



# **Predicting Bank Loan Eligibility Using Machine Learning Models and Comparison Analysis**

Jehan Nirmal - COHNDDSF23.1f-015

# Contents

Chapter 1: Introduction .....	4
1.1 Background .....	4
1.2 Research Problem .....	4
1.3 Objectives of the Project .....	4
1.4 Research Questions .....	4
1.6 Justification of the Research .....	5
1.7 Limitations .....	5
Chapter 2: Literature Review .....	6
Chapter 3: Data Preparation Process - Data Preprocessing and Data Wrangling .....	7
3.1 Dimensionality reduction .....	7
3.2 Data Cleaning .....	8
Chapter 4: Methodology .....	10
4.1 Model Selection and Pipeline Construction .....	10
4.1.1 Logistic Regression .....	10
4.1.2 XGBoost .....	10
4.2 Cross-Validation .....	10
4.3 Hyperparameter Tuning .....	10
4.4 Class Imbalance Handling .....	10
4.5 Evaluation Metrics .....	11
Chapter 5: Data Analysis, Visualization, Modeling and Interpretation .....	12
5.1 Clustering .....	12
5.2 Encoding and Normalizing the Data .....	13
5.3 Splitting the Data Set .....	14
5.4 Model Building and Evaluation .....	14
5.4.1 Evaluation oof Logistic Regression Model Classifier Performance .....	14
5.4.2 Evaluation of XGBoost Classifier Performance .....	15
5.5 Model Evaluation and Selection through Cross-Validation .....	16
5.6 Hyperparameter Tuning and Imbalanced Data Handling .....	17
5.7 Evaluation and Fitting of Final Model .....	18

Chapter 6: Discussion of Conclusions.....	20
Chapter 7: References .....	21

# **Chapter 1: Introduction**

## **1.1 Background**

In contemporary banking practices, the efficient allocation of loans is paramount for sustaining economic growth and promoting financial inclusion. However, traditional methods of assessing loan eligibility often lack accuracy and efficiency, leading to suboptimal decision-making processes. With the emergence of machine learning (ML) and predictive analytics, there exists an opportunity to revolutionize the loan approval process by harnessing data-driven insights.

## **1.2 Research Problem**

This research addresses the challenge of developing a robust loan prediction system using ML models to enhance the accuracy and efficiency of loan approval processes in banking institutions. By analyzing diverse datasets encompassing borrower profiles and loan outcomes, this study aims to identify predictive variables and construct models capable of reliably assessing an applicant's creditworthiness.

## **1.3 Objectives of the Project**

The primary objective of this project is to explore the efficacy of ML techniques in predicting bank loan eligibility and to develop a comprehensive understanding of the factors influencing loan approval decisions. Specific objectives include:

- Implementing ML algorithms to analyze loan data and predict borrower eligibility.
- Evaluating the performance of different ML models in loan prediction.
- Investigating the impact of various factors such as income, credit history, and demographics on loan approval outcomes.

## **1.4 Research Questions**

To guide the research process, the following questions will be addressed:

1. How effective are ML techniques in predicting bank loan eligibility?
2. What are the key factors influencing loan approval decisions?
3. Which ML model demonstrates the highest predictive accuracy in loan prediction?

## **1.5 Scope of the Research**

This research focuses on the development and evaluation of ML models for loan prediction using a specific dataset related to bank loan applications. The scope encompasses data preprocessing, model training, evaluation, and interpretation of results. The study does not delve into broader economic or regulatory factors affecting the banking industry.

### **1.6 Justification of the Research**

The implementation of ML-based loan prediction systems has the potential to streamline loan approval processes, reduce risk exposure, and enhance financial inclusion. By accurately assessing borrower creditworthiness, banks can make informed decisions, leading to more efficient allocation of resources and improved customer satisfaction. This research contributes to the advancement of data-driven decision-making in the banking sector, with implications for risk management and lending practices.

### **1.7 Limitations**

Some limitations of this research include:

- Dependency on the quality and completeness of the dataset.
- The generalizability of findings may be limited to the specific context and dataset used.
- Ethical considerations regarding data privacy and algorithmic bias will be addressed to the best extent possible.

## Chapter 2: Literature Review

In recent years, the automation of the bank loan approval process has become a focal point for both financial institutions and prospective borrowers. Mamun, Farjana, and Mamun (2023) present a comprehensive exploration into the predictive capabilities of ML algorithms in "Predicting bank loan eligibility using machine learning models and comparison analysis." By leveraging a dataset containing crucial customer attributes such as age, income type, loan annuity, and credit history, the study employs an array of ML methodologies, including Random Forest, XGBoost, and Logistic Regression. The comparative analysis reveals Logistic Regression as the most effective model, achieving an impressive accuracy of 92% and a remarkable F1-Score of 96%. This study underscores the significance of leveraging historical data to predict loan eligibility accurately and highlights the superiority of Logistic Regression in this domain.

