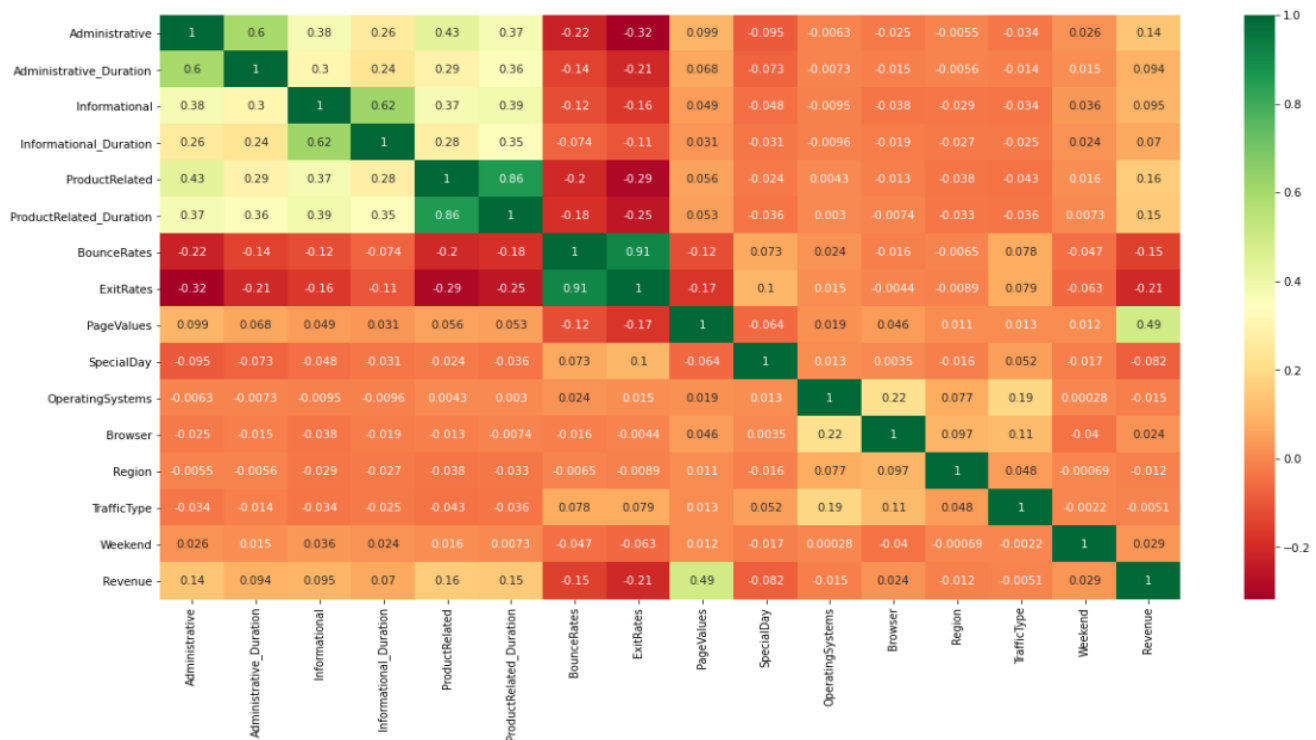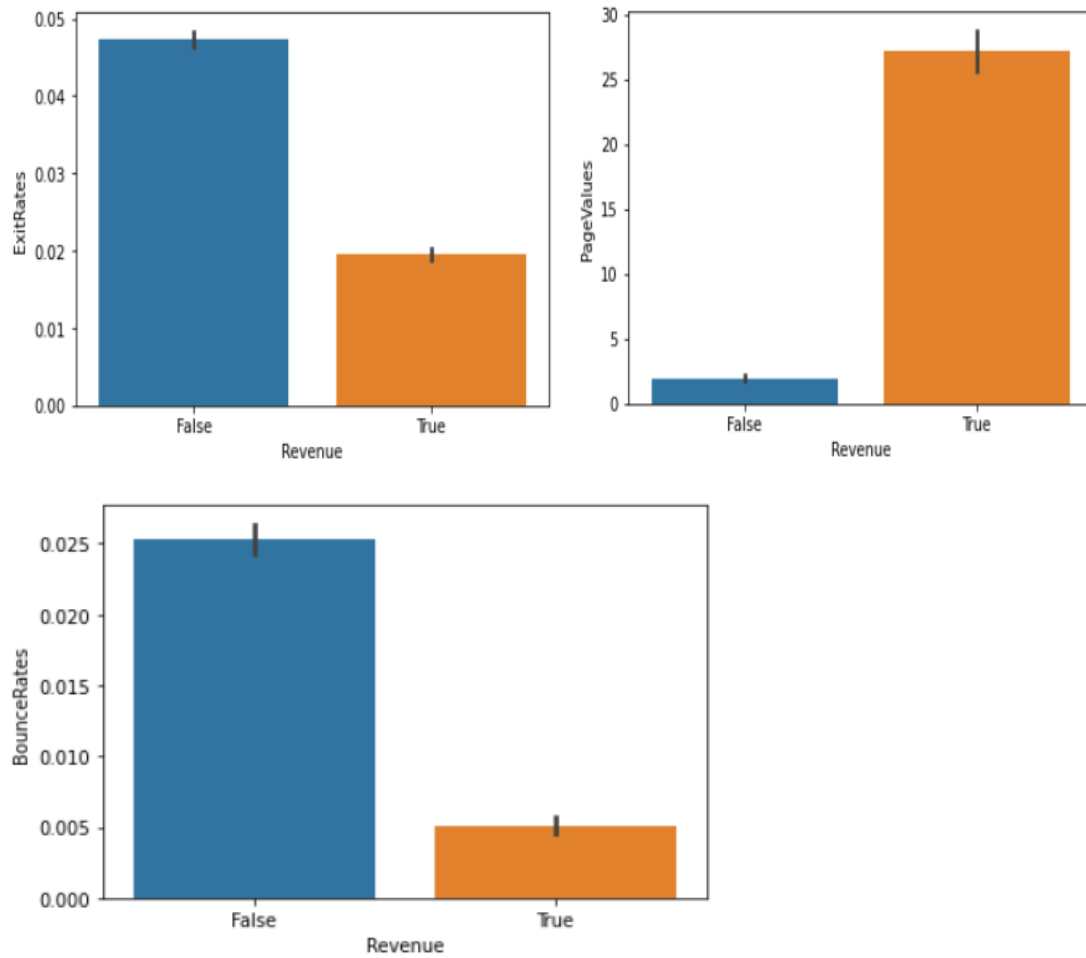# Data exploration

