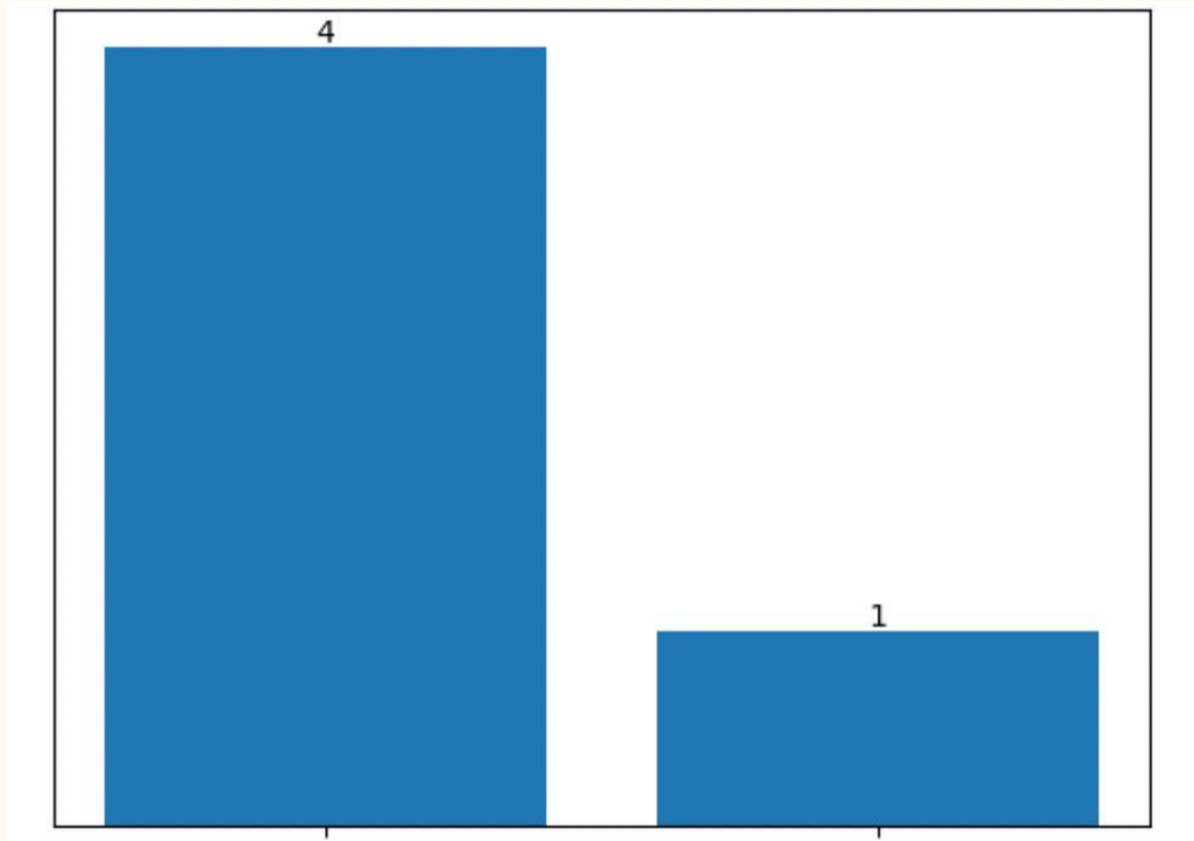# Statistics Practitioners' Challenge

*Provided by Allianz*

—

Presentation by "The Rebaggers"

# Class imbalances

# The team!



Distribution of our team

- MSc Data Science (4): Simone Moawad, Barath Raaj, Anushka Agarwal, Michele Bergami
- MSc Statistics (1): Elisabetta Sanasi

# Overview
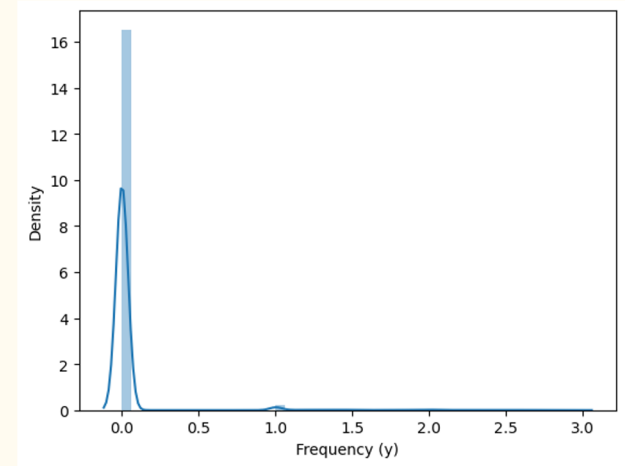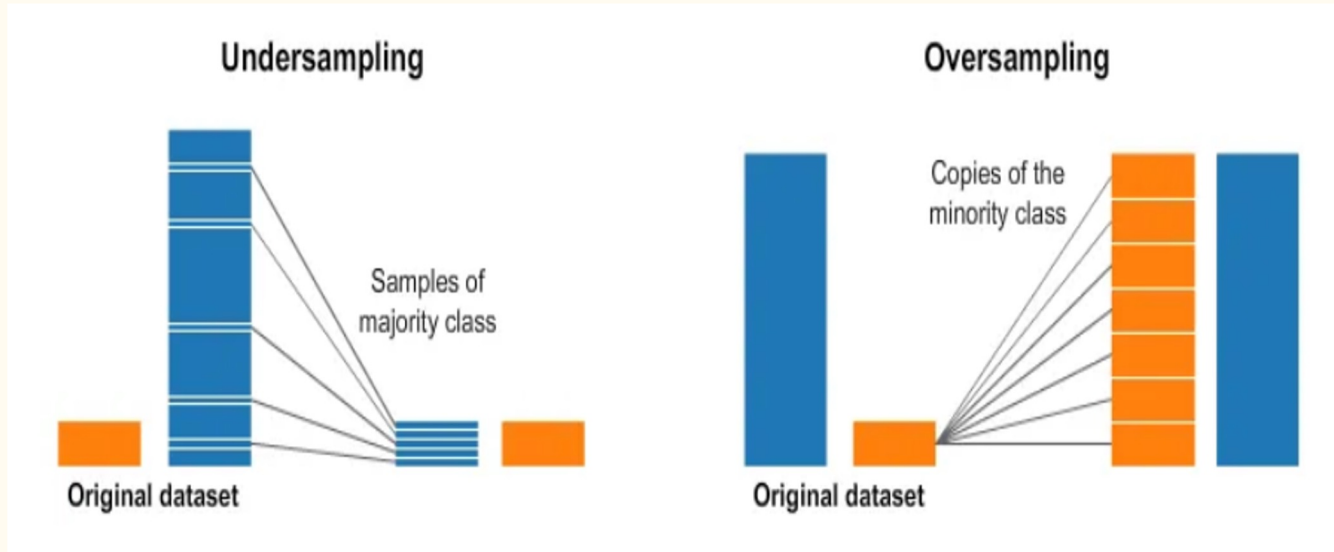
# Introduction

Importance of predicting claim frequency.

