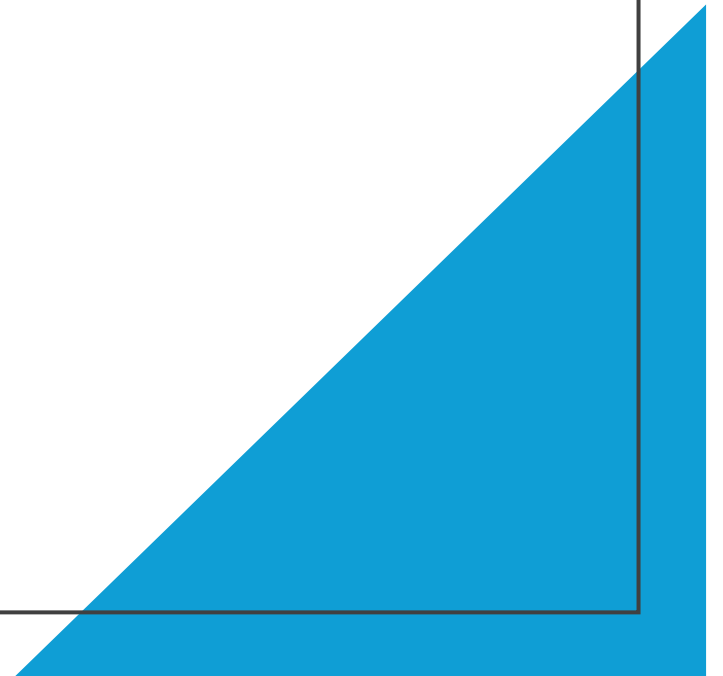


Data Science Advanced

Bank Fraud Classification Project

