



Final Project Report

Capstone2: Bank Marketing Campaign



Kevin C. (Springboard)

11/15/2025

Executive Summary

This project analyzes a direct marketing campaign conducted by a Portuguese banking institution. The data contains client information and details of past marketing contacts, with the goal of predicting whether a client will subscribe to a term deposit (y = yes/no). A variety of classification models were evaluated to identify the most effective approach for improving marketing performance and conversion rates.

Problem Statement

The objective of this project is to support a Portuguese bank in improving the efficiency of its direct marketing efforts by developing a model that predicts whether a client will subscribe to a term deposit (y = yes/no). The dataset includes demographic characteristics, interaction history, and economic indicators. A key challenge is that the target variable is highly imbalanced, as only a small proportion of clients choose to subscribe.

Addressing this imbalance and constructing a reliable classification model will enable the bank to better prioritize high-potential clients, optimize resource allocation, and ultimately increase campaign effectiveness.

Data Overview

This dataset was obtained from the [UCI Machine Learning Repository - Bank Marketing Dataset](#). It contains information collected from a Portuguese banking institution related to direct marketing campaigns.

The dataset includes:

- 1) Demographic data (age, job, marital status, etc),
- 2). Contact data (communication type, date, duration).
- 3). Campaign performance indicators (# of contacts, previous outcomes)
- 4). Economic indicators (employment variation rate, consumer confidence, Euribor, etc)

Methodology

We will clean and preprocess the dataset, including handling categorical variables and scaling numerical features. Next, we will conduct exploratory data analysis (EDA) to visualize patterns in both the dependent and independent variables and assess whether the binary target is balanced or imbalanced. If class imbalance is present, we will apply technique such as SMOTE or class weighting. We will then build classification models, including logistic regression, decision trees, and ensemble methods – and evaluate their performance using precision, recall, F1-score, and the confusion matrix. Finally, we will identify the key features that influence term deposit subscription.

Results & Evaluation

After preprocessing, exploratory data analysis, four models were trained and evaluated to address the imbalanced classification problem:

- XGBoost
- Random Forest
- Logistic Regression
- LightGBM

A consistent modeling pipeline was applied to all models using 5-fold Stratified Cross-Validation:

- Splitting the dataset into an 80/20 train test ratio using stratification
- Applying Stratified K-fold (K=5) only on the training set
- Using SMOTE within each training fold to oversample the minority class
- Scaling numerical features with StandardScaler
- Training an XGBoost classifier within each fold
- Evaluating on the validation fold using F1-score, Recall, and ROC AUC

Accuracy is not reliable for imbalanced classification.

Recall and F1 were prioritized because the business objective is to identify as many potential term-deposit subscribers as possible (minimizing false negatives).

Cross-Validation Performance Summary

Model	Avg F1	Avg Recall	Avg ROC AUC
LightBGM	0.6367	0.7928	0.9422

Random Forest	0.6247	0.6875	0.9411
XGBoost	0.6111	0.6479	0.9415
Logistic Regression	0.5821	0.7398	0.9163

Model Comparison

- LightBGM achieved the strongest overall performance, with the highest F1-score and the highest recall among all models.
- Logistic Regression produced strong recall but lower F1 and AUC
- Random Forest and XGBoost performed well but did not match LightBGM's balance across all metrics.
- ROC AUC values were high for all tree-based models, indicating consistent ranking ability, but LightBGM performed best overall.

Selected Final Model

Given its superior balance of recall, F1-score, and AUC, LightBGM was selected as the final model and retrained on the full training dataset (with SMOTE applied). The final model was then evaluated on the test set to validate generalization performance.

Retraining on the Full Training Dataset

The final LightBGM model was retrained on the entire training set (with SMOTE applied) and then evaluated on the held-out test set. This full data retraining ensures that the final deployed model has access to all available training information.