In the data provided, the first 10 columns are numerical values and the last 8 are categorical values. From these categorical features, the Month and Visitor Type are objects, Weekend and Revenue are Boolean, the other four features are all nominal variables which mean they represent each a respective category without having an order for the numbers.
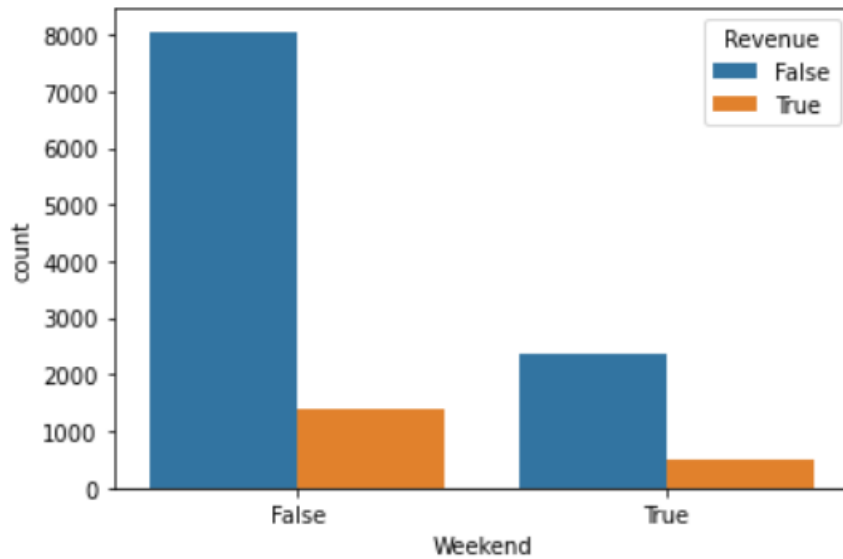
We need to analyze if there is any relation between the variables in the data using a correlation heatmap that will check the strength and direction of the correlation between the features.
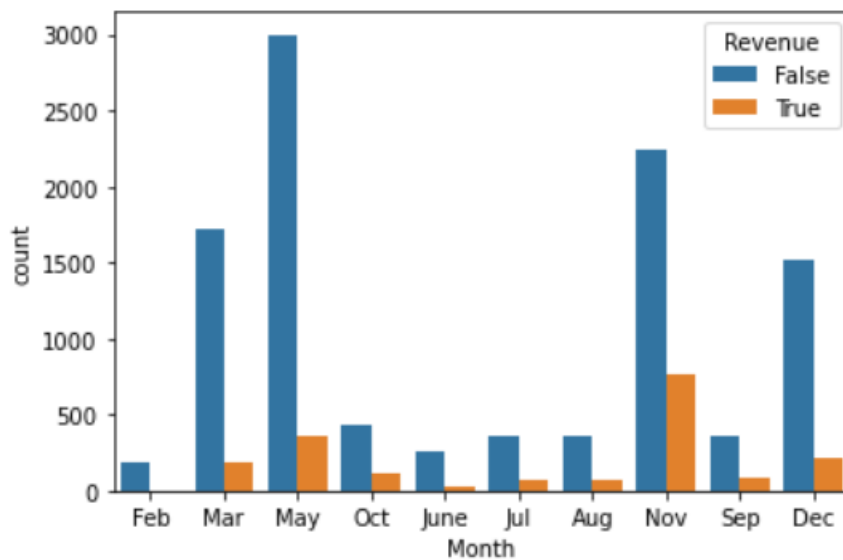