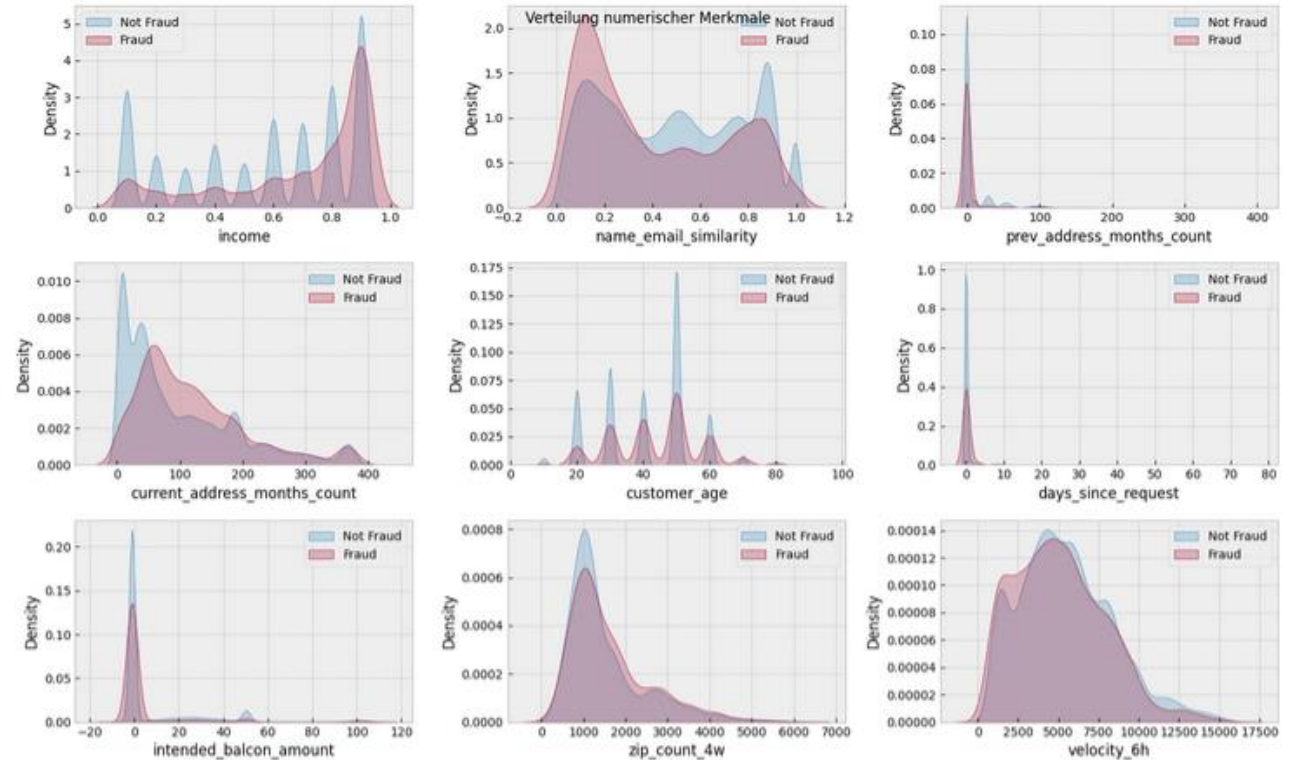
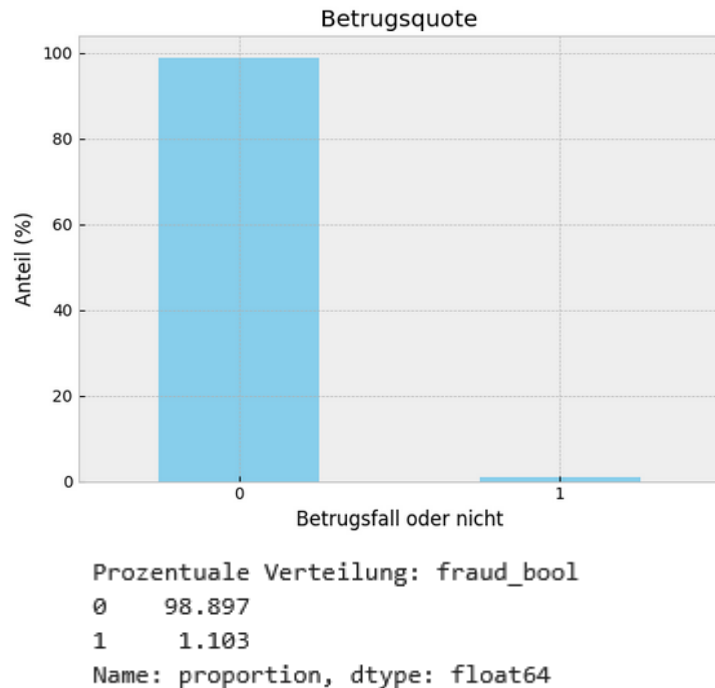
Marc-Vincent Müller

