BINDURA UNIVERSITY OF SCIENCE EDUCATION

DEPARTMENT OF MATHEMATICS AND STATISTICS

FACULTY OF SCIENCE



An Assessment On The Machine Learning Approaches For Credit Scoring And Default Prediction: A Case Study Of Cabs Bank

BY

Vuyo Nyembezi (B200666B)

A PROJECT SUBMITTED IN PARTIAL FULFILMENT OF THE REQUIREMENTS OF BACHELOR OF SCIENCE HONOURS DEGREE IN STATISTICS AND FINANCIAL MATHEMATICS

SUPERVISOR: MS.P.HLUPO

APPROVAL FORM

This is to certify, that this research project is the result of my own research work and has not been copied or extracted from past sources without acknowledgement. I hereby declare that no part of it has been presented for another degree in this University or elsewhere.



DEDICATION

I dedicate this project to my uncle Sifelani Tsiko,my father Vusunzi Walter Nyembezi, my mother Thokozile Nyembezi and also my younger brother Tulani Nyembezi for their support and love throughout the research process.

ACKNOWLEDGEMENTS

I would like to thank the Lord Almighty for making the dream come true. Without God's guidance, this research could not have been successful. I really appreciate Ms Hlupo my superviser for her support and guidance throughout the study. My token of appreciation goes to the Nyembezi family for their support and sacrifices throughout it was not easy.

Abstract

This research examines how machine learning methods can be used for credit scoring and default prediction at CABS Bank in Zimbabwe. Conventional credit evaluation techniques that mainly rely on historical financial data and set criteria often struggle to accurately predict credit risks, especially with the complexity of modern financial data. This study explores using advanced machine learning approaches to improve the precision and efficiency of assessing creditworthiness. The research involves creating and implementing machine learning models, comparing their effectiveness, and evaluating their fairness and ethical implications. The main findings demonstrate that machine learning models, like Random Forest and Neural Networks, outperform traditional methods in anticipating loan defaults, offering a more thorough and precise evaluation of credit risk. The research concludes by suggesting ways to incorporate these advanced methodologies into banking practices to enhance loan approval processes and promote wider financial inclusivity.

TABLE OF CONTENTS

APPROVAL FORM	ii
DEDICATION	iii
ACKNOWLEDGEMENTS	iv
Abstract	v
TABLE OF CONTENTS	vi
LIST OF TABLES	xi
LIST OF FIGURESx	ii
LIST OF ACRONYMS	iii
CHAPTER 1: INTRODUCTION	1
1.0 Introduction	1
1.1 Background of study	2
1.2 Problem statement	3
1.3 Research objectives	4
1.4 Research questions	4
1.5 Assumptions	4
1.6 Significance of study	5
1.7 Delimitation of the study	5
1.8 Limitations of the study	5
1.9 Research hypothesis	6
1.10 Scope	6
1.11 Definition of terms	7
1.12 Conclusion	8
CHAPTER 2: LITERATURE REVIEW	9
2.0 Introduction	9
2.1 Theoretical Literature	9

2.1.1 Credit Risk and Probability of Default
2.1.2 Creditworthiness
2.1.3 Conventional Credit Scoring Models
2.1.4 Machine Learning Approaches
2.1.5 Underlying Theories and Frameworks
2.1.5.1 Credit Risk Assessment Frameworks
2.1.5.2 Agency Theory 10
2.1.5.3 Moral Hazard Theory 11
2.1.5.4 Liquidity Preference Theory 11
2.1.5.5 Financial Theory 11
2.1.5.6 Asymmetric Information Theory 11
2.1.5.7 Signaling Theory
2.2 Empirical Literature
2.2.1 Traditional Models
2.2.2 Machine Learning Approaches
2.2.3 Hybrid Approaches
2.2.4 Credit Scoring for Specific Contexts
2.2.5 Limitations and Challenges
2.3 Research Gap
2.4 Proposed Conceptual Model
2.5 Conclusions
CHAPTER 3: RESEARCH METHODOLOGY 17
3.1 Introduction
3.2 Research Design
3.3 Data Sources
3.4 Research Instruments

3.5 Description of Variables and Expected Relationships	. 18
3.5.1 Independent variables: Demographic Variables	. 18
3.5.1.1 Age	. 18
3.5.1.2 Gender	. 19
3.5.1.3 Income	. 19
3.5.1.4 Educational Level	. 19
3.5.1.5 Employment Status	. 19
3.5.2.1 Credit Score	. 20
3.5.2.2 Payment History	. 20
3.5.3 Loan Application Variables	. 20
3.5.3.1 Loan Amount	. 20
3.5.3.2 Loan Term	. 21
3.5.3.3 Interest Rate	. 21
3.5.4 Dependent Variable	. 21
3.5.4.1 Default Status	. 21
3.5.5 Expected Relationships	. 21
3.6 Data Analysis Procedures	. 22
3.6.1 Logistic Regression	. 22
3.6.2 Machine learning models	. 23
3.6.2.1 Random Forest	. 23
3.6.2.2 Support Vector Machine	. 23
3.7 Ethical Considerations	. 24
3.8 Conclusion	. 25
CHAPTER 4: DATA PRESENTANTION, ANALYSIS AND INTERPRETATION	. 26
4.0 Introduction	. 26
4.1 Descriptive Statistics	. 26

4.1.2 Graphs showing descriptive statistics	
4.2 Pre-tests	
4.2.1 Correlations	
4.3 Model Output/Results	
4.3.1 Confusion Matrix for Logistic Regression:	
4.3.2 Confusion Matrix for SVM:	
4.3.3 Confusion Matrix for Random Forest:	
4.4 Model Validation Tests	
4.5 Discussion of Findings	
4.6 Conclusion	
CHAPTER 5: SUMMARY, CONCLUSIONS AND RECOMMENDATIONS	
5.0: INTRODUCTION	
5.1 Summary of Findings	
5.2 Conclusions	
5.3 Recommendations	50
5.4 Areas for Further Research	
5.5 Conclusion	
References	53
A1 code for ML models used in the thesis	
## Importing Dependencies	56
## Define Features and Target Variable	57
## Split the Dataset into Training and Testing Sets	57
## Standardize Features	57
## Standardize Features ## Logistic Regression Model	57 58
## Standardize Features ## Logistic Regression Model ## Random Forest Model	57 58 59

## ROC Curve and AUC for Random Forest	
## Plot ROC Curves	61
## Cross-Validation for Logistic Regression	
## Cross-Validation for SVM	
Appendix	
A1 code for ML models used in the thesis	

LIST OF TABLES

Table 4.1: Descriptive statistics	26
Table 4.2 Correlations	34
Table 4.3 Machine learning outputs results	40
Table 4.4 Cross-validation scores	44

LIST OF FIGURES

15
28
30
30
31
33
36
38
38
42

LIST OF ACRONYMS

- SVM Support Vector Machine
- LR Logistic Regression
- RF Random Forest
- AI Artificial intelligence
- ML Machine learning
- PD Probability of default
- LGD loss given default
- EL expected loss

CHAPTER 1: INTRODUCTION

1.0 Introduction

In the modern landscape of banking and finance, the use of advanced computational techniques has garnered significant attention in redefining credit assessment and default prediction strategies. Artificial Intelligence (AI), machine learning branch has become a potent instrument within banking institutions regarding evaluating creditworthiness and predicting loan defaults. This paradigm shift towards machine learning methodologies in credit assessment has been driven by the need for more accurate, efficient, and adaptive systems capable of handling vast amounts of data. According to Smith and Johnson (2018), the traditional methods of credit scoring and default prediction have faced limitations in handling the complexities of modern financial data and predicting credit risks accurately.

Hence, this study aims to delve into the realm of machine learning applications within banking institutions, evaluating their efficacy, limitations, and ethical implications in credit scoring and default prediction. By examining recent advancements and insights in the field, this research endeavors to contribute to the ongoing discourse, shedding light on the potential and challenges of integrating machine learning into credit assessment processes within contemporary banking institutions.

This section serves as a comprehensive introduction to the study, laying the groundwork by discussing the study's background, problem statement, objectives, research questions, assumptions, limitations, and definitions. It serves as a foundation for subsequent chapters that delve deeper into the study's various aspects. The forthcoming Chapter 2 conducts a thorough review of existing research concerning machine learning approaches for comparison. It aims to critically analyze and synthesize the insights provided by prior studies in this domain. Moving forward, Chapter 3 meticulously explores and elucidates the methods and models employed in the study, offering a detailed explanation of the applied methodology. This chapter focuses on the technical aspects and the procedures used to conduct the research. Chapter 4 is dedicated to presenting the analysis and findings derived from the conducted research. This section is pivotal

in comprehending the outcomes and implications of the study. Finally, Chapter 5 encapsulates the study, drawing conclusions and presenting recommendations based on the findings. This chapter serves as a culmination, offering insights and actionable suggestions derived from the study's results.

Overall, this section provides an in-depth overview and sets the context for the study's comprehensive exploration of machine learning approaches in credit assessment.

1.1 Background of study

In the world of finance, assessing an individual's creditworthiness has been a cornerstone of lending practices for centuries. This evaluation process is pivotal in determining whether an applicant is eligible for a loan and under what terms. Traditional credit scoring methods, once solely reliant on historical financial data and predetermined criteria, have evolved significantly with the advent of machine learning.

Within Zimbabwe's financial landscape, access to credit and reliable credit assessment mechanisms plays a pivotal role in economic development. As the country navigates through economic challenges, such as inflation and currency volatility, there's a growing emphasis on inclusive financial practices and innovative approaches to credit evaluation (The World Bank, 2018). With advancements in technology and the adoption of machine learning in various sectors, Zimbabwe's financial institutions are exploring novel methodologies to refine credit scoring systems, aiming for more accurate risk assessment and increased financial inclusivity within the population (Reserve Bank of Zimbabwe, 2015).

Historically, credit assessment models typically focused on basic factors like income, credit history, and job status. However, these traditional models had their limitations, often resulting in errors and the exclusion of potentially creditworthy individuals, especially those with limited credit histories. This led to a change in the financial industry.

The emergence of machine learning algorithms presented a transformative opportunity. Artificial Intelligence (AI) includes machine learning as a subset, that provided a innovative way to credit assessment by enabling systems to gain knowledge by analysing data, detect patterns and forecast

outcomes without the need for explicit programming. This technology's capacity to process vast amounts of data, including non-traditional variables such as social media behaviour or payment history, revolutionized the credit evaluation landscape.

Machine learning models provide numerous advantages over traditional credit scoring methods. They permit a more thorough analysis, giving financial institutions the ability to evaluate creditworthiness beyond the constraints of traditional criteria. By incorporating various data sources and factors, these models have the potential to offer a more holistic and precise assessment of an individual's creditworthiness.

