Master's degree in Artificial Intelligence and Data Engineering

Data Mining and Machine Learning project documentation

# CHURNPREDICT

Denny Meini

**INDICE**

# 1. INTRODUCTION

ChurnPredict is an application which permits the user, an Ecommerce holder, to know if a customer is likeable or not to leave the Ecommerce.
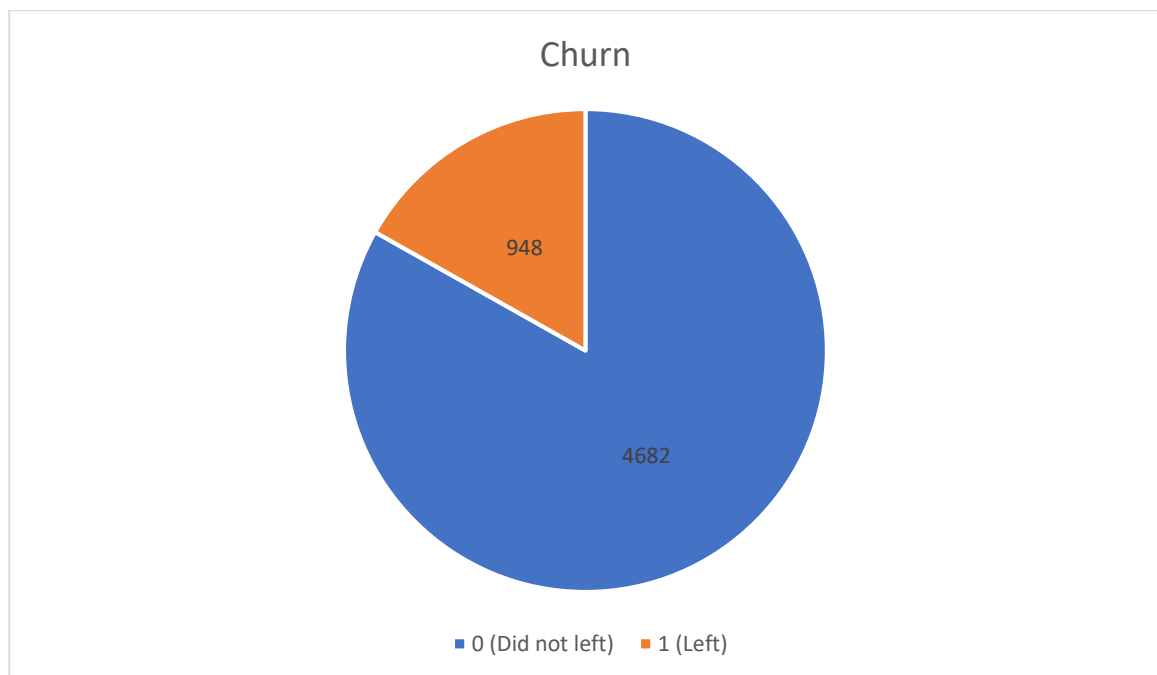
The main idea behind this project is to compare classifiers of different types, in order to evaluate their performance.

Before writing the application, I made some studies on the dataset, with the goal of preparing it to the classification and to use it with some classifiers. I compared the results I obtained, then, for my application, I used the classifier which gave me the best performances in the study phase.

The final application will give the user the possibility to insert customer's data in the specified field and will return on the screen whether there is a risk or not that the customer will leave the Ecommerce.

## 2. DATASET

The dataset I chose for this project comes from Kaggle, and it is named "Ecommerce Customer Churn Analysis and Prediction". It has 5630 rows, each of them represents a customer. Each row has 20 attributes, 19 of them are data relative to the customer, like tenure, preferred payment method, order count. The last feature is a binary attribute which tells us whether or not the customer has left the Ecommerce.

### Churn

948

4682

■ 0 (Did not left)    ■ 1 (Left)

As we can see from the graphic, the dataset is not balanced, with a division that is 83.2%/16.8%. In fact I had to apply a rebalancing method before the classification phase.