Checking the Revenue row as it is our goal to predict this outcome, we can notice the highest correlation of +0.49 with the Page Values, and a small negative relation with Exit Rates and Bounce Rates.
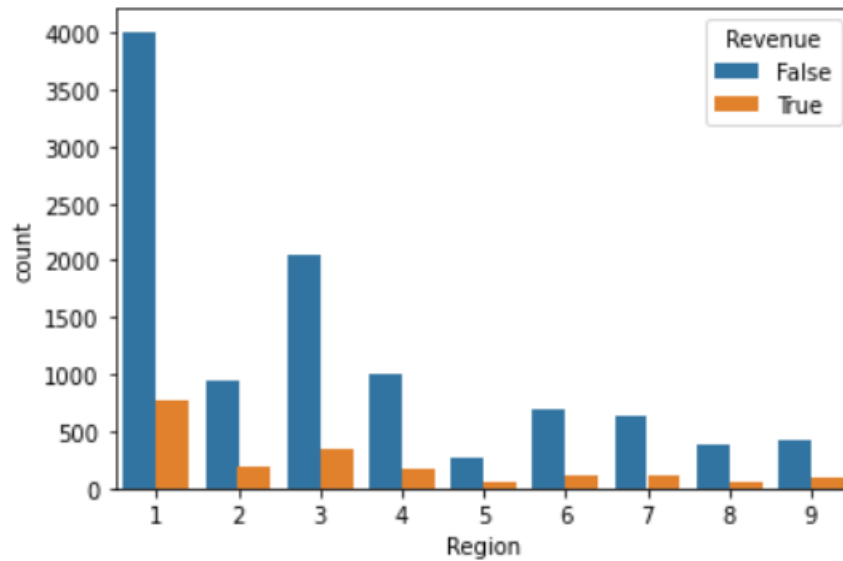
From the three figures, we might deduce that the higher the Page Values the more chance there is for a customer to buy a product, but the higher the Exit Rates and Bounce Rates the more chance the outcome will be to not buy.
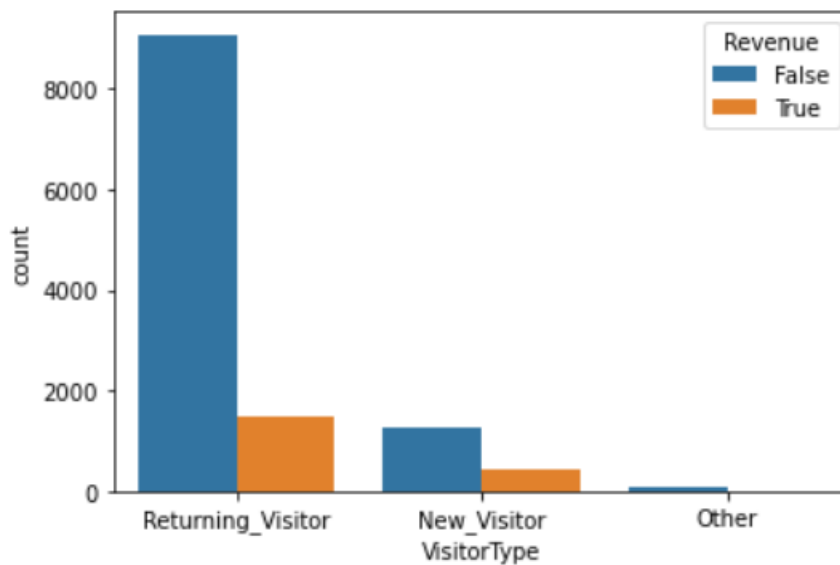
There is a lot of people visiting the sites during the week and the majority are not buying. During the weekend, the same thing happens but the number of total visits is much lower.



