

# Credit Card churn prediction

Quoc Frejter 430201, Katarzyna Dybcio 430189



# Problem description

- **Background**

The business manager, that is taking care of consumer credit card portfolio, is concerned about the increasing number of customers leaving their credit card services.

- **Problem**

With only 16.07% of customers exhibiting churn, predicting churn is challenging due to class imbalance.

- **Aim of the project**

Develop a predictive model to identify potential churners, empowering the bank to proactively engage and retain customers through personalized services.



# Dataset

- **Target Variable: Attrition\_customer**

Internal event (customer activity) variable - if the account is closed then 1 else 0

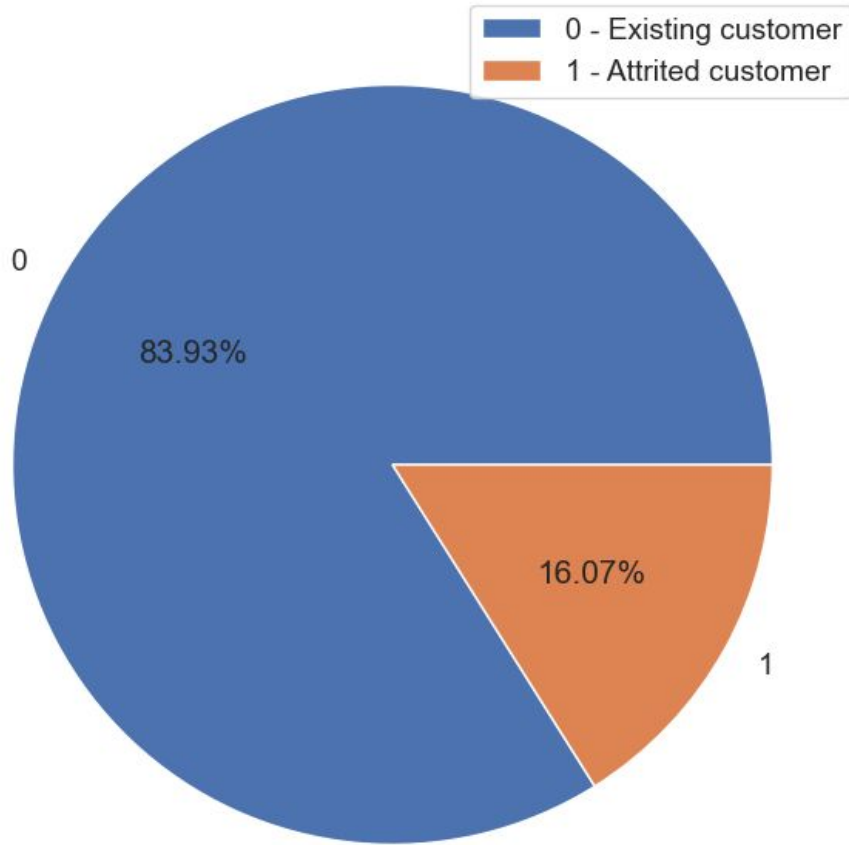
- **Independent Variables: 20 features**

'CLIENTNUM', 'Attrition\_Flag', 'Customer\_Age', 'Gender', 'Dependent\_count', 'Education\_Level', 'Marital\_Status', 'Income\_Category', 'Card\_Category', 'Months\_on\_book', 'Total\_Relationship\_Count', 'Months\_Inactive\_12\_mon', 'Contacts\_Count\_12\_mon', 'Credit\_Limit', 'Total\_Revolving\_Bal', 'Avg\_Open\_To\_Buy', 'Total\_Amt\_Chng\_Q4\_Q1', 'Total\_Trans\_Amt', 'Total\_Trans\_Ct', 'Total\_Ct\_Chng\_Q4\_Q1', 'Avg\_Utilization\_Ratio'

- **10127 observations**



Existing and attrited customers (%)