Final Test Set Results

- Precision: $TP / (TP + FP) = 802 / (802 + 828) \approx \mathbf{0.492}$
- Recall: $TP / (TP + FN) = 802 / (802 + 95) \approx \mathbf{0.8941}$
- F1-score: 0.6347
- ROC AUC: 0.9480

The confusion matrix demonstrated a significant reduction in false negatives, showing that the model successfully captures more clients who are likely to subscribe.

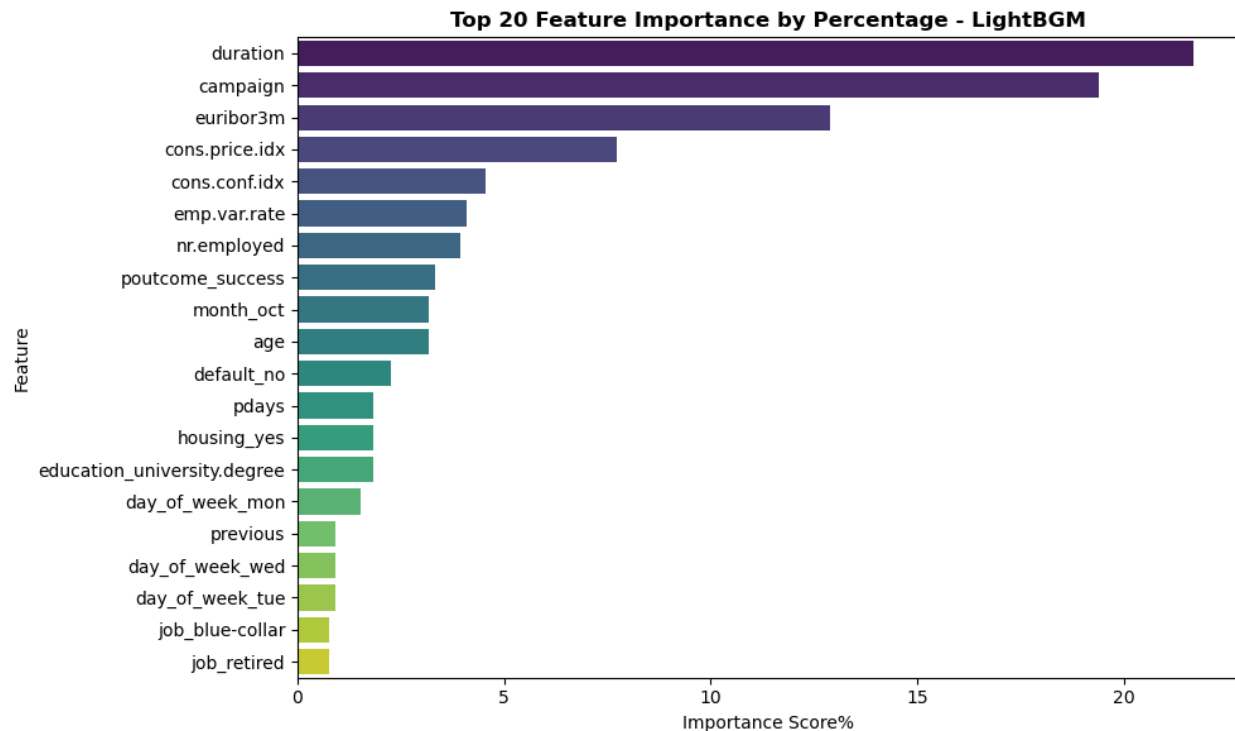
Feature Importance (LightGBM)