The evolution of machine learning in credit assessment marks a significant leap forward in the quest for more accurate, inclusive, and adaptive lending practices. Understanding the historical context of credit evaluation, the impact of machine learning, and the challenges it poses is essential for shaping a future where creditworthiness assessment is both effective and equitable.

This study aims to explore the matter more into the application of artificial intelligence approaches for predicting loan approval, exploring their potential, limitations, and ethical implications within the realm of creditworthiness assessment.

1.2 Problem statement

Enhancing accuracy in creditworthiness assessment through machine learning models remains a critical challenge in modern banking. Amidst a plethora of available algorithms, selecting the most effective model for credit evaluation poses a pivotal concern (Smith et al., 2021; Chen & Liu, 2019). The pressing need to optimize model selection for credit scoring and loan approval is evident, requiring a focused exploration into identifying and adopting the most robust and accurate machine learning models. This investigation seeks to determine the most suitable model that significantly improves accuracy and reliability in credit assessment within banking institutions.

1.3 Research objectives

These objectives aim to delve deeply into the practical application and comparative analysis of different machine learning models, providing a comprehensive understanding of their strengths, limitations, and adaptability in the intricate landscape of financial markets.

- 1. Develop and implement a machine learning-based credit scoring model within a banking institution to significantly enhance loan repayment prediction accuracy.
- 2. Conduct a comparative analysis of machine learning algorithms for default prediction within a banking institution.

1.4 Research questions

These research questions aim to probe into the viability, adaptability, and implications of employing advanced computational methodologies in the realm of portfolio optimization within Zimbabwe's financial context.

- 1. How can machine learning help banks decide if someone will repay a loan?
- 2. Which machine learning methods work best for predicting if someone might not repay a loan?
- 3. What data should bank use to make machine learning-based credit decisions?
- 4. Are machine learning models fair in predicting who might default on a loan?
- 5. How can banks use machine learning to improve their accuracy in lending decisions?

1.5 Assumptions

- 1. Assuming that using machine, learning can make loan decisions more accurate.
- 2. Assuming that more data sources will improve machine learning models for credit scoring.
- 3. Assuming that machine learning can identify patterns to predict loan defaults.
- 4. Assuming that fairness and lack of bias can be achieved in machine learning based credit assessments.
- 5. Assuming that banks can adopt machine learning to enhance their lending decisions' accuracy and efficiency.

1.6 Significance of study

- 1. Helps banks make better decisions: This study can assist banks in deciding who gets loans and who might not, making their decisions more accurate.
- 2. Improves fairness: By using machine learning, the study aims to make sure everyone gets a fair chance for a loan, regardless of biases.
- 3. Reduces loan defaults: It's important because it might help banks predict when people might not pay back their loans, reducing the number of loans that don't get repaid.
- 4. Enhances lending accuracy: This study might help banks make more accurate decisions about who should get a loan, making the process better for both banks and borrowers.
- 5. Shapes future banking: By using these new methods, the study might change how banks make decisions about loans, potentially making banking better for everyone involved.

1.7 Delimitation of the study

- 1. Focus on specific methods: This study looks at only certain machine learning methods used in credit scoring and default prediction, not all available methods.
- 2. Limited to banking: The study focuses on machine learning in credit decisions made by banks and doesn't include other financial institutions.
- 3. Based on available data: It's limited to using data that's already available and might not cover all possible sources relevant to credit scoring.
- 4. Time constraints: This study might not cover the most recent developments in machine learning for credit decisions due to time limitations.
- 5. Regional focus: The study might concentrate more on specific regions or countries' banking systems and might not include global perspectives.

1.8 Limitations of the study

- 1. Limited fairness: Sometimes, machine learning might make unfair decisions, leaving out certain people or groups.
- 2. Not always accurate: Machine learning isn't perfect and might make mistakes in predicting who might not repay a loan.

- 3. Needs a lot of data: Using machine learning requires a lot of information, and sometimes, not all the needed data is available.
- 4. Hard to understand: Machine learning models can be complicated, making it tough for people to know why a certain decision was made.
- 5. Costly and time-consuming: Implementing machine learning can be expensive and take a lot of time, making it tough for smaller banks to use.

1.9 Research hypothesis

In the context of assessing machine learning approaches for credit scoring and default prediction in banking institutions, the research hypothesis aims to investigate whether machine learning methods outperform or provide better predictions compared to traditional methods used in credit assessment within the banking sector. This hypothesis seeks to explore whether employing machine learning algorithms enhances accuracy, fairness, or overall effectiveness in determining creditworthiness and predicting defaults when compared to conventional approaches utilized by banks.

H0: Traditional methods used in banking for credit scoring and default prediction are equally effective as machine learning approaches.

H1: Machine learning approaches in banking for credit scoring and default prediction are more effective than traditional methods.

1.10 Scope

The study focuses on CABS bank to assess and introduce new machine learning models to facilitate fast and smooth flow of credit worthiness.

1.11 Definition of terms

Machine Learning: A subset of artificial intelligence that involves the development of algorithms enabling computers to learn from data, identify patterns, and make predictions or decisions without explicit programming.

Credit Scoring: A statistical method used by lenders to assess the creditworthiness of individuals or entities considering their credit background, financial behaviour, as well as various different elements in order to assess the risk of lending to them.

Default Prediction: The process of using historical data and predictive analytics, often employing machine learning algorithms, to forecast the likelihood of a borrower failing to repay a loan or meeting financial obligations.

Algorithm Bias refers to a systematic errors or biases within machine learning algorithms that can potentially result in unfair or discriminatory outcomes, often due to the data used to train the model.

Random Forest is a machine learning algorithm that consists of multiple decision trees used for tasks involving classification and regression by combining the forecasts of numerous individual trees.

Support Vector Machines (SVM): A supervised learning technique that identifies an optimal hyper plane to separate different classes of data by maximizing the margin between them.

Ethical Considerations in Machine Learning: The examination and incorporation of ethical principles, fairness, transparency, and accountability inside the design, development, as well as deployment from machine learning models to mitigate biases and verify responsible use.

Neural Networks: A computational model inspired by the composition and operation of the human brain, composed of networked nodes or neurons that perform and analyse complex data to learn patterns and make predictions.

1.12 Conclusion

The chapter has laid the groundwork for the entire research. It has presented the core concepts of the study, covering its background, the issue at hand, the queries driving the research, underlying presumptions, and the importance of the study. This section has paved the way for the following chapter, which will extensively explore the existing literature.

CHAPTER 2: LITERATURE REVIEW

2.0 Introduction

The purpose of this literature review is to provide a comprehensive overview of existing research on machine learning approaches in credit scoring and default prediction, with a specific focus on CABS Bank, the largest building society in Zimbabwe. CABS Bank offers a wide range of financial products and services, including transaction and savings accounts, mortgage loans, and mobile banking. As a leading financial institution, CABS Bank faces the challenge of credit risk management, where accurate credit scoring and default prediction are crucial for minimizing losses and maximizing profits. This literature review aims to synthesize the findings of previous research, identify gaps in the literature, and provide a foundation for the subsequent chapters of this study.

2.1 Theoretical Literature

Credit scoring and default prediction are built on a foundation of theoretical concepts and frameworks that have evolved over time. Understanding these theoretical underpinnings is essential for appreciating the development and application of credit scoring models.

2.1.1 Credit Risk and Probability of Default

Risk associated with credit is the likelihood of a debtor failing to make payments or credit obligation (Saunders & Allen, 2010). Probability of default (PD) is a measure of this risk, representing how likely it is that a borrower will default over a certain period (Bloomberg, 2020). Credit scoring models aim to estimate PD to determine creditworthiness.

2.1.2 Creditworthiness

Creditworthiness refers to a borrower's ability and willingness to repay debts (FICO, 2020). It encompasses various factors, including credit history, payment behavior, credit utilization, and credit depth (Experian, 2020). Credit scoring models assess these factors to determine a borrower's creditworthiness.

2.1.3 Conventional Credit Scoring Models

Conventional credit scoring models utilize statistical methods like logistic regression and discriminant analysis to forecast one's creditworthiness (Thomas, 2017). These models rely on a restricted amount of credit information, like credit reports and payment history, to calculate Probability of Default (PD).

2.1.4 Machine Learning Approaches

Machine learning approaches, such as neural networks and decision trees, have been increasingly applied to credit scoring (Santos et al., 2020). These models can handle large datasets, complex relationships, and non-linear interactions, enhancing predictive accuracy.

2.1.5 Underlying Theories and Frameworks

Several theories and frameworks underpin credit scoring models, providing a foundation for understanding credit risk assessment and default prediction. These theories and frameworks help explain the complex relationships between borrowers, lenders, and credit markets, and inform the development of credit scoring models that accurately assess credit risk and predict default probability.

2.1.5.1 Credit Risk Assessment Frameworks

Credit risk assessment frameworks, such as the Basel Accords, have been updated to incorporate advanced approaches for credit risk modeling (BCBS, 2019). These frameworks provide guidelines for managing credit risk in banks and other financial entities, including the use of probability of default (PD), loss given default (LGD), and expected loss (EL) models. The Basel Accords have been influential in shaping credit risk assessment practices globally, and have been adopted by many countries as a basis for their regulatory capital requirements.

2.1.5.2 Agency Theory

Agency Theory has been applied to credit risk assessment, highlighting the importance of information asymmetry and moral hazard (Minton et al., 2018). Agency Theory posits that a

principal-agent relationship exists between lenders and borrowers, where lenders delegate credit decisions to borrowers who may have different interests and risk preferences. This theory helps explain how credit scoring models can influence borrower behavior and credit risk.

2.1.5.3 Moral Hazard Theory

Moral Hazard Theory has been used to explain how credit scoring models can influence borrower behavior (Dimitrov et al., 2020). Moral Hazard Theory posits that borrowers may take on excessive risk or engage in risky behavior after obtaining credit, knowing that lenders bear the risk. This theory helps explain how credit scoring models can create incentives for borrowers to take on more credit risk than they can afford.

2.1.5.4 Liquidity Preference Theory

Liquidity Preference Theory has been incorporated into credit risk models to account for short term creditworthiness and liquidity metrics (Kashyap et al., 2017). Liquidity Preference Theory posits that lenders prioritize liquidity when making credit decisions, and that borrowers with higher liquidity preferences are more likely to default. This theory helps explain how credit scoring models can account for short-term creditworthiness and liquidity metrics.

2.1.5.5 Financial Theory

Financial Theory has been applied to credit risk assessment, including the use of expected loss and risk premia (Giesecke et al., 2018). Financial Theory posits that credit risk is a function of expected loss and risk premia, and that credit scoring models should account for these factors. This theory helps explain how credit scoring models can estimate expected loss and risk premia, and how these estimates can be used to predict default probability.

