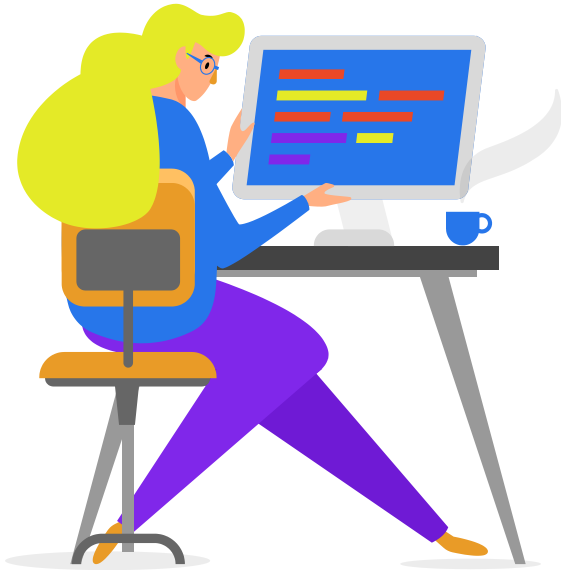




# ***Bank Account Fraud***

António Pedro Pinheiro  
up201704931

# Table of Contents



01

## Problem Definition

Business Understanding

02

## Data Understanding

Concise summary, with focus on main findings

03

## Data Preparation

Outline, with focus on the main operations

04

## Predictive Modelling

Experimental Results

05

## Conclusion

Limitations and future work

06

## Annexes

# Business Understanding: Problem Definition

## Goal

### Problem Statement

Detect fraudulent applications to bank accounts by training a machine learning model that accurately predicts if we're dealing with credit fraud

## Output

### Classification Problem

Binary: Fraud vs. No Fraud

## How?

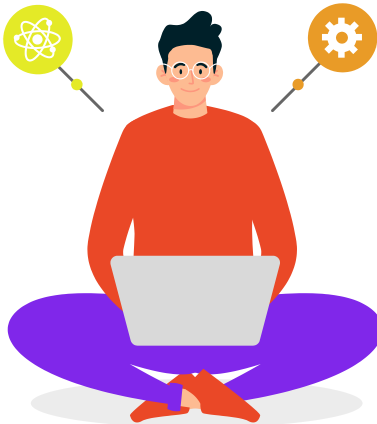
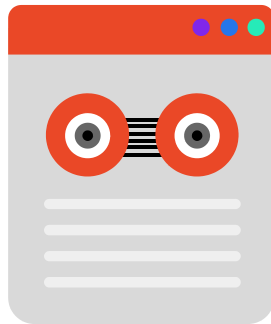
### ML Models

Predicting the outcome given a set of predictors

## Success

### Metrics

Area under Curve



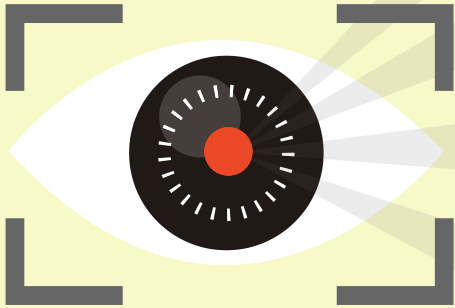
# Data Understanding: Summarization

## Dataset

36 predictors  
1 target variable



## Analysis & Summary



### Significant portion of missing values

prev\_address\_months\_count  
bank\_months\_count

### Device Fraud Count

All cases with value 0

### Customer Age

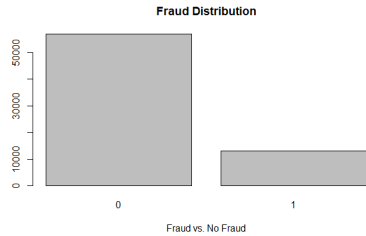
Cases <10 & > 90

### Imbalanced Domain Learning

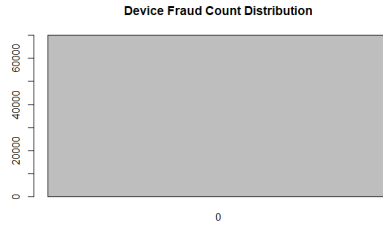
81% of cases on the train dataset have "no fraud" status

# Data Understanding: Visualization

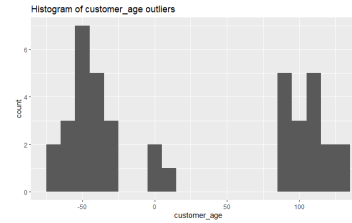
## Imbalanced Domain Learning



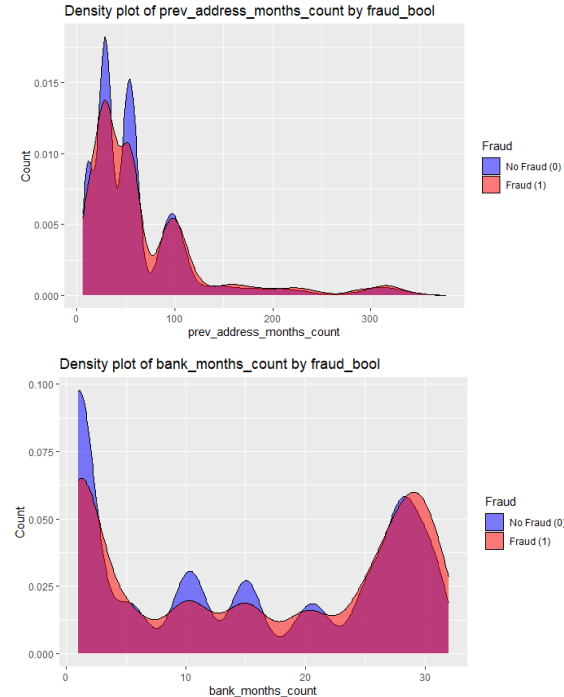
## Device Fraud Count



## Customer Age



## Significant portion of missing values



# Data Preparation: Data Quality & Transformation

## 01 Sampling Data

Stratified  
Train/Validation Sets  
(70/30)