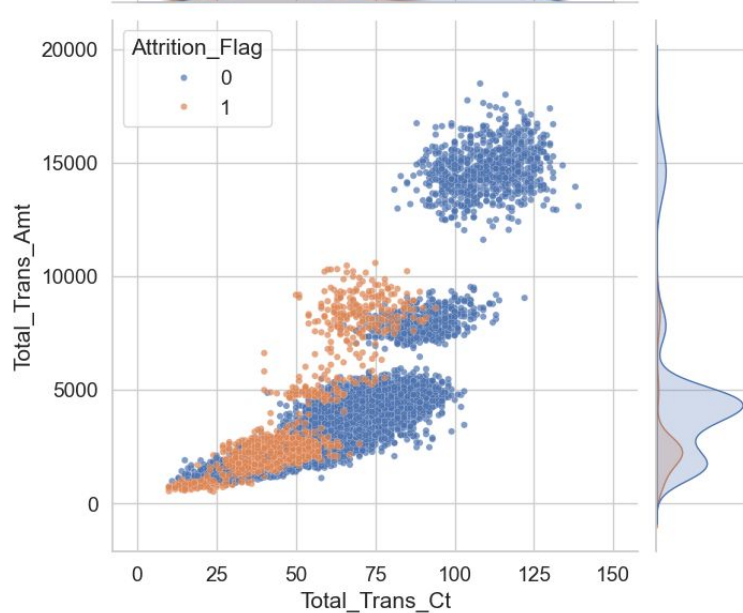


# Target Variable

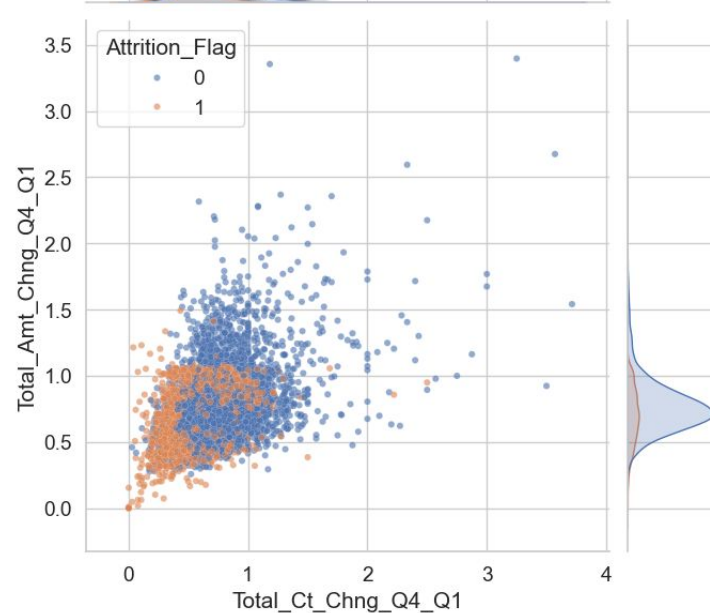
- Highly imbalanced dataset

Attrition_Flag	Count	Percent of all
0	8500	84%
1	1627	16%

Number of total transactions and total transaction amount based on customers attrition

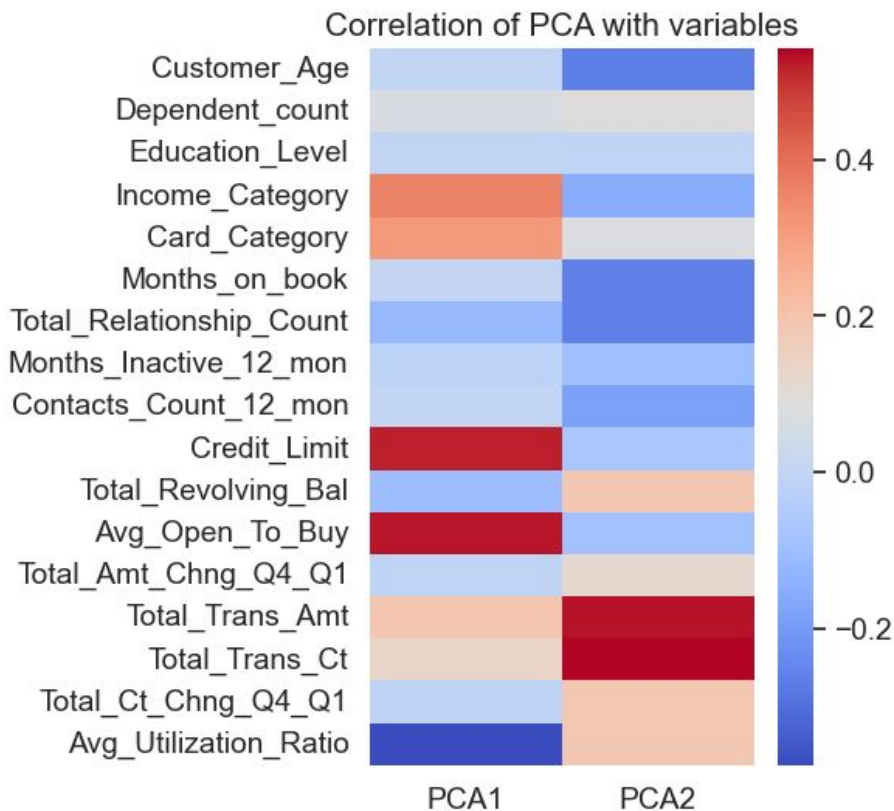


Total number and amount change in transaction Q4 over Q1 based on customer attrition



Customers who have opted out of the bank's services have significantly fewer completed transactions, and these transactions involve a lower average amount of funds.

