# Predicting Diabetes

DiabQuest

*Aleksandra Muciek & Magdalena Buszka*

# Research Question

- How well can we predict which patient has diabetes based on various predicting factors? (BMI, weight, age,…)

- Which variables are the most important?

# Challenges

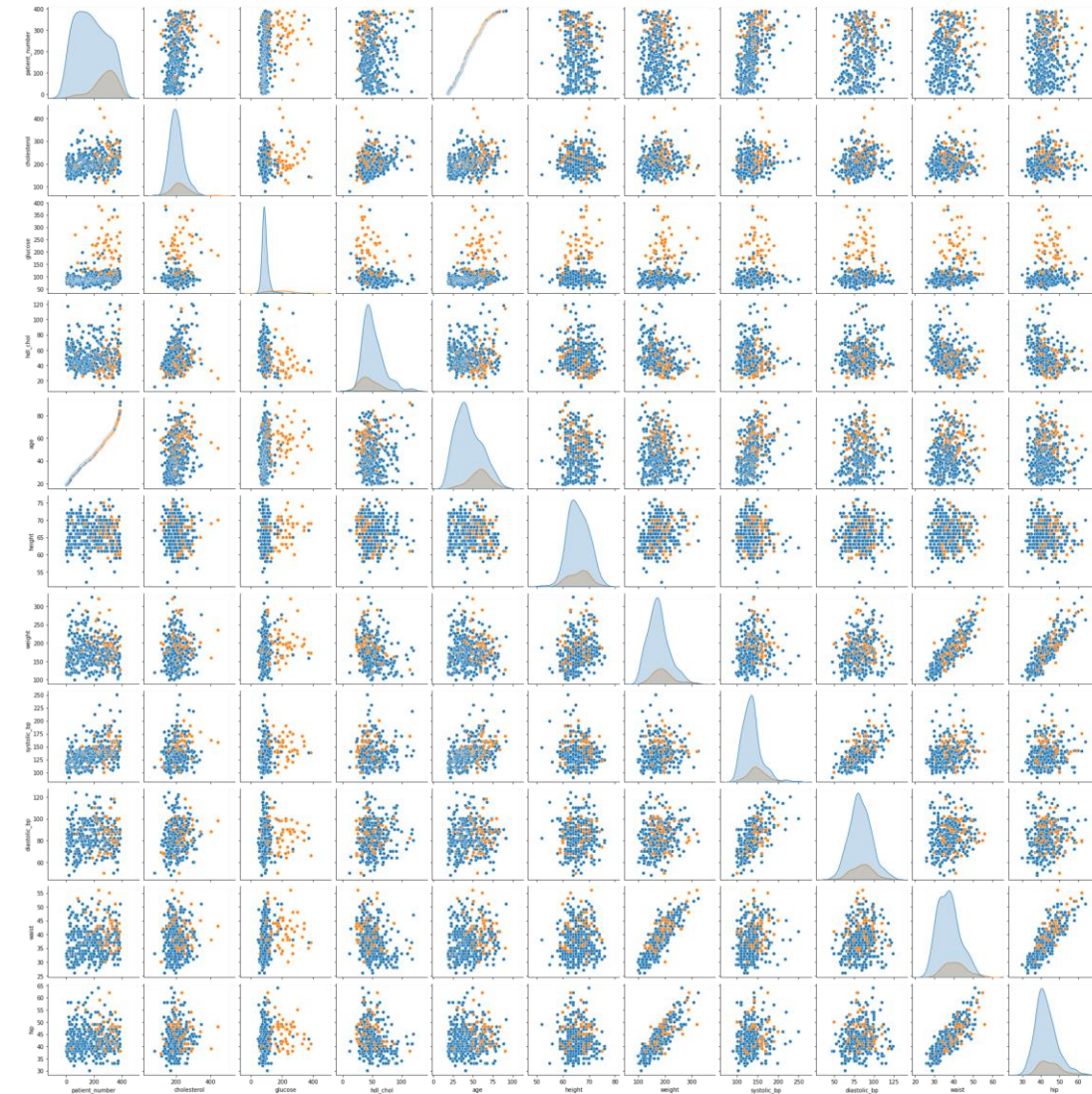- Picking the right metric for evaluation
- Dealing with imbalanced data