2.1.5.6 Asymmetric Information Theory

Asymmetric Information Theory has been used to develop credit scoring models that account for information imbalance between lenders and borrowers (Bolton et al., 2016). Asymmetric Information Theory posits that lenders and borrowers have different levels of information about credit risk, and that credit scoring models should account for this information imbalance. This

theory helps explain how credit scoring models can account for asymmetric information and predict default probability.

2.1.5.7 Signaling Theory

Signaling Theory has been applied to credit risk assessment, highlighting the importance of borrower signals and credit history (Berg et al., 2019). Signaling Theory posits that borrowers send signals to lenders about their creditworthiness through their actions and characteristics, and that credit scoring models should account for these signals. This theory helps explain how credit scoring models can use borrower signals and credit history to predict default probability.

These theories and frameworks provide a solid foundation for understanding credit scoring and default prediction, and inform the development of credit scoring models that accurately assess credit risk and predict default probability. By understanding the underlying theories and frameworks, credit scoring models can be developed that better account for credit risk and default probability, and that provide more accurate and reliable credit scores.

2.2 Empirical Literature

This section provides a comprehensive review of existing empirical studies on credit scoring and default prediction, encompassing their findings, methodologies, and limitations.

2.2.1 Traditional Models

Traditional credit scoring models have been extensively utilized in practice, and numerous studies have evaluated their performance. Wang et al. (2019) found that logistic regression models outperformed decision trees in predicting default probability, while Kim et al. (2018) demonstrated that credit scoring models using financial ratios and credit history data achieved high accuracy in predicting default. Lessor et al. (2019) showed that traditional credit scoring models performed well in predicting default but struggled with capturing complex relationships.

2.2.2 Machine Learning Approaches

Machine learning techniques have gained popularity in credit scoring due to their ability to handle complex data and non-linear relationships. . Research has shown that machine intelligence models, neural networks and random forests, for example, can outperform traditional models in predicting default probability. Santos et al. (2020) developed a neural network model that achieved higher accuracy than traditional models in predicting default, while Zhang et al. (2019) demonstrated that random forests outperformed logistic regression in predicting default probability.

2.2.3 Hybrid Approaches

Hybrid models combining traditional and machine learning techniques have also been explored. Chen et al. (2020) developed a hybrid model using logistic regression and decision trees, achieving higher accuracy than individual models. Similarly, Huang et al. (2020) showed that hybrid models combining traditional and machine learning techniques performed well in predicting default.

2.2.4 Credit Scoring for Specific Contexts

Several studies have focused on credit scoring for specific contexts, such as small businesses or online lending. Xu et al. (2020) developed a credit scoring model for small businesses using machine learning techniques, while Li et al. (2020) developed a credit scoring model for online lending using hybrid approaches.

2.2.5 Limitations and Challenges

While empirical studies have made significant contributions to credit scoring and default prediction, limitations and challenges remain. Common issues include data quality, sample bias, and model overfitting. Additionally, most studies focus on specific contexts or datasets, making generalizability a concern. Königstorfer and Thalmann (2020) highlighted the need for more research on credit scoring for specific industries.

2.3 Research Gap

Despite the significant contributions of existing literature to credit scoring and default prediction, several gaps and limitations remain. One of the main limitations is the lack of generalizability across different industries and populations. For instance, Wang et al. (2019) developed a credit scoring model using logistic regression, but their study was limited to a specific dataset and industry. Similarly, Santos et al. (2020) proposed a neural network model for credit scoring, but their study focused only on a specific type of loan.

Another gap in the existing literature is the limited exploration of hybrid approaches combining traditional and machine learning techniques. While machine learning techniques have shown promise in credit scoring, there is a need for further research on hybrid models that can leverage the strengths of both approaches. For example, Chen et al. (2020) developed a hybrid model using logistic regression and decision trees, but their study was limited to a specific dataset and did not explore other machine learning techniques.

Furthermore, existing literature often prioritizes accuracy over interpretability, making it difficult to understand the underlying factors driving default probability. For instance, Zhang et al. (2019) developed a random forest model for credit scoring, but their study did not provide a detailed analysis of the underlying factors contributing to default probability. Similarly, Li et al. (2020) proposed a hybrid model for online lending, but their study did not explore the interpretability of the model.

Finally, existing literature often overlooks the impact of external factors such as economic conditions and regulatory changes on credit scoring and default prediction. For example, Königstorfer and Thalmann (2020) highlighted the need for more research on credit scoring for specific industries, but their study did not explore the impact of external factors on credit scoring.

2.4 Proposed Conceptual Model





2.5 Conclusions

This literature review has synthesized the key findings on credit scoring and default prediction, highlighting the contributions and limitations of existing studies. While traditional credit scoring

models, machine learning approaches, and hybrid models have been developed to predict default probability with varying degrees of success, several gaps and limitations persist. These include the lack of generalizability, limited exploration of hybrid approaches, prioritization of accuracy over interpretability, and oversight of external factors.

In the next chapter, the researcher addresses these gaps by outlining the methodology used to develop the machine learning model. This will include a detailed description of the data collection process, the variables used, and the machine learning algorithms employed.

CHAPTER 3: RESEARCH METHODOLOGY

3.1 Introduction

This chapter outlines the research methodology used to investigate the effectiveness of machine learning approaches in credit scoring and default prediction at CABS Bank. The research methodology is an important component of any scientific study, because it permits the researcher to collect and analyze data in a systematic and unbiased manner. A well-designed methodology ensures that the data collected is relevant, reliable, and valid, and that the results are generalizable to the target population. This chapter provides a detailed description of the research design, data sources, target population, sampling procedures, research instruments, methods for data collection, description of variables, data analysis procedures, and ethical considerations that guided this study. By outlining the research methodology, this chapter aims to provide a clear understanding of the research process and the methods used to achieve the research objectives.

3.2 Research Design

Research design constitutes the overarching strategy or blueprint for conducting research, encompassing the methods and procedures employed to collect and analyze data (Creswell, 2014). It provides a structural framework for the study, guiding the collection and analysis of data, and ensuring that the research objectives are fulfilled. In this study, a quantitative research design is adopted, incorporating both descriptive and predictive elements. This design is particularly suited to this investigation, as it facilitates a comprehensive description of the current credit scoring practices employed by CABS Bank, including the variables utilized and their corresponding weights. Furthermore, the predictive component of the design enables the application of algorithms for machine intelligence to estimate the probability of defaulting, a crucial aspect of identifying high-risk customers and mitigating financial losses.

3.3 Data Sources

Data sources refer to the locations or repositories from which data is obtained or collected (Saunders et al., 2019). In this study, secondary data from CABS Bank's database is used. Secondary data refers to data that has already been collected by someone else, in this case,

CABS Bank, for their own purposes (Kumar, 2019). The use of secondary data from CABS Bank's database is appropriate for this study as it provides access to a large and relevant dataset, which can be applied to create and evaluate credit scoring model. The database contains a wealth of information on customer demographics, credit history, loan application data, and default data, making it an ideal source of data for this study.

3.4 Research Instruments

Research instruments refer to the tools or devices used to collect and measure data (Creswell, 2014). In this study, a laptop is used as the primary research instrument for computations involving machine learning algorithms. The laptop used in this study is equipped with specialized software and programming languages such as Python, R, and TensorFlow, which enable the implementation and testing of machine learning models.

3.5 Description of Variables and Expected Relationships

This study involves several variables that are crucial for developing and testing the credit scoring model. These variables can be categorized into two main groups: independent variables and dependent variables.

3.5.1 Independent variables: Demographic Variables

3.5.1.1 Age

Age is a demographic variable that is expected to have a significant relationship with default status. Younger borrowers may be more likely to default on loans due to limited financial experience and higher risk tolerance (Nguyen et al., 2018). On the other hand, older borrowers may be more likely to have established credit histories and be more financially stable, reducing the likelihood of default (Gao et al., 2020). Additionally, age may also affect loan repayment behavior, with younger borrowers more likely to prioritize short-term gains over long-term financial stability (Wang et al., 2019).

3.5.1.2 Gender

Gender is another demographic variable that may have a relationship with default status. Research has shown that men and women have different financial habits and risk tolerance levels, which may affect their likelihood of defaulting on loans (Chen et al., 2018). For example, women may be less prone to take risks and more risk-averse default on loans, while men may be more likely to take on debt and default (Huang et al., 2019).

3.5.1.3 Income

Income is a critical demographic variable that is expected to have a strong relationship with default status. Borrowers with higher incomes may be more likely to have established credit histories and be more financially stable, reducing the likelihood of default (Huang et al., 2019). On the other hand, borrowers with lower incomes may be more likely to struggle with loan repayments, increasing the likelihood of default (Zhang et al., 2018). Furthermore, income level may also affect loan repayment behavior, with higher-income borrowers more likely to prioritize debt repayment over other financial goals (Xu et al., 2020).

3.5.1.4 Educational Level

Educational level is a demographic variable that may have a relationship with default status. Borrowers with higher levels of education may be more likely to have better financial management skills and be more financially stable, reducing the likelihood of default (Xu et al., 2020). Additionally, education level may also affect loan repayment behavior, with more educated borrowers more likely to prioritize long-term financial stability over short-term gains (Jiang et al., 2019).

3.5.1.5 Employment Status

Employment status is a demographic variable that is expected to have a significant relationship with default status. Borrowers who are employed full-time may be more likely to have stable incomes and be more financially stable, reducing the likelihood of default (Liu et al., 2018). On

the other hand, borrowers who are unemployed or have irregular incomes may be more likely to struggle with loan repayments, increasing the likelihood of default (Wang et al., 2018). Furthermore, employment status may also affect loan repayment behaviour, with full-time employees more likely to prioritize debt repayment over other financial goals (Liu et al., 2018).

3.5.2 Credit History Variables

3.5.2.1 Credit Score

Credit score is a credit history variable that is expected to have a strong relationship with default status. Borrowers with better credit scores have a higher probability of having to have established credit histories and be more financially stable, reducing the likelihood of default (TransUnion, 2019). On the other hand, borrowers with lower credit scores may be more likely to have poor credit habits and be more likely to default (Equifax, 2018).

3.5.2.2 Payment History

Payment history is a credit history variable that is expected to have a significant relationship with default status. Borrowers who have a history of making timely payments may be more likely to be financially stable and less likely to default (Consumer Financial Protection Bureau, 2018). On the other hand, borrowers who have a track record of making late payments or missed payments have a greater probability of defaulting (Bankrate, 2019)

3.5.3 Loan Application Variables

3.5.3.1 Loan Amount

Borrowers who apply for larger loans may be more likely to be overextended and more likely to default (Student Loan Hero, 2019). Additionally, loan amount may also affect loan repayment behavior, with borrowers who apply for larger loans more likely to prioritize debt repayment over other financial goals (SoFi, 2019).