12.12.2025



Data Science Advanced – Bank Fraud

Data set & analysis



Highly unevenly distributed data, increase of fraud with higher income

Data Science Advanced – Bank Fraud

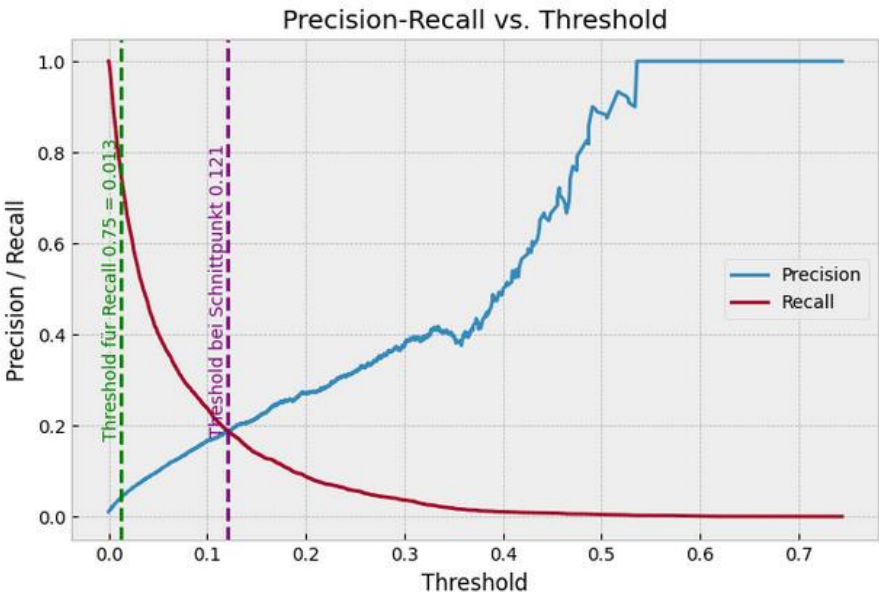
Steps

1. Drop x1 and x2
2. Search missing values
3. Drop columns or refill values with imputation
4. One hot encoding of categorical features
5. Scaling all features
6. Feature selection (only log regression)
7. Define models
8. Data reduction to 10 % and train-test-split
9. Gridsearch (only log regression and decision trees) → Optimization Score: ROC-AUC
10. Fit best model with best hyperparameters on all train data
11. Model evaluation → Precision-Recall-Curve: Threshold optimization on at least 75 % Recall on test data
12. Visualisation with Confusion Matrix
13. Comparison of results on train and test data

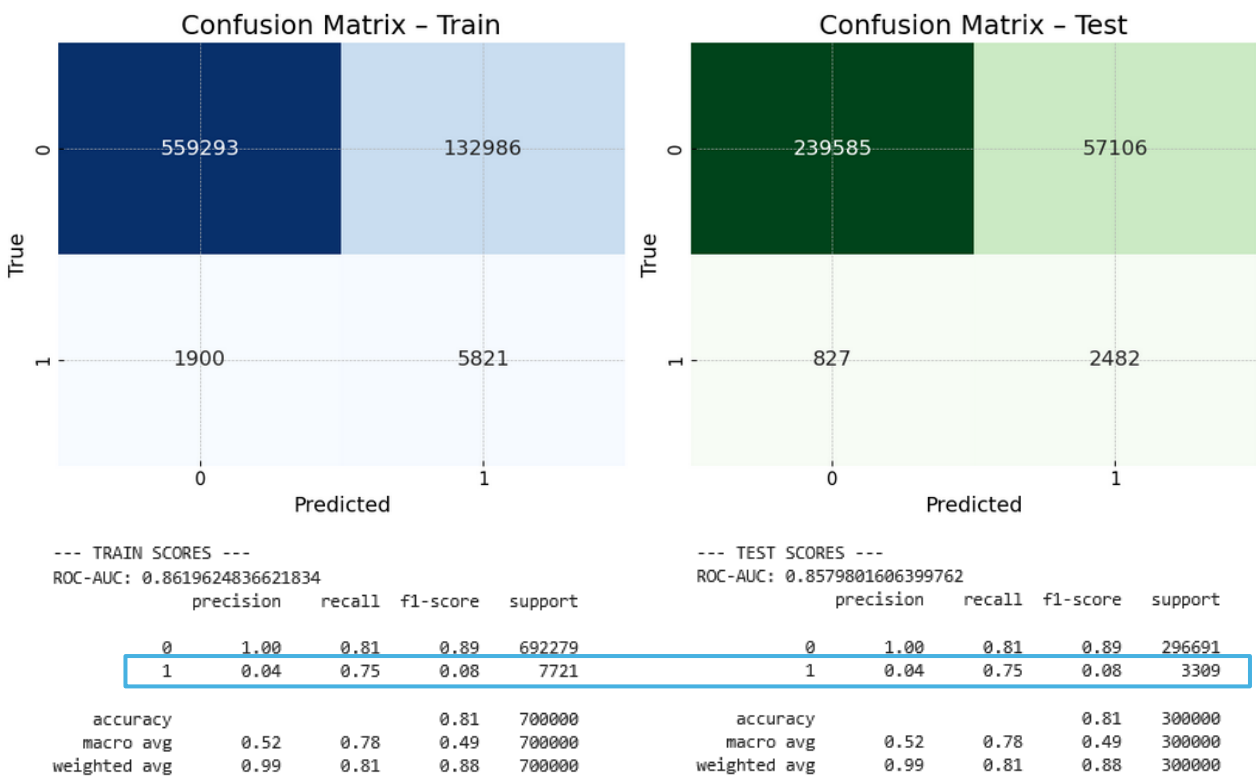
First optimization on ROC-AUC score, then threshold adjustment on at least 75 % Recall

