



Università  
di Genova

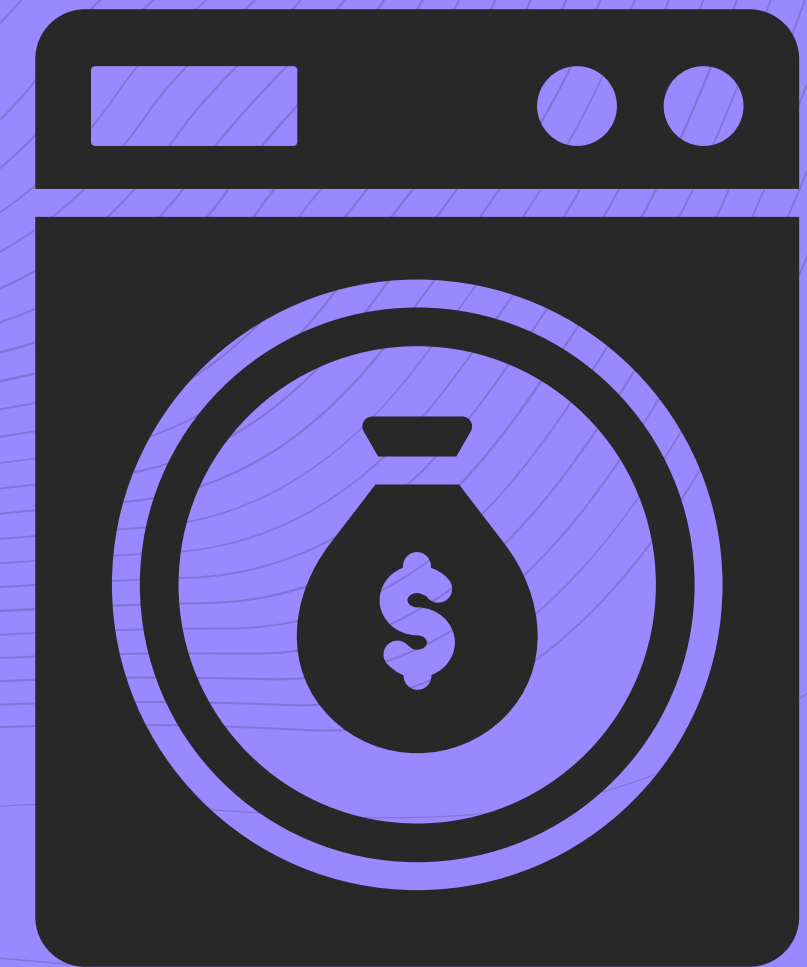
**DIBRIS** DIPARTIMENTO  
DI INFORMATICA, BIOINGEGNERIA,  
ROBOTICA E INGEGNERIA DEI SISTEMI

# BUSINESS ANALYSIS PROJECT : BANK ACCOUNT FRAUD

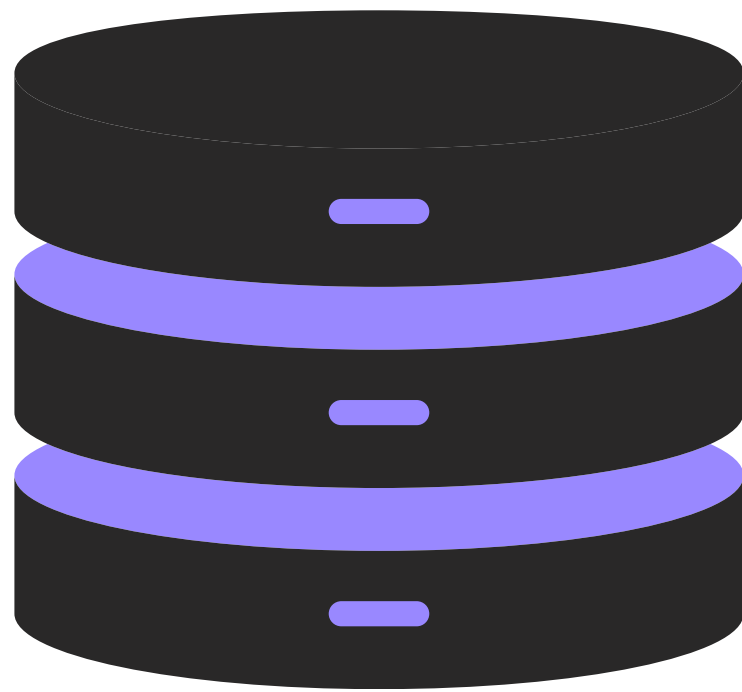
by  
**Nicolò Guainazzo**

# Understand the problem: money laundering

The problem lies in deciding whether a bank account opening request is legal or fraudulent. It is of paramount importance, in my opinion, for our domain (banking domain) to have the least number of false positives.



# Bank Account Fraud Dataset

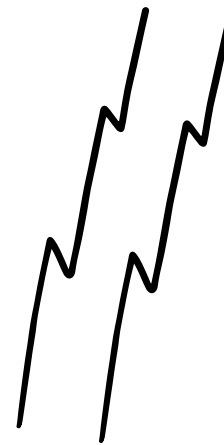


UniGe | DIBRIS



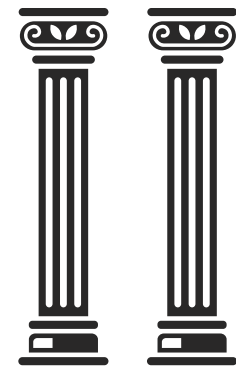
## OVERVIEW

- Size: 1 milion rows x 32 features
- No missing values (only columns with all values equal to zero)
- Highly imbalanced (approx. 1% fraud labels)
- Low correlation between the features
- Found on Kaggle



## ROWS

Each row in the dataset represents a request ,for a new bank account, labeled as good or fraudulent.