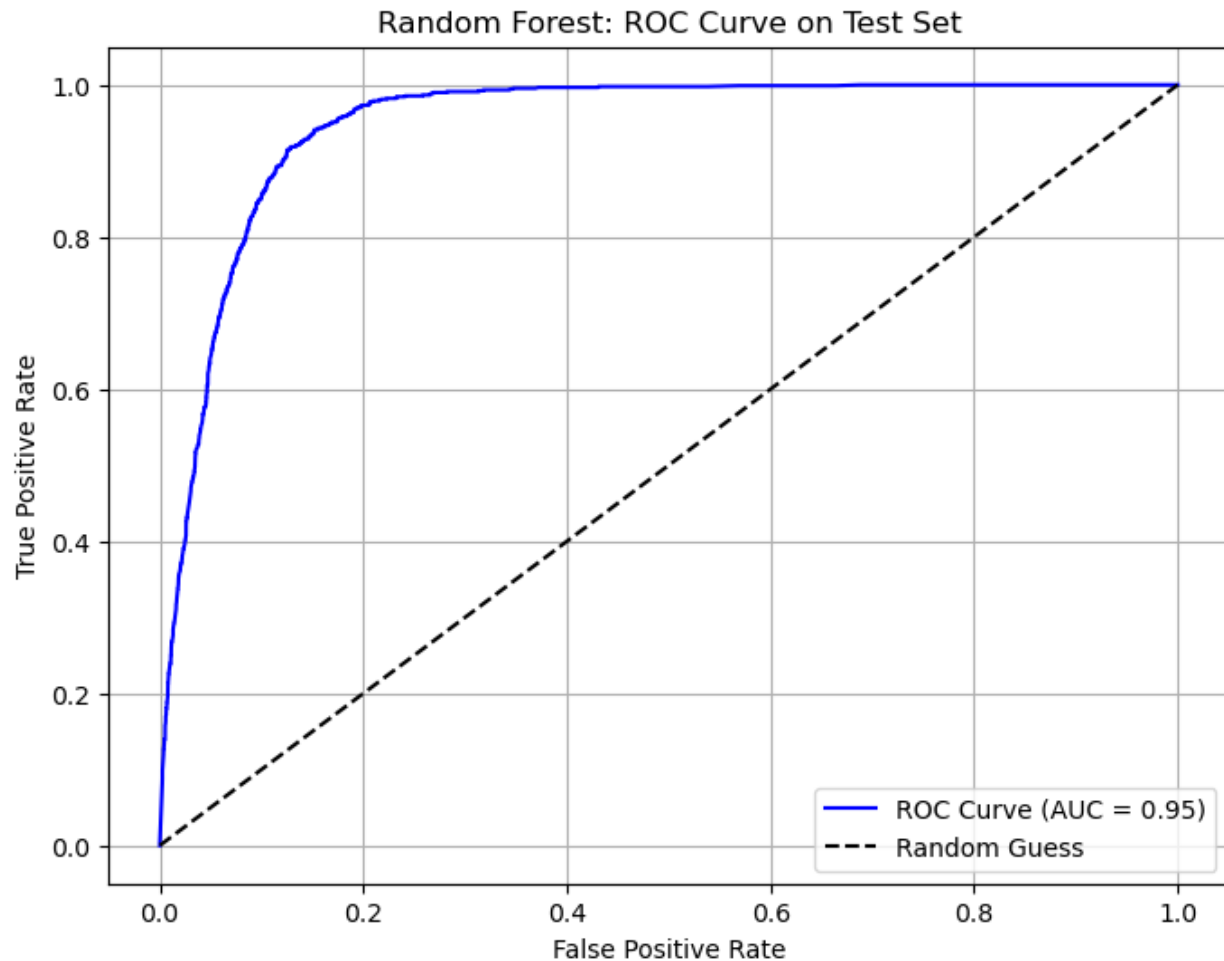
The most influential predictors included:

- Duration
- Previous campaign outcome
- Number of contacts

- Employment variation rate
- Euribor 3-month rate

These features align strongly with marketing intuition: longer, more successful past interactions and favorable economic conditions correlate with a higher likelihood of subscription.





