

Fraud Detection



WHY FRAUD DETECTION?

With our custom-built fraud detection model, you will:

- **Prevent revenue loss** by identifying fraudulent consumption early
- **Increase efficiency** by reducing time-consuming manual investigations
- **Ensure fair billing practices** to maintain customer trust and regulatory compliance





OVERVIEW



Part 1

Getting the best
out of the data

Part 2

Finding the best model for
prediction

Part 3

Conclusion &
Evaluation

Part 01

Getting the best out of the data

Merging Data



Client data

- Client ID
- Client Category – Classification of the customer
- Region
- Target (Fraud/Non-Fraud)

135,493

Invoice data

- Tarif Type – pricing category
- Consumption – in kWh (elec.) or m³ (gaz)
- Months – timespan of the invoice
- Reading Remark – annotations for counter readings
- Counter Type – gaz or electricity

4,476,749

Merging Data



Data

- ~~Client ID~~
- Client Category – Classification of the customer
- Region
 - Target (Fraud/Non-Fraud)
- Tarif Type – pricing category
 - Consumption – in kWh (elec.) or m³ (gaz)
 - Months – timespan of the invoice
- Reading Remark – annotations for counter readings
- Counter Type – gaz or electricity

53 Features

135,433

1. Imbalanced Data



1. Oversampling



1. SMOTE*

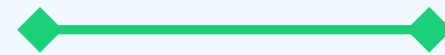


* Synthetic Minority
Oversampling
Technique

Part 02

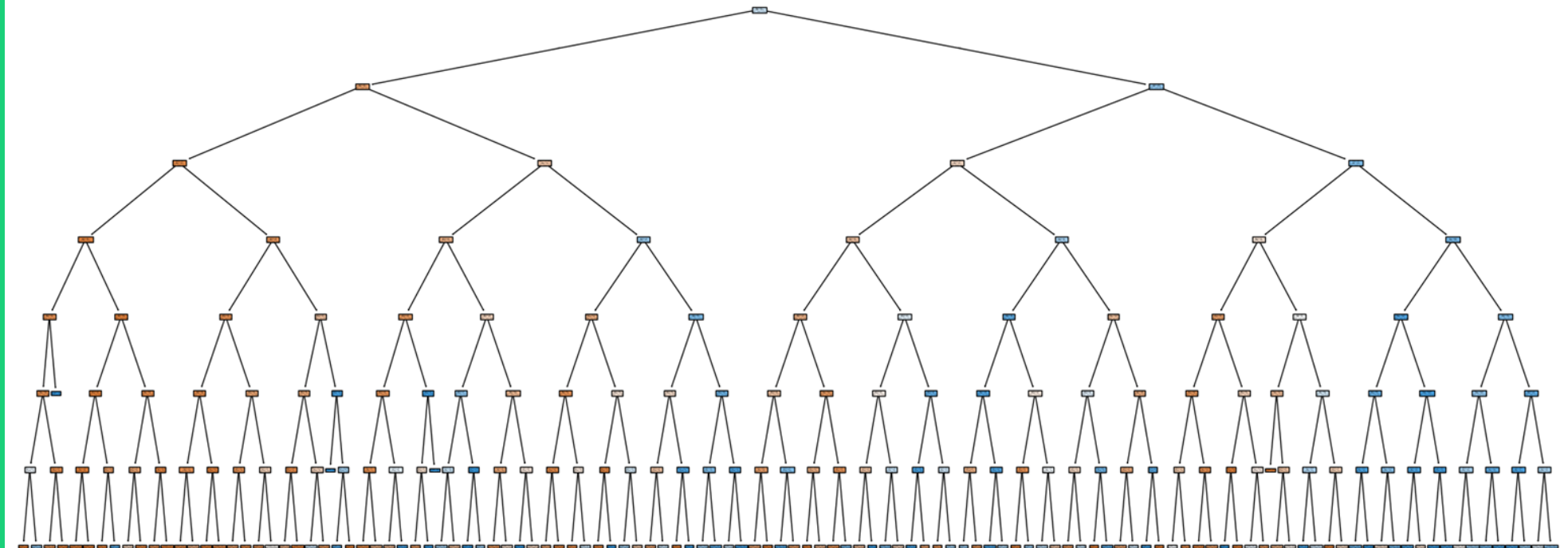
Finding the best model for prediction

Baseline Model



Decision Tree

A simple model with a single Decision Tree and
a maximum depth of 7, oversampling applied



Approaches to get the best model



Random Forest

A model with multiple decision trees that work together by averaging their outcomes

Stacking

A technique that combines multiple different models (base learners) and uses a meta-model to learn from their predictions (DecisionTree, KNN, SGDClassifier)

Boosting

A method that trains models sequentially, where each model corrects the errors of the previous one, gradually improving performance

ROC AUC

Measures how well the model separates fraud from non-fraud
→ used for overall assessment

Recall

The percentage of actual fraud cases correctly identified
→ important as the goal is to identify as much fraud as possible

Accuracy

The proportion of all correct predictions out of total cases.