## COLUMNS

5 categorical feature(*payment\_type, source, employment\_status, housing\_status, device\_os*) 28 binary or numerical feature (e.g. *income, foreign\_request, session\_lenght\_in\_minute*).

# Problem definition in Machine learning terms

- The problem is a **supervised problem** and in general is a **unbalanced binary classification** problem.
- Many types of model can be use (*Gradient boost, ANN, One-class, isolation forest*).



## GIVEN OUR PROBLEM:

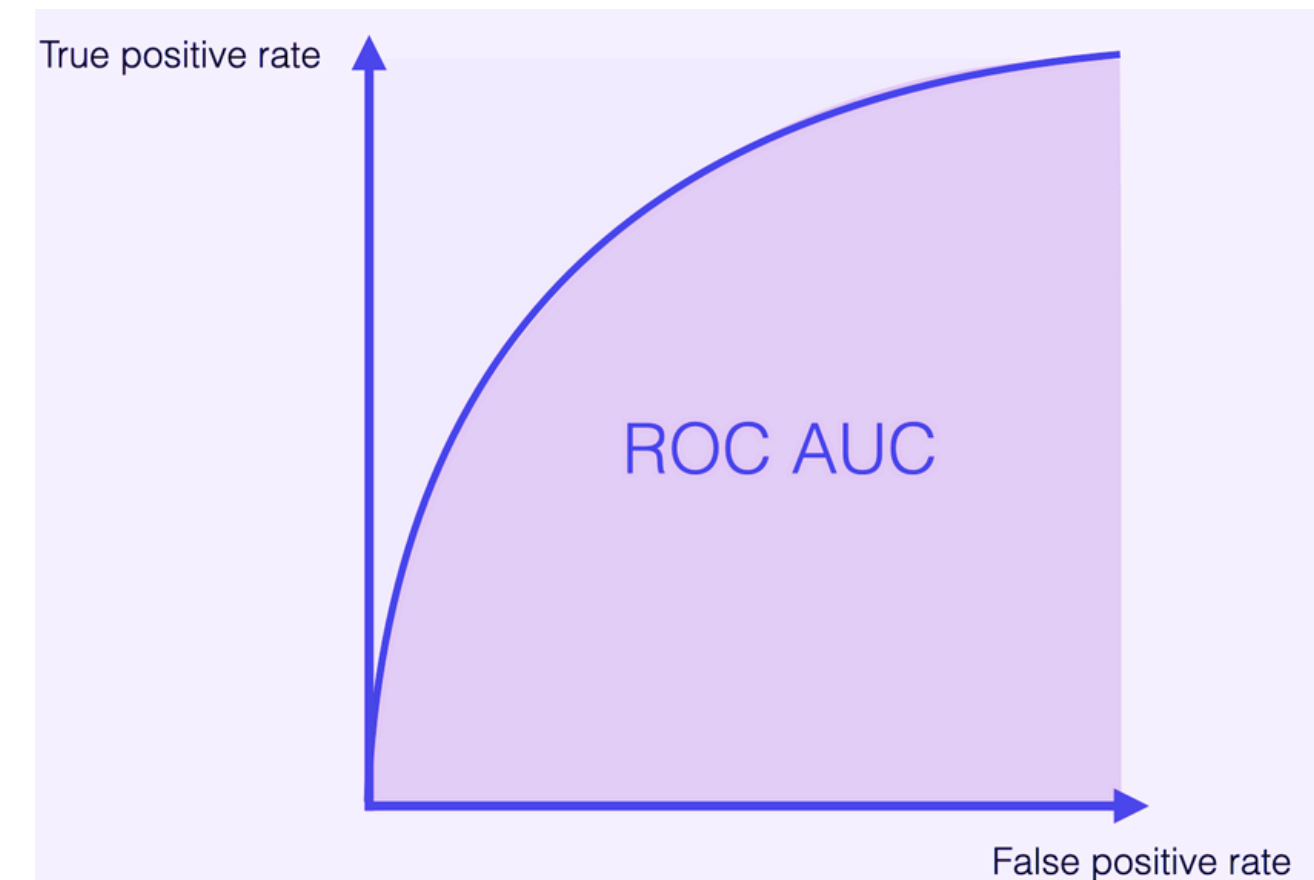
I chose the Random Forest model because, first of all, it does **not require data preparation** and also it is **robust and efficient**. In general, it is known to **work well with unbalanced data** as in our case.



# Metrics

## Confusion Matrix, Roc curve, Auc

Because of the imbalance of the data set, we cannot use the usual metrics (e.g. *accuracy*) to evaluate the model.





## 1<sup>st</sup> ATTEMPT

### Try using the raw data

Achieve Auc = 0.85 and a very large number of false positives

## 2<sup>nd</sup> ATTEMPT

### Try scaling, delete categorical features and weight of classes

Achieve Auc = 0.88 and a better balance between false and true positives

## 3<sup>rd</sup> ATTEMPT

### Try with resampling

Undersampling done using NearMiss method.

Achieve Auc = 0.94 and detect more true positives than false.



## 4<sup>th</sup> ATTEMPT

### Try with oversampling

Oversampling to rebalanced the dataset done with SMOTE method. Achieve Auc = 0.87 but not a good number of false positives

Are the results as I expected ?

In general, yes, I did not expect to  
have striking results.





# What's Next? How to develop my project

UniGe | DIBRIS

Finding other ways to decrease false positives.

---

Understand which features are more important for the prediction.

---

Increase the number of trees in the Random Forest model

---

Compare the results with other results obtained with different model.

---

# Time spend

for every main steps in hours