## 02 Training Set: Data Cleaning

Elimination of cases  
with customer\_age  
<10 & >90

## 03 Missing Values

Imputation of Mean  
and Mode in numeric  
and categorical  
variables, respectively

## 04 Data Normalization

Z-Score  
normalization

## 05 Training Set: Removal of Redundant & Unnecessary Attributes

- ID
- prev\_address\_months\_count
- bank\_month\_count
- device\_fraud\_count

# Data Preparation: Feature Selection

## Correlation Matrix

Numeric Variables

## Chi-Square Test

Categorical Variables

## ANOVA One-Way Test

Numeric/Categorical Variables

## Information Gain

Reduction in entropy

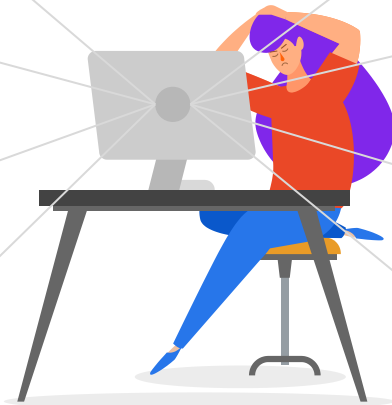
01

02

03

04

## Methods



05

06

07

08

**Variable Importance** with  
Random Forests  
Mean Decreasing Gini

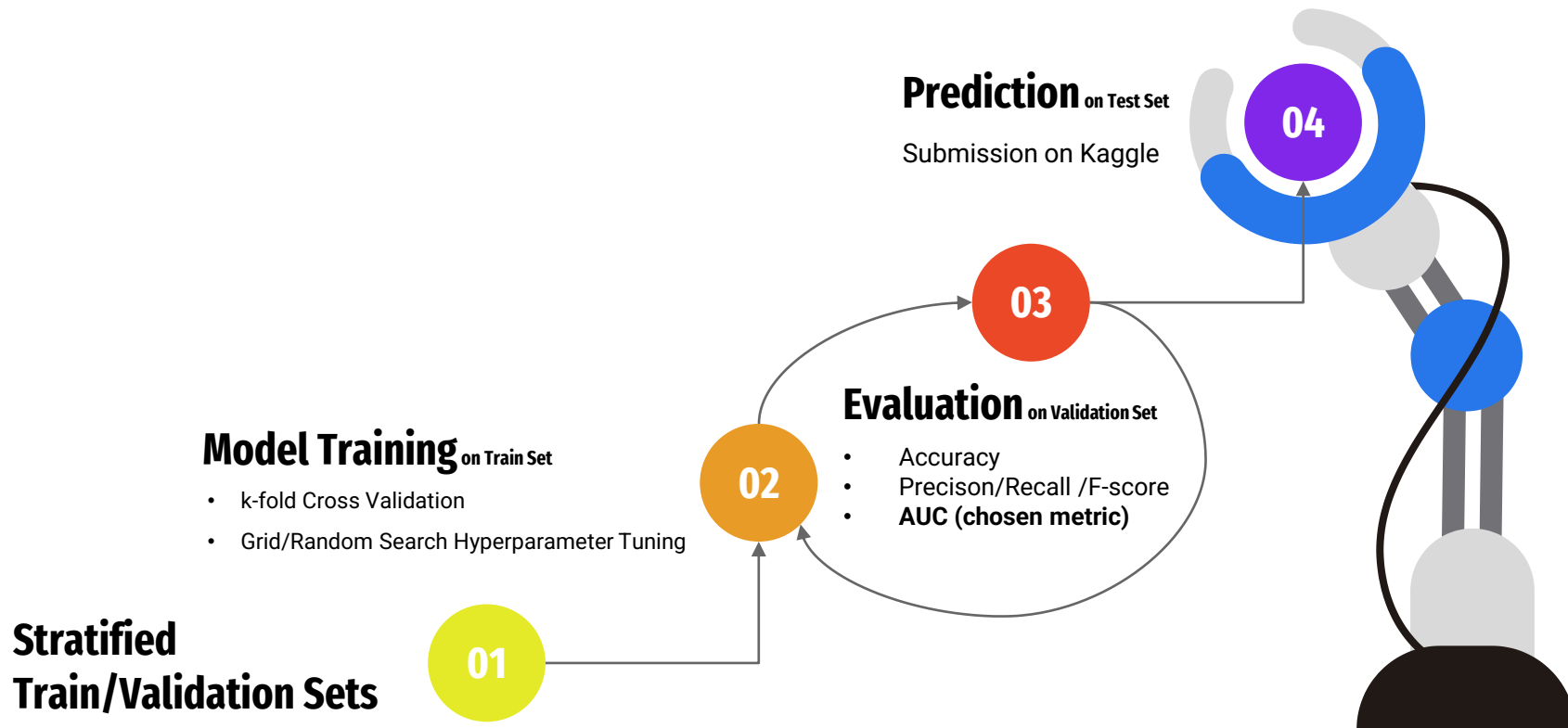
**Boruta Algorithm** (Wrapper Method)  
Copied shuffled features