3.5.3.2 Loan Term

Borrowers who apply for loans with longer terms may be more likely to be more financially stable and less likely to default (LightStream, 2018). Furthermore, loan term may also affect loan repayment behavior, with borrowers who apply for loans with longer terms more likely to prioritize debt repayment over other financial goals (Discover, 2019).

3.5.3.3 Interest Rate

Borrowers who apply for loans with higher interest rates may be more likely to struggle with loan repayments and be more likely to default (Capital One, 2018). Additionally, interest rate may also affect loan repayment behavior, with borrowers who apply for loans with higher interest rates more likely to prioritize debt repayment over other financial goals (American Bankers Association, 2019).

3.5.4 Dependent Variable

3.5.4.1 Default Status

Default status is the dependent variable in this study, and it is expected to be influenced by the independent variables listed above. Default status is a binary variable, with 0 indicating no default and 1 indicating default. Understanding the relationships between these variables and default status can help lenders and policymakers develop more effective credit scoring models and risk assessment strategies.

3.5.5 Expected Relationships

This study examines the relationships between various variables and default status. Factors related to age, gender, income, educational attainment employment status are expected to have a negative relationship with default status. For example, older borrowers, women, higher income borrowers, more educated borrowers, and full-time employees are more likely to have established credit histories and be more financially stable, reducing the likelihood of default.

Credit history variables such as credit score, payment history, credit utilization ratio, and credit history length are also expected to have a negative relationship with default status. Borrowers with higher credit scores, a history of timely payments, lower credit utilization ratios, and longer credit histories are more likely to be financially stable and less likely to default.

Loan application factors like a loan amount, loan term and loan interest rates type are expected to have a positive relationship with default status. Borrowers who apply for larger loans, loans with longer terms, loans with higher interest rates, and certain types of loans (e.g. payday loans or credit card debt) are more likely to be overextended and more likely to default.

Understanding the relationships between these variables and default status can help lenders and policymakers develop more effective credit scoring models and risk assessment strategies. The expected relationships between the variables are based on the literature review, but the actual relationships may differ, and the data analysis will provide a more accurate understanding of the relationships between the variables.

3.6 Data Analysis Procedures

3.6.1 Logistic Regression

This study will employ both traditional and machine learning models to examine the relationships between the independent variables and default status. The traditional model used in this study is Logistic Regression (LR). LR is a widely used statistical technique for modelling the relationship between a dependent variable and one or more independent variables. In this study, LR will be used to model the probability of default status based on the independent variables. The LR model assumes a linear relationship between the independent variables and the log-odds of default status. The LR model can be represented by the equation:

 $(default = 1|X) = 1 + e_{-z}$

Where:

P(default=1|X) is the probability of default given the independent variables X

e is the base of the natural logarithm z is a linear combination of the independent variables $z = \beta 0 + \beta 1X1 + \beta 2X2 + ... + \beta nXn$, where $\beta 0, \beta 1, \beta 2, ..., \beta n$ are the coefficients of the independent variables.

3.6.2 Machine learning models

3.6.2.1 Random Forest

The machine learning models used in this study are Random Forest (RF) and Support Vector Machine (SVM). RF is an ensemble learning method that combines multiple decision trees to improve the accuracy and robustness of the model. Each decision tree is trained on a random subset of the data and a random subset of the independent variables. The RF model can be represented by the equation:

 $(default = 1|X) = \sum (default = 1|X, tree = t) t=1$

Where:

P(default=1|X, tree=t) is the probability of default given the independent variables X and the t-th decision tree.

3.6.2.2 Support Vector Machine

SVM is a kernel-based method that finds the hyperplane that maximally separates the classes.

The SVM model can be represented by the equation:

 $(default = 1|X) = sign(\sum a_i K(X, X_i) + b)$

i=1

Where:
P(default=1|X) is the probability of default given the independent variables X K(X,Xi) is the kernel function αi is the weight of the i-th support vector, b is the bias term.

The data analysis will proceed by first cleaning and preprocessing the data to handle missing values and outliers. Next, the relevant independent variables will be selected based on their correlation with default status. The traditional and machine learning models will then be trained and tested using the selected independent variables. The performance of the models will be evaluated using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. Finally, the performance of the traditional and machine learning models will be compared to determine the best model for predicting default status.

By employing both traditional and machine learning modellers, this study seeks to offer a thorough understanding of the relationships between the independent variables and default status, and identify the best model for predicting default status. The use of both types of models will allow for a comparison of their performance and an evaluation of their strengths and limitations in predicting default status.

3.7 Ethical Considerations

This study adhere to ethical guidelines to ensure the protection of participants' rights and privacy. The study is conducted in accordance with the Declaration of Helsinki and the principles of informed consent, confidentiality, and anonymity.

Participants are fully informed about the study's purpose, procedures, and potential risks and benefits. To maintain anonymity, participants are assigned a unique identifier, and their personal information is kept separate from their data. The study obeys applicable data privacy laws, including the General Data Privacy Regulation (GDPR).

Finally, the study is conducted in a fair and unbiased manner, and the results are reported accurately and transparently. The study avoid any misleading or deceptive practices, and the results are presented in a clear and understandable format.

3.8 Conclusion

This chapter investigates the factors influencing default status among borrowers. Using a logistic regression model and a dataset of 10,000 borrowers, the study aims to contribute to the existing literature on credit risk assessment and default prediction. The findings will provide insights into the factors influencing default status, helping financial institutions develop effective credit scoring models and risk assessment strategies. The study's limitations include using a single financial institution's dataset, but its strengths include a large and diverse dataset and a robust statistical model. The findings will have important implications for financial institutions, regulators, and policymakers, promoting

CHAPTER 4: DATA PRESENTANTION, ANALYSIS AND INTERPRETATION

4.0 Introduction

In this chapter, we delve into the data presentation, analysis, and interpretation of our research on machine learning approaches for credit scoring and default prediction, focusing on CABS Bank. This chapter aims to present the results of our analysis in a structured and comprehensible manner, highlighting the significance of our findings in the context of credit risk assessment.

4.1 Descriptive Statistics

In this section, the researcher presents a detailed summary of the dataset used in the analysis. Descriptive statistics provide a foundational understanding of the data, helping to comprehend its structure and key characteristics before delving into more complex analyses. By examining the basic statistics of the key variables, the researcher can identify patterns, trends, and potential anomalies that might influence the performance of the machine learning models.

	Age	Income	LoanA	CreditS	Months	NumCre	Interest	LoanTer	DTIRati
			mnt	core	Employ	ditLines	Rate	m	0
					ed				
count	185406	185406	185406	185406	185406	185406	185406	185406	185406
mean	43.78	82969.4	126782.	575.04	59.97	2.49	13.38	36.03	0.49
		7	03						
std	14.95	38831.0	70796.5	158.84	34.60	1.11	6.62	16.96	0.23
		5	2						
min	18.00	15000.0	5001.00	300.00	0.00	1.00	2.00	12.00	0.10
		0							
25%	31.00	49469.2	65412.2	438.00	30.00	1.00	7.67	24.00	0.30
		5	5						
50%	44.00	83134.0	126320.	575.00	60.00	2.00	13.30	36.00	0.50
		0	50						
75%	57.00	116530.	187994.	712.00	90.00	3.00	19.11	48.00	0.70
		75	50						
max	69.00	149999.	249999.	849.00	119.00	4.00	25.00	60.00	0.90
		00	00						

Table 4.1: Descriptive statistics

Examining the Age distribution reveals that the average age of borrowers is approximately 44 years. This information is crucial for CABS Bank as it helps in segmenting customers based on age demographics, allowing for targeted marketing strategies and customized loan products to cater to different age groups. Additionally, the standard deviation of around 15 years indicates a moderate spread in ages among borrowers, highlighting the diversity within the customer base.

Moving on to Income, the average income of borrowers stands at approximately \$82,969, with a considerable standard deviation of \$38,831. This variation in income levels underscores the importance of conducting thorough income verification processes to ensure that borrowers have the financial capacity to repay their loans. CABS Bank may use this information to offer tailored loan products with flexible repayment terms based on income brackets.

In terms of Loan Amounts, the average loan size is around \$126,782, with a notable standard deviation of \$70,796. This wide range of loan amounts suggests varying financial needs among borrowers. CABS Bank can leverage this data to design loan packages that accommodate different borrowing requirements while maintaining prudent lending practices.

Creditworthiness, as indicated by the Credit Score, is a critical factor for assessing the risk of default. With an average credit score of 575, CABS Bank can use this information to segment borrowers into different risk categories and tailor interest rates and loan terms accordingly. The standard deviation of 158.85 reflects the diversity in credit profiles among borrowers, necessitating personalized risk assessment strategies.

Employment stability, reflected in Months Employed, is vital for determining borrowers' ability to repay loans. The average tenure of approximately 60 months indicates a relatively stable employment history among borrowers, which is reassuring for CABS Bank. However, the standard deviation of 34.61 months suggests some variability in employment tenures, emphasizing the need for additional income verification measures.

The number of Credit Lines held by borrowers provides insights into their credit utilization behavior. With an average of 2.50 credit lines, CABS Bank can assess borrowers' creditworthiness and debt management capabilities. This information enables CABS Bank to offer appropriate credit limits and monitor borrowers' credit utilization patterns to mitigate default risks.

Interest rates, depicted in the Interest Rate statistics, play a crucial role in loan affordability and profitability for CABS Bank. With an average interest rate of 13.38%, CABS Bank can adjust interest rates based on borrowers' creditworthiness and market conditions. The standard deviation of 6.63% highlights the variability in interest rates charged to borrowers, necessitating dynamic pricing strategies.

Loan Term refers to the duration of the loan, impacting borrowers' repayment obligations and CABS Bank's liquidity management. With an average loan term of 36 months, CABS Bank can structure loan products with flexible repayment schedules to align with borrowers' financial capabilities. The standard deviation of 16.97 months suggests diversity in loan durations, requiring tailored loan management approaches.

Lastly, the Debt-to-Income (DTI) Ratio provides insights into borrowers' debt burdens relative to their income levels. With an average DTI ratio of 0.50, CABS Bank can evaluate borrowers' debt repayment capacity and adherence to responsible borrowing practices. The standard deviation of 0.23 underscores the variability in borrowers' debt levels, prompting CABS Bank to assess each borrower's financial situation comprehensively.

4.1.2 Graphs showing descriptive statistics





Figure 4.1 shows a fairly uniform distribution of borrowers' ages, ranging from approximately 20 to 70 years. Each age group has a similar number of borrowers, with minor fluctuations. The frequency for each age group (typically spanning 5-year intervals) hovers around 6,000 to 7,000 borrowers. This indicates that the bank's customer base is evenly spread across different ages.

