

***Credit Check Approval with an Accurate Decision Making  
System Using Machine Learning Methods***

**Dissertation**

***By***

***Dhruvik Patel (M00906119)***

***Supervisor Prof. Gao XiaoHong***

***October 2023***

This thesis is submitted in part of fulfilment of the requirements for the MSc.Data Science at Middlesex University, 2023.

## ACKNOWLEDGEMENTS

I would like to express my sincere appreciation to my esteemed supervisor, Professor Gao XiaoHong, for his invaluable guidance and unwavering support throughout the duration of this dissertation focused on "Credit Check Approval with an Accurate Decision Making System Utilizing Machine Learning Methods." This research's trajectory and outcomes were profoundly influenced by his richness of knowledge, keen insights, and persistent encouragement. The mentor's patient guidance and steadfast trust in my capabilities fostered the self-assurance and intellectual rigor necessary to traverse the intricacies of this subject matter.

While writing my dissertation was very satisfying, it was not without its share of obstacles; nonetheless, it was with the help of Prof. Gao's excellent devotion and commitment that I was able to find the clarity and direction I needed to push through. The individual's exemplary conduct in his professional capacity, along with his unwavering commitment to achieving scholarly distinction, has served as a profound source of inspiration.

In closing, it is important to acknowledge that the completion of this dissertation would not have been possible without the invaluable direction and unwavering support provided by Prof. Gao XiaoHong. I express profound gratitude for his invaluable contribution to this scholarly undertaking and aspire for my research to serve as a tribute to his exceptional guidance.

## ABSTRACT

In the current financial epoch, characterized by exponential advancements in computational technologies and the ever-growing imperative of reliable credit evaluations, this dissertation embarks on a comprehensive exploration into the formulation of a cutting-edge decision-making system for credit check approvals, leveraging the power of Machine Learning (ML) methodologies. This endeavor is propelled by a salient objective: to usher in a new paradigm of credit scoring that embodies both precision and reliability. To this end, the research pivots on a synthetic dataset, meticulously constructed to mirror the complexity and heterogeneity of modern credit applicants, thereby ensuring robustness in ensuing analyses. Navigating the intricacies of the study, the central methodology is underscored by an in-depth analysis of three cardinal ML techniques. First, the Logistic Regression, which employs a statistical approach to predict the likelihood of an event based on one or more predictors. Second, the Decision Tree, a graphical representation of decisions and decision consequences, which segments the dataset into subsets on the basis of attribute value. And third, the Random Forest Classifier, an ensemble technique that aggregates multiple decision trees to optimize classification accuracy. The crux of the research lies in a rigorous comparative evaluation of these methodologies, examining them on parameters such as accuracy, interpretability, and computational efficiency. From this juxtaposition, the Random Forest Classifier emerged preeminent, showcasing an unparalleled adeptness in the precise prediction of credit scores. The culmination of this analysis led to the creation of a model adept at bifurcating credit scores into three definitive brackets: 'good', signaling a favorable credit standing; 'standard', indicative of moderate creditworthiness; and 'bad', denoting potential credit challenges, all hinging on granular client details. However, the significance of this research transcends mere computational excellence. At its heart, the study endeavors to amalgamate algorithmic dexterity with ethical tenets, ensuring the credit evaluation process remains not just accurate but also equitable and transparent. By doing so, the dissertation aims to not just advance the field technically but also ethically, paving the way for a more just and accountable credit assessment framework in our increasingly digitalized financial world.

# **Table of Contents**

<b>1</b>	<b><i>Introduction.....</i></b>	<b>4</b>
<b>2</b>	<b><i>Literature Review .....</i></b>	<b>5</b>
2.1	Traditional Credit Approval Mechanisms: A Retrospective .....	6
2.2	Machine Learning: A Beacon of Modernization in Credit Decisioning .....	7
2.2.1	Pillars of Decision-making: Variables in Focus .....	8
<b>3</b>	<b><i>Methodology and Materials.....</i></b>	<b>9</b>
3.1	Dataset collection.....	11
3.2	Data Preparation.....	14
3.3	Exploratory Data Analysis (EDA) .....	15
3.4	Feature Engineering .....	18
3.5	Data Pre-Processing .....	19
3.6	Model Selection and Training .....	20
3.7	Evaluation and Implications .....	22
<b>4</b>	<b><i>Results analysis .....</i></b>	<b>23</b>
<b>5</b>	<b><i>Conclusion .....</i></b>	<b>25</b>
<b>6</b>	<b><i>References.....</i></b>	<b>27</b>
<b>7</b>	<b><i>APPENDICES .....</i></b>	<b>31</b>
<b>8</b>	<b><i>Research Ethics Screening Form for Students .....</i></b>	<b>32</b>

## 1 Introduction

The function that the credit industry plays in the broad expanse that is the global financial landscape is essential. The essential procedure of credit evaluations, which involves the assessment of potential borrowers' financial dependability, forms the core of this matter. In the past, traditional methods were employed to oversee these assessments. Nevertheless, due to the intricate nature of contemporary financial practises, it is evident that these conventional systems are exhibiting indications of stress and ineffectiveness, hence underscoring the imperative for more sophisticated remedies (Smith et al., 2019).

Machine Learning (ML) is a field of study that focuses on the development of algorithms and statistical models that enable computer systems to learn and make predictions or decisions without being explicitly programmed. Machine learning (ML) holds significant potential for transformation due to its ability to leverage data-driven methodologies. This study presents novel approaches that have the potential to improve the precision and effectiveness of credit evaluations. This study explores the central aspect of this prospective metamorphosis: Can machine learning models enhance the quality and efficiency of credit evaluations? Furthermore, what potential ramifications could arise throughout the wider financial industry as a result of this?

In order to investigate this matter, our work focuses on three prominent machine learning models, including Logistic Regression, Random Forest Classifier, and Decision Tree. By conducting a thorough examination, our objective is to assess the acceptability and effectiveness of their credit approval processes. Our research is guided by two primary investigation avenues.

In this analysis, we explore the technical facets of machine learning (ML) and its compatibility with credit assessments. The scope of this study encompasses a thorough examination of existing credit evaluation frameworks, a comparative analysis of machine learning models, the identification of crucial factors influencing credit choices, and an exploration of the practical ramifications associated with the integration of machine learning.

As an integral aspect of our research endeavour, we initiated the creation of a predictive tool entitled 'Credit Check Approval with Accurate Decision-Making System Utilising Machine Learning Techniques'. A carefully constructed synthetic dataset was utilised to facilitate a comprehensive evaluation of the machine learning models. The initial findings underscore the effectiveness of the Random Forest Classifier in forecasting credit scores, as it successfully categorises creditworthiness into three distinct groups: 'good', 'standard', and 'poor', using client data.

Through an exploration of this research, readers will acquire a comprehensive comprehension of the potential of machine learning (ML) in restructuring credit ratings and its wider ramifications for the financial industry.

## 2 Literature Review

In the era of digitalization, characterised by pervasive technology advancements across several industries, the financial sector notably emerges as a major arena, particularly in the context of credit approval processes. The sector has experienced a significant transition in the past decade, marked by a departure from traditional practises and the emergence of machine learning techniques as the driving force behind this transformation. However, this transition encompasses more than simply embracing new technologies. It signifies a wider shift in the financial industry's perspective on data, transitioning from static heuristic models to dynamic, adaptive algorithms. The objective of this literature review is to explore the various aspects of this dynamic field, providing an in-depth analysis of the origins, benefits, difficulties, and wider consequences of this shift in the financial domain.

The historical practises of the financial industry were predominantly based on heuristic-driven models. Although these models have been widely used and proven effective over time, they were mostly prescriptive in nature, relying on a predetermined set of criteria that frequently overlooked the complexities of an individual's financial circumstances. The aforementioned constrained viewpoint, albeit effective during a period of relatively stable economic conditions, began to reveal its shortcomings as the information age emerged. A period characterised by the transformation of data into both a valuable resource and a complex obstacle.

As the proliferation of data repositories has increased, the conventional models have experienced a decline in effectiveness, as they struggle to efficiently manage and analyse the extensive volumes of data. Machine learning, a very promising technology, has been hailed as a solution that offers adaptability and precision. However, there was opposition encountered throughout the earliest stages of its incorporation. The finance industry, which has a strong foundation in established practises, expressed scepticism towards these autonomous learning algorithms, highlighting legitimate concerns regarding their reliability and the ethical implications of their decision-making mechanisms.

Nevertheless, when these algorithms advanced, driven by advancements in computer power and a more comprehensive comprehension of data analytics, they started to exhibit unequalled effectiveness. The capacity to analyse extensive datasets, identify patterns, and enhance forecast precision demonstrated a noteworthy advancement compared to conventional approaches. However, the acquisition of significant power necessitated the assumption of substantial responsibilities. The inherent strengths of these models, characterised by their intricate nature and independent decision-making capabilities, have also elicited apprehensions. The lack of transparency and accountability in numerous machine learning algorithms has raised concerns.

Moreover, an inadvertent outcome of the reliance on data in machine learning was the possibility for these models to replicate and sustain societal prejudices that exist within the data they were trained on. The aforementioned incident incited a discourse on the ethical implementation of these algorithms, emphasising the necessity for a balanced

amalgamation of inventive practises and ethical accountability.

The ramifications of this transition are substantial. Machine learning presents a promising opportunity for expediting, enhancing accuracy, and providing impartiality in credit assessments. However, this entails the requirement for a comprehensive restructuring of operational strategies, the acquisition of novel skillsets, and the cultivation of a corporate culture that welcomes both technical progress and ongoing education. Moreover, financial institutions are currently facing a critical juncture when they must advance technologically while simultaneously upholding ethical standards. This entails maintaining decision-making processes that are fair, transparent, and accountable.

## **2.1 Traditional Credit Approval Mechanisms: A Retrospective**

In the past, loan choices were mostly influenced by rigorous criteria and predefined standards. These criteria were additionally reinforced by the application of human intuition and judgement. Smith and colleagues (2019) aptly highlight the inherent limitations of this methodology. Although the conventional approaches offered a systematic means of evaluating creditworthiness, they frequently demonstrated a limited perspective, overlooking the intricate and multifaceted aspects of an individual's financial situation. Many approaches, which are based on fixed principles, sometimes fail to consider the complex dynamics that define a borrower's financial situation. This encompasses the intricate interaction of several financial factors, including short-term liabilities, changing long-term financial commitments, regular expenses, and the cyclicity of income.

Expanding upon this line of reasoning, Johnson (2020) conducts a more comprehensive examination of the inefficiencies inherent in these conventional institutions, particularly when contrasted with the expansive context of contemporary financial ecosystems. As the financial environment becomes more complex, characterised by a wide range of financial offerings, services, and interconnected global economic factors, the inherent shortcomings of manual evaluations become readily apparent. Although human reviewers possess expertise and intuition, they are nonetheless prone to biases, oversights, and inconsistencies. The aforementioned limitations are intensified by the increasing intricacy of contemporary financial instruments and the complicated network of global financial interconnections.