**Variable Importance** with  
Rpart  
Recursive Partitioning and Regression Trees

**PCA**  
Principal Component Analysis

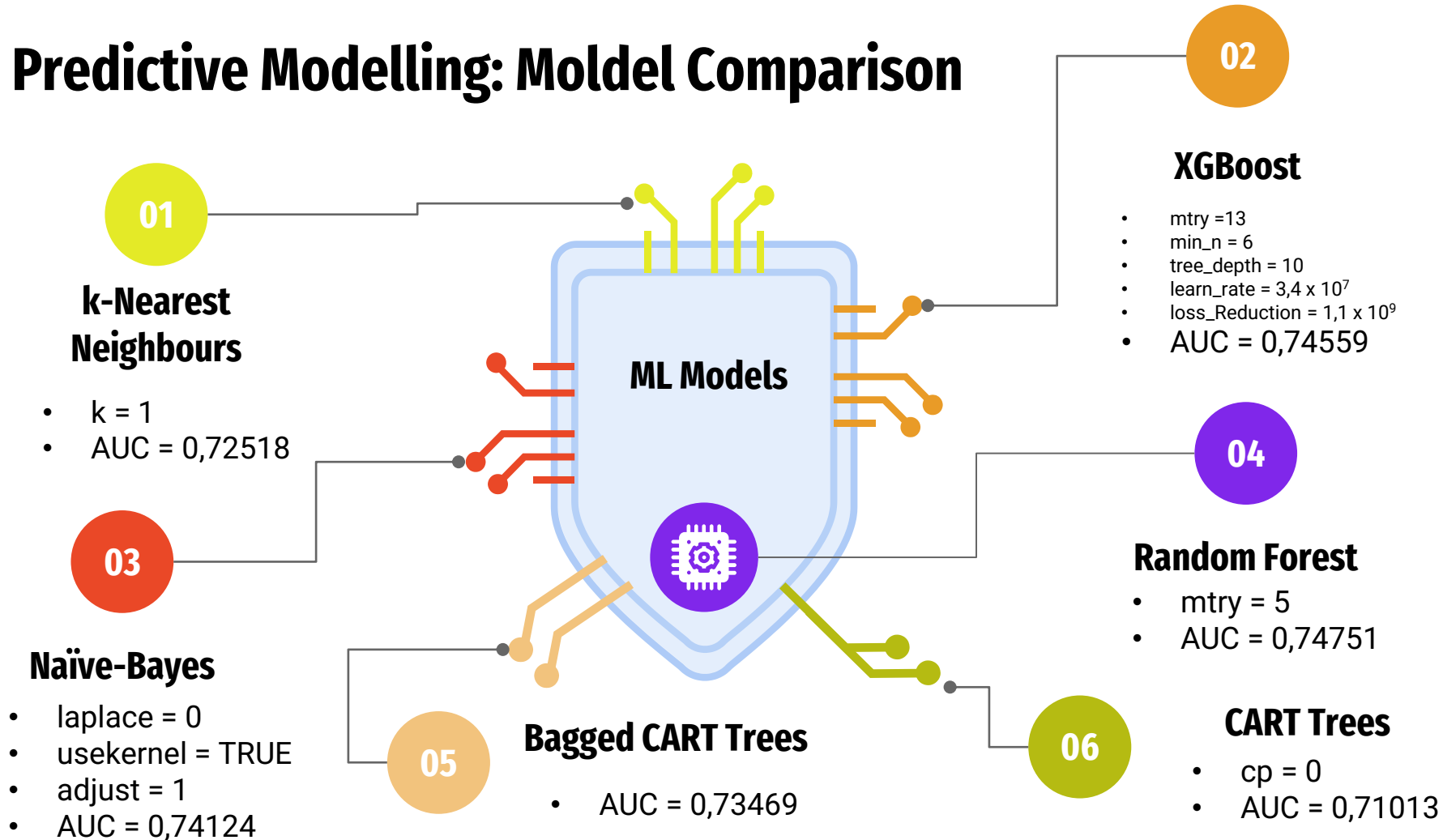
# Predictive Modelling:

Training + Hyperparameter Tuning + Evaluation Methodologies + Prediction





# Predictive Modelling: Model Comparison



# Conclusion



## Conclusions & Limitations

- Imbalanced Learning is a prevalent issue
- Recall is too low because the prevalent class is 0 in the target variable
- Accuracy is not a relevant metric
- AUC is more suited for this analysis
- Ensemble Models have better scores than primary models
- Final Chosen Model: Random Forest