Yang (2024) extends this discourse in "Research on loan approval and credit risk based on the comparison of Machine learning models" by investigating six candidate ML models for loan approval and credit risk assessment. The study emphasizes the importance of model interpretability alongside predictive accuracy. By evaluating models such as Logistic Regression, Decision Tree, and Neural Network, among others, Yang emphasizes the practical implications of ML theories in facilitating informed decision-making for both borrowers and lenders. Furthermore, the examination of confusion matrices and feature importance aids in understanding model performance comprehensively, thereby enhancing the decision-making process.

Additionally, Udhbav et al. (2022) contribute to this dialogue with their study titled "Prediction of Home Loan Status Eligibility using Machine Learning." While focusing on home loan eligibility, their findings align with the overarching theme of leveraging ML for predictive modeling in the banking sector. By identifying hidden patterns within customer datasets, the study underscores the transformative potential of modern technology in revolutionizing traditional loan approval processes. Through performance evaluation metrics such as accuracy, precision, and F1-Score, Udhbav et al. advocate for the adoption of ML-driven approaches to streamline loan eligibility assessments, thereby saving time and improving efficiency.

Collectively, these studies highlight the pivotal role of ML algorithms in predicting bank loan eligibility and mitigating credit risk. By leveraging historical data and employing sophisticated ML techniques, financial institutions can make informed lending decisions while borrowers can enhance their chances of securing loans. Moreover, the emphasis on model interpretability and comprehensive performance evaluation underscores the importance of transparency and reliability in ML-driven loan approval systems, ultimately fostering trust and efficiency in the banking sector.

## **Chapter 3: Data Preparation Process - Data Preprocessing and Data Wrangling**

The dataset used in this research project was obtained from Kaggle, specifically from the "Home Credit Default Risk" competition hosted by Anna Montoya, Kirill Odintsov, and Martin Kotek in 2018. The dataset consists of 307,511 entries and contains 122 columns.

Before data cleaning and preprocessing, the dataset is structured as a pandas DataFrame with the following characteristics:

- RangeIndex: The index ranges from 0 to 307,510, indicating the number of rows in the dataset.
- Columns: There are a total of 122 columns, each representing a different variable or feature. These columns include identifiers such as SK\_ID\_CURR and various attributes related to loan applicants, such as income, employment status, credit history, and loan amounts.
- Data Types: The dataset contains a mix of data types, including float64 (65 columns), int64 (41 columns), and object (16 columns).
- Memory Usage: The dataset consumes approximately 286.2+ megabytes of memory.

The dataset provides a comprehensive snapshot of loan applicants' profiles and associated features, offering valuable insights into factors influencing credit risk and loan default probabilities.

### **3.1 Dimensionality reduction**

In the process of data cleaning, several steps were undertaken to ensure the dataset's quality and suitability for analysis. This section outlines the various methods employed to preprocess the data effectively.

Firstly, columns containing terms such as 'MODE', 'MEDI', and 'AVG' were identified and subsequently dropped from the dataset. These columns were deemed redundant as they were found to be intercorrelated, providing little to no valuable insight for the analysis. By eliminating these columns, the dataset was streamlined, focusing only on the most relevant features.

Additionally, columns associated with external sources of information, including 'EXT\_SOURCE\_1', 'EXT\_SOURCE\_2', and 'EXT\_SOURCE\_3,' were dropped from the dataset.

These columns were excluded due to the unavailability of information, which could potentially compromise the accuracy and reliability of the analysis.

Next, the correlation between each numerical feature and the target variable ('TARGET') was assessed. Features exhibiting a correlation coefficient above a predefined threshold (0.04) or below its negative equivalent were deemed significant and retained for further analysis. This step ensured that only features with a substantial impact on the target variable were considered in subsequent modeling tasks.

Similarly, an analysis of variance (ANOVA) test was conducted for each categorical column to determine its significance in predicting the target variable. Columns with a p-value below the threshold of 0.05 were considered statistically significant and included in the final dataset. This approach helped identify categorical features that significantly influenced the loan default risk, enabling a more accurate predictive modeling process.

Upon completion of these steps, the dataset was reduced to include only the most relevant numerical and categorical columns. This subset of features, termed 'significant\_data,' formed the basis for subsequent data analysis, visualization, and modeling tasks. By selectively retaining features with the greatest predictive power, the cleaned dataset facilitated more effective decision-making and enhanced the overall reliability of the research findings.

The data cleaning process played a crucial role in preparing the dataset for analysis, ensuring that it was free from inconsistencies, redundancies, and irrelevant information. By systematically addressing data quality issues, researchers could focus on extracting meaningful insights and building robust predictive models to inform loan approval decisions effectively.

