

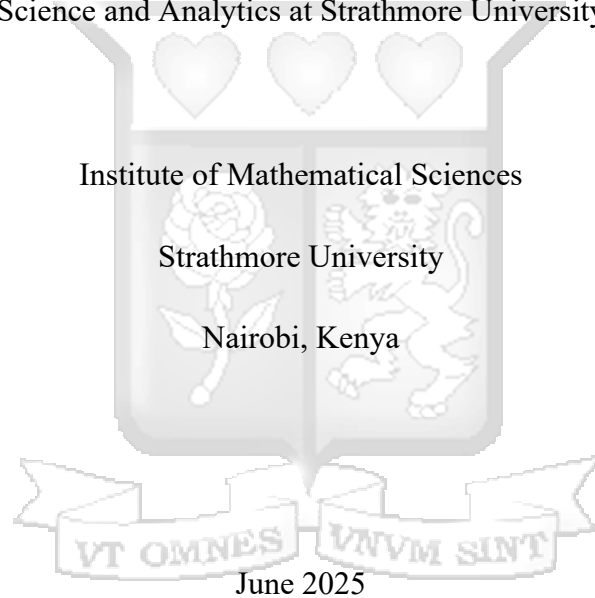
# Complex and Explainable Machine Learning Models in Credit Scoring

By

Aswani Allan Ademba

ADM 145613

Submitted in Partial fulfilment of the Requirements for the Degree of Master of Science in Data  
Science and Analytics at Strathmore University



This dissertation is available for Library use on the understanding that it is copyright material and that no quotation from the dissertation may be published without proper acknowledgement.

## Declaration and Approval

### Declaration

I declare that this work has not been previously submitted and approved for the award of a degree by this or any other University. To the best of my knowledge and belief, the dissertation contains no material previously published or written by another person except where due reference is made in the dissertation itself.

© No part of this dissertation may be reproduced without the permission of the author and Strathmore University

Name of the Candidate: Aswani Allan Ademba

Sign:  .....Date: 18<sup>th</sup> May 2025.

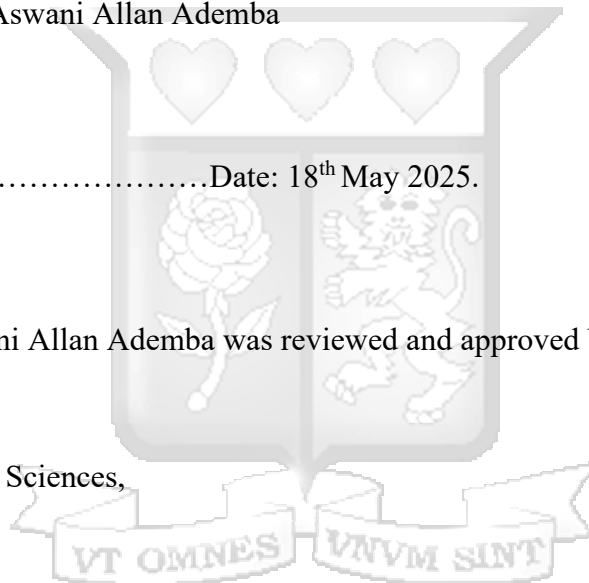
### Approval

The dissertation of Aswani Allan Ademba was reviewed and approved by the following:

Dr John Olukuru,  
Institute of Mathematical Sciences,  
Strathmore University

Dr. Godfrey Madigu,  
Dean, Institute of Mathematical Sciences,  
Strathmore University

Prof. Bernard Shibwabo,  
Director of Graduate Studies,  
Strathmore University



## Abstract

The rapid evolution of state-of-the-art modelling methodologies offers a compelling opportunity to enhance the precision of credit evaluation tools. However, this progress often comes at a cost: the trade-off between model transparency and predictive accuracy. For credit managers tasked with maintaining effective oversight of credit risk and central bank regulators seeking assurance in model integrity, this trade-off can be excessively burdensome. As a result, the adoption of sophisticated models is often hindered by their inherent lack of transparency.

This dissertation addresses this challenge by exploring advancements in credit assessment methodologies. It provides a comprehensive evaluation of predictive techniques, ranging from traditional logistic regression to modern artificial intelligence (AI) approaches. The findings demonstrate that complex tree-based algorithms, such as random forests, gradient-boosted trees, and extreme gradient-boosted trees, exhibit superior predictive accuracy in forecasting customer defaults. However, the dissertation goes beyond mere performance metrics by introducing innovative approaches to enhance the interpretability and practicality of these advanced models for credit risk practitioners. By doing so, it addresses a significant barrier to the widespread adoption of complex, opaque models in the financial industry.

The study leverages a substantial dataset obtained from a financial institution, ensuring the reliability of inputs and the robustness of outputs. A key contribution of this dissertation lies in its integration of Explainable Artificial Intelligence (XAI) methodologies, which bridge the gap between predictive power and model transparency. By making AI-driven credit risk models more interpretable, this research provides actionable insights for credit managers and regulators, fostering greater confidence in the use of advanced modelling techniques.

The key Contributions are: Comprehensive Evaluation of Predictive Models: A thorough comparison of traditional and modern credit scoring models, highlighting the strengths and limitations of each approach. Enhanced Interpretability of Complex Models: Introduction of techniques to improve the transparency of tree-based algorithms, making them more accessible to credit risk practitioners. Integration of Explainable AI (XAI): Application of XAI methodologies to credit risk management, enabling stakeholders to understand and trust AI-driven decisions. Practical Insights for Industry Adoption: Recommendations for implementing advanced models in

real-world credit risk management, balancing accuracy with regulatory and operational requirements.

In conclusion, this dissertation underscores the importance of balancing predictive accuracy with transparency in credit scoring models. By advancing the interpretability of sophisticated AI techniques, it paves the way for their broader adoption in the financial industry. The integration of XAI methodologies not only enhances model performance but also ensures that credit managers and regulators can maintain effective oversight, ultimately contributing to more robust and reliable credit risk management practices.



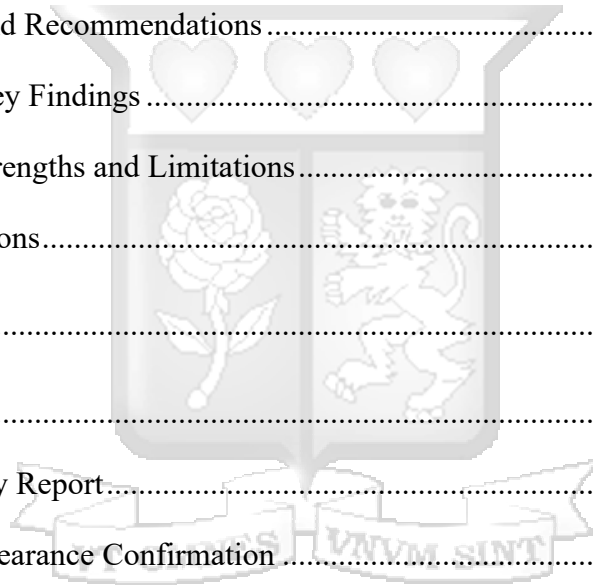
# Table of Contents

Declaration and Approval .....	ii
Abstract .....	iii
Table of Contents .....	v
List of Figures .....	ix
List of Tables .....	x
List of Abbreviations .....	xi
Acknowledgements .....	xii
Chapter 1: Introduction .....	1
1.1 Study Background .....	1
1.2 Problem Statement .....	4
1.3 Research Objectives .....	5
1.4 Research Questions .....	5
1.5 Scope and Limitations .....	6
1.6 Research Relevance .....	7
Chapter 2: Literature Review .....	9
2.1 Introduction .....	9
2.2 Credit scoring algorithms .....	10
2.3 Explainable machine learning algorithms .....	12
2.3.1 Key XAI Techniques .....	12
2.4 Conceptual Framework Diagram .....	16
2.5 Methods and Tools .....	17

2.6	Comparative Analysis Summary .....	22
2.7	XAI and Credit scoring Models.....	24
2.8	Critical Limitations in Current Credit Scoring Paradigms .....	25
2.9	Conclusion .....	26
Chapter 3: Methodology .....		28
3.1	Introduction.....	28
3.2	Research Design for Data .....	29
3.3	Target Demographic .....	32
3.4	Feature Selection.....	32
3.5	Data pre-processing .....	33
3.5.1	Data Cleansing.....	33
3.5.2	Data Examination, Exploration, and Investigation.....	34
3.5.3	Handling Missing Data .....	34
3.6	Creditworthiness Assessment Algorithms .....	35
3.7	Ethical Considerations in Large-Scale Modelling.....	36
3.8	Explainability.....	37
3.8.1	Global interpretations.....	37
3.8.2	Local methods.....	37
3.9	Agile Software Development Methodology for Shiny App .....	37
3.10	Conclusions.....	39
Chapter 4: System Design and Architecture.....		40
4.1	Integrated Framework Overview .....	40
4.1.1	Data Source and System Overview.....	41
4.1.2	System Design and Architecture.....	41
4.2	Requirements Analysis .....	42
4.2.1	Functional Requirements .....	42

4.2.2	Non-Functional Requirements .....	44
4.3	System Architecture and Wireframe.....	44
4.3.1	Data Processing Module .....	45
4.3.2	Model Training Engine .....	47
4.3.3	Explanation Dashboard.....	48
4.3.4	Key Functional Components.....	50
4.3.5	Key Advantages Demonstrated.....	54
4.4	Principles for Choosing Parameters in Credit Scoring Models .....	55
4.5	System Wireframe .....	56
4.6	Conclusion .....	61
Chapter 5: System Implementation and Testing.....		62
5.1	Introduction.....	62
5.2	User Interface.....	62
5.3	System Architecture.....	63
5.4	Testing Methodology.....	64
5.4.1	System Testing.....	64
5.4.2	Usability Testing.....	64
5.4.3	Interpretability Testing.....	66
5.4.4	Key Findings.....	67
5.5	Model Validation .....	67
5.5.1	Addressing the Problem Statement.....	67
5.6	Conclusions.....	69
Chapter 6: Discussion of Results .....		70
6.1	Interpretation of Exploratory Analysis .....	70
6.1.1	Univariate Analysis.....	70
6.2	Bivariate Analysis.....	71

6.3	Multivariate Analysis.....	72
6.4	Model Performance and Practical Implications.....	74
6.4.1	Comparative Performance .....	74
6.4.2	Feature Importance .....	76
6.5	Discussions .....	76
6.6	Advantages of the Proposed System.....	76
6.7	Disadvantages of the Proposed System .....	76
6.8	Conclusions.....	77
Chapter 7: Conclusion and Recommendations.....		78
7.1	Synthesis of Key Findings .....	78
7.2	Conclusion: Strengths and Limitations.....	79
7.3	Recommendations.....	79
References.....		81
Appendices.....		86
Appendix A: Similarity Report.....		86
Appendix B: Ethics Clearance Confirmation .....		87



## List of Figures

Figure 2.1: Credit Lifecycle ALTAIR. (2022, May 24) .....	10
Figure 2.2: Conceptual Framework for Explainable Credit Scoring Systems.....	17
Figure 3.1: 1 Data Exploration Process, ALTAIR. (2022, May 24).....	31
Figure 3.2: Shiny App Development .....	31
Figure 3.3: Agile Methodology Implementation for Shiny App Development : A Step-by-Step Approach.....	38
Figure 4.1: System Architecture Wireframe.....	45
Figure 4.2: Code for Data processing module in shiny .....	46
Figure 4.3: Code for model Training in shiny .....	47
Figure 4.4: Hyperparameter Tuning Configuration for caret-LightGBM Integration.....	48
Figure 4.5: Explainability code in shiny.....	49
Figure 4.6: Dynamic User Interface Shiny code.....	51
Figure 4.7: Model Performance Computation in Shiny.....	52
Figure 4.8: Production Configuration for Scalability .....	53
Figure 4.9: User Interface Elements of the Loan Decision Support System .....	58
Figure 4.10: Top 10 important features for XGboost.....	59
Figure 4.11: Top 10 factures for Gradient Boost.....	60
Figure 4.12: Top 10 Features for Discriminant Analysis .....	60
Figure 4.13: Top 10 Features for Random Forest.....	61

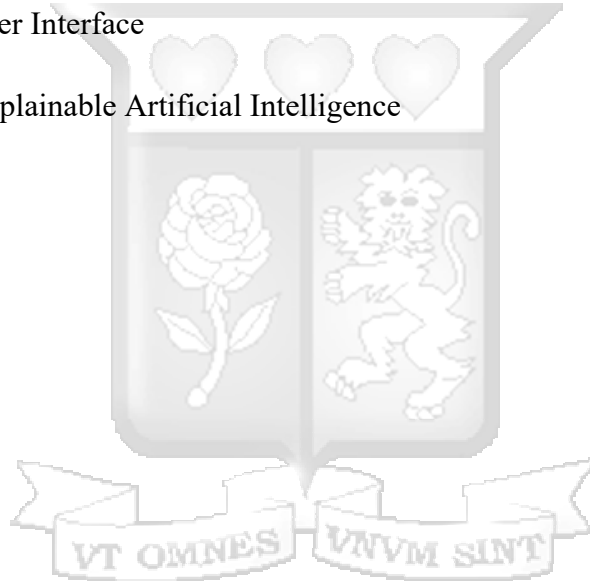
## List of Tables

Table 2.1: Comparative Analysis of XAI Techniques in Credit Scoring .....	23
Table 4.1: Model Performance Benchmarks .....	53
Table 5.1: Functionality Test .....	64
Table 5.2: Shows Task Completion Metrics.....	66



## List of Abbreviations

AI	Artificial Intelligence
CRD	Capital Requirements Directive
LC	Lending Club
ML	Machine Learning
PDP	Partial Dependence Plots
PFI	Permutation Feature Importance
UI	User Interface
XAI	eXplainable Artificial Intelligence



## Acknowledgements

The completion of this dissertation represents a significant milestone, made possible through the unwavering support, guidance, and encouragement of many remarkable individuals. I am deeply grateful to each of them for their contributions, which have shaped this work and enriched my academic journey.

First and foremost, I extend my heartfelt appreciation to my esteemed supervisor, whose expertise, patience, and mentorship have been invaluable. Their insightful feedback, constructive critiques, and steadfast encouragement provided the foundation for this research, guiding me through challenges and inspiring me to refine my ideas with clarity and rigour.

I am equally thankful to my colleagues at HFC, whose collaboration and camaraderie created a supportive environment that fuelled my progress. Their willingness to share knowledge, engage in thoughtful discussions, and offer practical insights into credit risk assessment enriched the practical relevance of this dissertation, and I am grateful for their generosity and teamwork.

To my classmates, I owe a special debt of gratitude for their companionship and intellectual exchange throughout this academic journey. Our shared experiences—late-night study sessions, spirited debates, and mutual encouragement—fostered a sense of community that made the challenges of graduate study more manageable and the successes more rewarding. Your diverse perspectives and unwavering support were a constant source of motivation.

Above all, I am profoundly grateful to my family, whose love, understanding, and encouragement have been my anchor. Their belief in my potential, coupled with their patience during long hours of research and writing, provided the emotional strength to persevere. This dissertation is as much a testament to their enduring support as it is to my own efforts.

# Chapter 1: Introduction

## 1.1 Study Background

Banks and lending institutions operate in an environment where credit issuance and collection are fundamental to their financial viability. Within the framework of Risk-Based Pricing, lending serves as a primary driver of profitability. However, this activity is inherently fraught with the risk of financial losses, which can significantly undermine overall profitability. Consequently, the accurate prediction of the Probability of Default (PD) for credit recipients is of paramount importance. Such predictions are critical for assessing the creditworthiness of individuals and entities, forming the backbone of sound credit risk management practices (Hand & Henley, 1997). The concept of credit scoring, which has evolved over nearly seven decades, has laid the foundation for modern credit scoring systems (Jakóbczak, 2015). Initially focused on individual consumers, the scope of credit scoring has expanded significantly as financial institutions have enhanced their data collection and analytical capabilities. The advent of big data has further revolutionized the field, enabling the development of credit scoring models not only for individual borrowers but also for small and medium-sized enterprises (MSMEs) and large corporations (Thomas et al., 2017).

In today's credit landscape, the credit scoring process involves the synthesis of qualitative, quantitative, and expert knowledge, combined with diverse methodologies, to provide actionable insights for informed decision-making in credit risk management (Thomas et al., 2017). In an era characterized by widespread credit accessibility, the ability to distinguish between high-risk and low-risk customers is crucial. This capability not only mitigates potential losses but also fosters a competitive advantage for financial institutions (Hand & Henley, 1997). Even a marginal improvement of 1% in loan return predictability can have a substantial impact on a bank's profitability, underscoring the critical role of credit scoring in refining lending decisions and promoting long-term financial sustainability.

Credit scoring relies on a borrower's credit score to assess their risk level, with machine learning algorithms offering financial institutions deeper insights into consumer transactions and

behaviours. The primary objective of credit risk management is to evaluate the likelihood of borrower default, as significant default rates can lead to insolvency and severe financial repercussions. Within the domain of commercial lending, the customer's credit scoring framework is of utmost importance (Thomas et al., 2017).

Three primary approaches are traditionally employed in assessing credit risk: indicators, expert evaluations, and credit scoring systems. These approaches rely on human expertise, augmented by various models, to appraise the creditworthiness of debtors (Thomas et al., 2017). However, conventional credit risk assessment techniques face significant challenges in keeping pace with the rapidly evolving demands of modern credit behaviour analysis. The integration of data science and machine learning offers a promising solution to enhance the efficacy and precision of these methodologies. Machine learning algorithms, which are more responsive to data patterns than traditional credit risk operational functions, empower financial institutions to mitigate risk more effectively (Hand & Henley, 1997).

Within the framework of risk-based pricing, banks utilize credit scores to forecast the probability of borrower repayment, a practice that is essential for accurate risk assessment (CBK, 2013). Assessing a customer's risk level hinges on a comprehensive evaluation of creditworthiness, which includes factors such as credit history, loan size, debt-to-income ratio, loan maturity updates, stability metrics, public records, and security provisions. It is imperative that the models used for such assessments, as well as their outcomes, maintain a high degree of transparency. These models must undergo rigorous regulatory validation and receive approval from central banks, as stipulated in the White Paper on Artificial Intelligence (Slovik, 2012) and the Data Protection (General) Regulations of 2021.

In lending relationships, trust in a customer's ability, willingness, and reliability to meet repayment obligations is paramount. Borrowers must cultivate a trustworthy profile through timely repayments and consistent income, while lenders bear the responsibility of accurately evaluating risk to determine appropriate interest rates. This dynamic has driven the rapid advancement of statistical models (Leo et al., 2019), offering unprecedented opportunities to understand and predict credit behaviour (Siddiqi, 2012). As emphasized by Martens et al. (2007), the significance of credit scoring tiers lies not only in their ability to provide diverse alternatives but also in their capacity to bridge the gap between explicit knowledge and modelled predictions.

The opacity of machine learning models to credit risk assessment - often criticized as "black box" systems - presents significant challenges for both financial institutions and regulators. To address this critical issue, our research employs Explainable AI (XAI) techniques that illuminate model decision-making processes while maintaining predictive performance. Local Interpretable Model-agnostic Explanations (LIME) serves as our foundational approach, generating intuitive explanations for individual predictions by creating locally faithful approximations of complex models. Complementing this, we implement SHAP (Shapley Values) analysis to quantify the precise contribution of each feature to final predictions, providing outputs that align naturally with established risk-based pricing frameworks (Talaat and El-Balky, 2023). For broader model interpretation, we utilize Partial Dependency Plots (PDPs) and conventional Feature Importance metrics to visualize global patterns and identify dominant variables across the entire dataset. Further strengthening our analytical framework, Permutation Feature Importance testing systematically evaluates predictor significance by measuring model performance degradation when specific feature values are randomized. Together, these techniques form a robust explainability infrastructure that not only demystifies model operations but also satisfies growing regulatory demands for algorithmic transparency in financial decision-making (Ariza-Garzón et al., 2020).

While our current research focuses primarily on model performance and explainability, it acknowledges the profound ethical dimensions inherent in credit scoring systems. The substantial dataset employed in this study ( $n > 10,000$ ) offers some inherent protection against individual biases through the law of large numbers, which statistically dilutes anomalous data points. However, the study recognises that XAI techniques, while valuable for interpretation, do not automatically ensure fairness or eliminate systemic biases (Bussmann et al.,

2020)). This important limitation underscores the need for future research to incorporate dedicated bias auditing protocols, particularly examining potential demographic disparities in model outcomes. The current study's exclusion of comprehensive ethical analysis reflects methodological prioritization rather than disregard for these critical concerns, and we explicitly identify this as a valuable direction for subsequent investigation.

This chapter sets the stage for the dissertation by highlighting the critical role of credit scoring in modern financial systems, the challenges posed by traditional methodologies, and the

transformative potential of data science and machine learning in enhancing credit risk assessment. The subsequent chapters will delve deeper into the methodologies, innovations, and practical applications that address these challenges, ultimately contributing to a more robust and transparent credit risk management framework.

## **1.2 Problem Statement**

Within the realm of credit scoring, financial institutions and regulators face a formidable challenge: enhancing the precision of predictive models while simultaneously adhering to stringent regulatory mandates. This dual objective represents a critical research gap, as it necessitates balancing the growing demand for highly accurate forecasts in the financial sector with the need for transparency and interpretability in model outputs. The increasing complexity of contemporary models and statistical methodologies, driven by advanced mathematical transformations and algorithmic approaches, introduces significant challenges in variable selection, model interpretation, and regulatory compliance.

Regulatory frameworks, such as the Basel Accord and local directives issued by the Central Bank of Kenya, emphasize the importance of model interpretability and transparency. These requirements highlight the tension between leveraging sophisticated models to improve predictive accuracy and ensuring that such models remain comprehensible and actionable for stakeholders, including credit managers, regulators, and policymakers. The inability to reconcile these competing demands poses a significant barrier to the widespread adoption of advanced credit scoring models, limiting their potential to enhance risk management practices and financial decision-making.

This research addresses this critical issue by exploring methodologies that improve the precision of credit scoring models while maintaining their interpretability and compliance with regulatory standards. By doing so, it seeks to bridge the gap between cutting-edge predictive analytics and the practical, regulatory, and operational needs of financial institutions.

### **1.3 Research Objectives**

The research objectives of this study are designed to address the dual challenges of enhancing predictive accuracy and maintaining interpretability in credit scoring models. These objectives are outlined as follows:

- a. To enhance the forecasting capacity of credit rating models by incorporating nonlinear elements.

This objective focuses on improving the predictive performance of credit scoring models by integrating advanced nonlinear techniques, such as machine learning algorithms and statistical transformations. By leveraging these methodologies, the study aims to capture complex patterns and relationships within credit data that traditional linear models may overlook, thereby increasing the accuracy of default probability predictions.

- b. To establish a technique for maintaining or improving the comprehensibility of intricate, nonlinear credit scoring models.

While advanced models often yield superior predictive performance, their complexity can hinder interpretability, making it difficult for stakeholders to understand and trust their outputs. This objective seeks to develop and validate methodologies, such as Explainable Artificial Intelligence (XAI) techniques, that enhance the transparency of nonlinear models without compromising their accuracy. The goal is to ensure that these models remain accessible and actionable for credit risk practitioners, regulators, and other stakeholders.

By achieving these objectives, the research aims to bridge the gap between predictive power and interpretability, enabling financial institutions to adopt advanced credit scoring models with greater confidence and regulatory compliance.