Recommendations

Based on the findings from the model and feature importance analysis, several actionable recommendations can help the bank improve the effectiveness of its future marketing campaigns:

1. Prioritize clients with higher predicted subscription probability
Use the LightGBM model to score customers before launching campaigns.
Focus marketing resources (calls, follow-ups, personalized messages) on high-probability segment to increase conversion rates and reduce unnecessary outreach.
2. Improve call quality and engagement duration
The feature duration is the strongest predictor of subscription.
This suggests that higher quality and longer conversations rather than more frequent calls lead to better outcomes.
The bank should train agents to:

- Improve conversation quality
- Ask more engaging questions
- Offer personalized product explanations

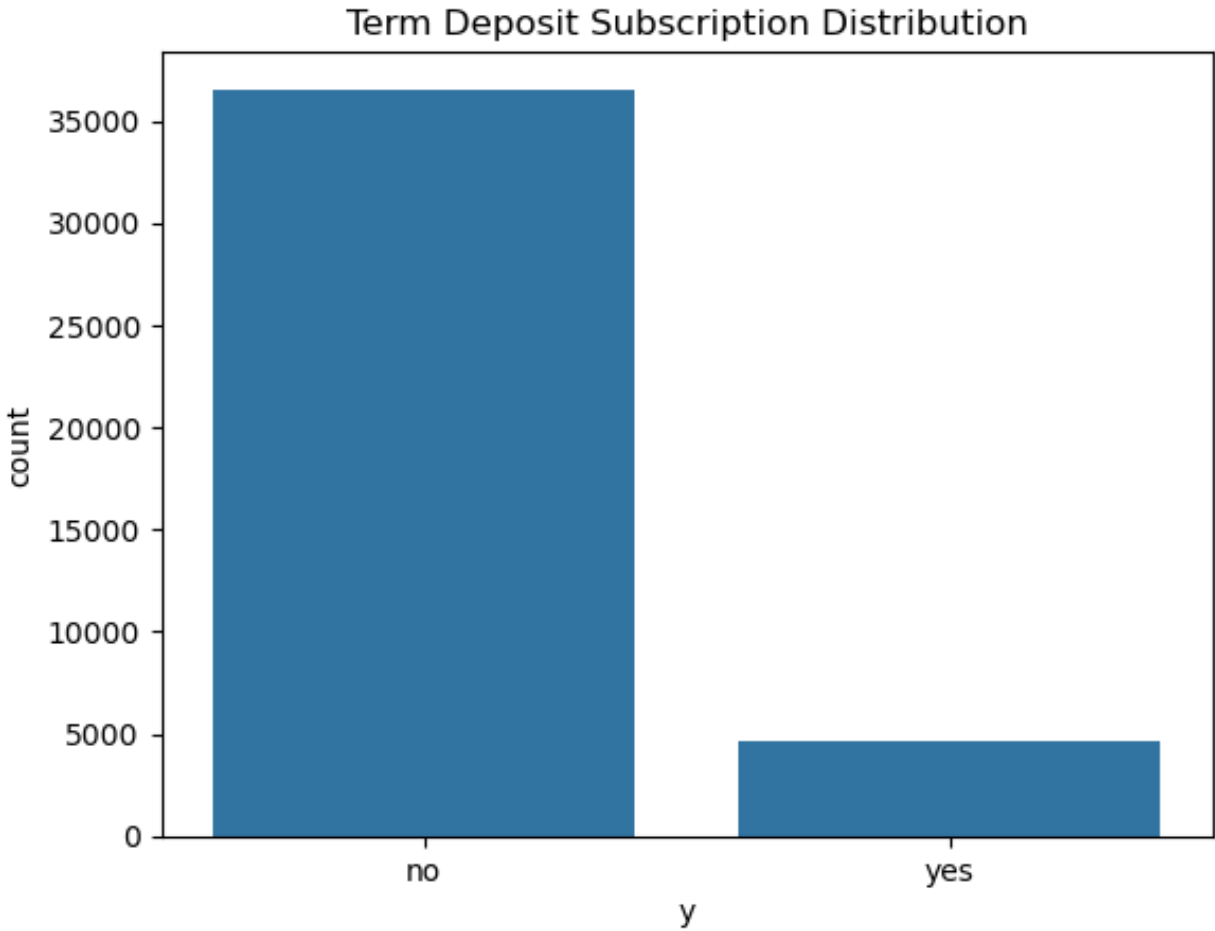
This creates more meaningful client interactions.

