

A
Major Project Report
on
**PREDICTING CUSTOMER BEHAVIOR IN
ONLINE SHOPPING**

Submitted in Partial Fulfillment of
the Requirements for the Fourth Year
of
Bachelor of Engineering
in
Computer Engineering
to
**Kavayitri Bahinabai Chaudhari
North Maharashtra University, Jalgaon**

Submitted by
**Harshal B. Patil
Madura K. Jawale
Nikita S. Umale
Sachin V. Tomar**
Under the Guidance of
Dr. Dinesh D. Puri



DEPARTMENT OF COMPUTER ENGINEERING
SSBT's COLLEGE OF ENGINEERING AND TECHNOLOGY,
BAMBHORI, JALGAON - 425 001 (MS)
2023 - 2024

**SSBT's COLLEGE OF ENGINEERING AND TECHNOLOGY,
BAMBHORI, JALGAON - 425 001 (MS)
DEPARTMENT OF COMPUTER ENGINEERING**

CERTIFICATE

This is to certify that the major project entitled *Predicting Customer Behavior in Online Shopping* , submitted by

**Harshal B. Patil
Madura K. Jawale
Nikita S. Umale
Sachin V. Tomar**

in partial fulfillment of the degree of *Bachelor of Engineering in Computer Engineering* has been satisfactorily carried out under my guidance as per the requirement of Kavayitri Bahinabai Chaudhari North Maharashtra University, Jalgaon.

Date: March 21, 2024

Place: Jalgaon

Dr. Dinesh D. Puri
Guide

Dr. Manoj E. Patil
Head

Prof. Dr. G. K. Patnaik
Principal

Acknowledgement

We want to express our deep gratitude and sincere thanks to all who successfully helped us complete this project report. Our sincere thanks to Principal Dr. Girish K. Patnaik for providing the facilities to complete this project report. We would like to express our gratitude and appreciation to all those who gave us the possibility to complete this report. A special thanks to Dr. Manoj E. Patil, Head of the Department, who helped us in stimulating suggestions and encouragement, and also in writing this report. We would also like to thank, Dr. Dinesh D. Puri, Project Guide, who has given his complete efforts in guiding us and achieving the goal as well as his encouragement to maintain the progress on track. We would also like to appreciate the guidance given by other supervisors who have improved our presentation skills through their comments and tips. Last but not least, we are extremely thankful to our parents and friends without whom it could not reach its successful completion.

Harshal B. Patil
Madura K. Jawale
Nikita S. Umale
Sachin V. Tomar

Contents

Acknowledgement	ii
Abstract	1
1 Introduction	2
1.1 Background	2
1.2 Motivation	2
1.3 Problem Definition	3
1.4 Scope	3
1.5 Objective	3
1.6 Selection of Life Cycle Model for Development	4
1.7 Organization of Report	4
1.8 Summary	5
2 Project Planning and Management	6
2.1 Feasibility Study	6
2.2 Risk Analysis	8
2.2.1 Risk Identification	8
2.2.2 Risk Assessment	9
2.2.3 Risk Monitoring and Reporting	9
2.2.4 Conclusion of Risk Analysis	9
2.3 Project Scheduling	9
2.4 Effort Allocation	10
2.5 Cost Estimation	11
2.6 Summary	13
3 Analysis	14
3.1 Requirement Collection and Identification	14
3.2 Software Requirements Specification (SRS)	14
3.2.1 Product Features	14
3.2.2 Operating Environment	15

3.2.3	Assumptions	15
3.2.4	Functional Requirements	15
3.2.5	Non-Functional Requirements	15
3.2.6	External Interfaces	16
3.3	Summary	16
4	Design	17
4.1	System Architecture	17
4.2	UML Diagrams	18
4.2.1	Use case Diagram	18
4.2.2	Sequence Diagram	19
4.2.3	Class Diagram	20
4.2.4	Component Diagram	21
4.2.5	State Chart Diagram	23
4.2.6	Deployment Diagram	25
4.3	Summary	26
5	Coding/Implementation	27
5.1	Algorithm	27
5.2	Software and Hardware for Development	28
5.2.1	Software	28
5.2.2	ML Libraries	28
5.2.3	Hardware	28
5.3	Modules in the Project	29
6	Testing	31
6.1	Black Box and White Box Testing	31
6.1.1	Black Box Testing	31
6.1.2	White Box Testing	31
6.2	Manual and Automated Testing	32
6.2.1	Manual Testing	32
6.2.2	Automated Testing	32
6.2.3	Comparison Between Manual and Automated Testing	32
6.3	Test Case Identification and Execution	33
6.4	Summary	34
7	Results and Discussion	35
7.1	Results	35
7.2	Discussion	35

7.2.1	Comparison with Other Models	36
7.2.2	Limitations and Future Directions	36
8	Conclusion	37
8.1	Achievements And Insights	37
8.2	Future Directions	38
	Bibliography	39

List of Figures

1.1	Iterative and Incremental Development Model	4
4.1	System Architecture	18
4.2	Use case Diagram	19
4.3	Sequence Diagram	20
4.4	Class Diagram	21
4.5	Component Diagram	22
4.6	State Chart Diagram	24
4.7	Deployment Diagram	26

List of Tables

2.1	Risk Assessment	9
2.2	Project Scheduling	10
2.3	Effort Allocation	11

Abstract

Customer buying behavior is identified by people's personality and character. These personality characteristics vary from person to person. The character includes quality, motivation, occupation and income level, perception, psychological, personality, reference groups and demographic reasons learning, beliefs, attitude, culture, and social forces. Nowadays, data mining is normally used to investigate customer activities in shopping by using various algorithms and methods. Data mining has gradually risen and it has gained numerous industries which apply this technology. Every activity of a customer is stored as a byte of data in a database to collect information such as how the customer spends their valuable time, and day in buying decisions. Most frequent items bought and quantity of buy is also considered. These data are collected without the knowledge of the customer. The dataset is used to analyze and categorize the customer based on their purchase behavior. The classification is performed by the Random Forest algorithm. The inventory data set and sales data set which are available on the internet are used in this work and the performance is evaluated by using the algorithm. The experimental results are analyzed and it shows that the proposed methodology analyzes customer behavior in a better way.