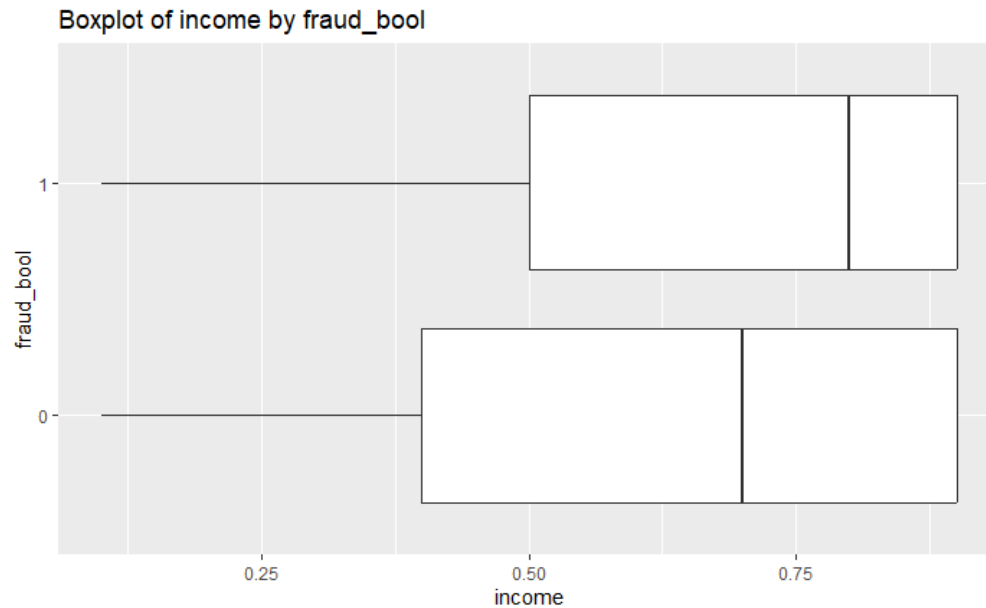


## Future Work

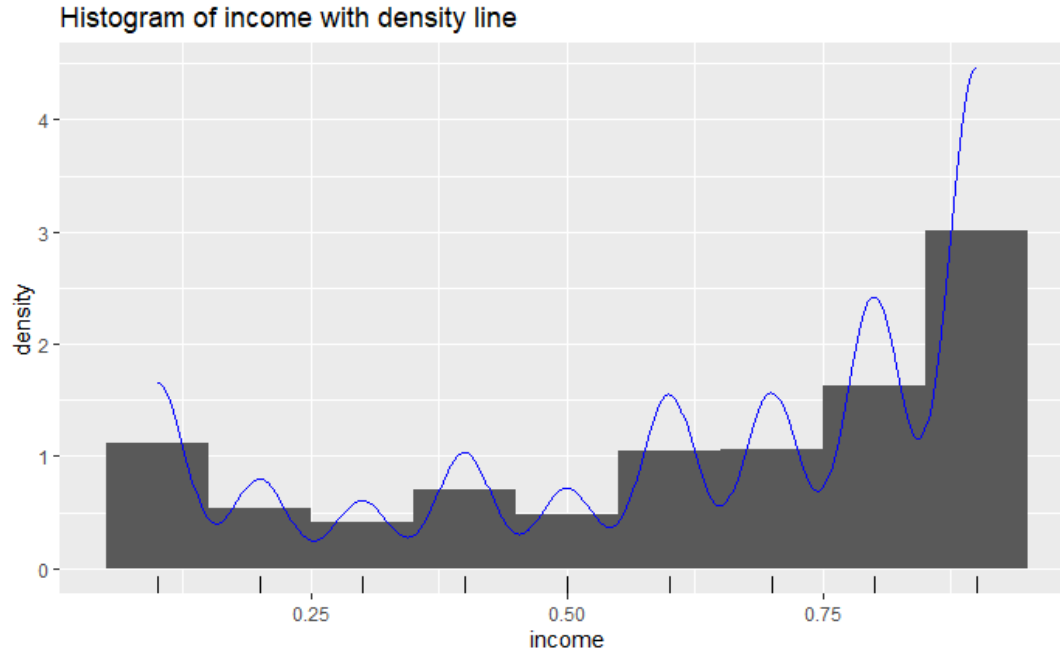
- Explore other strategies for imbalanced domains
- Train with more diverse models
- In-Depth Study to advance the feature selection techniques
- Try extracting insightful information from other clustering techniques

# **Annexes**

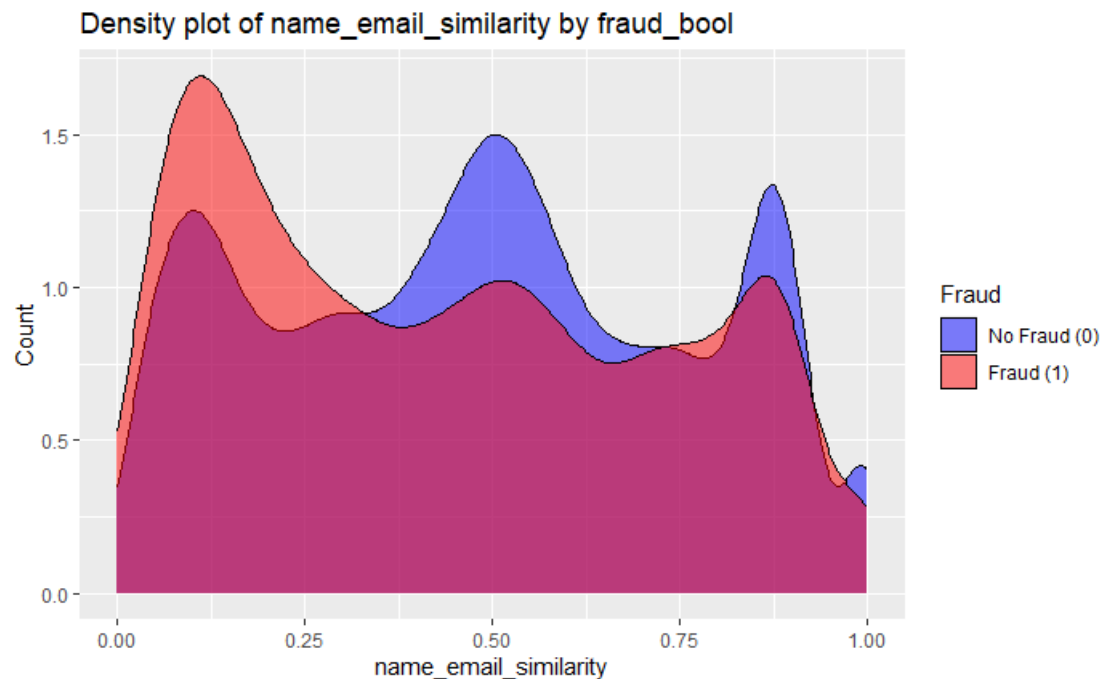
# Annex A: Income Boxplot by fraud\_bool