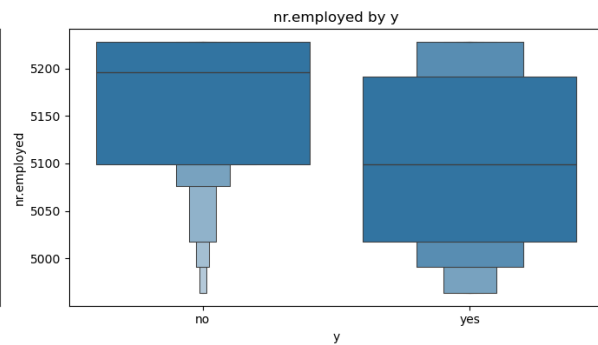
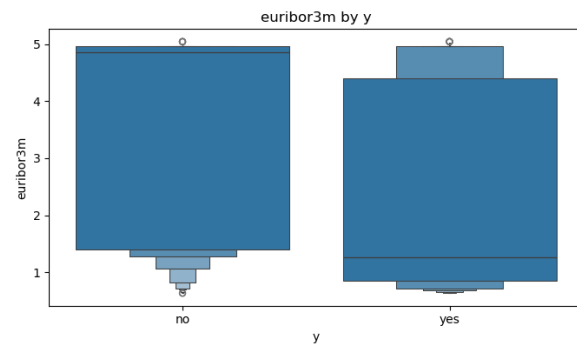
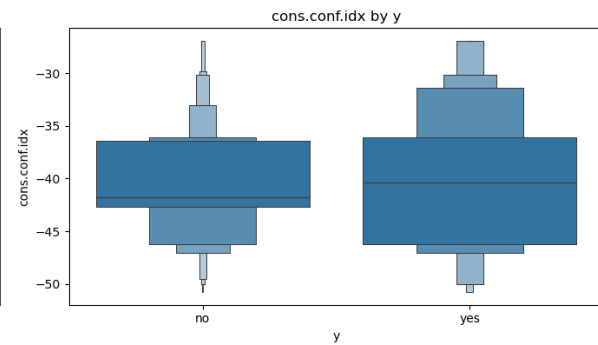
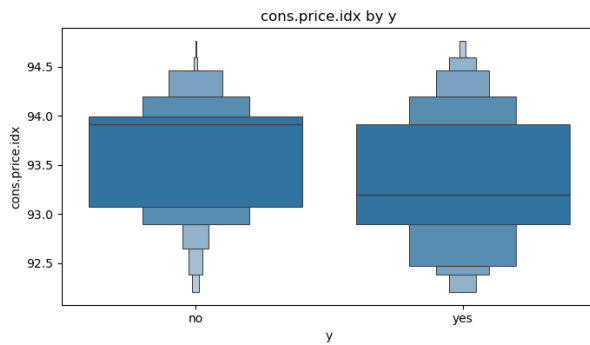
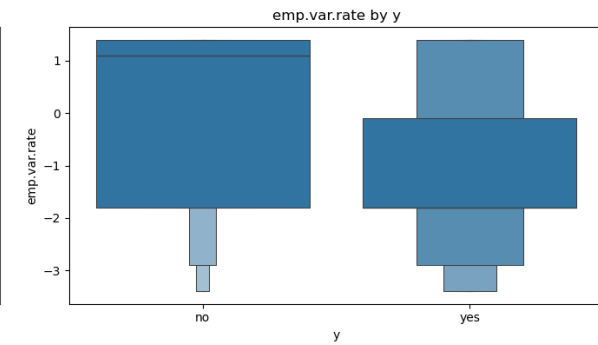