- *Fraud detection and prevention*
- *Financial planning and reserving*
- *Pricing of the policies.*
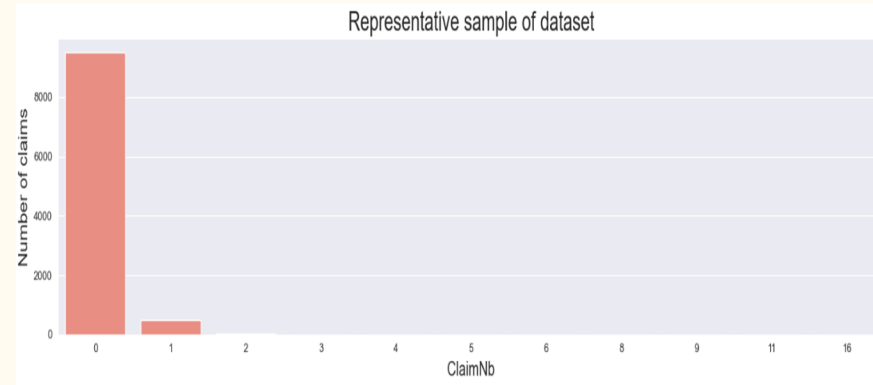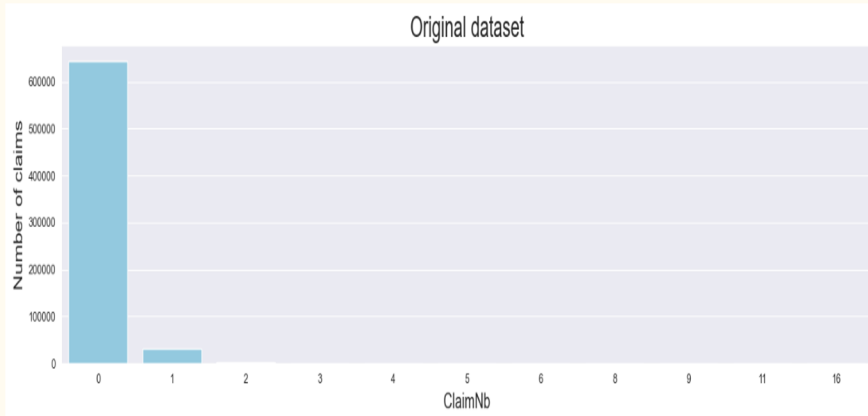
# Challenges faced

## 1. Bias and poor performance metrics



Suggested solution: Application of different resampling techniques to capture the best of both the classes and comparing it within different performance metrics.

*Performance metrics - (D2 explained, Mean absolute error, Time taken to train the model, F1 score)*

# 2. Computational complexity and resource constraints



Suggested solution: Representative sampling of the original data to replicate similar results before actually implementing the model on the actual dataset.