# Data

- Database from Kaggle: link

- Each row contains data from one patient

- For each patient we have several medical measurements, personal and anthropometric data

- Imbalanced data-set: 390 rows with only 15% from people with diabetes
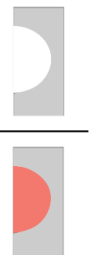
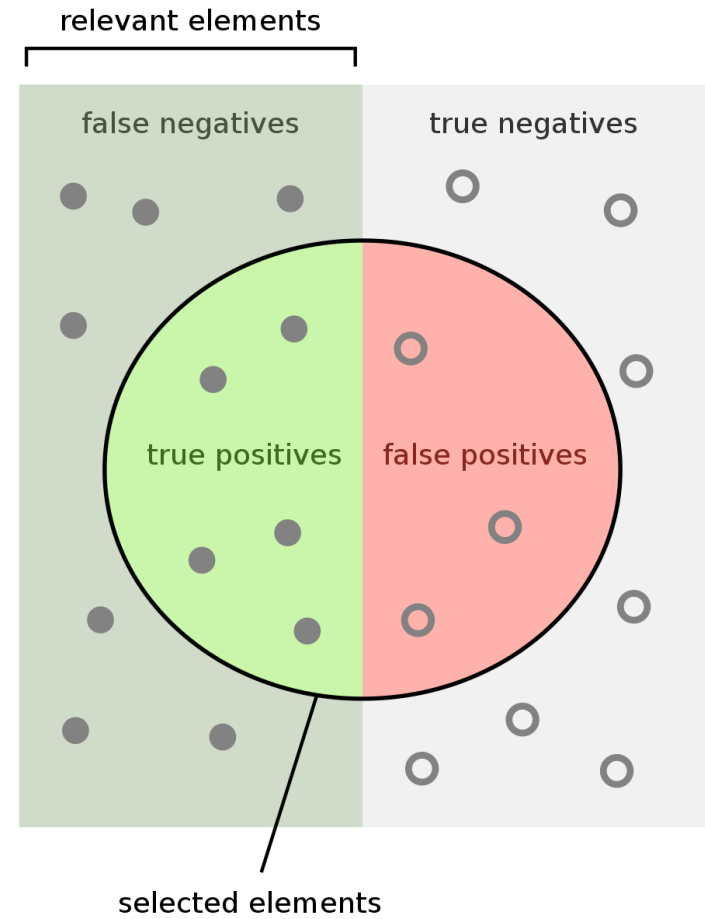| No diabetes | 85% |
| Diabetes | 15% |

# Tools and Techniques

- Logistic regression
- Naïve Bayes
- Decision trees
- Random forest
- Explainable Boosting Machine
- K nearest Neighbours
- SVC
- Ada-Boost
- XG Boost