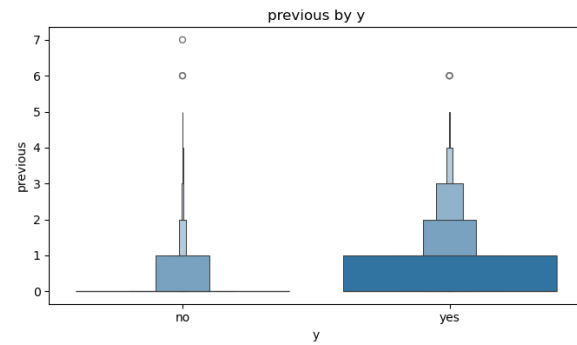
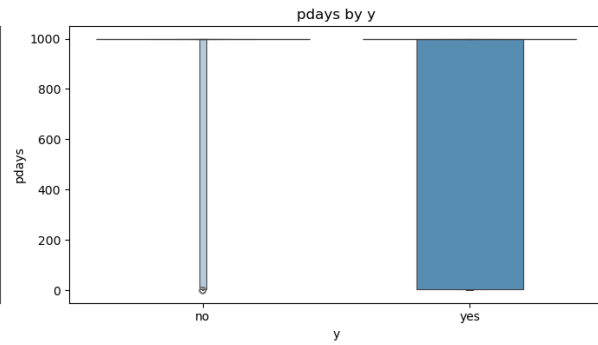
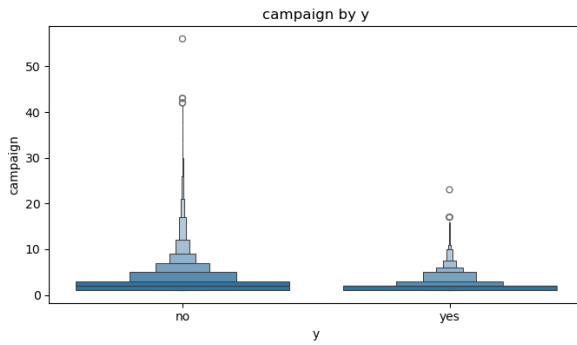
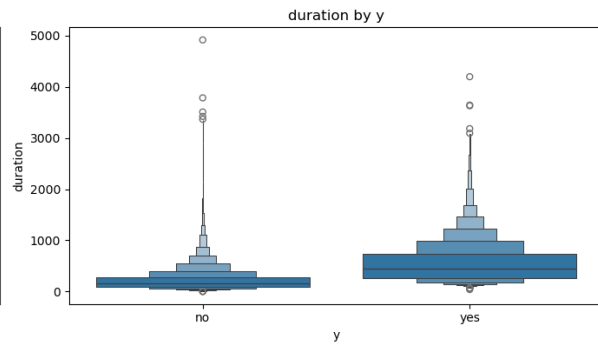
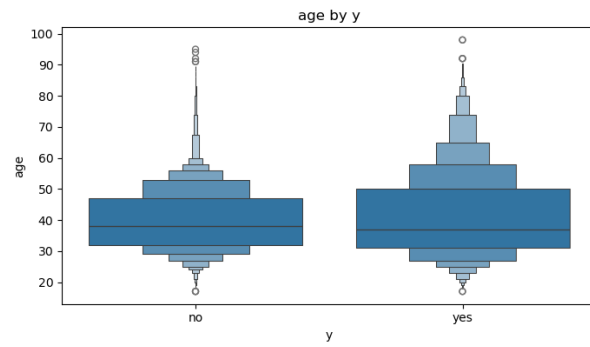
Future Work