# K-means, PCA - clustering of customers



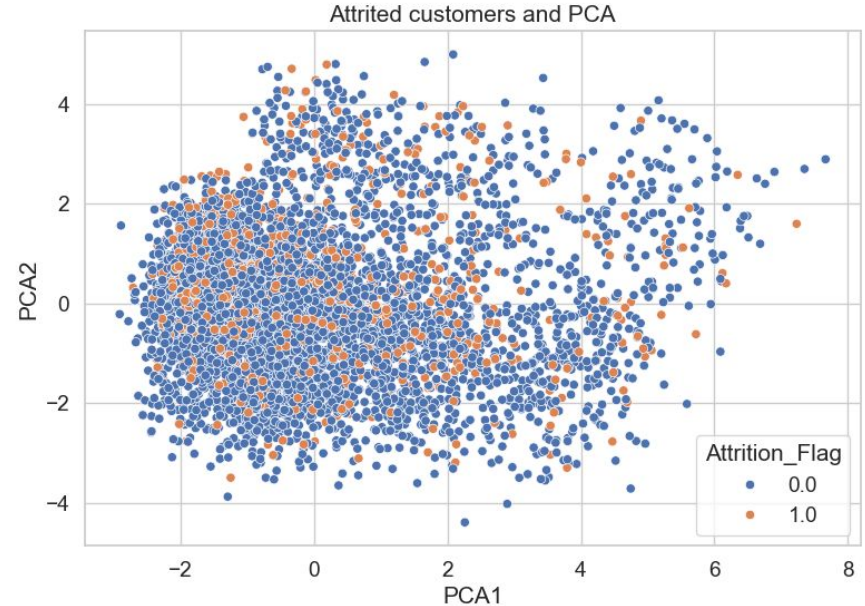
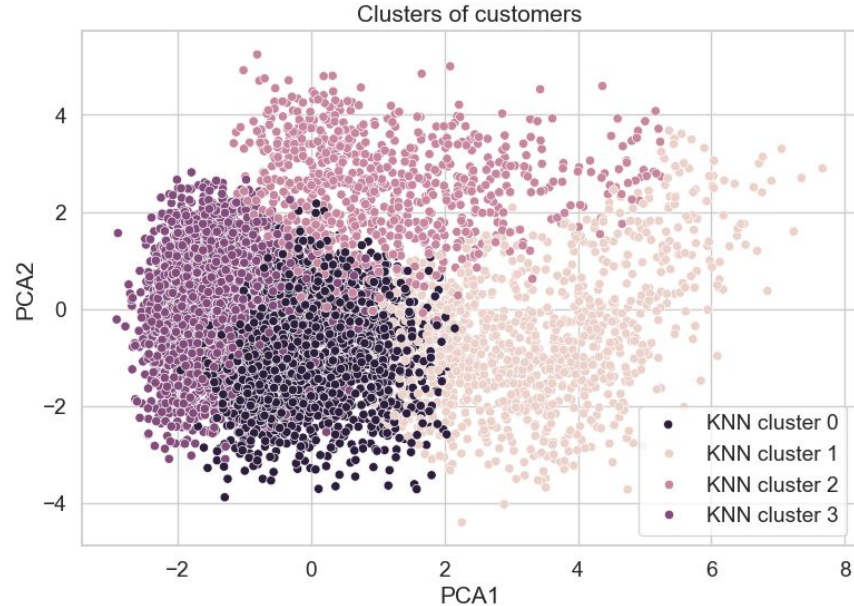
## High PCA1

- High income, "premium" card, high credit limit but low utilization
- Wealthy customers who exhibit a low utilization ratio.