relevant elements

false negatives

true negatives

true positives

false positives

selected elements

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}}$$

Precision =
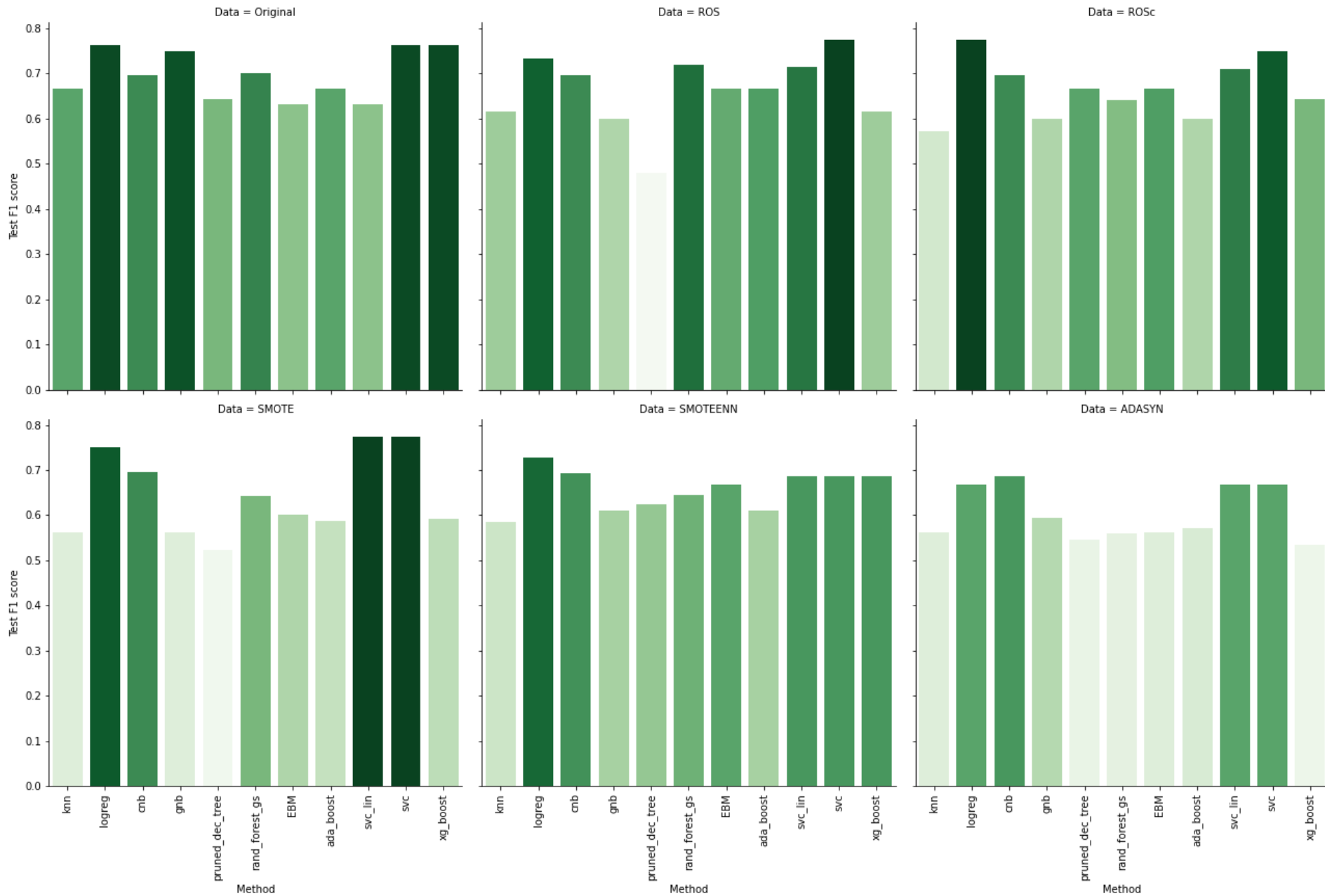
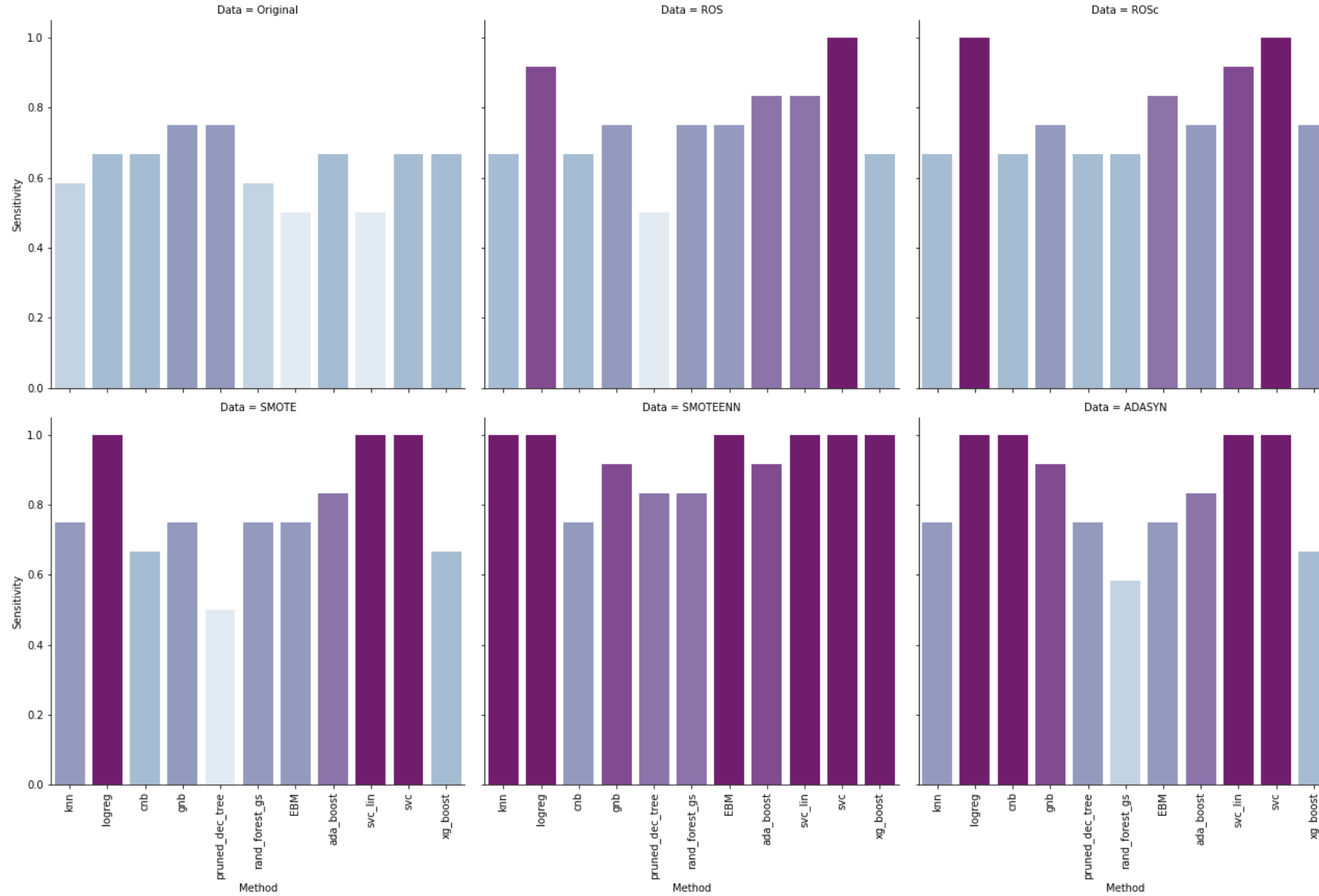Recall = Sensitivity=