# Overview

# Proposed solutions

# Overview
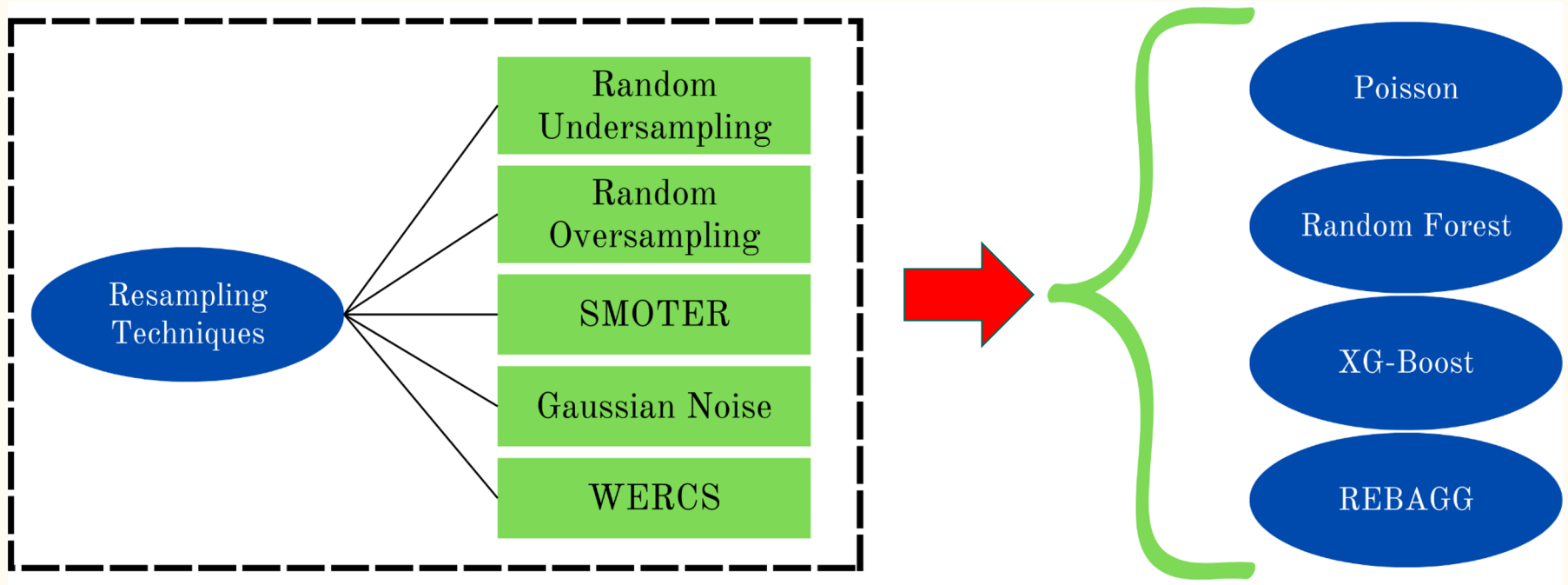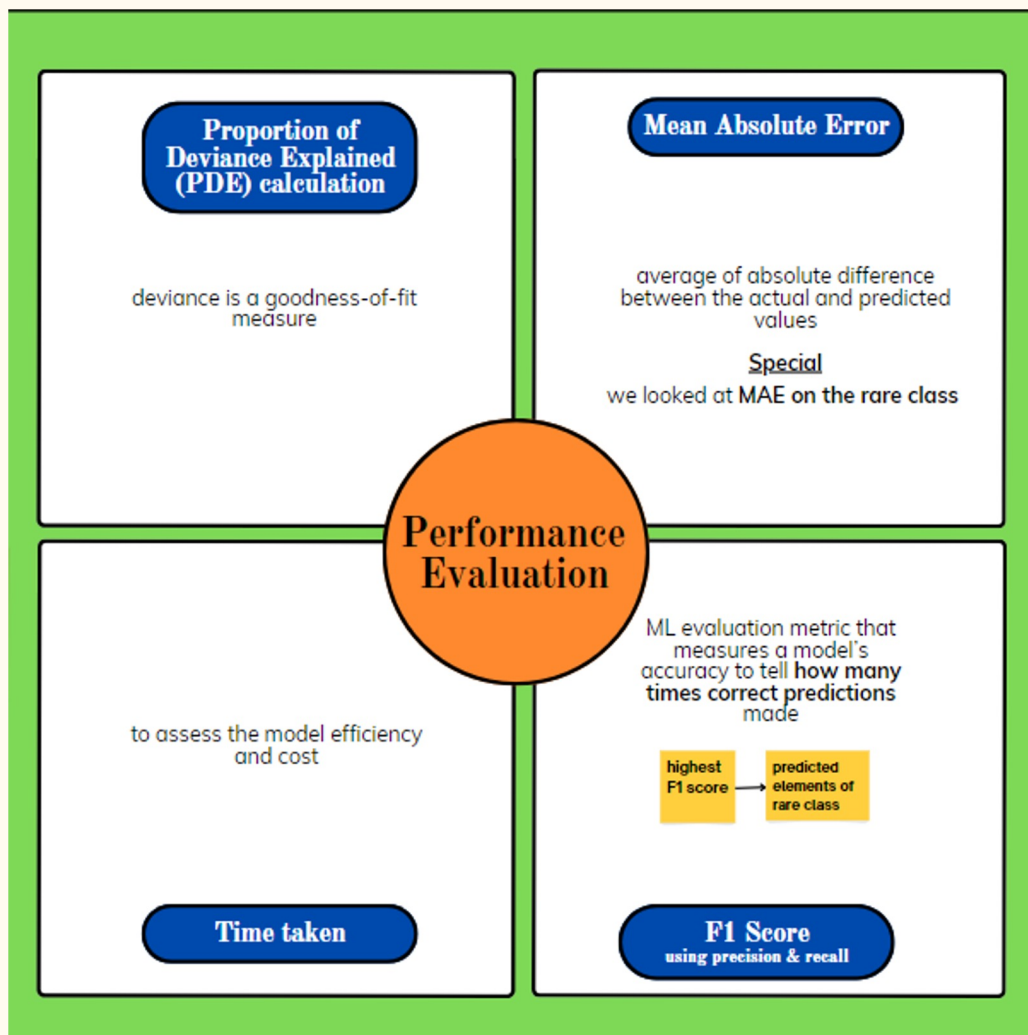
# Resampling techniques



- Random Undersampling
- Random Oversampling
- SMOTER
- Gaussian Noise
- WERCS
