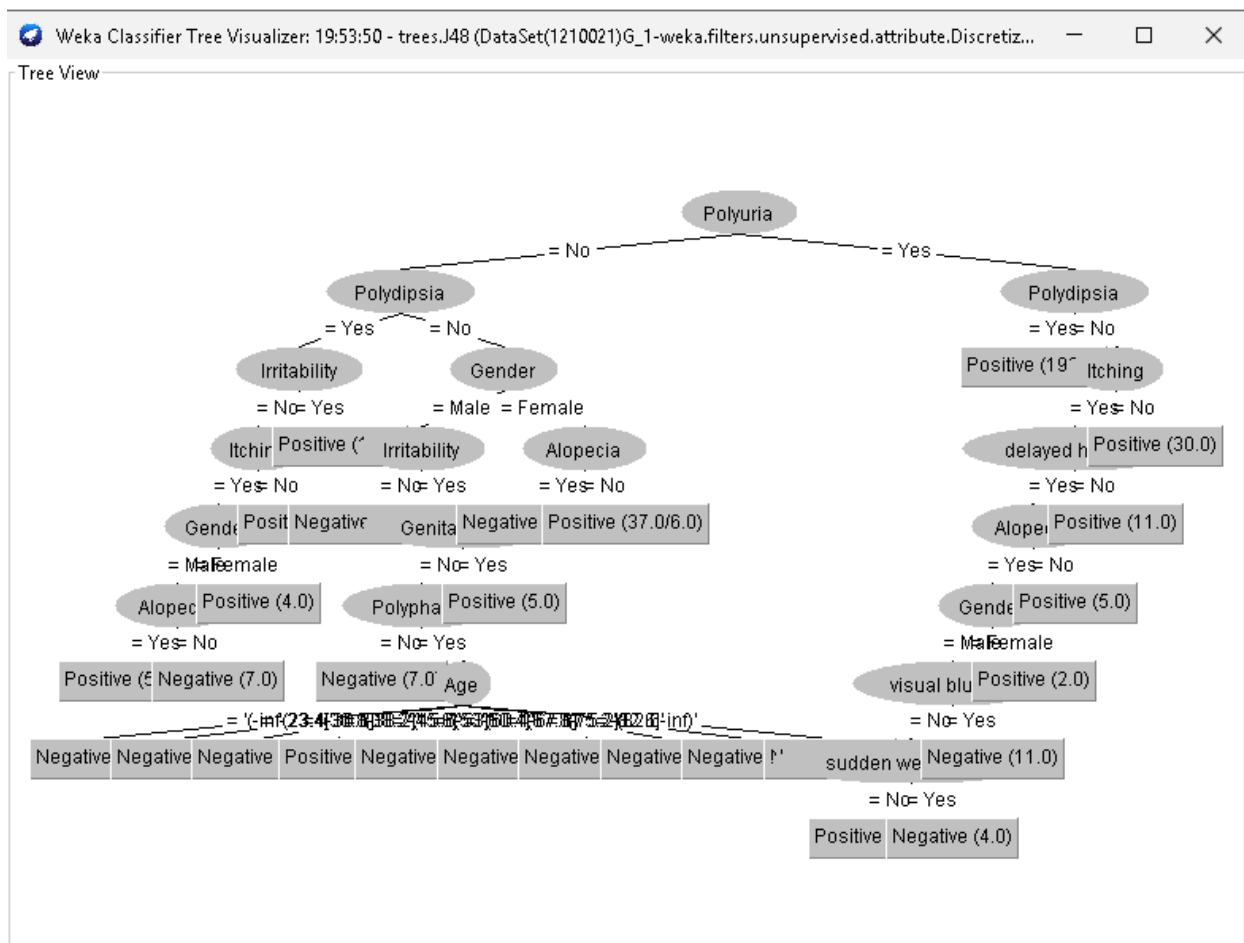
## Tree View:



Figure 6: Tree, bins=10

# Changing hyper-parameters.

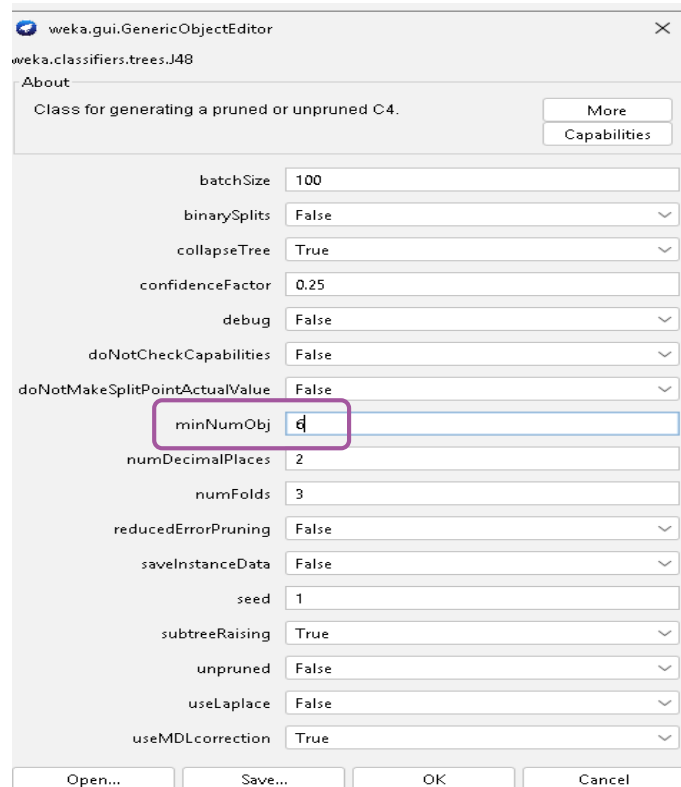Changing was done on 10 bins data output, and the MinNumObj was 2 Then set 6



Figure 7:Change minNumObj to 6.



Figure 8: Result After Change minNumObj in Decision Tree.

**Notice the changes:**

- The Number of Leaves: 11, Size of the tree: 21
- From the figure, it is shown that the correctly classified instances are 90.5769 %
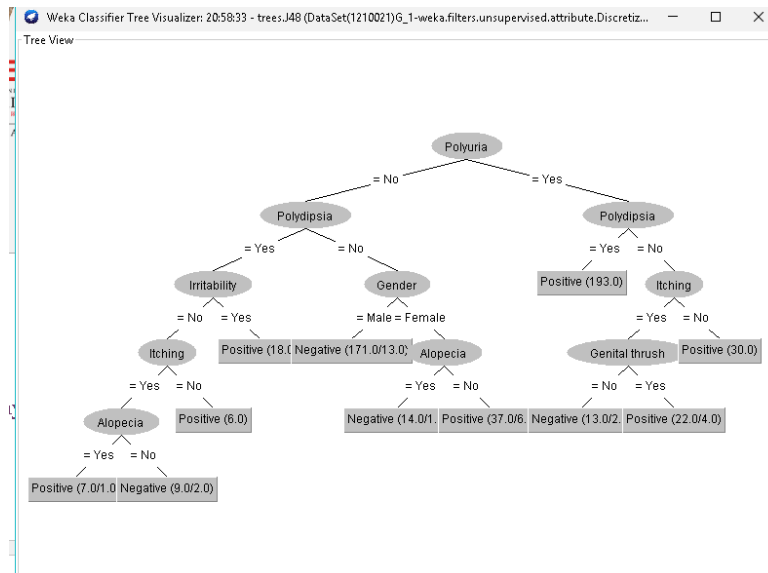- The Result Found from the Confusion Matrix: TP =291, FN =29, FP = 20 and TN = 180.



Figure 9: Tree after changing minNumObj.

As the precision is 0.936, recall 0.909, F-Measure = 0.922.

These results could be acceptable.

Notice when increase the number of (minNumObj) the correctly classified instances are will decrease and this will effect on the Confusion Matrix.