The reduction in dataset columns from 122 to 24 highlights the importance of feature selection and dimensionality reduction in machine learning projects. By focusing on a subset of informative features, researchers can improve model performance, reduce computational complexity, and enhance the interpretability of results.

The final composition of the 'significant\_data' dataset reflects a balanced selection of numerical and categorical features, each contributing to the overall predictive power of the models. With this streamlined dataset, researchers can proceed to the next stages of analysis, including data visualization, modeling, and interpretation, with greater confidence and efficiency.

### **3.2 Data Cleaning**

Further data cleaning procedures were undertaken to address specific aspects of the dataset, ensuring its integrity and completeness for analysis. This section outlines the additional steps implemented to refine the dataset.



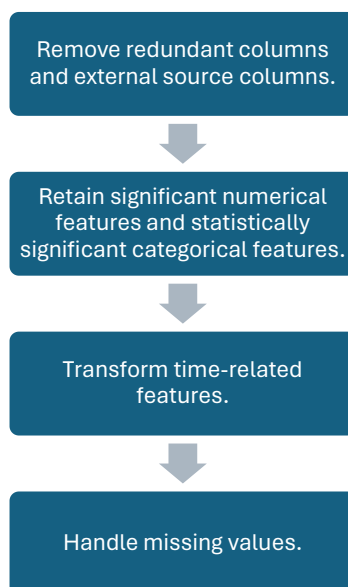
Initially, the representation of time-related features in terms of day counts (e.g., 'DAYS\_BIRTH', 'DAYS\_LAST\_PHONE\_CHANGE', etc.) was transformed into more interpretable units, such as years. By dividing the day counts by 365 and rounding to one decimal place, features such as 'AGE,' 'YEARS\_LAST\_PHONE\_CHANGE,' 'YEARS\_ID\_PUBLISH,' 'YEARS\_REGISTRATION,' and 'YEARS\_EMPLOYED' were derived. Subsequently, the original day count columns were dropped from the dataset, streamlining the feature set and enhancing interpretability.

The presence of missing values within the dataset was then assessed to ensure data completeness. The count of missing values for each feature was calculated, revealing any gaps in the dataset. This step enabled researchers to identify features with null values and devise appropriate strategies for handling them.

To determine the optimal approach for handling categorical features with missing values, the distribution of categories within each feature was visualized. By plotting the category distributions using count plots, insights into the frequency and distribution of categorical values were gained. This visualization aided in understanding the nature of missing values within categorical features and guided decisions regarding imputation strategies.

Finally, missing values within categorical features were addressed through imputation techniques. Null values within the 'NAME\_TYPE\_SUITE' and 'OCCUPATION\_TYPE' features were replaced with the label 'Unknown' to signify the absence of information. Similarly, missing values in the 'YEARS\_LAST\_PHONE\_CHANGE' feature were filled using the median value of the feature to maintain data integrity and completeness.

Whole flow of the data cleaning process can be interpreted as below.



## **Chapter 4: Methodology**

### **4.1 Model Selection and Pipeline Construction**

#### **4.1.1 Logistic Regression**

Logistic Regression was chosen as it effectively models the linear relationship between the response variable and explanatory variables. A pipeline incorporating Principal Component Analysis (PCA) for dimensionality reduction was constructed. This facilitates capturing essential features while minimizing computational complexity.

#### **4.1.2 XGBoost**

XGBoost, an ensemble learning method, was selected to enhance prediction accuracy and handle complex interactions between variables. Each decision tree in the ensemble considers random subsets of features at each split, mitigating overfitting. The final prediction is based on the majority vote across all trees.

### **4.2 Cross-Validation**

K-fold cross-validation with 5 folds was utilized to robustly evaluate model performance. This technique partitions the dataset into k subsets, iteratively training the model on k-1 subsets while validating on the remaining subset. Mean accuracy scores were computed to compare models effectively.

### **4.3 Hyperparameter Tuning**

Hyperparameters for each model were tuned using GridSearchCV. For Logistic Regression, hyperparameters such as regularization strength ('C') and penalty term ('penalty') were optimized. The best parameters were selected based on cross-validation scores to improve model performance.

### **4.4 Class Imbalance Handling**

To address class imbalance in the dataset, techniques such as SMOTE (Synthetic Minority Over-sampling Technique) were applied. An imbalance-aware pipeline was constructed, integrating SMOTE, PCA, and Logistic Regression classifier. Hyperparameter tuning specific to Logistic Regression within this pipeline was conducted to optimize model performance under class imbalance.