The 20 features of the dataset are:

1. **CustomerID**
   The ID of the customer.

2. **Churn**
   Output class [0: did not left, 1 left].

3. **Tenure**
   The number of years the customer has been a customer.

**4. PreferredLoginDevice**

The device the customer has used the most to access the Ecommerce.

**5. CityTier**

The tier of the city (Chinese model) [1: big, 2: medium, 3: small].

**6. WarehouseToHome**

Distance between the warehouse and the customer's house.

**7. PreferredPaymentMode**

Customer's preferred method of payment.

**8. Gender**

The gender of the customer [Female/Male].

**9. HourSpendOnApp**

Number of hours the customer has spent on the app.

**10.NumberOfDeviceRegistered**

Number of devices the customer registered.

**11.PreferredOrderCat**

The preferred item category of the customer based on his/her orders

**12.SatisfactionScore**

Grade of satisfaction of the customer from 1 (not satisfied) to 5 (very satisfied).

**13.MaritalStatus**

Marital status of the customer [Married/Single/Divorced].

**14.NumberOfAddress**

Number of address added by a customer.

**15.Complain**

Whether the customer complained or not during the last month [0: no, 1: yes].

**16.OrderAmountHikeFromlastYear**

Percentage increase regarding the orders from the previous year.

### 17. CouponUsed
Number of coupon the customer used in last month.

### 18. OrderCount
Number of order the customer placed during last month.

### 19. DaySinceLastOrder
Number of days since the customer's last order.

### 20. CashbackAmount
Average cashback of the customer during the last month.

What is interesting to see is that the dataset presents some null values.

# 3. PREPROCESSING

The preprocessing phase is composed by several operations, which have the goal of cleaning and adjusting the dataset in order to make it ready to the classification phase.

### I. First analysis and division of the dataset

I started with an analysis of the values of the column. I found that there were no evident outliers in the dataset (eg. OrderCount<1), but there were some values in the categorical featured that was written in different ways even if they represent the same value:

- 'Mobile phone' and 'Phone' in PreferredLoginDevice
- 'CC' and 'Credit Card' in PreferredPaymentMode
- 'COD' and 'Cash on Delivery' in PreferredPaymentMode
- 'Mobile' and 'Mobile Phone' in PreferedOrderCat

As a consequence I had to fix this problem changing the values "Phone" and "Mobile" in "Mobile Phone", "CC" in "Credit Card" and "COD" in "Cash on Delivery".
Then I faced the problem of the null values and I saw that there were two possible ways:

- Dropping the records that contained a null value
- Filling the null value

I decided to use both ways and to compare them in the classification phase. By dropping the record with null values the dataset was reduced to 3774 rows, distributed in the following way:
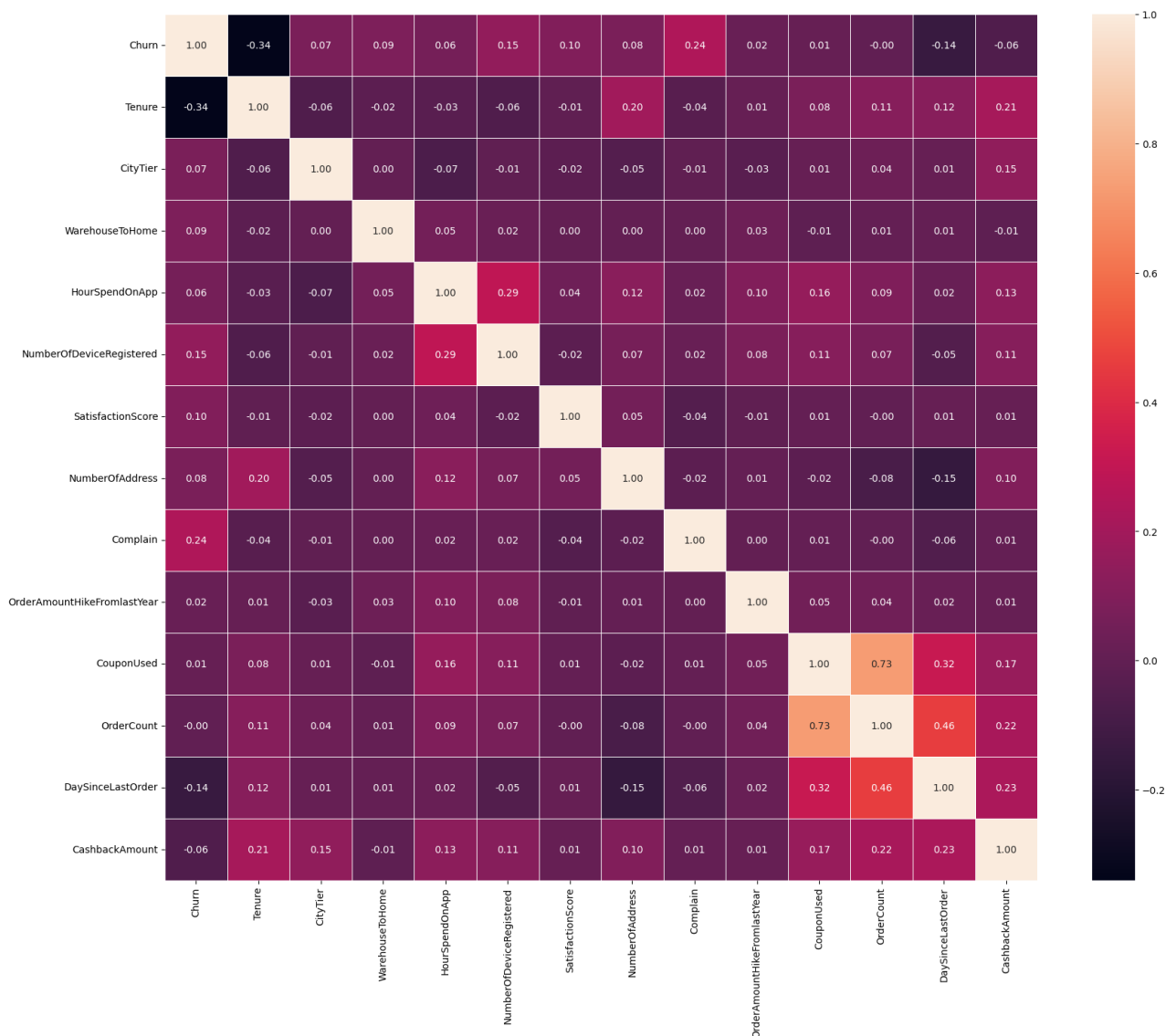


Churn

631

3143

■ 0 (Did not left)    ■ 1 (Left)

It's clear that the proportion between the two classess remains the same. For what concerns the dataset with the null values filled, I decided to replace every null value with the mean of the feature (all the null values was about a numerical feature), however I considered only the records belonging to the same class. The proportion between the two classes obviously remains the same.

Finally I dropped from both datasets the column "CustomerID" because the ID of the customer is unique and does not have effects on the classification.

## II.  Correlation analysis and feature extraction

I performed a correlation analysis in order to see if there were some features that were strongly correlated (at least 0.9).

I obtained this two correlation tables:

Dropped values:

## Filled values



In both cases the correlation between the output feature (Churn) and the others is really weak (at maximum -0.36); in general the strongest correlation in both datasets is between OrderCount and CouponUsed (0.73 in dropped, 0.66 in filled). Since they were not greater than 0.9 I decided not to drop one of them and I maintained them both.

Then I used the get_dummies in order to transform the categorical values in a set of numerical features (0/1). At that point I had only numerical features, so I gave another look at the correlation table. I did not obtain any relevant difference with respect to the two values I highlighted in the previous case, as a matter of fact I could simply extract the two datasets and end the preprocessing phase (I did not report the last two correlation tables because

they would not be clearly readable in this document. However, they are still present in the "Preprocessing" Notebook).