There is mainly one month that stands out between them all which is November with the highest number of purchases. May on the other hand has the highest number of visitors that does not buy.
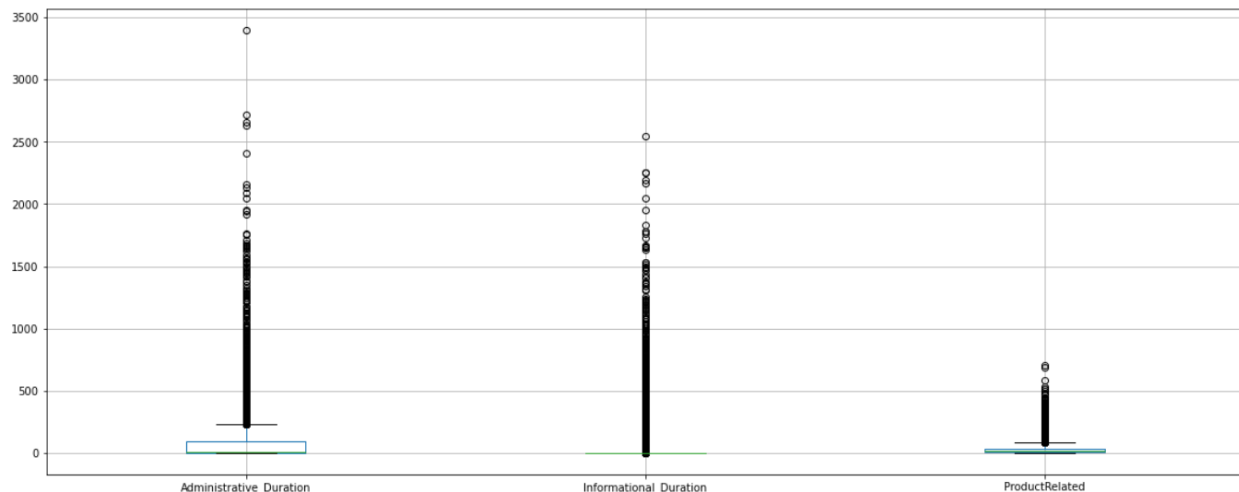
The highest number of values is in the region 1 that dominates the chart with region 3 , 4 and 2 following respectively.