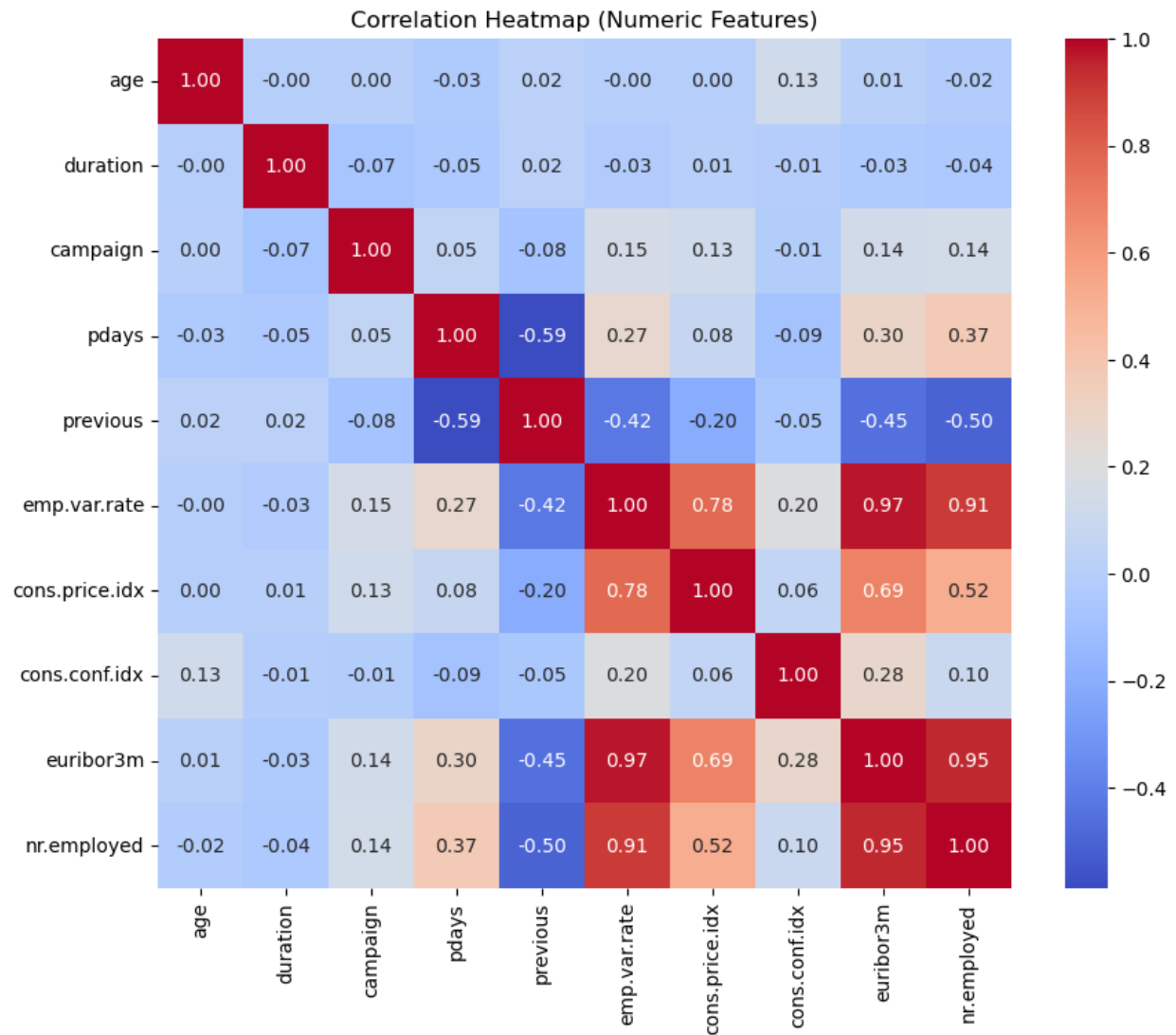
Hyperparameter tuning for improved model performance

The current LightGBM model uses a basic configuration. Applying grid search, random search, or Bayesian optimization could further enhance Recall, F1-score, and model stability.

Appendix







XGBoost Classification Average Static Results

Avg F1 Score: 0.6111
Avg Recall Score: 0.6479
Avg ROC AUC: 0.9415

Random Forest - Average Static Results

Avg F1 Score: 0.6247
Avg Recall Score: 0.6875
Avg ROC AUC: 0.9411

Logistic Regression - Average Static Results

Avg F1 Score: 0.5821
Avg Recall Score: 0.7398
Avg ROC AUC: 0.9163

LightBGM - Average Static Results

Avg F1 Score: 0.6367
Avg Recall Score: 0.7928
Avg ROC AUC: 0.9422

	Avg F1 Score	Avg Recall Score	Avg ROC AUC
Model			
LightGBM	0.6367	0.7928	0.9422
Logistic Regression	0.5821	0.7398	0.9163
Random Forest	0.6247	0.6875	0.9411
XGBoost	0.6111	0.6479	0.9415