# DECISION TREE AND RANDOM FOREST FOR THE ADVENTUREWORKDW DATASET

**Name:** Oscar Hernandez Mata

In this project I wanted to predict whether a costumer would be classified as a potential bike buyer or a no-bike buyer. Hence, I used the AdventureWorkDW dataset from MySQL, which I modified by selecting the necessary columns for this project.

We converted some of the integer variables to factors in order to create a training set that contains 67% of the data and a training set that contains the rest 33% by using the hold out method.

Next, we created a subset of the training data (TrainTree) to limit the number of columns in the decision tree. Then we repeated the process for the testing data (See Figure 1).
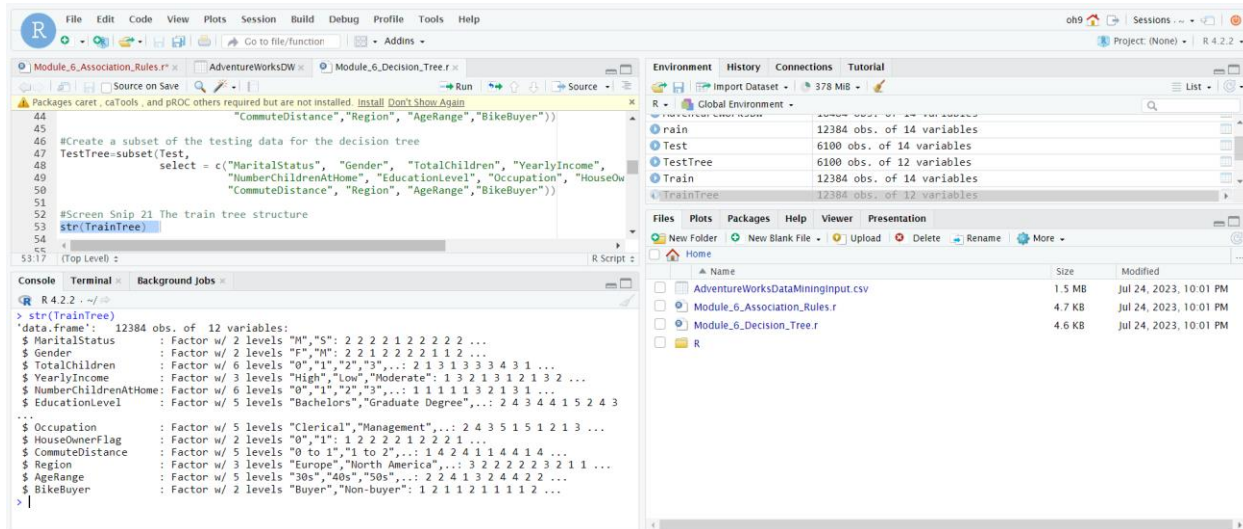


**Figure 1**

The subsequent phase in the classification procedure involved constructing the decision tree model using the training subset, with "Bike Buyer" designated as the class label. Following this, the model was evaluated by presenting it with the testing data.

To assess the effectiveness of our model, a confusion matrix was generated. This matrix provides a detailed breakdown of the true positives, true negatives, false positives, and false negatives. Additionally, measures of sensitivity and specificity were calculated. The confusion matrix is crucial as it offers insights into the robustness and accuracy of our model, as illustrated in Figure 2.
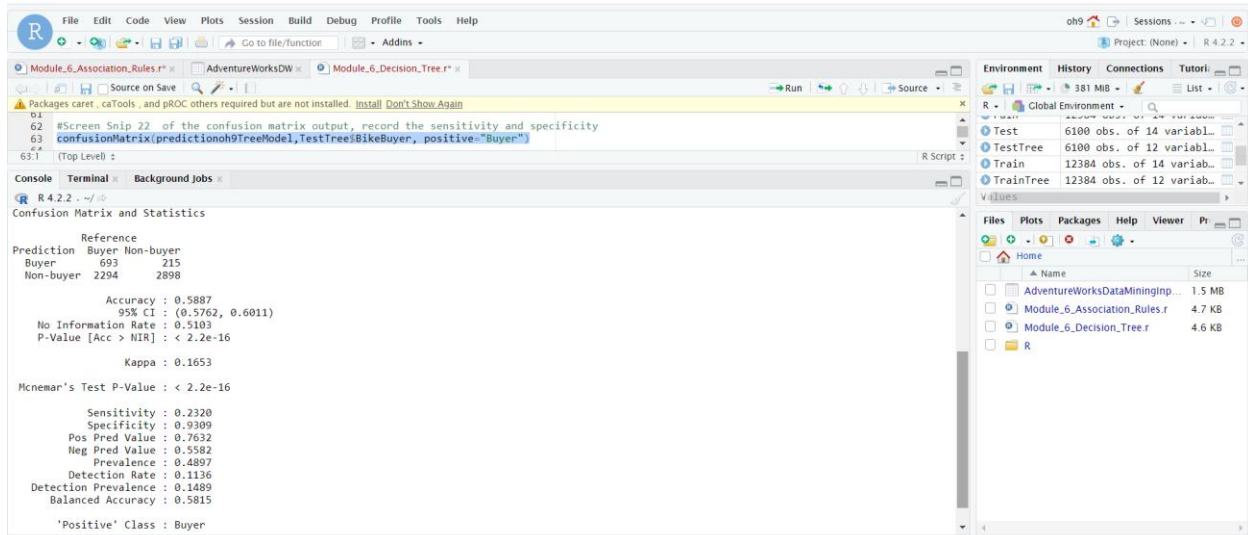
**Figure 2**