# Naïve Bayes

## 5-folds cross validation.

Using the 5-fold cross validation to train the model, and applying Naïve Bayes as a classifier, the result of cross matrix, accuracy, recall, precision, and f-source were calculated as shown below.
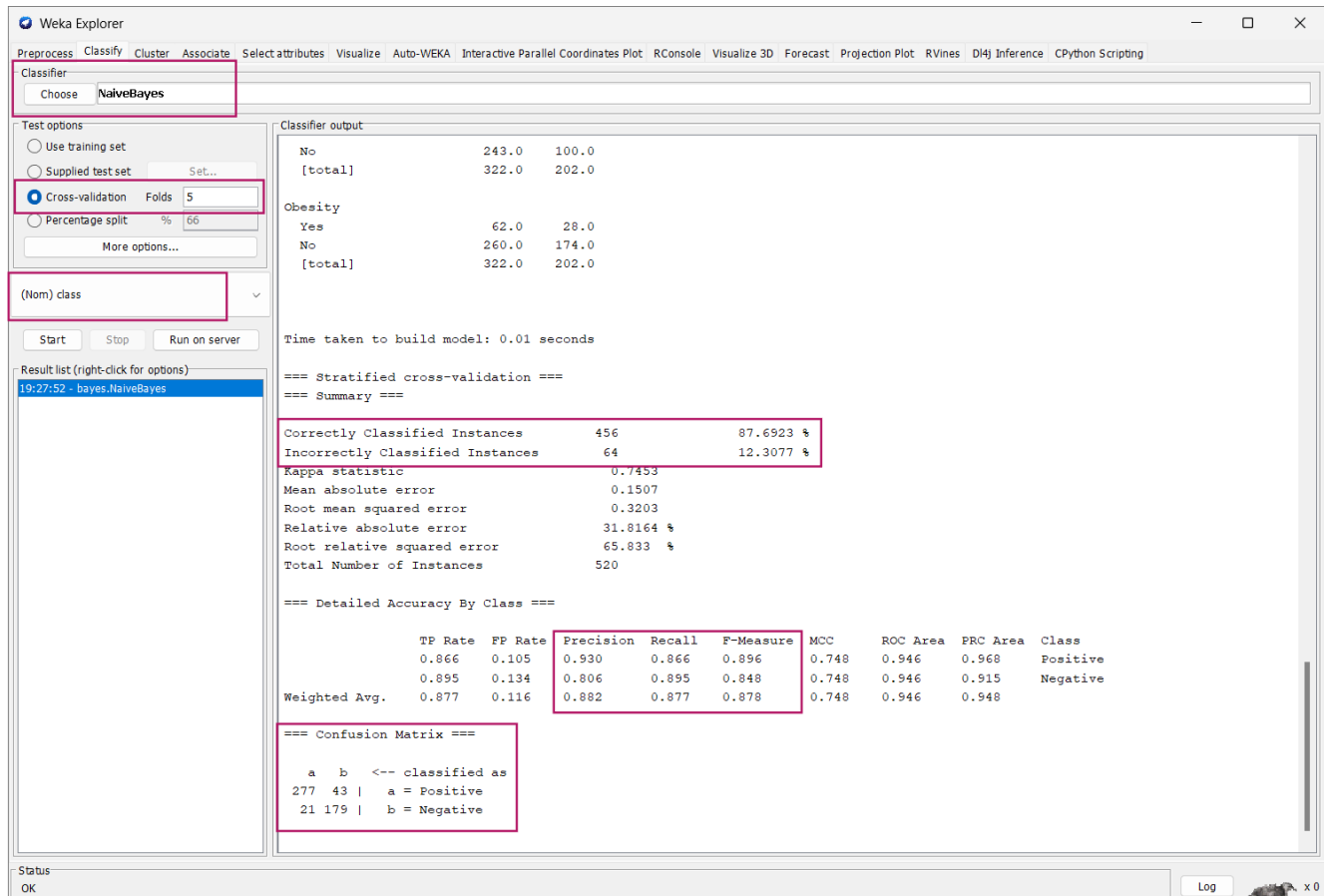


Figure 10: The result of using Naive Bayes as Classifier.