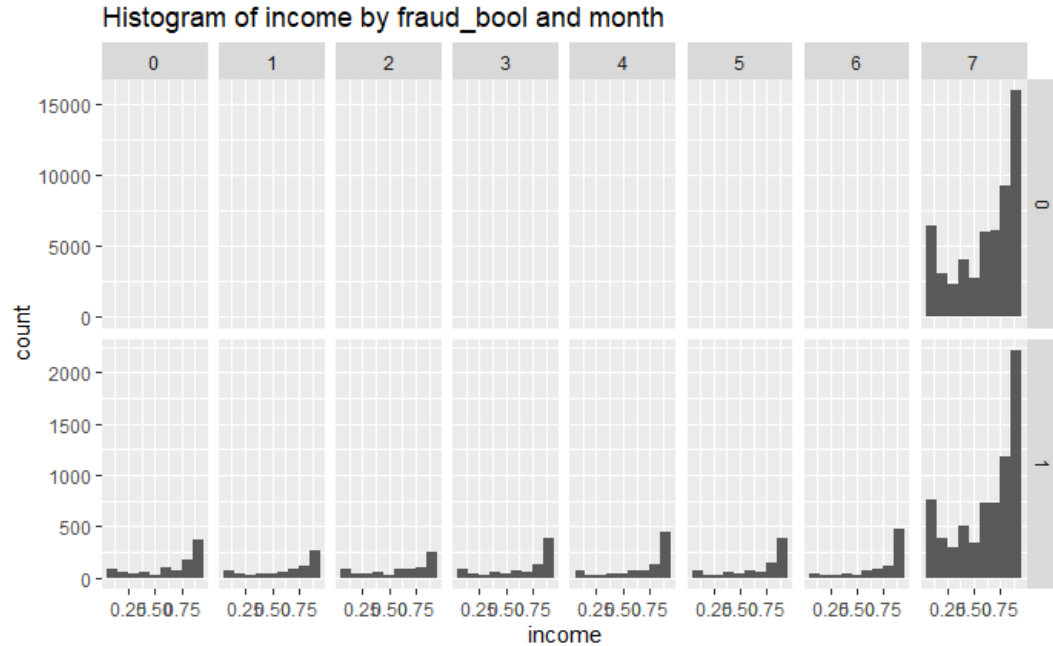
# Annex B: Income Histogram with Density Line



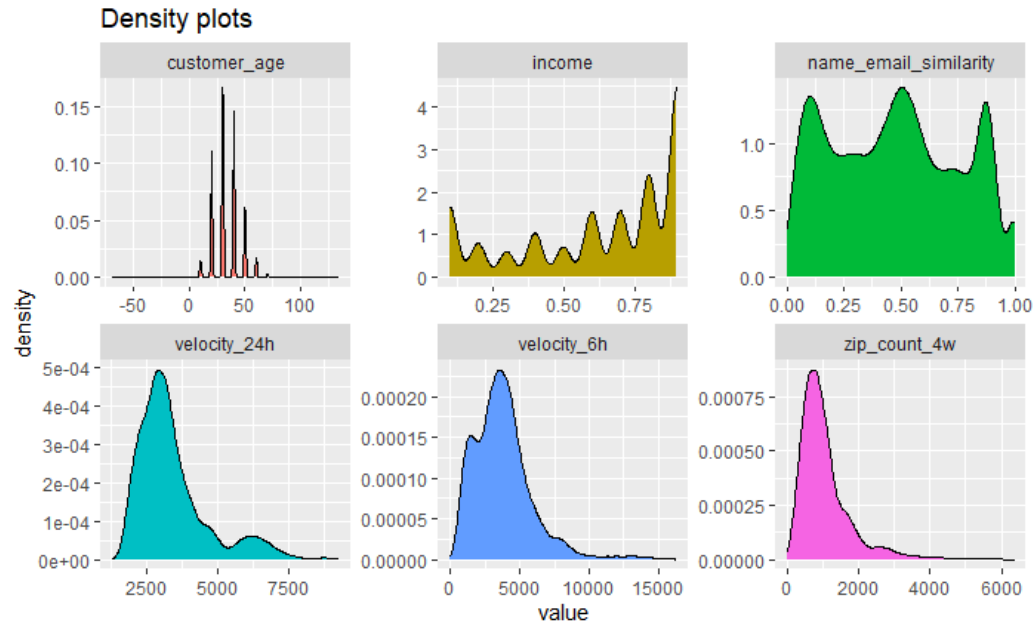
# Annex C: Density Plot of name\_email\_similarity by fraud\_bool