Returning Visitors are the main type of visitors with a very high number of customers not buying with a very low percentage of them buying. New visitors consist of customers that have higher chance to buy than returning ones.

Next, we should check for outliers and decide what to do with them. The simplest way is to plot boxplots and see whether there are points that are situated outside of the whiskers of the box plot.

This is one example of the boxplots, they all seem to have outliers. For this project, the outliers were not removed as they did not change or affect the outcome of the models used.
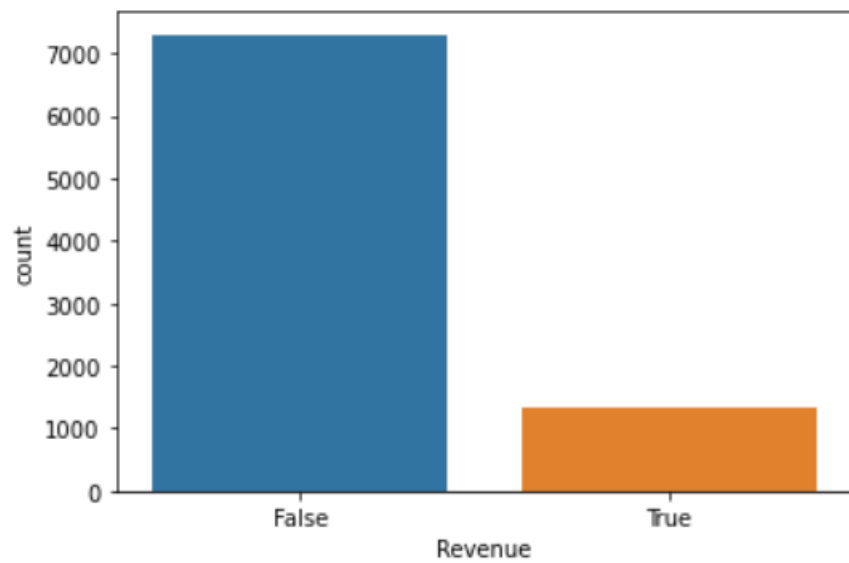
## Data preprocessing

There are no missing values in this data and all the variables are in the correct ranges such as values like the Durations do not have negative numbers with a minimum of 0.

We split the data between the Revenue and the other features with a 70% for training the data and 30% for the testing part. All the data preprocessing techniques will be applied only on the first part of the split data, that is the X train and X test in this case, so that we insure there will be no data leakage to the results when we will test our prediction, because then we may have wrong answers and very high accuracy and an f1 score that is too good to be true.
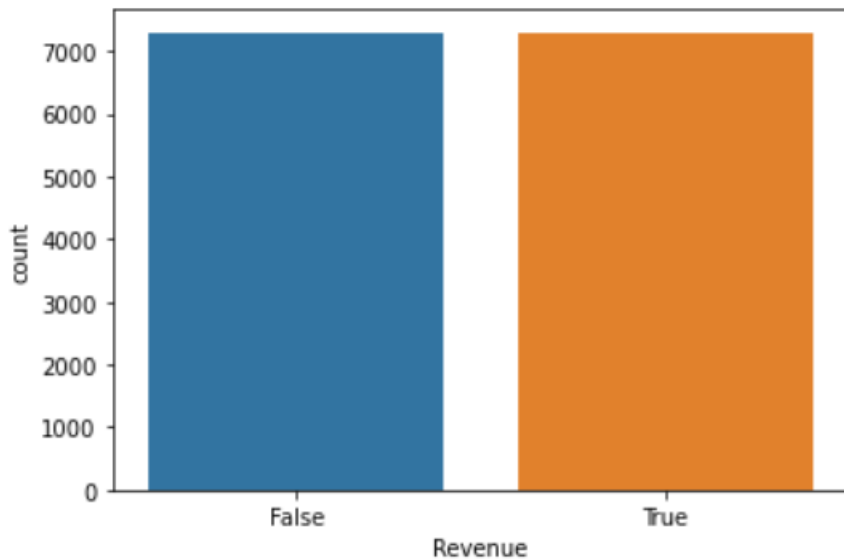
We are going to apply label encoding which will transform all categorical values to numerical values. For each month, a number will be assigned to represent it, from 1 to 9 as there is 9 months in this data set. Weekend will be 0 or 1 depending on if it is false or true and the Visitor Type will be from 0 to 2 as there is three type of visitors. The other categorical features do not need to be transformed as they are already in numerical form.