Data Science Advanced – Bank Fraud

Model – Logistic Regression



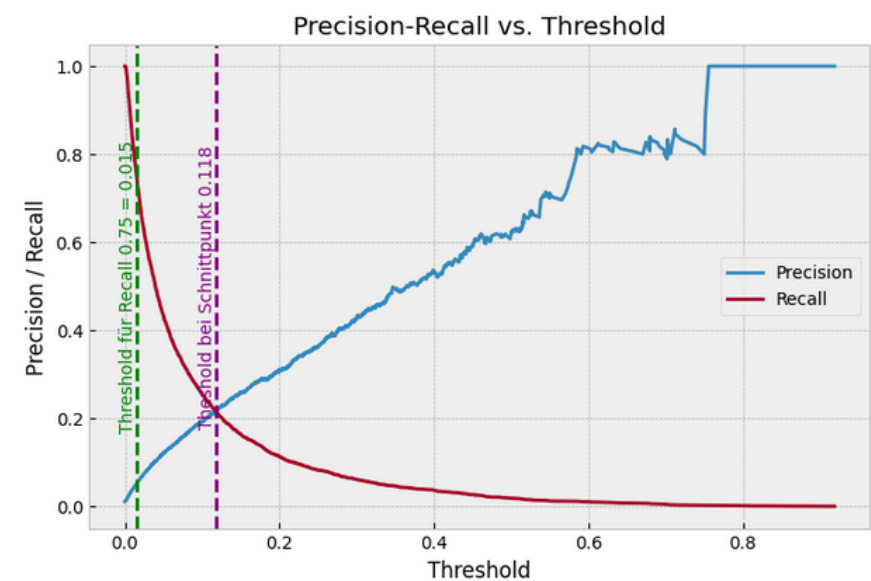
Threshold bei Recall ~75 %: 0.01256743519921646
Threshold bei Schnittpunkt: 0.12129834878162572