# 4. CLASSIFICATION

For the classification phase the objective was to find the classifier which will correctly predict the output of an instance. In order to do this, I dropped the Churn column from the dataset and I used it as target column.
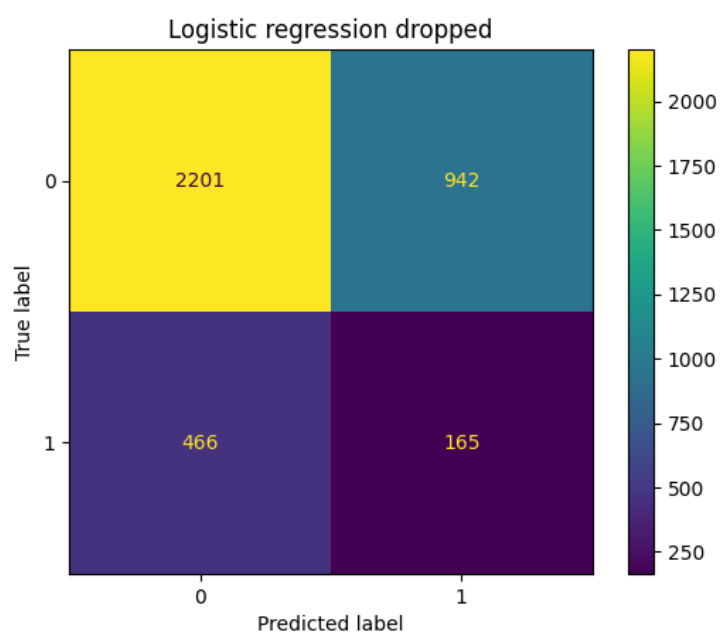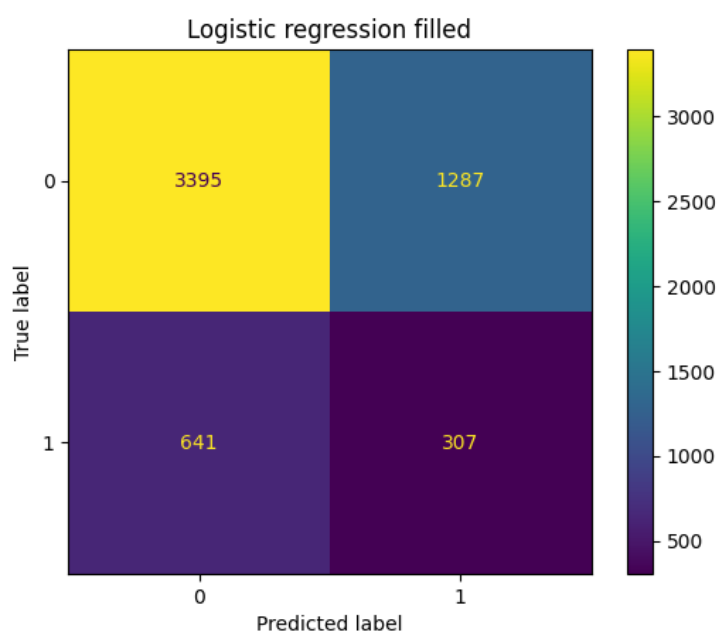