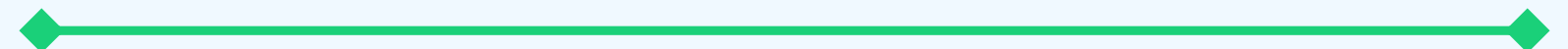
Precision

The percentage of predicted fraud cases that are actually fraud.

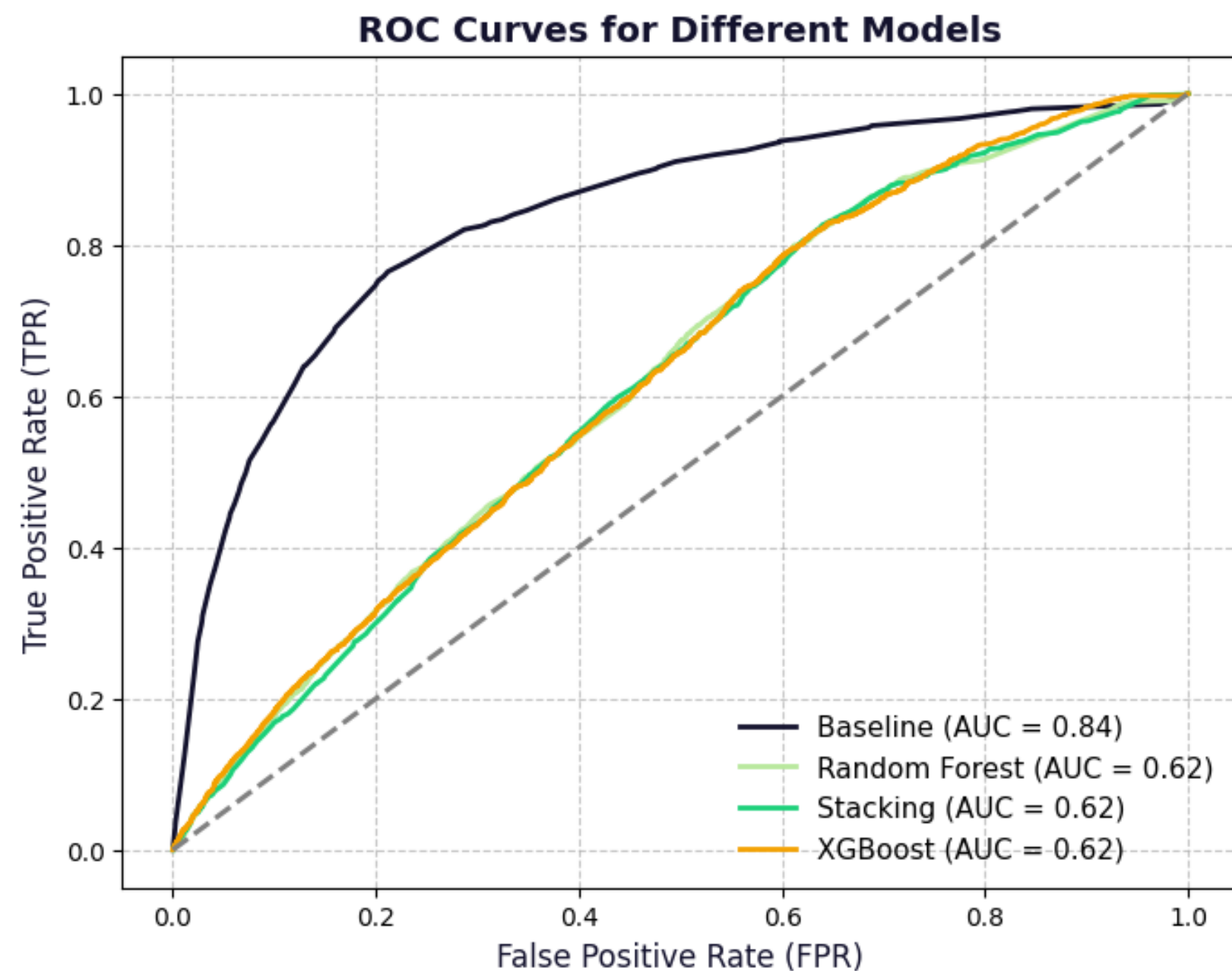
F1 Score

The balance between Precision and Recall.