We will normalize the first 10 features, which are non-categorical. This will ensure that all the values are consistent with the same format and range, while keeping the information intact. Also this data does not follow a gaussian distribution.

We can clearly see that the two classes are not represented equally so this data is highly imbalanced, thus making prediction harder with this kind of data as we do not have a lot of information about the positive outcome we would like to predict.

We will be using an over sampling method, which will increase the number of the minority class of the Revenue to the number of the majority class so we can have a balanced data.

# Model Implementation

We used three different kind of decision trees implementations as they are robust to outliers and are used for supervised classification problems and do not require a large amount of computing power to perform.

the random forest is a random grouping of many other decision trees. Random forest combines the output of many decision trees as an average of the scores.
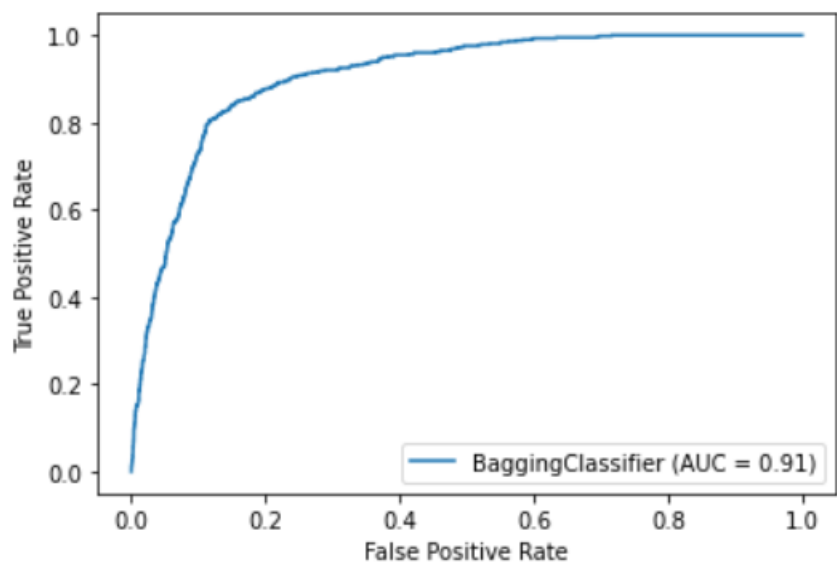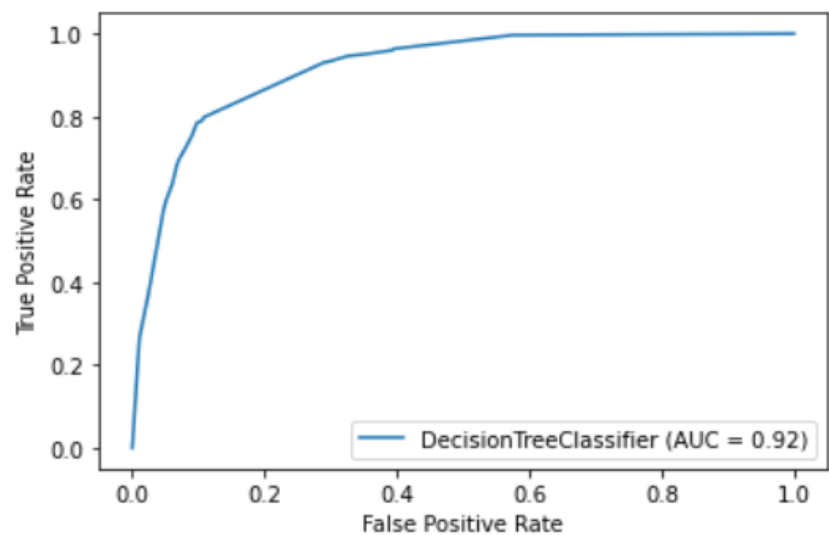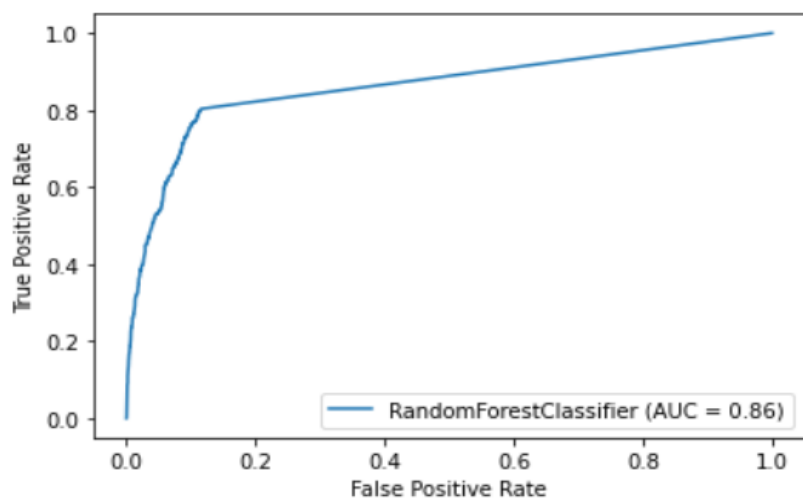
Bagging ensembles is another kind of classifier that is similar to the random forest, the difference is that the bagging considers all the features rather than a subset of feature for splitting a node. Usually, random forest is a better enhanced version of bagging.

After the implementation of each model, they were all configured with hyperparameter optimization to try and get a better performance.