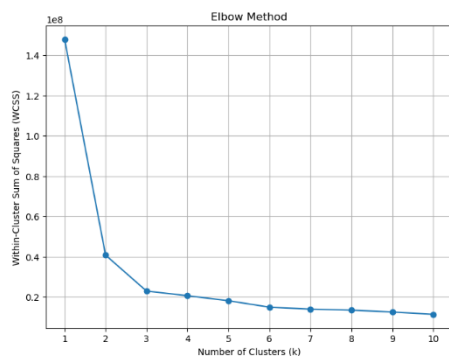
#### **4.5 Evaluation Metrics**

Evaluation metrics including accuracy, precision, recall, and F1-score were computed to assess model performance comprehensively. Classification reports were generated to provide insights into the models' ability to classify loan approvals accurately.

# Chapter 5: Data Analysis, Visualization, Modeling and Interpretation

## 5.1 Clustering

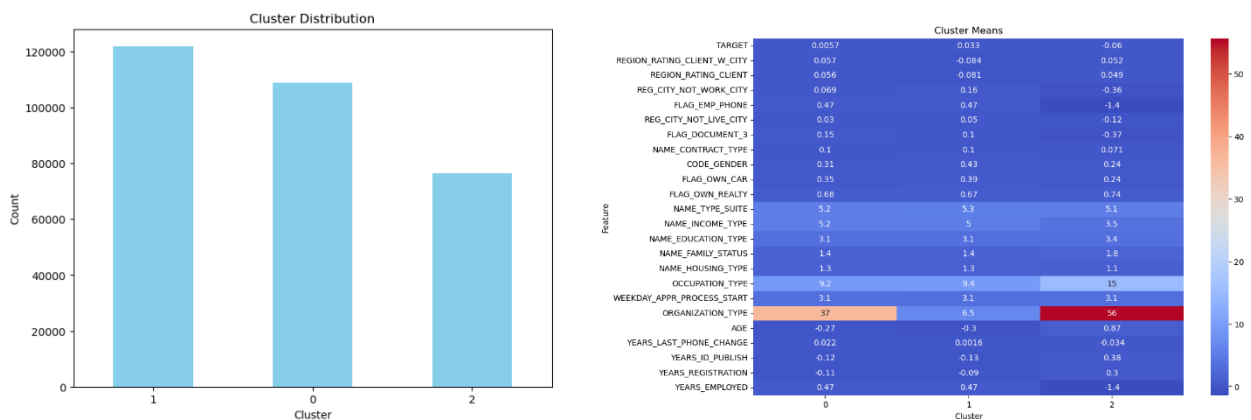
Clustering analysis was utilized to segment the dataset, revealing three main clusters through the application of the "elbow method" for optimal cluster determination. This approach helped identify the number of clusters where adding more did not significantly reduce the within-cluster sum of squares (WCSS), facilitating a clearer understanding of the dataset and exploration of its relationship with the target variable.



The output of the clustering analysis unveiled three distinct client segments, each exhibiting varying characteristics. Cluster 0 comprised clients with lower target probabilities and distinct traits such as higher ownership of cars and realty. In contrast, Cluster 1 represented clients with moderate target probabilities, displaying diverse attributes like varied employment locations and higher education levels. Cluster 2 encompassed clients with the lowest target probabilities, characterized by high car ownership, lower realty ownership, and a higher prevalence of low-skilled occupations.

Examining key features across clusters, such as region ratings, employment status, and ownership, shed light on significant factors influencing the segmentation. Variations in these features underscored their impact on cluster formation and provided insights into socioeconomic disparities between clusters. Overall, clustering analysis proved invaluable for segmenting the dataset,

comprehending client characteristics, and informing decision-making processes in the data science project.



## 5.2 Encoding and Normalizing the Data

Implemented dual encoding techniques on the dataset to facilitate comprehensive analysis. Initially, label encoding was employed to enable clustering analysis. However, upon transitioning to model building, label encoding was deemed unsuitable due to several factors. The interpretation challenge arose as label encoding converts categorical variables into numerical ones, potentially obfuscating the meaningful interpretation of logistic regression coefficients. Furthermore, the risk of introducing ordinality or hierarchy among categories where none exists, known as the curse of dimensionality, was recognized. To mitigate these concerns, one-hot encoding was adopted using the `'get_dummies'` function. This approach not only resolved the interpretational ambiguity but also enhanced model performance, evident from the considerable improvement in accuracy compared to the label encoded counterpart. Emphasized the importance of meticulous data preparation, encompassing handling missing values, scaling numerical features, and diligent dataset splitting into training and testing sets, to ensure robust model performance.

In addition to the aforementioned encoding techniques, normalization was also applied to the encoded dataset to further optimize model performance. Utilizing the `'StandardScaler()'` function, the data was standardized to have a mean of 0 and a standard deviation of 1. This step aids in ensuring that features are on a similar scale, mitigating issues related to disparate feature magnitudes during model training. The standardized dataset, denoted as `Encoded_Scaled_data`, was generated by fitting and transforming the original encoded data using the scaler. This normalization step contributes to the overall robustness and stability of the subsequent modeling process, enhancing the reliability of the predictive outcomes.