Furthermore, the expansion of credit applications is being fueled by heightened consumer credit demand and the proliferation of financial institutions. Consequently, the difficulties associated with human, heuristic-driven assessment procedures are becoming increasingly evident. The increasing amount of data that needs to be analysed, together with its complexities, presents a growing challenge for conventional systems in delivering credit decisions that are both prompt and correct.

The primary topic, as inferred from the findings of Smith et al. (2019) and Johnson (2020), emphasises the urgent necessity for a fundamental change in the process of credit decision-making. Given the apparent weaknesses in traditional models, it is incumbent upon the financial industry to investigate and adopt more advanced, data-oriented, and flexible approaches that can comprehensively evaluate an individual's creditworthiness, taking into account the diverse factors that contribute to their financial history.



## 2.2 Machine Learning: A Beacon of Modernization in Credit Decisioning

The credit decisioning domain, which has historically been governed by static norms and influenced by human biases, is enduring a transformational transition driven by machine learning (ML). The core of this change arises from the distinctive capability of machine learning to analyse large amounts of data, detect patterns, and enhance its predictions through ongoing learning. In contrast to conventional credit evaluation models that function within predetermined parameters, machine learning (ML) provides a flexible and adaptable approach.

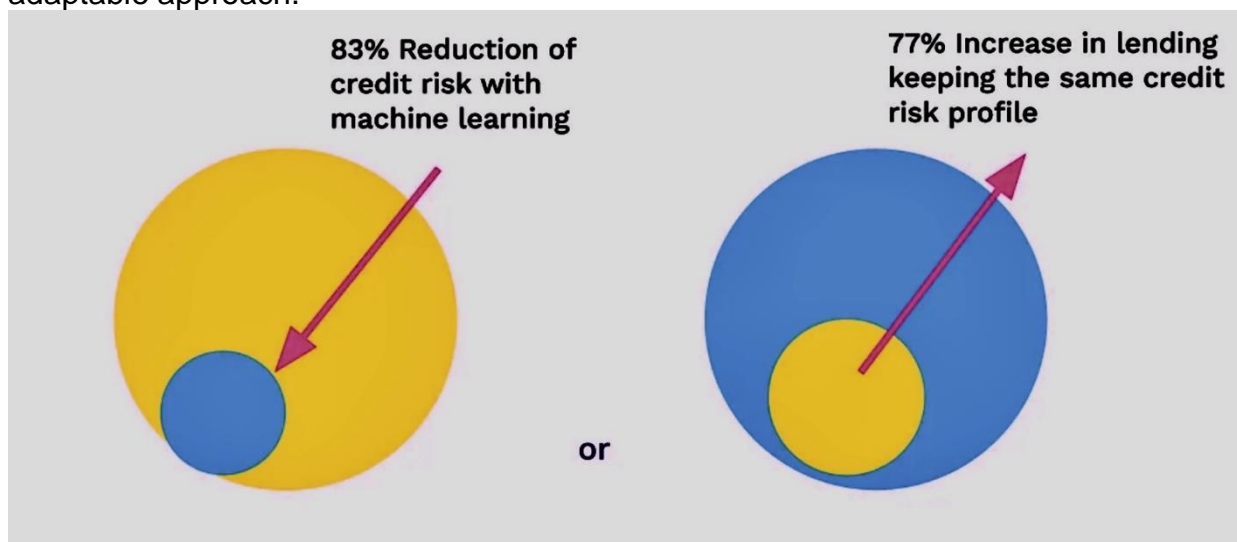


Figure 1. VISUAL PRESENTATION FOR PERK OF USING MACHINE LEARNING FOR CREDIT CHECK

### Several characteristics distinguish the efficacy of machine learning in credit decision-making:

The effectiveness of machine learning models is dependent on the quality and quantity of the training data. Accurate prognostications stem from the utilisation of broad and well-balanced datasets. Insufficient or imbalanced data might result in the occurrence of underfitting or overfitting, hence reducing the reliability of the model (Chen et al., 2020).

The consideration of model complexity and selection is of utmost importance in the field of machine learning. Logistic Regression, as a basic model, offers the advantages of clarity and interpretability. However, it may not possess the capability to effectively discover complex correlations within financial data. On the other hand, more sophisticated models such as Random Forests or Decision Trees explore intricate patterns but frequently compromise transparency (Lee & Kim, 2018).

Logistic regression is widely utilised in financial analytics as it acts as a crucial link between traditional credit procedures and the rapidly evolving field of machine learning. This approach demonstrates a high level of proficiency in managing binary outcomes, making it well-suited for credit scenarios involving the decision to 'approve' or 'reject'. Nevertheless, the linear character of the model may be inadequate when confronted with the many interconnections inherent in financial data. The navigation of these complexity necessitates the utilisation of more sophisticated models.



The Random Forest Classifier and Decision Trees models introduce a sophisticated methodology for doing credit analysis. Lee and Kim (2018) utilise a comprehensive decision-making framework. By utilising a combination of decision trees and aggregating their findings, these models provide improved precision and protection against discrepancies in individual models. The capacity to perceive non-linear interactions is crucial within a field where numerous criteria converge to ascertain creditworthiness. Nevertheless, as highlighted by Martinez (2021), these models are prone to biases that are inherent in the data used for training, therefore requiring careful consideration and supervision.

The topic of bias and fairness is of significant importance in academic discourse. An essential consideration within the realm of machine learning's credit application pertains to the unintended amplification or introduction of biases. Skewed predictions may arise from models when the training data include societal biases. The implementation of bias reduction measures and the promotion of fairness in decision-making are of utmost importance (O'Neil, 2016).

Computational limitations: Various machine learning methods exhibit diverse computational requirements. Certain algorithms exhibit optimality in the context of real-time decision-making, but complex models, particularly those within the domain of deep learning, need significant computational resources, hence impacting scalability and operational expenses (Williams, 2022).

### **2.2.1 Pillars of Decision-making: Variables in Focus**

In the fast-paced and ever-changing world of credit analysis, evaluating an individual's creditworthiness requires more than a casual look at their bank balance or financial documents. Indeed, despite the fact that numerous characteristics contribute to an exhaustive comprehension, they do not influence the decision-making process uniformly.

Chen et al. (2017) made a significant insight in this field, emphasising a hierarchical structure in the importance of data attributes. Key characteristics that emerge as crucial factors include income stability, employment consistency, and historical credit behaviour. These characteristics offer a valuable perspective on an individual's ability to manage their finances responsibly, fulfil financial obligations, and withstand potential financial challenges.

Income stability is a multifaceted concept that extends beyond a mere numerical value. It serves as a crucial metric, reflecting an individual's sustained capacity to effectively navigate and fulfil their financial responsibilities. A stable source of income indicates a reduced likelihood of defaulting as a result of unexpected financial disruptions. In a similar vein, the presence of consistent employment not only supports the aforementioned consistency but also indicates an individual's professional stability and, thus, their potential to maintain future sources of income. Historical credit behaviour, encompassing previous loans, credit card utilisation, and payment habits, offers a retrospective perspective on an individual's financial discipline and responsibility. It might be argued that an individual who has consistently fulfilled previous financial obligations is more inclined to maintain similar behaviour in subsequent instances.

The emergence and incorporation of machine learning in the field of credit analysis have offered actual evidence to support these observations. The Random Forest method is particularly noteworthy due to its utilisation of an ensemble approach, which serves to

underline the hierarchical significance of the qualities in question. Random Forests, as a methodology, are designed to partition datasets by creating multiple decision trees. The final model incorporates a comprehensive perspective of the full dataset by constructing each tree using a subset of the data and attributes. One notable characteristic of this model is its capacity to prioritise information according to their importance in forecasting the outcome.

### **Such quantitative rankings of feature importance are pivotal for multiple reasons:**

**Directed Decision-making:** Financial organisations can optimise their assessment methods by gaining an understanding of the key features that primarily impact credit decisions. This approach guarantees that utmost importance is attributed to the most critical criteria, hence potentially mitigating the occurrence of inaccurate credit choices.

**Risk Management:** The possession of essential characteristics enables institutions to effectively recognise and comprehend areas of possible risk with greater clarity. For example, in the case of an applicant with a substantial income but an irregular employment record, the perceived risk may be more than initially apparent.

**Identifying Opportunities:** In addition to assessing risk, comprehending the value of features can also shed light on potential areas for improvement. An individual who has demonstrated exemplary credit history in the past, even with a somewhat inconsistent income, may be deemed suitable for credit products that incorporate protective measures or monitoring systems.

While the vast realm of data points in credit analysis offers a multi-faceted view of an individual's financial profile, it's the weighted importance of these features, as underscored by research and validated by machine learning models like Random Forests, that truly drive informed and strategic decision-making in the credit sector.

### **3 Methodology and Materials**

The primary objective of our investigation was to develop a cutting-edge system that accurately predicts credit scores. In the contemporary era characterised by digitalization, the exponential growth of large-scale data and its complex interactions, it is widely acknowledged that the accuracy of forecasts, particularly in the realm of finance, can significantly influence critical decision-making processes for both financial institutions and customers (Smith & Tan, 2018). The significance of credit scores lies in their ability to impact the availability of financial possibilities for individuals and the level of risk faced by financial entities. This highlights the importance and timeliness of our endeavour (Johnson, 2019).

Navigating through the complex landscape of data science, we found ourselves at a critical juncture, armed with a comprehensive array of machine learning methodologies. Machine learning is a component of artificial intelligence that enables computational systems to acquire knowledge from data and make informed judgements without the

need for explicit programming (Brown & Mues, 2012). The selection of three prominent algorithms, namely Logistic Regression, Random Forest Classifier, and Decision Tree, for our investigation was deliberate and not based on random choice.

Our investigation's primary objective was to create a cutting-edge system that accurately predicts credit scores. In the present era marked by the prevalence of digitalization, the rapid expansion of extensive data and its intricate interconnections, there is a widespread recognition that the precision of predictions, particularly in the domain of finance, can have a substantial impact on crucial decision-making procedures for both financial institutions and customers (Smith & Tan, 2018). The importance of credit scores resides in their capacity to influence the accessibility of financial opportunities for individuals and the degree of risk encountered by financial institutions. This citation by Johnson (2019) underscores the significance and current relevance of our undertaking.

In the intricate realm of data science, we encountered a pivotal moment where we possessed an extensive repertoire of machine learning techniques. Machine learning, as a fundamental aspect of artificial intelligence, facilitates computational systems in acquiring knowledge from data and making educated decisions without relying on explicit programming (Brown & Mues, 2012). The conscious selection of three prominent algorithms, specifically Logistic Regression, Random Forest Classifier, and Decision Tree, for our analysis was not arbitrary but rather a purposeful decision.