## High PCA2

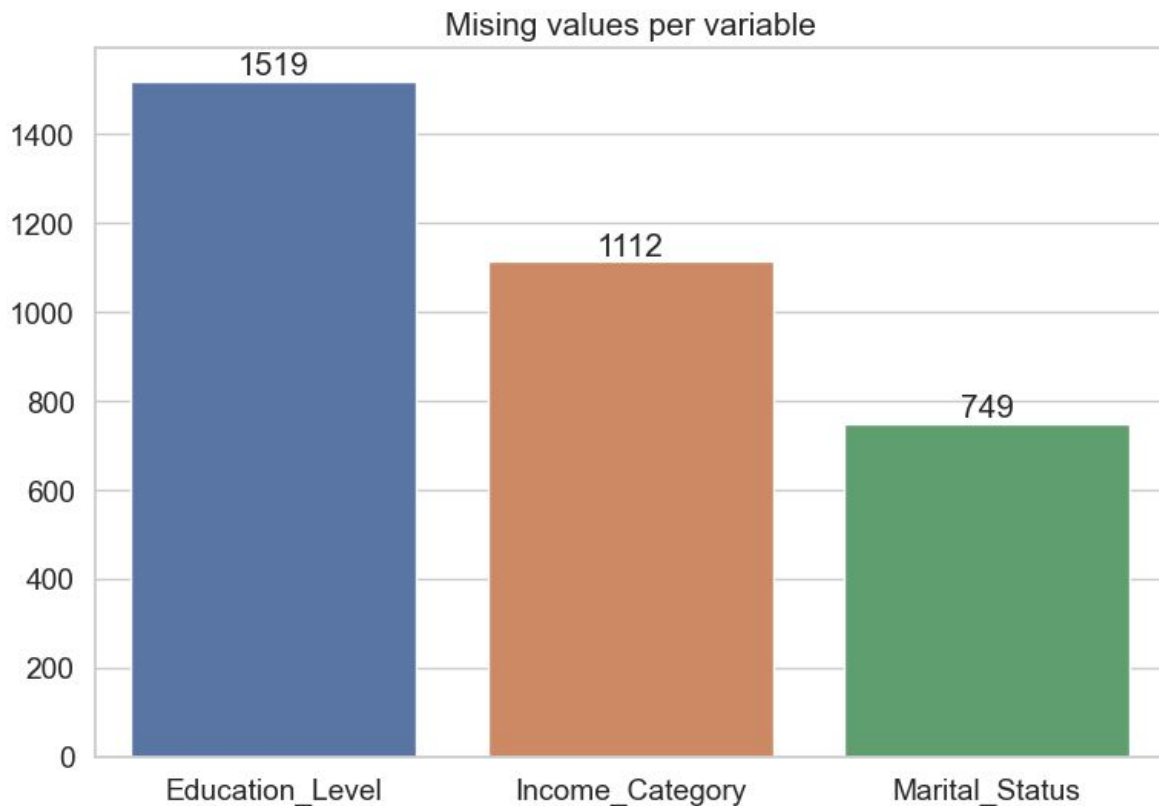
- Young, low income, low number of month on book, low number of products, high transaction count and amount
- Young/new customers who actively use their cards

# Plot of clusters



- On this plot we can see, that employing K-Nearest Neighbors (KNN) and Principal Component Analysis (PCA) for modeling the attrition flag could pose challenges. We can see that customers are quite mixed up and our algorithm would have problems dividing the group.