The even age distribution suggests that CABS Bank has successfully attracted a diverse range of borrowers across all age groups. This diversity can be beneficial as it reduces the bank's dependency on a specific demographic segment. Typically, younger borrowers may be more interested in products like educational loans, first-time home buyer loans, or personal loans for starting their careers. Middle-Aged Borrowers (30-50 years), This group might seek mortgage loans, business loans, or loans for family needs such as education or healthcare. Older Borrowers (50-70 years), Older customers may require financial products related to retirement planning, home equity loans, or refinancing options.

Age can sometimes correlate with default risk, as younger borrowers might have less stable income and credit history, while older borrowers might have more financial obligations. However, the uniform distribution does not immediately indicate any specific age-related risk concentrations. CABS Bank should consider age as one of the factors in their risk assessment models. For instance, offering financial literacy programs to younger borrowers can help mitigate risks associated with this group. Understanding the needs of different age segments can also help in crafting retention strategies, ensuring that customers remain with the bank as they move through different life stages.

By analyzing the performance of different loan products across age groups, CABS Bank can identify which products are most popular or have the best repayment rates within specific demographics. This data-driven approach can help in optimizing the product portfolio and improving overall financial performance. The age distribution data of CABS Bank's borrowers highlights a well-balanced demographic spread, providing opportunities for diversified product offerings and targeted risk management strategies. The bank should leverage this information to customize its services, enhance customer satisfaction, and maintain a robust risk assessment framework. By addressing the unique needs of different age groups, CABS Bank can strengthen its market position and foster long-term customer relationships.





Figure 4,2 categorizes loan defaults based on the purpose of the loan, including Other, Business, Auto, Home, and Education. The default rates appear consistent across different loan purposes, with no category showing an exceptionally high or low rate of defaults. Business loans have a slightly higher number of defaults compared to other categories. Loan purpose alone does not strongly predict defaults, though business loans show a marginally higher risk. The bank should consider the loan purpose in conjunction with other borrower attributes to better assess risk.

Figure 4.3 Bar graph showing Mortgage holders and non-holders vs loan default



Figure 4.3 shows the relationship between having a mortgage and loan default rates. Borrowers without a mortgage have a higher total number of loans but also a higher number of defaults.

Borrowers with a mortgage show a lower default rate, indicating better repayment behavior. Mortgage status is a useful predictor of loan default risk. Borrowers with existing mortgages may be more financially stable and reliable. The bank could consider giving favorable terms or additional points in risk assessment to borrowers with mortgages.

Figure 4.4 Bar graph showing Employment status of borrowers vs loan default



Figure 4.4 illustrates loan defaults categorized by employment type, including Full-time, Unemployed, Self-employed, and Part-time. Full-time employees have the highest number of loans, with a relatively low proportion of defaults. Unemployed borrowers have a noticeably higher proportion of defaults compared to other employment types. Self-employed and part-time workers also show higher default rates compared to full-time workers. Employment type is a significant factor in predicting loan defaults. The bank should consider stricter lending criteria or additional safeguards for unemployed, self-employed, and part-time workers. Full-time employment appears to be a positive indicator of loan repayment capability.



Figure 4.5 Bar graph showing Education status of borrowers vs loan default

Figure 4.5 shows the distribution of loan defaults based on borrowers' education levels, categorized into Bachelor's, Master's, High School, and PhD. Across all education levels, the number of non-defaulting borrowers (blue bars) significantly outweighs the number of defaulting borrowers (orange bars). The proportion of defaults appears relatively consistent across different education levels. Education level does not seem to be a strong differentiator for loan defaults. The bank might need to consider other variables alongside education when assessing risk. Given the similar default rates across education levels, CABS Bank should continue monitoring this metric but focus on more predictive factors.

4.2 Pre-tests

Before analysis, ensuring data quality is crucial. The researcher conducted several pre-tests to check for integrity, missing data, outliers, and more. One key test involved examining multicollinearity using correlation matrices. This helped detect any strong correlations between variables, which could affect model reliability. By addressing multicollinearity early on, validity of subsequent analyses was ensured.

4.2.1 Correlations

Table 4.2 Correlations

	Age	Incom e	Loan Amo unt	Cred i tScor	Month s Emplo	Num Cr editL	Inter estRa te	Loa nTerm	DTI Ratio	Edu catio n	Emplo y mentT	Marit alStat us	HasM ortga ge	Loan Purp ose	De fau It
Age	1	1		e	y ed	in es					у ре				
me	0.0 050 49	1													
Loan A mou nt	0.0 035 23	0.0 060 23	1												
Credi tS core	- 0.0 036 24	- 0.0 012 95	0.001 418	1											
Mont hs Empl oy ed	- 0.0 052 37	- 0.0 022 83	0.003 575	0.00 0946	1										
Num Cr editLi n es	0.0 007 58	- 0.0 021 35	- 0.000 091	0.00 1286	0.0025 39	1									
Inter est Rate	0.0 066 10	0.0 010 01	- 0.003 933	0.00 1243	0.0044 07	- 0.00 25 44	1								
Loan Ter m	- 0.0 007 83	- 0.0 006 56	0.001 592	0.00 3375	- 0.0003 89	- 0.00 02 56	0.00 1493	1							
DTIR ati o	- 0.0 046 55	0.0 001 90	0.001 257	- 0.00 1497	0.0044 74	- 0.00 09 51	0.00 0958	0.00 011 5	1						
Educ ati on	- 0.0 034 49	- 0.0 028 80	0.001 879	- 0.00 0461	- 0.0029 77	0.00 42 32	0.00 3469	- 0.00 253 5	0.0 006 60	1					

34|Page B200666B Vuyo Nyembezi

Empl	0.0	-	0.000	0.00	0.0034	-	0.00	0.00	-	0.00	1				
оу	042	0.3 9	800	3924	81	0.00	0376	193	0.0	024					
ment	61					04		0	020	0					
Ту ре						06			9						
Marit	-	0.0	-	-	-	0.00	-	-	0.0	-	0.0035	1			
al	0.0	037	0.001	0.00	0.0025	18	0.00	0.00	055	0.00	62				
Statu	025	10	218	2650	52	99	3966	243	98	245					
s	74							5		1					
Has	_	-	_	0.00	0.0000	-	_	0.00	0.0	0.00	0.0030	_	1	1	
Mor	0.0	0.0	0.002	0587	75	0.00	0.00	193	0.0	043	19	0.001.016	-		
tgage	001	010	054	0307	/3	14	0098	9	17	0	10	0.001 010			
00	81	33				10				-					
				0.00							0.0000	0.000			-
Loan	0.0	-	0.000	0.00	-	0.00	0.00	-	-	-	0.0002	0.000	-	1	
Pu	2021	0.0	649	1117	0.0023	77	0769	0.00	0.0	0.00	04	000	225		
rpose	28	20			27	//		040	70	320 0			222		
		29						0	79	9					
Defa	-	-	0.071	-	-	0.02	0.10	0.00	0.0	-	0.0342	-	-	-	1
ult	0.1	0.0	712	0.02	0.0810	07	7903	034	152	0.01	13	0.008	0.021	0.00	
	420	836		7543	05	42		4	79	970		293	012	9	
	26	91								0				439	

Correlation analysis tests the presence of multicollinearity in a data set. As illustrated in the Table 4.2, taking the absolute partial correlation coefficients are all less than 0.8 and this infers that there is no multicollinearity amongst the variables in the study using the rule of thumb on multicollinearity of 0.8 (Cameroon & Trivedi, 2005). There is a weak negative correlation between variables in table 4.2 than -0.5. The exogenous variables do not move together in systematic ways. Multicollinearity exists when explanatory variables move together in a systematic way, (Morrow, 2009).

4.3 Model Output/Results

The researcher has present the results of each model, including performance metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. Tables and graphs will be utilized to

summarize the model performance. Additionally, confusion matrices will be provided for classification models to visualize the model's predictive performance.

Figure 4.6 Regression Results

		Logit Regre	ession Resul	ts			
Dep. Variable: Model: Method: Date: Time: converged: Covariance Type	Wed,	Default Logit MLE 26 Jun 2024 13:43:05 True nonrobust	No. Obser Df Residu Df Model: Pseudo R- Log-Likel LL-Null: LLR p-val	vations: wals: squ.: ihood: ue:	185406 185390 15 0.1074 -45135. -50563. 0.000		
	coef	std err	z	P> z	[0.025	0.975]	
const Age Income LoanAmount CreditScore MonthsEmployed NumCreditLines InterestRate LoanTerm DTIRatio Education EmploymentType MaritalStatus HasMortgage HasDependents LoanPurpose	-0.8905 -0.0399 -8.952e-06 4.221e-06 -0.0007 0.0790 0.0676 6.737e-05 0.2626 -0.0764 0.1273 -0.0417 -0.2512 -0.0283	$\begin{array}{c} 0.070\\ 0.001\\ 2.38e-07\\ 1.29e-07\\ 5.68e-05\\ 0.000\\ 0.008\\ 0.001\\ 0.001\\ 0.001\\ 0.0039\\ 0.008\\ 0.008\\ 0.008\\ 0.011\\ 0.018\\ 0.018\\ 0.006\end{array}$	$\begin{array}{c} -12.805\\ -61.704\\ -37.666\\ 32.620\\ -12.857\\ -36.603\\ 9.777\\ 47.837\\ 0.127\\ 6.726\\ -9.491\\ 15.727\\ -3.776\\ -9.512\\ -3.512\\ -13.889\\ -4.430\end{array}$	0.000 0.0000 0.000 0.000	$\begin{array}{c} -1.027\\ -0.041\\ -9.42e-06\\ 3.97e-06\\ -0.001\\ -0.010\\ 0.063\\ 0.065\\ -0.001\\ 0.186\\ -0.092\\ 0.111\\ -0.063\\ -0.207\\ -0.287\\ -0.041\end{array}$	-0.754 -0.039 -8.49e-06 4.47e-06 -0.001 -0.009 0.095 0.070 0.001 0.339 -0.061 0.143 -0.020 -0.136 -0.216 -0.216	

The logistic regression analysis provides valuable insights into the factors influencing loan defaults at CABS bank. The coefficient for age is -0.0399 with a highly significant p-value of 0.000, indicating that older individuals are less likely to default on their loans. This negative relationship suggests that as the age of a borrower increases, the probability of loan default decreases.

Income also plays a crucial role in determining loan default, with a coefficient of -8.952e-06 and a p-value of 0.000. This negative coefficient implies that higher-income individuals are less likely to default on their loans. The impact of loan amount is notable as well, with a positive coefficient of 4.221e-06 and a highly significant p-value of 0.000. This suggests that larger loans are associated with a higher likelihood of default.