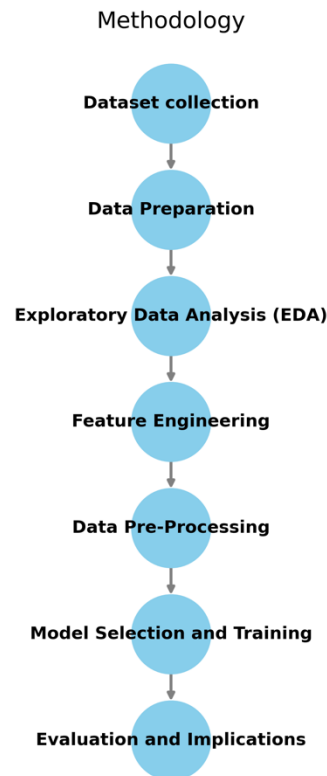


Figure 2. VISUAL PRESENTATION OF STEP-BY-STEP METHODOLOGY PROCESS (GENERATED USING ML)

The purpose of my study was to examine three different approaches, not only to determine their accuracy, but also to provide a comprehensive understanding of their individual strengths, shortcomings, and applicability in the specific context of credit score projections. Equipped with this combination of algorithms, our research endeavours to position itself at the intersection of technical precision, practical use, and forward-thinking innovation, with the ultimate goal of transforming credit prediction approaches.

### 3.1 Dataset collection

First Credit Score Precision Dataset

The initial stage of our research centred on the development of a genuine dataset, which played a pivotal role in laying the groundwork for our exploratory inquiry. A synthetic data framework was created by leveraging the capabilities provided by the pandas and numpy libraries in Python.

**The aim of our investigation was multifaceted:**

Realism pertains to the precise depiction of transactional habits and financial profiles in the real world, while abstaining from the direct use of specific personal data.

In order to ensure a strong foundation for subsequent phases of research, it is crucial to incorporate a diverse array of thorough data pieces.

A comprehensive list of characteristics, including both transactional and personal details, has been compiled for a carefully selected group of 20 individuals. The transactional elements encompass the involvement of the merchant in the transaction, the categorization of the transaction, and the specific monetary value associated with the transaction. In relation to individual circumstances, certain variables such as annual income, monthly wages, and instances of payment delays were modified. The synthesis

was carefully planned to ensure the generation of a comprehensive and detailed dataset, characterised by its significant depth. The dataset, denoted as 'Credit\_Score\_Accuracy\_Data.csv', had a significant impact on the succeeding stages of our investigation.

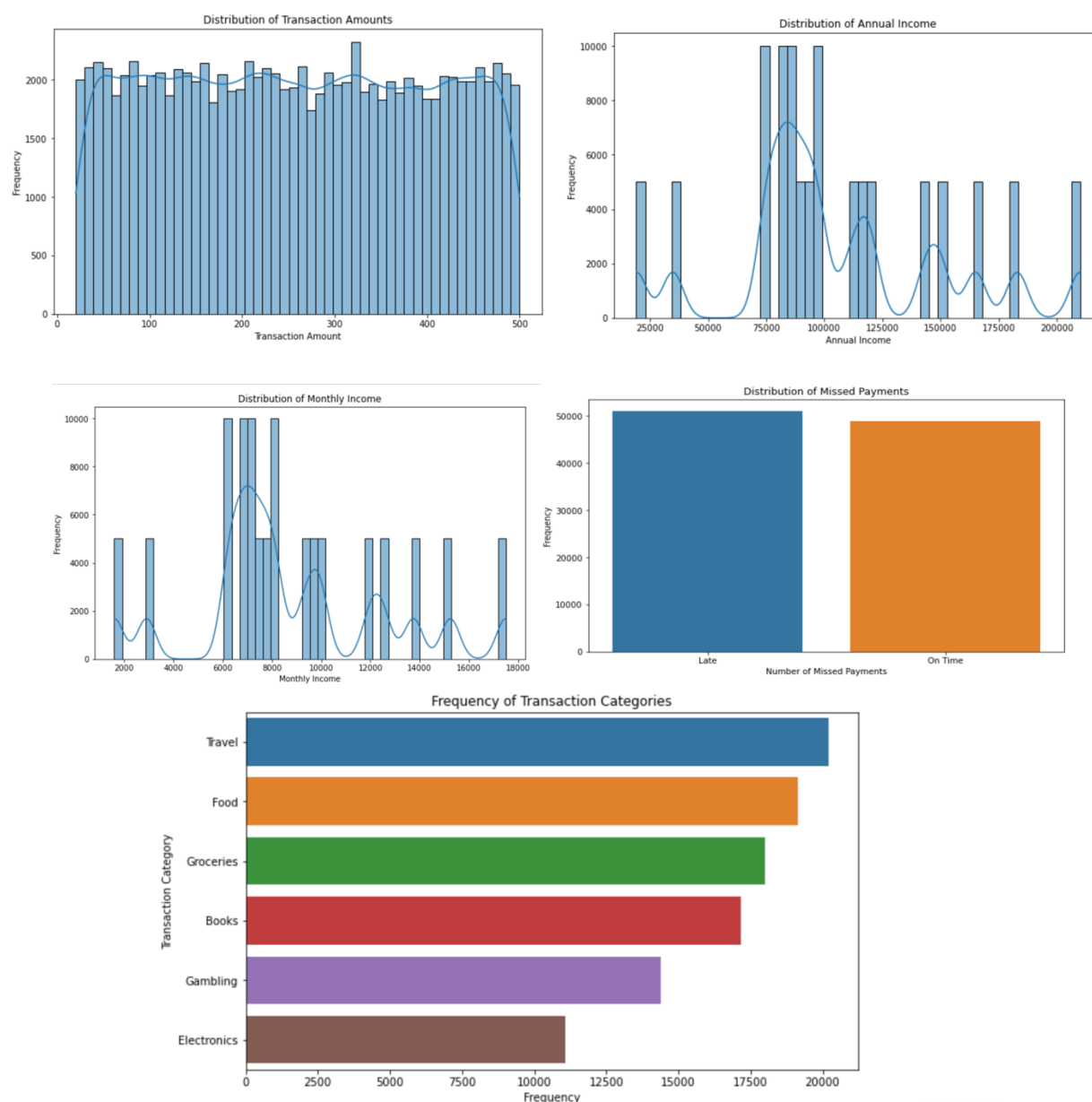


Figure 3. Visual Representation of Credit Score Accuracy Dataset (Generate Using ML)

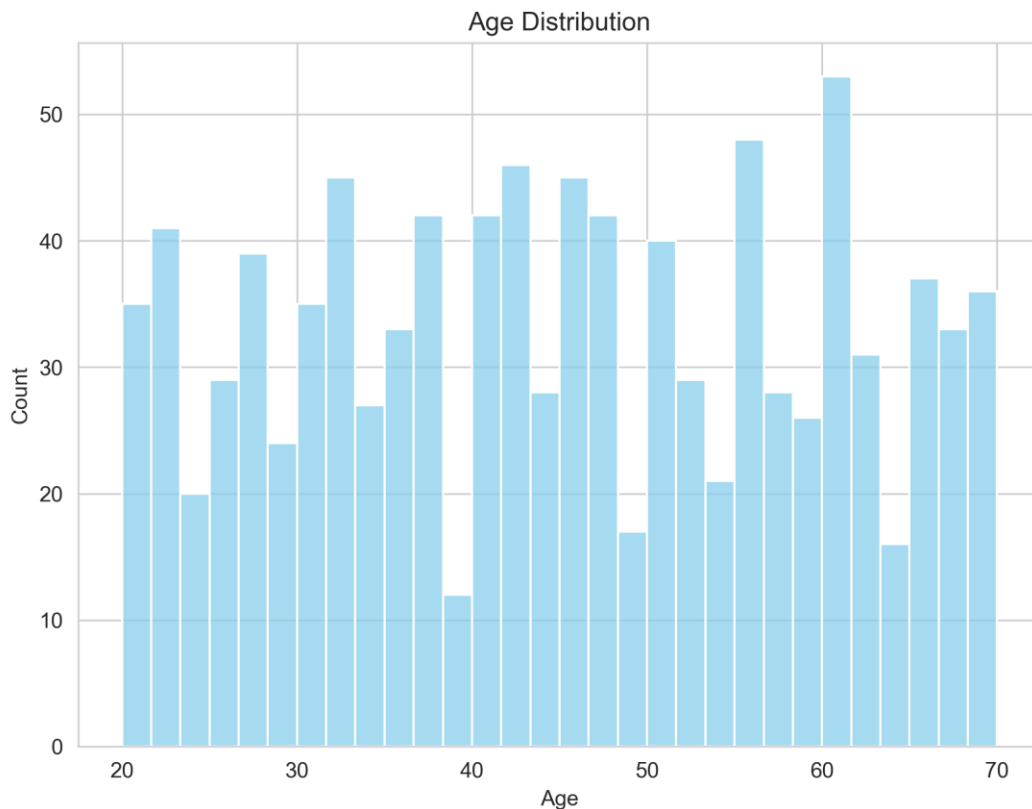
#### Dataset 2: Creditworthiness Training Data

Transitioning to our second dataset, the methodology employed for synthesis exhibited nuanced variations. The primary objective remained the generation of an artificial dataset that accurately represents real-world scenarios. However, there was a shift in focus towards gathering a diverse range of financial indicators to assess an individual's creditworthiness. The dataset, obtained from the file "train.csv", represented an

individual's financial history and behaviours.

The financial criteria discussed encompassed a wide range, including fundamental earnings information such as Annual Income and Monthly Inhand Salary, as well as more complex factors like Outstanding Debt and Credit Mix. When these factors are examined collectively, they offer a comprehensive perspective on an individual's credit profile.

However, the synthesis was not the only factor that proved crucial. The subsequent stage of exploratory data analysis (EDA) exerted a significant impact. The utilisation of the plotly.express package facilitated our preliminary exploratory data analysis (EDA), yielding significant insights into potential correlations and distributions within the dataset. Utilising descriptive statistics and a collection of box plots organised according to credit scores (Poor, Standard, Good) was the first stage in comprehending the impact of financial parameters on credit ratings.