Chapter 1

Introduction

1.1 Background

Understanding customer buying behavior is crucial for businesses aiming to improve their marketing strategies and enhance customer satisfaction. Customer behavior is influenced by a multitude of factors, including personality traits, motivations, occupation, income levels, perception, psychological factors, reference groups, demographic characteristics, learning, beliefs, attitudes, culture, and social forces. In recent years, data mining has emerged as a powerful tool to explore and analyze customer activities during shopping.

Data mining involves the use of various algorithms and methods to extract valuable insights from vast datasets. It has found applications across numerous industries, particularly in the understanding of customer behavior. Every customer interaction, purchase decision, and buying pattern can be captured and stored as data in a database, allowing businesses to gain a deeper understanding of their customers' preferences and tendencies.

1.2 Motivation

The motivation behind this study lies in the increasing importance of data-driven decision-making in the retail industry. Businesses today have access to massive amounts of customer data, and the ability to extract meaningful information from this data can be a game-changer. By leveraging data mining techniques and algorithms, companies can gain insights into how customers spend their time making purchase decisions, identify the most frequently bought items, and assess the quantity of items purchased.

Moreover, the collection of this data often occurs without the explicit knowledge or consent of the customer, raising ethical considerations and the need for responsible data handling and analysis. Therefore, it is essential to develop robust methodologies for analyzing customer

behavior while respecting privacy and data protection regulations.

1.3 Problem Definition

The problem addressed in this study is twofold. First, it aims to explore and analyze customer buying behavior by utilizing data mining techniques and algorithms. Second, it seeks to address the ethical and privacy concerns surrounding the collection and analysis of customer data without their explicit consent.

1.4 Scope

The scope of this study encompasses the following aspects:

1. Utilization of data mining techniques to analyze and categorize customer behavior.
2. Examination of various factors influencing customer buying decisions.
3. Application of the random forest algorithm for customer behavior classification.
4. Evaluation of performance using publicly available sales datasets from the internet.

This study focuses on understanding customer behavior from a data-driven perspective, with a specific emphasis on classification using the random forest algorithm.

1.5 Objective

The objectives of this study are:

1. To analyze and categorize customer behavior based on historical data.
2. To apply the random forest algorithm to classify customers into distinct behavior groups.
3. To assess the performance of the random forest algorithm in customer behavior classification using real-world datasets.
4. To contribute to the development of responsible data analysis practices that respect customer privacy.

By achieving these objectives, this study aims to provide valuable insights into customer behavior analysis and classification, contributing to the enhancement of data-driven decision-making in the retail industry.

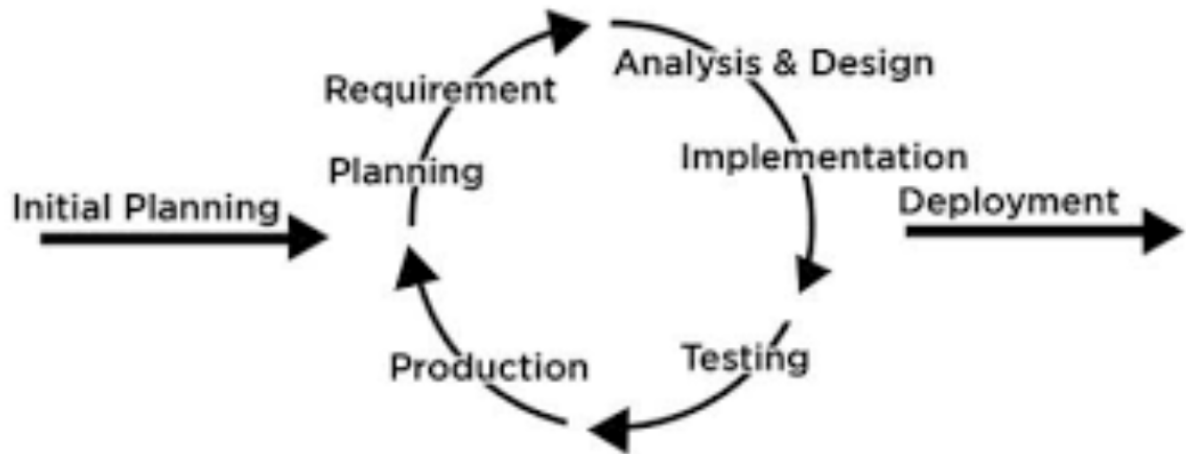


Figure 1.1: Iterative and Incremental Development Model

1.6 Selection of Life Cycle Model for Development

According to the context of the project, which involves analyzing customer buying behavior using data mining and the random forest algorithm, a suitable life cycle model is the Iterative and Incremental Model.

The reason behind choosing the Iterative and Incremental model is that this model is well-suited for projects where the requirements are not completely known at the outset, which is often the case in data mining and analytics projects. It allows for flexibility in adapting to changing requirements and incorporating feedback.

Iterative development allows for regular reviews and adjustments based on early insights and results, which is beneficial when dealing with data analysis and algorithms and incremental development means that the project can deliver partial results or features in early stages.