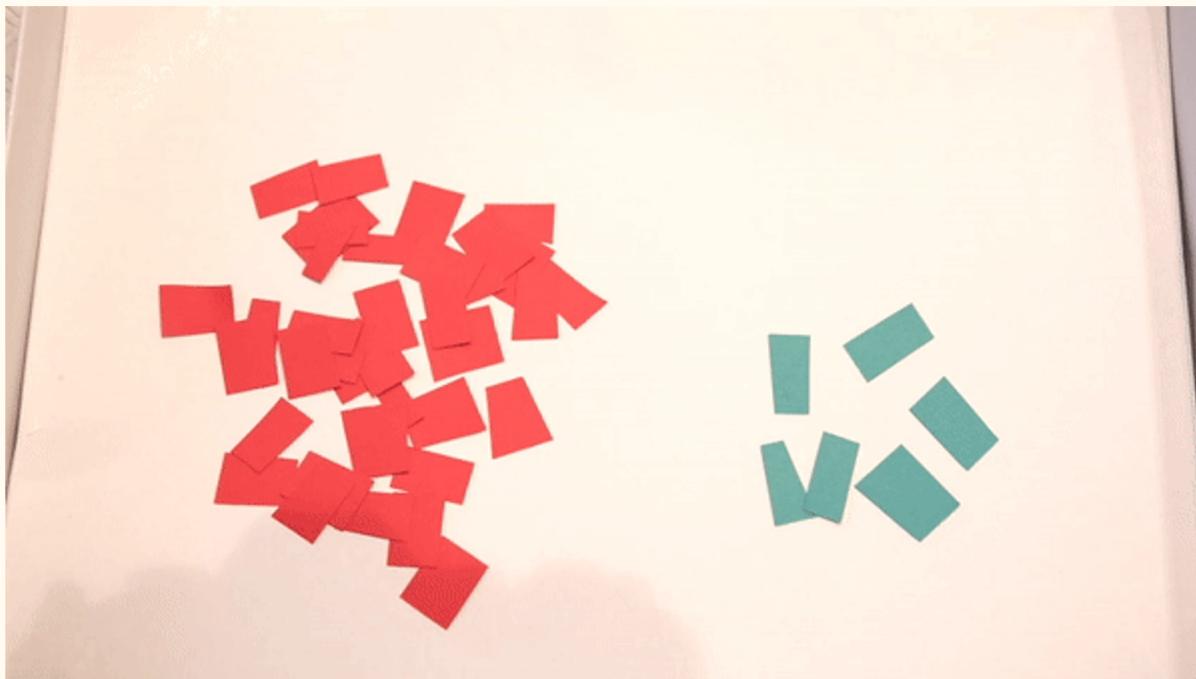
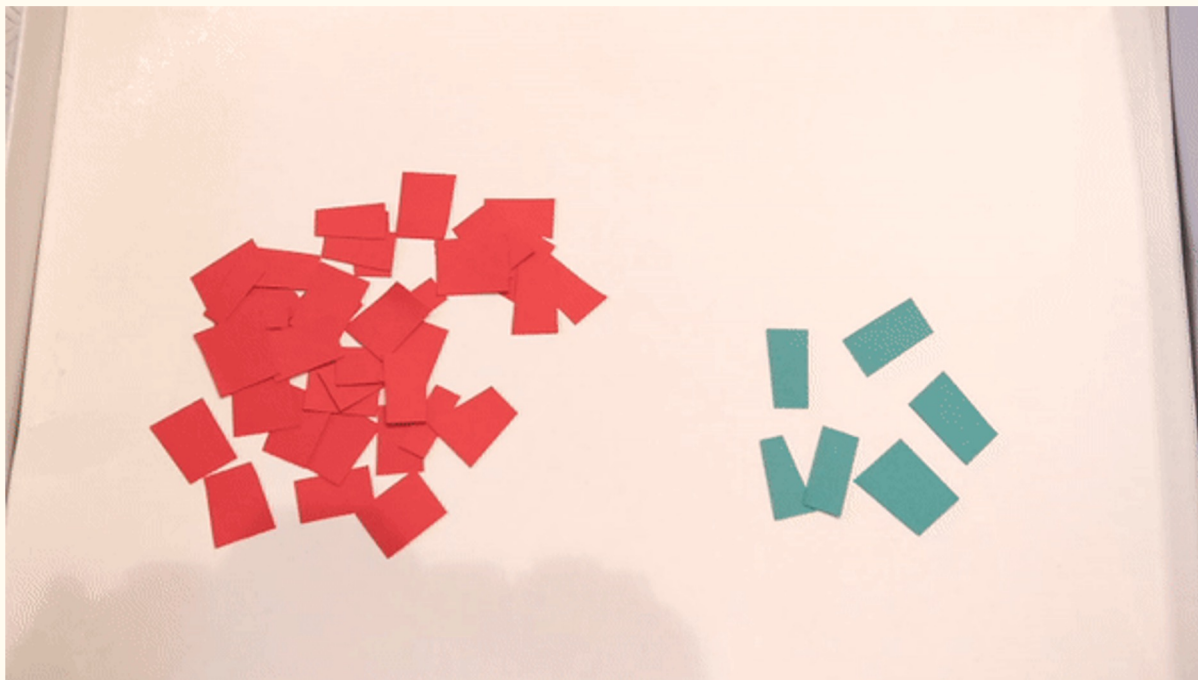**What do these techniques mean?**

# Random undersampling



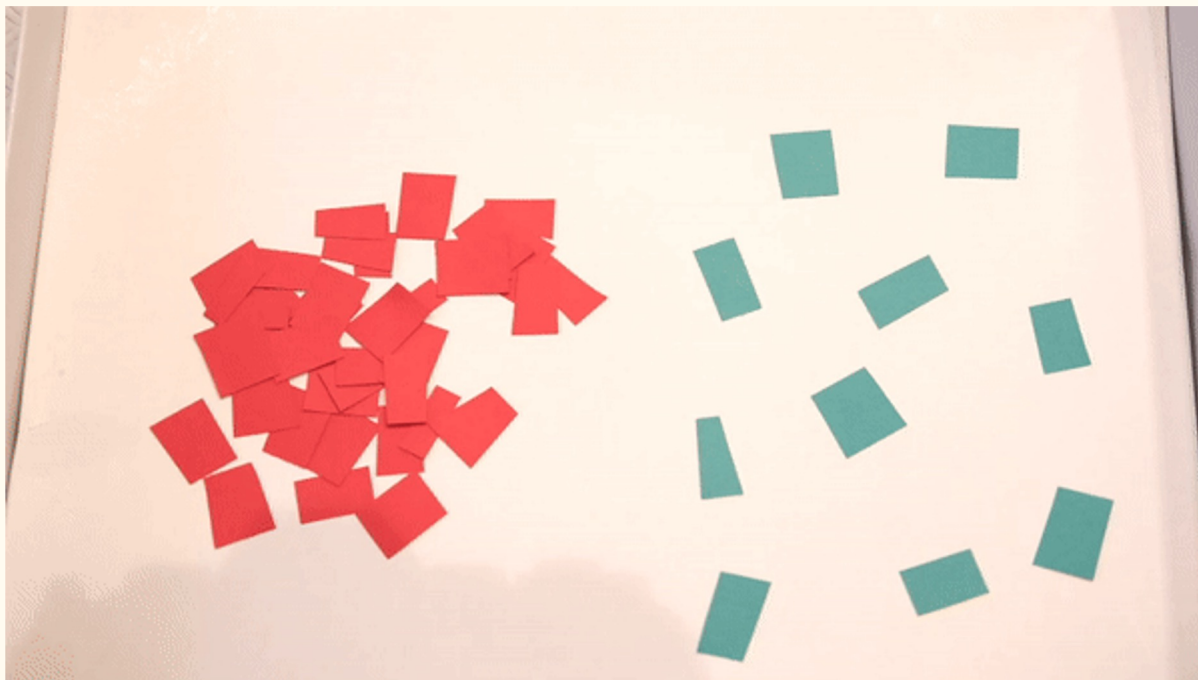| undersampled (randomly) | left untouched |