Main KPIs and their priority for evaluating the models



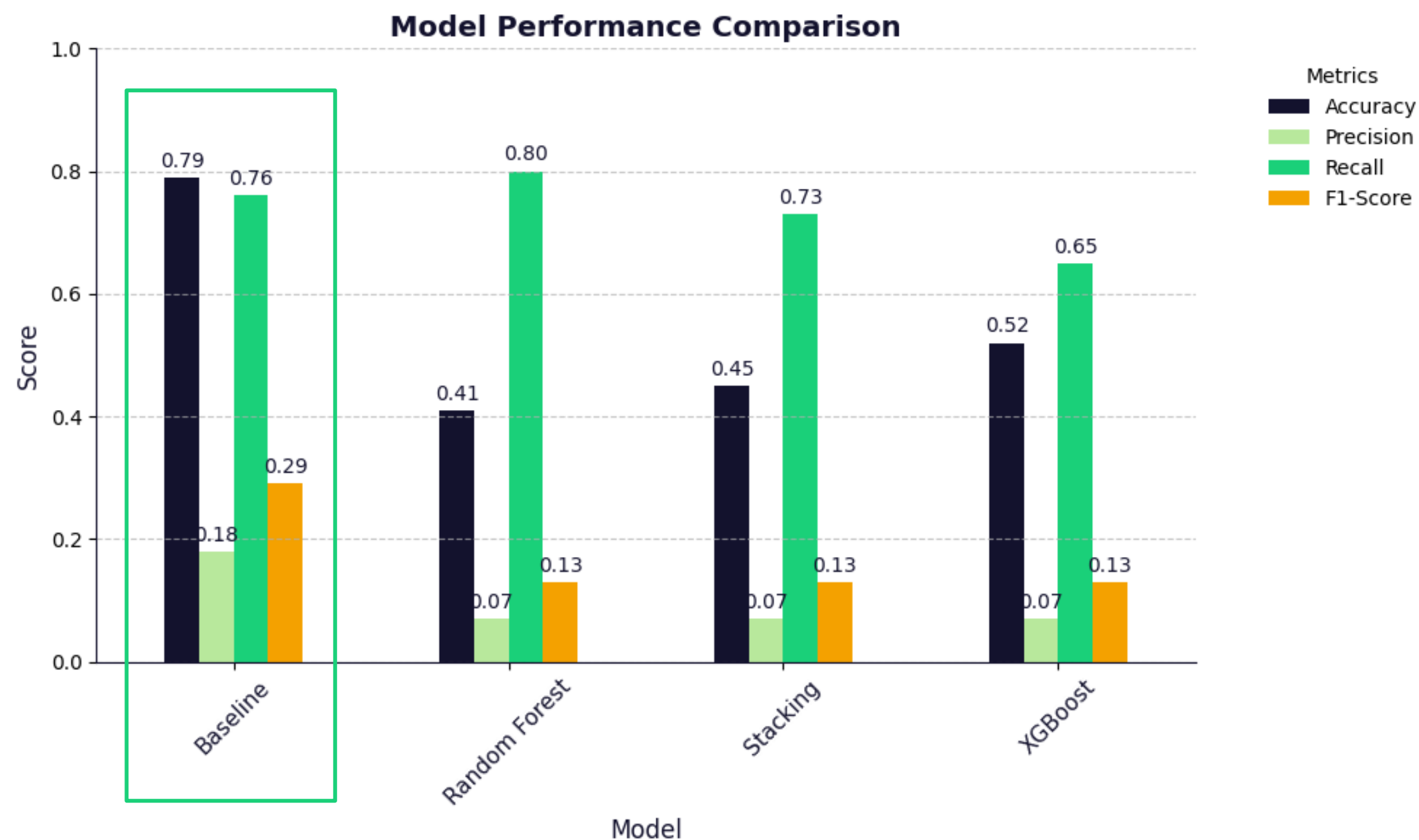
Model Comparison - ROC Curves



- The Baseline model is by far outperforming the complex models
- all complex model achieve the same ROC AUC values of 0.62 while the Baseline model achieves 0.84

Model Comparison - Other KPIs

Also looking at the other KPIs Accuracy, Precision, Recall and the F1 Score, the Baseline Model achieves the overall best results



Decision Tree

A simple model with a single Decision Tree and a maximum depth of 7, oversampling applied



Random Forest

A model with multiple decision trees that work together by averaging their outcomes

Stacking

A technique that combines multiple different models (base learners) and uses a meta-model to learn from their predictions

Boosting

A method that trains models sequentially, where each model corrects the errors of the previous one, gradually improving performance

Part 03

Conclusion & Evaluation

CONCLUSION

- A Simple solution can be the best
- Still need to check and compare with different models
- Importance of EDA and Feature Engineering
- Fraud detection is difficult to detect and compound correlations hard to find
- Working as a team is essential



LOOK FURTHER

- Optimize the base model again
- More Data Engineering
- Read the decision tree as good as possible