The analysis indicated that the decision tree model demonstrated a sensitivity of 0.23 and a specificity of 0.93. These results suggest that while the model is highly effective at identifying true negatives, it struggles considerably with accurately predicting true positives, as detailed in Figures 2 and 3.
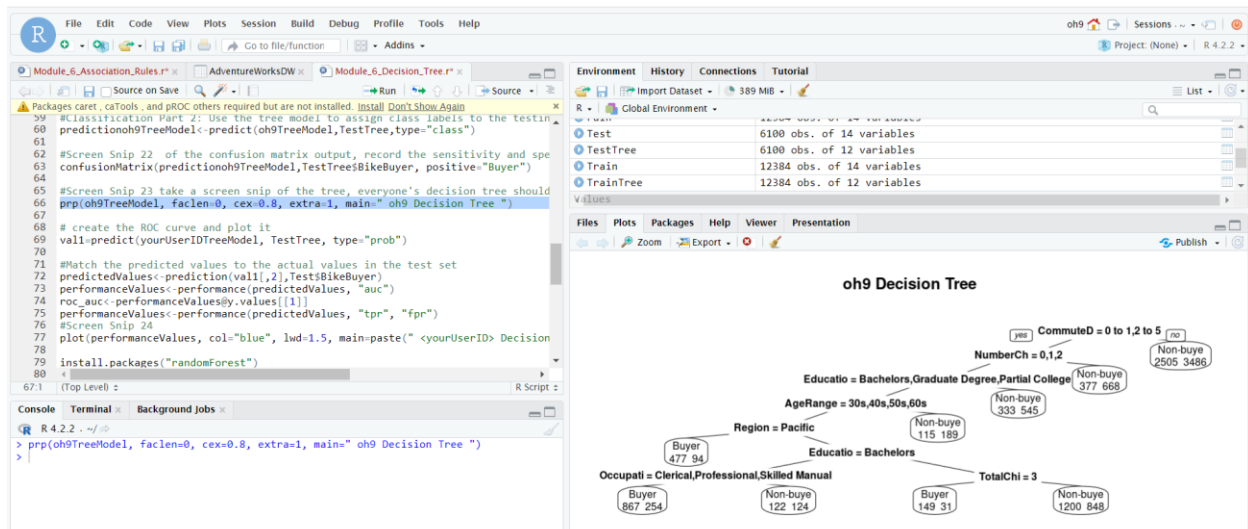


**Figure 3**

The ROC curve graphic provides an assessment of the model's predictive accuracy by comparing predicted values against actual outcomes. The AUC, or Area Under the Curve, for the decision tree model is 0.637, as shown in Figure 4. This value indicates that the model's performance in distinguishing between positive and negative instances is suboptimal, given that the AUC is significantly lower than the ideal value of 1.