Specificity =

# Balancing data

- Naïve random over-sampling

- Naïve random over-sampling with shrinkage
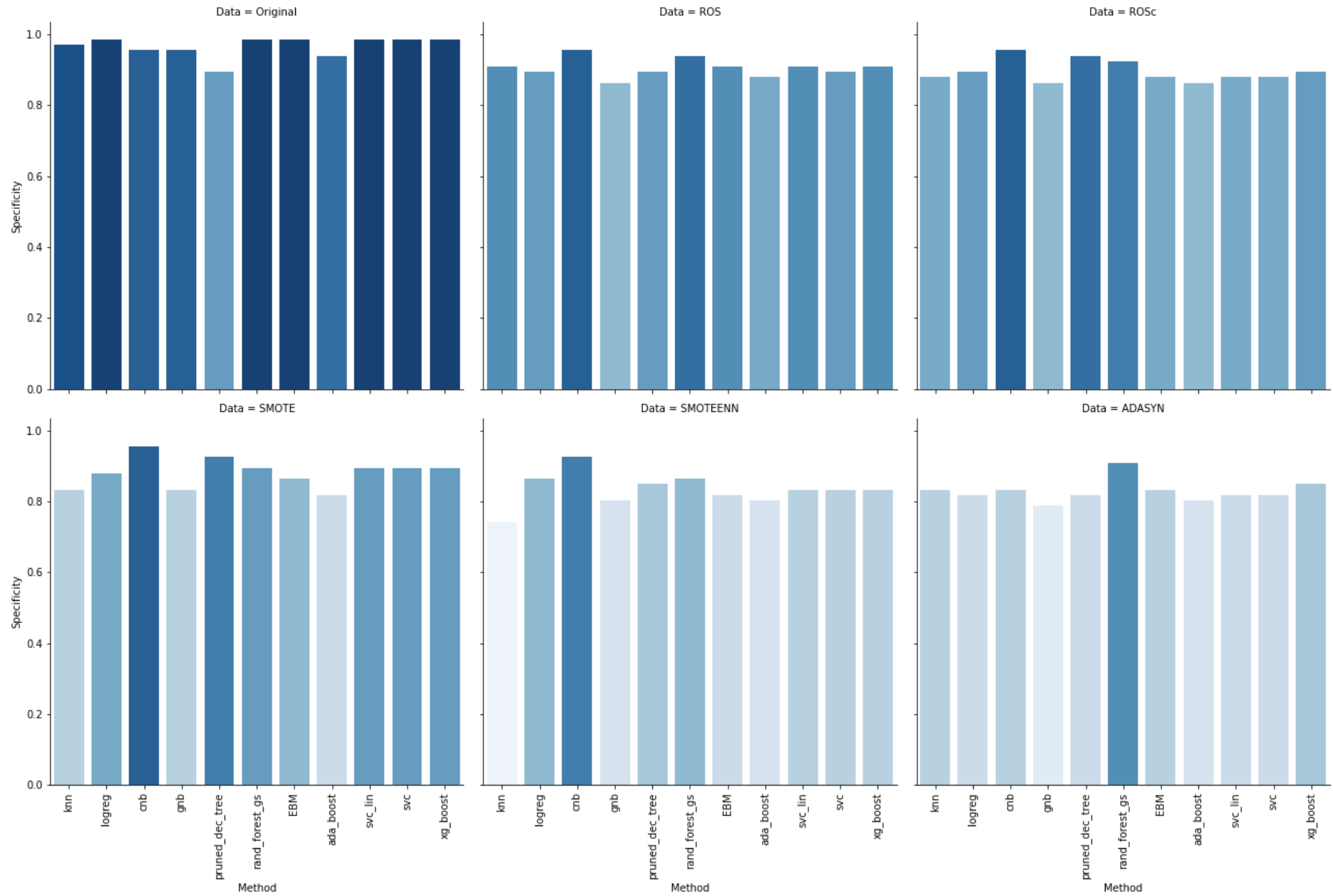
- SMOTE