# Random oversampling



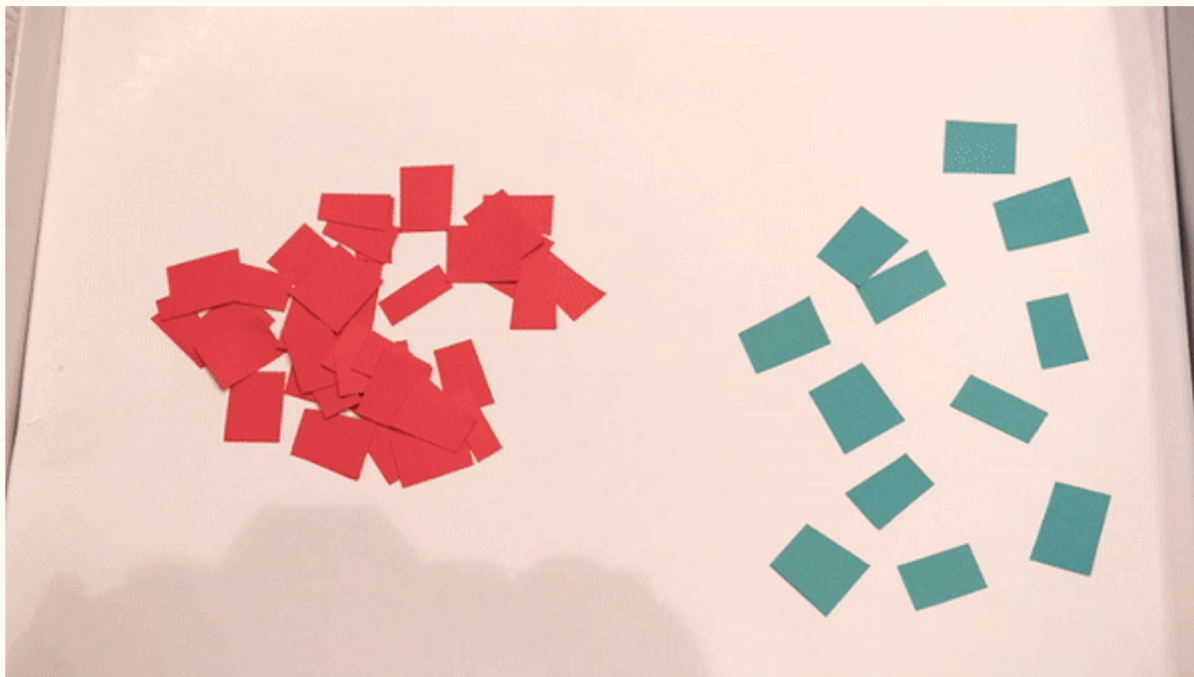left untouched

oversampled
(by duplication)

# SMOTER



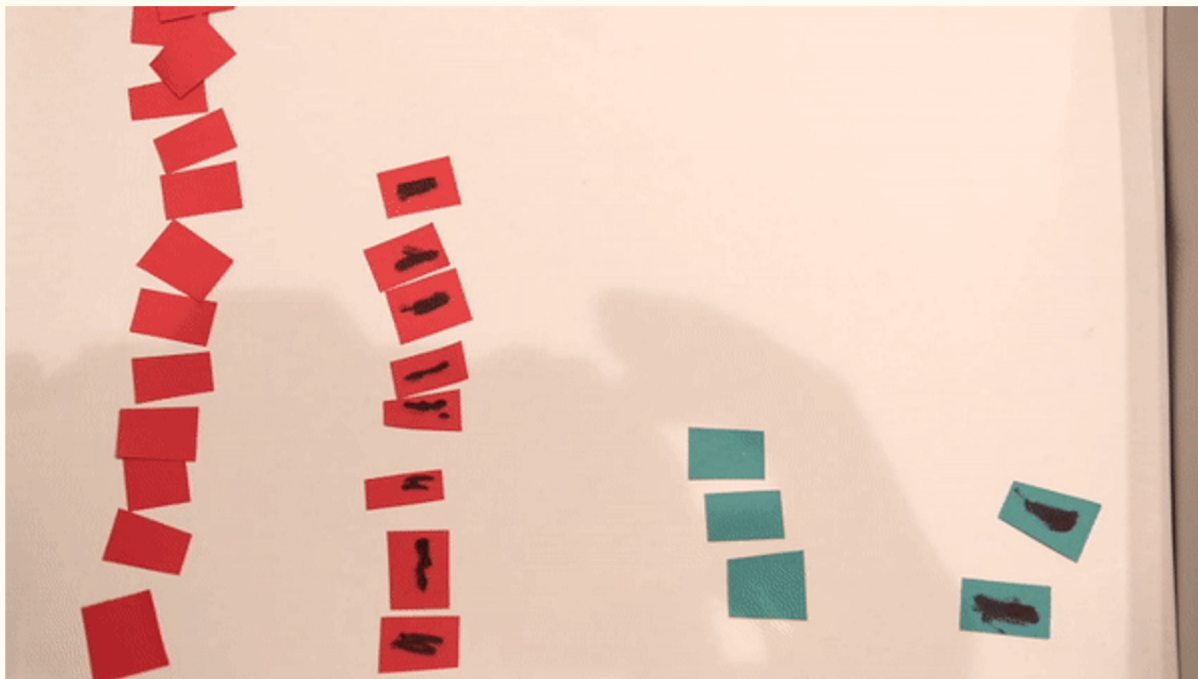| undersampled (randomly) | oversampled (using kNN) |

# Gaussian noise



| undersampled | oversampled |
|:---:|:---:|
| (randomly) | (adding gaussian noise) |