**Figure 4**

Since the decision tree did not accurately predict the true positives, I first installed the randomForest and caTools packages, then created a random forest model using the randomForest command for the training data (See Figure 5). Next, I used the random forest model to assign class labels to the testing data: **predict(rf, newdata=TestTree[-12]) and able(TestTree[,12], pred)**.
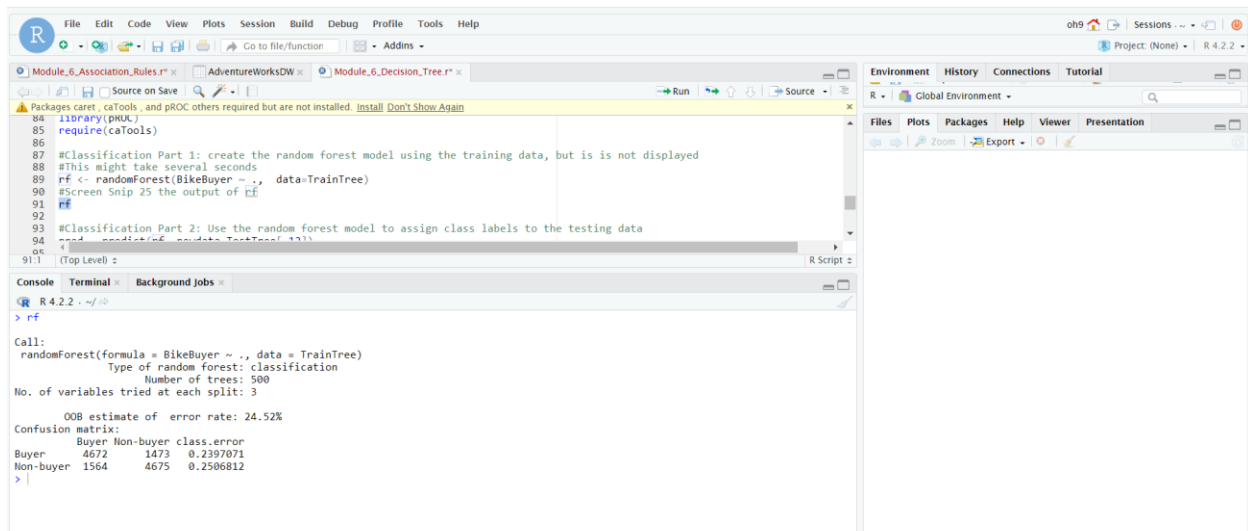