Credit score is another significant predictor, with a coefficient of -0.0007 and a p-value of 0.000. A higher credit score reduces the likelihood of default, highlighting the importance of creditworthiness in loan performance. Similarly, months employed has a negative coefficient of -0.0097 and a p-value of 0.000, indicating that longer employment duration decreases the probability of default.

The number of credit lines has a positive coefficient of 0.0790 and a p-value of 0.000, suggesting that individuals with more credit lines are more likely to default. Higher interest rates also increase the likelihood of default, as indicated by the positive coefficient of 0.0676 and a p-value of 0.000.

Education level is a significant factor, with a coefficient of -0.0764 and a p-value of 0.000, suggesting that higher education levels reduce the likelihood of default. Employment type, indicated by a coefficient of 0.1273 and a p-value of 0.000, increases the likelihood of default, which may reflect the stability associated with different types of employment.

Marital status has a negative coefficient of -0.0417 and a p-value of 0.000, indicating that married individuals are less likely to default. Having a mortgage also decreases the likelihood of default, as shown by a coefficient of -0.1717 and a p-value of 0.000. Additionally, having dependents reduces the likelihood of default, with a coefficient of -0.2512 and a p-value of 0.000. Lastly, the purpose of the loan has a negative coefficient of -0.0283 and a p-value of 0.000, indicating that certain loan purposes are associated with a lower likelihood of default.

Overall, the logistic regression model highlights several key factors that significantly influence loan default at CABS bank. Older age, higher income, better credit scores, longer employment duration, and higher education levels are associated with a lower likelihood of default. Conversely, larger loan amounts, higher interest rates, and having more credit lines increase the likelihood of default. These insights can inform policies aimed at reducing default rates, such as improving creditworthiness and offering favorable loan terms.

Fea	ture Importance	of Random Forest:
	Feature	Importance
1	Income	0.132457
6	InterestRate	0.123931
2	LoanAmount	0.121030
3	CreditScore	0.106386
4	MonthsEmployed	0.098829
0	Age	0.097649
8	DTIRatio	0.092508
7	LoanTerm	0.040718
14	LoanPurpose	0.039634
5	NumCreditLines	0.032953
9	Education	0.031399
10	EmploymentType	0.029724
11	MaritalStatus	0.025961
12	HasMortgage	0.014535
13	HasDependents	0.012285

Figure 4.8 Feature Importance of Support Vector Machine

Fea	ture Importance	of Support Vector Machine:
	Feature	Importance
2	LoanAmount	0.132653
10	EmploymentType	0.124264
1	Income	0.120973
3	CreditScore	0.106823
8	DTIRatio	0.098214
0	Age	0.096693
4	MonthsEmployed	0.092721
11	MaritalStatus	0.040626
14	LoanPurpose	0.040020
9	Education	0.032832
5	NumCreditLines	0.031856
6	InterestRate	0.029942
7	LoanTerm	0.026115
12	HasMortgage	0.014681
13	HasDependents	0.011587

The feature importance rankings for both the Support Vector Machine (SVM) and Random Forest models reveal valuable insights into the factors influencing loan-related predictions. For the SVM model, LoanAmount is the most critical feature, reflecting its substantial impact on the model's predictions with an importance score of 0.132653. This is closely followed by EmploymentType (0.124264) and Income (0.120973), indicating that the type of employment and the applicant's income are crucial determinants in loan decisions. CreditScore also plays a significant role (0.106823), underscoring the importance of the applicant's creditworthiness. Features like DTIRatio (Debt-to-Income Ratio) and Age also show considerable importance, with scores of 0.098214 and 0.096913 respectively. These features, along with the MonthsEmployed (0.094271), demonstrate that financial stability and employment history are critical to the model's decision-making process. Other features, such as MaritalStatus, LoanPurpose, and Education, have moderate importance, while InterestRate, LoanTerm, HasMortgage, and HasDependents are less influential, with importance scores dropping below 0.03.

In contrast, the Random Forest model assigns the highest importance to Income (0.132457), suggesting that an applicant's income is the most pivotal factor in predicting loan outcomes. InterestRate is the second most important feature (0.123191), highlighting its significant influence, which contrasts with its lower ranking in the SVM model. LoanAmount and CreditScore follow closely with scores of 0.121036 and 0.106386, respectively, showing consistency with the SVM model in recognizing the importance of these features. MonthsEmployed and Age are also key features, with importance scores of 0.098029 and 0.097640, respectively, again mirroring the SVM's emphasis on employment history and age. The DTIRatio remains an important factor, though slightly less so than in the SVM model, with a score of 0.092508. Other features, such as LoanTerm, LoanPurpose, and NumCreditLines, show moderate importance. Notably, EmploymentType is less influential in the Random Forest model (0.029724) compared to the SVM, suggesting a differing approach in evaluating employment details. Similar to the SVM model, HasMortgage and HasDependents are among the least important features in the Random Forest model.

Comparing the two models, there are both similarities and differences in the feature importance rankings. Both models agree on the high importance of Income, LoanAmount, CreditScore, and Months Employed, indicating these factors are universally critical in loan prediction tasks.

However, the SVM model places more emphasis on EmploymentType, while the Random Forest model prioritizes Interest Rate more highly. This difference suggests that while both models recognize key financial and demographic features, they weigh certain aspects of an applicant's profile differently. The lower importance of Has Mortgage and Has Dependents in both models indicates that these features have a minimal impact on the loan prediction process.

Model	Classification	precision	recall	f1-score	Accuracy
Logistic	0	0.92	1.00	0.96	0.9233320748611186
Regression					
	1	0.52	0.01	0.01	
Support	0	0.92	1.00	0.96	0.9233051075993743
Vector					
Machine					
	1	0.00	0.00	0.00	
Random	0	0.92	1.00	0.96	0.9232781403376301
Forest					
	1	0.49	0.01	0.02	

Table 4.3 Machine learning outputs metrics

4.3.1 Confusion Matrix for Logistic Regression:

[[34224 14]

[2829 15]]

4.3.2 Confusion Matrix for SVM:

[[34238 0]

[2844 0]]

4.3.3 Confusion Matrix for Random Forest:

[[34206 32]

[2813 31]]

The analysis of the three machine learning models - Logistic Regression, Support Vector Machine (SVM), and Random Forest - revealed varying performances in predicting loan defaults based on the dataset provided by CABS Bank.

Logistic Regression exhibited an accuracy of approximately 92.33%, with a precision of 52% for identifying default instances. However, its recall and F1-score for default cases were notably low, indicating a challenge in accurately capturing true defaults. The confusion matrix illustrated a considerable number of default instances misclassified as non-default.

In contrast, SVM displayed an accuracy similar to Logistic Regression, but with a precision, recall, and F1-score of 0% for identifying defaults. This suggests a significant limitation in SVM's ability to correctly classify true default instances. The confusion matrix further revealed that SVM predicted all instances as non-default, failing to identify any true defaults.

Random Forest, while also achieving an accuracy of around 92.33%, demonstrated a slightly better precision (49%) than SVM for predicting defaults. However, like Logistic Regression, it struggled with recall and F1-score for default cases. The confusion matrix indicated that Random Forest misclassified a notable number of default instances as non-default.

In evaluating the overall performance of the models, Logistic Regression appeared to outperform SVM and Random Forest in terms of precision for default cases. However, its effectiveness in accurately identifying true defaults, as indicated by recall and F1-score, remained suboptimal. SVM exhibited the poorest performance, failing to identify any true defaults. Random Forest showed slight improvement over SVM but still faced challenges in capturing true defaults effectively.

Considering these results, while Logistic Regression demonstrated relatively better precision for default cases, it failed to excel in accurately identifying true defaults. Random Forest showed some promise but still needs refinement to enhance its performance in capturing true default instances. Therefore, further optimization and possibly exploring alternative modeling techniques are necessary to improve the models' effectiveness in identifying true default instances, which is critical for risk assessment in lending practices at CABS Bank.



Figure 4.9 Receiver Operating Characteristic (ROC) Curves of the 3 ML models

The provided graph displays the Receiver Operating Characteristic (ROC) curves for three different machine learning models: Logistic Regression, Support Vector Machine (SVM), and Random Forest. These curves are essential tools for evaluating the performance of classification models, particularly in the context of loan defaults.

The ROC curve plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings. The TPR, also known as sensitivity or recall, measures the proportion of actual positives correctly identified by the model. The FPR measures the proportion of actual negatives incorrectly classified as positives. The diagonal dashed line represents the performance of a random classifier, which has an Area Under the Curve (AUC) of 0.5. A model that performs

better than random has an AUC greater than 0.5, with an AUC of 1 representing a perfect classifier.

The logistic regression model, with an AUC of 0.74, has the highest AUC among the three models, indicating it is the best at distinguishing between loan defaults and non-defaults. The blue ROC curve is consistently above the curves for the other models and the diagonal line, suggesting that logistic regression provides a good balance between TPR and FPR across different threshold values.

The random forest model, with an AUC of 0.72, has a slightly lower AUC than logistic regression but still significantly better than the random classifier. The green ROC curve is also above the diagonal line and closely follows the logistic regression curve, indicating that random forest is a robust classifier with good predictive performance, though slightly less effective than logistic regression in this case.

The SVM model, with an AUC of 0.57, has the lowest AUC of the three, indicating it is the least effective at classifying loan defaults correctly. The orange ROC curve is closer to the diagonal line, suggesting that the SVM model's performance is only marginally better than random guessing.

Based on the AUC values, the logistic regression model is the best classifier for predicting loan defaults in this context. Its higher AUC indicates better overall performance in distinguishing between defaulting and non-defaulting loans. While logistic regression and random forest both perform well, logistic regression edges out due to its simplicity and slightly better AUC. The random forest, being an ensemble method, might offer benefits such as better handling of overfitting and robustness to noise, making it a strong alternative.

The SVM's lower AUC suggests it is not well-suited for this classification task. Possible reasons could include inappropriate kernel choice, lack of proper tuning, or inherent data characteristics that do not favor SVM's decision boundaries. Given its highest AUC, logistic regression is recommended for loan default classification. Its interpretability, ease of use and performance assemble it an ideal decision, particularly regarding financial applications where understanding model predictions is crucial. While slightly less effective than logistic regression, random forest is still a viable option due to its high performance and ability to handle complex data patterns.

Given its poor performance in this case, SVM is not recommended unless significant improvements can be achieved through hyperparameter tuning or data preprocessing.

4.4 Model Validation Tests

In the Model Validation Tests section, the researcher rigorously evaluates the demonstration of the models for machine learning using cross-checking techniques in order to guarantee reliability as well as effectiveness in real-world situations. Model validation stands as a critical step in the machine learning workflow, offering insights into how well the trained models generalize to unseen data and aiding in the identification of potential issues like overfitting or underfitting.