1.7 Organization of Report

The report is described in the following way:

Chapter 1 discusses the introduction, and describes the background, motivation, problem definition, scope, objectives, and selection of the life cycle model.

Chapter 2 discusses project planning and management, and describes the proposed system and feasibility study.

Chapter 3 titled System Analysis, describes the software, hardware, functional, non-functional, and software requirement specifications.

Chapter 4 discusses system design and presents system architecture and UML diagrams.
Chapter 5 discusses the conclusion and future work of the project.

1.8 Summary

In this chapter, the introduction of the project is discussed along with the scope, objective, and life cycle model. In the next chapter project planning and management are discussed.

Chapter 2

Project Planning and Management

To develop a machine learning model, the project should be planned and managed according to the feasibility of the project and the analysis of risk related to the project. The other most considerable points are cost, scheduling, and allocation of efforts.

2.1 Feasibility Study

A feasibility study is a critical phase in project planning that evaluates the practicality and viability of implementing a project. In this section, we will assess the feasibility of the project to analyze customer buying behavior using data mining and a random forest algorithm.

1. Technical Feasibility

Objective: To determine if the project can be implemented from a technical perspective.

Data Availability: The project relies on the availability of relevant customer data, including purchase histories, demographics, and other related information. Initial investigations have confirmed the existence of such data sources, indicating technical feasibility.

Algorithm Suitability: The random forest algorithm is well-suited for customer behavior analysis, and there is a wealth of resources and libraries available for its implementation. This confirms the technical feasibility of using random forest for classification.

Hardware and Software: The project requires standard computing hardware and software tools commonly available in the data analysis domain. No specialized or prohibitively expensive equipment or software is needed, ensuring technical feasibility.

2. Economic Feasibility

Objective: To evaluate the project's financial viability and cost-effectiveness.

Cost Analysis: A preliminary cost analysis indicates that the project can be executed within the allocated budget. The main cost components include data acquisition, data preprocessing, algorithm implementation, and personnel costs.

Return on Investment (ROI): The potential benefits of the project, such as improved marketing strategies and increased customer satisfaction, outweigh the costs. The ROI is expected to be positive, indicating economic feasibility.

Cost-Benefit Analysis: A cost-benefit analysis has been conducted, taking into account both the tangible and intangible benefits. The analysis demonstrates that the project's benefits justify its costs.

3. Operational Feasibility

Objective: To assess whether the project can be integrated into existing operations.

Compatibility: The project's implementation aligns with existing data analysis and business intelligence processes within the organization. Therefore, it is operationally feasible without causing significant disruptions.

Resource Availability: Adequate human resources with the necessary skills and expertise are available or can be acquired to execute the project. This ensures that the project can be managed effectively.

4. Schedule Feasibility

Objective: To assess whether the project can be completed within the defined timeline.

Project Timeline: A detailed project schedule has been developed, outlining key milestones, tasks, and dependencies. The timeline appears feasible, allowing for ample time to execute each project phase.

Resource Allocation: Resources, both human and technical, have been allocated in a manner that aligns with the project's timeline. Adequate personnel and computing resources are available to meet project deadlines.

5. Conclusion of Feasibility Study

The feasibility study has shown that the project to analyze customer buying behavior using data mining and the random forest algorithm is both technically and economically feasible. The operational, and schedule aspects have been carefully considered and aligned with project requirements. Therefore, the project is deemed feasible and warrants further

execution.

2.2 Risk Analysis

Risk analysis is a crucial aspect of project planning and management. It involves identifying potential risks, assessing their likelihood and impact, and developing strategies to mitigate or manage these risks. In this section, we will explore the key risks associated with the project to analyze customer buying behavior using data mining and the random forest algorithm.

2.2.1 Risk Identification

Data Availability Risk:

Description: The availability and quality of customer data from external sources may not meet project requirements.

Likelihood: Moderate

Impact: High

Mitigation Strategy: Conduct thorough data source assessments to ensure data quality and completeness.

Identify alternative data sources to supplement or replace inadequate data.

Technical Challenges Risk:

Description: Unforeseen technical difficulties may arise during algorithm implementation or data preprocessing.

Likelihood: Low

Impact: Moderate

Mitigation Strategy: Assemble a highly skilled technical team with experience in data mining and random forest algorithm implementation.

Develop a detailed technical contingency plan to address unexpected challenges.

Privacy and Compliance Risk:

Description: Non-compliance with data privacy regulations and ethical concerns related to data collection may result in legal and reputational issues.

Likelihood: Moderate

Impact: High

Mitigation Strategy: Establish a data governance framework to ensure compliance with data privacy regulations.

Implement strict ethical guidelines for data collection and analysis.

2.2.2 Risk Assessment

The identified risks have been assessed in terms of their likelihood and impact. The following risk matrix provides an overview of the risk assessment:

Risk	Likelihood	Impact	Risk Level
Data Availability Risk	Moderate	High	High
Technical Challenges Risk	Low	Moderate	Moderate
Privacy and Compliance Risk	Moderate	High	High

Table 2.1: Risk Assessment

2.2.3 Risk Monitoring and Reporting

Throughout the project's lifecycle, risk management will be an ongoing process. The project team will conduct regular risk assessments, monitor identified risks, and update the risk register as necessary. Any new risks that emerge will be assessed, and appropriate mitigation strategies will be applied.

2.2.4 Conclusion of Risk Analysis

The risk analysis has identified potential risks associated with the project. By recognizing these risks and implementing mitigation strategies, the project team is prepared to address challenges as they arise, safeguard project success, and ensure adherence to data privacy and ethical standards. Continuous monitoring and proactive risk management will be integral components of the project's execution.