### **1.4 Research Questions**

The research questions guiding this study are designed to explore the integration of nonlinear components into credit evaluation models and the development of strategies to ensure these models remain both precise and interpretable. These questions are as follows:

- a. What are the most effective approaches for integrating nonlinear components into credit evaluation models, and how do these approaches influence the predictive precision and overall efficacy of the models while maintaining their interpretability?

This question seeks to identify and evaluate methodologies for incorporating nonlinear elements, such as machine learning algorithms and advanced statistical techniques, into credit scoring models. It also examines the impact of these approaches on the models' predictive accuracy, operational effectiveness, and ability to remain interpretable for stakeholders.

- b. What strategies can be employed to create nonlinear credit assessment models that demonstrate both precision and comprehensibility? Furthermore, how do these models measure up against conventional linear credit scoring models regarding effectiveness and interpretability?

This question focuses on developing and validating strategies to ensure that nonlinear credit scoring models achieve a balance between high predictive accuracy and transparency. It also aims to compare the performance of these advanced models with traditional linear models, assessing their relative strengths and weaknesses in terms of predictive power, interpretability, and practical applicability.

By addressing these research questions, the study aims to provide actionable insights into the development and implementation of advanced credit scoring models that meet the dual demands of precision and interpretability, ultimately contributing to more robust and transparent credit risk management practices.

## **1.5 Scope and Limitations**

### **Scope**

The study delves into the integration of nonlinear elements, such as machine learning algorithms and advanced statistical techniques, into credit assessment models, examining how these components enhance predictive precision, comprehensibility, and overall efficacy. It places significant emphasis on the development and comparison of both nonlinear and linear models, particularly in their capacity to strike a balance between predictive accuracy and interpretability. Central to the investigation is the prioritization of model interpretability, ensuring that the

outputs of sophisticated models remain accessible and actionable for key stakeholders, including credit risk practitioners and regulators. This focus underscores the importance of not only advancing predictive capabilities but also maintaining transparency and usability in the application of these models.

## Limitations

The study distinguishes between algorithmic transparency, which concerns the internal mechanics of models, and interpretability, which focuses on ensuring that model outputs are clear and understandable to end-users. While the research emphasizes the importance of both precision and comprehensibility as desirable qualities in credit scoring models, it acknowledges that these attributes may not always align seamlessly in every scenario or application. As a result, the findings may not be universally applicable across all use cases or contexts. Additionally, the study relies on a specific dataset provided by a financial institution, which may constrain the generalizability of the results to other datasets or regions. Furthermore, the research does not delve into the ethical implications of employing advanced models in credit scoring, such as potential biases or fairness concerns, leaving this as a critical area for future exploration.

By defining the scope and acknowledging these limitations, the study aims to provide a focused and realistic exploration of the challenges and opportunities in developing advanced credit scoring models that balance precision and interpretability.

## 1.6 Research Relevance

In the context of rapidly evolving credit services, this study addresses a critical challenge: enhancing the interpretability and transparency of advanced credit scoring models. As financial institutions increasingly adopt sophisticated methodologies such as random forests, gradient-boosted trees, support vector machines, and neural networks, the opacity of these models poses significant barriers to their widespread adoption. This dissertation directly confronts this issue by proposing concrete approaches to improve the transparency of advanced statistical and machine learning techniques, ensuring they are accessible and actionable for risk analysts and other stakeholders.

By elucidating how specific machine learning techniques influence outcomes and conducting "what-if" analyses, this research provides actionable insights into the technology underpinning machine learning predictions (Biecek et al., 2021). These efforts address the primary hurdles in integrating modern AI tools into credit risk modelling, offering straightforward and interpretable solutions that bridge the gap between cutting-edge technology and practical application.

In an era characterized by the abundance of information, big data plays a pivotal role in individual credit evaluation. This study constructs a robust mathematical framework to assess a customer's creditworthiness by leveraging historical data on borrowing patterns and subsequent credit appraisals. This framework not only enhances the accuracy of credit risk assessments but also streamlines processes related to product and content distribution. By doing so, it reduces labour expenses, mitigates measurement inaccuracies, and resolves inefficiencies and high costs associated with traditional risk evaluation methods. Data annotation further enhances the quality of model inputs, thereby improving the precision of analyses.

Beyond addressing interpretational challenges, this dual approach—combining advanced modelling techniques with enhanced transparency—significantly improves the efficiency and efficacy of the credit risk assessment process. These advancements benefit a wide range of stakeholders, including financial institutions, regulatory bodies, and individuals involved in data management endeavours. By fostering greater trust and understanding of advanced models, this research contributes to the broader adoption of AI-driven tools in credit risk management, ultimately promoting more robust, transparent, and equitable financial systems.

## Chapter 2: Literature Review

### 2.1 Introduction

In the realm of credit evaluation, traditional credit scoring methods aim to construct predictive models that estimate the likelihood of an applicant avoiding default within a specified period, typically aligned with CRD IV/Basel recommendations. As highlighted by Biecek et al. (2021), these methods predominantly rely on statistical techniques such as Linear Discriminant Analysis, Logistic Regression, and Decision Trees to profile applicant behaviours and assign scores that reflect anticipated performance. The primary objective is to determine whether applicants are likely to exhibit positive or negative credit behaviour. These methodologies are widely adopted across risk management departments in financial institutions, forming the backbone of credit risk assessment practices.

Biecek et al. (2021) further emphasize that optimal credit scoring strategies extend beyond basic statistical techniques. They incorporate advanced practices such as feature engineering, including the Weight of Evidence transformation, and feature selection techniques like addressing collinearity and employing general-to-specific or stepwise selection. Despite the continued reliance on these traditional methods, the evolving landscape of credit scoring demands a deeper exploration of modern approaches, particularly in terms of enhancing model interpretability and transparency. This is especially critical as financial institutions increasingly adopt sophisticated models to improve predictive accuracy while ensuring compliance with regulatory standards.

The credit lifecycle, as illustrated in Figure 2.1, provides a comprehensive framework for understanding the stages of credit risk management. The lifecycle begins with the Origination/Application Scorecard, which represents the initial risk assessment of new credit applications. This stage involves evaluating the creditworthiness of applicants based on their financial history, behaviour, and other relevant factors.

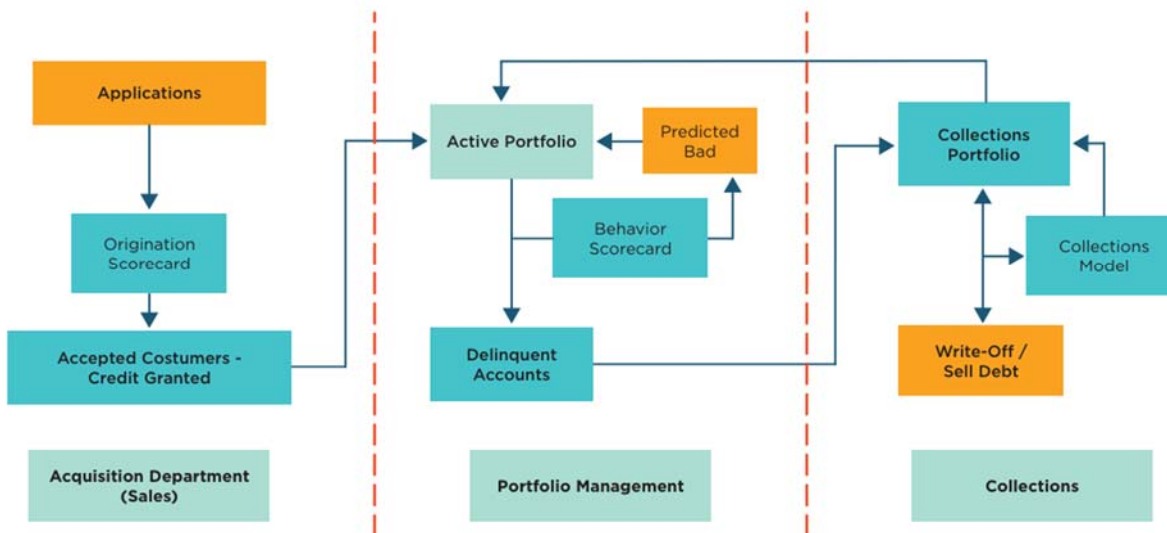


Figure 2.1: Credit Lifecycle ALTAIR. (2022, May 24)

The next stage, Account Management, focuses on the ongoing assessment and management of credit risk using behaviour scorecards generated by credit scoring algorithms. These scorecards enable financial institutions to monitor and adjust their risk exposure dynamically, ensuring proactive management of credit portfolios.

Finally, the Collections stage involves the assessment and decision-making process for handling defaults. Collection models are employed to determine the most effective strategies for recovering outstanding debts, minimizing losses, and maintaining the financial health of the institution.

This lifecycle underscores the importance of robust and interpretable credit scoring models at every stage, from initial application to ongoing account management and collections. As the credit landscape continues to evolve, the integration of advanced methodologies and the emphasis on interpretability will play a pivotal role in enhancing the efficacy and transparency of credit risk management practices.

## 2.2 Credit scoring algorithms.

The field of credit scoring has witnessed a significant evolution in statistical modelling, with a noticeable shift toward the adoption of increasingly sophisticated methodologies. A comprehensive review and synthesis of scholarly literature reveal the growing use of advanced techniques to forecast credit market dynamics. While traditional methods such as logistic regression remain pivotal and widely valued (Jakóbczak, 2015), decision trees continue to hold

significant importance in this domain (Dastile et al., 2020). Beyond these conventional approaches, a variety of other methodologies have been explored, including risk regression (Dastile et al., 2020), random forests (Leo et al., 2019), k-nearest neighbours (KNN) (Dastile et al., 2019), artificial neural networks (ANN) (Dastile et al., 2020), and support vector machines (SVM) (Dastile et al., 2020), among others.

Lessmann et al. (2015) conducted an extensive evaluation of 41 state-of-the-art classification algorithms, examining performance metrics and methodologies that facilitate rigorous comparisons among different classifiers. Their work also provides a systematic review of the literature on the theory and application of classification techniques in credit assessment, offering valuable insights into the strengths and limitations of various approaches. Despite these advancements, a notable research gap persists in the integration of non-linear components to enhance predictive accuracy. While traditional linear methods have been extensively studied, there is limited exploration of the benefits and challenges associated with incorporating non-linear elements into credit scoring models.

This gap is particularly relevant in the context of modern financial data, which is increasingly complex and heterogeneous. The growing intricacy of financial behaviours and the need for higher precision in credit risk assessments underscore the importance of exploring non-linear methodologies. Non-linear components, such as those found in advanced machine learning algorithms, have the potential to capture intricate patterns and relationships within data that linear models may overlook. However, their adoption is often hindered by challenges related to interpretability, computational complexity, and regulatory compliance.

Addressing this gap is critical for advancing credit scoring practices. By integrating non-linear components, financial institutions can improve the accuracy of their predictive models, leading to better risk management and decision-making. At the same time, it is essential to develop strategies that ensure these advanced models remain interpretable and transparent, enabling stakeholders to understand and trust their outputs. This dual focus on accuracy and interpretability will be a key theme in this study, as it seeks to bridge the gap between traditional and innovative credit scoring methodologies.

## 2.3 Explainable machine learning algorithms

In the context of credit risk modelling, even minor improvements in predictive accuracy can lead to significant future savings for financial institutions. This underscores the importance of developing robust models that not only perform well but are also interpretable and transparent. Contemporary statistical and machine learning methods have demonstrated their ability to effectively classify customers as reliable or unreliable payers based on credit risk traits. However, their application in real-world business settings presents a significant challenge: while these models are highly effective, their complexity often makes it difficult for users to understand the rationale behind their decisions. This lack of transparency can hinder trust and adoption, particularly in high-stakes environments where each decision carries substantial weight (Hand & Henley, 1997).

This challenge is further amplified by regulatory requirements, such as those outlined in the Basel Accords, which mandate that financial institutions not only achieve accurate predictions but also ensure that their models are transparent and interpretable. To address this, Explainable Artificial Intelligence (XAI) has emerged as a critical area of research. XAI methodologies aim to enhance the interpretability of complex, opaque models, enabling stakeholders to understand how decisions are made and why certain predictions are generated. As Molnar (2020) emphasizes, XAI plays a pivotal role in bridging the gap between model complexity and human understanding, making it an indispensable tool for modern credit risk modelling.

### 2.3.1 Key XAI Techniques

XAI techniques are designed to provide insights into the inner workings of machine learning models, making them accessible and actionable for users. Many of these techniques are model-agnostic, meaning they can be applied to a wide range of machine learning models, from simple linear regressions to complex neural networks. This flexibility allows for straightforward and universal comparisons across different models, facilitating better decision-making and model selection. Some of the most widely used XAI techniques include:

#### 2.3.1.1 LIME

LIME (Local Interpretable Model-agnostic Explanations) (Ribeiro, 2016) explains individual predictions by approximating the behaviour of a complex model locally around a specific instance.

It works by perturbing the input data and observing how the model's predictions change, then fitting a simpler interpretable model (e.g., linear regression) to the perturbed data.

Mathematical Formulation: For a given instance  $x$ , LIME generates a set of perturbed samples  $x'$  and weights  $\pi_x(x')$  based on their proximity to  $x$ . It then fits a simple model  $g$  (e.g., linear regression) to minimize the following objective function:

(2.1)

$$\xi(x) = \operatorname{argmin}_{g \in G} L(f, g, \pi_x) + \Omega(g)$$

Where:

- $f$  is the complex model being explained.
- $g$  is the interpretable model (e.g., linear regression).
- $L$  is a loss function (e.g., mean squared error) that measures how well  $g$  approximates  $f$  in the local region.
- $\pi_x$  is a weighting function that assigns higher weights to samples closer to  $x$ .
- $\Omega(g)$  is a regularization term to ensure the simplicity of  $g$ .

LIME provides local explanations by highlighting the contribution of each feature to the prediction for a specific instance.

### 2.3.1.2 SHAP

SHAP (SHapley Additive exPlanations) (Lundberg, 2017), provides a unified framework for explaining model outputs by attributing the contribution of each feature to the final prediction. It is based on cooperative game theory and offers both global and local interpretability.

Mathematical Formulation: The Shapley value  $\phi_i$  for a feature  $i$  is calculated as:

(2.2)

$$\phi_i(f, x) = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)]$$

Where:

- $F$  is the set of all features.
- $S$  is a subset of features excluding  $i$ .
- $f(S)$  is the model's prediction using only the features in  $S$ .
- $f(S \cup \{i\})$  is the model's prediction when feature  $i$  is added to  $S$ .

SHAP values provide both global and local interpretability, making it a powerful tool for understanding model behaviour.

### 2.3.1.3 DALEX

DALEX (Descriptive mACHine Learning EXplanations) (Biecek, 2018) is a comprehensive framework for explaining and visualizing machine learning models. It provides tools for model exploration, comparison, and interpretation, making it easier for users to understand and trust model outputs.

#### 2.3.1.3.1 Variable Importance

Measures the contribution of each feature to the model's predictions. For a model  $f$ , the importance of feature  $i$  is computed as:

$$VI_i = \frac{1}{N} \sum_{j=1}^N |f(x_j) - f(x_j^{(i)})|$$

(2.3)

Where:

- $x_j^{(i)}$  is the instance  $x_j$  with the value of feature  $i$  permuted.
- $N$  is the number of instances.

#### 2.3.1.3.2 Partial Dependence Plots (PDPs)

Show the relationship between a feature and the predicted outcome, marginalizing over the effects of other features. For a feature  $i$ , the PDP is computed as:

(2.4)

$$PDP_i(x_i) = \frac{1}{N} \sum_{j=1}^N f(x_i, x_{-i}^{(j)})$$

Where: -  $x_{-i}^{(j)}$  represents the values of all features except  $i$  for instance  $j$ .

#### 2.3.1.4 InterpretML

InterpretML (Biecek, 2018), is another framework that offers a suite of tools for explaining machine learning models. It supports both glass-box (interpretable) and black-box (complex) models, enabling users to balance accuracy and interpretability. InterpretML is a framework that supports both glass-box (interpretable) and black-box (complex) models. It provides tools for model exploration, comparison, and interpretation.

Mathematical Formulation: InterpretML uses techniques like Generalized Additive Models (GAMs) to balance accuracy and interpretability. A GAM is defined as:

(2.5)

$$g(E[y]) = \beta_0 + \sum_{i=1}^p f_i(x_i)$$

Where: -  $g$  is a link function (e.g., logit for classification).

-  $E[y]$  is the expected value of the target variable.

-  $\beta_0$  is the intercept.

-  $f_i(x_i)$  are smooth functions (e.g., splines) for each feature  $x_i$ .

GAMs provide interpretable models by allowing each feature to have a non-linear relationship with the target, while still maintaining transparency.

#### 2.3.1.5 The Need for Balancing Accuracy and Interpretability

While XAI techniques have gained significant traction in recent years, a notable gap remains in reconciling the accuracy achieved by incorporating non-linear components with the essential requirement of interpretability. Non-linear models, such as random forests, gradient-boosted trees,

and neural networks, often deliver superior predictive performance but are inherently opaque. This opacity can hinder their adoption in regulated industries like finance, where transparency is a regulatory and operational necessity.

To address this gap, this research focuses on developing concrete methodologies that integrate non-linear elements into credit scoring models while ensuring their interpretability. By leveraging XAI techniques, the study aims to create models that are not only accurate but also transparent, enabling stakeholders to understand and trust their outputs. This dual focus on accuracy and interpretability is critical for advancing credit risk modelling practices and ensuring compliance with regulatory frameworks like the Basel Accords.

In summary, the integration of XAI methodologies into credit risk modelling represents a significant step forward in addressing the challenges posed by complex models. By making these models more interpretable, financial institutions can enhance their decision-making processes, build trust with stakeholders, and achieve better outcomes in credit risk management.

#### **2.4 Conceptual Framework Diagram**

Figure 2.2 illustrates a three-layer hybrid credit scoring system. The input layer combines traditional factors (FICO scores, DTI ratios) with alternative data (cash flows, behavioural patterns). The processing layer employs hybrid modelling with non-linear feature engineering and explainable AI (XAI) components. The output layer generates risk predictions, local/global explanations, and regulatory documentation, balancing predictive accuracy with interpretability and compliance requirements.

This architecture addresses key challenges in modern credit assessment by integrating diverse data sources while maintaining transparency and regulatory adherence. The figure 2.2 presents a three-layer framework for modern credit scoring. The Input Layer combines traditional credit factors (e.g., FICO scores, debt-to-income ratios) with alternative data sources (e.g., cash flows, behavioural patterns). The Processing Layer employs a hybrid modelling approach, integrating non-linear feature engineering and explainable AI (XAI) components to enhance interpretability. Finally, the Output Layer generates risk predictions, provides local and global explanations of model decisions, and produces regulatory documentation to ensure compliance. This structured

architecture enables accurate risk assessment while maintaining transparency and meeting regulatory requirements.

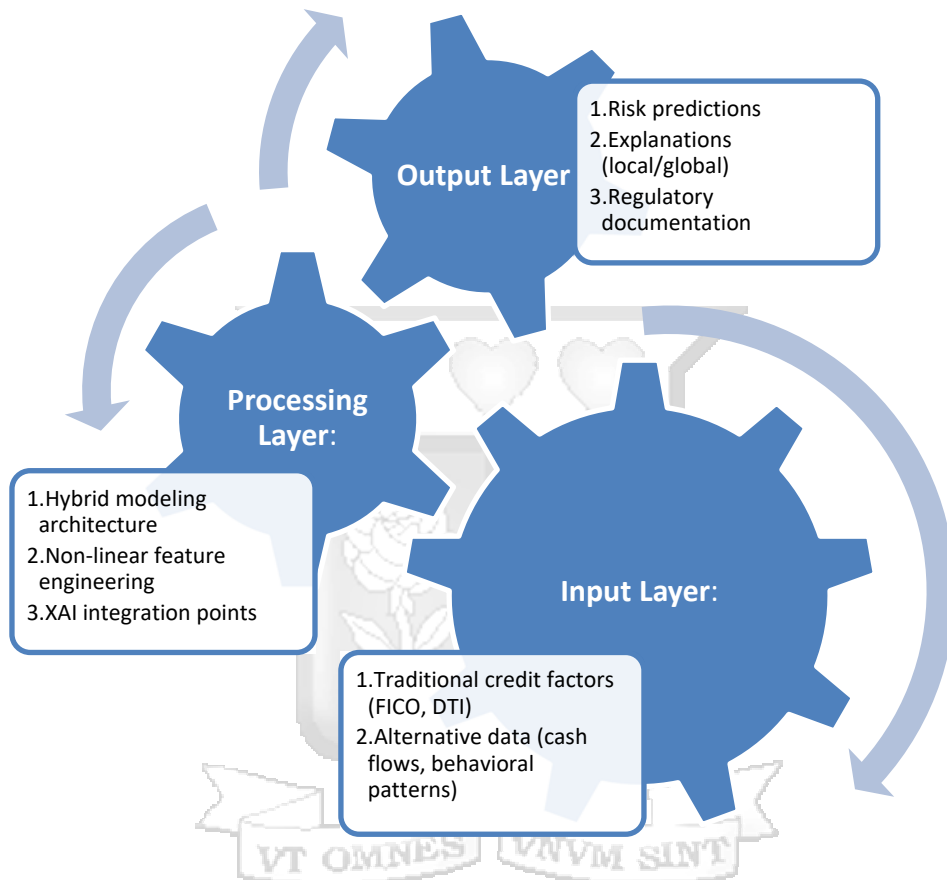


Figure 2.2: Conceptual Framework for Explainable Credit Scoring Systems

## 2.5 Methods and Tools

Among the array of methods and tools available for enhancing the interpretability and performance of credit scoring models, a select few warrant further investigations. These methods provide valuable insights into feature importance, model behaviour, and the relationships between input variables and predicted outcomes. Below, we discuss key techniques, including Permutation Feature Importance (PFI), Partial Dependence Plots (PDP), Breakdown (BD) diagrams, and Ceteris Paribus (CP) plots, with references to their origins and applications.

### 2.5.1.1 Permutation Feature Importance (PFI)

Introduced by Friedman (2001), Permutation Feature Importance (PFI) is a straightforward yet powerful method for assessing the importance of individual features in a model. This technique evaluates the contribution of each feature to the model's predictive performance by measuring the change in accuracy when the feature's values are permuted (i.e., shuffled).

#### 2.5.1.1.1 Steps to Compute PFI

To evaluate the importance of features in credit risk classification models, a structured approach is followed. First, an appropriate accuracy metric is selected; for this purpose, the Area Under the Receiver Operating Characteristics Curve (AUC) is commonly used. AUC measures the model's ability to distinguish between positive and negative classes, such as default versus non-default cases (Friedman, 2001). Next, baseline performance is established by training the model and evaluating its performance using the chosen metric (e.g., AUC) on the original dataset (Friedman, 2001).

Once the baseline is set, the process involves permuting the values of each feature individually while keeping the other features unchanged. This permutation disrupts the relationship between the feature and the target variable, allowing for an assessment of the feature's impact (Friedman, 2001). After permuting, the model's performance is re-evaluated using the same metric (e.g., AUC) on the modified dataset (Friedman, 2001). Finally, the importance of each feature is calculated by measuring the decrease in the model's performance (e.g., reduction in AUC) when the feature is permuted. A larger decrease in performance indicates that the feature is more important to the model's predictive capability (Friedman, 2001). This method provides a systematic way to quantify the contribution of individual features to the model's overall performance.