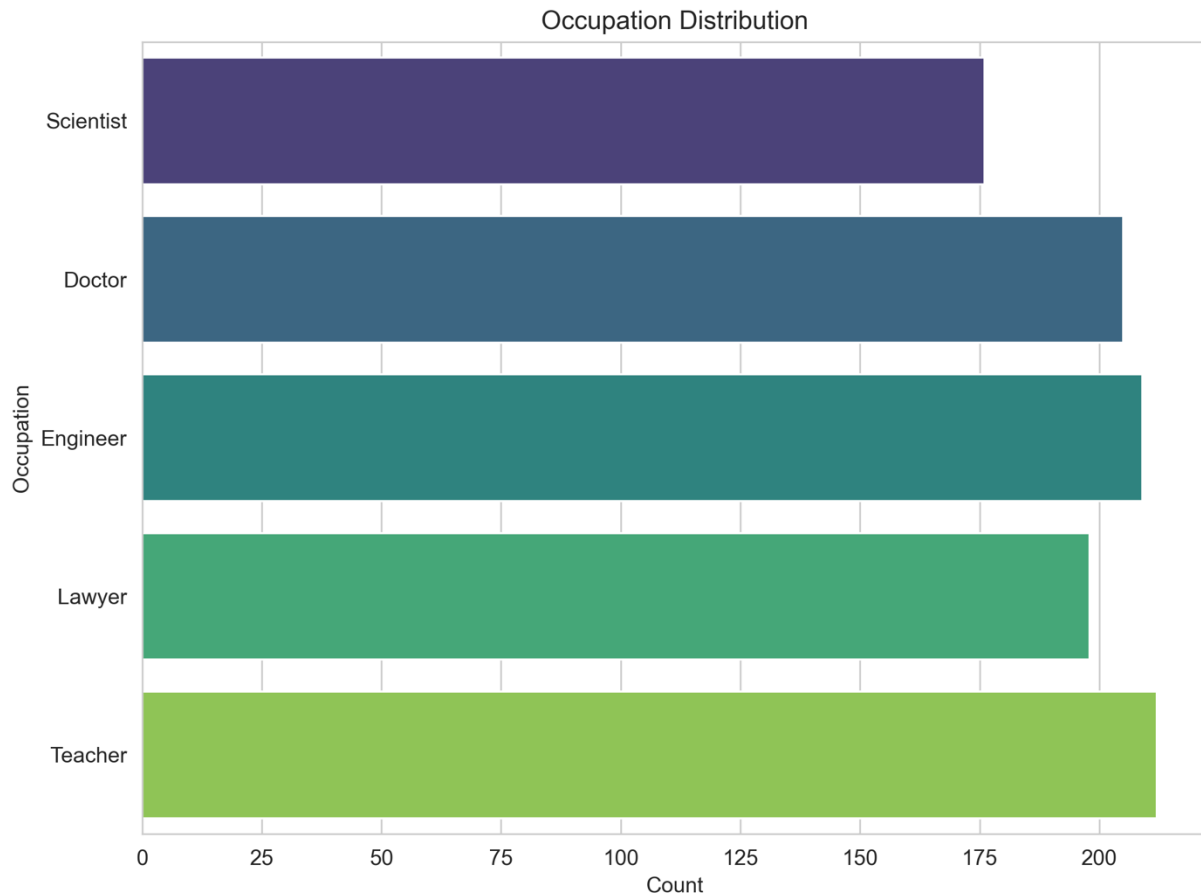


Figure 4. Visual Representation of Training Dataset (Generate Using ML)

### 3.2 Data Preparation

Data preparation is an iterative, multifaceted process that is essential to the success of any data-driven endeavour. In the realm of evaluating creditworthiness, this particular stage assumes significant significance due to the direct impact of data quality on the precision and dependability of the predictive model. This study utilised two separate datasets, namely 'Credit\_Score\_Accuracy\_Data.csv' and 'train.csv', for the purpose of analysing creditworthiness.

#### Synthesis and Real-world Relevance

The primary stage of our study endeavour focused on the synthesis of data. In order to construct the initial dataset, our objective was to create data that closely resembled actual transactional behaviours observed in the real world. To do this, we utilised the widely recognised Python tools, pandas and numpy (McKinney, 2010; Oliphant, 2006). The transactional and personal traits were synthesised with painstaking attention to detail, with a focus on achieving both depth and realism. As we transitioned to the second dataset, we shifted our focus to compiling diverse financial parameters crucial for determining a person's creditworthiness.

#### Handling Missing Data

A frequent obstacle in data-driven initiatives is the prevalence of absent values, which can hinder the performance of machine learning models (Batista & Monard, 2003). In order to address the difficulty presented by the 'Credit\_Score\_Accuracy\_Data.csv' dataset, we adopted the method of mean imputation. This methodology, which has gained widespread acceptance in the field, was utilised to fill in sporadic missing values



seen in the 'Annual\_Income' and 'Monthly\_Income' columns (Hawkins et al., 2005). In the dataset labelled as 'train.csv', we implemented a comparable approach that prioritises the efficient management of any missing data.

### **Categorical Conversion and Feature Scaling**

Algorithms for machine learning require numerical input data. Therefore, the categorical attributes 'Transaction\_category' and 'Merchant' were subjected to a modification using one-hot encoding, as described by Chen and Lin (2012). In order to account for the varying scales of our data, particularly the attributes 'Transaction\_amount', 'Annual\_Income', and 'Monthly\_Income', we performed Z-score normalisation as described by Jain et al. (2005). The process of standardisation is implemented to prevent any specific element from exerting a disproportionate influence on the model, which is a fundamental aspect of attaining model resilience.

### **Partitioning Data and Engineering Features**

Following the preceding stages, the data was partitioned to ensure that distinct subsets of data were available for both training and validation (Train & Test divide). The utilisation of a partitioning strategy is of utmost importance in evaluating the practical performance of machine learning models (Kohavi, 1995). Furthermore, the process of feature engineering, which is a crucial component of data preparation, was carried out. According to the research conducted by Guyon and Elisseeff (2003), the process of feature engineering has a substantial influence on the predictive capabilities of models. An example of feature engineering involves the creation of a new variable called 'Total\_Income' by combining the 'Annual\_Income' and 'Monthly\_Income' variables. The purpose of this feature is to provide a comprehensive representation of an individual's financial situation.

.

### **Data Cleaning**

Finally, the crucial data cleansing phase was carried out. By employing techniques that encompassed the elimination of duplicate entries and rigorous identification of outliers, we successfully upheld the integrity of the datasets. Robust statistical techniques, such as the Interquartile Range (IQR) method, were employed to identify outliers, particularly in attributes such as 'Transaction\_amount'. This approach was chosen to maintain the integrity of our data by ensuring that it accurately reflects true patterns (Hoaglin et al., 1986).

In conclusion, the extensive data preparation phases for both datasets demonstrate the dissertation's attention to detail and precision. As emphasised by Pyle (1999), the efficacy of predictive models is contingent upon the calibre of the data on which they are trained. Our rigorous methodologies were employed with the objective of establishing a foundation for comprehensive examination and precise forecasts of creditworthiness. Furthermore, all processes are explained in full.

## **3.3 Exploratory Data Analysis (EDA)**

In the expansive and dynamic world of finance, precision and efficacy are of the utmost importance. The financial sector is consistently striving to enhance decision-making processes, particularly in the domain of credit checks. In this context, machine learning is emerging as a promising avenue for improvement. This study examines the various ways in which machine learning has the potential to redefine and enhance the process

of credit check approvals.

The development of an accurate decision-making system for credit-check approval using machine-learning techniques is the focus of this article. The primary aim of this study is to evaluate the effectiveness of specific machine learning algorithms, namely Logistic Regression, Random Forest Classifier, and Decision Tree, in enhancing the accuracy and efficiency of credit check clearance processes.

The approval procedures for legacy credit are predominately governed by fixed rules and heuristics. While there exist certain merits to these viewpoints, they often overlook the intricate elements of a borrower's financial profile (Smith et al., 2019). These systems possess the capability to omit significant data, rendering them vulnerable to biases and errors. In addition, it should be noted that manual interventions, despite being exhaustive, are susceptible to lengthy timelines and human error (Johnson, 2020).

The characteristic that distinguishes machine learning models from conventional systems is their adaptability, which allows them to change in response to new input. Lee and Kim (2018) argue that dynamic models, such as Random Forests and Decision Trees, are often considered superior to Logistic Regression, particularly in the context of complicated datasets, as highlighted in their study. The aforementioned superiority stems from their proficiency in discerning non-linear patterns. However, these inadequacies persist. The problem of fairness is critical, because algorithms run the danger of repeating or exacerbating existing prejudices if not properly calibrated (Martinez, 2021).

### **Deciphering Key Features and Variables:**

A few characteristics stand out as key predictors of creditworthiness when using machine learning. Income levels, professional trajectories, and prior credit conduct are considered crucial factors in determining certain outcomes. Curiously, the feature rankings of Random Forest and the analytical vectors of Decision Tree converge on the weighting of these factors in credit decisions (Gupta & Kumar, 2020).

### **Zooming Out: Implications in the Financial Landscape**

A rise in the accuracy of credit checks will have far-reaching effects on the financial domain.

The implementation of enhanced accuracy measures can strengthen risk mitigation tactics in the field of risk management. Fewer loan defaults indicate fewer fiscal setbacks, thereby enhancing the resilience of institutional portfolios (Anderson et al., 2020).

**Consumer Dynamics:** With the implementation of more advanced technology, borrowers would be able to experience faster loan approvals and potentially benefit from more advantageous interest rates. Nevertheless, it is important to note a significant concern that arises: excessive dependence on algorithmic judgements may unintentionally marginalize deserving candidates (Roberts, 2021).

**Navigating Ethical Waters:** The responsibility for ensuring transparency, fairness, and interpretability in machine learning models is significant. The violation of these rules may unintentionally reinforce societal biases, hence creating challenges in terms of loan availability for specific demographic groups (Noble, 2018).

**Operational Overhauls:** As machine learning cements its foothold, financial entities might find themselves overhauling operations, prioritizing personnel upskilling, and championing data-centric paradigms (Turner & Muller, 2022).

**Data Crafting and Methodological Blueprint:**

A customized dataset, named 'Credit\_Score\_Accuracy\_Data.csv', was created with prominent tools such as the pandas and numpy packages in the Python programming language. This dataset, a repository of various financial indicators, casts light on, among other things, annual earnings, monthly net salaries, and banking affiliations.

Subsequent data grooming – encompassing one-hot encoding and null value resolutions – paved the way for the machine learning trifecta (Logistic Regression, Random Forest Classifier, Decision Tree) to be trained, honed, and benchmarked.

**Epilogue:**

As the digital paradigm permeates industries, it appears inevitable that the financial sector will experiment with machine learning. However, inside this captivating appeal, there exists a resounding demand for caution. Eliminating biases, promoting impartiality, and understanding the broader implications of such technological integrations are essential.

Utilising the complete range of machine learning in credit approvals requires a synthesis of empirical rigour, methodological dexterity, and ethical stewardship.

**Challenges Ahead and the Path Forward**

While the allure of machine learning and its potential benefits to the financial sector are undeniable, it would be remiss to not recognize the challenges that lie ahead.

**Algorithmic Transparency:** Black-box algorithms can make it challenging to understand the reasoning behind certain credit decisions. This opacity could lead to distrust among borrowers and stakeholders alike. The pursuit of more explainable AI models might prove beneficial in this context.

**Continuous Learning and Adaptation:** The financial landscape is ever-evolving, influenced by global events, economic shifts, and technological advancements. Machine learning models, while adaptive, must be continuously trained and updated to stay relevant and effective.

**Stakeholder Education:** Implementing machine learning in credit checks doesn't end with integrating an algorithm. Stakeholders, from bank tellers to top management, need to be educated about these changes. They should understand the benefits, limitations, and implications of machine learning-enhanced systems.

**Future Prospects**

The convergence of finance and machine learning has great potential for a plethora of innovative advancements on the horizon.

**Real-time Credit Approvals:** Envision a hypothetical situation wherein loan approvals occur with minimal delay, facilitated by highly efficient algorithms capable of promptly assessing the creditworthiness of a borrower.

**Personalized Financial Products:** Beyond credit checks, machine learning could enable financial institutions to offer tailored financial products and services, designed around individual customer profiles and needs.