# Annex D: Income Histogram by fraud\_bool + month

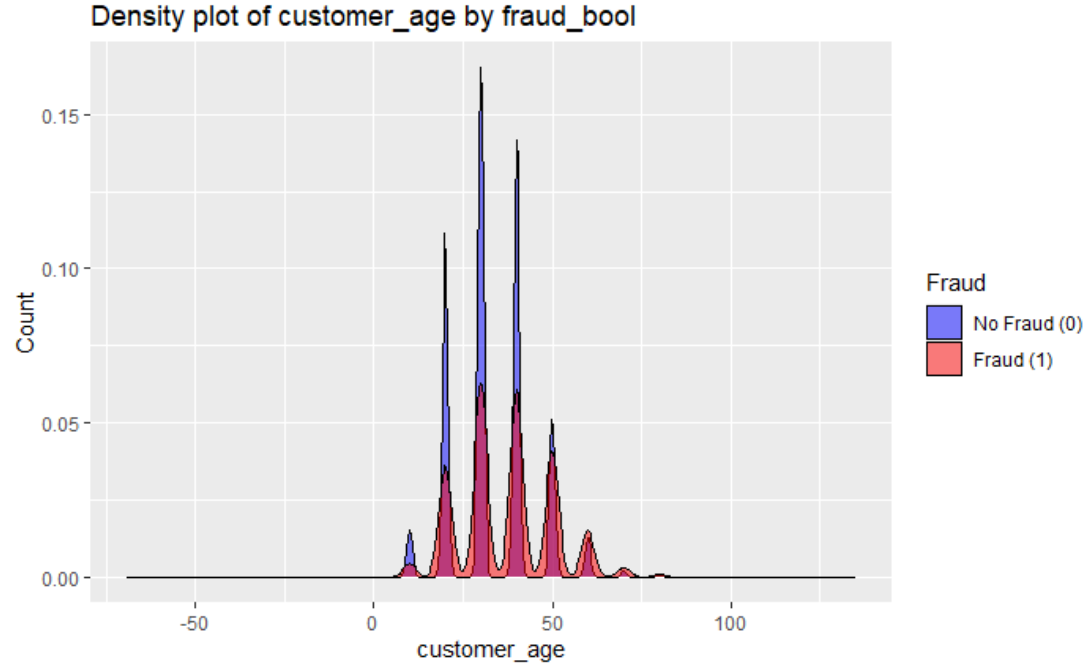


# Annex E: Density Plots

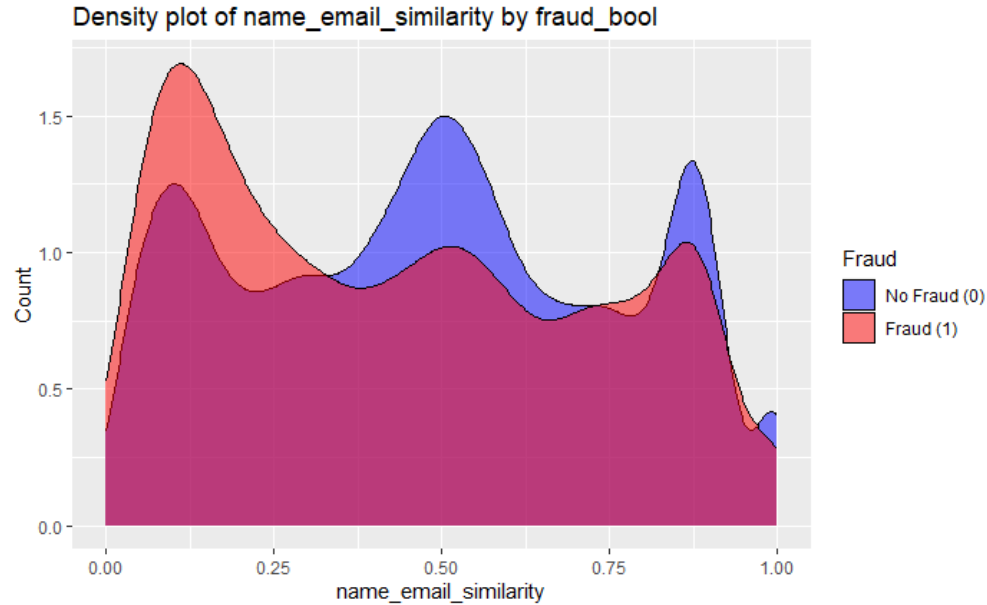




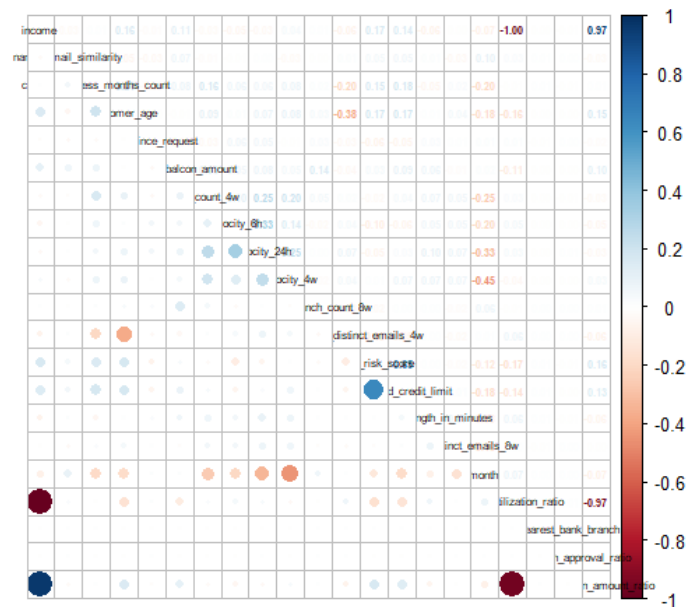
# Annex F: Density Plot of customer\_age by fraud\_bool



# Annex G: Density Plot of customer\_age by fraud\_bool



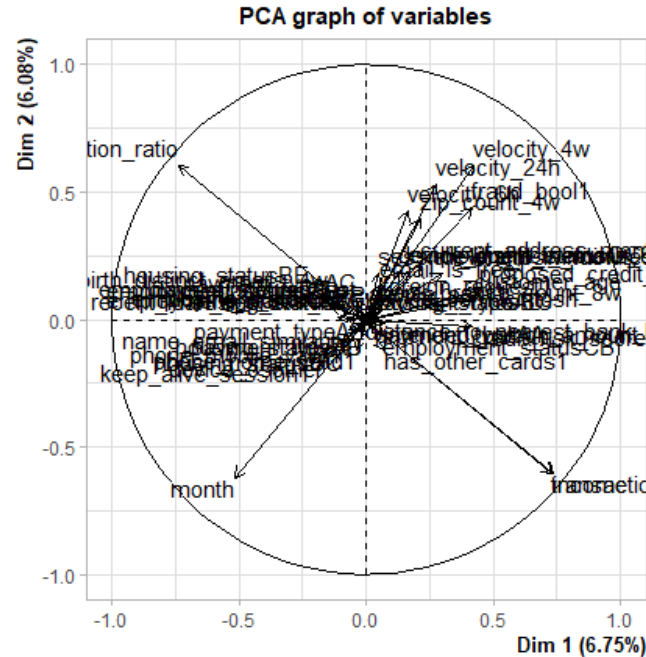
## Annex H: Correlation Matrix (Spearman Coefficients)