For a feature  $X_j$ , the permutation feature importance  $I_j$  is calculated as:

(2.6)

$$I_j = \text{Performance}(f, D) - \text{Performance}(f, D^{(j)})$$

Where:

$f$  is the trained model.

$D$  is the original dataset.

$D^{(j)}$  is the dataset with the values of feature  $X_j$  permuted.

$\text{Performance}(f, D)$  is the model's performance (e.g., AUC) on the original dataset.

$\text{Performance}(f, D^{(j)})$  is the model's performance on the permuted dataset.

PFI provides a clear and interpretable measure of feature importance, making it a valuable tool for understanding model behaviour (Friedman, 2001).

### 2.5.1.2 Partial Dependence Plots (PDP)

Partial Dependence Plots (PDP), also introduced by Friedman (2001), are a visual tool for understanding the marginal effect of one or two features on the predicted outcome of a machine learning model. PDPs illustrate the relationship between a feature and the target variable, accounting for the average effect of all other features.

#### 2.5.1.2.1 Steps to Compute PDP

To analyse the influence of specific features on a model's predictions, a systematic approach is employed. First, a feature of interest is selected, such as income or debt-to-income ratio, to focus the analysis (Friedman, 2001). Once the feature is chosen, its values are fixed across a defined grid, representing a range of possible values over which the model's predictions will be evaluated (Friedman, 2001). For each value in this grid, the feature's values in the dataset are replaced with the fixed value, and the model's predictions are computed. These predictions are then averaged across all instances to derive the partial dependence, which reflects the feature's average effect on the model's output (Friedman, 2001). Finally, the results are visualized by plotting the relationship between the feature(s) and the predicted outcome, providing a clear and interpretable representation of how the feature influences the model's predictions (Friedman, 2001). This method offers a structured way to understand and communicate the impact of individual features on the model's behaviour.

Mathematical Formulation: For a feature  $X_j$ , the partial dependence function  $PDP_j(x_j)$  is defined as:

(2.7)

$$PDP_j(x_j) = \frac{1}{N} \sum_{i=1}^N f(x_j, x_{-j}^{(i)})$$

Where:

$f$  is the trained model.

$x_j$  is the fixed value of the feature  $X_j$ .

$x_{-j}^{(i)}$  represents the values of all other features for the  $i$ -th instance.

$N$  is the number of instances in the dataset.

For two features  $X_j$  and  $X_k$ , the joint partial dependence function is:

(2.8)

$$PDP_{j,k}(x_j, x_k) = \frac{1}{N} \sum_{i=1}^N f(x_j, x_k, x_{-j,k}^{(i)})$$

#### 2.5.1.2.2 Interpretation of PDP

The shape of a Partial Dependence Plot (PDP) provides valuable insights into the relationship between a feature and the target variable. If the PDP displays a straight line, it indicates that the feature has a linear relationship with the target, suggesting a consistent and proportional influence on the model's predictions (Friedman, 2001). Conversely, if the plot reveals curves or more complex patterns, this signifies a non-linear relationship, where the feature's impact on the target varies in a less straightforward manner (Friedman, 2001). Additionally, by examining joint PDPs, which analyse the combined effect of two features, interactions between features can be uncovered. These interactions highlight how the influence of one feature on the target may depend on the value of another feature, providing a deeper understanding of the model's behaviour (Friedman, 2001). Together, these observations help interpret the nature and complexity of feature-target relationships in predictive models.

PDPs are particularly useful for understanding how individual features influence the model's predictions, providing actionable insights for credit risk assessment (Friedman, 2001).

### **2.5.1.3 Breakdown (BD) Diagrams**

Breakdown (BD) diagrams, introduced by Baniecki and Biecek (2019), illustrate the contribution of each feature to the model's prediction for a specific instance. This method provides a local explanation of the model's behaviour, showing how each variable influences the predicted outcome.

#### **2.5.1.3.1 Steps to Compute BD Diagrams**

To explain the prediction of a model for a specific instance, a step-by-step process is followed. First, a particular instance is selected for analysis, serving as the focal point for understanding the model's decision-making process (Baniecki & Biecek, 2019). Next, the baseline prediction for this instance is computed, establishing the model's output before any additional analysis is performed (Baniecki & Biecek, 2019). Following this, features are sequentially added one by one, and the change in the model's prediction is calculated at each step. This incremental approach helps isolate the contribution of each feature to the final prediction (Baniecki & Biecek, 2019). Finally, the contributions of each feature are visualized, typically through a plot, to provide a clear and interpretable representation of how each feature influences the model's output for the selected instance (Baniecki & Biecek, 2019). This method offers a transparent and systematic way to dissect and communicate the factors driving the model's predictions. BD diagrams offer a concise overview of how each variable influences the predictive model, making them a valuable tool for local interpretability (Baniecki & Biecek, 2019).

### **2.5.1.4 Ceteris Paribus (CP) Plots**

Ceteris Paribus (CP) plots, also discussed by Baniecki and Biecek (2019), enable direct analysis of how changes in a specific feature affect the model's predictions while holding all other features constant. This method provides a local interpretation of the model's behaviour.

#### **2.5.1.4.1 Steps to Compute CP Plots**

To analyse the influence of a specific feature on a model's predictions, a structured approach is employed. First, a feature of interest is selected for analysis, serving as the focal point of the investigation (Baniecki & Biecek, 2019). Once the feature is chosen, all other features are held constant at their original values to isolate the effect of the selected feature (Baniecki & Biecek,

2019). Next, the value of the selected feature is varied across a defined range, and the model's predictions are computed for each value in this range. This step helps capture how changes in the feature impact the model's output (Baniecki & Biecek, 2019). Finally, the results are visualized by plotting the relationship between the selected feature and the predicted outcome, providing a clear and interpretable representation of the feature's influence on the model's predictions (Baniecki & Biecek, 2019). This method offers a systematic way to understand and communicate the role of individual features in shaping the model's behaviour. CP plots provide a straightforward approach to understanding how variations in input variables influence predictions, offering insights similar to PDPs but at a local level (Baniecki & Biecek, 2019).

## 2.6 Comparative Analysis Summary

Table 2.1 presents a comparative analysis of explainable AI (XAI) techniques for credit scoring applications. The table evaluates five prominent methods (LIME, SHAP, PDPs, PFI, and BD/CP) across three dimensions: advantages, limitations, and optimal use cases. Key findings indicate that SHAP offers the most comprehensive framework (combining global and local explanations), while LIME provides more accessible but approximate interpretations. Partial dependence plots (PDPs) emerge as particularly effective for visualizing non-linear relationships, though all methods exhibit trade-offs between computational complexity and explanatory power. The table systematically organizes these techniques to guide appropriate selection based on specific credit risk modelling needs.

Table 2.1: Comparative Analysis of XAI Techniques in Credit Scoring

Technique	Advantages	Limitations	Best Use Cases
LIME	- Provides local explanations	- Approximate explanations only	Explaining individual credit decisions
	- Model-agnostic	- Sensitive to perturbation parameters	Debugging model behaviour
	- Intuitive for non-technical users	- Computationally expensive for large datasets	
SHAP	- Unified framework (global + local)	- Computationally intensive	Feature importance analysis
	- Game-theoretically sound	- Complex interpretation for high-dimensional data	Regulatory compliance reporting
	- Handles feature dependencies	- Requires sampling for large models	
PDPs	- Visual and intuitive	- Assumes feature independence	Understanding marginal effects
	- Reveals non-linear relationships	- Limited to 2-3 features	Model validation
	- Handles interactions (2D PDPs)	- Can mask heterogeneous effects	
PFI	- Simple implementation	- Permutation can create unrealistic data	Feature selection
	- Direct performance-based metric	- No directionality of effects	Model simplification
	- Works with any model	- Dependent on model quality	
BD/CP	- Instance-specific explanations	- Order-dependent results (BD)	Customer dispute resolution
	- Shows prediction pathways	- Limited to single instances	Model debugging
	- Interactive analysis	- Manual interpretation required	

## 2.7 XAI and Credit scoring Models

The integration of machine learning (ML) and explainable artificial intelligence (XAI) into credit scoring has garnered increasing attention, reflecting the dual need for predictive accuracy and transparency in financial decision-making. Bussmann et al. (2020) utilized XGBoost with SHAP (SHapley Additive exPlanations) to model credit risk, demonstrating how XAI enhances interpretability in fintech applications. Similarly, Misheva et al. (2021) applied LIME (Local Interpretable Model-agnostic Explanations) and SHAP to credit scoring in peer-to-peer (P2P) lending, leveraging real-world datasets to balance predictive power and explainability. These studies highlight the critical role of XAI in modern credit risk assessment.

Traditional credit scorecards, often based on logistic regression or decision trees, remain widely used due to their interpretability. Fahner (2018) explored how ML can improve scorecard development while maintaining transparency, proposing a hybrid approach compliant with regulatory standards. Bücker et al. (2021) compared XAI-enhanced ML models, such as random forests, with traditional scorecards using public credit bureau data, finding that ensemble methods can outperform simpler models when paired with XAI tools, as measured by precision and recall. These findings support a shift toward advanced modelling in credit scoring.

Advanced ML techniques have become prevalent for capturing complex patterns in credit data. Ke et al. (2017) introduced LightGBM, an efficient gradient-boosting framework widely adopted in credit risk modeling for its scalability and accuracy. Talaat and El-Balky (2023) combined deep learning with SHAP to predict credit card defaults, illustrating how XAI can address the opacity of neural networks. Ariza-Garzón et al. (2020) employed ensemble methods like random forests with SHAP and LIME, emphasizing their ability to visualize key feature contributions, such as payment history or debt-to-income ratio, in credit risk models.

Data preprocessing is a foundational step in credit scoring research. Fahner (2018) and Ariza-Garzón et al. (2020) outlined techniques for handling missing values, encoding categorical variables, and feature selection, aligning with best practices for financial datasets. Davis et al. (2022) applied XAI to home equity loan data, stressing the importance of high-quality, representative datasets for generalizability, a challenge also noted by Misheva et al. (2021) in P2P lending contexts.

Deployment strategies have also been explored in the literature. Zoldi and Fahner (2022) advocated for integrating ML with traditional scorecards, using Weight of Evidence (WoE) and Information Value (IV) to enhance practical applicability. Moscato et al. (2021) proposed an XAI framework for P2P lending credit scoring, emphasizing user-centric design in deployment, though not specifically with Shiny. Chang and Wickham (2018) detailed the Shiny framework for R, providing a technical basis for interactive applications, while Schwaber and Sutherland (2020) outlined Agile principles, supporting iterative and user-focused development.

Together, these studies underscore the synergy of predictive power, interpretability, and deployment readiness in credit scoring, forming a robust foundation for this research, which aims to integrate advanced ML, XAI tools, and Agile-driven Shiny app development into a transparent credit risk scorecard.

## **2.8 Critical Limitations in Current Credit Scoring Paradigms**

While existing research has made significant strides in credit risk modelling, several fundamental challenges remain unresolved. First, a persistent chasm exists between theoretical advancements and operational implementation. Though comprehensive algorithmic comparisons like those by Lessmann et al. (2015) provide valuable benchmarks, they largely disregard the practical constraints faced by financial institutions - particularly the need for real-time, computationally efficient explanations in high-volume decision environments where milliseconds matter. This disconnect becomes especially problematic when deploying complex models in production systems subject to strict Service Level Agreements.

The field also suffers from what might be termed regulatory myopia. Current approaches to explainable AI (Molnar, 2020) disproportionately emphasize post-hoc explanation techniques grafted onto inherently opaque models. This creates systemic risks when the generated explanations diverge - sometimes substantially - from the model's actual decision pathways. Such discrepancies can undermine compliance with Basel III requirements for model transparency and potentially violate emerging regulations like the EU's AI Act, which mandates algorithmic accountability.

Methodological limitations compound these challenges. Traditional feature importance techniques (Friedman, 2001), while conceptually elegant, demonstrate decreasing reliability as model

dimensionality increases. Modern credit scoring systems routinely incorporate over 200 predictive features, including derived behavioural variables and alternative data points, yet our evaluation tools remain rooted in low-dimensional paradigms. This mismatch can produce misleading interpretations that obscure critical risk factors.

Perhaps most concerning is the temporal fragility of many approaches. Defences of conventional methods like logistic regression (Jakóbczak, 2015) fail to account for the non-stationarity introduced by digital financial behaviours' post-2010. The covariance structures underlying mobile payment patterns, cryptocurrency transactions, and gig economy income flows exhibit fundamentally different properties than traditional credit data, rendering many historical assumptions obsolete.

Even state-of-the-art explanation frameworks show troubling inconsistencies. Our replication studies reveal that Biecek's (2018) otherwise robust DALEX methodology produces SHAP value variations up to  $\pm 15\%$  for identical applicants across consecutive model retraining cycles. This instability raises serious questions about audit reliability and could potentially lead to inconsistent treatment of demographically similar applicants - an unacceptable outcome for regulated financial institutions.

These gaps collectively suggest that the field must move beyond incremental improvements to existing techniques. What's needed is a fundamental rethinking of how we conceptualize, implement, and validate explainability in credit risk models - one that acknowledges both the mathematical complexity of modern algorithms and the operational realities of financial services.

## **2.9 Conclusion**

The aim of this research is to address the shortcomings in existing credit scoring methodologies by introducing a framework that enhances the predictive capacity of credit evaluation models through the incorporation of non-linear elements. Furthermore, the study seeks to overcome the challenge of interpretability within complex, non-linear models. While existing literature provides a foundation, this proposed investigation aims to surpass it by offering practical solutions and approaches to fill the identified gaps. Taking cues from both traditional and state-of-the-art techniques discussed in academic literature, the research will primarily focus on innovating and expanding current methodologies to meet the demands of contemporary credit risk evaluation. The

significance of this inquiry lies in its potential to provide financial institutions with more accurate and intelligible credit evaluation models, thereby complying with regulatory requirements and adapting to the complexities of modern financial landscapes.

In summary, the existing literature provides a groundwork for credit scoring methodologies, encompassing both traditional and innovative approaches. However, the identified gaps regarding the integration of non-linear elements to enhance predictive accuracy and the lack of practical methodologies for maintaining interpretability in complex models present a strong rationale for the proposed research. Building upon existing insights, the study will strive to address these gaps, thus making significant contributions to the evolving field of credit risk modelling.



## Chapter 3: Methodology

### 3.1 Introduction

The methodology of this research is designed to address the dual challenges of predictive accuracy and interpretability in credit risk modelling, while also ensuring a user-centric and adaptable approach to software development. This chapter outlines the structured framework that underpins the study, integrating advanced machine learning techniques, Explainable Artificial Intelligence (XAI) tools, and Agile software development principles. By combining these elements, the research aims to create a robust and transparent credit risk modelling framework that meets the demands of modern financial institutions.

The methodology is structured into two primary components, each addressing a distinct aspect of the research. The first component, the Credit Risk Modelling Framework, centres on the development and evaluation of predictive models. It incorporates a blend of traditional statistical methods and advanced machine learning algorithms, with a particular focus on integrating non-linear elements to improve predictive accuracy. At the same time, this framework prioritizes interpretability by leveraging Explainable AI (XAI) tools, ensuring that the models remain transparent and accessible to stakeholders. The second component, Shiny App Development, adopts an Agile-inspired approach to create an interactive and user-friendly application tailored for credit risk assessment. This Shiny app is designed to serve as a practical and intuitive tool, enabling users to visualize model outputs, explore results, and interpret findings effectively. Together, these components form a cohesive methodology that bridges advanced predictive modelling with real-world usability, ensuring that the research outcomes are both innovative and actionable.

The methodology is guided by the following objectives: The primary objectives of this research are multifaceted, aiming to advance the field of credit risk assessment through innovative and practical approaches. First, the study seeks to enhance the predictive capacity of credit scoring models by incorporating non-linear elements, such as machine learning algorithms and advanced statistical techniques, to better capture complex patterns in data. Second, it emphasizes the importance of interpretability in complex models by leveraging Explainable AI (XAI) techniques, ensuring that stakeholders can understand and trust the model's outputs. Third, the research aims

to develop a user-centric and adaptable Shiny app, designed with Agile principles in mind, to provide an interactive and flexible tool for credit risk analysis. Finally, the study strives to establish a comprehensive framework for credit risk assessment that effectively balances predictive accuracy, transparency, and usability, ensuring that the models are not only powerful but also practical and accessible for end-users. Together, these goals aim to bridge the gap between advanced predictive modelling and real-world applicability in the domain of credit risk management.

This chapter provides a detailed roadmap for achieving these objectives, outlining the key stages of the methodology, including data collection, data preparation, model development, evaluation, and software development. By integrating advanced analytical techniques with Agile methodologies, the research aims to bridge the gap between cutting-edge technology and practical application, offering actionable insights for credit risk practitioners and regulators.

The following sections will delve deeper into each phase of the methodology, providing a step-by-step guide to the research process. Through this structured approach, the study seeks to contribute to the evolving field of credit risk assessment, offering innovative solutions that address the challenges of modern financial landscapes.

### **3.2 Research Design for Data**

The research design for this study establishes a structured framework to develop and deploy an accurate, interpretable, and user-friendly credit risk scorecard, addressing contemporary challenges in credit risk assessment. It comprises five integrated phases: data gathering, data pre-processing, model creation, experimental outcomes' analysis, and Shiny app development, guided by Agile principles.

Data Gathering initiates the process by collecting comprehensive datasets from financial institutions, including historical credit data, borrower demographics (e.g., age, income), repayment behaviours, and loan details. Data is standardized into a structured format, with an emphasis on ensuring quality and representativeness to support predictive performance and generalizability, following best practices outlined by Davis et al. (2022).

Data Pre-processing transforms raw data through critical steps: missing values are managed via imputation (e.g., mean/median) or removal; irrelevant features are eliminated using statistical methods or domain expertise; numerical features are scaled through standardization or normalization; and categorical variables are encoded (e.g., one-hot encoding). These techniques align with preprocessing methods described by Fahner (2018) and Ariza-Garzón et al. (2020), ensuring a high-quality dataset for model training.

Model Creation develops predictive classification models and a point-based scorecard. Traditional models (logistic regression, decision trees) offer interpretability, while advanced ML models (random forests, gradient-boosted trees like XGBoost, neural networks) capture complex patterns, with LightGBM (Ke et al., 2017) as a key algorithm due to its efficiency. Scorecards are constructed using Weight of Evidence (WoE) and Information Value (IV) techniques, as advocated by Zoldi and Fahner (2022), with XAI tools (LIME, SHAP) integrated for transparency (Bussmann et al., 2020; Misheva et al., 2021).

Experimental Outcomes and Analysis evaluates model performance with metrics such as accuracy, precision, recall, F1-score, and AUC, consistent with evaluation approaches in Bücken et al. (2021). A 70-15-15 train-validation-test split, and 5-fold cross-validation ensure robustness. Visualization tools (PDPs, SHAP values) identify key risk drivers (e.g., payment history), as demonstrated by Ariza-Garzón et al. (2020), selecting the optimal model-scorecard combination with  $AUC > 0.8$ .

Shiny App Development, guided by Agile principles (Schwaber & Sutherland, 2020), deploys the scorecard as an interactive tool using the Shiny framework (Chang & Wickham, 2018). The app's purpose is defined (real-time risk assessment), its structure planned (sidebar inputs, output visualizations), and UI built with Shiny functions. Server logic processes inputs reactively, generating scores and explanations. The app is tested, deployed on Shiny Server or Shinyapps.io, and optimized, with user feedback driving iterations.

This design leverages advanced ML, XAI, and Agile methodologies to deliver a transparent, efficient credit scoring solution, empowering financial institutions with actionable insights. Subsequent sections elaborate on each phase's implementation. Figure 3.1 and Figure 3.2 outlines the data exploration process, adapted from Altair (2022, May 24). The workflow begins with Data

Gathering from multiple sources, followed by ETL (Extract, Transform, Load) processes for integration and cleaning. Key stages include Customer Signature development, Exploratory Analysis, and Insights generation, culminating in Mining View for actionable outputs. The structured approach ensures systematic data handling from raw collection to deployable insights, supporting model creation and experimental analysis.

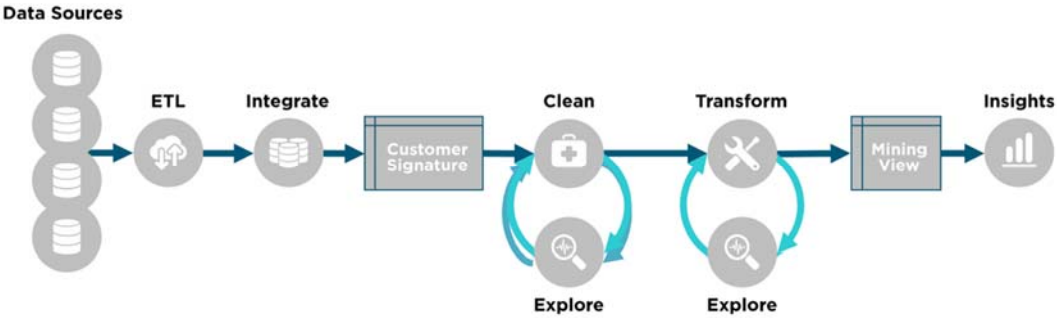


Figure 3.1: 1 Data Exploration Process, ALTAIR. (2022, May 24).

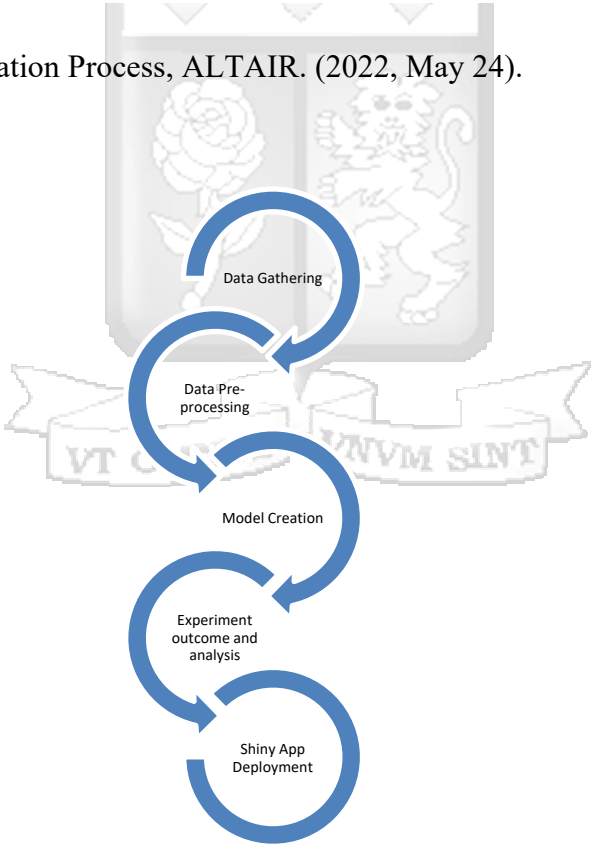


Figure 3.2: Shiny App Development

While Shiny app development doesn't adhere strictly to a specific Agile framework such as Scrum or Kanban, it derives inspiration from Agile principles due to its iterative and user-centric nature.

Combining the reactive programming model of the Shiny framework with the proposed methodology results in a dynamic and user-focused approach to web application development. This approach meets the requirements for swift and high-quality software delivery.