# WERCS



**combination**
of the other methods

**+**

**weights**
to the data by relevance

# Overview

# The models

| Poisson regression | Random forest | XGboost | REBAGG |
|---|---|---|---|

**Poisson regression**

💚Well-suited for count data

💔Strong assumption
💔Linear relationships

**Random forest**

💚Robust method against overfitting

💔Computationally complex

**XGboost**

💚Efficient and scalable

💔Tuning of hyperparameters

**REBAGG**

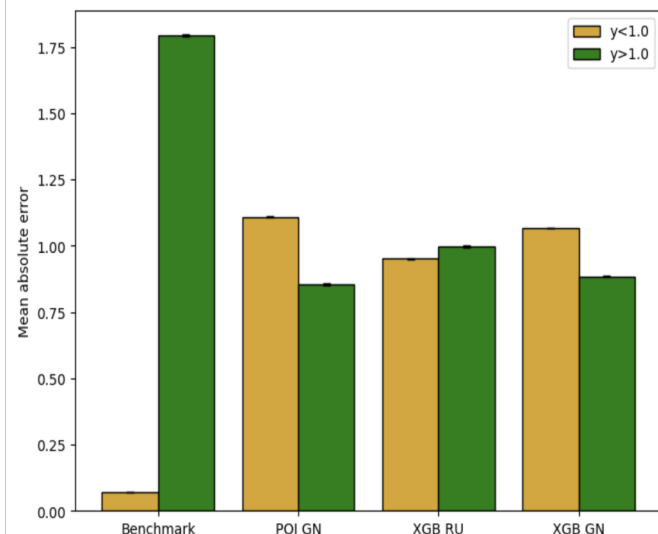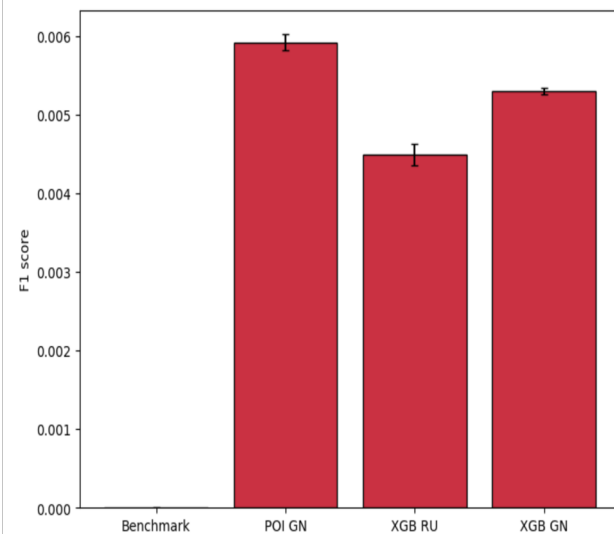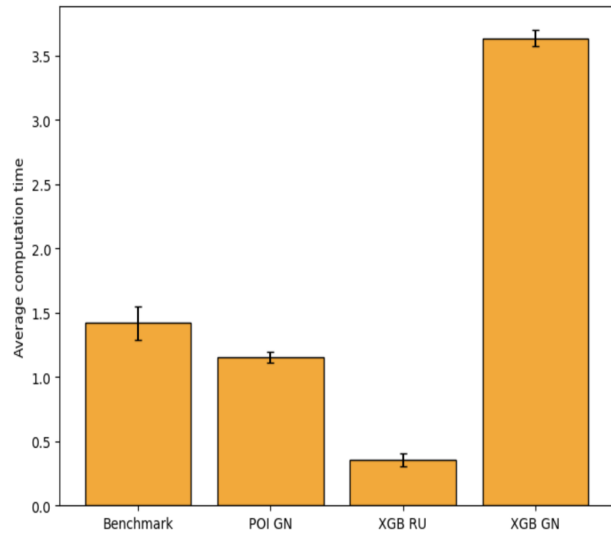💚Stability and accuracy
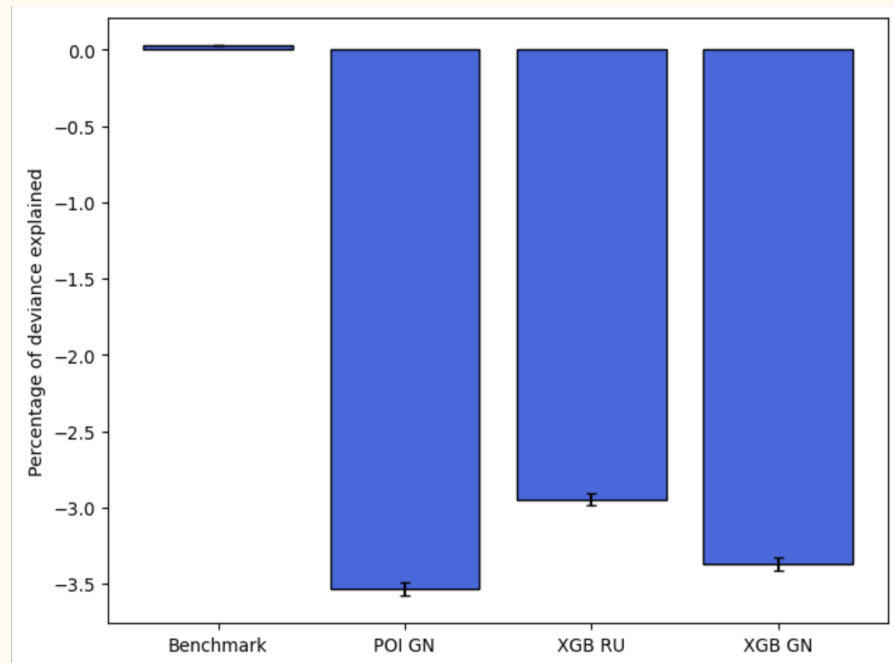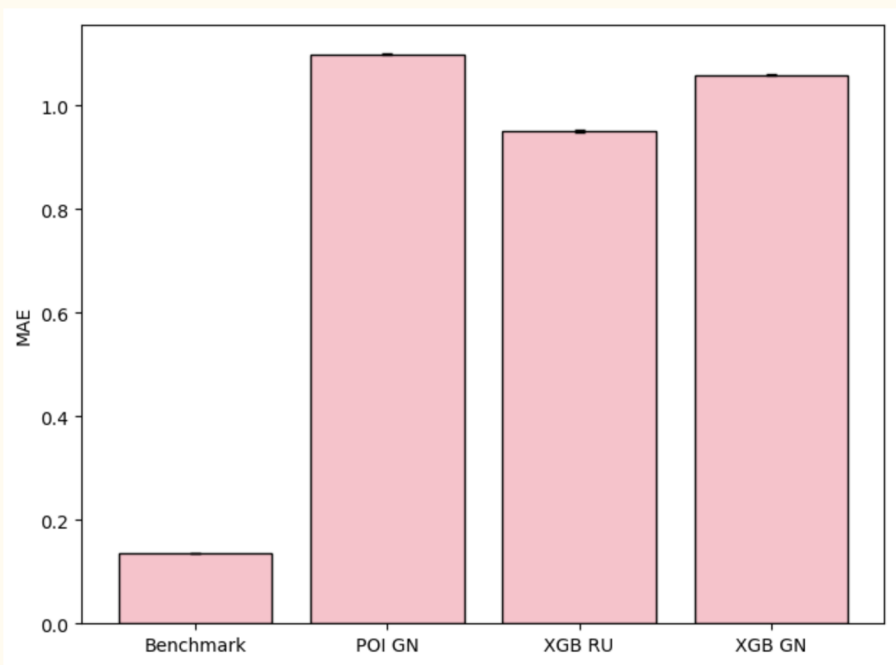
💔Computational resources and time for training