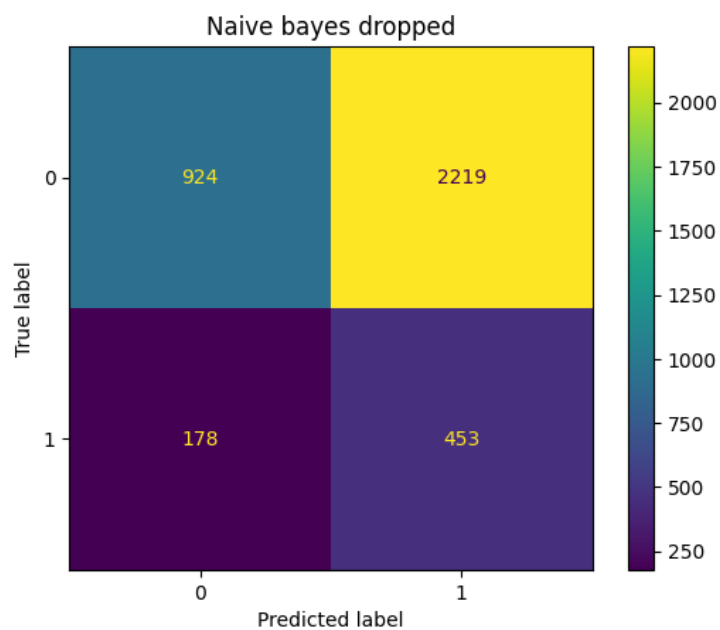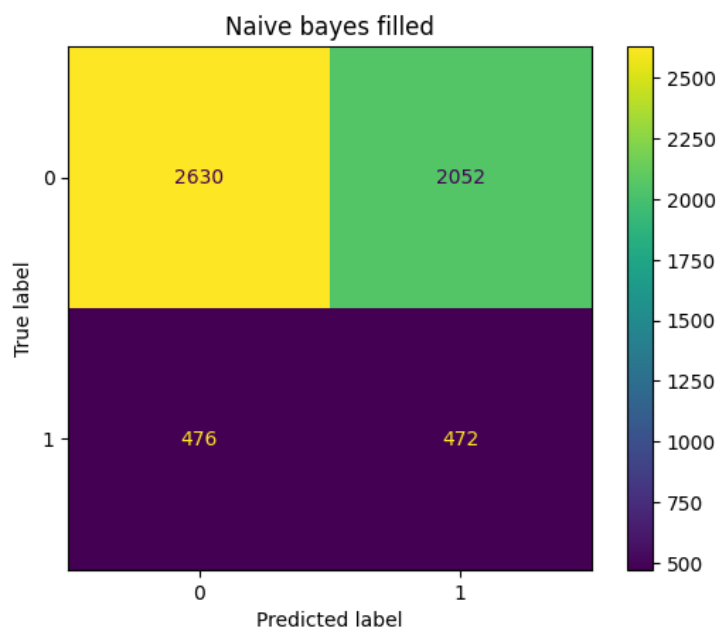
I decided to evaluate the accuracy of the classifiers using a k-fold cross validation (with k=5) in order to obtain results that were the least dependent by the choice of test and training set; specifically I used Stratified K fold. The parameters that I used to evaluate the models were f-score and average accuracy.