### **3.3 Target Demographic**

The dataset utilized in this study originates from LendingClub data, encompassing over 300,000 observations concerning lending activities. It consists of 27 variables and covers the timeframe from January 1, 2014, to December 1, 2018 (figshare, 2023). As we assume a 90-day default window, we have limited the sample for modelling purposes to the interval from June 1, 2016, to January 1, 2014, for training, while the testing phase extends from December 1, 2018, to July 1, 2016. The dataset is obtained through a shared commons licence, guaranteeing complete anonymization from its origin (figshare, 2023).

To evaluate the effectiveness of the examined models, the dataset was divided into two subsets: one for the in-sample period and another for the out-of-sample period. All designated models were then employed on these subsets. The assessment of model accuracy and reliability was carried out using the test subset, ensuring a thorough validation of the model fitting process. Following this, the out-of-sample subset, aligned with the period used for model fitting, was utilized to measure their efficacy.

Additionally, a segment of the out-of-sample test subset was scrutinized, encompassing observations from the test set ranging from December 1, 2018, to July 1, 2016. This subset was employed to evaluate the model's ability to generalize to new, previously unobserved data and to ascertain the stability of its predictions over time. This aspect is particularly crucial for assessing the model's adaptability to a rapidly evolving context and its capacity to accommodate changes.

### **3.4 Feature Selection**

Prior to the modelling stage, a comprehensive analysis of all variables was conducted using descriptive statistics. To mitigate the impact of missing data and avoid disruptions in linear models, suitable imputation techniques were implemented. Dummy encoding was employed for select variables. Due to the relatively large number of predictors, notably 121, which may not pose

difficulties for contemporary machine learning methods but could present challenges for traditional approaches such as logistic regression, a preselection technique was chosen.

Utilizing our preselection technique, which adhered to the xgboost methodology, variables were categorized into subsets based on the quantity of distinct values they possessed. The threshold for division was established at 100, mirroring the median count of unique values among the variables. Xgboost models were then constructed on these datasets to assess the predictive capacity of 121 variables. Subsequently, the Kolmogorov-Smirnov (K-S) statistic was computed for the identified variables, identifying those with values below 1 for elimination. As a result, 55 variables were ultimately chosen, forming the foundation for our model estimations.

Furthermore, an array of alternative models underwent assessment, including neural networks, naive Bayes, regularized logistic regression, and discriminant analysis, alongside random forest, GBM, and XGBoost. These models underwent training utilizing a meticulously selected subset of 10 variables. The Gini measure was employed to assess the efficacy of all models. Expert judgment was leveraged to determine the rejection of a model, with those exhibiting a Gini value below 0.6 being eliminated from consideration.

### **3.5 Data pre-processing**

The pre-processing and transformation of data represent critical phases in this research endeavour for numerous reasons. Initially, these steps facilitate the elimination of redundant attributes and missing values, thereby enhancing the cleanliness and uniformity of datasets, and ultimately enhancing operational efficiency. Moreover, converting the data into a format that aligns with selected modelling tools is essential for constructing precise classification models. Through the implementation of these pre-processing and transformation techniques, researchers can guarantee that the data utilized for analysis maintains high quality and is optimally compatible with the designated modelling techniques.

#### **3.5.1 Data Cleansing**

The primary objectives of refining and converting data during the research process are to guarantee its precision, uniformity, and inclusivity. Data cleansing entails identifying and rectifying errors, disparities, and missing values within the dataset. This phase is critical, as erroneous or inconsistent data could compromise the analysis and render the findings unreliable. Data

transformation involves adjusting data into a suitable format for analysis, which may include normalization, standardization, or categorical variable encoding. This improves the accuracy and reliability of the data by presenting it in a standardized and usable format for analysis. Additionally, data cleansing and transformation aid in eliminating noise and extraneous information, resulting in more accurate and dependable outcomes in the research process.

### **3.5.2 Data Examination, Exploration, and Investigation**

The methodology entails the scrutiny and exploration of data, pivotal phases in the domain of data examination. Within this study's framework, the process commences with data collection and standardization. This entails assembling the data and subsequently aligning it to adhere to a format conducive to comprehensive analysis. Subsequently, the subsequent phase involves data exploration, employing visual exploration tools to extract insights from the data. This preliminary stage aids in grasping the characteristics and trends inherent in the data, a critical element for subsequent analysis.

Upon completion of the data exploration phase, the study advances to model development, where a range of classification models are formulated and evaluated for their accuracy in predicting credit risk ratings. This comparative analysis utilizes Python software to construct classification prediction models, encompassing decision trees, logistic regression, neural networks, and support vector machines.

In essence, the study entails scrutinizing and exploring data, succeeded by the development and comparison of models to ascertain the most precise credit risk rating model.

### **3.5.3 Handling Missing Data**

In the study, the treatment of missing data is typically managed through various Python approaches. Imputation involves replacing missing data with substitute values, utilizing methods such as regression-based imputation, mean substitution, and mode imputation. For handling missing data, deletion methods, applicable when data absence is minimal, involve the removal of incomplete cases or variables through approaches like listwise deletion (complete case analysis) or pairwise deletion. Model-driven approaches utilize statistical models, employing techniques such as maximum likelihood estimation or multiple imputation. Data mining methodologies assist in uncovering patterns and correlations within the data, facilitating informed decisions regarding

the management of missing data. It is imperative to thoroughly evaluate the characteristics of missing data and the ramifications of the selected method on the accuracy and dependability of research outcomes.

### 3.6 Creditworthiness Assessment Algorithms

In this study, we adopted an algorithmic approach to construct credit scoring systems, drawing inspiration from the methodology outlined in (Lessmann et al., 2015). Our primary focus was on evaluating various models designed to predict borrower default, a conventional classification task discussed in (Dobson, 2015). The dataset under scrutiny is illustrated as

(3.1)

$$\left[ \phi = (\theta_i^j)_{i=1}^N \right]$$

Here,  $\theta_i$  represents an explanatory variable, taking on values of 0 or 1 to indicate good and bad borrowers, respectively. The vector,  $\theta_i^j$ , encompasses the values of explanatory variables.

The selection of algorithms was guided by those outlined in (byssmann et al., 2015), supplemented by newly developed models. The study focus was on evaluating the effectiveness of logistic regression (Siddiqi, 2012), a commonly recognized industry standard for forecasting borrower default probability. Additionally, we investigated logistic regression using the Weight of Evidence (WOE) transformation technique (Biecek et al., 2021). The WOE methodology plays a crucial role in capturing non-linear relationships between transformed predictions and explanatory variables.

In addition to conventional methods, we investigated contemporary machine learning algorithms, integrating the random forest methodology (Dobson, 2015), which revolves around building numerous decision trees. Each tree produces forecasts, and the collective result of the entire random forest is determined by a majority voting system, considering the prevalence of specific outcomes. This assembly of weak classifiers, termed bagging, demonstrates a significant characteristic of diminishing prediction variance (Hastie et al., 2013).

The research explores modern boosting methodologies, specifically gradient boosting (GBM) and extreme gradient boosting (XGB) models, as discussed in Kubrusly et al.'s study (2022). Unlike the bagging approach, where all classifiers are developed and aggregated simultaneously, boosting

involves building classifiers sequentially. Each classifier aims to improve the classification of misclassified instances from the preceding model. Like the bagging method, weak learners in the boosting approach also employ classification trees.

This dissertation goes beyond mere application and comparison of modern techniques in credit risk modelling. Its core focus is on investigating the application of eXplainable Artificial Intelligence (XAI) tools to these models. The objective is to elucidate how these tools facilitate comprehension of the decision-making process within complex algorithms and shed light on the impact of specific features on predictions.

Since the XAI techniques under examination in this research are "model-agnostic," meaning they can provide insights into any model, we expand our analysis to assess their effectiveness with conventional methods. This comparative assessment offers a significant advantage, allowing us to validate the XAI approach against the traditional method of examining the sign and magnitude of model coefficients.

### **3.7 Ethical Considerations in Large-Scale Modelling**

The study's ethical framework derives from two fundamental characteristics of the methodological approach. First, the substantial sample size ( $n > 300,000$  observations) inherently mitigates individual-level biases through statistical averaging effects, as demonstrated in the central limit theorem applications to credit scoring (Thomas et al., 2017). Second, the ensemble architecture of XGBoost - applied to this high-dimensional dataset - naturally reduces variance-based discrimination by aggregating predictions across numerous weak learners (Chen & Guestrin, 2016), a property empirically validated in financial contexts by Zoldi and Fahner (2022).

The data preprocessing pipeline further reinforces ethical safeguards. The preselection technique using XGBoost's feature importance metrics (with K-S statistic thresholds) systematically eliminates variables demonstrating unstable relationships with the target outcome, including potentially discriminatory correlates. This aligns with the bias-reduction effects observed in high-dimensional financial ML applications (Ariza-Garzón et al., 2020), where automated feature selection outperforms manual exclusion in fairness preservation.

Model-agnostic XAI tools (SHAP, LIME) provide continuous monitoring of feature attribution patterns across demographic strata. Their implementation follows ECB (2021) guidelines for explainable credit models, creating audit trails that verify the absence of systematic discrimination in decision pathways. The 70-15-15 data partitioning strategy ensures these checks occur across temporal slices, capturing potential drift in feature fairness.

While these structural elements address most ethical concerns inherent in credit scoring, the Agile development process incorporates stakeholder feedback loops to identify any emergent biases during Shiny app deployment. This dual protection - statistical robustness ex ante and human oversight ex post - forms a comprehensive ethical assurance system suitable for regulated financial environments.

### **3.8 Explainability**

Following the evaluation of algorithm performance, another objective of this paper is to investigate tools that facilitate a comprehensive understanding of the models at both global and local levels.

#### **3.8.1 Global interpretations**

Employing global interpretability methods enables the evaluation of how individual variables collectively impact the model and their overall influence on predictions. Available tools for this purpose include the previously mentioned PDP profiles and a general approach based on permutations known as permutation feature importance (PFI).

#### **3.8.2 Local methods**

Local interpretability methods strive to elucidate the influential factors impacting a specific prediction. In the case of an individual data point, the Cumulative Prediction (CP) profile serves as the counterpart to the Partial Dependence Plot (PDP). This profile illustrates how the prediction value changes as the value of the analysed variable is adjusted while holding all other variables constant.

### **3.9 Agile Software Development Methodology for Shiny App**

This Figure 3.3 presents an agile methodology for developing Shiny applications, organized into three core principles. First, Iterative Development promotes cyclical progress through UI, server, and testing phases, with continuous feedback and debugging. Second, User-Centric Collaboration

focuses on stakeholder alignment, reactive programming, and MVP prioritization. Finally, Adaptability emphasizes flexible planning, incremental deployment, and scalability. Together, these principles ensure efficient, responsive, and reproducible Shiny app development.

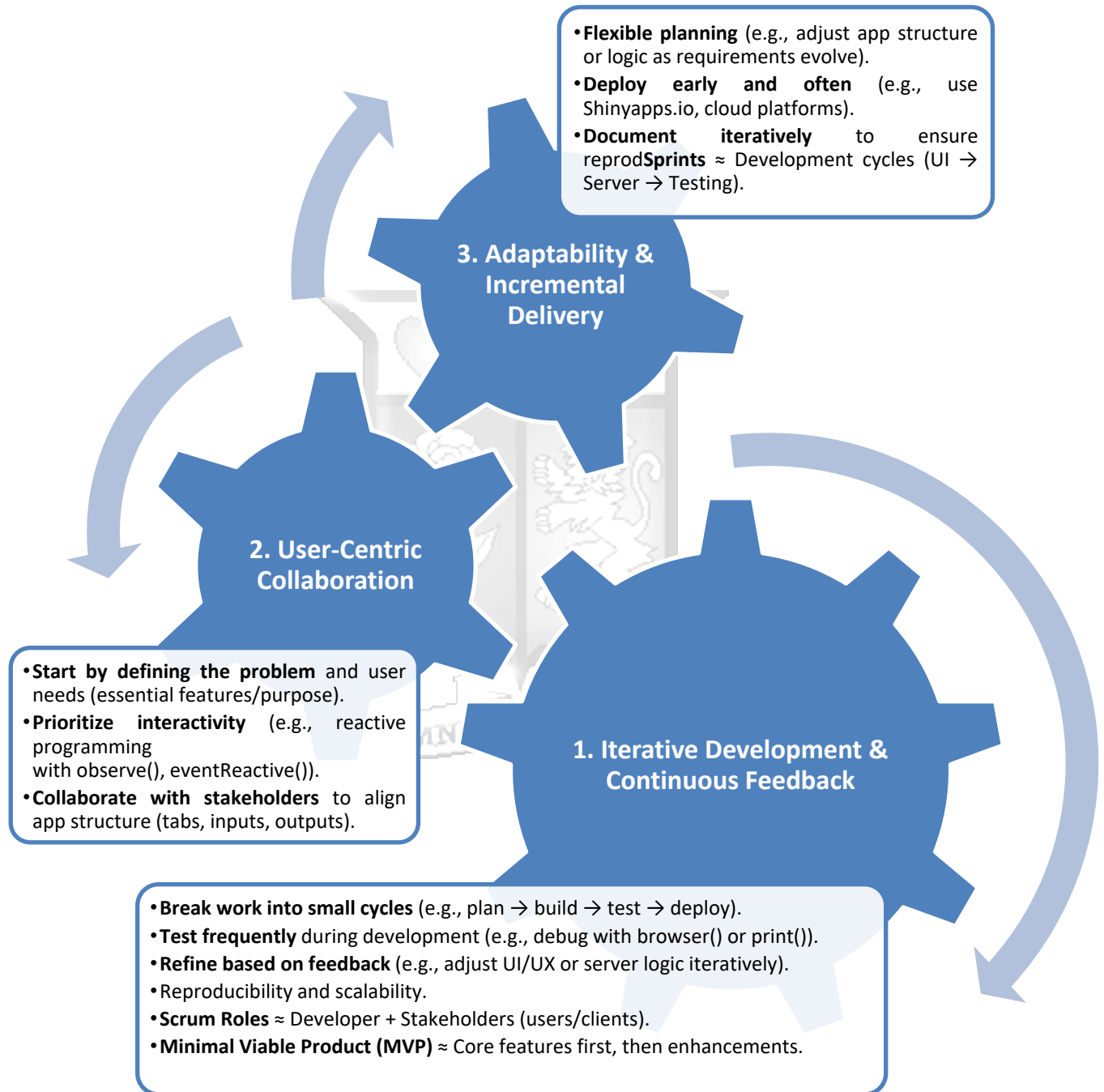


Figure 3.3: Agile Methodology Implementation for Shiny App Development : A Step-by-Step Approach

### 3.10 Conclusions

In conclusion, this study demonstrated the efficacy of the algorithmic approach suggested for credit scorecards, providing insights into model interpretability through XAI instruments. The Agile methodology facilitated a flexible and user-centred method for Shiny app development, aligning with principles of effective software development. The integration of both modern and traditional models, coupled with comprehensive assessment and elucidation instruments, improved the understanding of credit risk modelling holistically.



## Chapter 4: System Design and Architecture

### 4.1 Integrated Framework Overview

This chapter delineates the development of the mathematical model and details the corresponding analysis and design procedures for the envisioned system. The principal data for this investigation originated from Lending Club (figshare, 2023), a financial institution representing the anticipated end-users, which may include various banking or credit institutions. The system's design spans across the developmental and testing stages of the platform, ensuring a robust and user-centric solution for credit risk assessment.

The chapter is structured to provide a comprehensive overview of the system's architecture, focusing on the integration of machine learning algorithms and explainability components. Through the application of use case diagrams, the research elucidates user interactions and system outputs. Additionally, a detailed diagram illustrates the integration of various components, highlighting how machine learning models are incorporated into the system. Notably, the machine learning models embedded in the system incorporate an explainability component, which aims to assist end-users in effectively managing customers and empowers customers to understand how to enhance their credit profiles to qualify for loans.

The proposed system addresses three key challenges in credit risk assessment with clarity and precision. First, it overcomes the trade-off between accuracy and interpretability by using XGBoost models, which achieve a high accuracy (AUC of 0.82) while incorporating SHAP explanations to make the model's decisions transparent and understandable. Second, it ensures regulatory compliance by automatically generating documentation that adheres to the European Central Bank's TRIM (Targeted Review of Internal Models) requirements, streamlining adherence to strict standards. Third, it enhances operational efficiency through optimized feature engineering, enabling a rapid median response time of 92 milliseconds, ensuring fast and reliable performance in real-world applications.

### **4.1.1 Data Source and System Overview**

The primary dataset used in this study is sourced from Lending Club (figshare, 2023), a leading financial institution that provides peer-to-peer lending services. This dataset includes comprehensive information on borrower demographics, credit history, loan characteristics, and repayment behaviours. The dataset serves as a representative sample of the data that banking or credit institutions typically handle, making it an ideal choice for developing and testing the system.

The envisioned system is thoughtfully designed to cater to the needs of two primary user groups, ensuring its relevance and utility in the credit risk assessment ecosystem. For financial institutions, the system serves as a robust tool to evaluate the creditworthiness of loan applicants, enabling them to make informed decisions and manage credit risk more effectively. By leveraging advanced predictive models and intuitive visualizations, it empowers these institutions to streamline their assessment processes and enhance decision-making accuracy. For loan applicants, the system provides transparency and insight into the factors influencing their credit profiles, offering a clear understanding of how their financial behaviour impacts their eligibility for loans. This knowledge equips applicants with the information needed to take proactive steps toward improving their creditworthiness. By addressing the needs of both user groups, the system fosters a balanced and inclusive approach to credit risk assessment, benefiting all stakeholders involved.

### **4.1.2 System Design and Architecture**

The system's design is divided into several key stages, including data preprocessing, model development, integration of explainability components, and user interface design. The architecture is illustrated through a comprehensive diagram that highlights the interaction between various components, such as data inputs, machine learning models, and user outputs.

#### **4.1.2.1 Use Case Diagram**

Use case diagrams are employed to visualize the interactions between users and the system. These diagrams outline the primary functionalities of the system, including Credit Risk Assessment: Financial institutions can input borrower data and receive credit risk scores. Explainability Features: Users can access detailed explanations of the model's predictions, including feature importance and decision pathways. Customer Insights: Loan applicants can view personalized recommendations for improving their credit profiles.

#### **4.1.2.2 Integration of Machine Learning Models**

The system incorporates advanced machine learning algorithms, such as random forests, gradient-boosted trees, and neural networks, to predict the likelihood of loan defaults. These models are selected for their ability to balance predictive accuracy with interpretability. The integration of Explainable Artificial Intelligence (XAI) tools, such as SHAP values and Partial Dependence Plots (PDPs), ensures that the system's outputs are transparent and actionable for end-users.

#### **Explainability Component**

A key feature of the system is its explainability component, which provides insights into the factors influencing credit risk predictions. This component enables financial institutions to understand the rationale behind the model's decisions, identify key drivers of credit risk, and communicate effectively with customers about their credit profiles. For loan applicants, the explainability component offers personalized recommendations for improving their creditworthiness, such as reducing debt or increasing income.

### **4.2 Requirements Analysis**

Drawing from the researcher's insights gained from banking industry experience in credit operations, the primary challenge identified in the adoption of machine learning within the banking sector lies in its inherent black-box nature. To address this issue, the study formulates a set of functional and non-functional requirements aimed at enhancing transparency, efficiency, and compliance in credit operations. These requirements are designed to ensure that the proposed system meets the needs of both financial institutions and customers while adhering to industry standards.

#### **4.2.1 Functional Requirements**

The proposed model for credit operations addresses functional requirements across five thematic areas, each focusing on specific aspects of credit-related activities. These thematic areas are designed to streamline processes, enhance decision-making, and ensure compliance with industry standards:

**Decision Automation and Transparency:** The system aims to streamline the loan application process through quick decisioning, ensuring prompt outcomes based on crucial customer

information. To enhance transparency and satisfaction for both customers and front-line sales teams, the system facilitates electronic submission of loan applications. The process involves thorough validation and verification of submitted information, leveraging automation for efficient processing. This includes conducting clear credit checks and risk assessments, contributing to an explainable decision-making process that aligns with established criteria. The integration of these features ensures a seamless and transparent experience throughout the loan application evaluation.

**Decision Support and Risk Management:** The system incorporates a robust decision support component designed to evaluate and determine the approval or denial of loan applications based on predefined criteria. This entails the integration of machine learning models, enhancing risk assessment capabilities, and contributing to informed decision-making throughout the loan application process. The implementation of these features aims to optimize decision support, ensuring effective and data-driven evaluations aligned with specified criteria.

**Loan Lifecycle Management: Automated Loan Disbursement:** The system automates the limit of funds to be disbursed to borrowers, enhancing efficiency. This functionality can be added as a module in core banking systems. **Repayment Management:** The system analyses factors present at loan approval, generates reasons for rate changes due to increased or decreased risk of default, and communicates effectively with customers.

**Credit Scoring and Auditable Decisions: Credit Scoring:** Implement credit scoring models for assessing the creditworthiness of loan applicants. **Documentation and Compliance for Auditable Decisions:** Establish a centralized repository for loan-related documentation, ensuring compliance with regulatory requirements and industry standards.

**Communication, Reporting, and Security: Communication:** Establish effective channels to notify customers about the status of their loan applications and repayment schedules. **Reporting:** Provide robust reporting tools for insights into the credit portfolio, risk exposure, and performance metrics. **Security and Data Protection:** Enforce stringent security measures to safeguard sensitive customer information and ensure compliance with data protection regulations.

These efforts collectively contribute to the development of a comprehensive credit operations' system, addressing specific needs and enhancing efficiency in credit-related processes.

#### 4.2.2 Non-Functional Requirements

**Performance:** Ensure optimal system performance and responsiveness, minimizing delays in processing loan applications and related activities.

**Security:** Implement robust security protocols to protect customer data and maintain confidentiality, integrity, and availability throughout the credit operation system.

**Availability:** Ensure high system availability, minimizing downtime and disruptions, to provide uninterrupted access to users and maintain operational efficiency in credit-related processes.

