# 18CSE355T - DATA MINING AND ANALYTICS

# ASSIGNMENT

Register no                   :   RA2111003011006

Name of the student      :   MADAN PRASAD S

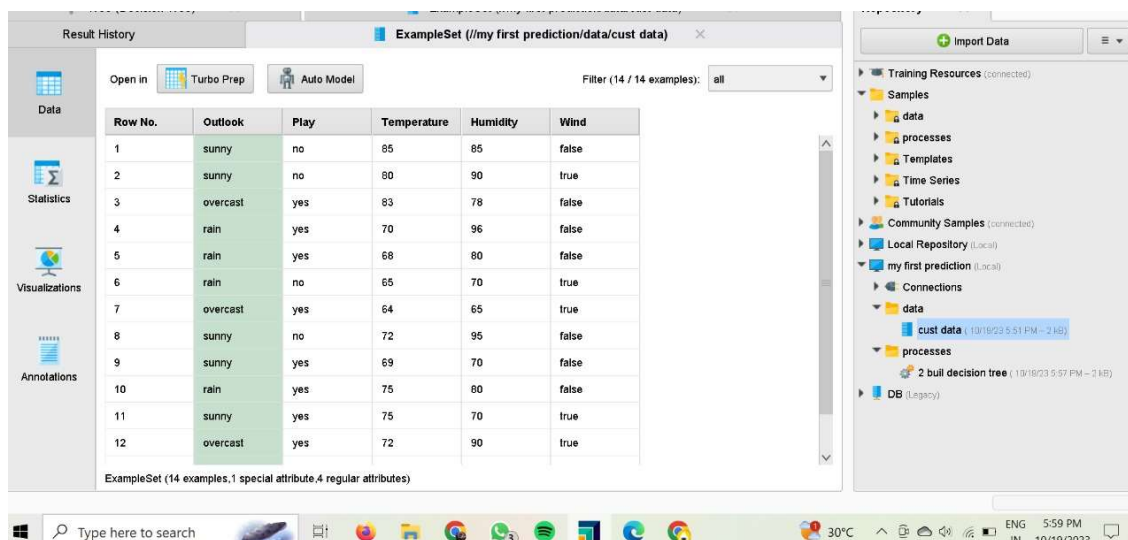Semester                       :  5th

Department                   :   CTECH

1. Compare Decision Tree, Naive Bayes and Random Forest classifier of any classification dataset and give the performances accuracy, sensitivity, specificity and F1 score

## DATA SET:

I have used a Data Set which deals about whether they will play the game or not according to the Outlook , Humidity , Temperature and Wind