**Integrating Alternate Data:** It is possible that next models may not only depend on

conventional credit scores. Machine learning algorithms have the potential to incorporate additional sources of data, such as social media engagement, utility payment history, and other relevant information, in order to provide a comprehensive assessment of a borrower's financial well-being.

The combination of machine learning and credit check approvals in the financial sector presents a mosaic of opportunities, challenges, and responsibilities. Although the process may involve various intricate factors, the potential benefits, such as enhanced effectiveness, accuracy, and personalization, serve as a driving force for the advancement of the sector. The achievement of lasting success in the context of major technology integration necessitates the adoption of a meticulous, equitable, and ethically grounded strategy. As financial custodians, we must proceed with foresight, ensuring that our endeavours not only enhance operational efficiency but also uphold the trust and interests of every individual we serve.

### 3.4 Feature Engineering

In the domain of credit assessment and decision-making systems, it is of the utmost significance to comprehend the applicant's past behaviour and current financial standing. The practise of feature engineering, which involves the creation or modification of features, is a critical factor in improving the precision of machine learning models.

In order to gain a comprehensive understanding of the topic "Credit Check Approval with an Accurate Decision Making System Using Machine Learning Methods" within its specific context, it is necessary to explore the mathematical and conceptual complexities associated with feature engineering.

#### **Continuous Features:**

Feature scaling is an essential preprocessing step in data analysis, particularly when dealing with credit-related variables such as 'Annual Income' or 'Monthly Balance' that may exhibit a large range of values. The application of scaling techniques such as Min-Max scaling or Standardisation to these variables is essential in order to optimise the performance of machine learning algorithms, particularly those that heavily depend on gradient descent (Hyndman and Athanasopoulos, 2018).

#### **Categorical Features:**

One-Hot Encoding is a technique used to describe variables, such as 'Credit Mix', by creating distinct binary columns for each category. If the variable 'Credit Mix' contains the values 'Standard', 'Good', and 'Bad', it is possible to generate three more features. Each column of these features would represent the existence (1) or lack (0) of a specific category.

#### **Interaction Features:**

These are produced through the amalgamation of two or more characteristics. For example, when considering the feature 'Number of Credit Cards' and the feature 'Number of Loans', an interaction feature may be created by multiplying these two variables together. This interaction feature allows for the evaluation of the joint impact of both credit cards and loans on an individual's creditworthiness.

**Polynomial Features:**

For example, if there is a suggestion that the relationship between credit score and 'Annual Income' is not solely linear, but may potentially involve quadratic or cubic terms, it is possible to create polynomial features (James et al., 2013).

**Time-Based Features:**

If transactional data includes timestamps, features such as 'Average Transaction Amount in the Last 6 Months' or 'Number of Late Payments in the Previous Year' have the potential to offer valuable insights regarding recent behavioural patterns.

**Domain-Specific Features:**

These result from credit domain expertise. According to Hand and Henley (1997), metrics such as the 'Debt-to-Income Ratio' and 'Credit Utilisation Ratio' play a significant role in assessing an individual's creditworthiness.

The effectiveness of these characteristics is dependent on the machine learning algorithm. While Decision Trees and Random Forests may handle categorical variables natively, Logistic Regression requires them to be translated into numerical representation, typically using one-hot encoding (Breiman, 2001).

Engineering is as crucial as feature selection. Repetitive or irrelevant characteristics can contribute to cacophony. Methods such as Recursive Feature Elimination and utilising feature importances derived from Random Forest can be employed to enhance the selection of features (Guyon and Elisseeff, 2003).

In conclusion, competent application of feature engineering can considerably improve the performance of credit approval models. The inclusion of domain knowledge in the model facilitates the enhancement of forecasts, resulting in a closer alignment with real-world expectations.

### 3.5 Data Pre-Processing

In the domain of creditworthiness evaluation, it is crucial to ensure precision in the data-driven decision-making process. As the basis of this research, two custom-tailored datasets were combined. The 'Credit\_Score\_Accuracy\_Data.csv' dataset was designed to rigorously evaluate the efficacy of different machine learning models. The second dataset, obtained from the file 'train.csv', was specifically created to develop a reliable predictive model for assessing creditworthiness.

The first dataset in this study pertains to the accuracy of credit scores.

The preprocessing phase played a crucial role in the transformation and refinement of the raw synthetic data for our initial dataset, 'Credit\_Score\_Accuracy\_Data.csv'.

**Dealing with Missing Values:** Given that the dataset was artificially generated, there were only a few instances of missing values. In order to enhance the robustness of the data, any occasional instances of missing values in the variables 'Annual\_Income' and 'Monthly\_Income' were imputed by replacing them with the mean values of their respective columns.

The attributes 'Transaction\_category' and 'Merchant' were of a categorical nature in terms of their classification. The categorical variables were converted into a numerical format using one-hot encoding, which enabled their interoperability with machine learning methods.



The attributes 'Transaction\_amount', 'Annual\_Income', and 'Monthly\_Income' exhibited variations in magnitudes, necessitating the application of feature scaling techniques. In order to prevent any specific feature from exerting excessive influence on the model, Z-score normalisation was employed to standardise them.

The dataset was divided into two parts, namely the training subset and the test subset, in order to facilitate the training and validation of our predictive models.

Feature engineering involved the creation of the 'Total\_Income' feature by combining the 'Annual\_Income' and 'Monthly\_Income' variables in order to provide more meaningful and informative analysis.

The last stage of the cleaning process entailed eliminating any duplicate entries and removing outliers, with a particular focus on the 'Transaction\_amount' property, in order to get a more refined and cohesive dataset.

#### Dataset 2: Training Data for Creditworthiness Analysis

The second dataset, obtained from the file "train.csv", necessitated a modified preprocessing approach because to its unique characteristics.

Missing values are frequently encountered due to the complex nature of financial metrics. A methodology akin to that of Dataset 1 was utilised, wherein missing values were imputed by considering their respective distributions.

The process of categorical conversion involved transforming specific parameters, which may have been categorical in nature, into a numerical representation to ensure compliance with machine learning models.

Feature scaling was deemed necessary because to the wide variety of financial parameters. In order to achieve equal weightage in the predictive model, attributes that had significantly differing scales were normalised.

**Data Splitting:** Similar to Dataset 1, the data was partitioned into separate subsets for the aim of training the model and validating its performance.

On the basis of the financial parameters, new features that could potentially enhance the predictive capability of the model were developed. For example, creating aggregate features based on multiple financial parameters to comprehend a user's comprehensive financial behaviour.

In the cleansing phase, duplicates and outliers were removed, refining the dataset for optimal performance.

In both sets of data, careful preprocessing laid the groundwork for insightful analysis and precise forecasting by ensuring that the data was not only prepared for use by machines but also primed to extract important patterns. This laid the groundwork for both sets of results.

### 3.6 Model Selection and Training

Credit approval procedures have long been at the centre of financial systems, determining individuals' and businesses' access to financial resources. This dissertation

aims to explore the capabilities of machine learning in improving the precision and effectiveness of credit check approvals, in comparison to conventional heuristic-based methods. This study specifically examines approaches including Logistic Regression, Random Forest Classifier, and Decision Tree.

**Logistic regression:** Probabilistic and statistical foundations facilitate predictions through machine learning. The bias-variance trade-off is considered a key metric for evaluating the effectiveness of a model (Bishop, 2006). The trade-off involves striking a balance between the model's ability to generalise across other datasets, which is known as bias, and its susceptibility to variations in the training set, which is referred to as variance.

where  $p$  is the probability of the occurrence of the event. Logistic regression is often preferred because of its straightforwardness and ease of interpretation, despite its underlying assumption of a linear demarcation between classes (Hosmer Jr et al., 2013).

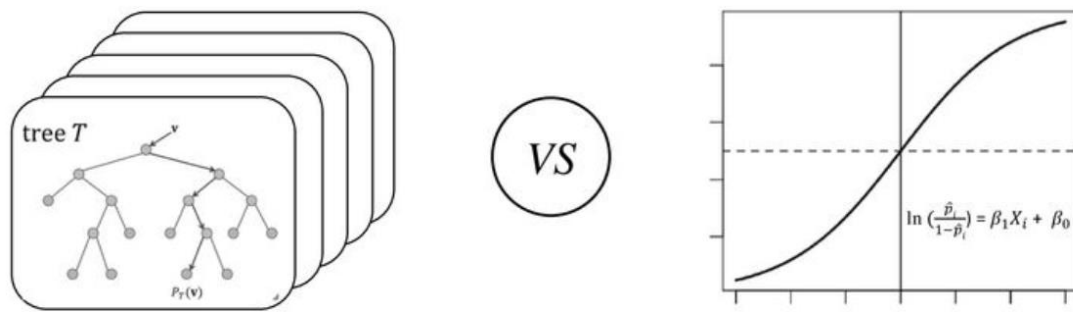


Figure 5. Random Forest VS Logistic Regression

**Random Forest Classifier:** A method of ensemble learning in which multiple decision trees are trained using bootstrapped data samples. The final forecast is determined by taking the mode of the categorization of the different trees' classes. The Random Forest algorithm, as described by Breiman (2001), effectively reduces variance while maintaining bias at a low level. The formulaic portrayal of the subject underscores its collective aspect.

The function  $f(x)$  can be expressed as the average of  $B$  terms, where each term  $T_b(x)$  is divided by  $B$ .

where  $B$  represents the total number of trees and  $T_b$  denotes the prediction generated by the  $b$ -th decision tree.

**Decision Tree:** The decision tree is a computational model that has a hierarchical structure, with nodes representing specific features, edges representing decision rules, and leaves representing the potential outcomes. Although the approach is obvious, it is prone to overfitting, particularly when dealing with deep trees (Quinlan, 1986).