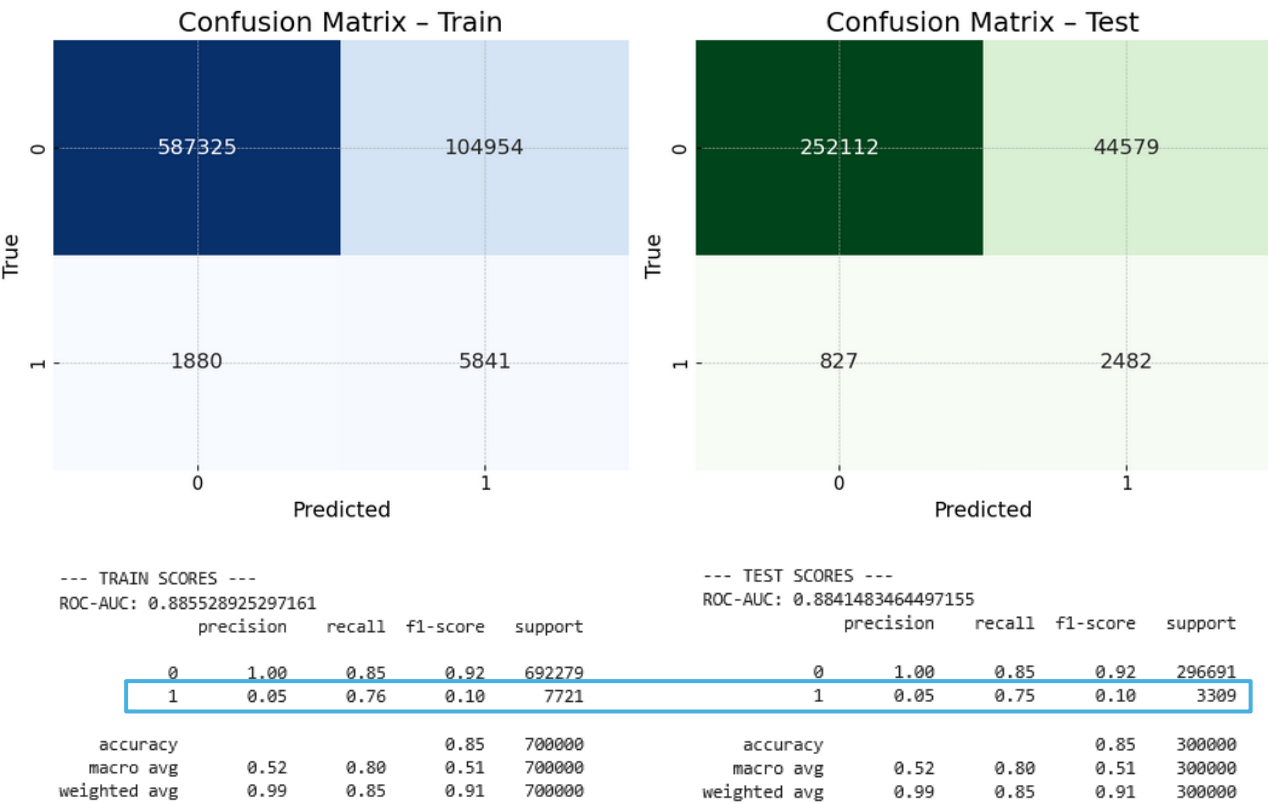
Without SMOTE, ridge regression, class weight: None, C: 3