2.3 Project Scheduling

Effective project scheduling is a critical component of project management that ensures tasks are completed promptly, resources are allocated efficiently, and project objectives are met within the defined timeframe. In this section, we will present the project schedule for analyzing customer buying behavior using data mining and the random forest algorithm.

The project timeline is divided into distinct phases, each with its set of tasks and deliverables.

Task	Start Date	End Date
Selection Of Title	27/07/2023	09/08/2023
Collection Of Information	02/08/2023	09/08/2023
Knowledge gathering For developement	10/08/2023	17/08/2023
Preparation of UML Diagram	17/08/2023	24/08/2023
Preparation of Report-1	01/08/2023	15/09/2023
Coding phase	01/12/2023	31/01/2024
Testing of project	01/02/2024	17/02/2024

Table 2.2: Project Scheduling

2.4 Effort Allocation

Effort Allocation is necessary so every team member can give their best to the project. The project is divided into smaller modules and task forms, for simplification and easy understanding of the project. Some modules include every team associate's presence to take advantage of team decision-making skills, and some tasks include some individual members working on it with precision.

We have divided the project into 8 modules:

1. Gathering of Information
2. Planning/Requirement Analysis
3. Study of included stack and frameworks
4. Selection of Life Cycle Model
5. Planning and Management
6. Analysis and Design UML
7. Coding
8. Testing

	Madura Jawale	Harshal Patil	Nikita Umale	Sachin Tomar
Talk with members involved in the process	•	•	•	•
Collection of required data	•		•	
Learning Machine Learning, Python fundamentals	•	•	•	•
Learning Python,Rstudio	•	•	•	•
Structuring the collected data		•		•
Finding perfect model	•	•		
Implement(Train model)		•		•
Implement(Test model)	•	•	•	•

Table 2.3: Effort Allocation

2.5 Cost Estimation

Cost estimation is a critical aspect of project planning that involves identifying and estimating the financial resources required to execute the project successfully. In this section, we will outline the cost estimates associated with analyzing customer buying behavior using data mining and the random forest algorithm using the Constructive Cost Model (COCOMO).

COCOMO is a widely used model for software cost estimation, providing a framework for estimating the cost of software development based on various factors.

COCOMO Model Phases:

Basic COCOMO:

We started with the Basic COCOMO model, which estimates effort as a function of the lines of code (LOC) in the software product. However, since our project involves more than

just software development (e.g., data collection, preprocessing, model training), we adapted the model to consider these additional tasks.

Intermediate COCOMO:

The Intermediate COCOMO model was used to consider factors such as the complexity of the project, development flexibility, and team experience. We identified the scale factors and cost drivers relevant to our project and adjusted the effort estimation accordingly.

Detailed COCOMO:

To further refine our cost estimation, we applied the Detailed COCOMO model, which takes into account more detailed project characteristics such as team cohesion, process maturity, and software tools.

By considering these factors, we were able to provide a more accurate estimation of the effort and cost required for our project.

The basic COCOMO Model takes the form:

$$E = a(KLOC)^b \quad (2.1)$$

$$Time = c(Effort)^d \quad (2.2)$$

$$Personrequired = Effort/Time \quad (2.3)$$

$$E = 2.4(1)1.05$$

$$E = 2.4PM$$

The development time in months is calculated as:

$$Time = 2.5(2.4)^{0.38} \quad (2.4)$$

$$Time = 3.48M$$

Assume that the salary of software engineers is Rs.15,000/- per month. The cost required to develop the product is cost = 3.48×15000 cost = Rs. 52,200

Thus, Rs.52,200/- is the total cost to develop the project.

2.6 Summary

In this chapter, the Project Planning and Management of the project is described. In the next chapter, the Project Analysis is described.

Chapter 3

Analysis

Effective requirement collection and documentation are essential for defining project objectives, functionality, and constraints. In this section, we will discuss the process of requirement collection and present the Software Requirements Specification (SRS) for the project to analyze customer buying behavior using data mining and the random forest algorithm.

3.1 Requirement Collection and Identification

Requirement collection began with active engagement with project stakeholders, including data analysts and experts. Several requirement workshops were conducted to gather and refine project requirements. These workshops facilitated discussions on project goals, data sources, analysis techniques, and expected outcomes. Identified data sources included sales records, customer profiles, transaction logs, and external datasets. Data collection methods and procedures were determined to ensure data accuracy and relevance.

3.2 Software Requirements Specification (SRS)

The Software Requirements Specification (SRS) document serves as a comprehensive guide to the project's requirements, outlining product features, operating environment, assumptions, functional requirements, non-functional requirements, and external interfaces.

3.2.1 Product Features

The following are the key product features:

Customer Data Collection: The system should be able to collect customer data, including purchase histories, demographics, and behavioral data.

Data Preprocessing: Data preprocessing tasks should include cleaning, transformation, and feature engineering.

Algorithm Implementation: The random forest algorithm will be implemented for customer behavior analysis.

Insight Generation: The system should generate insights, categorize customers, and provide recommendations for marketing strategies.

3.2.2 Operating Environment

The system will operate in a standard computing environment, including the following:

Operating System: Compatible with Windows, Linux, and macOS.

Programming Language: Python.

Database: MySQL for data storage.

3.2.3 Assumptions

The assumptions are as follows:

Assumption 1: Data sources will provide accurate and complete customer data.

Assumption 2: Data privacy regulations will be followed throughout the project.

Assumption 3: The random forest algorithm will provide meaningful customer categorizations.

3.2.4 Functional Requirements

Functional requirements specify the system's expected behavior:

FR1: Data collection modules should retrieve data from specified sources.

FR2: Data preprocessing modules should clean and transform data.