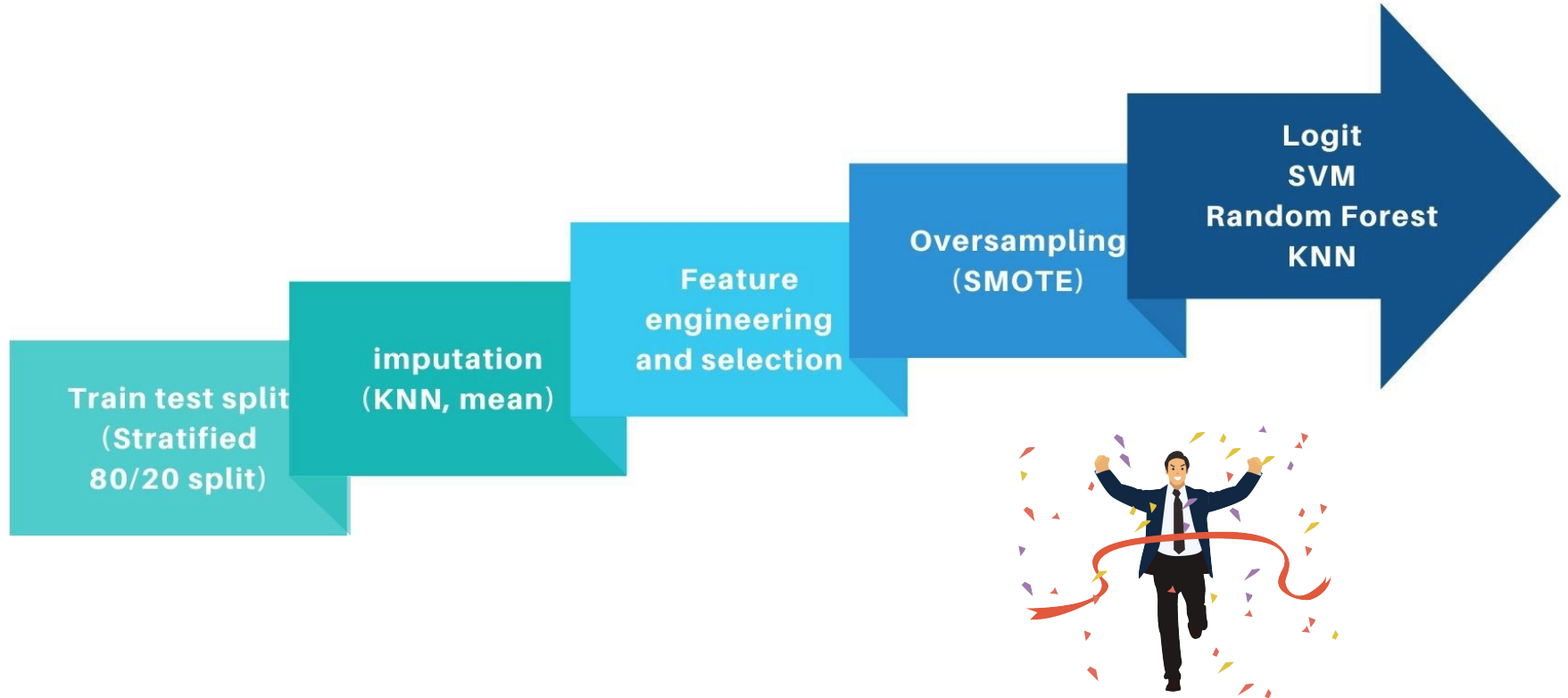
# Missing values



- We have missing values only for Education level, Marital status and income category
- Imputation methods
  - Education\_level - imputed using KNN
  - Income\_category - imputed using KNN
  - Marital\_Status - Filled missing values with most frequent



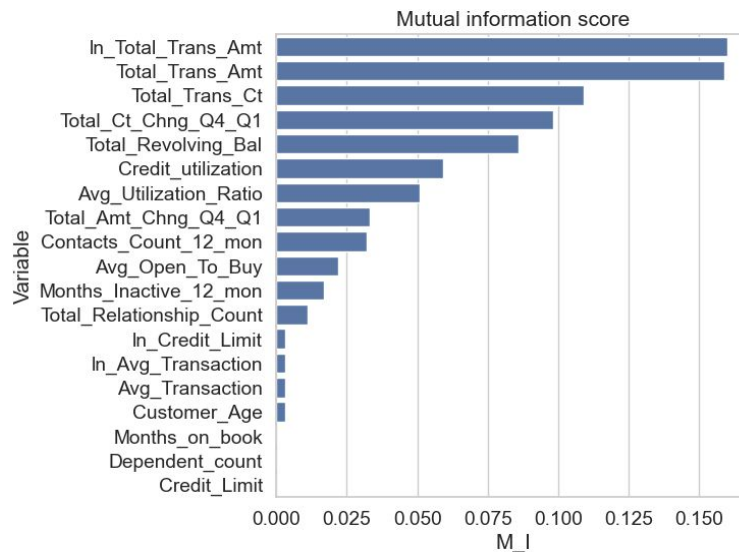
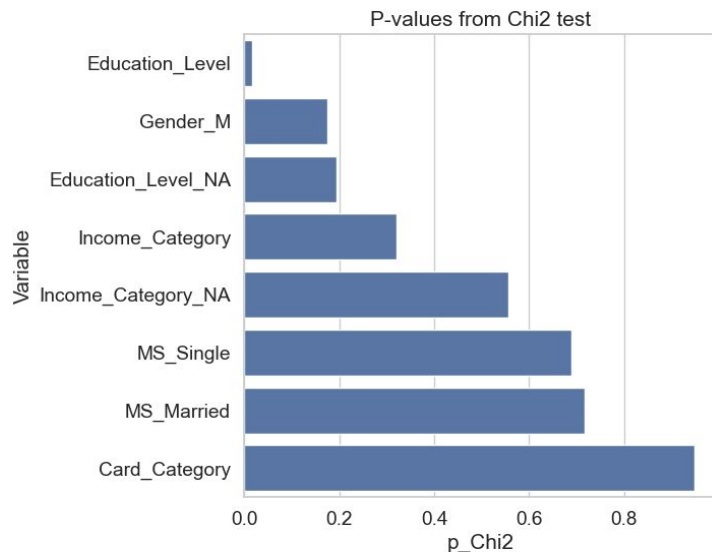
# Steps for success