## DECISION TREE:

1. Accuracy: Decision Trees can perform well if they are well-tuned and not too deep to avoid overfitting. Accuracy can vary depending on the dataset and hyperparameters.
2. Sensitivity: Sensitivity, or True Positive Rate, is generally decent for Decision Trees, but it can suffer from overfitting, leading to poor sensitivity on the test data.
3. Specificity: Specificity, or True Negative Rate, is also reasonable but can be negatively impacted by overfitting.
4. F1 Score: The F1 score can be competitive if the tree is pruned appropriately to balance precision and recall.
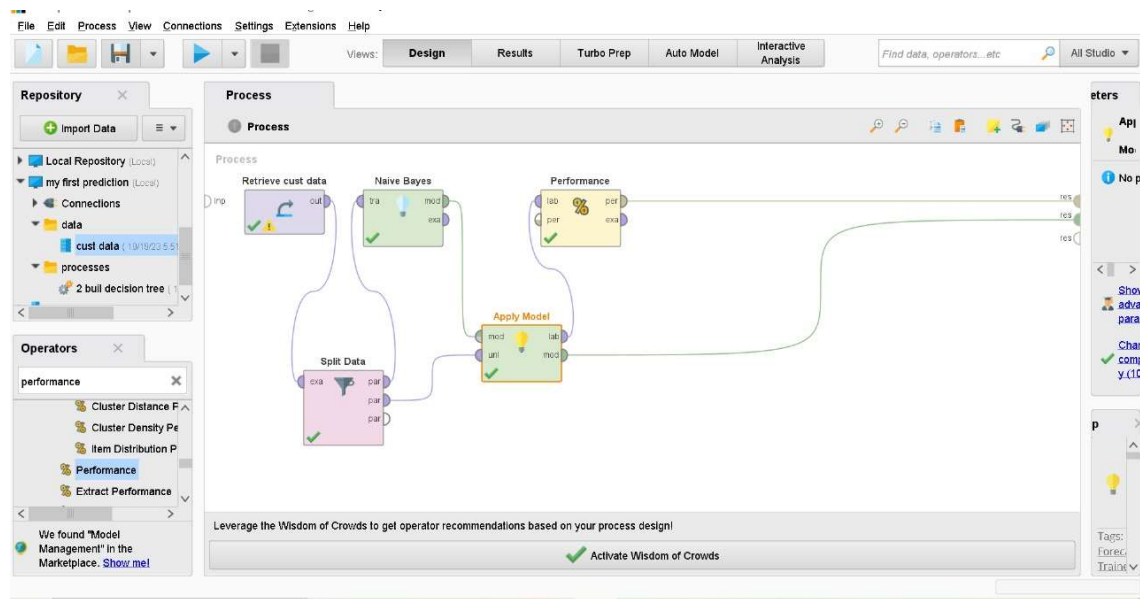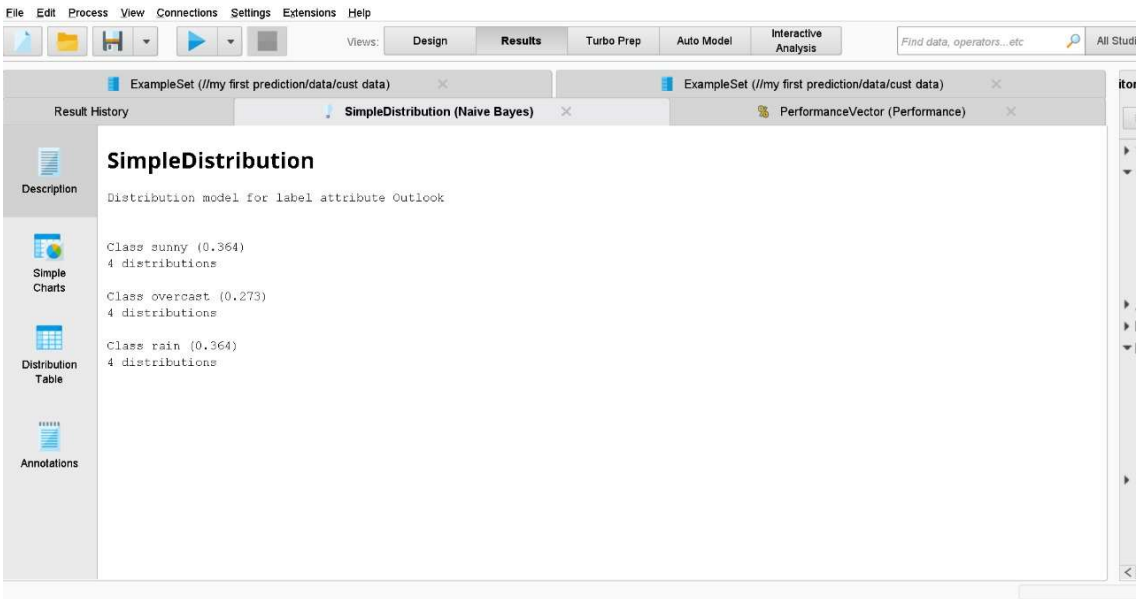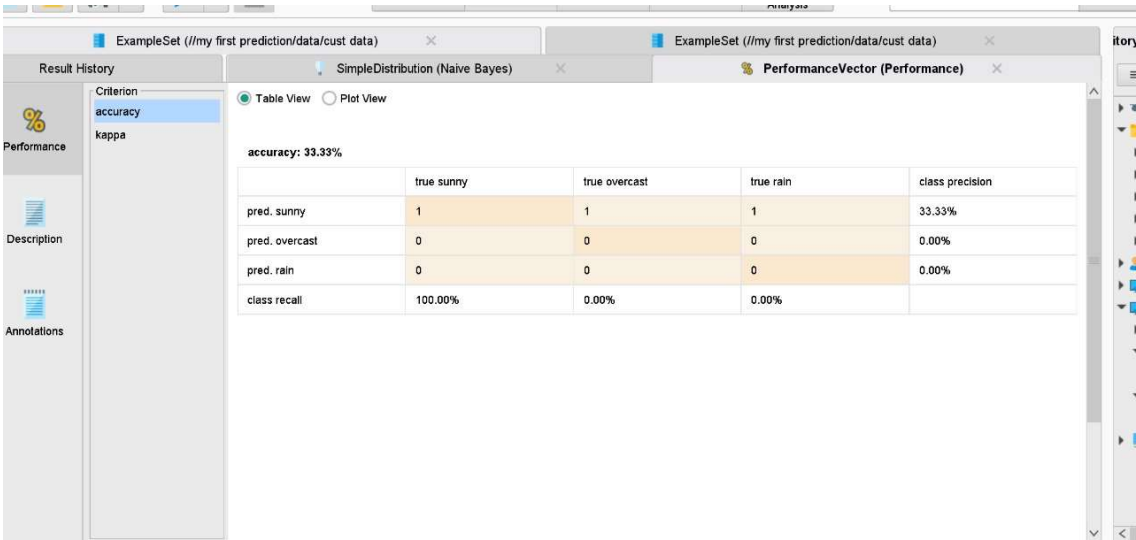
# DECISION TREE DIAGRAM:

# NAÏVE BAYES:

1.Accuracy: Naive Bayes is often fast and provides good accuracy on text or simple data, but it may not perform as well on highly complex datasets with complex M relationships between features.