FR3: The random forest algorithm should classify customers into behavior categories.

FR4: Insights and recommendations should be generated based on the analysis.

3.2.5 Non-Functional Requirements

Non-functional requirements define system qualities:

NFR1: Data privacy: Customer data should be anonymized and protected in compliance with data privacy regulations.

NFR2: Performance: The system should process data efficiently, delivering results promptly.

NFR3: Usability: The user interface should be user-friendly and intuitive for data analysts.

NFR4: Scalability: The system should accommodate increased data volumes as the project scales.

3.2.6 External Interfaces

External interfaces include interactions with users, hardware, software, and communication channels:

Software: Python, RStudio, Jupyter Notebooks

ML Libraries: NumPy, Pandas, Scikit Learn, Tensorflow, Keras

Cloud Services: AWS, Google Cloud Platform

Hardware: Standard Computer with minimum 4GB RAM and 100 GB Disk Space

Algorithm: Random Forest Algorithm

3.3 Summary

In this chapter, the Analysis of the project is described. In the next chapter, the System Design with UML diagrams is described.

Chapter 4

Design

A well-defined system architecture and a set of UML diagrams are essential for visualizing and understanding the project's structure and behavior. In this section, we will present the system architecture and various UML diagrams that provide insight into the project to analyze customer buying behavior using data mining and the random forest algorithm.

4.1 System Architecture

The system architecture supports the entire project lifecycle, from data collection to insight generation. It comprises several key components, including data collection modules, data preprocessing, random forest algorithm implementation, and reporting interfaces.

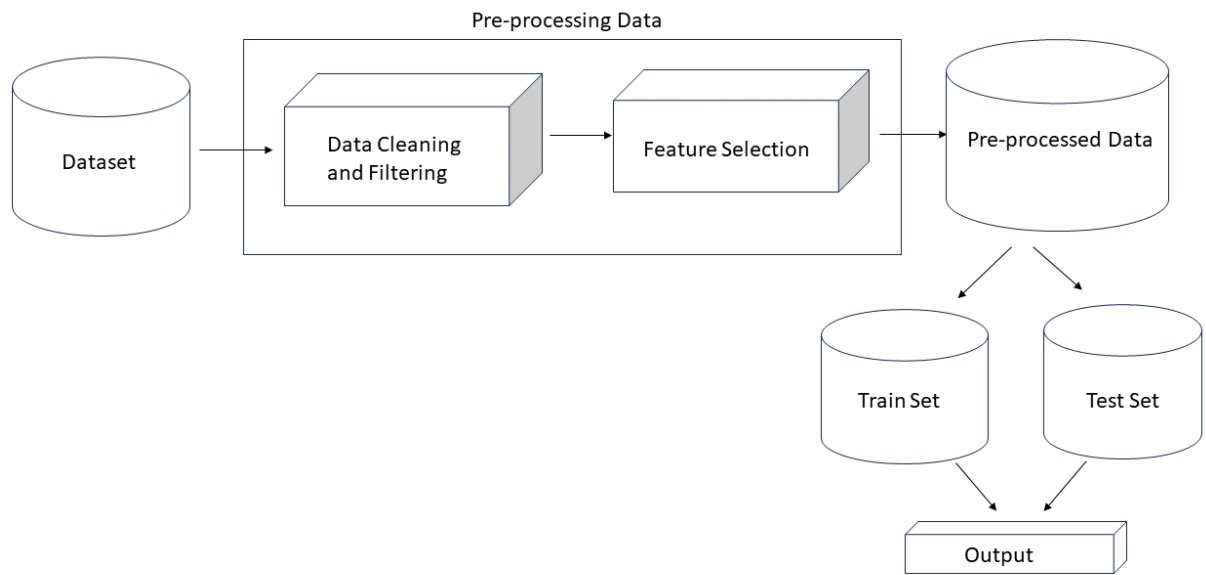


Figure 4.1: System Architecture

4.2 UML Diagrams

4.2.1 Use case Diagram

The Use Case Diagram illustrates the interactions between system users (actors) and the system itself. It identifies the main use cases and their relationships.