Then we chose what features are the most important for each model (except for the bagging because it should use all features) and implement them with the new features selected.

We will use f1 score because this was an imbalanced target, it is a way to combine recall and precision under one score. The values are before and after we implemented the selected features and optimized parameters which improved the performance and results. At the end we also compared results with AUC score, which is unbiased towards a minority or majority class, it is a better score to use in this case rather than the accuracy.

|  | Bagging ensemble | Decision Tree | Random Forest |
|---|---|---|---|
| F1 score before | 0.65 | 0.65 | 0.65 |
| F1 score after | 0.66 | 0.67 | 0.66 |
| Time (seconds) before | 0.4 | 0.06 | 11.1 |
| Time (seconds) after | 5.07 | 0.01 | 5.7 |

RandomForestClassifier (AUC = 0.86)



DecisionTreeClassifier (AUC = 0.92)



BaggingClassifier (AUC = 0.91)

# Conclusion

After seeing all the results and the performances of each model, we can conclude that the decision tree is the most effective between them as it is the fastest (0.01 seconds) and with better results regarding f1 (0.67) and AUC (0.92) score.

  If we are going to use Decision Tree as the main predictive model, by looking at the selected features we can assume that the customers that wants to buy are affected by the high number of page values and the month they are visiting the site, especially the month of November as it had the highest number of visits and purchases as seen before in the figure.