# Feature engineering and selection

## Transformations:

- Categorical variables: one-hot encoding, ordinal encoder
- Numerical variables: log transformation, standard scaler
- Interactions: average transaction, account utilization



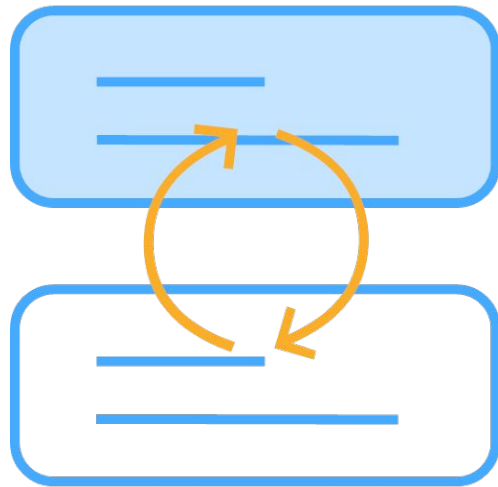
# Hyperparameter tuning

Method: **Grid search**

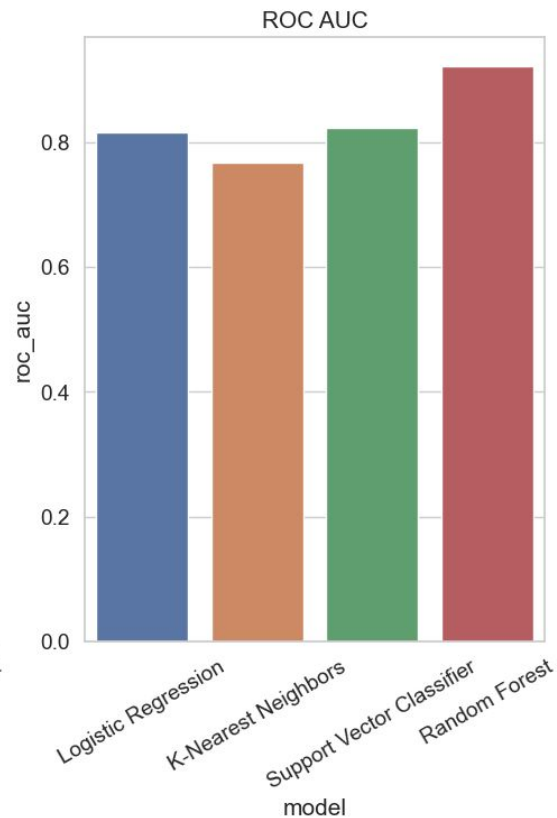
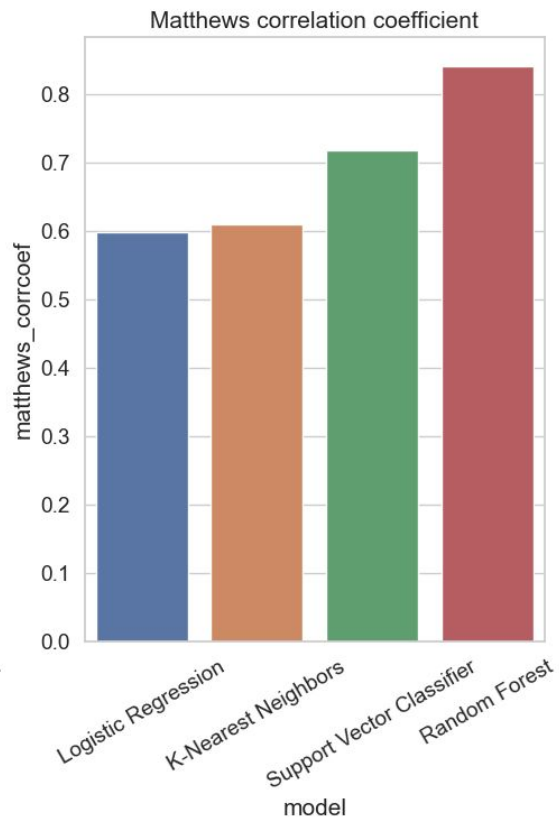
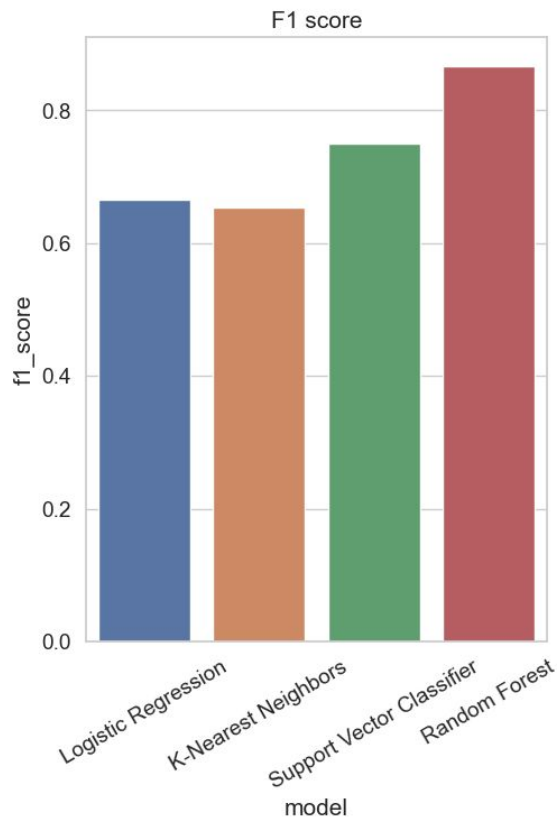
Cross validation: **Stratified K-fold** (8 splits)