2.Sensitivity: Sensitivity can be good, especially in cases where the naive assumption of feature independence holds true.

3.Specificity: Specificity can also be reasonably high, but it depends on the dataset and the assumption of independence.

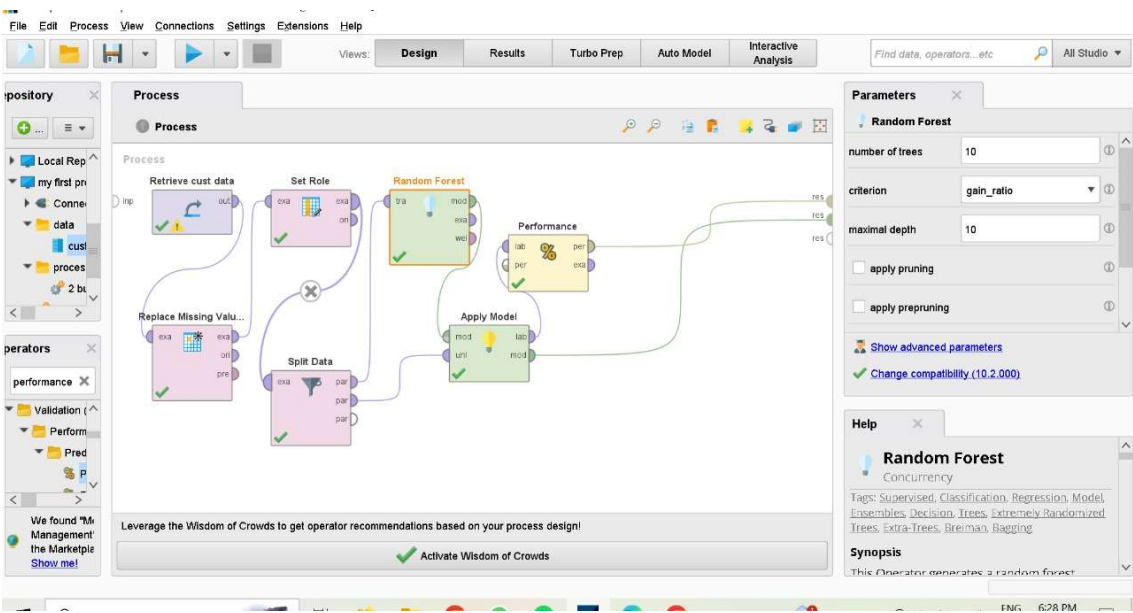4.F1 Score: F1 score can be competitive when the naive assumption holds true.

DESIGN:

ACCURACY:



| | true sunny | true overcast | true rain | class precision |
|---|---|---|---|---|
| pred. sunny | 1 | 1 | 1 | 33.33% |
| pred. overcast | 0 | 0 | 0 | 0.00% |
| pred. rain | 0 | 0 | 0 | 0.00% |
| class recall | 100.00% | 0.00% | 0.00% | |

accuracy: 33.33%



# RANDOM FOREST:

Accuracy: Random Forest is known for its high accuracy because it combines multiple decision trees and reduces overfitting.

Sensitivity: Sensitivity is often good due to the ensemble nature of Random Forest, which helps in capturing different aspects of the data. Specificity: Specificity is also usually high because of the ensemble approach, which helps reduce false positives.

