

LOAN DEFAULT PREDICTION

. . .

Tatiana Meshkova



Overview of the problem

Content



Approach for the solution



Key findings & insights



Evaluation & choosing models



Recommendations & next steps



The potential benefits





Overview of the problem

Loan default is a significant problem for the lending industry, resulting in financial losses and other negative consequences



Accurately predicting loan defaults can help banks minimize risk, allocate resources more effectively, comply with regulations, and promote fair lending practices



Developing an unbiased and data-driven loan approval process can enhance customer service quality and provide a competitive advantage

The goal is to build a classification model to predict loan defaults, considering important features and ensuring reliable performance and interpretability

Approach for the solution

Exploratory Data Analysis:

. . .

Analyzing the data to understand relationships and patterns between features and the target variable.

Model Evaluation:

Assessing the model's performance using evaluation metrics such as accuracy, precision, recall, and F1-score.

3

Model Selection and Training:

Comparing different machine learning models and selecting the best performing one, and fine-tuning the hyperparameters to improve its performance. Choosing the optimal model for deployment and providing recommendations

4

•••

Key findings & insights

LGB model Feature Importance



The most important features for predicting mortgage default are debt-to-income ratio (DEBTINC), age of oldest credit line in months (CLAGE), number of credit lines (CLNO), amount of mortgage owed (MORTDUE) and loan amount (LOAN),.

All TOP5 features were consistently identified across multiple models and the correlation matrix, indicating their strong impact on loan default prediction.

Key findings & insights

Analysis confirms the importance of delinquencies and derogatory reports as predictors of loan default.



The average loan amount for customers who default on their loans is lower than for customers who pay back their loans.



Older credit lines are more reliable and may be a useful factor for lenders to consider when assessing loan applications.





Evaluation models

- Recall was chosen to evaluate the model's performance, as it tells us how many actual positive cases the model identified correctly, which is important in minimizing false negatives.
- Precision and F1-score are also considered to compare similar models, with Precision indicating the number of correctly identified positive cases and F1-score providing a balance between Precision and Recall.
- Considering all three metrics allows for a more informed decision on the best model for predicting loan default.
- In addition, the model's training time and ease of implementation were also taken into account when selecting the best model for deployment.

•••



Choosing model

- Based on the analysis and evaluation of several models, I propose to adopt the **LGBM Classifier** as the final solution.
- This model outperformed the other models in terms of metric scores, as well as having a reasonable training time. It also provides feature importance ranking, which can help to identify the most important factors affecting the target variable.



8



Next step - model implementation



Collect data on borrowers, including both numerical and categorical variables, and split the data into training and testing sets.



Train a predictive model using the LGBM Classifier on the training set, using the identified important predictors of loan default.



Integrate the model into the bank's loan approval process, using the model's predictions as one of many factors to consider when making lending decisions.



Set a threshold for the model's predicted probability of default and only approve loans for borrowers whose predicted probability of default falls below that threshold.



Recommendations

Implement the LGBM Classifier into the loan approval process to more accurately assess the risk of loan default and take appropriate measures to mitigate that risk.



Continue monitoring the model's performance in production and updating it as needed to ensure that it remains accurate and effective in identifying borrowers at risk of loan default.



Establish a process for updating the model as needed, taking into account changes in economic conditions, regulatory requirements, and borrower profiles.



Communicate the use of the model to customers to promote transparency and build trust, providing clear explanations of how the model works and how it informs lending decisions.



The potential benefits

Increase revenue: Our model can accurately predict loan defaults, allowing us to avoid risky borrowers and focus on profitable ones.



Mitigate risk: By avoiding loans that are likely to default, we can reduce the risk of financial losses and improve the overall health of our loan portfolio.



Better customer experience:

With a faster and more accurate loan approval process, customers will have a better experience and are more likely to return for future loans.



Scalability: The solution can be easily scaled to handle an increasing volume of loan applications without requiring significant additional resources or manual labor.



. . .

. . .

12

· • •

Correlation matrix



Test performance comparison

Model	Recall	Precision	Accuracy	Total Rank
Tuned Catboost classifier oversampling without				
one-hot encoding	0.916319	0.916319	0.916319	4
lgbm classifier	0.910015	0.910015	0.910015	6
Lgbm classifier without one-hot encoding	0.909665	0.909665	0.909665	12
Tuned Catboost classifier oversampling	0.912819	0.912819	0.912819	13
Tuned Catboost classifier	0.908265	0.908265	0.908265	16
Tuned lgbm classifier oversampling	0.909668	0.909668	0.909668	20
Tuned lgbm classifier	0.907566	0.907566	0.907566	22
Catboost classifier	0.904063	0.904063	0.904063	24
XGBoost classifier	0.905114	0.905114	0.905114	25
Decision Tree classifier	0.907589	0.907589	0.907589	32
Tuned Catboost classifier undersampling	0.908996	0.908996	0.908996	33
Random Forest classifier	0.896361	0.896361	0.896361	33
Tuned Random Forest classifier	0.872200	0.872200	0.872200	37
Tuned lgbm classifier undersampling	0.902369	0.902369	0.902369	42
Tuned Random Forest classifier 2	0.861016	0.861016	0.861016	43
Tuned Decision Tree classifier	0.883830	0.883830	0.883830	46



Test performance comparison

Model	Time	Rank
Lgbm classifier without one-hot encoding	0.35	1
Lgbm classifier	0.39	2
Tuned Catboost classifier oversampling without one-hot encoding	01.01	3
Tuned Catboost classifier oversampling	1.32	4



Set a threshold 20%

	yl_test	yl_pred_proba_lgbm	approved
4394	0	0.002379	1
5000	0	0.018923	1
2786	0	0.00344	1
2256	0	0.005746	1
114	0	0.307778	0
3787	0	0.000594	1
1289	0	0.00525	1
1189	1	0.997221	0
4715	0	0.002677	1
70	1	0.894823	0
1268	1	0.989448	0