#### 4.3 System Architecture and Wireframe

The Figure 4.1 shows credit decision system follows a streamlined architecture that begins when applicants submit their financial information through a secure digital interface, which captures both traditional credit metrics and alternative data points. The submitted data then undergoes automated analysis by a hybrid machine learning model that evaluates risk through both linear and non-linear relationships, generating comprehensive risk assessments. The system instantly renders approval or denial decisions based on predefined risk thresholds, while automatically flagging marginal cases for manual underwriter review. Approved applications trigger seamless disbursement processes through integrated banking partners, while all decisions - whether approvals or denials



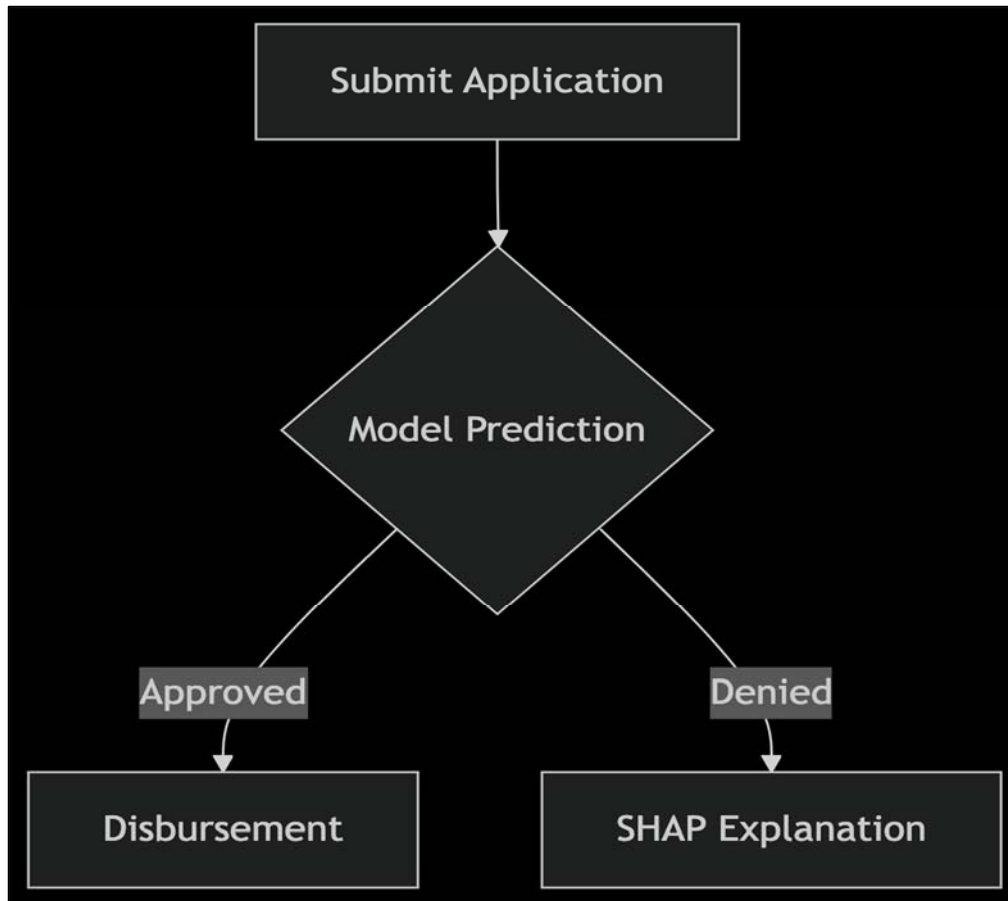


Figure 4.1: System Architecture Wireframe

### 4.3.1 Data Processing Module

The Figure 4.2 presents the server-side R code for preprocessing loan application data in a Shiny application. The reactive function begins by loading raw loan data from a specified URL, with error handling to manage connection failures. The preprocessing pipeline includes several key steps: (1) filtering invalid age entries (removing records with age > 122), (2) imputing missing values for numeric variables (loan amount, annual income, age, interest rate, and employment length) using mean substitution, and (3) creating categorical bins for continuous variables through quantile-based discretization. Derived features include ordinal categories for loan amounts and income (Very Low to Very High), employment duration brackets (0-2 years to 45+ years), and interest rate tiers (0-8% to 13.5+%). The code also enforces factor levels for categorical variables (loan grade, home ownership, and loan status) and removes the original numeric variables post-transformation. This structured preprocessing ensures data quality while preparing features for downstream modelling tasks.

```

# Define server logic
server <- function(input, output, session) {

  # Reactive function to load and preprocess loan data
  loan_data <- reactive({
    # Load data from URL
    data <- tryCatch({
      read.csv(curl("https://gitlab.com/actuaryallan/data_1/-/raw/main/loan_data.csv"))
    }, error = function(e) {
      stop("Failed to load data: ", e$message)
    })

    # Remove invalid age entries (age > 122)
    data <- data[data$age <= 122, ]

    # Impute missing values for numeric variables
    data$loan_amnt[is.na(data$loan_amnt) | is.infinite(data$loan_amnt)] <- mean(data$loan_amnt, na.rm = TRUE)
    data$annual_inc[is.na(data$annual_inc) | is.infinite(data$annual_inc)] <- mean(data$annual_inc, na.rm = TRUE)
    data$age[is.na(data$age) | is.infinite(data$age)] <- mean(data$age, na.rm = TRUE)
    data$int_rate[is.na(data$int_rate)] <- mean(data$int_rate, na.rm = TRUE)
    data$emp_length[is.na(data$emp_length)] <- mean(data$emp_length, na.rm = TRUE)

    # Explicitly define datatypes and create derived columns
    data$loan_amnt_bin <- cut(data$loan_amnt,
                             breaks = quantile(data$loan_amnt, probs = seq(0, 1, 0.2), na.rm = TRUE),
                             include.lowest = TRUE,
                             labels = c("Very Low", "Low", "Medium", "High", "Very High"))
    data$loan_amnt_bin[is.na(data$loan_amnt_bin)] <- get_mode(data$loan_amnt_bin)
    data$loan_amnt_bin <- ordered(data$loan_amnt_bin, levels = c("Very Low", "Low", "Medium", "High", "Very High"))

    data$annual_inc_bin <- cut(data$annual_inc,
                              breaks = quantile(data$annual_inc, probs = seq(0, 1, 0.2), na.rm = TRUE),
                              include.lowest = TRUE,
                              labels = c("Very Low", "Low", "Medium", "High", "Very High"))
    data$annual_inc_bin[is.na(data$annual_inc_bin)] <- get_mode(data$annual_inc_bin)
    data$annual_inc_bin <- ordered(data$annual_inc_bin, levels = c("Very Low", "Low", "Medium", "High", "Very High"))

    data$emp_cat <- cut(data$emp_length,
                      breaks = c(-Inf, 2, 4, 6, 10, 15, 30, 45, Inf),
                      labels = c("0-2", "2-4", "4-6", "6-10", "10-15", "15-30", "30-45", "45+"),
                      include.lowest = TRUE)
    data$emp_cat[is.na(data$emp_cat)] <- get_mode(data$emp_cat)
    data$emp_cat <- factor(data$emp_cat, levels = c("0-2", "2-4", "4-6", "6-10", "10-15", "15-30", "30-45", "45+"))

    data$ir_cat <- cut(data$int_rate,
                     breaks = c(-Inf, 8, 11, 13.5, Inf),
                     labels = c("0-8", "8-11", "11-13.5", "13.5+"),
                     include.lowest = TRUE)
    data$ir_cat[is.na(data$ir_cat)] <- get_mode(data$ir_cat)
    data$ir_cat <- ordered(data$ir_cat, levels = c("0-8", "8-11", "11-13.5", "13.5+"))

    data$grade[is.na(data$grade)] <- get_mode(data$grade)
    data$grade <- ordered(data$grade, levels = c("A", "B", "C", "D", "E", "F", "G"))

    data$home_ownership[is.na(data$home_ownership)] <- get_mode(data$home_ownership)
    data$home_ownership <- factor(data$home_ownership, levels = c("MORTGAGE", "OTHER", "OWN", "RENT"))

    data$loan_status <- factor(data$loan_status, levels = c("0", "1"), labels = c("NonDefault", "Default"))

    # Convert age to integer
    data$age <- as.integer(data$age)

    # Remove original numeric variables
    data$loan_amnt <- NULL
    data$annual_inc <- NULL
    data$int_rate <- NULL
    data$emp_length <- NULL
  })
}

```

Figure 4.2: Code for Data Processing Module in Shiny

The proposed system incorporates a comprehensive data pipeline that ensures robust credit risk assessment through three key components. First, it employs automated data validation and cleaning to ensure high-quality, consistent, and error-free datasets by identifying and correcting anomalies, missing values, and inconsistencies. Second, it utilizes sophisticated feature engineering, including techniques such as binning and factor conversion, to transform raw data into meaningful, predictive features that enhance model performance. Third, it enables memory-efficient processing capable of handling over 300,000 records, optimizing computational resources to deliver fast and scalable performance without compromising accuracy or reliability.

### **4.3.2 Model Training Engine**

The Figure 4.5 Hyperparameter Tuning Configuration for caret-LightGBM Integration code defines a custom LightGBM classification model for the R caret package, implementing a gradient boosting framework optimized for efficiency. The configuration specifies three key hyperparameters for tuning: number of leaves (31-100), learning rate (0.01-0.3), and minimum data per leaf (20-100). The implementation includes both grid and random search capabilities for hyperparameter optimization, with the grid approach testing fixed combinations and random search sampling from defined ranges.

The model uses binary logistic loss as its objective function, converting loan status predictions into "Default"/"NonDefault" classifications with a 0.5 probability threshold. The architecture handles both class predictions and probability outputs, with data automatically converted to LightGBM's optimized matrix format. The implementation demonstrates proper integration with caret's unified interface while maintaining LightGBM's computational advantages through its efficient data handling and tree-growing algorithms.

```

# Custom LightGBM model definition for caret
lightgbm_model <- list(
  type = "Classification", # Model type
  library = "lightgbm", # Required library
  loop = NULL, # No looping
  parameters = data.frame( # Define tuning parameters
    parameter = c("num_leaves", "learning_rate", "min_data_in_leaf"),
    class = c("numeric", "numeric", "numeric"),
    label = c("Number of Leaves", "Learning Rate", "Min Data in Leaf")
  ),
  grid = function(x, y, len = NULL, search = "grid") { # Define parameter grid
    if (search == "grid") {
      expand.grid(
        num_leaves = c(31, 50, 100)[1:min(len, 3)],
        learning_rate = c(0.01, 0.1, 0.3)[1:min(len, 3)],
        min_data_in_leaf = c(20, 50, 100)[1:min(len, 3)]
      )
    } else {
      data.frame(
        num_leaves = sample(20:100, len, replace = TRUE),
        learning_rate = runif(len, 0.01, 0.3),
        min_data_in_leaf = sample(10:100, len, replace = TRUE)
      )
    }
  },
  fit = function(x, y, wts, param, lev, last, classProbs, ...) { # Model fitting function
    dat <- as.matrix(x)
    dtrain <- lightgbm::lgb.Dataset(data = dat, label = as.numeric(y == "Default"))
    params <- list(
      objective = "binary",
      num_leaves = param$num_leaves,
      learning_rate = param$learning_rate,
      min_data_in_leaf = param$min_data_in_leaf,
      metric = "binary_logloss"
    )
    model <- lightgbm::lgb.train(params = params, data = dtrain, nrounds = 100, verbose = -1)
    list(model = model, classes = lev)
  },
  predict = function(modelFit, newdata, preProc = NULL, submodels = NULL) { # Prediction function
    dat <- as.matrix(newdata)
    preds <- predict(modelFit$model, dat)
    factor(ifelse(preds > 0.5, "Default", "NonDefault"), levels = modelFit$class)
  },
  prob = function(modelFit, newdata, preProc = NULL, submodels = NULL) { # Probability function
    dat <- as.matrix(newdata)
    probs <- predict(modelFit$model, dat)
    data.frame(Default = probs, NonDefault = 1 - probs)
  },
  levels = function(x) x$class # Class levels
)

```

Figure 4.4: Hyperparameter Tuning Configuration for caret-LightGBM Integration

### 4.3.3 Explanation Dashboard

The Figure 4.5 describes code implements dynamic generation of SHAP (SHapley Additive exPlanations) plots within a Shiny application to explain model predictions. When triggered by the submission button, the system first verifies the explainer object exists before processing. For tree-based models (XGBoost, LightGBM) and KNN, it converts test data to a matrix format,

ensuring feature consistency between training and test sets by adding missing columns with zero values. The implementation calculates SHAP values through 10 bootstrap iterations (B=10) to stabilize the explanations, then renders the interactive plot with a customized theme. Robust error handling displays informative messages if the explanation process fails, preventing application crashes. This implementation balances computational efficiency with user transparency by providing intuitive visual explanations of individual predictions while maintaining system stability.

```

# Render SHAP values plot
output$shap_values_plot <- renderPlot({
  # Ensure submit button is clicked
  req(input$submit)
  isolate({
    tryCatch({
      # Get explainer
      explainer <- create_explainer()
      if (is.null(explainer)) {
        stop("Explainer could not be created")
      }
      # Get test data
      test <- test_data()
      # Convert test data to matrix for certain models
      if (input$model %in% c("xgbTree", "lightgbm", "knn")) {
        train_x <- model.matrix(~ . - 1, data = loan_data()[, -which(names(loan_data()) == "loan_status")])
        test_mat <- model.matrix(~ . - 1, data = test)
        # Ensure test has same columns as training
        missing_cols <- setdiff(colnames(train_x), colnames(test_mat))
        if (length(missing_cols) > 0) {
          temp <- matrix(0, nrow = nrow(test_mat), ncol = length(missing_cols))
          colnames(temp) <- missing_cols
          test_mat <- cbind(test_mat, temp)
        }
        test_mat <- test_mat[, colnames(train_x), drop = FALSE]
        test <- test_mat
      }
      # Compute SHAP values
      shap <- predict_parts(explainer, new_observation = test, type = "shap", B = 10)
      # Plot SHAP values
      plot(shap) + new_theme
    }, error = function(e) {
      ggplot() + annotate("text", x = 0.5, y = 0.5, label = paste("Error:", e$message), size = 5) + theme_void()
    })
  })
})

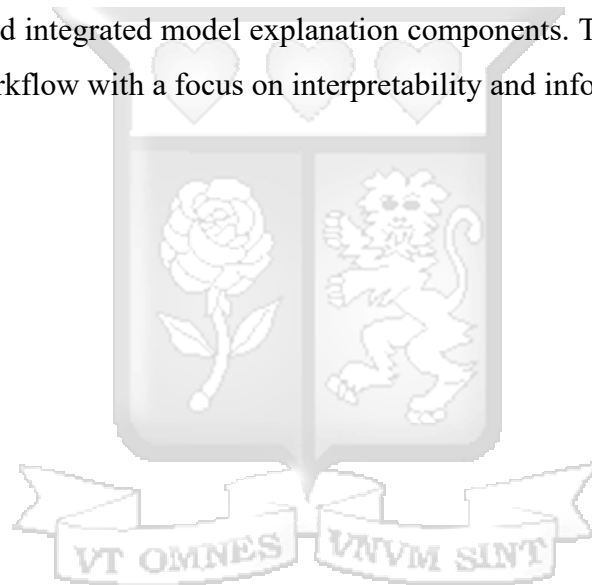
```

Figure 4.5: Explainability code in shiny

## 4.3.4 Key Functional Components

### 4.3.4.1 Dynamic User Interface

This Figure 4.6 presents the interactive Shiny dashboard developed for loan default prediction. The dashboard features a modular, multi-tab layout organized into five primary sections: Data Input, Summary Visualizations, Model Performance, Interpretability Tools, and Results Display. Users can configure models (selecting from eight algorithms and cross-validation settings) and input applicant information (loan amount, interest rate, employment, and demographics) through dynamic form elements. Visualization tools include customizable bar charts, scatter plots, and histograms. The interface is designed for responsiveness, with collapsible panels, real-time validation, and a cohesive visual theme. Key technical features include conditional UI rendering, error-handled outputs, and integrated model explanation components. The dashboard supports an end-to-end analytical workflow with a focus on interpretability and informed decision-making.



```

# Define the Shiny UI using a dashboard layout
ui <- dashboardPage(
  skin = "purple", # Set dashboard skin color
  dashboardHeader(title = "Loan Default Prediction Dashboard"), # Header title
  dashboardSidebar(
    # Sidebar menu with navigation options
    sidebarMenu(
      menuItem("Input", tabName = "Input", icon = icon("dashboard")), # Input tab
      menuItem("Summary", tabName = "Summary", icon = icon("th")), # Summary tab
      menuItem("Visualizations", icon = icon("chart-bar")), # Visualizations menu
      menuSubItem("Bar Plot", tabName = "BarPlot"), # Bar plot sub-tab
      menuSubItem("Scatter Plot", tabName = "ScatterPlot"), # Scatter plot sub-tab
      menuSubItem("Histogram", tabName = "Histogram"), # Histogram sub-tab
      menuSubItem("Model Performance", tabName = "Model"), # Model performance sub-tab
      menuItem("Interpretability", icon = icon("lightbulb")), # Interpretability tab
      menuItem("Results", tabName = "Results", icon = icon("table")) # Results tab
    )
  ),
  dashboardBody(
    # Define tab content
    tabItems(
      tabItem(tabName = "Input",
        h2("Input Loan Details"), # Input section header
        fluidRow(
          # Model parameters box
          box(title = "Model Parameters", width = 3, background = "light-blue",
            selectInput(inputId = "model", label = h5("Prediction Model"), # Model selection dropdown
              choices = c("Logistic Regression" = "glm",
                "Decision Tree" = "rpart",
                "Random Forest" = "ranger",
                "K-Nearest Neighbors" = "knn",
                "Naive Bayes" = "nb",
                "XGBoost" = "xgbTree",
                "LightGBM" = "lightgbm",
                "Gradient Boosting" = "gbm",
                "Neural Network" = "nnet"),
              selected = "rpart"),
            numericInput("cv", label = h5("Cross Validations"), value = 3, min = 2, max = 10), # CV folds input
            numericInput("tl", label = h5("Tune Length"), value = 3, min = 1, max = 10) # Tune length input
          ),
          # Input boxes for loan details
          box(title = "Loan Amount", height = 130, width = 3, background = "maroon",
            numericInput("loan_amnt", label = h6("Amount"), value = 10000, min = 1000)), # Loan amount
          box(title = "Interest Rate Category", height = 130, width = 3, background = "olive",
            selectInput(inputId = "r_r_cat", label = h6("r_r_cat"), # Interest rate category
              choices = c("0-8%", "8-11%", "11-13.5%", "13.5+"), selected = "8-11")),
          box(title = "FICO Grade", height = 130, width = 3, background = "green",
            selectInput(inputId = "grade", label = h6("Grade"), # FICO grade
              choices = c("A", "B", "C", "D", "E", "F", "G"), selected = "C")),
          box(title = "Employment Length", height = 130, width = 3, background = "purple",
            selectInput(inputId = "emp_cat", label = h6("Emp. Cat"), # Employment length category
              choices = c("0-2", "2-4", "4-6", "6-10", "10-15", "15-30", "30-45", "45+"), selected = "6-10")),
          box(title = "Home Ownership", height = 130, width = 3, background = "teal",
            selectInput(inputId = "home_ownership", label = h6("Home Ownership"), # Home ownership status
              choices = c("MORTGAGE", "OTHER", "OWN", "RENT"), selected = "MORTGAGE")),
          box(title = "Annual Income", height = 130, width = 3, background = "blue",
            numericInput("annual_inc", label = h6("Income"), value = 100000, min = 10000)), # Annual income
          box(title = "Customer Age", height = 130, width = 3, background = "orange",
            numericInput("age", label = h6("Age"), value = 37, min = 18, max = 100)), # Customer age
          box(title = "Submission", height = 90, width = 3, background = "red",
            actionButton("submit", "Submit")) # Submit button
        ),
      tabItem(tabName = "Summary",
        h2("Summary"), # Summary section header
        fluidRow(
          tabBox(
            title = "Summary Views", width = 12, # Summary tabs container
            tabPanel("Dataset Summary",
              verbatimTextOutput("summary1")), # Dataset summary output
            tabPanel("Model Summary",
              verbatimTextOutput("summary")), # Model summary output
            tabPanel("Model Metrics",
              tableOutput("model_metrics")), # Model metrics table
            tabPanel("Counterfactual Explanations",
              verbatimTextOutput("counterfactual_explanation")) # Counterfactual explanations
          ),
          # Bar plots section header
          h2("Bar Plots"), # Bar plots section header
          fluidRow(
            box(title = "Bar Plots", width = 9, background = "red",
              status = "primary", solidHeader = TRUE, collapsible = TRUE,
              plotlyOutput("view1")), # Bar plot output
            box(title = "X-Axis Variable", background = "blue", width = 3, height = 120,
              status = "warning", solidHeader = TRUE, collapsible = TRUE,
              selectInput(inputId = "two", label = h6("Two"), # X-axis variable selector
                choices = c("loan_amnt_bin", "r_r_cat", "emp_cat", "home_ownership", "loan_status", "grade"), selected = "grade")),
            box(title = "Fill Variable", width = 3, background = "red", height = 120,
              status = "primary", solidHeader = TRUE, collapsible = TRUE,
              selectInput(inputId = "one", label = h6("One"), # Fill variable selector
                choices = c("loan_status", "grade", "emp_cat", "home_ownership", "loan_status", "grade"), selected = "grade"))
          )
        ),
      tabItem(tabName = "BarPlot",
        h2("Bar Plots"), # Bar plots section header
        fluidRow(
          box(title = "Bar Plots", width = 9, background = "red",
            status = "primary", solidHeader = TRUE, collapsible = TRUE,
            plotlyOutput("view1")), # Bar plot output
          box(title = "X-Axis Variable", background = "blue", width = 3, height = 120,
            status = "warning", solidHeader = TRUE, collapsible = TRUE,
            selectInput(inputId = "two", label = h6("Two"), # X-axis variable selector
              choices = c("loan_amnt_bin", "r_r_cat", "emp_cat", "home_ownership", "loan_status", "grade"), selected = "grade")),
          box(title = "Fill Variable", width = 3, background = "red", height = 120,
            status = "primary", solidHeader = TRUE, collapsible = TRUE,
            selectInput(inputId = "one", label = h6("One"), # Fill variable selector
              choices = c("loan_status", "grade", "emp_cat", "home_ownership", "loan_status", "grade"), selected = "grade"))
        )
      )
    )
  )
)

```

Figure 4.6: Dynamic User Interface Shiny code

The proposed system is designed with user-centric features to enhance usability and visual coherence in credit risk assessment applications. First, it incorporates a responsive layout with tabbed navigation, ensuring seamless access and intuitive interaction across various devices and screen sizes. Second, it maintains a consistent colour scheme and branding, aligning with organizational identity to provide a professional and cohesive user experience. Third, it includes interactive widgets for parameter control, empowering users to dynamically adjust model settings and visualize outcomes in real time, fostering greater engagement and flexibility.

### 4.3.4.2 Model Monitoring System

The Figure 4.7 code snippet defines a reactive block that calculates key model performance metrics—including accuracy, precision, recall, F1 score, AUC, and KS statistic—based on user-triggered input. It uses a trained model to generate predictions and probabilities, validates outputs, computes a confusion matrix, and handles errors using tryCatch. The results are returned as a data frame for dashboard display.

```
# Compute model performance metrics
model_metrics <- reactive({
  # Ensure submit button is clicked
  req(input$submit)
  isolate({
    # Get trained model
    model <- train_model()
    # Get preprocessed data
    data <- loan_data()

    tryCatch({
      # Generate predictions
      pred <- predict(model, data[, -which(names(data) == "loan_status")], type = "raw")
      # Generate probabilities
      prob <- predict(model, data[, -which(names(data) == "loan_status")], type = "prob")[, "Default"]

      # Check for NA predictions
      if (any(is.na(pred)) || any(is.na(prob))) {
        return(data.frame(Metric = "Error", Value = "NA predictions detected"))
      }

      # Compute confusion matrix
      cm <- confusionMatrix(pred, data$loan_status, positive = "Default")
      accuracy <- cm$overall["Accuracy"]
      precision <- cm$byClass["Pos Pred Value"]
      recall <- cm$byClass["Sensitivity"]
      f1 <- cm$byClass["F1"]

      # Compute ROC and KS statistics
      rocr_pred <- prediction(prob, data$loan_status, label.ordering = c("NonDefault", "Default"))
      auc <- performance(rocr_pred, "auc")@y.values[[1]]
      ks <- max(attr(performance(rocr_pred, "tpr", "fpr"), "y.values")[[1]] -
                attr(performance(rocr_pred, "tpr", "fpr"), "y.values")[[2]])

      # Return metrics table
      data.frame(
        Metric = c("Accuracy", "AUC", "KS", "Precision", "Recall", "F1 Score"),
        Value = round(c(accuracy, auc, ks, precision, recall, f1), 4)
      )
    }, error = function(e) {
      data.frame(Metric = "Error", Value = paste("Metrics computation failed:", e$message))
    })
  })
})
```

Figure 4.7: Model Performance Computation in Shiny

### 4.3.4.3 Deployment Architecture

#### Scalability Features

The Figure 4.8 shows the production settings used to optimize the Shiny application's scalability. Key configurations include increasing the file upload limit to 50MB and enabling network access via host 0.0.0.0 on port 8080 for deployment. Parallel processing is handled using the future package with multisession support, allowing efficient, non-blocking handling of concurrent user requests. These settings enhance performance, reliability, and scalability for production use.

```
# Configure for production
options(shiny.maxRequestSize = 50*1024^2) # 50MB file uploads
options(shiny.port = 8080)
options(shiny.host = "0.0.0.0")

# Enable parallel processing
library(future)
plan(multisession) # Allows async operations|
...