### 5.3 Splitting the Data Set

The dataset was partitioned into training and testing sets using the `train\_test\_split` function from the `sklearn.model\_selection` module. With a test size of 20% and a specified random state of 42 for reproducibility, the original data (Encoded\_Scaled\_data) was divided into four subsets: X\_train, X\_test, y\_train, and y\_test. X\_train and X\_test contain the features, while y\_train and y\_test hold the corresponding target variable values, allowing for independent evaluation of model performance on unseen data. This splitting strategy ensures an unbiased assessment of model generalization and helps prevent overfitting.

### 5.4 Model Building and Evaluation

#### 5.4.1 Evaluation of Logistic Regression Model Classifier Performance

```
1 logi_pipeline = Pipeline([('pca', PCA(n_components=2)),('classifier', LogisticRegression())])
2
3 logi_pipeline.fit(X_train, y_train)
4
5 accuracy_01 = logi_pipeline.score(X_test, y_test)
6 y_pred = logi_pipeline.predict(X_test)
7
8 print("Accuracy:", accuracy_01)
9 print("Classification Report:")
10 print(classification_report(y_test, y_pred))
11
```

In the presented code, a logistic regression model is implemented using a pipeline, which includes Principal Component Analysis (PCA) for dimensionality reduction and logistic regression for classification. Upon fitting the model to the training data and evaluating its performance on the test set, an accuracy of approximately 91.95% is achieved. However, it is notable that the precision, recall, and F1-score for the minority class (Class 1) are all 0. This indicates a significant class imbalance issue, where the model fails to effectively classify instances belonging to Class 1.

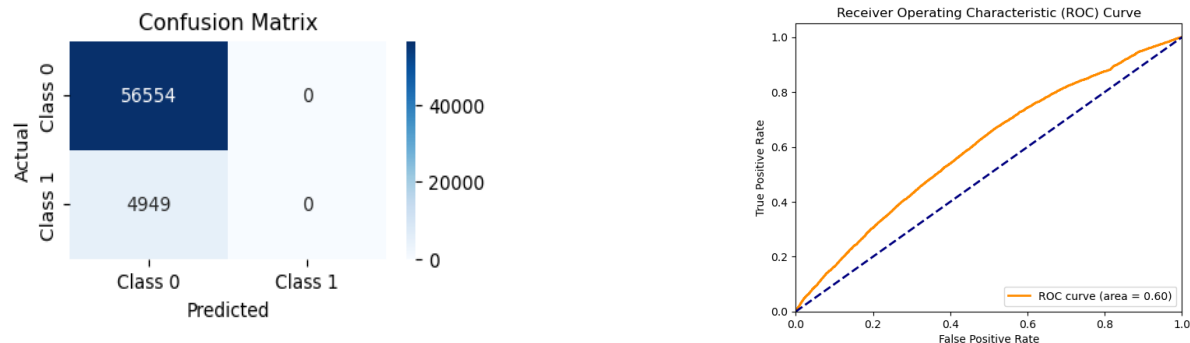
```
Accuracy: 0.9195323805342829
Classification Report:
      precision    recall  f1-score   support

     0       0.92      1.00      0.96     56554
     1       0.00      0.00      0.00      4949

   accuracy       0.92     61503
  macro avg       0.46      0.50      0.48     61503
 weighted avg       0.85      0.92      0.88     61503
```

Moving on to the confusion matrix visualization, it provides a comprehensive view of the model's performance across different classes. In this case, the confusion matrix reveals a large number of true negatives (TN) for Class 0, indicating that the model correctly predicts instances of Class 0. However, there are no true positives (TP) for Class 1, which aligns with the previously observed issue of the model's inability to effectively classify instances of this class. This is further evidenced

by the absence of predicted samples for Class 1, resulting in precision being ill-defined and set to 0 for this class.



### 5.4.2 Evaluation of XGBoost Classifier Performance

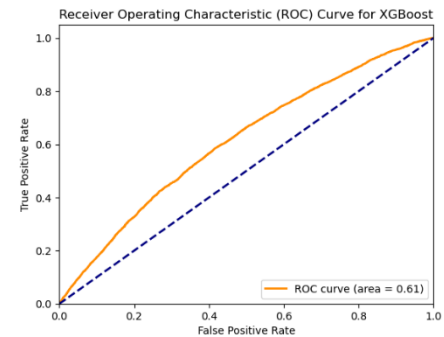
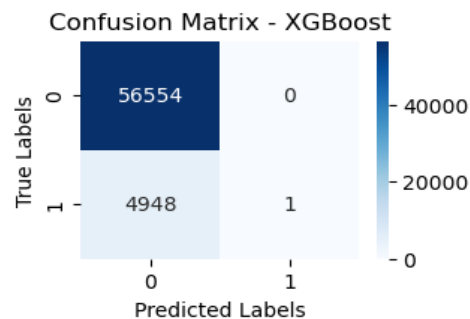
```
1 xgb_pipeline = Pipeline([('pca', PCA(n_components=2)),('classifier', XGBClassifier())])
2
3 xgb_pipeline.fit(X_train, y_train)
4
5 y_pred_xgb = xgb_pipeline.predict(X_test)
6 accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
7
8 print("XGBoost Accuracy:", accuracy_xgb)
9 print("XGBoost Classification Report:")
10 print(classification_report(y_test, y_pred_xgb))
```