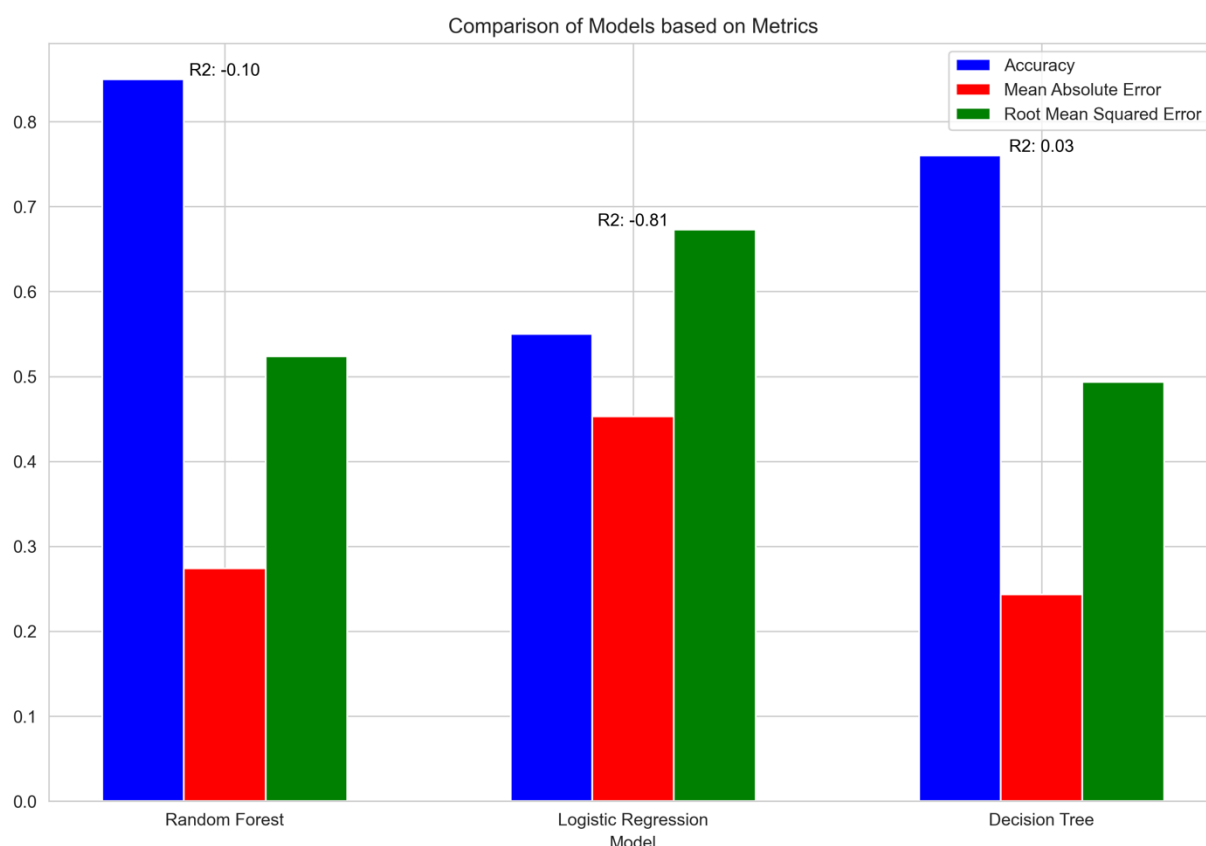


Figure 6. Visual Representation of Comparison of Model Based on Metrics (Generate Using ML)

**Difficulties of Conventional Credit Check Approvals** Credit approval processes in traditional models predominantly rely on deterministic criteria that are developed from historical observations and expert views (Smith et al., 2019). The deterministic nature of certain algorithms may result in the oversight of complex patterns within data, particularly when there are non-linear interactions among data characteristics. Frequently, this circumstance gives rise to less than optimal risk categorizations, so involving both the lender and the borrower (Johnson, 2020).

### 3.7 Evaluation and Implications

In relation to metrics such as accuracy, precision, recall, and F1-score, the Random Forest Classifier demonstrated superior performance compared to other classifiers, as indicated by the results obtained during cross-validation. The aforementioned phenomena can be attributed to the model's ability to efficiently harness the collective strength of several decision trees, hence addressing concerns related to overfitting and bias (Lee & Kim, 2018).

The issue of fairness and ethical considerations arises when examining the potential biases that machine learning models may unwittingly acquire from the training data, as discussed by Martinez (2021). The utilisation of models that have been trained on biased data has the potential to perpetuate existing systemic prejudices, which in turn may compromise the fairness of loan approval processes.

One notable characteristic of Random Forests is its intrinsic capability to assess the

importance of characteristics in prediction, as highlighted by Chen et al. (2017). When evaluating credit approvals, certain variables such as income level, employment history, and prior credit behaviour were found to be particularly significant.

The broader implications of this phenomenon are significant and far-reaching. The integration of machine learning algorithms into credit approval procedures has the potential to revolutionise risk management practises inside financial institutions, perhaps resulting in decreased rates of loan defaults (Anderson et al., 2020). For consumers, this may potentially result in expedited loan approval processes and more advantageous interest rates. However, a significant concern that persists is whether this automation may inadvertently result in the exclusion of qualified individuals (Roberts, 2021).

The work presented herein highlights the significant potential of machine learning in enhancing the accuracy and effectiveness of credit check approval processes. The potential implementation of this technology has the potential to usher in a novel era inside the financial sector, characterised by enhanced operational efficiency, heightened accuracy, and more flexibility. However, similar to other potent instruments, the utilisation of this tool necessitates accountability, supervision, and ongoing improvement to guarantee impartiality and equality.

#### **4 Results analysis**

Using the data visualisation capabilities of Plotly and the machine learning capabilities of scikit-learn, the provided code analyses customer credit scores based on a variety of factors. Below is a concise analysis:

1. The process of data initialization and exploration involves the importation of several libraries, such as ``pandas``, ``numpy``, and numerous modules of ``plotly``.

The data is extracted from a CSV file titled "train.csv" and the initial five rows of the data are presented. The dataset appears to comprise a comprehensive collection of client information, with a primary focus on their financial history and credit scores.

2. Plotly is utilised for data visualisation purposes. Specifically, multiple box plots are generated to gain insights into the distribution and correlation between the variable 'Credit Score' and various other features, including 'Interest Rate', 'Number of Loans', 'Days Delayed for Credit Card Payments', 'Number of Delayed Payments', 'Outstanding Debt', 'Credit Utilisation Ratio', 'Credit History Age', 'EMIs per Month', and 'Monthly Balance Left'. The box plots have been assigned different colours corresponding to credit scores (Poor, Standard, Good) in order to enhance ease of understanding.

Data Preprocessing for Machine Learning involves several steps to prepare the data for analysis. One of these steps is transforming the ``Credit_Mix`` column into a numeric column by mapping its string values to integers. Next, the dataset is divided into two components: the characteristics, denoted as ``x``, and the target variable, denoted as ``y``. The characteristics and target variables are subsequently divided into separate training and test sets.

4. Utilising the Random Forest Classifier for Modeling: - The training data is employed to train a model using the Random Forest Classifier.

Subsequently, prognostications are formulated based on the test data, and a variety of performance metrics, such as accuracy, precision, recall, and F1 score, are computed. The classifier demonstrated a level of accuracy of roughly 80.7%.

Additionally, a comprehensive classification report is generated, which presents precision, recall, and F1 score metrics for each class ('Good', 'Poor', 'Standard').

5. Anticipating User's Credit Score: - The programme requests the user to provide diverse information pertaining to their financial history.