Scoring measures:

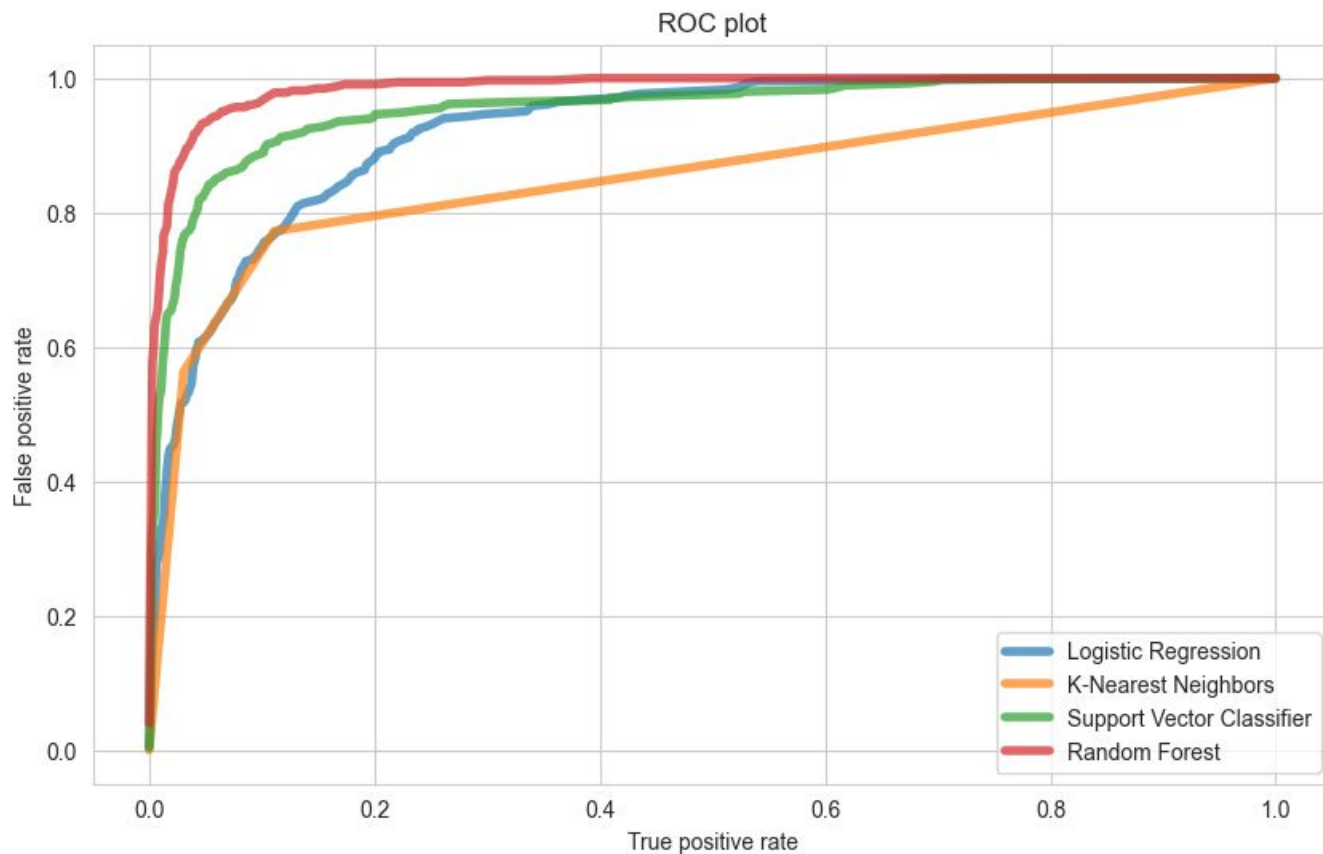
- **F1 score**
- **Matthews correlation coefficient**
- ROC AUC (additional)