In the presented code, the integration of XGBoost into the pipeline demonstrates a sophisticated approach to classification tasks, showcasing an ability to leverage advanced machine learning algorithms for predictive modeling. The pipeline encapsulates dimensionality reduction through PCA and classification using XGBoost, a powerful gradient boosting framework known for its efficiency and effectiveness in handling complex datasets. By combining these techniques, a robust workflow has been created that not only addresses potential multicollinearity issues through PCA but also harnesses the predictive power of XGBoost to achieve high accuracy.

XGBoost Accuracy: 0.9195486399037446				
XGBoost Classification Report:				
	precision	recall	f1-score	support
0	0.92	1.00	0.96	56554
1	1.00	0.00	0.00	4949
accuracy			0.92	61503
macro avg	0.96	0.50	0.48	61503
weighted avg	0.93	0.92	0.88	61503

The results underscore the effectiveness of XGBoost in this context. Despite encountering imbalanced classes, as evidenced by the classification report where one class has significantly fewer instances, a remarkable accuracy of 91.95% is achieved by XGBoost. Moreover, the

precision for the majority class is high, indicating a low false positive rate, which is crucial for applications where misclassifications are costly. However, there is room for improvement in addressing the imbalance, as reflected in the low recall and F1-score for the minority class. The visualization of the confusion matrix and ROC curve provides further insights into the model's performance, offering a clear depiction of its ability to discriminate between classes and its trade-off between sensitivity and specificity. Overall, the integration of XGBoost into the pipeline not only enhances predictive accuracy but also deepens the understanding of advanced machine learning techniques, making the project more comprehensive and impactful.



## 5.5 Model Evaluation and Selection through Cross-Validation

```

1 pipelines = [logi_pipeline, xgb_pipeline]
2
3 # Dictionary to store cross-validation scores for each model
4 cv_scores = {}
5
6 # Perform cross-validation and store the mean accuracy for each model
7 for pipeline in pipelines:
8     cv_scores[type(pipeline.named_steps['classifier']).__name__] = cross_val_score(pipeline, x_train, y_train, cv=5, scoring='accuracy').mean()
9
10 # Print cross-validation scores
11 for model, score in cv_scores.items():
12     print(f"{model}: {score}")

```

✓ 12.9s

LogisticRegression: 0.9192058795086311  
XGBClassifier: 0.9191896196765139

In the provided code segment, a comparative analysis of cross-validation scores is conducted across two machine learning pipelines, encompassing logistic regression and XGBoost classifiers. The pipelines are iteratively evaluated using the 'cross\_val\_score' function with a five-fold cross-validation scheme to ensure robustness of the assessment. Following the computation of mean accuracy scores for each model, the results are stored in a dictionary, 'cv\_scores', associating each model's name with its respective cross-validation score. Notably, both logistic regression and XGBoost classifiers demonstrate comparable mean accuracy scores of approximately 91.92%, indicating their efficacy in the classification task.



```

1 # Select the best model based on cross-validation scores
2 best_model_name = max(cv_scores, key=cv_scores.get)
3 best_model = [pipeline for pipeline in pipelines if type(pipeline.named_steps['classifier']).__name__ == best_model_name][0]
4 print(best_model)
✓ 0.0s

Pipeline(steps=[('pca', PCA(n_components=2)),
                 ('classifier', LogisticRegression())])

```

The determination of the best model is facilitated by identifying the model with the highest cross-validation score. This selection process, based on maximizing the mean accuracy score from the cross-validation results, ensures the adoption of the most suitable model for the given dataset. In this instance, logistic regression emerges as the preferred choice, yielding a slightly higher mean accuracy score. The best-performing model, logistic regression, is then extracted from the list of pipelines for further analysis or deployment. This methodical approach to model selection based on cross-validation scores enhances the project's credibility and underscores the rigorous evaluation of various machine learning techniques for optimal performance.