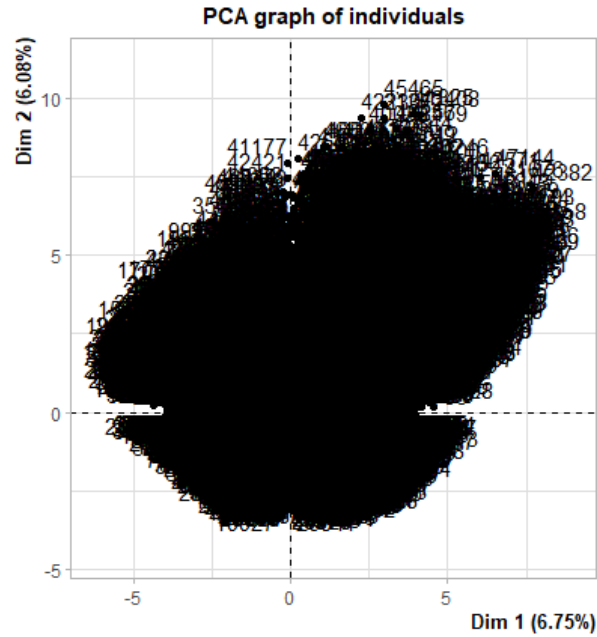
# Annex I: Chi-Square Test between payment\_type & employment\_status



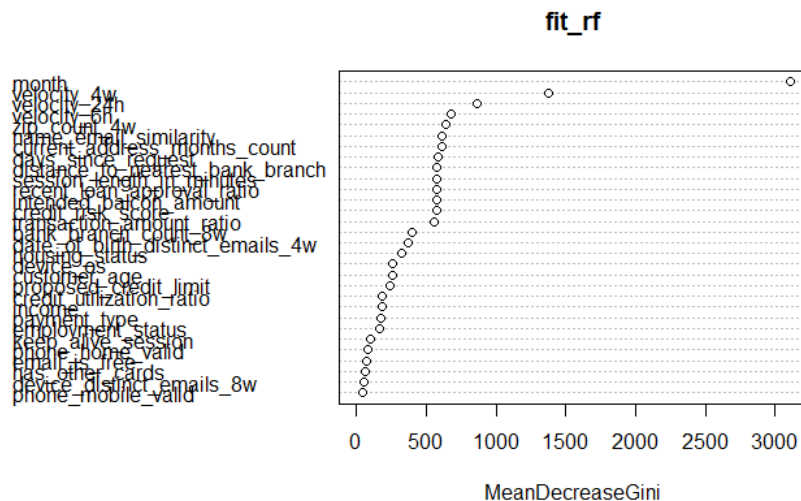
## Annex J: PCA Graph of Variables/Dimensions