# Which is the best model?

# Every model!

Apart from WERCS + Random Forest, that's terrible :(

Performance on the rare class

Performance on the rare class

# General performance

General performance

The compromise

The compromise

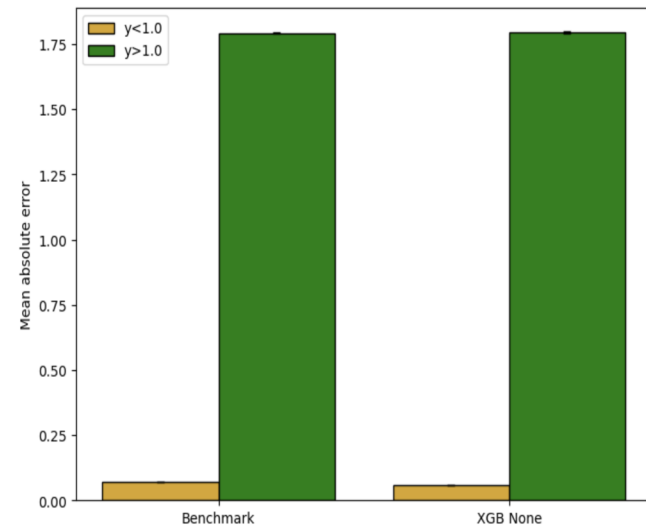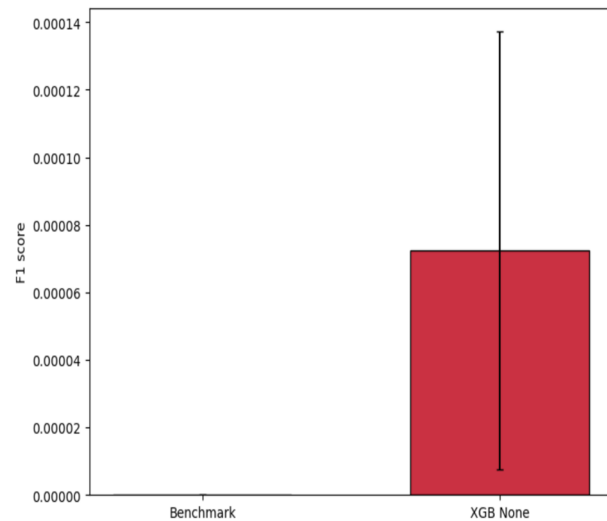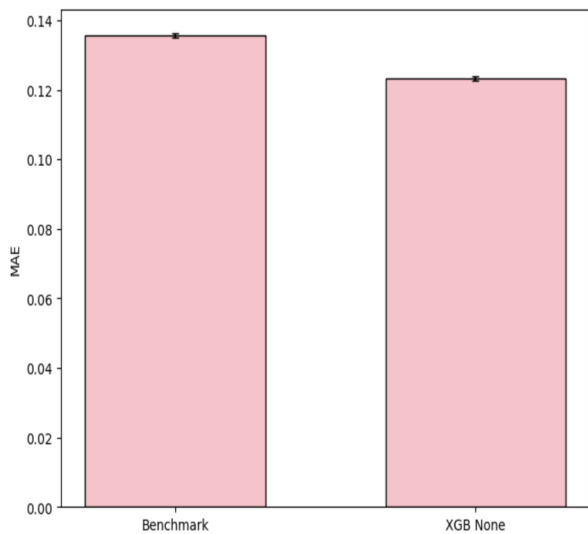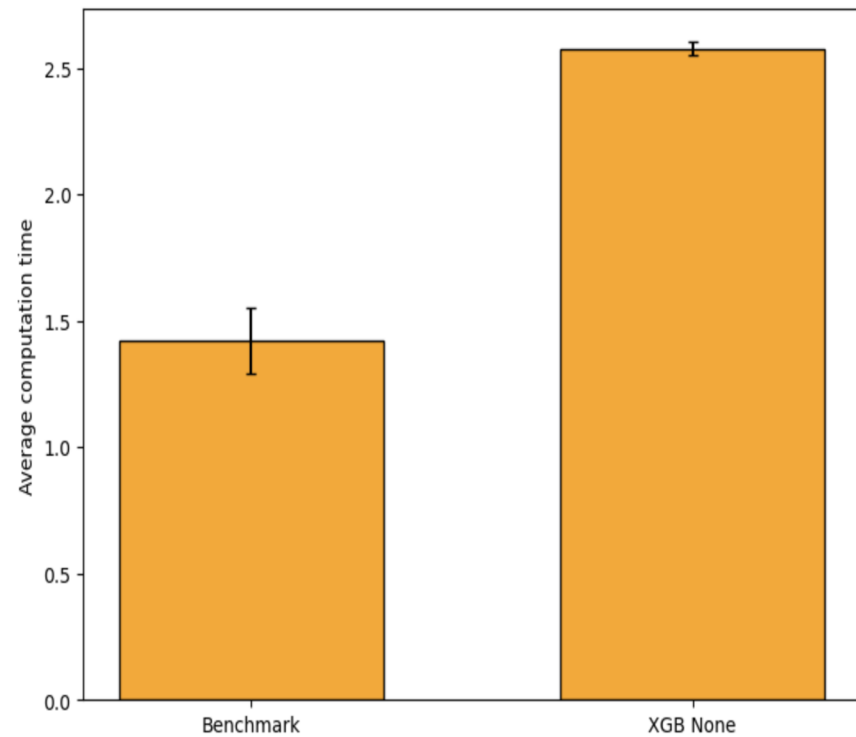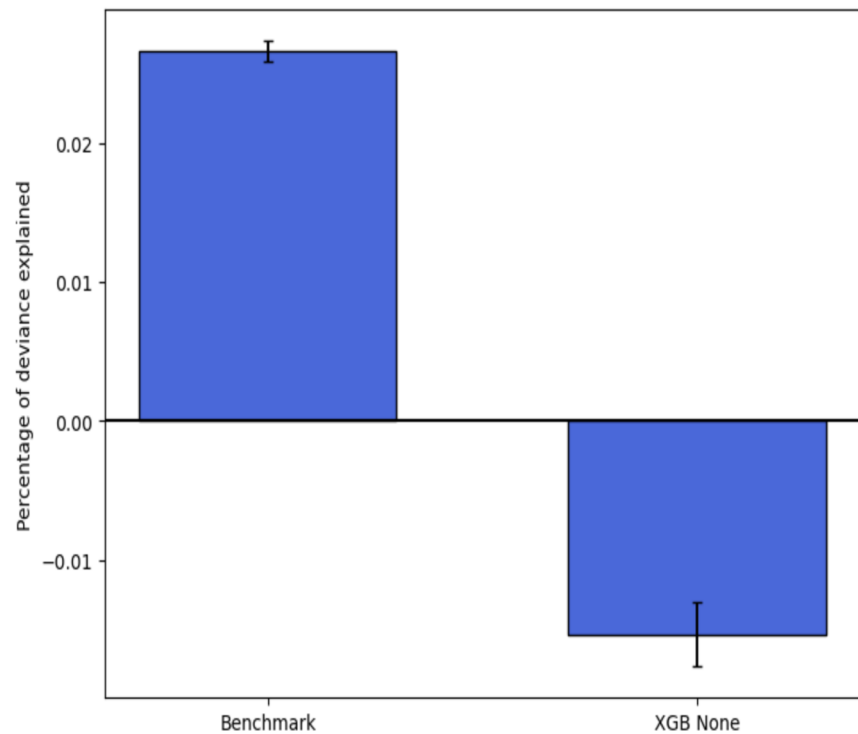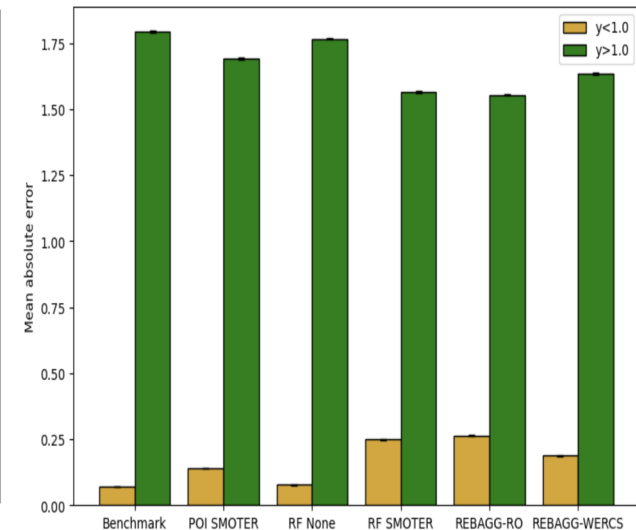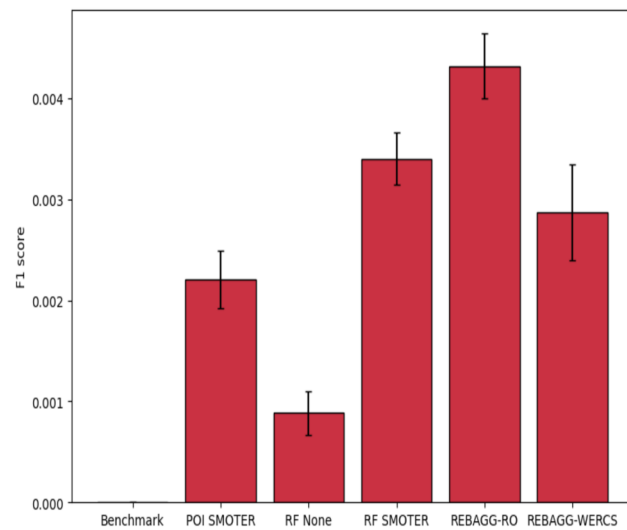|  | PDE | F1 Score | MAE on rare class | MAE | Time |
|---|---|---|---|---|---|
| **General performance** | | | | | |
| Benchmark | 0.02660 | 0.00000 | 1.79269 | 0.13565 | 1.42199 |
| XGB None | -0.01537 | 0.00007 | 1.79514 | 0.12326 | 2.57665 |
| | | | | | |
| **Performance on the rare class** | | | | | |
| Benchmark | 0.02660 | 0.00000 | 1.79269 | 0.13565 | 1.42199 |
| POI GN | -3.53659 | 0.00592 | 0.85354 | 1.09977 | 1.15724 |
| XGB RU | -2.94768 | 0.00449 | 0.99613 | 0.95204 | 0.35723 |
| XGB GN | -3.36994 | 0.00531 | 0.88476 | 1.05970 | 3.63604 |
| **Compromise** | | | | | |
| Benchmark | 0.02660 | 0.00000 | 1.79269 | 0.13565 | 1.42199 |
| POI SMOTER | -1.24956 | 0.00220 | 1.69023 | 0.19723 | 31.18506 |
| RF None | -1.69685 | 0.00088 | 1.76638 | 0.13885 | 16.67995 |
| RF SMOTER | -0.93917 | 0.00340 | 1.56522 | 0.29707 | 51.94275 |
| REBAGG-RO | -0.94582 | 0.00432 | 1.55349 | 0.31093 | 43.82252 |
| REBAGG-WERCS | -1.14715 | 0.00287 | 1.63469 | 0.24052 | 58.96560 |

The results above are obtained on a device with Apple M1 Chip (8-core CPU, 7-core GPu and 16-core Neural Engine)

# Overview

## Pros

Differentiation of strategies with improvements on MAE in general and on the rare class

Overall computational intensity of selected models that is comparable to the benchmark

## Cons

Very low PDE, probably due to the resampling itself

There is not a unique model capable of addressing this model without downsides

## Further analysis

Gridsearch on models hyperparameters

Simultaneous gridsearch of resampling techniques and model hyperparameters to possibly achieve an optimal model

# Thank you for your attention!