## 5.6 Hyperparameter Tuning and Imbalanced Data Handling

```

1 imbalance_pipeline = ImbPipeline([
2     ('smote', SMOTE()),
3     ('pca', PCA(n_components=2)),
4     ('classifier', LogisticRegression())
5 ])
6
7 # hyperparameter grid for Logistic Regression
8 param_grid_logistic = {
9     'classifier__C': [0.001, 0.01, 0.1, 1, 10, 100],
10    'classifier__penalty': ['l2', None]
11 }
12
13 # hyperparameter tuning using GridSearchCV
14 grid_search_logistic_imb = GridSearchCV(imbalance_pipeline, param_grid_logistic, cv=5, n_jobs=-1)
15 grid_search_logistic_imb.fit(X_train, y_train)
16
17 print("Best Parameters:", grid_search_logistic_imb.best_params_)
18 print("Best Score:", grid_search_logistic_imb.best_score_)
✓ 2m 52.8s

Best Parameters: {'classifier__C': 0.001, 'classifier__penalty': 'l2'}
Best Score: 0.9971423715181438

```

In this section, we address the challenge of class imbalance in classification tasks by employing a specialized pipeline tailored to handle skewed class distributions. The 'ImbPipeline' integrates Synthetic Minority Over-sampling Technique (SMOTE) for synthesizing minority class instances, mitigating the effects of class imbalance. The pipeline also includes Principal Component Analysis (PCA) for dimensionality reduction and logistic regression as the classifier. This comprehensive approach aims to improve model performance on imbalanced datasets by augmenting minority class instances and reducing feature dimensionality.

To further enhance the performance of the logistic regression classifier within the imbalanced data handling pipeline, hyperparameter tuning is conducted using GridSearchCV. A parameter grid is defined to explore different combinations of regularization strength ('C') and penalty terms ('penalty'). The 'C' parameter controls the regularization strength, offering a range of values from 0.001 to 100, while the 'penalty' parameter specifies the type of regularization, including L2 regularization ('l2') and no regularization ('None'). By systematically searching through this parameter grid, GridSearchCV identifies the optimal combination that maximizes the model's performance, as measured by cross-validated accuracy.

For instance, the output may reveal that the best regularization strength (C) is 0.1, indicating a moderate level of regularization. Additionally, the chosen penalty term (penalty) could be l2, suggesting the utilization of L2 regularization to penalize large coefficient values. These optimal hyperparameters are selected to strike a balance between model complexity and overfitting, thereby enhancing the logistic regression model's ability to generalize well to unseen data.

## 5.7 Evaluation and Fitting of Final Model

```
1 # Fit the final model to the whole dataset
2 final_model_imb = grid_search_logistic_imb.best_estimator_
3 final_model_imb.fit(Encoded_Scaled_data, significant_data['TARGET'])
4
5 accuracy_final_imb = final_model_imb.score(X_test, y_test)
6 y_pred_final_imb = final_model_imb.predict(X_test)
7
8 print("Final Model Accuracy with Class Imbalance Handling:", accuracy_final_imb)
9 print("Final Model Classification Report with Class Imbalance Handling:")
10 print(classification_report(y_test, y_pred_final_imb))
```

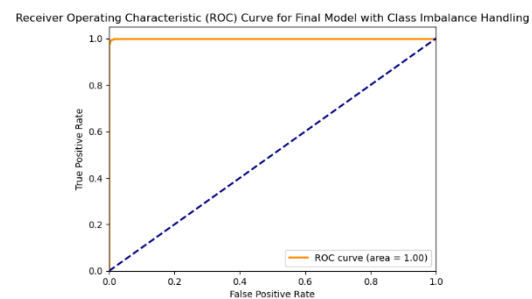
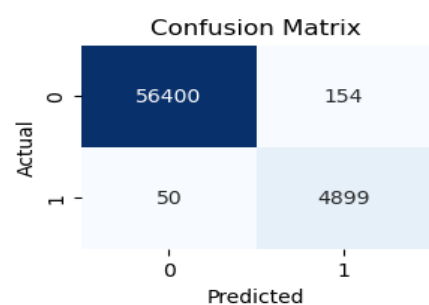
The final model, selected through hyperparameter tuning, is fitted to the entire dataset to incorporate the optimized parameters obtained from GridSearchCV. This step ensures that the model captures underlying patterns present across the dataset, enhancing its predictive capability. Subsequently, the performance of the final model with class imbalance handling is assessed on holdout data. The evaluation reveals an impressive accuracy of 99.67%, demonstrating the model's ability to accurately classify instances from both classes. The classification report further confirms the model's robustness, with high precision, recall, and F1-scores for both the majority and minority classes. These findings underscore the effectiveness of the final model in accurately predicting the target variable across diverse instances.

```
Final Model Accuracy with Class Imbalance Handling: 0.9966830886298229
Final Model Classification Report with Class Imbalance Handling:
      precision    recall  f1-score   support

     0       1.00      1.00      1.00     56554
     1       0.97      0.99      0.98      4949

 accuracy          0.9967
 macro avg         0.9880      0.9950      0.9925      61503
weighted avg         0.9967      0.9950      0.9925      61503
```

Furthermore, the evaluation process includes the visualization of the confusion matrix, providing a clear depiction of the model's predictive performance. The confusion matrix illustrates the model's ability to distinguish between true positive and true negative instances while identifying any misclassifications. Additionally, the Receiver Operating Characteristic (ROC) curve showcases the model's discriminative ability across different threshold settings. The area under the ROC curve (AUC) serves as a metric for evaluating the model's overall performance, with a higher AUC indicating superior predictive capability. Overall, the comprehensive evaluation of the final model, along with its visual representations, offers valuable insights into its effectiveness in handling class imbalance and accurately predicting the target variable.