Data Science Advanced – Bank Fraud

Model – Decision Tree



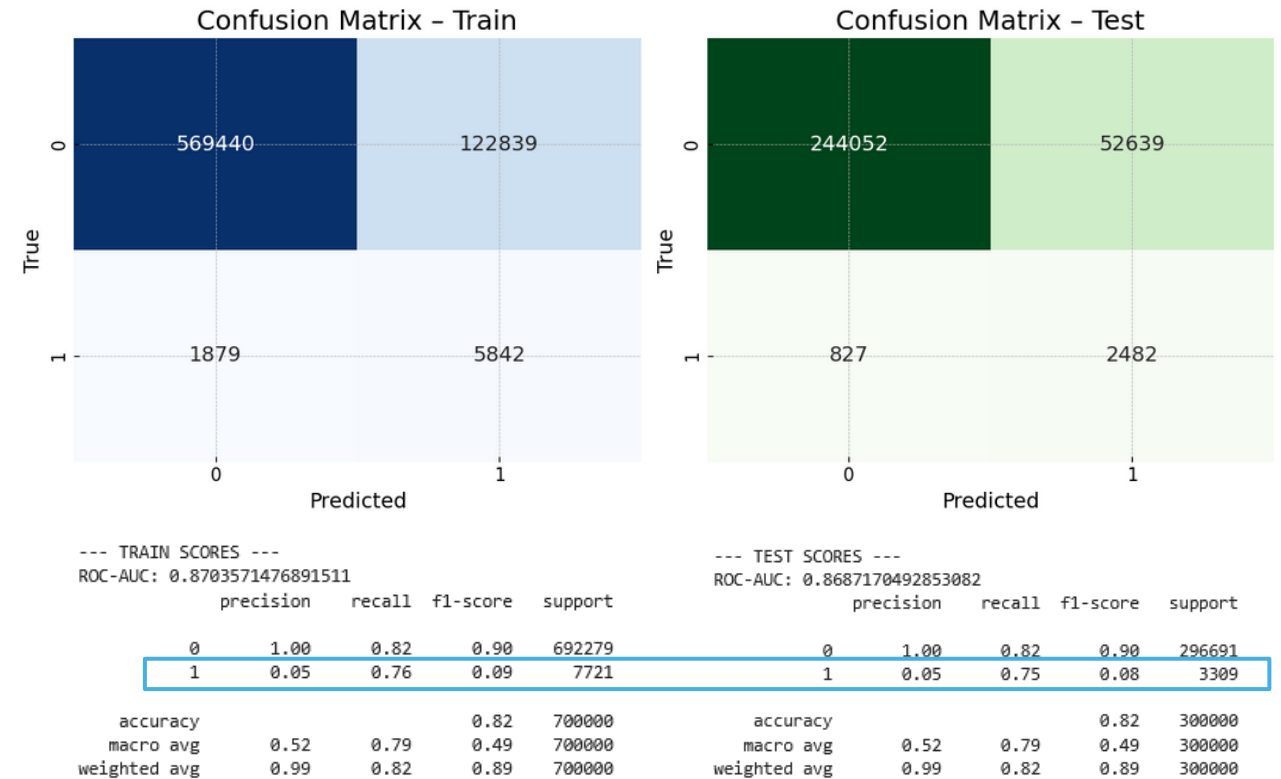
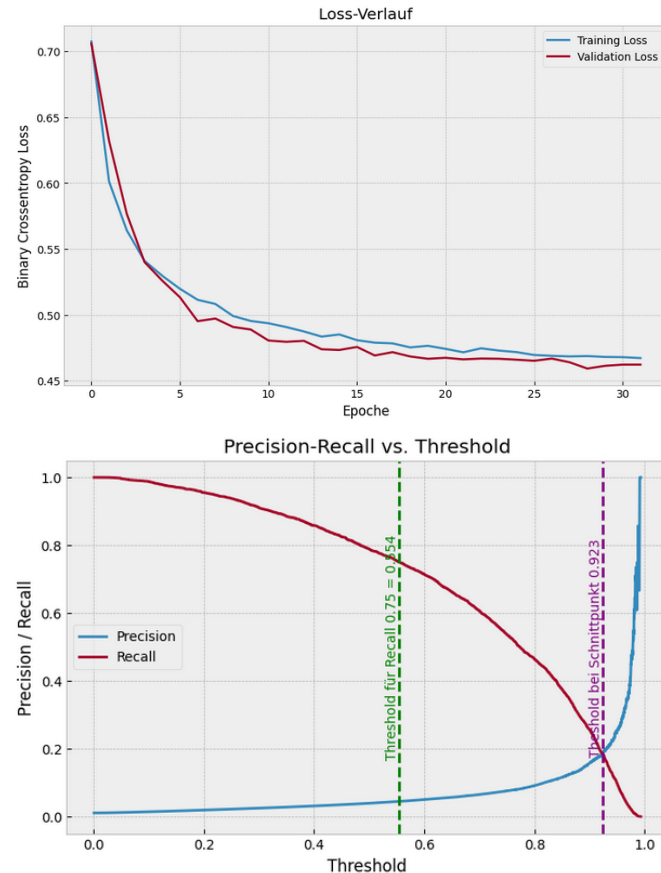
Threshold bei Recall ~75 %: 0.01521433988097543
Threshold bei Schnittpunkt: 0.11760161799233325



Best tree: LightGBM, max depth: 3, n estimators: 300, num leaves: 3, scale pos weight: 1, learning rate: 1