From figure above, from confusion matrix part. The matrix illustrates that, TP =277, FN =43, FP =21, and TN = 179. From these values, precision, recall and f- measure values can be calculated. The value of the precision is 0.882, Recall is 0.877 and f-measure = 0.878 all in average, and they specify for both positive and negative class above. Also, from the figure the accuracy was equals 87.69%, Which is considered good.

# Changing hyper-parameters.

Applying changing on batchSize from 100 to 10000, and the doNotCheckCapabilities to true, and the useSupervisedDiscretization to true.



Figure 11:The result after change hyper-parameters.

The values after changing hyper-parameters does not change from the values before changing hyper-parameters as shown in the figure above.

# MLP

## 5-folds cross validation.

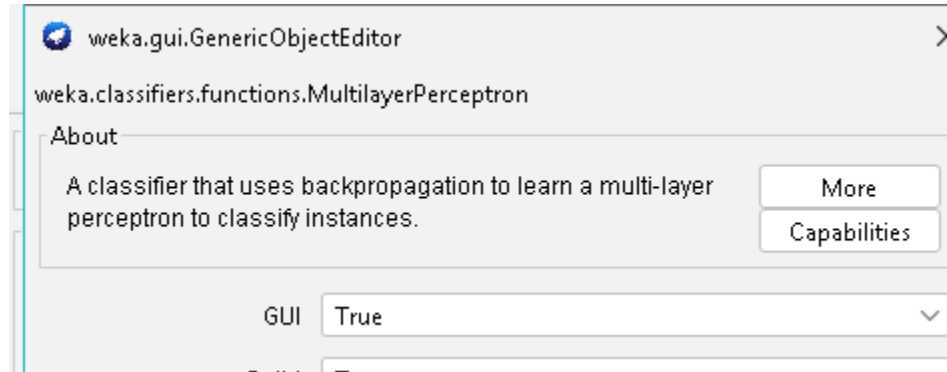To apply the 5 –fold cross validation on the MLP the setup will be:



Figure 12: Setup for applying 5-fold cross validation.

When starting the classify, the following is the result:



Figure 13: MLP 5-folds cross validation result part-1.

Figure 14: MLP 5-folds cross validation part-2.

From this result:

The classification model achieved an overall accuracy of 94.62%, correctly classifying 492 instances out of 520. Specifically, it accurately identified 302 instances of the 'Positive' class, while 10 instances were incorrectly classified as 'Negative'. Conversely, for the 'Negative' class, the model correctly classified 190 instances, with 18 instances misclassified as 'Positive'. These results reflect a high level of precision and recall for both classes, indicating robust performance in distinguishing between the classes based on the given matrix:       TP =302, FN =18, FP = 10 and TN = 190.

So The Result Could Be acceptable.

**To show the NN set the GUI True the Neural Network:**



Figure 15: Neural Network.

# Changing hyper-parameters

By Changing the hidden layers From **a ➔2**



Figure 16: Before changing the parameter.



Figure 17: After changing the parameter.

## The NN become:



Figure 18: Neural Network  After Changing.