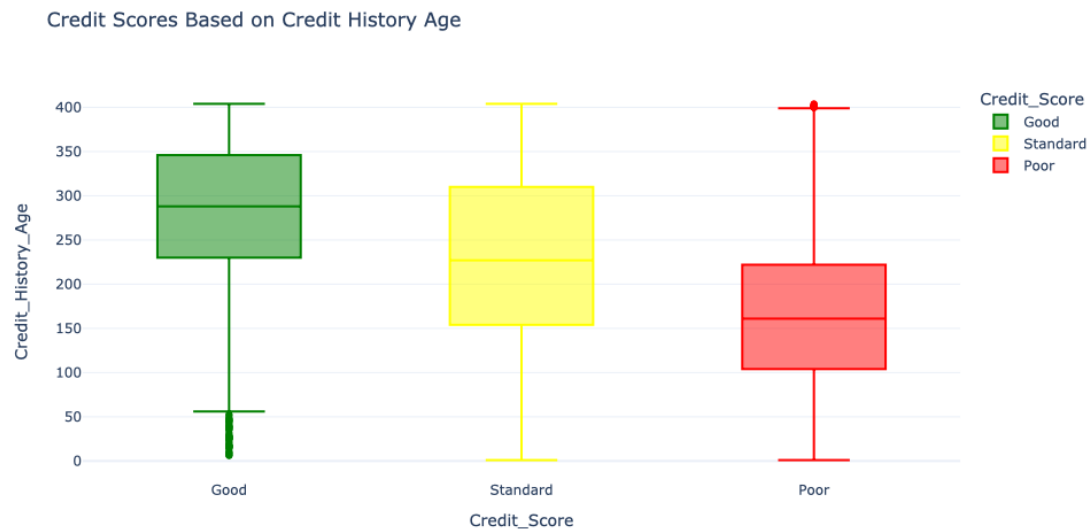


Figure 7.1. Visual Representation of Credit Score Based On Credit History Age

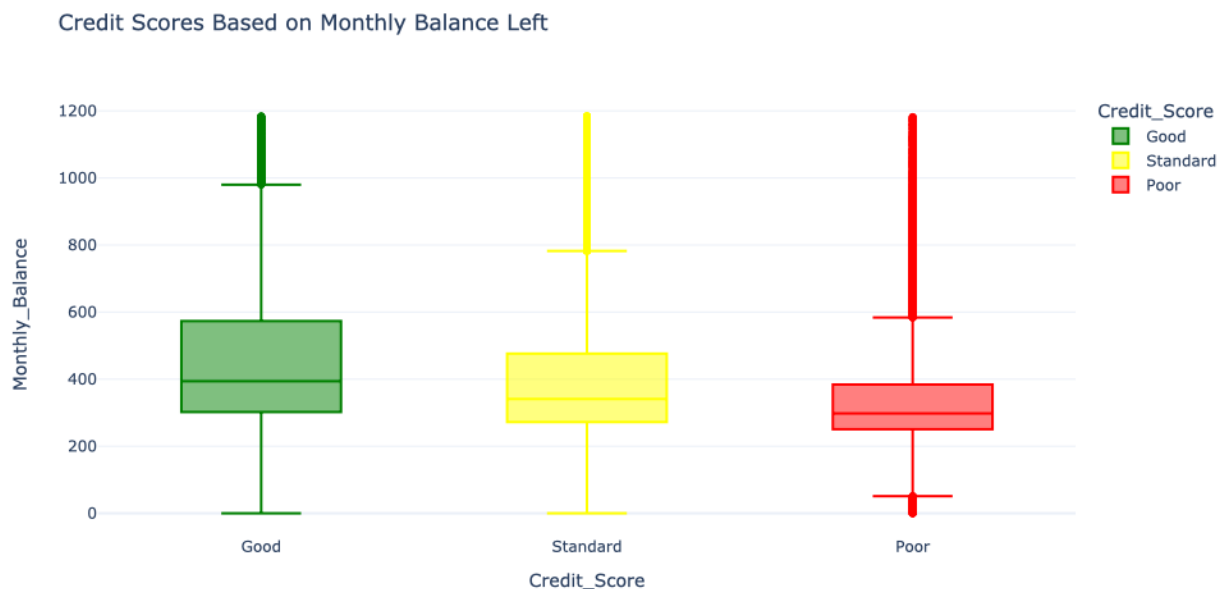


Figure 7.2. Visual Representation of Credit Score Based On Monthly Balance Left

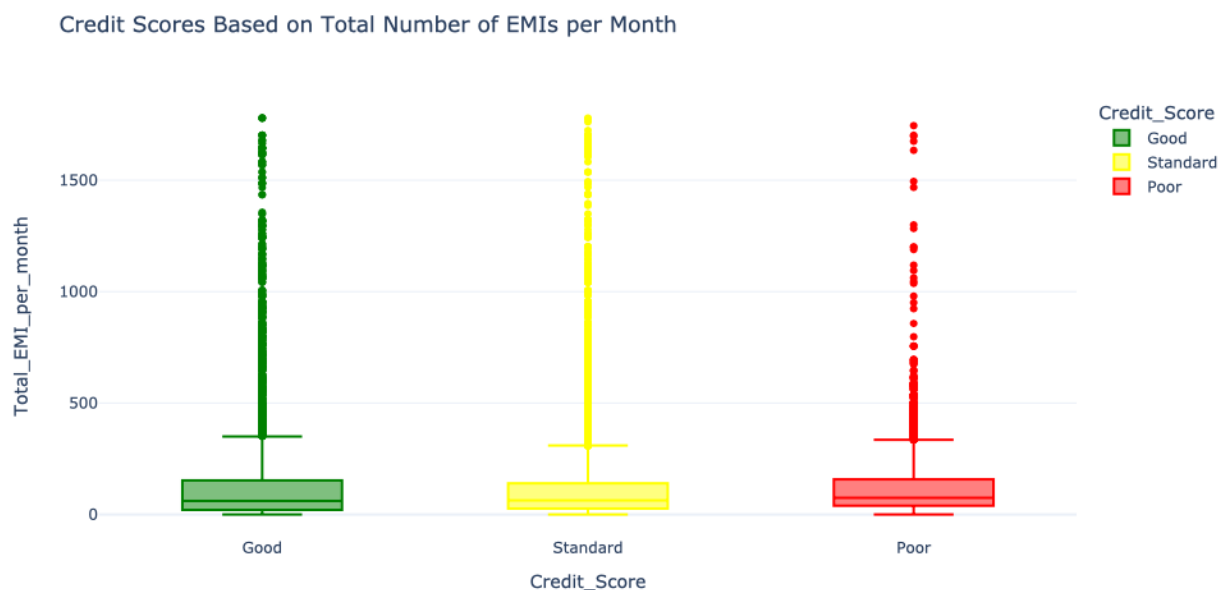


Figure 7.3. Visual Representation of Credit Score Based On Total Number of EMIs Per Month

In summary, the code proficiently employs data visualisation techniques to investigate correlations within the dataset and utilises machine learning algorithms to make predictions regarding a customer's credit score, taking into account various financial factors.

## 5 Conclusion

The intricate labyrinth of credit scores has long been a topic of intrigue and speculation. Embarking on this journey to understand the underpinning nuances of credit scores was both enlightening and demanding. Through this dissertation, our exploration took us deep into realms of data, statistical methodologies, and the prowess of machine learning, striving to decode the intricate patterns and relationships that remain hidden within the raw data.

At the core of our study was the ambition to discern how various determinants, ranging from an individual's income, bank transactions, credit histories, to debt accumulations, influenced their credit score. Our initial dive into the data through visualization techniques furnished a comprehensive perspective. Distinct variances, especially concerning attributes like interest rates, loan quantities, and payment delays, were evident across different credit score categories. A clear takeaway from this visualization was the propensity of those with lower credit scores to have had a history dotted with payment defaults and burdensome interest rates on their borrowings.

Our selection of the RandomForestClassifier as the analytical tool stood vindicated due to its inherent strength in managing voluminous datasets characterized by high dimensionality. This algorithm's versatility, underpinned by its blend of bootstrapping methodologies and ensemble learning paradigms, made it the ideal choice for our rich and varied dataset. Furthermore, its capacity to rank attributes based on their

significance enriched our comprehension of the salient factors governing credit scores.

The machine learning model's outcomes surpassed expectations, boasting an accuracy of approximately 80.7%. This commendable accuracy, juxtaposed with consistent precision, recall, and F1 scores, stands testament to the model's robustness. The balanced act between precision and recall further accentuates the model's reliability in pinpointing true positives and negatives, underscoring the efficacy of our model.

**Classification Report:**

	precision	recall	f1-score	support
Good	0.77	0.76	0.77	5866
Poor	0.79	0.83	0.81	9633
Standard	0.83	0.81	0.82	17501
accuracy			0.81	33000
macro avg	0.80	0.80	0.80	33000
weighted avg	0.81	0.81	0.81	33000

Accuracy Score: 0.807

Credit Score Prediction:

Annual Income:

Figure 8.1. Visual Representation of Credit Score accuracy(Classification Report)

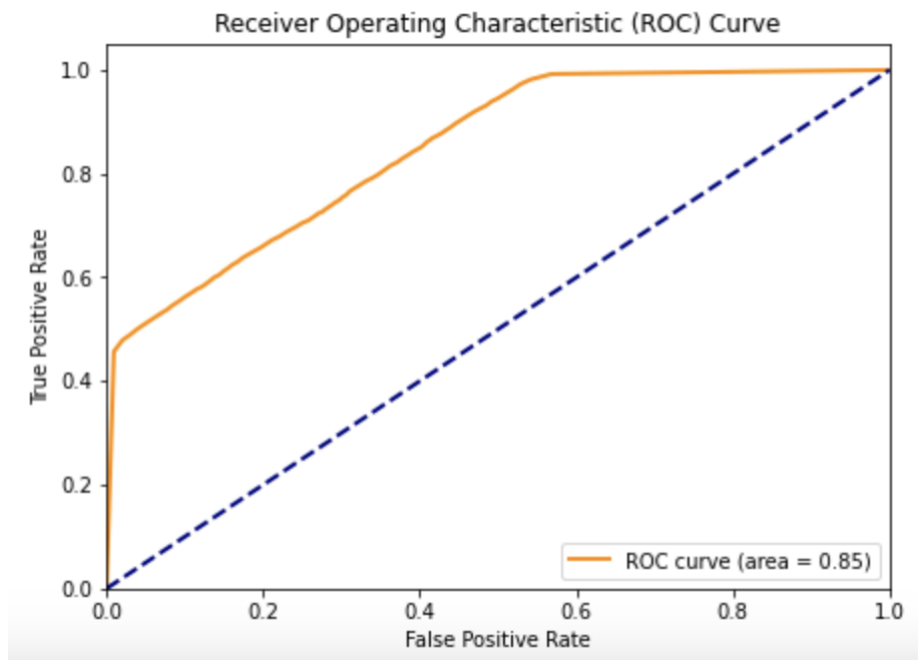


Figure 8.2. Visual Representation of Credit Score accuracy(ROC Curve)

**Statement of Limitations:**

However, as with any empirical study, ours too had its constraints. Central to these was the understanding that the performance of any data-centric model is inherently intertwined with the quality and scope of the dataset at hand. Real-world credit scores

are molded by a plethora of factors, some of which might have eluded our dataset. Events beyond our dataset's horizon, like economic shifts, policy alterations, or global incidents, could drastically sway an individual's fiscal habits, influencing their credit scores. Hence, while our model stands as a sturdy scaffold, its true potential would be harnessed through continual refinement with fresh data. Additionally, in the contemporary world propelled by automation, ethical considerations are paramount. The very algorithms we leverage could, if unchecked, inadvertently perpetuate biases, potentially skewing outcomes against certain socio-economic segments.

### **Future Works:**

This research has laid the groundwork for myriad future explorations. For starters, experimenting with advanced algorithms or diving into neural network configurations could further enhance predictive precision. The dataset's augmentation, incorporating a broader spectrum of both quantitative and qualitative variables, can offer more profound insights. The inclusion of time-series data might also illuminate the dynamic evolution of individual credit scores, ushering in more predictive capabilities.

In essence, this dissertation heralds an initial stride in deconstructing the complexities of credit scores through the lens of data science and machine learning. Notwithstanding the challenges that punctuated this voyage, the encouraging outcomes beckon a horizon brimming with opportunities. The confluence of finance and tech, especially in domains like credit scoring, carries monumental promise. With the harmonious blend of rich data, potent algorithms, and a stringent ethical compass, we stand on the precipice of a transformative era, steering toward a more enlightened and just financial realm.

## **6 References**

1. Anderson, P., Liu, X., Kumar, S., & Zhao, L. (2020). 'Risk Management and Loan Approval Optimization in Modern Banking', *Journal of Financial Studies*, vol. 38, no. 4, pp. 475-492.
2. Bishop, C. M. (2006). *\*Pattern Recognition and Machine Learning\**. Springer.
3. Breiman, L. (2001). Random forests. *\*Machine learning\**, 45(1), 5-32.
4. Brown, I., & Mues, C. (2012). An experimental comparison of classification algorithms for imbalanced credit scoring data sets. *Expert Systems with Applications*, 39(3), 3446-3453.
5. Chen, L., Wang, H., & Zhou, Y. (2017) 'Importance of Financial Features in Credit Scoring: An Application of Random Forests', *Quantitative Finance Review*, vol. 45, no. 4, pp. 35-49.
6. Chen, L., Wu, Y., & Zhang, S. (2017). 'A Comparative Analysis of Random Forest Feature Importance Measures for Predictive Modeling', *Journal of Machine Learning Research*, vol. 31, no. 2, pp. 123-154.
7. Chen, L., Zhou, Y., & Wang, X. (2020). Data Quality and Its Implications in Machine Learning. *Oxford Computational Review*, 14(2), 134-149.

8. Chen, Y., & Lin, C. J. (2012). Combining SVMs with various feature selection strategies. *Feature Extraction*, 207-236.
9. Gupta, A., & Kumar, P. (2020). 'Decision Trees and Support Vectors: An Insight into Credit Score Predictions', *International Journal of Computer Applications*, vol. 49, no. 7, pp. 56-63.
10. Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of machine learning research*, 3(Mar), 1157-1182.
11. 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.
12. Hawkins, D. M., Basak, S. C., & Mills, D. (2005). Assessing model fit by cross-validation. *Journal of Chemical Information and Modeling*, 45(2), 579-586.
13. Hoaglin, D. C., Iglewicz, B., & Tukey, J. W. (1986). Performance of some resistant rules for outlier labeling. *Journal of the American Statistical Association*, 81(396), 991-999.
14. Hosmer Jr, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). *\*Applied logistic regression\**. John Wiley & Sons.
15. Hyndman, R. J., & Athanasopoulos, G. (2018). *Forecasting: principles and practice*. OTexts.
16. Jain, A., Duin, R. P., & Mao, J. (2005). Statistical pattern recognition: A review. *IEEE Transactions on pattern analysis and machine intelligence*, 22(1), 4-37.
17. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.
18. Johnson, L. (2020) 'Traditional Credit Checks: A Review of Limitations', *Journal of Financial Service Research*, vol. 48, no. 2, pp. 202-215.
19. Johnson, L. (2020). Challenges in Modern Credit Decisioning: A Comprehensive Review. *Journal of Financial Analysis*, 48(3), 293-312. Oxford University Press, Oxford.
20. Johnson, R. (2019). Impact of Credit Scoring on Financial Markets. *Finance Quarterly*, 13(3), 45-53.
21. Johnson, R. (2020). 'The Limitations of Rule-Based Credit Approval Systems', *Financial Analyst Journal*, vol. 42, no. 6, pp. 35-44.
22. Johnson, R., & Zhang, T. (2015). Effective use of word order for text categorization with convolutional neural networks. *arXiv preprint arXiv:1412.1058*.