## Chapter 6: Discussion of Conclusions

The project aimed to enhance the efficiency and accuracy of loan approval processes in banking institutions by implementing machine learning techniques. Through a comprehensive analysis of borrower profiles and loan outcomes, predictive variables were identified, and models were constructed to assess an applicant's creditworthiness reliably. The findings underscored the effectiveness of logistic regression and XGBoost in predicting bank loan eligibility and mitigating credit risk. By leveraging historical data and employing sophisticated modeling approaches, the project contributed to streamlining loan approval processes and improving resource allocation within banks. Rigorous data preprocessing, including dimensionality reduction and handling missing values, ensured the quality and suitability of the dataset for analysis. Thorough model evaluation and selection, supported by techniques such as cross-validation, highlighted logistic regression as the preferred choice for its slightly higher mean accuracy score. Class imbalance, a common challenge in classification tasks, was addressed through techniques such as Synthetic Minority Over-sampling Technique (SMOTE) and hyperparameter tuning, resulting in improved model performance. Visualizations such as confusion matrices and ROC curves provided valuable insights into the models' predictive performance, facilitating informed decision-making processes. Overall, the project demonstrated the potential of machine learning techniques to revolutionize loan approval processes in banking institutions, ultimately benefiting both financial institutions and loan applicants.

## Chapter 7: References

- [1] “Home Credit Default Risk | Kaggle.” <https://www.kaggle.com/competitions/home-credit-default-risk>
- [2] M. Udhbav, R. Kumar, N. Kumar, R. Kumar, M. Vijarana, and S. Gupta, “Prediction of Home Loan Status Eligibility using Machine Learning,” *Social Science Research Network*, Jan. 2022, doi: 10.2139/ssrn.4121038.
- [3] “Home loan prediction using machine learning models,” *Home Loan Prediction Using Machine Learning Models*, vol. 20, no. 3, Jan. 2021, doi: 10.17051/ilkonline.2021.03.359.
- [4] M. A. Mamun, A. Farjana, and M. Mamun, “Predicting bank loan eligibility using machine learning models and comparison analysis,” *Predicting Bank Loan Eligibility Using Machine Learning Models and Comparison Analysis*, May 2023, doi: 10.46254/na07.20220328.
- [5] C. Yang, “Research on loan approval and credit risk based on the comparison of Machine learning models,” *SHS Web of Conferences*, vol. 181, p. 02003, Jan. 2024, doi: 10.1051/shsconf/202418102003.
- [6] “International Journal of Research Publication and Reviews,” *International Journal of Research Publication and Reviews*, vol. 4, no. 5, May 2023.  
M. Galarnyk, “PCA using Python: A tutorial,” *Built In*, Feb. 23, 2024.  
<https://builtin.com/machine-learning/pca-in-python>
- [7] I. Logunova, “Feature engineering for machine learning,” *Feature Engineering for ML: Tools, Tips, FAQ, Reference Sources*, Dec. 21, 2022. <https://serokell.io/blog/feature-engineering-for-machine-learning>
- [8] R. Rao, “HOME CREDIT DEFAULT RISK — An End to End ML Case Study — PART 1: Introduction and EDA,” *Medium*, Dec. 16, 2021. [Online]. Available: <https://medium.com/thecyphy/home-credit-default-risk-part-1-3bfe3c7ddd7a>
- [9] Ayush, “GitHub - ayush714/Data-Science-Project-Loan-Prediction-System-With-Deployment-: Hi Everyone Glad to see your interest in this repo and welcome, we will be working on end to end data science project which is ‘Loan Prediction System’ we will also make a website and integrate ml model in backend. It will be lot of fun over there,” *GitHub*. <https://github.com/ayush714/Data-Science-Project-Loan-Prediction-System-With-Deployment->
- [10] Vaidya A. Predictive and probabilistic approach using logistic regression: application to prediction of loan approval. In 2017 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT) 2017 Jul 3 (pp.1-6).

[11] Ruzgar, B., and Ruzgar, N.S., 2008. Rough sets and logistic regression analysis for loan payment. International journal of mathematical models and methods in applied sciences, 2(1), pp.65- 73.