Data Science Advanced – Bank Fraud

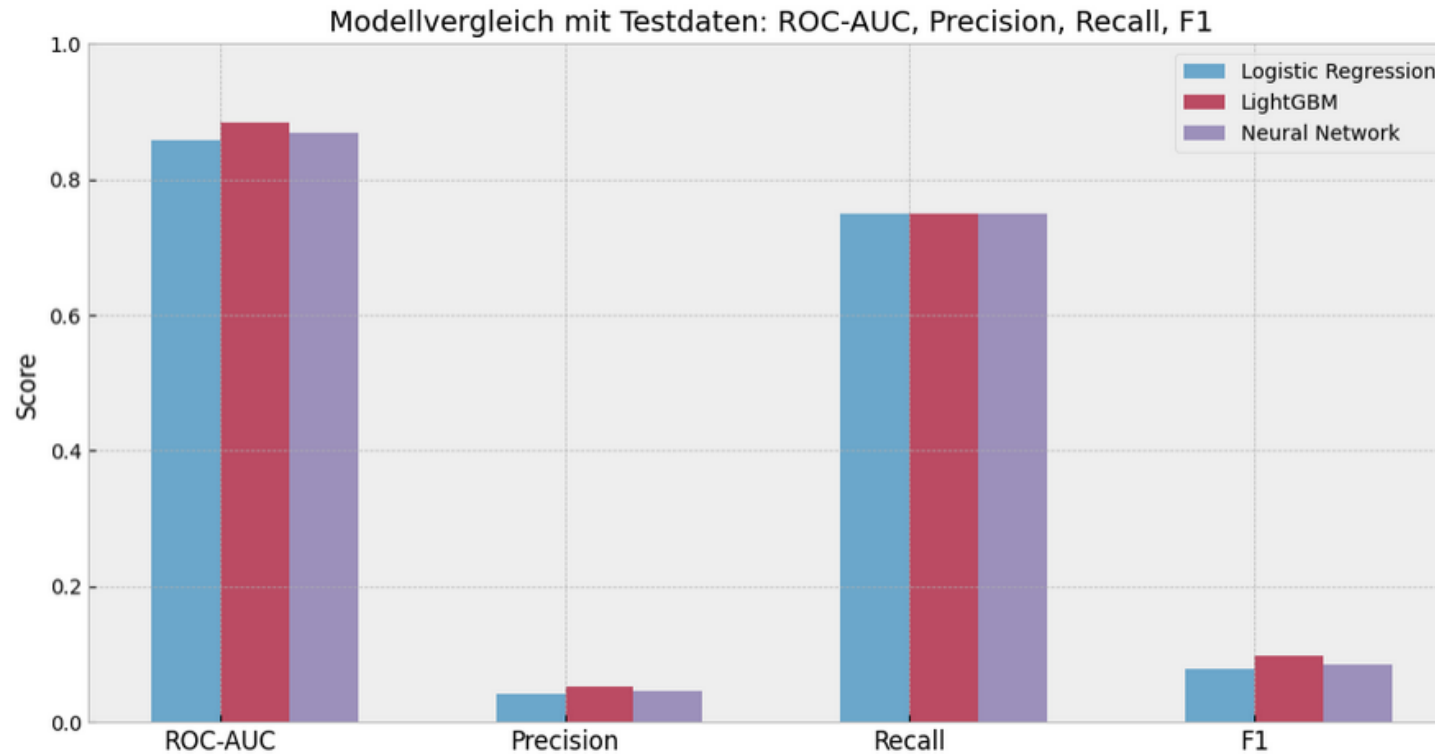
Model – Neural Network



2 intermediate layers, ReLU activation , regularisation: Dropout against overfitting, early stopping

Data Science Advanced – Bank Fraud

Model – Comparison and summary



Very similar scores regardless of the model, best results with LightGBM, difficult fraud detection due to uneven distribution, Additional features required to improve model (for example x1 and x2)