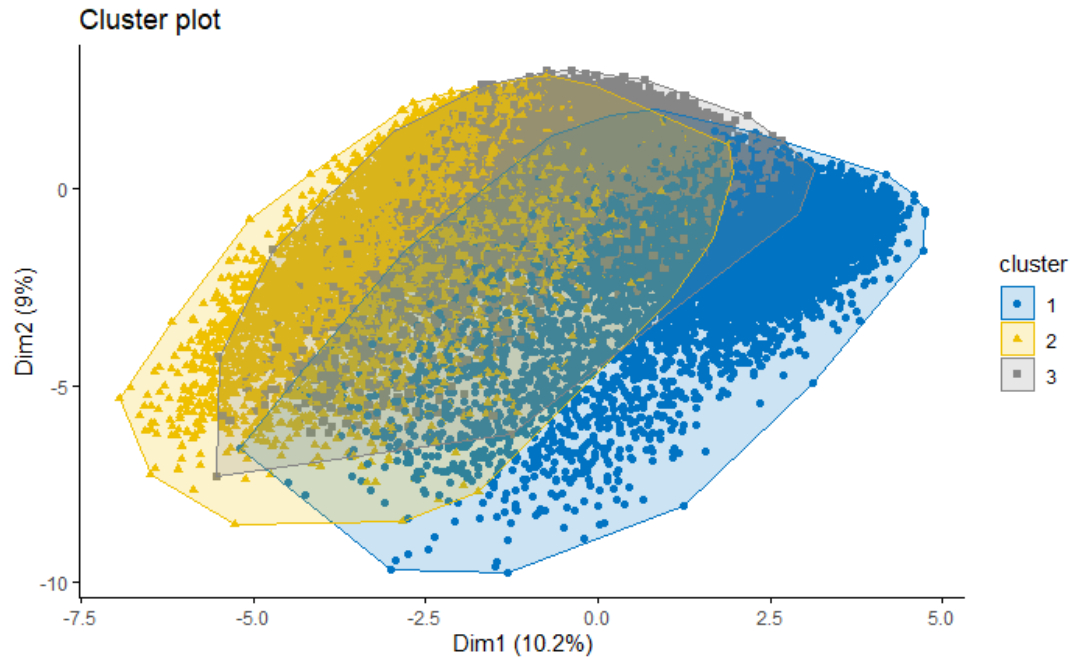
# Annex K: PCA Graph of Individuals



# Annex L: Variable Importance in a Random Forest Model

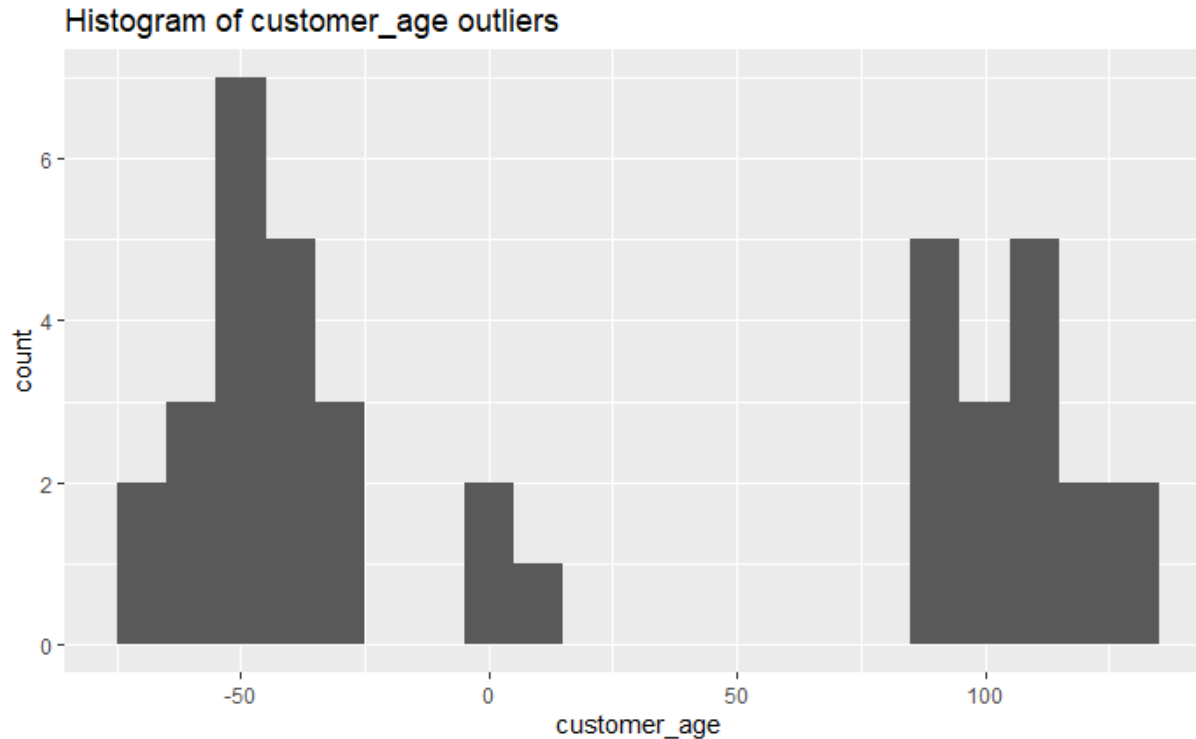


# Annex M: Cluster Plot for CLARA Clustering Algorithm with $k=3$





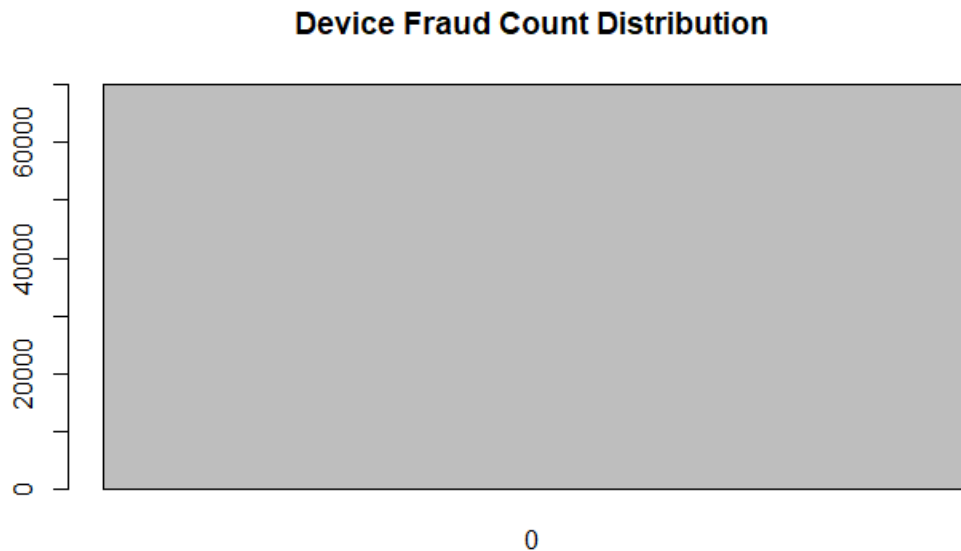
# Annex N: Histogram of customer\_age outliers



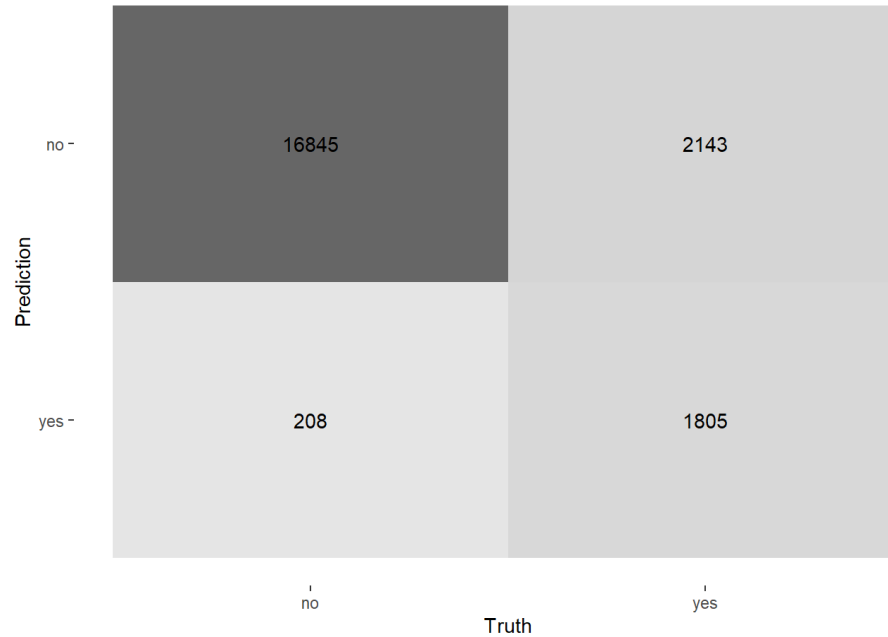
## Annex 0: fraud\_bool dataset distribution



# Annex P: device\_fraud\_count dataset distribution



# Annex Q: Confusion Matrix – Random Forest



# Annex Q: ROC Curve – Random Forest