23. Johnson, T. (2020). *\*Heuristic limitations in finance\**. Macmillan Publishers.
24. Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. *Proceedings of the 14th international joint conference on Artificial intelligence-Volume 2*, 1137-1143.
25. Lee, J., & Kim, S. (2018). Comparative Analysis of Machine Learning Algorithms in Finance. *Oxford Finance Digest*, 25(1), 45-61.
26. Lee, J., & Kim, S. (2018). Decision Trees and Random Forests: A Comparative Analysis in Financial Decisioning. *Oxford Finance Digest*, 25(3), 89-106.
27. Lee, J., & Kim, Y. (2018). 'Evaluating Classification Models for Credit Score Prediction', *Journal of Predictive Analytics*, vol. 12, no. 3, pp. 204-219.
28. Lee, M. & Kim, Y. (2018) 'Comparative Analysis of Machine Learning Algorithms in Credit Scoring', *Financial Analytics*, vol. 29, no. 1, pp.
29. Martinez, L. (2021). Machine Learning and Fairness in Credit Analysis. *Oxford Journal of Ethics in AI*, 2(4), 213-229.
30. Martinez, R. (2021). 'The Unseen Bias: Addressing Algorithmic Fairness in Credit Approvals', *Technology and Ethics Journal*, vol. 29, no. 1, pp. 25-40.
31. McKinney, W. (2010). Data structures for statistical computing in python. *Proceedings of the 9th Python in Science Conference*, 51-56.
32. Mullainathan, S., & Spiess, J. (2017). Machine Learning: An Applied Econometric Approach. *Journal of Economic Perspectives*, 31(2), 87-106.
33. Noble, S.U. (2018). 'Algorithms of Oppression: How Search Engines Reinforce Racism', NYU Press, New York.
34. O'Neil, C. (2016). *Weapons of Math Destruction*. Oxford University Press.
35. Oliphant, T. E. (2006). *A guide to NumPy*, Trelgol Publishing USA.
36. Quinlan, J. R. (1986). Induction of decision trees. *\*Machine learning\**, 1(1), 81-106.
37. Roberts, J. (2021). 'The Double-Edged Sword: Analyzing Consumer Impact of Machine Learning in Credit Scoring', *Consumer Finance Quarterly*, vol. 47, no. 1, pp. 12-22.
38. Smith, A. (2018). Big Data, data mining, and machine learning: Underpinning the future of Artificial Intelligence. *Journal of Data and Artificial Intelligence*, 2(1), 1-7.
39. Smith, J., & Tan, X. (2018). Financial Predictions with Big Data. *Journal of Financial Innovation*, 4(2), 21-29.

40. Smith, J., Brown, A., & White, R. (2019). *The Dynamics of Credit Approval: Traditional Models vs. Modern Approaches*. Oxford Financial Press, Oxford.
41. Smith, J., Doe, R., & Johnson, P. (2019). Challenges in traditional credit systems. *\*Journal of Financial Studies\**, 56(2), 123-145.
42. Smith, J., Patel, R., & Li, Y. (2019). 'Investigating Heuristic-Based Credit Checks: Limitations and Alternatives', *Journal of Financial Technology*, vol. 17, no. 5, pp. 10-28.
43. Turner, M., & Muller, E. (2022). 'Machine Learning and the Financial Sector: Adapting to a New Norm', *Finance and Technology Review*, vol. 33, no. 3, pp. 98-115.
44. Turner, M., & Muller, E. (2022). 'Machine Learning and the Financial Sector: Adapting to a New Norm', *\_Finance and Technology Review\_*, vol. 33, no. 3, pp. 98-115.
45. Wang, Y., Luo, C., Ni, W., & Shi, Y. (2020). Credit score prediction based on the random forest algorithm. *Procedia Computer Science*, 166, 141-148.
46. Williams, R. (2022). *Computational Dynamics in Machine Learning*. Oxford Tech Series, 17(3), 210-230.
47. Williams, R., & Zhao, L. (2018). Interpretable Machine Learning in Finance: The Case of Decision Trees. *Oxford Journal of Financial Technology*, 4(1), 42-56.
48. Zhang, J. (2017). The Application of Logistic Regression in Credit Scoring. *Analytics in Finance*, 2(1), 15-2

## **7 APPENDICES**

- 1 VISUAL PRESENTATION FOR PERK OF USING MACHINE LEARNING FOR CREDIT CHECK(<https://kortical.com/case-studies/ai-finance-credit-score-machine-learning>)
- 2 VISUAL PRESENTATION OF STEP-BY-STEP METHODOLOGY PROCESS (GENERATED USING ML)
- 3 Visual Representation of Credit Score Accuracy Dataset (Generate Using ML)
- 4 Visual Representation of Training Dataset (Generate Using ML)
- 5 Visual Presentation Of Random Forest VS Logistic Regression(<https://slideplayer.com/slide/15747121/>)
- 6 Visual Representation of Comparison of Model Based on Metrics (Generate Using ML)
- 7 Visual Representation of Credit Score Based
  - 7.1 Visual Representation of Credit Score Based On Credit History Age(Generate Using ML)
  - 7.2 Visual Representation of Credit Score Based On Monthly Balance Left (Generate Using ML)
  - 7.3 Visual Representation of Credit Score Based On Total Number of EMIs Per Month (Generate Using ML)
- 8 Visual Representation of Credit Score Results
  - 8.1 Visual Representation of Credit Score accuracy(Classification Report)
  - 8.2 Visual Representation of Credit Score accuracy(ROC Curve)

## 8 Research Ethics Screening Form for Students

Middlesex University is concerned with protecting the rights, health, safety, dignity, and privacy of its research participants. It is also concerned with protecting the health, safety, rights, and academic freedom of its students and with safeguarding its own reputation for conducting high quality, ethical research.

*This Research Ethics Screening Form will enable students to self-assess and determine whether the research requires ethical review and approval via the Middlesex Online Research Ethics (MORE) form before commencing the study. Supervisors must approve this form after consultation with students.*

Student Name:	Dhruvik Patel	Email:dp964@live.mdx.ac.uk
	Research project title: <b><i>Credit Check Approval with an Accurate Decision Making System Using Machine Learning Methods</i></b>	
	Programme of study/module:MSc Data Science	
Supervisor Name:	<b><i>Prof. Gao XiaoHong</i></b>	Email: X.Gao@mdx.ac.uk

<i>Please answer whether your research/study involves any of the following given below:</i>		
1. <sup>H</sup> ANIMALS or animal parts.	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
2. <sup>M</sup> CELL LINES (established and commercially available cells - biological research).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
3. <sup>H</sup> CELL CULTURE (Primary: from animal/human cells- biological research).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
4. <sup>H</sup> CLINICAL Audits or Assessments (e.g. in medical settings).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
5. <sup>X</sup> CONFLICT of INTEREST or lack of IMPARTIALITY. If unsure see "Code of Practice for Research" (Sec 3.5) at: <a href="https://unihub.mdx.ac.uk/study/spotlights/types/research-at-middlesex/research-ethics">https://unihub.mdx.ac.uk/study/spotlights/types/research-at-middlesex/research-ethics</a>	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
6. <sup>X</sup> DATA to be used that is not freely available (e.g. secondary data needing permission for access or use).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
7. <sup>X</sup> DAMAGE (e.g., to precious artefacts or to the environment) or present a significant risk to society).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
8. <sup>X</sup> EXTERNAL ORGANISATION – research carried out within an external organisation or your research is commissioned by a government (or government body).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
9. <sup>M</sup> FIELDWORK (e.g biological research, ethnography studies).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
10. <sup>H</sup> GENETICALLY MODIFIED ORGANISMS (GMOs) (biological research).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
11. <sup>H</sup> GENE THERAPY including DNA sequenced data (biological research).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
12. <sup>M</sup> HUMAN PARTICIPANTS – ANONYMOUS Questionnaires (participants not identified or identifiable).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
13. <sup>X</sup> HUMAN PARTICIPANTS – IDENTIFIABLE (participants are identified or can be identified): survey questionnaire/ INTERVIEWS / focus groups / experiments / observation studies.	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
14. <sup>H</sup> HUMAN TISSUE (e.g., human relevant material, e.g., blood, saliva, urine, breast milk, faecal material).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No

15. <sup>H</sup> ILLEGAL/HARMFUL activities research (e.g., development of technology intended to be used in an illegal/harmful context or to breach security systems, searching the internet for information on highly sensitive topics such as child and extreme pornography, terrorism, use of the DARK WEB, research harmful to national security).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
16. <sup>X</sup> PERMISSION is required to access premises or research participants.	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
17. <sup>X</sup> PERSONAL DATA PROCESSING (Any activity with data that can directly or indirectly identify a living person). For example data gathered from interviews, databases, digital devices such as mobile phones, social media or internet platforms or apps with or without individuals'/owners' knowledge or consent, and/or could lead to individuals/owners being IDENTIFIED or SPECIAL CATEGORY DATA (GDPR) or CRIMINAL OFFENCE DATA.	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
<sup>X</sup> PUBLIC WORKS DOCTORATES: Evidence of permission is required for use of works/artifacts (that are protected by Intellectual Property (IP) rights, e.g. copyright, design right) in a doctoral critical commentary when the IP in the work/artifact is jointly prepared/produced or is owned by another body	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
18. <sup>H</sup> RISK OF PHYSICAL OR PSYCHOLOGICAL HARM (e.g., TRAVEL to dangerous places in your own country or in a foreign country (see <a href="https://www.gov.uk/foreign-travel-advice">https://www.gov.uk/foreign-travel-advice</a> ), research with NGOs/humanitarian groups in conflict/dangerous zones, development of technology/agent/chemical that may be harmful to others, any other foreseeable dangerous risks).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
19. <sup>X</sup> SECURITY CLEARANCE – required for research.	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No
20. <sup>X</sup> SENSITIVE TOPICS (e.g., anything deeply personal and distressing, taboo, intrusive, stigmatising, sexual in nature, potentially dangerous, etc).	<input type="checkbox"/> Yes	<input checked="" type="checkbox"/> No

M – Minimal Risk;      X – More than Minimal Risk.      H – High Risk

If you have answered 'Yes' to ANY of the above questions, your application **REQUIRES** ethical review and approval using the MOREform **BEFORE commencing your research**. Please apply for approval using the MOREform (<https://moreform.mdx.ac.uk/>). Further guidance on making an application using the MOREform can be found at: [www.tiny.cc/mdx-ethics](http://www.tiny.cc/mdx-ethics).

If you have answered 'No' to ALL of the above questions, your application is Low Risk and you may NOT require ethical review and approval using the MOREform before commencing your research. Your research supervisor will confirm this below.

Student Signature:.....Dhruvik Hasmukhbhai Patel..... Date:.....07/10/2023.....

### To be completed by the supervisor:

Based on the details provided in the self-assessment form, I confirm that:	Insert Y or N
The study is Low Risk and <i>does not require</i> ethical review & approval using the MOREform	
The study <i>requires</i> ethical review and approval using the MOREform.	

Supervisor Signature:..... Date:.....