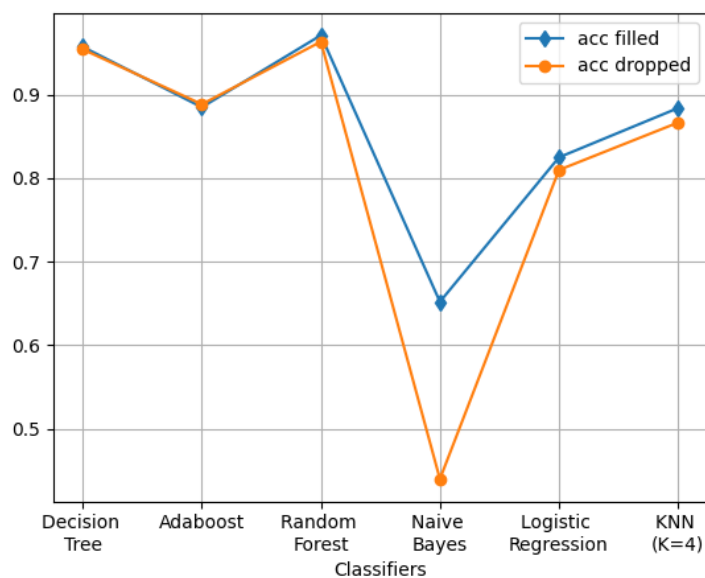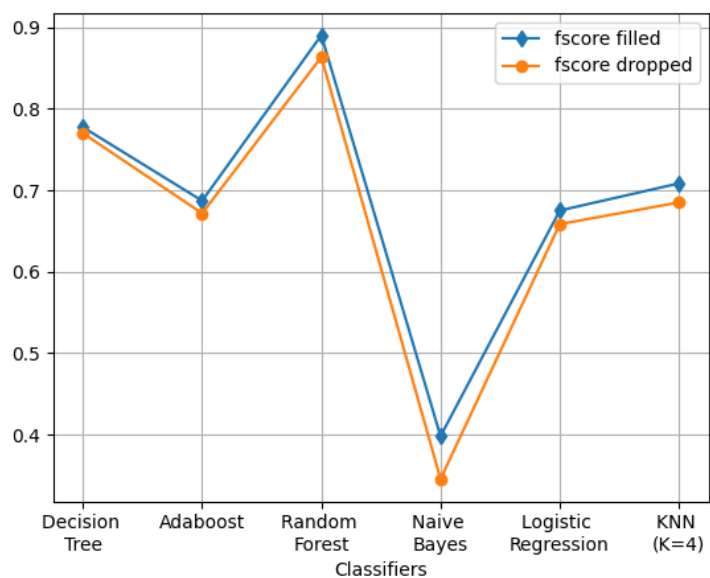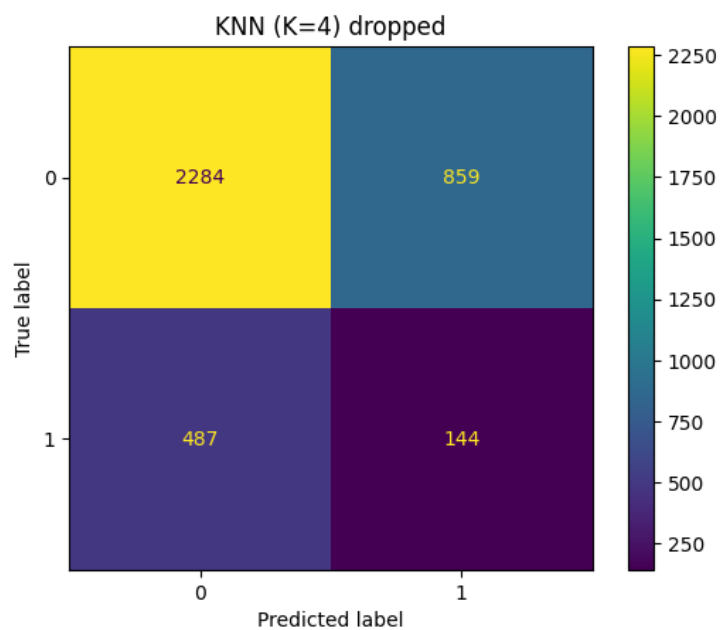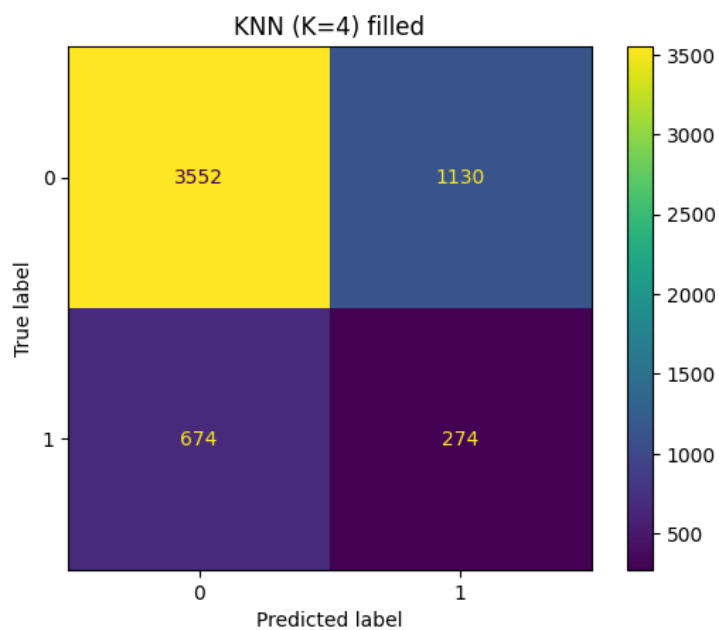
Due the fact that the two classes are imbalanced, I used SMOTE before every classification in order to rebalance the dataset.