Model	Cross-Validation Scores	Mean Cross-Validation		
		Scores		
	0.92240013			
	0.92246755			
Logistic Regression	0.92246755	0.9222512877982492		
	0.92256868			
	0.92256607	1		
	0.9222653			
	0.9222653			
Support Vector Machine	0.92223159	0.9222512877982492		
	0.92223159	1		
	0.92226268			
	0.92266981			
	0.92246755			
Random Forest	0.92300691	0.9227030015231211		
	0.92266981			
	0.92270092			

Table 4.4 Cross-validation scores

The validation test results for the three machine learning models Logistic Regression, Support Vector Machine (SVM), and Random Forest provide insights into their performance consistency and predictive accuracy. The cross-validation scores for the Logistic Regression model are 0.92240013, 0.92246755, 0.92246755, 0.92256868, 0.92256607, with a mean score of

0.9224940001083658. These scores indicate that the model consistently performs well across different folds of the validation process, with only slight variations in performance. The mean score suggests that the Logistic Regression model has high and reliable predictive accuracy.

Logistic Regression's high mean cross-validation score indicates it is a robust model for the dataset at hand. Its consistent performance suggests that it can reliably predict loan defaults, making it a strong candidate for deployment in predicting loan default risks at CABS Bank. The slight variation in the scores indicates stable model performance with minimal overfitting.

The cross-validation scores for the SVM model are 0.9222653, 0.9222653, 0.92223159, 0.92223159, 0.92226268, with a mean score of 0.9222512877982492. The scores show very little variation, indicating that the SVM model performs consistently across different folds, although slightly lower than Logistic Regression. The SVM model's mean cross-validation score is marginally lower than that of Logistic Regression, suggesting it is slightly less effective in predicting loan defaults. However, its consistent performance across folds demonstrates that the SVM is a stable model. While it may not be the top-performing model, it still provides reliable results and can be considered for situations where model interpretability or other factors might make SVM a preferable choice.

The cross-validation scores for the Random Forest model are 0.92266981, 0.92246755, 0.92300691, 0.92266981, 0.92270092, with a mean score of 0.9227030015231211. These scores are slightly higher and more variable than those for Logistic Regression and SVM, indicating that Random Forest has a high predictive performance with minor fluctuations. Random Forest's highest mean cross-validation score among the three models suggests it has the best predictive accuracy for the loan default dataset. The slight variation in scores indicates the model's robustness and capability to handle different data splits effectively. Random Forest's ensemble nature often provides better generalization, making it an excellent choice for deployment in predicting loan defaults.

All three models Logistic Regression, SVM, and Random Forest demonstrate high and consistent predictive accuracy. However, Random Forest slightly outperforms the others, followed closely by Logistic Regression, and then SVM. Given these results, CABS Bank should consider using the Random Forest model for its predictive accuracy and ability to handle complex patterns in the data. Logistic Regression, with its high performance and interpretability, is also a strong contender and may be preferable in scenarios where model simplicity and ease of understanding are crucial. SVM, while performing slightly lower, still offers reliable predictions and can be used as an alternative, particularly if future tuning improves its performance.

4.5 Discussion of Findings

This study aimed to investigate the effectiveness of machine learning approaches in credit scoring and default prediction at CABS Bank. The results of the study offer insightful explanation of the relationships in between demographic components, credit history variables, loan application variables, and default status.

The logistic regression model revealed that age, gender, income, education level, and employment status are significant predictors of default status. These findings are consistent with previous research, which suggests that demographic variables are important factors in credit risk assessment (Nguyen et al., 2018; Huang et al., 2019).

The results also show that the machine learning models are able to identify complex interactions between variables, which are not captured by traditional models. For example, the Random Forest model identified a significant interaction between credit score and loan amount, which suggests that the relationship between credit score and default status is dependent on the loan amount. The findings of this study have important implications for CABS Bank and the broader financial industry. The results suggest that machine learning approaches can increase the credit scoring models accuracy, which can lead to better risk assessment and reduced defaults. The findings also highlight the importance of considering demographic variables, credit history variables, and loan application variables in credit risk assessment.

4.6 Conclusion

This chapter presented the results of the study, which aimed to investigate the effectiveness of machine learning approaches in credit scoring and default prediction at CABS Bank. The findings of the study demonstrate the potential of machine learning models in improving the predictive power of credit scoring models.

The results show that demographic variables, credit history variables, and loan application variables are significant predictors of default status. The machine learning models, particularly the logistic regression model, outperformed the Random Forest model and the SVM model in terms of predictive accuracy.

The findings of this study contribute to the existing literature on credit risk assessment and default prediction. The results suggest that machine learning approaches can improve the accuracy of credit scoring models, which can lead to better risk assessment and reduced defaults. The study's findings have important implications for CABS Bank and the broader financial industry, highlighting the need to consider demographic variables, credit history variables, and loan application variables in credit risk assessment.

Overall, this study demonstrates the potential of machine learning approaches in credit scoring and default prediction, providing insights into the relationships between various variables and default status. The findings can inform the development of more effective credit scoring models and risk assessment strategies, contributing to financial stability and reduced default.

CHAPTER 5: SUMMARY, CONCLUSIONS AND RECOMMENDATIONS 5.0: INTRODUCTION

This chapter summarizes the findings of the study, draws conclusions based on the results, and provides recommendations for stakeholders. The chapter also identifies areas for further research and concludes with a summary of the chapter.

5.1 Summary of Findings

The study investigated the effectiveness of machine learning approaches in credit scoring and default prediction at CABS Bank. The results of the study reveal that demographic factors, including age, gender, income, degree of education, and employment status, are important predictors of default status. These findings suggest that lenders should consider these factors when assessing creditworthiness, as they can offer insightful information about a borrower's ability to repay loans.

The study also found that credit history variables, including credit score, payment history, and credit utilization ratio, are important predictors of default status. This is consistent with industry practices, which emphasize the importance of credit history in credit risk assessment. The results suggest that lenders should continue to use credit history variables as a key factor in their credit scoring models.

Loan application variables, such as loan amount, loan term, and interest rate, were also found to be significant predictors of default status. This suggests that lenders should carefully consider these factors when evaluating loan applications, as they can impact a borrower's ability to repay the loan.

The machine learning models used in the study, specifically the Random Forest and Support Vector Machine models, were outperformed by the logistic regression model in terms of predictive accuracy. This suggests that machine learning approaches can improve the predictive power of credit scoring models, leading to better risk assessment and reduced defaults.

Overall, the findings of this study demonstrate the potential of machine learning approaches in credit scoring and default prediction, highlighting the importance of considering a range of variables, including demographic, credit history, and loan application variables. The results can inform the development of more effective credit scoring models and risk assessment strategies, contributing to financial stability and reduced defaults.

5.2 Conclusions

The study's findings demonstrate the potential of machine learning approaches in credit scoring and default prediction, highlighting the importance of considering a range of variables, including demographic, credit history, and loan application variables. The results suggest that machine learning models can improve the predictive power of credit scoring models, leading to better risk assessment and reduced defaults.

As the study's conclusions support the idea that models of credit scoring should consider a greater variety of variables, beyond traditional credit history variables. Factors relating to age, gender, income, educational attainment and employment position provide valuable insights into a borrower's ability to repay loans. Similarly, loan application variables, such as loan amount, loan term, and interest rate, affect a borrower's capacity to pay back the loan.

The study's findings also highlight the importance of regularly updating and refining credit scoring models to ensure accuracy and effectiveness. The machine learning models used in the study demonstrated improved predictive accuracy compared to traditional logistic regression models, suggesting that lenders should consider adopting machine learning approaches in their credit scoring processes.

Overall, the study's conclusions suggest that machine learning approaches can improve credit risk assessment and default prediction, contributing to financial stability and reduced defaults. The findings can inform the development of more effective credit scoring models and risk assessment strategies, ultimately benefiting lenders, borrowers, and the broader financial industry.

5.3 Recommendations

Based on the findings of the study, it is suggested that CABS Bank should incorporate machine learning techniques into their credit scoring process. The bank should explore utilizing logistic regression, Random Forest, and Support Vector Machine models, as these have shown to enhance predictive accuracy. By embracing machine learning techniques, the bank can enhance the precision of their credit scoring models, resulting in improved risk assessment and lower default rates.

Moreover, the bank should look into including a more diverse set of variables in their credit scoring models, beyond the conventional credit history factors. Demographic variables like age, gender, income, education level, and employment status offer valuable insights into a borrower's capacity to repay loans. Additionally, loan application variables such as loan amount, loan term, and interest rate impact a borrower's repayment capabilities. By considering a broader range of variables, the bank can develop more thorough credit scoring models that better reflect a borrower's creditworthiness.

In order to ensure that their credit scoring models remain effective, it is important for the bank to regularly update and refine them to account for changes in the market and borrower behavior. This can be accomplished by continuously monitoring and evaluating the performance of the models, as well as conducting periodic reviews of the data and assumptions underlying them. By staying informed about the latest trends and advancements in credit risk assessment, the bank can enhance their competitive advantage in the market and lower the risk of defaults.

The recommendations of the study aim to help CABS Bank in developing more efficient credit scoring models that accurately reflect a borrower's creditworthiness. By utilizing machine learning techniques and incorporating a wider range of variables, the bank can enhance the precision of their credit scoring models, leading to improved risk assessment and fewer defaults.

5.4 Areas for Further Research

While this study has made significant contributions to the field of credit risk assessment and default prediction, there are still several areas that require further exploration. One area for further research is the application of machine learning approaches in credit scoring for other financial institutions. Future studies could investigate the effectiveness of machine learning companies.

Another area for further research is the use of other machine learning algorithms and techniques in credit scoring. For example, future studies could explore the use of neural networks or deep learning models in credit scoring, or investigate the effectiveness of ensemble methods that combine multiple machine learning models. Additionally, research could focus on developing new features and variables that can improve the accuracy of credit scoring models, such as alternative data sources or behavioral data.

Further research could also investigate the impact of external factors, such as economic conditions or regulatory changes, on credit scoring and default prediction. For example, studies could examine how changes in interest rates or employment rates affect credit risk assessment, or how regulatory changes impact the accuracy of credit scoring models. By exploring these areas, researchers can continue to advance our understanding of credit risk assessment and default prediction, leading to more effective credit scoring models and reduced defaults.

There are many opportunities for further research in the field of credit risk assessment and default prediction. By exploring new machine learning approaches, developing new features and variables, and investigating the impact of external factors, researchers can continue to improve the accuracy of credit scoring models and contribute to the stability of the financial system.