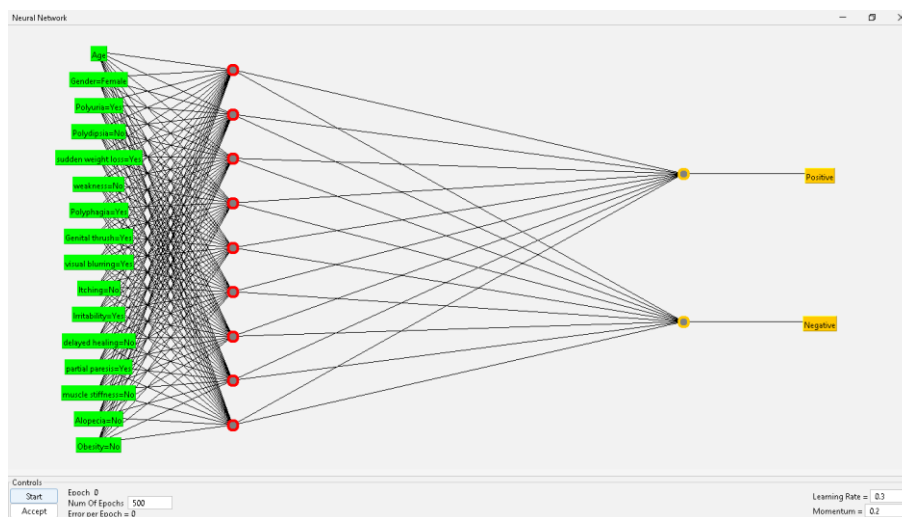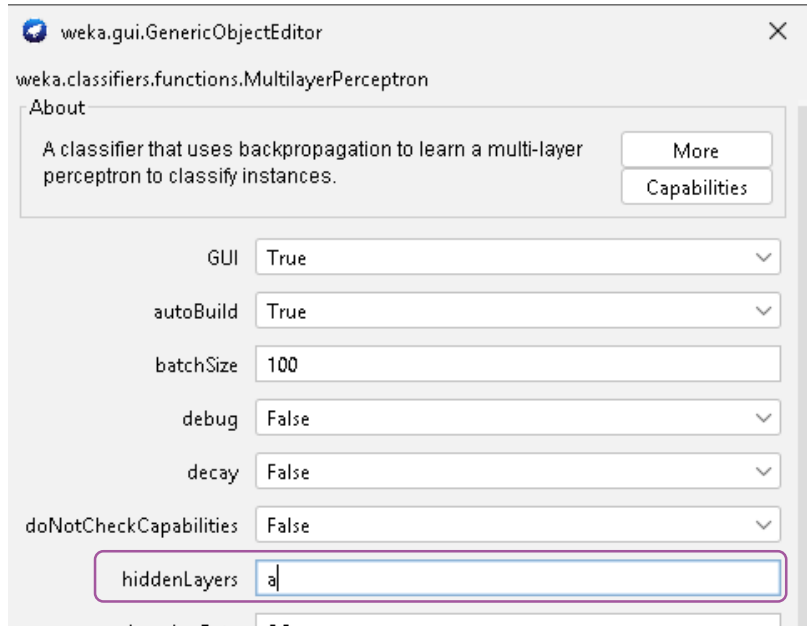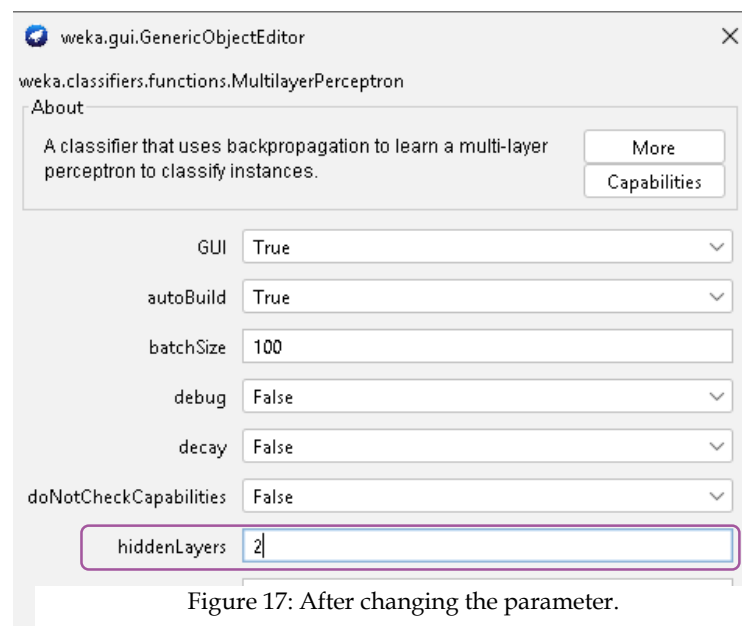