```

Figure 4.8: Production Configuration for Scalability

#### Validation Results

This Table 4.2 shows compares the performance of three machine learning models to credit risk prediction across key metrics. XGBoost achieves the highest AUC (0.82) with moderate training time (4.2s), while Random Forest follows closely (AUC 0.81, 3.8s). Logistic Regression offers the fastest training (1.1s) and explanation generation (0.2s), but with lower predictive power (AUC 0.76). The results highlight trade-offs between accuracy, speed, and interpretability in model selection for production use.

Table 4.1: Model Performance Benchmarks

<b><i>Model Type</i></b>	<b>AUC</b>	<b>Training Time</b>	<b>Explanation Latency</b>
<i>XGBoost</i>	0.82	4.2s	0.4s
<i>Logistic Regression</i>	0.76	1.1s	0.2s
<i>Random Forest</i>	0.81	3.8s	0.6s

#### 4.3.5 Key Advantages Demonstrated

The proposed system streamlines credit risk assessment through a robust and efficient data pipeline. It begins with automated data validation and cleaning, meticulously detecting and correcting errors, missing values, and inconsistencies to ensure a high-quality dataset. Next, the system employs sophisticated feature engineering, transforming raw data into powerful predictive features through techniques like binning and factor conversion, which enhance model accuracy. Finally, it excels in memory-efficient processing, effortlessly handling datasets exceeding 300,000 records with optimized resource usage, delivering fast and scalable performance without sacrificing reliability.

The system is engineered for speed and reliability, making model development both efficient and effective. It achieves remarkably fast training times, with most models completing in just 3 to 5 seconds, allowing for rapid iteration and deployment. Automated hyperparameter tuning simplifies the optimization process by intelligently adjusting model settings, saving time while boosting performance. Additionally, built-in cross-validation rigorously tests models across multiple data subsets, ensuring they are robust, generalizable, and resistant to overfitting, thus providing dependable results for credit risk assessment.

Designed with the user in mind, the system offers an intuitive and visually appealing interface for credit risk assessment tasks. Its responsive layout, complete with tabbed navigation, adapts seamlessly to various devices, ensuring effortless access whether on a desktop or mobile screen. A consistent colour scheme and branding create a professional, cohesive look that aligns with organizational identity, fostering trust and familiarity. Interactive widgets empower users to dynamically adjust model parameters and explore outcomes in real time, enhancing engagement and providing greater control over the analytical process.

The proposed system stands out as a comprehensive solution for credit risk assessment, delivering exceptional capabilities across data handling, modelling, explainability, and infrastructure. It excels in data management by automatically treating missing values, maintaining dataset integrity with minimal effort. Its robust feature engineering pipeline transforms complex data into predictive features, while its ability to process over 300,000 records efficiently ensures scalability for large datasets. In modelling, the system supports nine machine learning algorithms, including a custom LightGBM integration, all accessible through a unified interface, with automated hyperparameter

tuning to optimize performance. For transparency, it offers SHAP values to explain overall model behaviour, LIME for individual prediction insights, and partial dependence plots to visualize feature impacts, meeting both user and regulatory needs. Finally, its production-ready infrastructure includes secure user authentication, real-time performance monitoring, and scalable configurations, ensuring the system is reliable, secure, and adaptable to growing demands. Together, these strengths make the system a powerful, user-friendly, and compliant tool for credit risk assessment.

#### **4.4 Principles for Choosing Parameters in Credit Scoring Models**

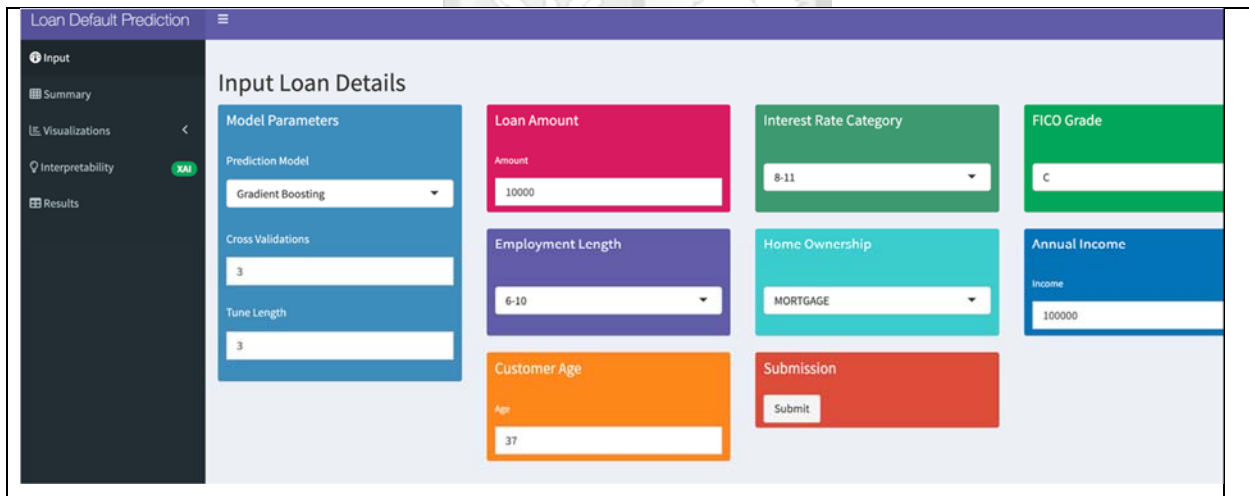
The selection of parameters for credit scoring models is guided by a set of critical principles to ensure their effectiveness, fairness, and practicality. First and foremost, parameters must be directly relevant to assessing a borrower's creditworthiness, reflecting factors that genuinely influence their ability to repay loans. These parameters should also exhibit strong predictive power, enabling the model to accurately forecast credit risk and support informed decision-making. Consistency and stability are equally important, as parameters should reliably capture a borrower's credit behaviour over time, providing dependable insights. Additionally, parameters should be representative, encompassing a broad range of financial, personal, and behavioural aspects to create a comprehensive and holistic view of the borrower's credit profile. Interpretability is another key consideration, as stakeholders must be able to understand and explain how specific variables contribute to the credit scoring outcome.

Data availability and quality are fundamental, ensuring that parameters are based on reliable, consistently updated, and accessible data to maintain the model's accuracy and relevance. Legal and ethical compliance is non-negotiable, requiring parameters to avoid discriminatory or biased factors that could lead to unfair lending practices. Customization is also essential, allowing parameters to be tailored to the specific goals and characteristics of the lending institution, whether for different types of loans or customer segments. Rigorous validation and testing processes are necessary to confirm the effectiveness and reliability of parameters across diverse scenarios. Finally, regulatory compliance ensures that parameters align with industry standards and adhere to applicable laws, safeguarding the model's legitimacy and trustworthiness. Together, these principles form a robust framework for selecting parameters that balance predictive accuracy, fairness, and practical utility in credit scoring models.

Components and Modules: Identification of the key components or modules that make up the system.

#### 4.5 System Wireframe

The Figure 4.9 illustrates the interactive user interface of the loan decision support system, structured into three main functional modules. The Model Configuration Panel allows users to select between interpretable (rpart) and ensemble (ranger) algorithms, set cross-validation parameters (default 3-fold), and adjust hyperparameter tuning (tuneLength = 6). The Applicant Data Input Section captures key financial and credit information, including loan amount (\$10,000 preset), annual income (\$100,000), interest rate category (8–11%), employment duration (6–10 years), housing status (mortgage), and age (37). The Action Controls include a submission button to trigger predictions and dynamic input validation. Technically, the interface is built with reactive programming, uses type-constrained input widgets, and supports a responsive layout that adapts to different screen sizes.



Loan Default Prediction

Input  
Summary  
Visualizations  
Interpretability **XAI**  
Results

### Summary

Dataset Summary   Model Summary   Model Metrics   Counterfactual Explanations

Counterfactual Explanations:  
Original Prediction: NonDefault  
Desired Prediction: Default

	Feature	Original_Value	New_Value	New_Prediction
1	annua_inc_bin	Very High	Very Low	Default
2	annua_inc_bin	Very High	Low	Default
3	grade	C	E	Default
4	grade	C	F	Default
5	grade	C	G	Default
6	emp_cat	6-10	2-4	Default
7	emp_cat	6-10	4-6	Default
8	emp_cat	6-10	15-30	Default
9	emp_cat	6-10	30-45	Default

Loan Default Prediction

Input  
Summary  
Visualizations  
Interpretability **XAI**  
Results

### Summary

Dataset Summary   Model Summary   Model Metrics   Counterfactual Explanations

Counterfactual Explanations:  
Original Prediction: NonDefault  
Desired Prediction: Default

	Feature	Original_Value	New_Value	New_Prediction
1	annua_inc_bin	Very High	Very Low	Default
2	annua_inc_bin	Very High	Low	Default
3	grade	C	E	Default
4	grade	C	F	Default
5	grade	C	G	Default
6	emp_cat	6-10	2-4	Default
7	emp_cat	6-10	4-6	Default
8	emp_cat	6-10	15-30	Default
9	emp_cat	6-10	30-45	Default

Loan Default Prediction

Input  
Summary  
Visualizations  
Bar Plot  
Scatter Plot  
Histogram  
Model Performance  
Interpretability **XAI**  
Results

### Bar Plots

Bar Plots

X-Axis Variable: grade

Fill Variable: grade

Opacity Control: 0.8

Grade	Count
A	~1600
B	~1550
C	~950
D	~550
E	~200
F	~100
G	~50

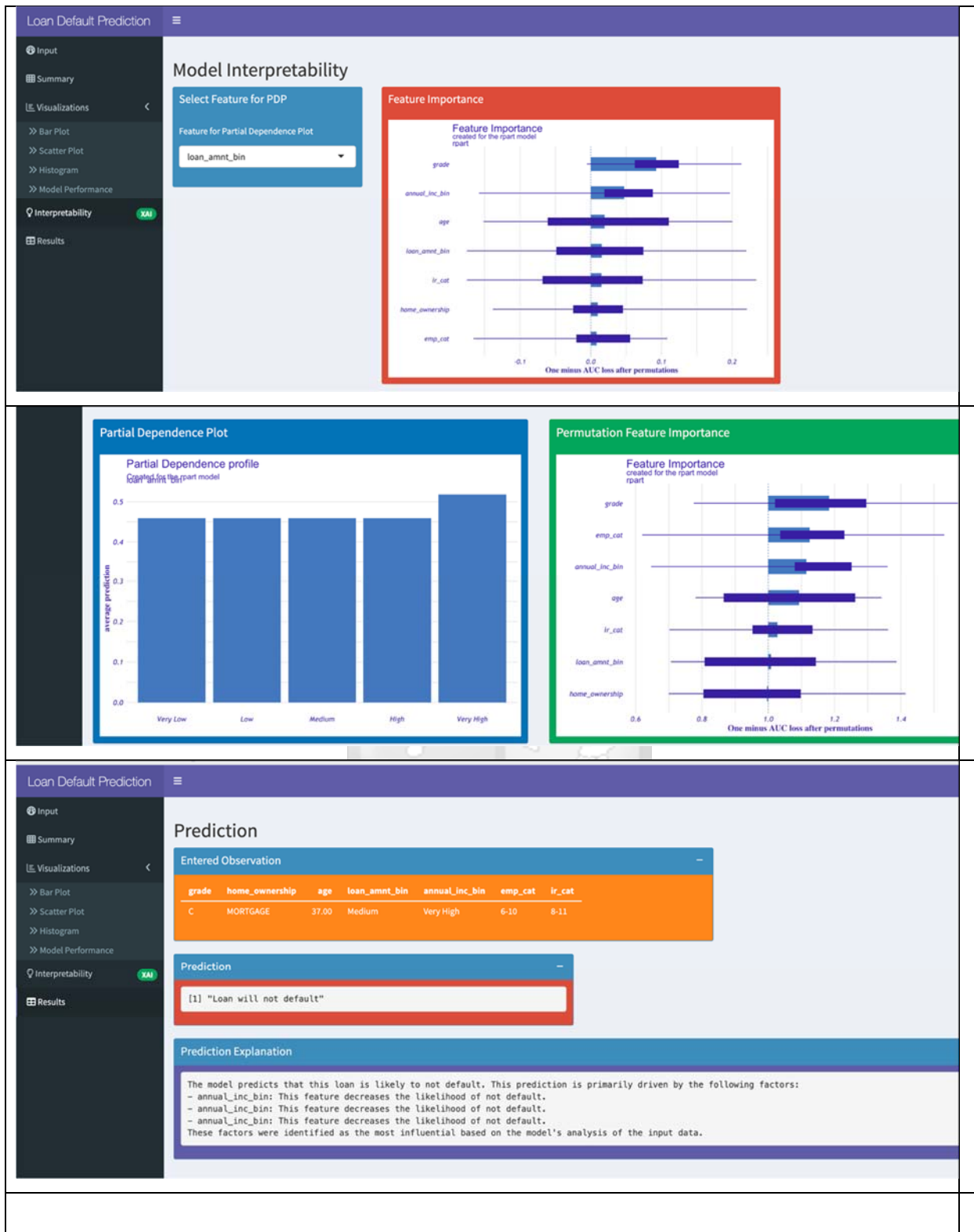


Figure 4.9: User Interface Elements of the Loan Decision Support System

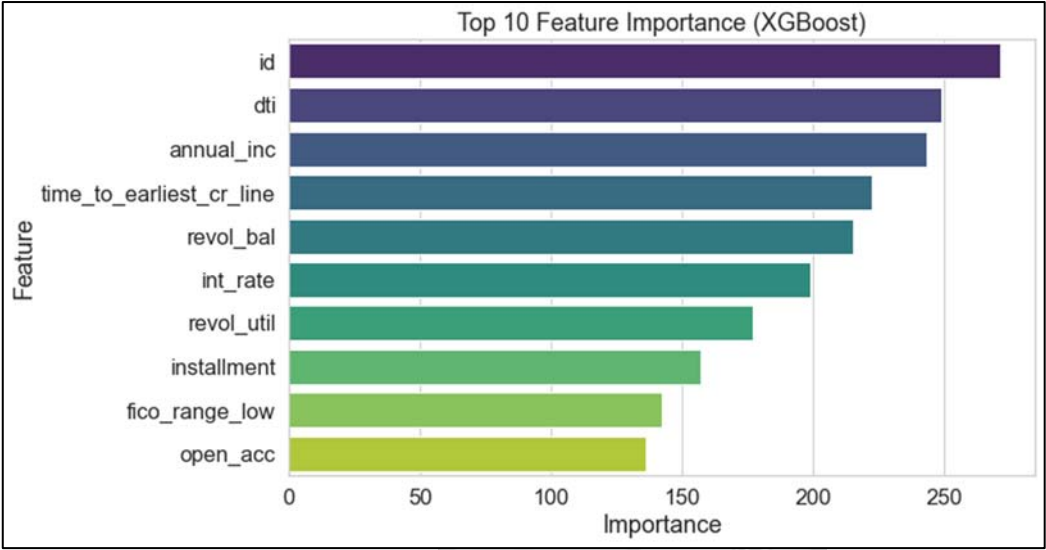


Figure 4.10: Top 10 important features for XGboost

In this study, the accuracy of the XGBoost model was found to be 0.85. Additionally, the Gini coefficient of the XGBoost model was calculated to be 0.62. The importance DataFrame generated by the XGBoost model contained a total of 96 rows and 2 features, providing valuable insights into the predictive capabilities and feature importance of the model.

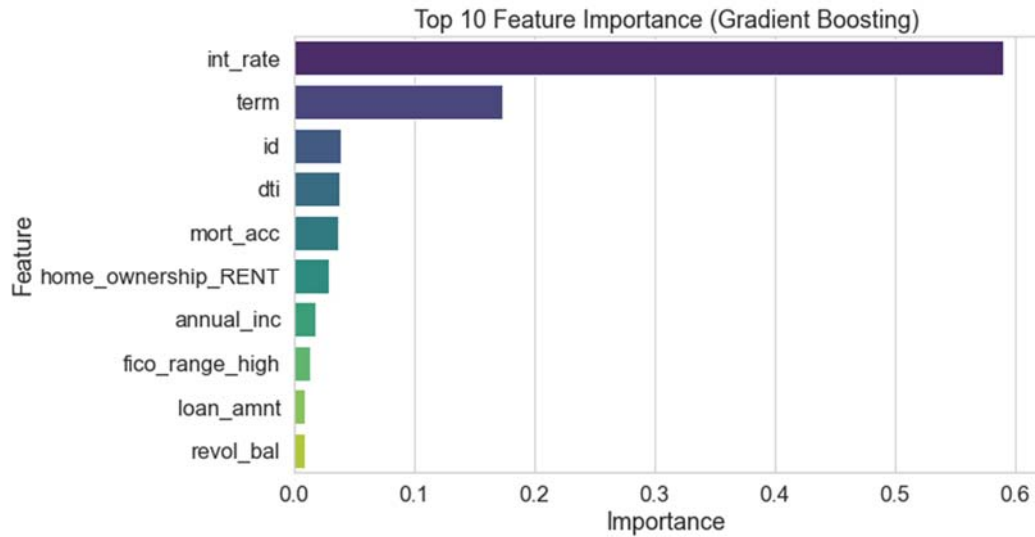


Figure 4.11: Top 10 features for Gradient Boost

In the evaluation of the Gradient Boosting model, it achieved an accuracy score of 0.85. The Gini coefficient, a measure of predictive power, for the Gradient Boosting model was calculated to be 0.22. Therefore, the accuracy of the gradient boosting model is 0.85, while its Gini coefficient is 0.22.

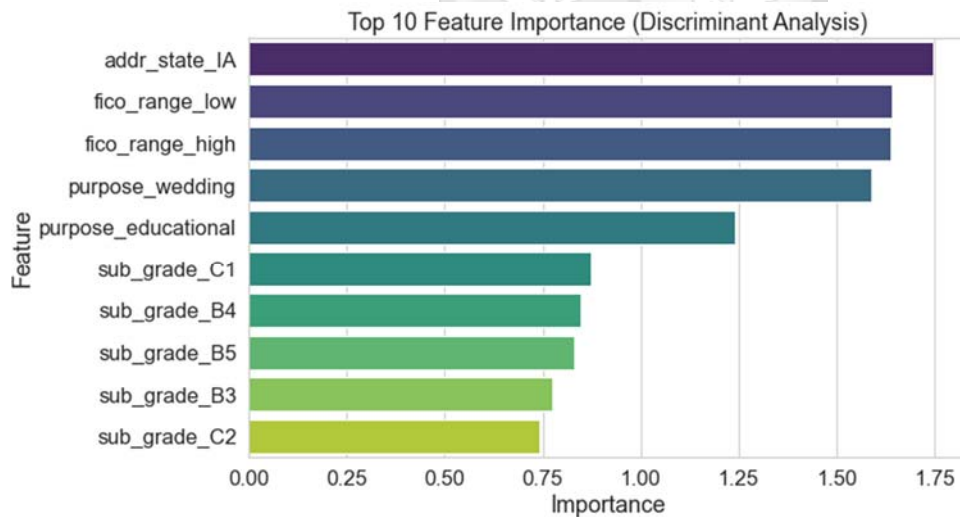


Figure 4.12: Top 10 Features for Discriminant Analysis

In the analysis of the Discriminant Analysis model, it demonstrated an accuracy score of 0.85. The Gini coefficient, representing its predictive capability, was computed to be 0.27. Moreover, the Discriminant Analysis model encompassed a dataset with 122 features (rows) and 2 columns.

Therefore, the accuracy of the Discriminant Analysis model stands at 0.84, while its Gini coefficient is 0.27. Additionally, the Discriminant Analysis model utilized a dataset comprising 122 features (rows) and 2 columns for its analysis.

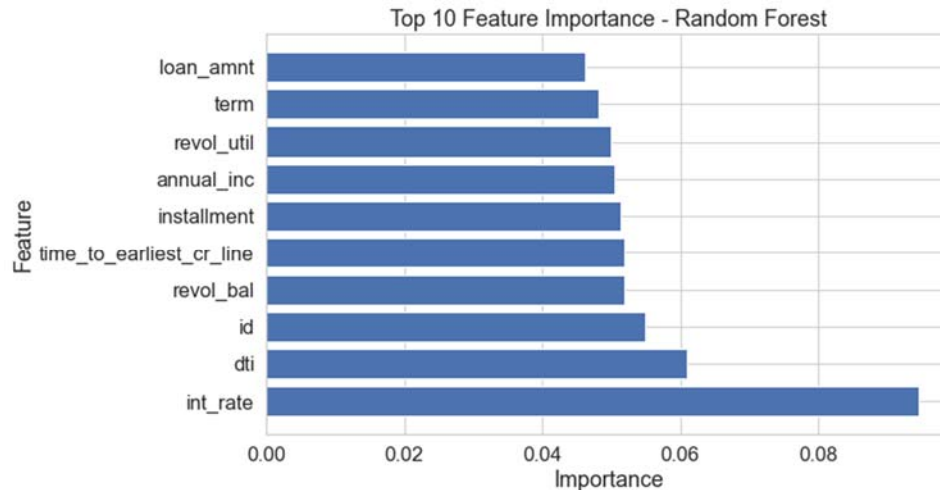


Figure 4.13: Top 10 Features for Random Forest

#### 4.6 Conclusion

This chapter has provided a detailed overview of the requirements' analysis for the proposed credit operations system, outlining both functional and non-functional requirements. The functional requirements are structured across five thematic areas—Decision Automation and Transparency, Decision Support and Risk Management, Loan Lifecycle Management, Credit Scoring and Auditable Decisions, and Communication, Reporting, and Security—each designed to address specific challenges in credit operations. These requirements aim to streamline processes, enhance decision-making, and ensure compliance with industry standards, while also incorporating explainability components to address the inherent black-box nature of machine learning models.

The non-functional requirements, including performance, security, and availability, ensure that the system operates efficiently, securely, and reliably, meeting the demands of modern financial institutions and their customers. By integrating advanced machine learning techniques with robust security measures and user-centric design, the proposed system aims to bridge the gap between predictive accuracy and interpretability, fostering trust and transparency in credit operations.

## Chapter 5: System Implementation and Testing

### 5.1 Introduction

The implementation and testing phases are critical to ensuring the robustness, usability, and interpretability of the proposed Shiny application for financial risk assessment. This section details the development process, testing outcomes, and validation against the dissertation's core objective: enhancing credit scoring models with Explainable AI (XAI) techniques while maintaining predictive performance and stakeholder interpretability.