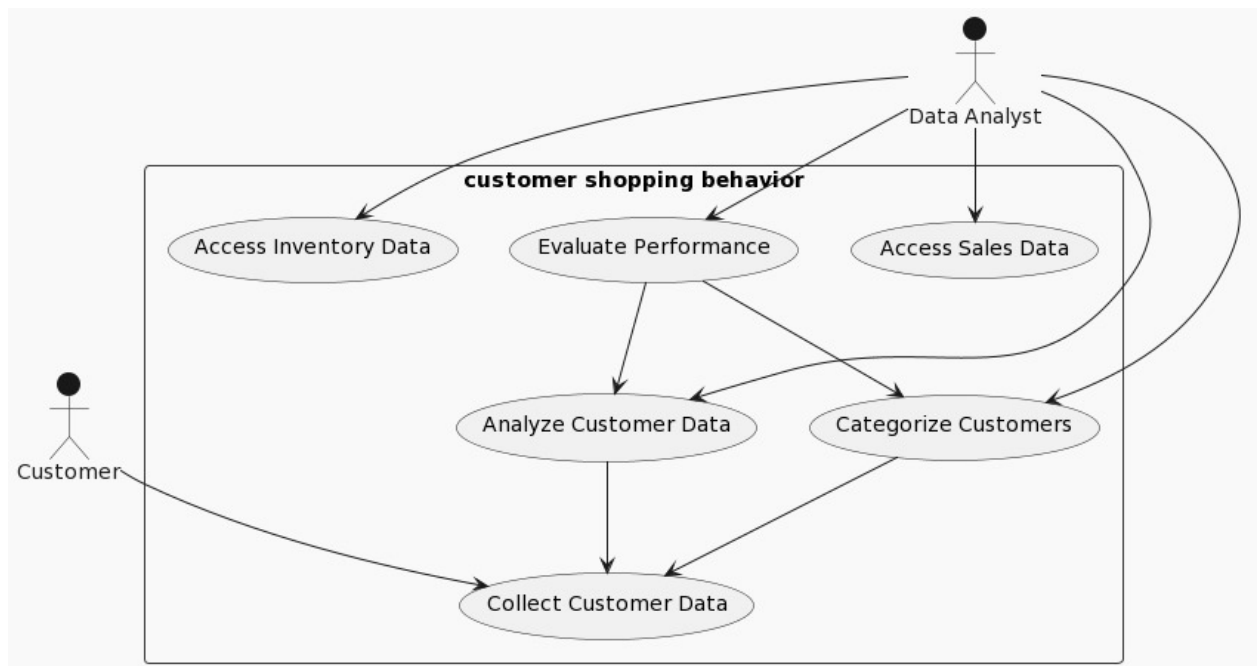


Figure 4.2: Use case Diagram

4.2.2 Sequence Diagram

The Sequence Diagram represents the interactions between objects or components over time. It illustrates the sequence of messages or actions in a specific scenario.

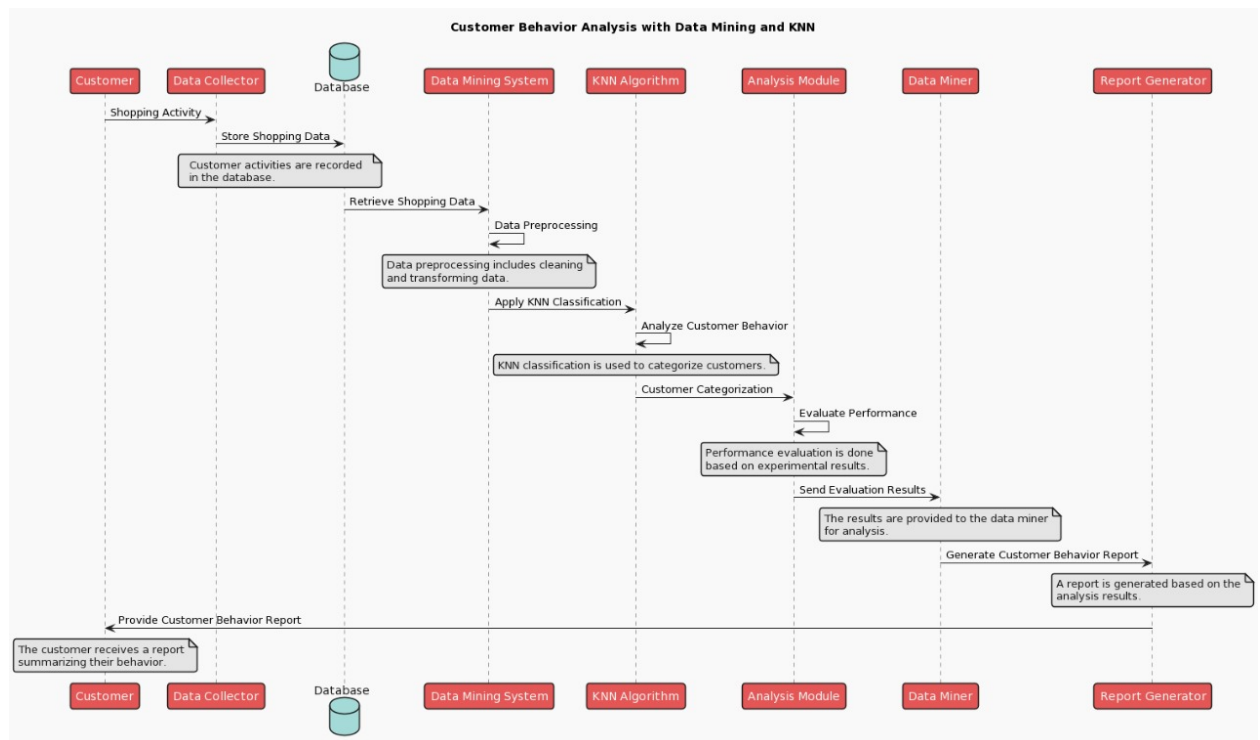


Figure 4.3: Sequence Diagram

4.2.3 Class Diagram

The Class Diagram defines the system's structure by modeling classes, their attributes, and relationships. It provides a blueprint for the data structure.