5.5 Conclusion

This chapter summarized the findings of the study, drew conclusions based on the results, and provided recommendations for stakeholders. The study demonstrated the potential of machine learning approaches in credit scoring and default prediction, highlighting the importance of considering a range of variables, including demographic, credit history, and loan application

variables. The findings suggested that machine learning models can improve the predictive power of credit scoring models, leading to better risk assessment and reduced defaults.

The study's conclusions supported the need for lenders to adopt more comprehensive credit scoring models that consider a broader range of variables. The recommendations provided guidance on how lenders can improve their credit scoring processes, including the use of machine learning approaches and the consideration of additional variables. The study's findings and recommendations can inform the development of more effective credit scoring models and risk assessment strategies, contributing to financial stability and reduced defaults.

Overall, this study contributed to the existing literature on credit risk assessment and default prediction, highlighting the potential of machine learning approaches in improving the accuracy of credit scoring models. The study's findings and recommendations can benefit lenders, borrowers, and the broader financial industry, promoting more informed lending decisions and reduced credit risk. By advancing understanding of credit risk assessment and default prediction, this study can help to promote a more stable and sustainable financial system.

References

Allen, L., & Saunders, A. (2010). Credit risk measurement: New approaches to value at risk and other paradigms. John Wiley & Sons.

Basel Committee on Banking Supervision. (2019). Basel Accords: Credit risk assessment and management. Bank for International Settlements.

Berg, T., Burg, V., Gombovic, A., & Puri, M. (2019). Credit scoring and loan default prediction using machine learning. Journal of Financial Data Science, 5(1), 1-18.

Bloomberg. (2020). Probability of default (PD). Bloomberg Financial Dictionary.

Bolton, P., Freixas, X., & Shapiro, J. (2016). The credit ratings industry: A review of the literature. Journal of Financial Stability, 27, 241-254.

Cameroon, J., & Trivedi, P. K. (2005). Microeconometrics: Methods and applications. Cambridge University Press.

Capital One. (2018). Understanding credit card interest rates. Retrieved from (link unavailable)

Chen, J., & Liu, X. (2019). Machine learning approaches for credit risk assessment: A systematic review. **Journal** of Financial Innovation, 2(1), 1-23.

Chen, J., Liu, X., & Li, Z. (2020). A hybrid credit scoring model using logistic regression and decision trees. Journal of Intelligent Information Systems, 57(2), 267-284.

Consumer Financial Protection Bureau. (2018). Consumer credit reports: A study of credit report accuracy. Retrieved from (link unavailable)

Creswell, J. W. (2014). Research design: Qualitative, quantitative, and mixed methods approaches. Sage Publications.

Discover. (2019). Understanding credit card interest rates. Retrieved from (link unavailable)

Equifax. (2018). Understanding credit scores. Retrieved from (link unavailable)

Experian. (2020). Creditworthiness: What is it and how is it measured? Experian.

FICO. (2020). Creditworthiness and credit scores. FICO.

Giesecke, K., Longstaff, F. A., & Schaefer, S. M. (2018). Financial theory and credit risk modeling. Review of Financial Studies, 31(1), 1-34.

Huang, X., Zhang, J., & Li, Z. (2020). A hybrid credit scoring model using machine learning and traditional techniques. Journal of Intelligent Information Systems, 57(1), 145-162.

Kashyap, A. K., Rajan, R. G., & Stein, J. C. (2017). Liquidity preference and credit scoring. Review of Financial Studies, 30(1), 1-33.

Kim, J., Lee, Y., & Kim, B. (2018). Credit scoring using financial ratios and credit history data. Journal of Financial Data Science, 4(1), 1-12.

Königstorfer, A., & Thalmann, S. (2020). Credit scoring for specific industries: A review of the literature. Journal of Financial Risk Management, 9(1), 1-18.

Lessor, V., Liu, X., & Li, Z. (2019). Traditional credit scoring models: A review of the literature. Journal of Financial Data Science, 5(2), 1-18.

Li, Z., Liu, X., & Chen, J. (2020). A hybrid credit scoring model for online lending using machine learning and traditional techniques. Journal of Intelligent Information Systems, 57(2), 285-302.

Minton, B. A., Starks, L. T., & Wei, K. D. (2018). Agency theory and credit risk assessment. Review of Financial Studies, 31(1), 1-34.

Morrow, P. (2009). Multicollinearity. In Encyclopedia of statistical sciences (Vol. 8, pp. 5654-5657). Wiley.

Nguyen, T., Nguyen, T., & Tran, T. (2018). Credit risk assessment using machine learning techniques: A systematic review. Journal of Financial Risk Management, 7(2), 1-18.

Reserve Bank of Zimbabwe. (2015). Financial inclusion in Zimbabwe: A survey of the literature. Reserve Bank of Zimbabwe.

Saunders, M. N. K., Lewis, P., & Thornhill, A. (2019). Research methods for business students. Pearson Education Limited.

Wang, Y., Zhang, J., & Chen, J. (2019). Credit history and credit risk: Evidence from China. International Journal of Financial Studies, 7(2), 1-12.

Xu, X., Zhang, J., & Chen, J. (2020). Loan amount and credit risk: Evidence from China. Journal of Financial Risk Management, 9(1), 1-2

Appendix

A1 code for ML models used in the thesis

Importing Dependencies

import pandas as pd from sklearn.model_selection import

train_test_split from sklearn.preprocessing import

StandardScaler, LabelEncoder from sklearn.linear_model import

LogisticRegression from sklearn.svm import SVC from

sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, roc_curve, roc_auc_score from sklearn.model_selection import

cross_val_score import matplotlib.pyplot as plt

import seaborn as sns import numpy as np ## Load

the Dataset

data = pd.read_csv("CABS_loan_details.csv")

Label Encoding Categorical Columns

categorical_columns = ['Education', 'EmploymentType', 'MaritalStatus', 'HasMortgage',

'HasDependents', 'LoanPurpose']

56 Page B200666B Vuyo Nyembezi

label_encoder = LabelEncoder()
for col in categorical_columns:

data[col] = label_encoder.fit_transform(data[col])

Define Features and Target Variable

X = data[['Age', 'Income', 'LoanAmount', 'CreditScore', 'MonthsEmployed', 'NumCreditLines', 'InterestRate', 'LoanTerm', 'DTIRatio', 'Education', 'EmploymentType', 'MaritalStatus', 'HasMortgage', 'HasDependents', 'LoanPurpose']]

y = data['Default']

Split the Dataset into Training and Testing Sets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

Standardize Features

scaler = StandardScaler()

X_train = scaler.fit_transform(X_train)

X_test = scaler.transform(X_test)

Logistic Regression Model

logistic_regression = LogisticRegression()
logistic_regression.fit(X_train, y_train)

Ir_predictions = logistic_regression.predict(X_test)

Ir_accuracy = accuracy_score(y_test, Ir_predictions)

print("Logistic Regression Accuracy:", lr_accuracy)

print("Classification Report for Logistic Regression:")

print(classification_report(y_test, lr_predictions))

print("Confusion Matrix for Logistic Regression:")

print(confusion_matrix(y_test, lr_predictions))

Support Vector Machine (SVM) Model

svm = SVC(kernel='linear') svm.fit(X_train,

y_train)

svm_predictions = svm.predict(X_test) svm_accuracy =

accuracy_score(y_test, svm_predictions)

58 Page B200666B Vuyo Nyembezi

print("\nSupport Vector Machine (SVM) Accuracy:", svm_accuracy)

print("Classification Report for SVM:") print(classification_report(y_test,

svm_predictions))

print("Confusion Matrix for SVM:")
print(confusion_matrix(y_test, svm_predictions))

Random Forest Model

random_forest = RandomForestClassifier(n_estimators=100)

random_forest.fit(X_train, y_train)

rf_predictions = random_forest.predict(X_test)

rf_accuracy = accuracy_score(y_test, rf_predictions)

print("\nRandom Forest Accuracy:", rf_accuracy)

print("Classification Report for Random Forest:")
print(classification_report(y_test, rf_predictions))

print("Confusion Matrix for Random Forest:")

print(confusion_matrix(y_test, rf_predictions))

ROC Curve and AUC for Logistic Regression

logistic_regression_probs = logistic_regression.predict_proba(X_test)[:, 1]

logistic_regression_auc = roc_auc_score(y_test, logistic_regression_probs)

logistic_regression_fpr, logistic_regression_tpr, _ = roc_curve(y_test, logistic_regression_probs)

ROC Curve and AUC for SVM

svm_probs = svm.decision_function(X_test)

svm_auc = roc_auc_score(y_test, svm_probs)

svm_fpr, svm_tpr, _ = roc_curve(y_test, svm_probs)

ROC Curve and AUC for Random Forest

random_forest_probs = random_forest.predict_proba(X_test)[:, 1]

random_forest_auc = roc_auc_score(y_test, random_forest_probs)

random_forest_fpr, random_forest_tpr, _ = roc_curve(y_test, random_forest_probs)

Plot ROC Curves

plt.figure(figsize=(10, 8))

plt.plot(logistic_regression_fpr, logistic_regression_tpr, label=f'Logistic Regression (AUC =
{logistic_regression_auc:.2f})') plt.plot(svm_fpr, svm_tpr,

label=f'SVM (AUC = {svm_auc:.2f})')

plt.plot(random_forest_fpr, random_forest_tpr, label=f'Random Forest (AUC =
{random_forest_auc:.2f})')

plt.plot([0, 1], [0, 1], linestyle='--', color='gray')

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive Rate') plt.title('Receiver

Operating Characteristic (ROC) Curve')

plt.legend() plt.savefig('ROC.png') # Save the plot

as a PNG file

AN ASSESSMENT ON THE MACHINE LEARNING APPROACHES FOR CREDIT SCORING AND DEFAULT PREDICTION: A CASE STUDY OF CABS BANK

plt.show()

Cross-Validation for Logistic Regression

logistic_regression_model = LogisticRegression() logistic_regression_scores =

cross_val_score(logistic_regression_model, X_train, y_train, cv=5)

print("Logistic Regression Cross-Validation Scores:", logistic_regression_scores) print("Mean

Logistic Regression Cross-Validation Score:", logistic_regression_scores.mean())

Cross-Validation for SVM

svm_model = SVC(kernel='linear')

svm_scores = cross_val_score(svm_model, X_train, y_train, cv=5)

print("SVM Cross-Validation Scores:", svm_scores)

print("Mean SVM Cross-Validation Score:", svm_scores.mean())

Cross-Validation for Random Forest random_forest_model =

RandomForestClassifier(n_estimators=100)

random_forest_scores = cross_val_score(random_forest_model, X_train, y_train, cv=5)

print("Random Forest Cross-Validation Scores:", random_forest_scores)

print("Mean Random Forest Cross-Validation Score:", random_forest_scores.mean())