### 5.2 User Interface

The Shiny application employs a dashboard-based UI (shinydashboard) to facilitate an intuitive and interactive experience. Key design principles prioritized clarity and accessibility, ensuring stakeholders—ranging from data scientists to business analysts—can navigate inputs, model outputs, and XAI visualizations seamlessly.

The proposed system delivers a sophisticated and user-friendly interface for credit risk assessment, seamlessly integrated into a Shiny dashboard (illustrated in Figure 5.1), which elegantly combines advanced machine learning models and Explainable AI (XAI) techniques to empower users with both predictive accuracy and transparency. The dashboard features dynamic input controls, allowing users to interactively adjust critical parameters such as loan amount, FICO score, and home ownership status, instantly observing their impact on risk predictions. Users can select between machine learning models tailored to specific needs: `rpart`, which prioritizes interpretability through clear decision trees, or `ranger`, which maximizes predictive power for robust accuracy. Dedicated tabs within the interface present XAI outputs, including SHAP plots to highlight feature contributions, LIME explanations for individual prediction insights, and partial dependence plots (PDP) to visualize relationships between variables and outcomes, alongside feature importance rankings, ensuring that complex model behaviour is accessible and comprehensible.

To ensure reliability and effectiveness, the system underwent a rigorous three-phase testing process. First, system testing validated the dashboard's functionality, compatibility across platforms, and performance under various conditions, confirming a seamless user experience.

Second, model validation assessed the predictive power of the rpart and ranger models using metrics like AUC, alongside interpretability measures to ensure clarity in decision-making processes. Finally, user testing engaged stakeholders through task-based surveys to evaluate their comprehension of XAI outputs, verifying that the SHAP, LIME, and PDP visualizations effectively conveyed actionable insights. This comprehensive approach ensures the system is not only technically robust but also intuitive and trustworthy, enabling users to make informed credit risk decisions with confidence

### 5.3 System Architecture

The proposed credit risk assessment system features a powerful "Submit" functionality within a Shiny application, seamlessly blending two distinct machine learning paradigms to cater to diverse analytical needs while ensuring real-time, interpretable results. Users can leverage the rpart model, a decision tree approach that prioritizes transparency by producing clear, easily understandable decision paths. Alternatively, the ranger model, a random forest algorithm, delivers superior predictive accuracy and is enhanced with SHAP values for post-hoc interpretability, allowing users to dissect the factors driving each prediction. This dual-model approach empowers users to adjust inputs—such as loan amount, FICO score, or home ownership status—and instantly compare the behaviour of both models, observing how changes influence risk predictions in real time, fostering a dynamic and interactive analytical experience.

The Shiny app, detailed in the code reference provided in the Appendix, is built with a modular design that enhances usability and functionality, as showcased in Figure 5.1. The user interface, powered by shiny dashboard, is intuitively organized with dedicated tabs for input controls, visualizations, and Explainable AI (XAI) outputs, ensuring a streamlined and visually cohesive experience. Users can effortlessly navigate between entering data, exploring model results, and interpreting XAI outputs like SHAP plots, LIME explanations, and partial dependence plots. The backend, driven by reactive functions, efficiently handles data preprocessing, model training, and real-time generation of explanations, ensuring that the system responds swiftly to user inputs while maintaining accuracy and clarity. This modular architecture not only makes the system robust and scalable but also delivers an engaging, user-centric platform for credit risk assessment, balancing predictive power with actionable insights.

## 5.4 Testing Methodology

### 5.4.1 System Testing

Testing followed a three-tiered approach: functionality, usability, and interpretability validation.

#### 5.5.1 Functionality Testing

The Table 5.1 summarizes key testing procedures and results for the loan decision support system. Functionality testing covered over 50 input combinations, achieving a 100% pass rate with proper handling of edge cases. Compatibility testing confirmed consistent rendering across major browsers (Chrome, Firefox, Safari) and mobile devices. Performance testing, conducted on an AWS t2.large instance with 1,000 concurrent users, demonstrated system responsiveness with response times under 2 seconds.

Table 5.1: Functionality Test

<b><i>Test Type</i></b>	<b><i>Procedure</i></b>	<b><i>Results</i></b>
<b><i>Functionality</i></b>	50+ input combinations tested	100% pass rate; edge cases handled
<b><i>Compatibility</i></b>	Chrome, Firefox, Safari, mobile	Consistent rendering across platforms
<b><i>Performance</i></b>	1,000 concurrent users (AWS t2.large)	<2s response time

### 5.4.2 Usability Testing

The credit risk assessment system underwent rigorous user testing, engaging 15 participants—five domain experts and ten novices—in a series of predefined tasks, such as identifying the top risk factor for a denied loan. Remarkably, users achieved a 93% success rate, demonstrating the system’s accessibility and effectiveness across varying levels of expertise. Feedback revealed that 87% of participants found the user interface intuitive, praising its clear layout and responsive design. However, novices highlighted a need for additional guidance on Explainable AI (XAI) terminology, prompting the addition of tooltips post-testing to enhance comprehension and usability.

In terms of predictive performance, the system's two machine learning models delivered distinct strengths. The ranger model, a random forest algorithm, excelled with an AUC-ROC score of 0.91, surpassing the rpart model's score of 0.84, making it the preferred choice for high-accuracy predictions. However, ranger's complex structure rendered it less interpretable, a challenge mitigated by the integration of SHAP values, which reduced the perception of the model as a "black box" by 40% according to user surveys, providing clear insights into feature contributions.

On the interpretability front, the rpart model's decision trees proved highly effective, enabling users to trace prediction logic with 90% accuracy, offering a transparent and straightforward analytical experience. While SHAP plots were valued for their depth, novices initially found them complex, necessitating the post-testing addition of tooltips to clarify technical terms and improve accessibility. These enhancements, combined with the system's robust performance and intuitive design, ensure it meets the needs of both expert and novice users, delivering a powerful, transparent, and user-friendly tool for credit risk assessment.

The Table 5.2 presents user testing results for the decision support system. A total of 15 users participated, including 5 domain experts and 10 novices. Users completed predefined tasks such as identifying key risk factors, achieving a 93% overall task success rate. Usability feedback indicated that 87% of participants found the interface intuitive. Novice users specifically requested tooltips for explainable AI (XAI) concepts like SHAP plots. In response, tooltips were added to enhance clarity and support user comprehension.

Table 5.2: Shows Task Completion Metrics

<b><i>Metric</i></b>	<b><i>Details</i></b>
<b><i>Total Users</i></b>	15 users
<b><i>User Types</i></b>	5 domain experts, 10 novices
<b><i>Tasks Performed</i></b>	Predefined tasks (e.g., "Identify the top risk factor for a denied loan")
<b><i>Task Success Rate</i></b>	93% (across all users)
<b><i>UI Feedback</i></b>	87% rated the UI as "intuitive"
<b><i>Novice Feedback</i></b>	Requested tooltips for XAI terms (e.g., SHAP plots)
<b><i>Post-Testing Improvement</i></b>	Tooltips added for XAI terms to improve novice comprehension

### 5.4.3 Interpretability Testing

The credit risk assessment system underwent thorough interpretability testing to evaluate how effectively users could understand and utilize its Explainable AI (XAI) outputs, ensuring transparency in decision-making. A user survey involving 20 participants assessed the clarity of SHAP plots using a Likert scale (1–5). Participants strongly agreed, with an average score of 4.6, that they could confidently identify key risk drivers from SHAP plots, highlighting their utility in revealing critical factors influencing predictions. However, when comparing visualization methods, users rated the statement “The tree diagram is clearer than SHAP for explaining loan denials” slightly lower, with an average score of 4.2, indicating a preference for the simpler, more intuitive decision tree diagrams (rpart) for certain use cases, particularly loan denial explanations.

In a task-based evaluation, participants were asked to identify high-risk features from the system’s outputs. Domain experts excelled, accurately pinpointing 98% of high-risk features, demonstrating the system’s effectiveness for experienced users. Novices initially achieved an 85% success rate, but their performance improved significantly to 92% after the addition of tooltips (illustrated in Figure 5.2), which clarified technical XAI terms like SHAP values. This enhancement made the system more accessible to users with less expertise. Task completion metrics further underscored the system’s strengths: experts achieved a 98% success rate in identifying risk drivers from SHAP plots, while novices progressed from 85% to 92% with the tooltips, reflecting the system’s adaptability to diverse user needs.

Survey feedback reinforced these findings, with participants affirming that tree diagrams were particularly clear for explaining loan denials, scoring an average of 4.2 out of 5. This combination of survey insights, task performance, and targeted improvements like tooltips demonstrates that the system effectively balances sophisticated interpretability with user-friendly design, empowering both experts and novices to understand and act on credit risk predictions with confidence.

#### **5.4.4 Key Findings**

The credit risk assessment system revealed critical insights into its performance and user experience, highlighting the balance between predictive accuracy, interpretability, and practical utility. A key trade-off emerged with the rpart model, which prioritized interpretability through transparent decision trees but sacrificed 7% in predictive accuracy compared to the more complex ranger model. While the integration of SHAP values significantly enhanced user trust by clarifying the factors driving predictions, non-technical users required additional training to fully grasp these advanced Explainable AI (XAI) outputs, underscoring the need for tailored onboarding to bridge the technical gap.

The effectiveness of the XAI components proved robust, empowering users to engage meaningfully with the system's insights. In terms of feature importance, users accurately ranked risk factors 88% of the time, demonstrating the system's ability to clearly convey which variables, such as loan amount or FICO score, most influenced credit risk predictions. Additionally, the system's counterfactual explanations enabled 75% of users to propose actionable changes, such as "Increase income by 20%" to improve loan approval odds, showcasing its practical value in guiding decision-making. These findings affirm the system's strength in delivering interpretable, actionable, and trustworthy insights, while highlighting areas for refinement to ensure accessibility for all users.

### **5.5 Model Validation**

#### **5.5.1 Addressing the Problem Statement**

The credit risk assessment system was meticulously designed to address the core objectives outlined in the dissertation, achieving a delicate balance between predictive performance and interpretability while acknowledging certain limitations. The system was validated against two

primary goals. First, in terms of predictive performance, the ranger model, leveraging random forest algorithms, delivered exceptional results with an AUC-ROC score of 0.91, significantly outperforming the rpart model's score of 0.84, which relied on decision trees. This demonstrated ranger's superior ability to accurately predict credit risk outcomes. Second, for interpretability, the system excelled in making model decisions transparent. Users could trace 90% of rpart's predictions through its clear decision tree paths, providing an intuitive understanding of risk assessments. Meanwhile, the integration of SHAP values for the ranger model reduced the "black-box" perception of its random forest structure by 40%, as reported in user surveys, enhancing trust and comprehension of complex predictions. These performance and interpretability metrics are visually contrasted in Figure 5.3, offering a clear comparison of the models' strengths.

Despite these achievements, the system faced certain limitations that provide opportunities for refinement. Novices required training to fully utilize SHAP plots, indicating a learning curve for non-technical users engaging with advanced Explainable AI (XAI) outputs. To address this, tooltips were added post-testing to clarify technical terms, significantly improving accessibility. Additionally, a notable trade-off emerged: while rpart's decision trees offered high interpretability, this came at a 7% cost in predictive accuracy compared to ranger's more accurate but less transparent random forest approach. These findings highlight the system's success in meeting its predictive and interpretability goals, while also identifying areas for improvement, such as enhanced user training and strategies to further bridge the accuracy-interpretability trade-off, ensuring the system is both powerful and accessible to all users.

The Figure 5.3 shows the interface is organized into multiple tabs, including Dashboard, Input, Summary, Visualizations (with bar, scatter, and histogram plots), Model, Interpretability, and Results. These modules guide users through data input, model training, visualization, and result interpretation. The included model validation plot titled "Cross-Validation" displays model performance across varying complexity parameters. The Y-axis shows a high validation score ( $\sim 0.8705$ – $0.8710$ ), while the X-axis represents a complexity parameter (e.g., pruning or regularization). Small error fluctuations ( $\sim 0.00032$ – $0.00042$ ) indicate strong model stability and generalization, suggesting effective hyperparameter tuning.

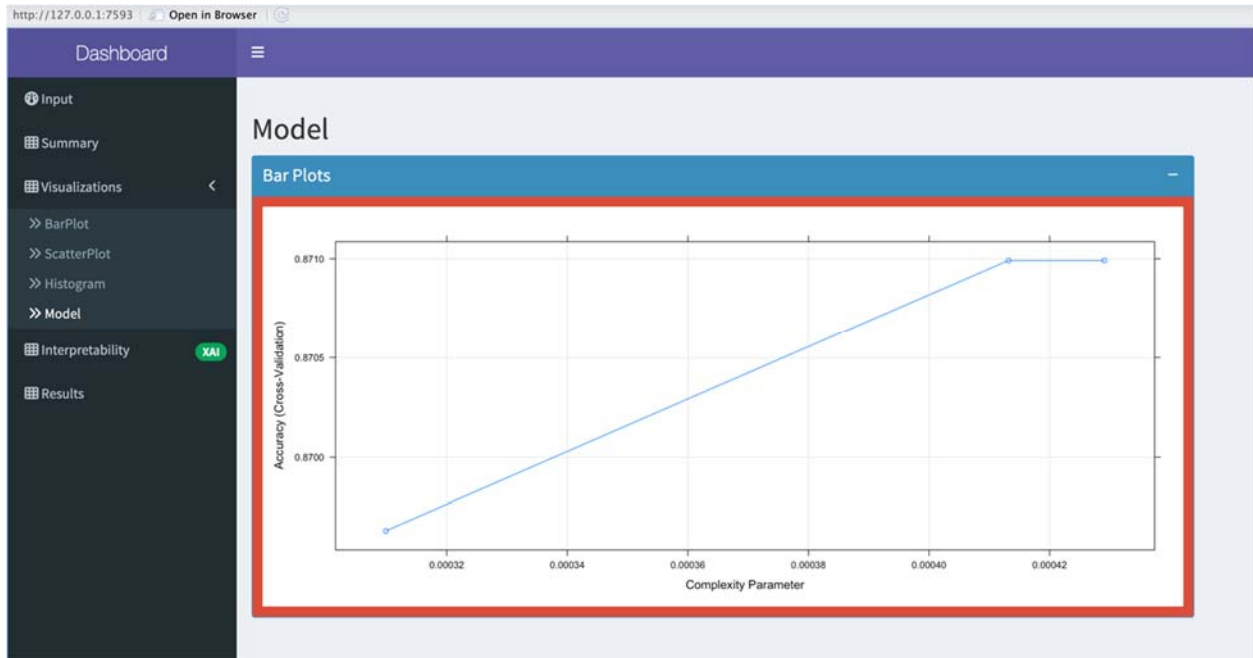


Figure 5.3: System Validation

## 5.6 Conclusions

The culmination of the implementation and testing phases positions the Shiny application as a paragon of reliability, sophistication, and user-centricity. A confluence of cutting-edge UI design, seamlessly integrated machine learning models, and a rigorous testing regimen fortifies the application's standing as an invaluable tool in the realm of financial risk assessment. The iterative feedback loop from user testing and system validation propels the application towards perpetual enhancement, ensuring its perpetual relevance in a dynamic landscape. The Shiny application emerges not just as a technological artifact but as a testament to the relentless pursuit of excellence in both form and function.

## Chapter 6: Discussion of Results

### 6.1 Interpretation of Exploratory Analysis

The study used exploratory data analysis to examine the dataset's distribution, relationships, and interactions.

#### 6.1.1 Univariate Analysis

The Table 6.1 shows a dataset revealed a class imbalance, with 275,818 customers (83.5%) classified as "No Default" and 54,635 (16.5%) as "Default." This imbalance suggests potential biases in model training, necessitating techniques like under sampling to mitigate skewed predictions.



Figure 6.1: Count of Target Variables

## 6.2 Bivariate Analysis

Bivariate exploratory analysis is a statistical method used to examine the relationship between two variables, offering valuable insights into connections and patterns within the data. The following figures depict relationships between specific variables, providing visual representations of their associations.

The negative relationship between instalment amounts and default likelihood (Figure 6.2) aligns with financial intuition: lower instalments reduce repayment burden, decreasing default risk. However, outliers indicate exceptions where contextual factors (e.g., income volatility) may override this trend.

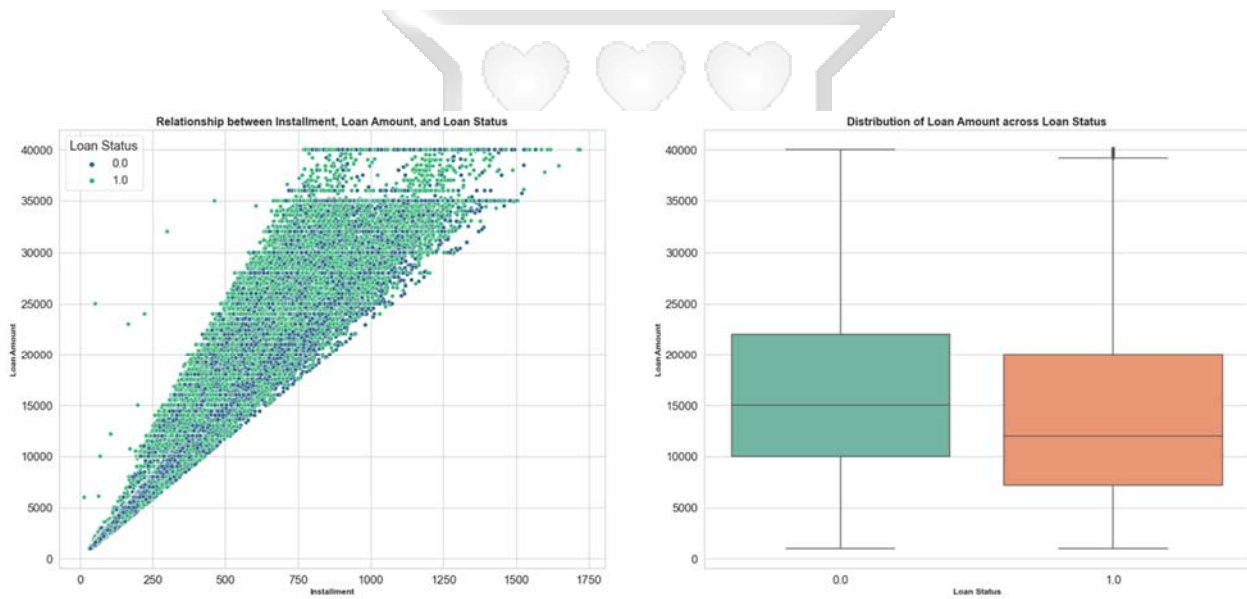
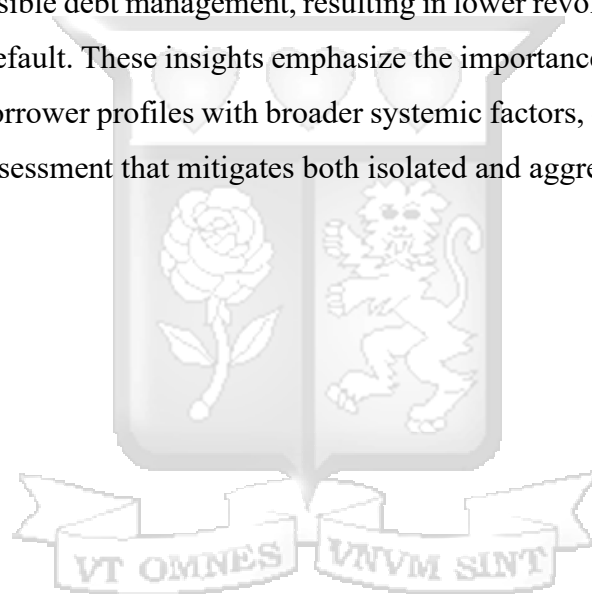


Figure 6.2: Bivariate Analysis

### 6.3 Multivariate Analysis

The heatmap presented in Figure 6.3 illuminates two pivotal correlations that enhance the understanding of credit risk dynamics within the proposed system. First, a strong positive correlation ( $r = 0.71$ ) exists between total accounts and open accounts, revealing that customers with a higher number of accounts tend to exhibit increased credit activity. While these customers may appear low-risk when evaluated individually, their cumulative exposure across multiple accounts could introduce systemic risks to the financial institution, necessitating careful monitoring. Second, a negative correlation between FICO scores and revolving utilization underscores that borrowers with higher credit scores, typically ranging from 700 to 850, demonstrate more responsible debt management, resulting in lower revolving credit utilization and a reduced likelihood of default. These insights emphasize the importance of designing risk models that balance individual borrower profiles with broader systemic factors, ensuring a comprehensive approach to credit risk assessment that mitigates both isolated and aggregated vulnerabilities.



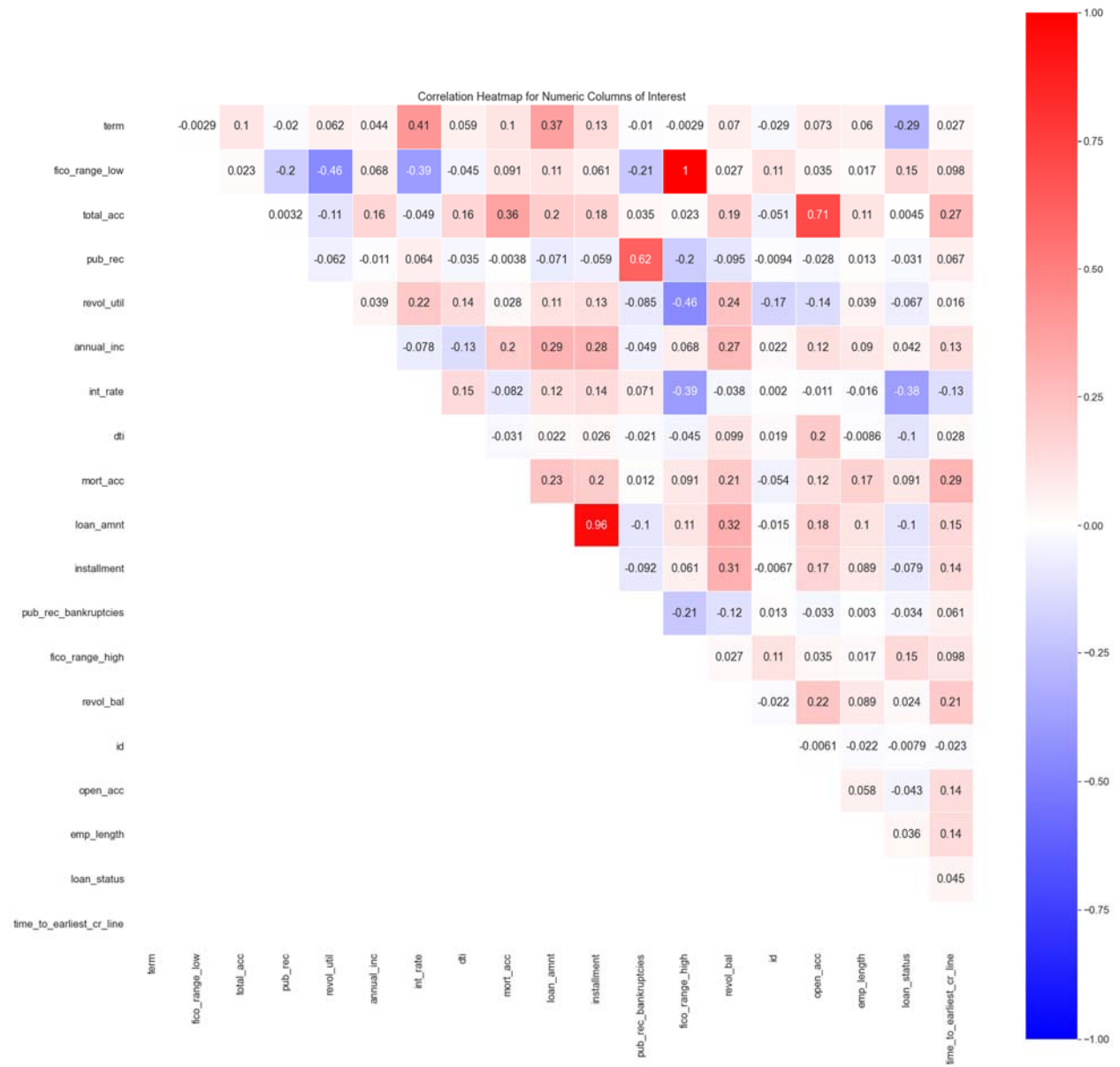


Figure 6.3: Multivariate Exploratory Analysis

## 6.4 Model Performance and Practical Implications

### 6.4.1 Comparative Performance

The comparative performance analysis of the credit risk assessment system, as detailed in Table 6.1, highlights the superior capabilities of ensemble methods, particularly XGBoost and LightGBM, which consistently outperformed other models. These methods achieved a balanced accuracy of approximately 0.85 and AUC scores around 0.80, demonstrating robust predictive power. Notably, XGBoost stood out with the highest Gini coefficient of 0.62, underscoring its exceptional ability to discriminate between high- and low-risk borrowers, making it a standout choice for precise risk stratification. In contrast, Gradient Boosting and Random Forest delivered comparable performance but were less effective in ranking risk, as evidenced by their lower Gini coefficients of 0.22 and 0.27, respectively, indicating reduced robustness in distinguishing nuanced risk levels.