I used each classifier two times: first with the filled dataset and then with the dropped one. For K Nearest Neighbors I tried with K berween 3 and 10 and I obtained the best results for K=4 (I do not report all the tries for simplicity). Here I show the confusion matrixes of the classification and then I will report the results in terms of f-score and average accuracy for each try.

Naive bayes filled

| | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 2630 | 2052 |
| True 1 | 476 | 472 |

Naive bayes dropped

| | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 924 | 2219 |
| True 1 | 178 | 453 |

Logistic regression filled

| | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 3395 | 1287 |
| True 1 | 641 | 307 |

Logistic regression dropped

| | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 2201 | 942 |
| True 1 | 466 | 165 |

| Classifiers | | f-score | Avg accuracy |
|---|---|---|---|
| Decision Tree | Filled | 0.777 | 0.957 |
| | Dropped | 0.770 | 0.954 |
| AdaBoost | Filled | 0.687 | 0.885 |
| | Dropped | 0.672 | 0.888 |
| Random Forest | Filled | 0.890 | 0.971 |
| | Dropped | 0.864 | 0.963 |
| Gaussian Naive Bayes | Filled | 0.398 | 0.652 |
| | Dropped | 0.345 | 0.440 |
| Logistic Regression | Filled | 0.675 | 0.825 |
| | Dropped | 0.658 | 0.810 |
| KNN (K=4) | Filled | 0.708 | 0.884 |
| | Dropped | 0.685 | 0.866 |

From these results I could draw two conclusions:

- The dataset with the null values filled always gave best results (the only exception is avg accuracy with AdaBoost, but the difference is just of 0.003)
- Random Forest is globally the best classifier with respect to both f-score and avg accuracy

The last analysis I performed was the Wilcoxon test because I wanted to know if the difference between the filled and the dropped results was statistically relevant. The test for Random Forest gave a pvalue of 0.0625. Since this value is greater than 0.05 I could say that the difference had no statistical relevance, therefore I cannot conclude that to fill the null value is globally a better choice than drop the relative records. Due to the fact that filling the value is not so complicated, compared to drop the records, and since the difference of f-score for Random Forest is not so small (0.026), I could conclude that Random Forest with the null values filled was the best option in order to obtain a better classification.

# 5. IMPLEMENTATION OF THE APPLICATION

The application I created, in order to use the dataset and Random Forest, is called "ChurnPredict". The goal of the application is to give the user the possibility to know if a specific customer is likable to leave the Ecommerce, thanks to a list of data asked by the application itself. The application is written in Python and makes use of the following libraries:

- Tkinter, for the GUI.
- Sklearn, to perform Random Forest.
- Pandas, to read the dataset.
- Imblearn, for using the pipeline and SMOTE
- Numpy, in order to transform the list of inputs in an array.

The homepage looks like this:

The user has to insert the values into the fields; some of them give the user a list of possible choices (like Preferred order category), others are just text fields.

Finally, the user will click on "Confirm" and the app will show a message on the screen telling whether there is a risk or not.

## 6. CONCLUSIONS

In this study I compared a list of classifiers using a dataset with a binary classification. The tests were done using a Stratified K Fold as cross validator and SMOTE, in order to deal with the fact that the dataset was imbalaced. Due to the fact that there were some null values, I created a dataset without the records that had a null value and a dataset with the null values filled with the mean of the feature (only considering the same class). All the classifiers, with the exception of Gaussian Naive Bayes, had a better performance on the filled dataset with respect to f-score and average accuracy. Particularly good was Random Forest, considering that it gave me a very high performance, even if Wilcoxon test told us that there is no statistical relevance in the performance difference. In conclusion, the last step was building a simple application in Python which uses the filled dataset and Random Forest.