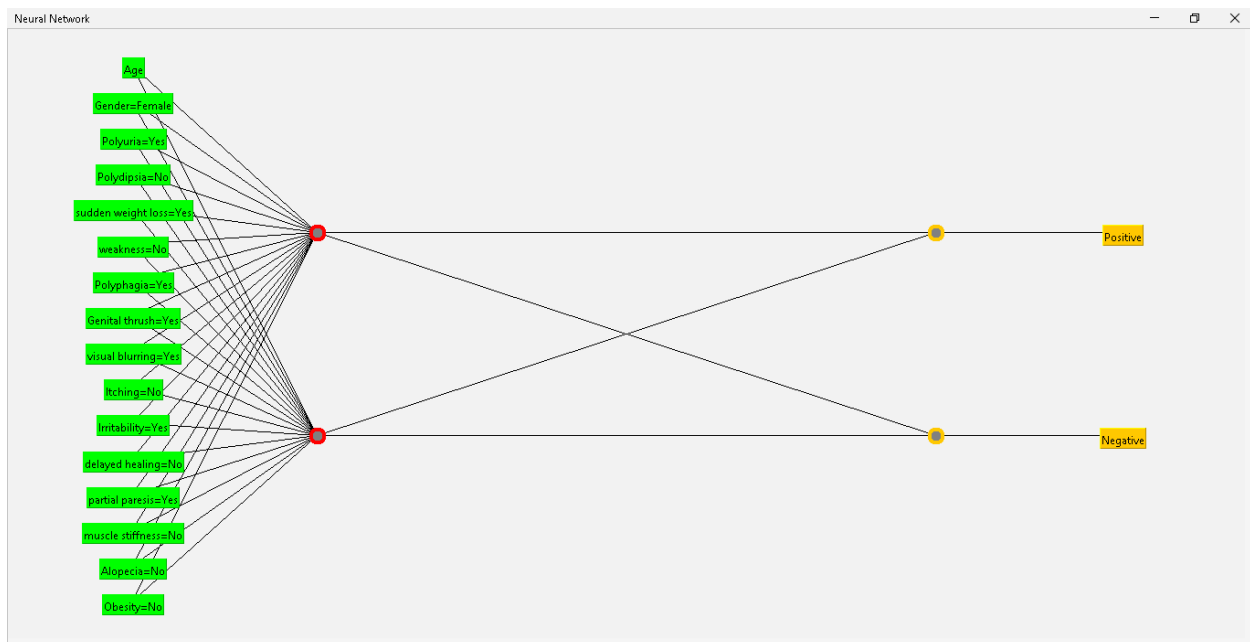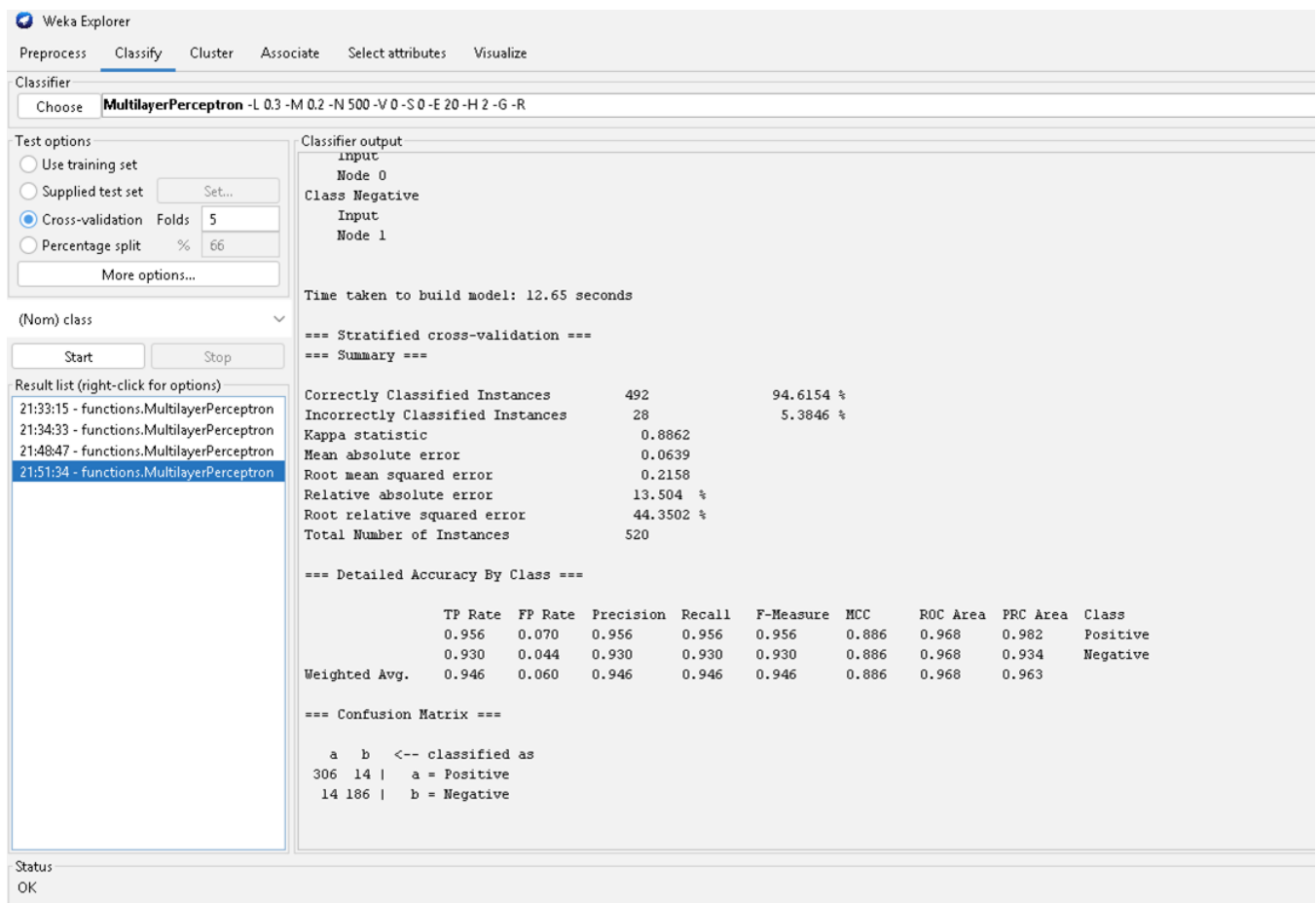Figure 19: The result of using MLP as classifier after change.