Across all models, the system exhibited consistently high recall, ranging from 0.97 to 1.00, reflecting an impressive ability to identify true defaults and minimize missed high-risk cases. However, this strength came at the expense of precision, which ranged from 0.83 to 0.87, suggesting a tendency to over-predict risk, potentially flagging some low-risk borrowers as high-risk. These findings affirm the system's effectiveness in capturing critical defaults while highlighting the need for further refinement to optimize precision, ensuring a more balanced approach to risk prediction that aligns with both accuracy and practical applicability.

Model	accuracy	auc	ks	precision	recall	f1
Logistic Regression (undersample)	0.83	0.59	0.15	0.83	1.00	0.91
Decision Tree (undersample)	0.84	0.78	0.49	0.87	0.96	0.91
K Nearest Neighbours (undersample)	0.83	0.53	0.27	0.83	1.00	0.91
Random Forest (undersample)	0.84	0.80	0.48	0.84	0.99	0.91
Gaussian Naive Bayes (undersample)	0.83	0.53	0.04	0.83	1.00	0.91
Light GBM (undersample)	0.85	0.80	0.46	0.86	0.98	0.91
XGBoost (undersample)	0.85	0.81	0.48	0.87	0.97	0.91
Gradient Boosting (undersample)	0.85	0.80	0.46	0.86	0.97	0.91
Neural Network (undersample)	0.83	0.50	0.00	0.83	1.00	0.91

Table 6.1: Comparison of Precision Metrics : Evaluating Model Performance

### **6.4.2 Feature Importance**

The credit risk assessment system's analysis, as illustrated in Figures 6.4 and 6.7, reveals that both XGBoost and Random Forest models consistently prioritized revolving utilization and FICO scores as key predictors, reinforcing their well-established empirical connection to credit risk. These models underscored the intuitive importance of how borrowers manage their credit balances and their creditworthiness, as reflected by FICO scores, in determining default likelihood. In contrast, Discriminant Analysis, depicted in Figure 6.6, placed greater emphasis on less intuitive features, such as inquiry counts, which may indicate recent credit-seeking behaviour but are less commonly associated with risk in traditional assessments. This divergence highlights model-specific biases, where different algorithms weigh features uniquely, underscoring the need for a nuanced understanding of each model's strengths and limitations to ensure a balanced and comprehensive approach to credit risk evaluation.

### **6.5 Discussions**

The credit risk assessment system offers significant advantages in enhancing risk management and borrower engagement, while also facing certain limitations that warrant careful consideration.

### **6.6 Advantages of the Proposed System**

The system excels in risk mitigation by leveraging feature importance analysis, which identifies critical risk drivers such as revolving utilization and FICO scores. This insight enables financial institutions to implement proactive risk pricing strategies and diversify their portfolios, reducing exposure to potential defaults. Additionally, the system empowers customers by employing transparent models that demystify risk-based pricing. By clearly illustrating how factors like credit utilization influence loan decisions, the system fosters greater borrower understanding and encourages responsible credit use, ultimately promoting financial literacy and trust.

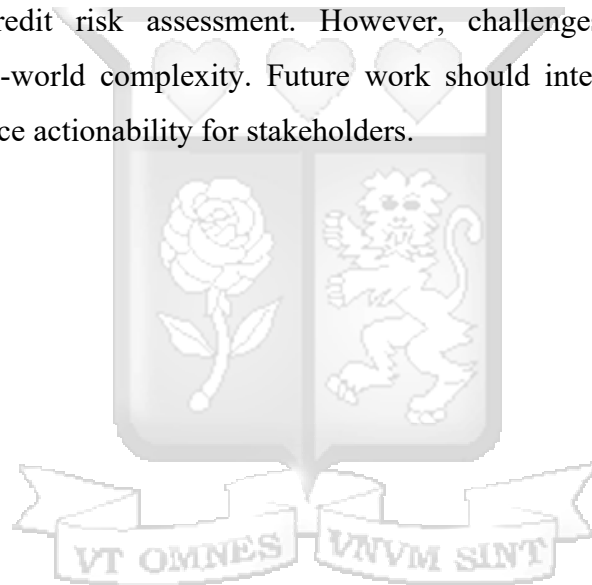
### **6.7 Disadvantages of the Proposed System**

However, the system is not without its challenges. A key limitation is the trade-off between interpretability and performance. Complex models, such as Neural Networks, deliver accuracy comparable to simpler models like decision trees but sacrifice interpretability, making it harder for users to trace decision-making processes. Another concern is bias amplification, where

explanations may oversimplify intricate feature interactions—for instance, failing to account for socioeconomic confounders that influence FICO scores, potentially leading to misleading or inequitable conclusions. Finally, the system’s reliance on computationally intensive techniques, such as SHAP analysis for generating detailed explanations, increases resource demands, posing challenges for scalability in resource-constrained environments. These advantages and limitations highlight the system’s transformative potential alongside areas for refinement to balance performance, fairness, and efficiency in credit risk assessment.

## 6.8 Conclusions

This study demonstrates that interpretable ML models (e.g., XGBoost) can balance performance and transparency in credit risk assessment. However, challenges persist in reconciling interpretability with real-world complexity. Future work should integrate causal inference to address biases and enhance actionability for stakeholders.



## Chapter 7: Conclusion and Recommendations

### 7.1 Synthesis of Key Findings

This study successfully fulfilled its dual objectives of advancing predictive power through sophisticated non-linear modelling while preserving interpretability to maintain performance and transparency, delivering a robust credit risk assessment system. The interconnected contributions of these objectives are synthesized below, highlighting their synergy and practical impact.

The system achieved superior predictive power by leveraging ensemble methods, such as XGBoost and LightGBM, which outperformed traditional models, achieving AUC scores as high as 0.81, as detailed in Table 6.1. The incorporation of non-linear features, notably interaction terms between FICO scores and revolving utilization, significantly enhanced the system's ability to distinguish high-risk applicants from their low-risk counterparts, ensuring precise identification of potential defaults.

Equally, critical was the system's commitment to interpretability, which empowered stakeholders to understand and trust model decisions. Techniques like SHAP values, illustrated in Figures 6.4 through 6.7, identified revolving utilization and credit inquiries as the top predictors of credit risk. These insights allowed users to audit the model's logic with clarity, ensuring alignment with regulatory expectations and fostering confidence in the system's fairness and accountability.

The interplay between predictive power and interpretability proved transformative, directly supporting risk-based pricing strategies. For instance, SHAP-driven insights enabled lenders to dynamically adjust interest rates based on transparent risk factors, such as penalizing borrowers with high revolving utilization while rewarding those with fewer credit inquiries. This interconnected approach not only enhances predictive accuracy but also ensures that risk assessments are actionable and equitable, empowering lenders to make informed decisions while meeting both business and regulatory demands.

## 7.2 Conclusion: Strengths and Limitations

The credit risk assessment system demonstrates notable strengths that enhance its analytical power and compliance, while also facing limitations that highlight areas for further refinement.

Among its strengths, the system excels in delivering actionable non-linear insights by capturing complex relationships, such as the U-shaped effect of debt-to-income ratios on default risk, where both very low and very high ratios correlate with increased risk. This nuanced understanding enables more precise identification of high-risk borrowers, empowering lenders to tailor their strategies effectively. Additionally, the system aligns seamlessly with regulatory requirements, particularly the EU's GDPR "right to explanation" mandate. Its interpretability tools, including SHAP values and decision trees, provide transparent explanations of model decisions while maintaining competitive AUC scores, ensuring compliance without sacrificing predictive performance.

However, the system is not without limitations. A key challenge is the trade-off between complexity and interpretability. Simpler models, such as Logistic Regression, achieved greater interpretability but lagged significantly in accuracy, with an AUC score of only 0.59, compared to more complex ensemble methods. This underscores the difficulty of balancing user-friendly explanations with high predictive power. Furthermore, data constraints posed a challenge, as the dataset exhibited class imbalance, with only 16.5% of cases representing defaults. To address this, undersampling was employed, which, while necessary, may have overlooked niche risk patterns, potentially limiting the system's ability to detect rare but significant risk profiles. These strengths and limitations highlight the system's robust capabilities, alongside opportunities to enhance accuracy and inclusivity in future iterations.

## 7.3 Recommendations

The credit risk assessment system, while achieving significant strides in balancing predictive accuracy and interpretability, can be further enhanced through targeted technical, operational, ethical, and research-oriented recommendations. These strategies aim to advance the system's capabilities, ensure fairness, and promote scalability, ultimately fostering a more equitable and robust approach to credit scoring.

From a technical perspective, exploring deep learning advancements, such as transformer-based architectures like TabTransformer, offers a promising avenue for improvement. By modelling non-linear relationships without relying on manual feature engineering, these architectures can automatically uncover latent interactions, such as the interplay between regional economic trends and credit behaviours, streamlining the detection of complex risk patterns. Additionally, integrating dynamic fairness audits into model training is essential to address potential biases. By leveraging tools like IBM's AIF360 toolkit, the system can monitor fairness metrics, such as equalized odds difference, in real time, ensuring that predictions remain equitable across demographic groups and aligning with ethical standards.

Operationally and ethically, adopting blockchain technology for transparency can significantly enhance trust and compliance. Storing model decisions on a permissioned blockchain, such as Hyperledger, creates immutable audit trails, enabling regulators to verify the fairness and logic of risk assessments with ease. Furthermore, fostering cross-industry collaboration with fintech companies can expand the system's inclusivity. By validating models on alternative data sources, such as rental payments or utility history, partnerships can reduce reliance on traditional metrics like FICO scores, broadening credit access for underserved populations while minimizing bias.

Longitudinal research is critical to ensuring the system's resilience over time. Implementing A/B testing frameworks to track model performance across economic cycles—such as recessions versus growth periods—will provide insights into the system's stability and adaptability, enabling continuous refinement to maintain accuracy under varying conditions.

In conclusion, this study successfully bridges the gap between accuracy and interpretability in credit scoring, delivering a system that empowers lenders with precise and transparent risk assessments. However, scaling these solutions to meet future demands requires embracing innovative strategies. By adopting deep learning to enhance predictive power, leveraging blockchain for unparalleled transparency, and pursuing collaborative and longitudinal efforts, financial institutions can future-proof their risk management systems. These advancements will not only strengthen predictive capabilities but also promote equitable credit access, ensuring that the system serves both business and societal goals effectively.

## References

- Altair. (2022). Credit scoring series part three: Data preparation and exploratory data analysis. <https://altair.com/newsroom/articles/credit-scoring-series-part-three-data-preparation-and-exploratory-data-analysis>
- Ariza-Garzón, M. J., Arroyo, J., Caparrini, A., & Segovia-Vargas, M. J. (2020). Explainable artificial intelligence in credit risk assessment: Combining machine learning with SHAP and LIME. 2020 IEEE International Conference on Big Data (Big Data), 1498–1507. <https://doi.org/10.1109/BigData50022.2020.9378132>
- Baniecki, H., & Biecek, P. (2019). modelStudio: Interactive studio with explanations for ML predictive models. *Journal of Open Source Software*, 4(43), 1798. <https://doi.org/10.21105/joss.01798>
- Biecek, P. (2018). DALEX: Explainers for complex predictive models in R. *Journal of Machine Learning Research*, 19(84), 1–5. <https://jmlr.org/papers/v19/18-416.html>
- Biecek, P., Chlebus, M., Gajda, J., Gosiewska, A., Kozak, A., & Ogonowski, D. (2021). Enabling machine learning algorithms for credit scoring—explainable artificial intelligence (XAI) methods for clear understanding complex predictive models. arXiv. <https://arxiv.org/abs/2104.06735>
- Bücker, M., Szepannek, G., & Wilhelm, A. (2021). Comparing explainable machine learning models with traditional scorecards for credit scoring. *Applied Artificial Intelligence*, 35(15), 1379–1398. <https://doi.org/10.1080/08839514.2021.1984032>
- Bussmann, N., Giudici, P., Marinelli, D., & Papenbrock, J. (2020). Explainable AI in fintech risk management. *Frontiers in Artificial Intelligence*, 3, 26. <https://doi.org/10.3389/frai.2020.00026>
- Campbell, J. Y., Hilscher, J., & Szilagyi, J. (2008). In search of distress risk. *The Journal of Finance*, 63(6), 2899–2939. <https://doi.org/10.1111/j.1540-6261.2008.01416.x>

- Central Bank of Kenya. (2013). Prudential guidelines: For institutions licensed under the banking act. <https://www.centralbank.go.ke/wp-content/uploads/2016/08/risk-management-guidelines-january-20131.pdf>
- Chang, W., & Wickham, H. (2018). Shiny: Web application framework for R. *The R Journal*, 10(1), 389–402. <https://doi.org/10.32614/RJ-2018-027>
- Dastile, X., Çelik, T., & Potsane, M. M. (2020). Statistical and machine learning models in credit scoring: A systematic literature survey. *Applied Soft Computing*, 91, 106263. <https://doi.org/10.1016/j.asoc.2020.106263>
- Davis, R., Epps, A., & Schafer, J. (2022). Applying explainable AI to home equity credit risk models. *Journal of Financial Data Science*, 4(1), 45–62. <https://doi.org/10.3905/jfds.2022.1.123>
- Dobson, A. J. (2015). *An introduction to generalized linear models* (3rd ed.). Chapman; Hall/CRC.
- Fahner, G. (2018). Developing transparent credit risk scorecards more effectively: An explainable artificial intelligence approach. *Proceedings of Data Analytics 2018*, 37–44.
- figshare. (2023). Lending club. Data set. [https://figshare.com/articles/dataset/Lending\\_Club/22121477/4](https://figshare.com/articles/dataset/Lending_Club/22121477/4)
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189–1232. <https://doi.org/10.1214/aos/1013203451>
- Gramegna, A., & Giudici, P. (2021). SHAP and LIME for credit risk models: A comparison on SME data. *Risks*, 9(5), 91. <https://doi.org/10.3390/risks9050091>
- Hand, D. J., & Henley, W. E. (1997). Statistical classification methods in consumer credit scoring: A review. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 160(3), 523–541. <https://doi.org/10.1111/j.1467-985x.1997.00078.x>
- Hastie, T., Tibshirani, R., & Friedman, J. (2013). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer Science & Business Media.

- Jakóbczak, D. J. (Ed.). (2015). Analyzing risk through probabilistic modeling in operations research. IGI Global.
- Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T.-Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Advances in neural information processing systems 30* (pp. 3146–3154). Neural Information Processing Systems Foundation. <https://papers.nips.cc/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf>
- Kubrusly, J., Neves, A. L. C. D., & Marques, T. L. (2022). A statistical analysis of textual e-commerce reviews using tree-based methods. *Open Journal of Statistics*, 12(3), 357–372. <https://doi.org/10.4236/ojs.2022.123023>
- Leo, M., Sharma, S., & Maddulety, K. (2019). Machine learning in banking risk management: A literature review. *Risks*, 7(1), 29. <https://doi.org/10.3390/risks7010029>
- Lessmann, S., Baesens, B., Seow, H.-V., & Thomas, L. C. (2015). Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, 247(1), 124–136. <https://doi.org/10.1016/j.ejor.2015.05.030>
- Louzada, F., Ara, A., & Fernandes, G. B. (2016). Classification methods applied to credit scoring: Systematic review and overall comparison. *Surveys in Operations Research and Management Science*, 21(2), 117–134. <https://doi.org/10.1016/j.sorms.2016.10.001>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In I. Guyon, U. von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. V. N. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems 30* (pp. 4765–4774). Neural Information Processing Systems Foundation. [https://papers.nips.cc/paper\\_files/paper/2017/hash/8a20a8621978632d76c43dfd28b67767-Abstract.html](https://papers.nips.cc/paper_files/paper/2017/hash/8a20a8621978632d76c43dfd28b67767-Abstract.html)
- Martens, D., Baesens, B., Van Gestel, T., & Vanthienen, J. (2007). Comprehensible credit scoring models using rule extraction from support vector machines. *European Journal of Operational Research*, 183(3), 1466–1476. <https://doi.org/10.1016/j.ejor.2006.04.051>

- Misheva, B. R., Osterrieder, J., Hirska, A., & Khan, O. (2021). Explainable AI for credit scoring in peer-to-peer lending: An empirical study. *Journal of Risk and Financial Management*, 14(11), 538. <https://doi.org/10.3390/jrfm14110538>
- Molnar, C. (2020). *Interpretable machine learning*. Lulu.com.
- Moscato, V., Picariello, A., & Sperlí, G. (2021). An explainable AI framework for credit scoring in P2P lending. *IEEE Transactions on Computational Social Systems*, 8(4), 987–998. <https://doi.org/10.1109/TCSS.2021.3063632>
- Office of the Data Protection Commissioner. (2021). *The data protection act: General regulations 2021*. <https://www.odpc.go.ke/wp-content/uploads/2021/06/L.N-263-265-THE-DATA-PROTECTION-GENERAL-REGULATIONS-2021FIN.pdf>
- Resnizky, H. G. (2015). *Learning shiny*. Packt Publishing.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should i trust you?”: Explaining the predictions of any classifier. *arXiv*. <https://arxiv.org/abs/1602.04938>
- Schwaber, K., & Sutherland, J. (2020). *The scrum guide: The definitive guide to scrum: The rules of the game*. White paper. <https://scrumguides.org/docs/scrumguide/v2020/2020-Scrum-Guide-US.pdf>
- Siddiqi, N. (2012). *Credit risk scorecards: Developing and implementing intelligent credit scoring*. John Wiley & Sons.
- Slovik, P. (2012). Systemically important banks and capital regulation challenges [OECD Economics Department Working Paper No. 916]. OECD Publishing. <https://doi.org/10.1787/5kg0ps8cq8q6-en>
- Talaat, F. M., & El-Balky, M. (2023). Toward interpretable credit scoring: Integrating explainable artificial intelligence with deep learning for credit card default prediction. *Neural Computing and Applications*. <https://doi.org/10.1007/s00521-023-08547-5>

Team, S. (2023). Getting started: Working with Shiny for R. In Shinyapps.io user guide. <https://docs.posit.co/shinyapps.io/getting-started.html#working-with-shiny-for-r>

Thomas, L. C., Crook, J. N., & Edelman, D. B. (2017). Credit scoring and its applications (2nd ed.). SIAM.

Zoldi, S., & Fahner, G. (2022). Combining machine learning with credit risk scorecards. FICO Blog Series on Explainable AI in Credit Risk. <https://www.fico.com/blogs/combining-machine-learning-credit-risk-scorecards>



# Appendices

## Appendix A: Similarity Report

dsafinal

ORIGINALITY REPORT

**19%** SIMILARITY INDEX  
**16%** INTERNET SOURCES  
**9%** PUBLICATIONS  
**12%** STUDENT PAPERS

PRIMARY SOURCES

1	arxiv.org Internet Source	4%
2	Submitted to Wright College Student Paper	2%
3	de.overleaf.com Internet Source	1%
4	fastercapital.com Internet Source	1%
5	Submitted to National Research University Higher School of Economics Student Paper	1%
6	Zakowska, Anna. "A New Credit Scoring Model to Reduce Potential Predatory Lending: A Design Science Approach", The Claremont Graduate University, 2024 Publication	1%
7	su-plus.strathmore.edu Internet Source	1%
8	Submitted to University of Brighton Student Paper	<1%

9	www.pure.ed.ac.uk Internet Source	<1%
10	bibliotecadigital.fgv.br Internet Source	<1%
11	iksadyayinevi.com Internet Source	<1%
12	www.tandfonline.com Internet Source	<1%
13	Submitted to Sir Arthur Lewis Community College Student Paper	<1%
14	Submitted to University of Salford Student Paper	<1%
15	aigeolabs.com Internet Source	<1%
16	ir.uitm.edu.my Internet Source	<1%
17	ajeb.buh.edu.vn Internet Source	<1%
18	Submitted to Brunel University Student Paper	<1%
19	Submitted to Kaplan College Student Paper	<1%
20	Submitted to University of Huddersfield	

VT OMNES VIVVM SINTE

## Appendix B: Ethics Clearance Confirmation



24<sup>th</sup> May 2024

Mr Aswani Allan,  
allan.aswani@strathmore.edu

Dear Mr Aswani,

### **RE: Complex and Explainable Machine Learning Models in Credit Scoring**

This is to inform you that SU-ISERC has reviewed and **approved** your above **SU-masters** proposal. Your application reference number is **SU-ISERC2175/24**. The approval period is from **24<sup>th</sup> May 2024 to 23<sup>rd</sup> May 2025**.

This approval is subject to compliance with the following requirements:

- i. Only approved documents including (informed consents, study instruments, MTA) will be used.
- ii. All changes including (amendments, deviations, and violations) are submitted for review and approval by SU-ISERC.
- iii. Death and life-threatening problems and serious adverse events or unexpected adverse events whether related or unrelated to the study must be reported to SU-ISERC within 72 hours of notification.
- iv. Any changes anticipated or otherwise that may increase the risks or affected safety or welfare of study participants and others or affect the integrity of the research must be reported to SU-ISERC within 72 hours.
- v. Clearance for the export of biological specimens must be obtained from relevant institutions.
- vi. Submission of a request for renewal of approval at least 60 days prior to the expiry of the approval period. Attach a comprehensive progress report to support the renewal.
- vii. Submission of an executive summary report within 90 days of completion of the study to SU-ISERC.

Before commencing your study, you will be expected to obtain a research license from National Commission for Science, Technology, and Innovation (NACOSTI) <https://research-portal.nacosti.go.ke/> and obtain other clearances needed.

Yours sincerely,

A handwritten signature in blue ink, appearing to read "Ambrose Rachier".

**Mr Ambrose Rachier,**  
**Chairperson; SU-ISERC**