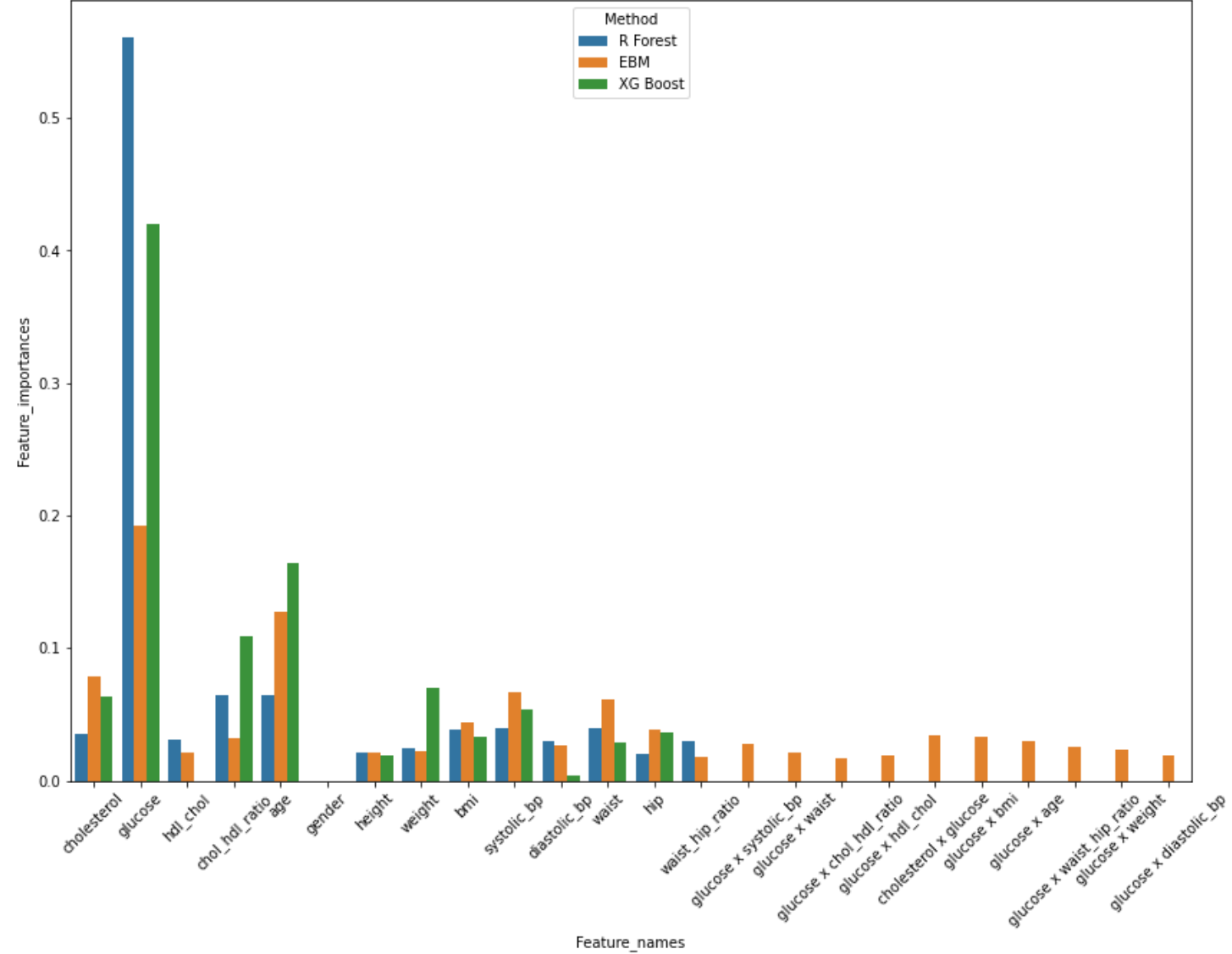