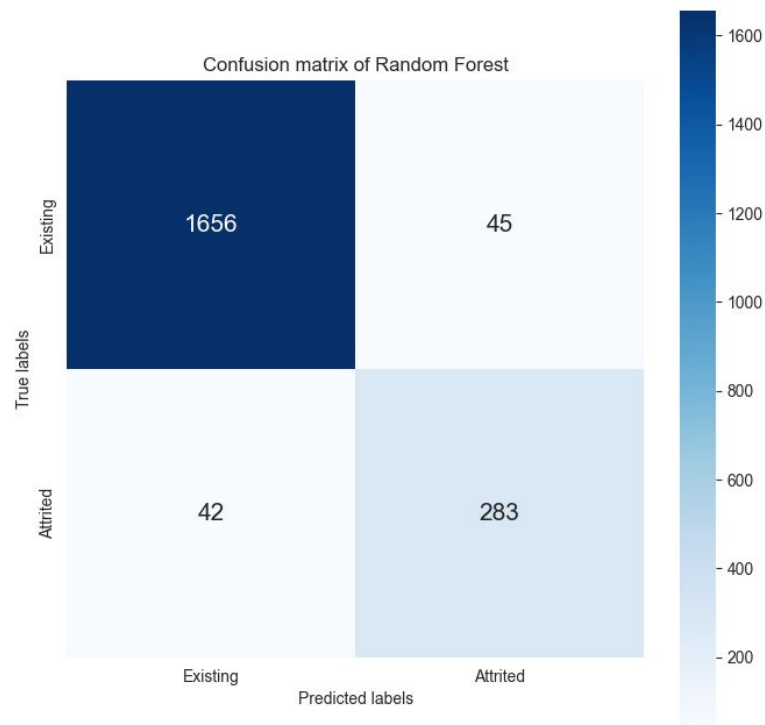
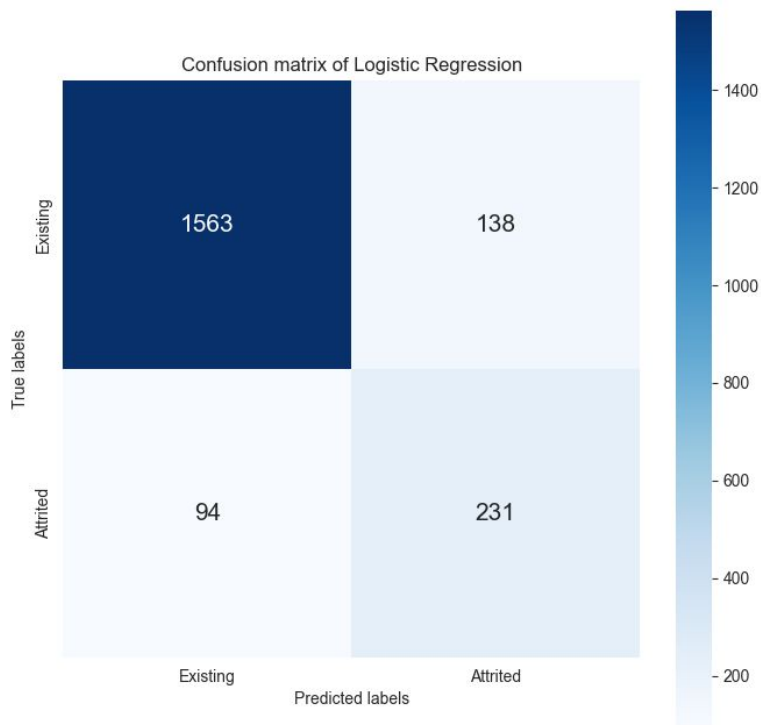
# Models performance



# Models performance



# Logistic Regression vs Random Forest



# Summary

- Additionally, Random Forest correctly identified the highest number of customers(238) who actually churned from the bank services
- In almost every model, the best results were achieved with data augmented using the SMOTE method. That is because our data was highly imbalanced.
- For the Random Forest, it was beneficial to reduce the number of variables when training the model what was confirmed by the test statistics.
- Based on this, we can conclude that the Random Forest model is the most effective method for our database.

# Thank you!