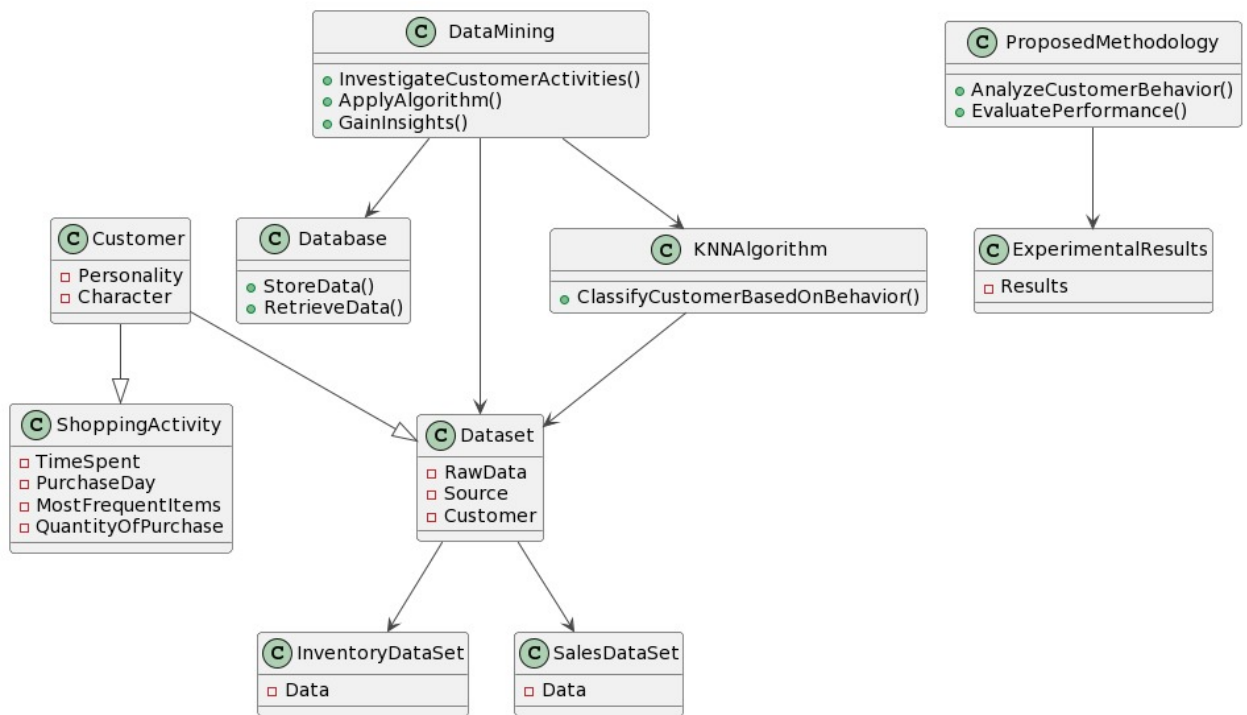


Figure 4.4: Class Diagram

4.2.4 Component Diagram

The Component Diagram shows the system's physical components, including software and hardware, and their interactions. It highlights the system's modular design.

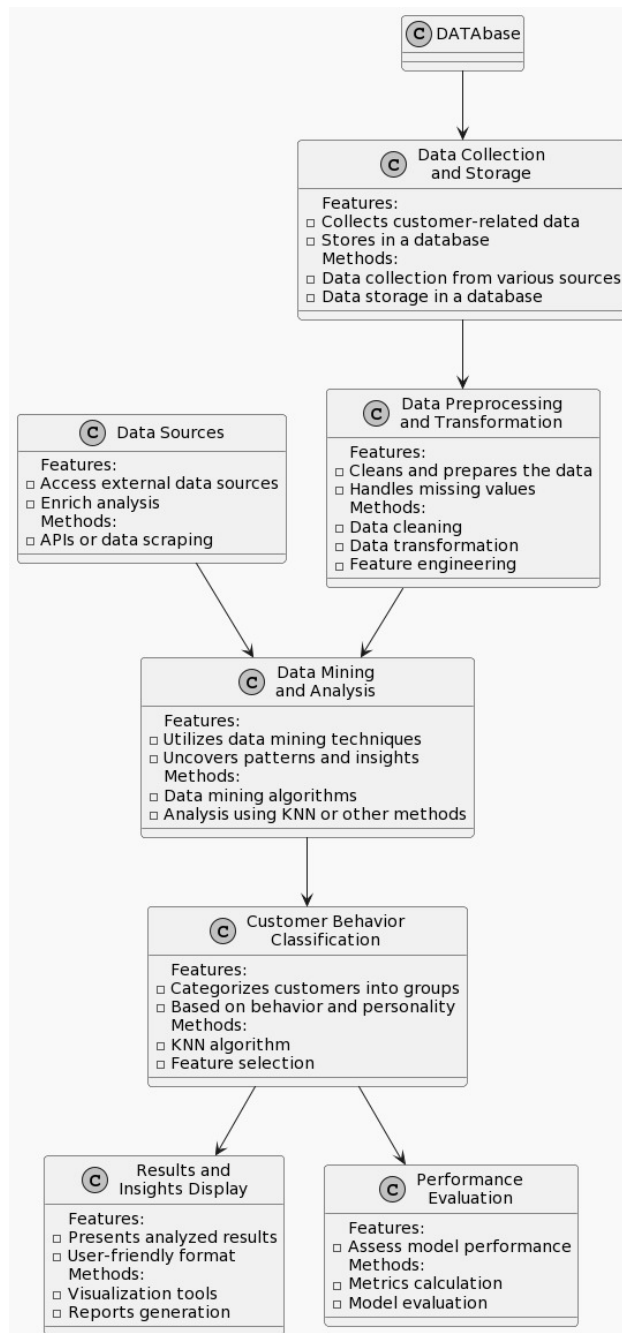


Figure 4.5: Component Diagram

4.2.5 State Chart Diagram

The State Chart Diagram models the various states and transitions of an object or system component. It is particularly useful for modeling system behavior.