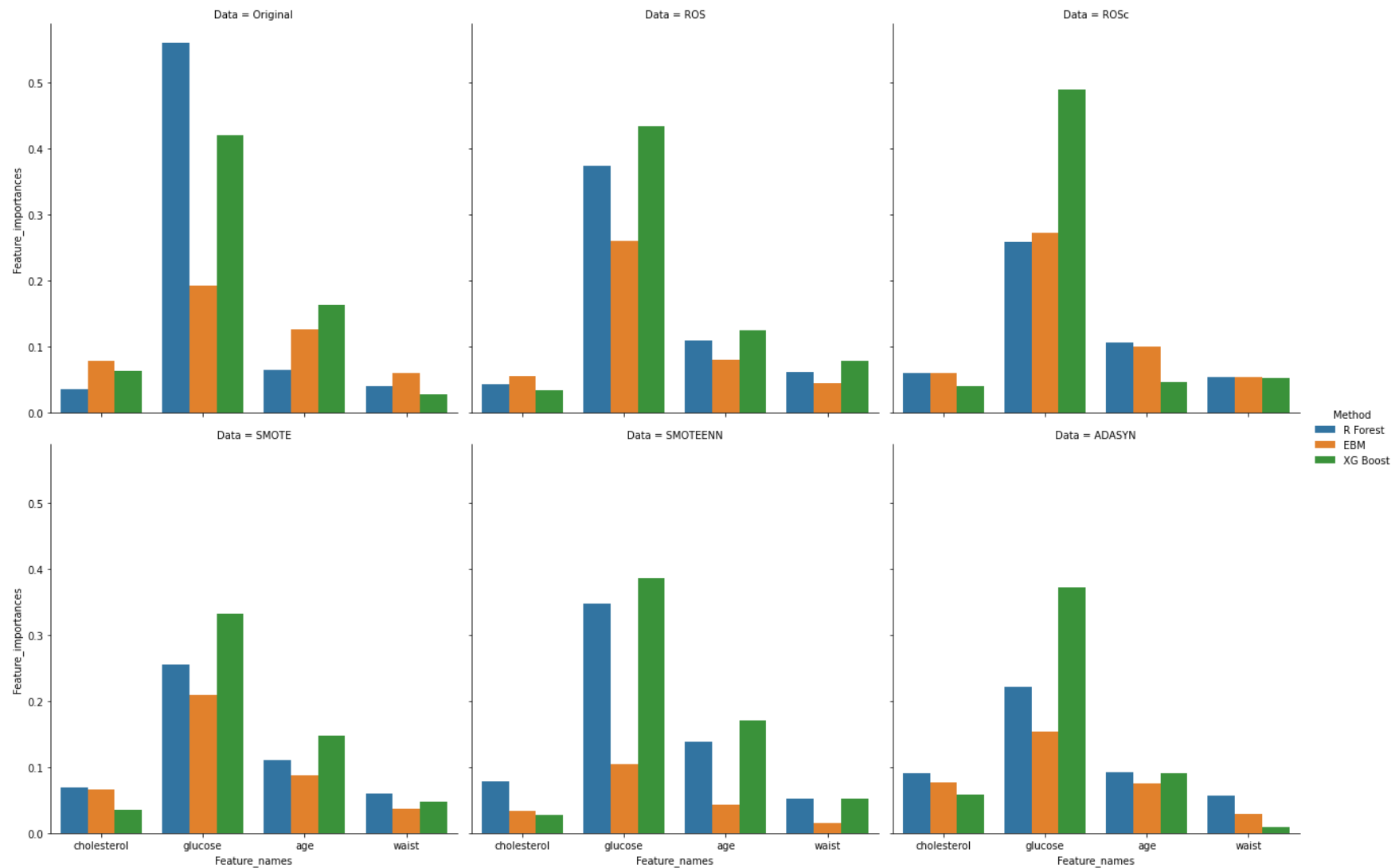- ADASYN