**From the result of classify:**

The model achieved an overall accuracy of 94.62%, correctly classifying 492 instances. Specifically, it accurately identified 306 instances of the 'Positive' class (indicating diabetes) and 186 instances of the 'Negative' class (indicating no diabetes). The model's performance metrics, including precision, recall, and F-measure, were consistently high for both classes, reflecting its effectiveness in distinguishing between individuals with and without diabetes based on the provided attributes.

TP =306, FN =14, FP = 14 and TN = 186.

As the precision is 0.956, recall 0.956, F-Measure = 0.956. These results could be acceptable.

**The effect of changing:**

Both results show good accuracy in identifying diabetes. Result Two found slightly more cases of diabetes and had fewer misses, but result one made fewer mistakes by falsely saying someone had diabetes when they didn't. Choosing between them would depend on whether it's more important to find all cases of diabetes or to avoid saying someone has diabetes when they don't.

## conclusion

Among the three classifiers, the optimized Multi-Layer Perceptron (MLP) demonstrated the highest performance in terms of accuracy, precision, recall, and F-measure, making it the most effective model for this diabetes classification task. The Decision Tree also showed strong performance, especially in terms of precision, but was more sensitive to hyper-parameter adjustments. Naïve Bayes, while stable, had slightly lower performance metrics compared to the other two methods. Overall, the choice of classifier can depend on the specific requirements of the task at hand. If achieving the highest possible accuracy and balanced performance is crucial, the MLP is the preferred choice. For scenarios where interpretability and stability are important, the Decision Tree and Naïve Bayes classifiers offer viable alternatives.

confidence factor is equal to 0.75. That test in this project has the best accuracy by 87.11% and with F-measure equals 0.874. In the Naïve Bayes algorithm, the results were different. They were the worst, as the accuracy equals 67.77% and the F-measure equals 0.753. In the last algorithm 'The Hoeffding Tree algorithm", the results were good, since the accuracy was equal to 85.33% and the F-measure to 0.858, even after playing with the parameters, the results were still good.