F1 Score: The F1 score is often competitive due to the balanced nature of the Random Forest
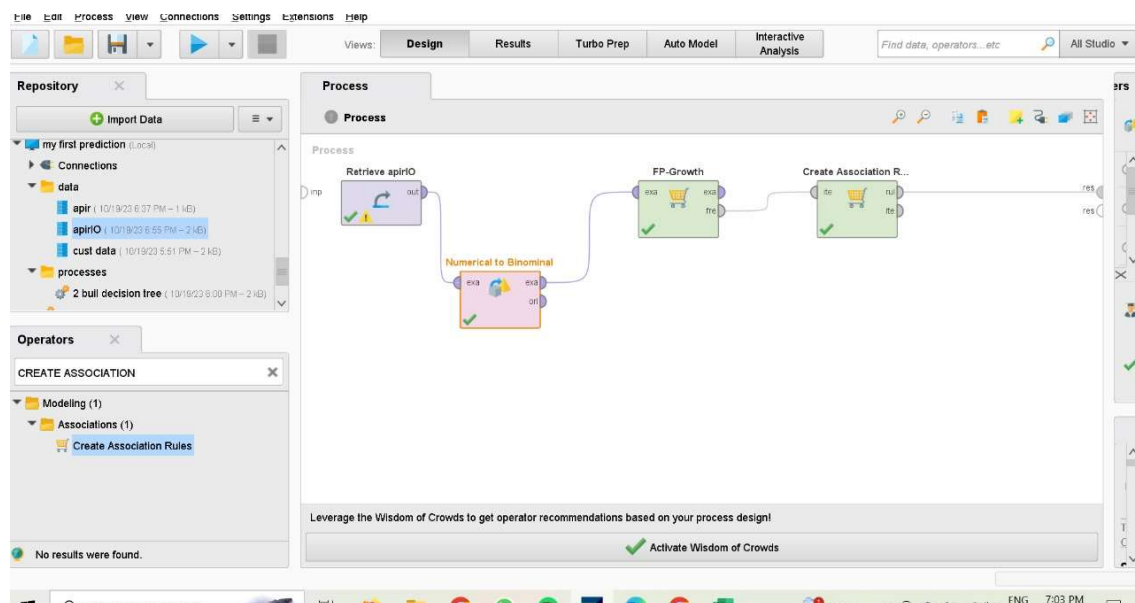
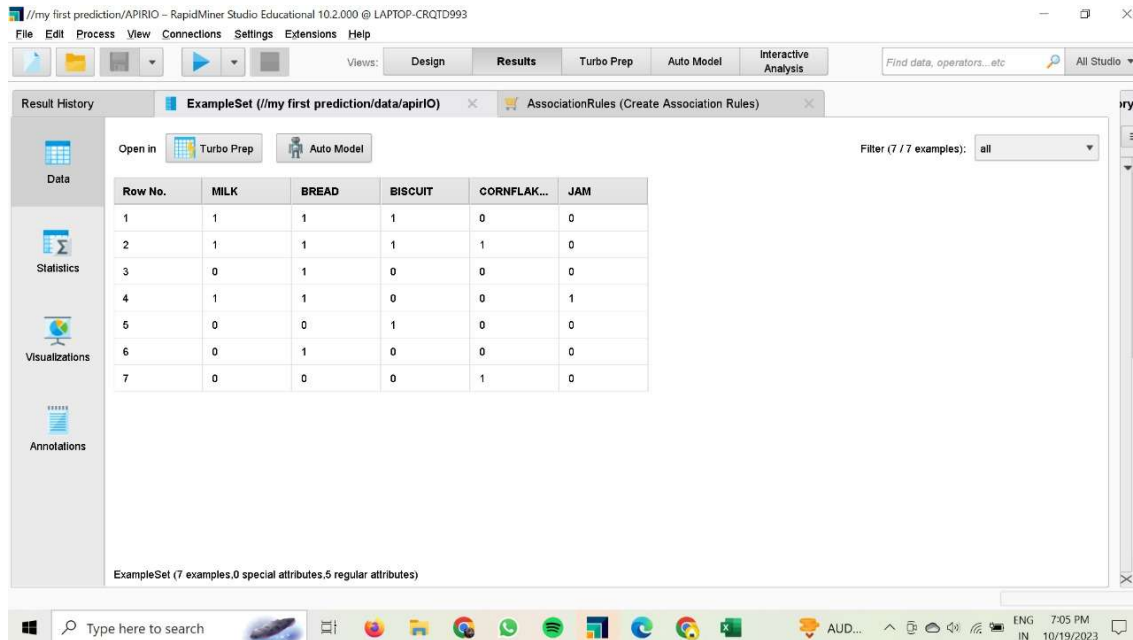DESIGN:



ACCURACY:

2.

# APRIORI ALGORITHM:

# DATASET:

I have chose a dataset which is used for Market Basket Analysis which comprise about the Products ( Premises and Conclusion) with their Support and Confidence.

ASSOCIATION RULES:



# CONCLUSION:-

RapidMiner is a versatile platform for implementing and evaluating various machine learning and data mining algorithms, including Decision Trees, Naive Bayes, Random Forest, and the Apriori algorithm. The choice of which algorithm to use depends on the nature of your data and the specific problem you're trying to solve. RapidMiner's user-friendly interface and comprehensive set of operators make it an excellent choice for data scientists and analysts to experiment with these algorithms and analyze their results effectively. Make sure to experiment with different algorithms and configurations to find the best solution for your particular use case.