- SMOTEENN

Test F1 score

Sensitivity

Specificity

Feature importances - original data

# Conclusions

**Outcome**:

- logistic regression did surprisingly well
- Balancing data resulted in higher sensitivity and lower specificity for most methods
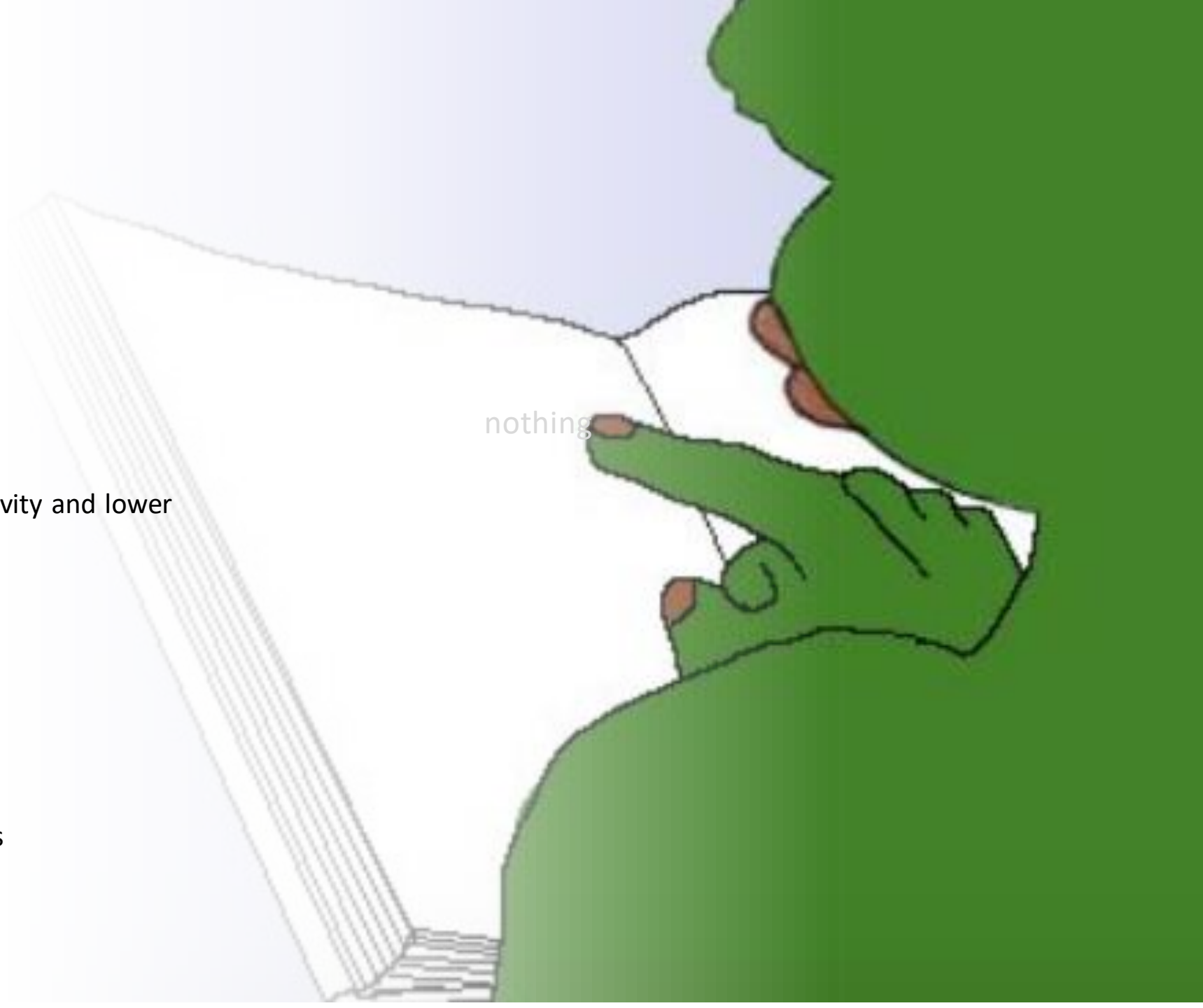
**What was good?**

- The data was quite easy to work with

**What was bad?**

- The data was quite easy to work with

**What have we learned?**

- How to balance data
- How to use scikit-learn built-in methods

nothing

# Report

The report can be found here:
[report](report)