**Figure 5**

Finally I matched the predicted values to the actual values in the test set using the follwing commands: **predictions <- as.data.frame(predict(rf, TestTree, type = "prob")), predictions$predict <- names(predictions)[1:2][apply(predictions[,1:2], 1, which.max)], predictions$observed <- TestTree$BikeBuyer** (See figure 6).
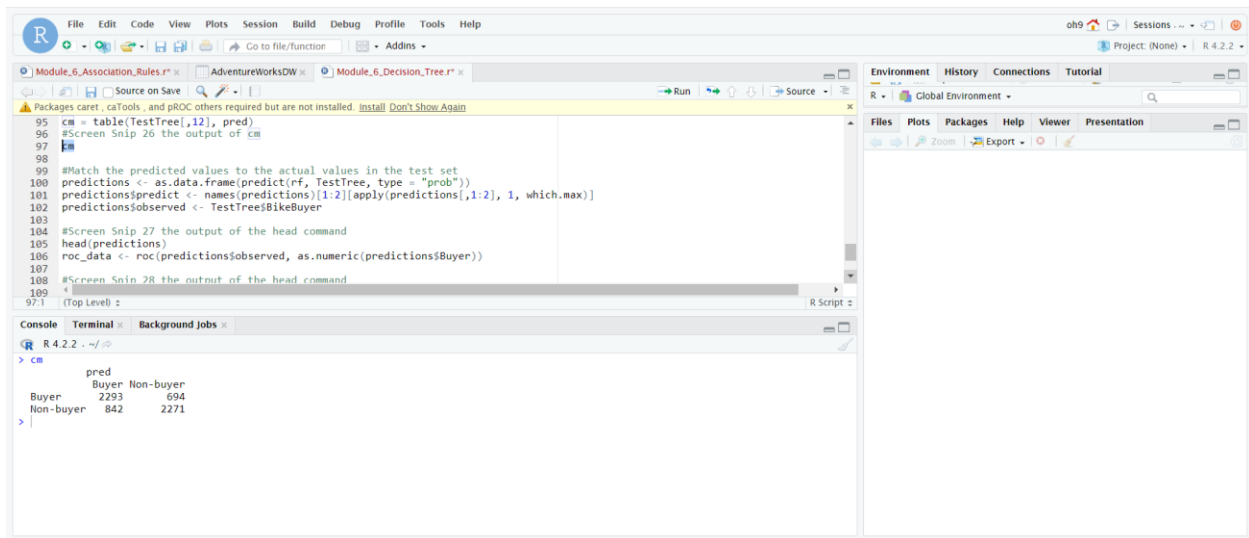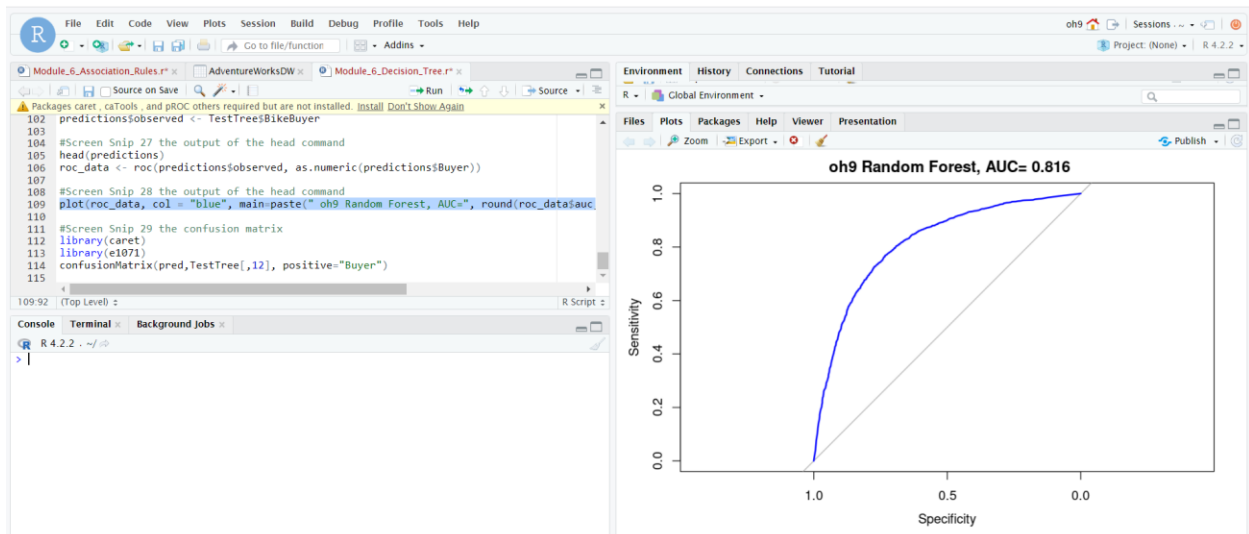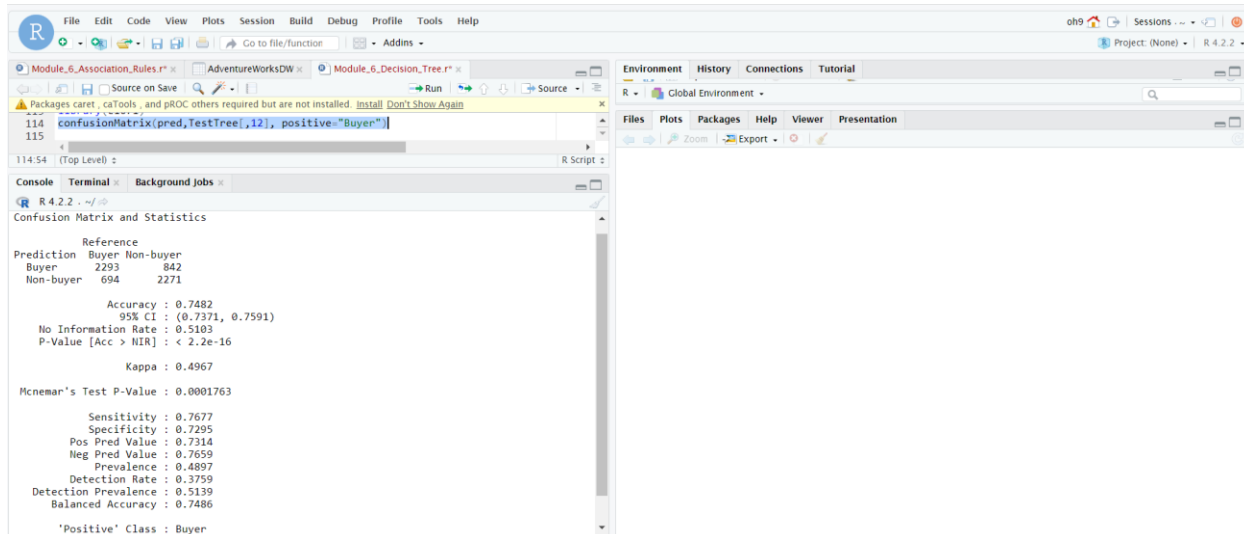
3

**Figure 6**



**Figure 7**

**Figure 8**

In the recent analysis, the performance of the random forest model was evaluated using the confusionMatrix function, yielding promising results. The model demonstrated a sensitivity of 0.76 and a specificity of 0.72, indicating a robust ability to accurately identify both true positive and true negative outcomes, as depicted in Figure 9. Further, the area under the Receiver Operating Characteristic (ROC) curve, or AUC, was calculated at 0.816. This value, illustrated in Graphic 7, underscores the model's effectiveness in distinguishing between positive and negative instances, approaching the optimal AUC value of 1.