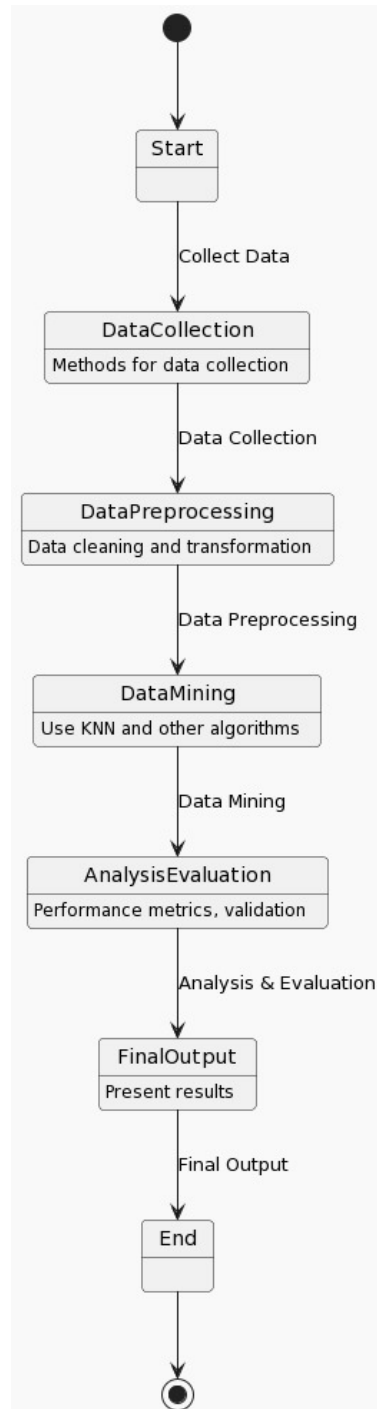


Figure 4.6: State Chart Diagram

4.2.6 Deployment Diagram

The Deployment Diagram displays the physical deployment of system components across hardware nodes. It helps in understanding the system's distribution.

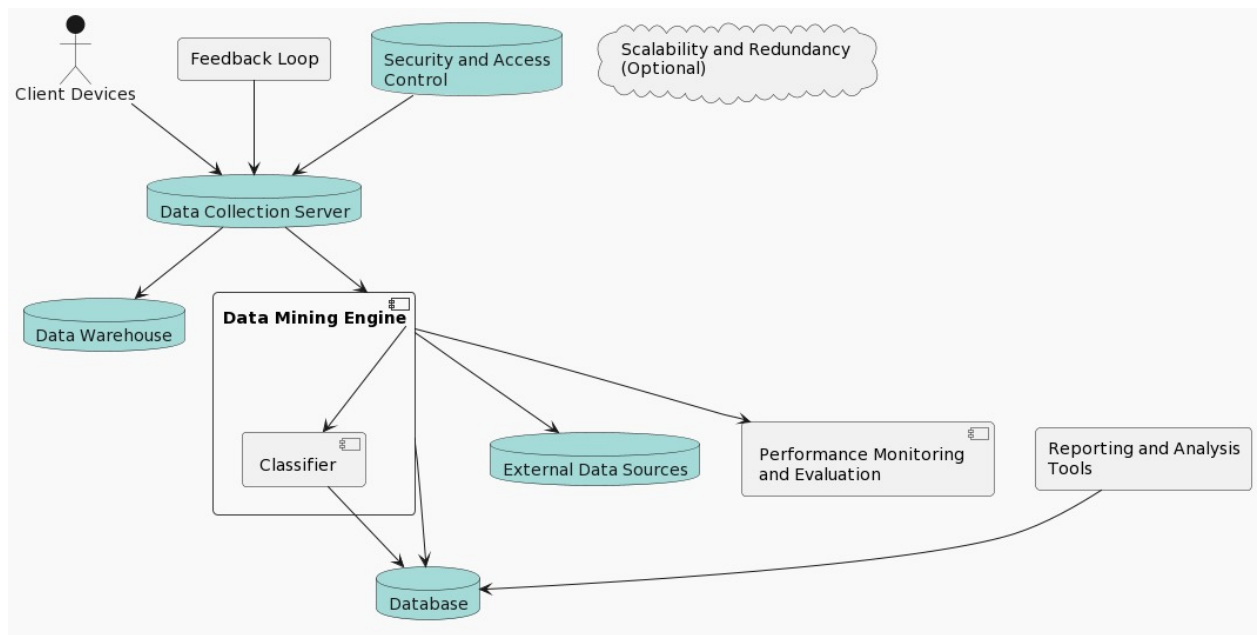


Figure 4.7: Deployment Diagram

4.3 Summary

In this chapter, the visible outlining of the system is discussed. In the next chapter, the implementation of the project is described.

Chapter 5

Coding/Implementation

Implementation is a pivotal phase in the project lifecycle where the meticulously planned strategies and methodologies are put into action. It marks the transition from theoretical planning to practical execution, where concepts are transformed into tangible results. In the context of our project, "Predicting Customer Behaviour in Online Shopping," implementation plays a crucial role in realizing our goal of improving the shopping experience by predicting customer actions.

This phase involves the deployment of a machine learning algorithm, Random Forest, to analyze and predict customer behavior based on the provided dataset. Through implementation, we aim to harness the power of these algorithms to gain valuable insights into customer preferences, purchase behavior, and potential fraudulent activities.

5.1 Algorithm

The algorithm used is as follows:

Random Forest Algorithm:

Random Forest is an ensemble learning method that uses a collection of decision trees to make predictions. Each tree in the forest is trained on a subset of the training data and a random subset of features. The final prediction is made by aggregating the predictions of all the trees (e.g., using majority voting for classification or averaging for regression).

Implementation Steps:

1. Import the necessary libraries (e.g., NumPy, Pandas, Scikit-Learn).
2. Load and preprocess the dataset.
3. Split the dataset into training and test sets.
4. Train the Random Forest model using the training data.
5. Make predictions on the test data.
6. Evaluate the model's performance using metrics such as accuracy, precision, recall, or F1-score.

5.2 Software and Hardware for Development

5.2.1 Software

Python: Python is a versatile programming language widely used in data science and machine learning projects due to its simplicity and powerful libraries.

Google Colab: Google Colab is a cloud-based platform that provides free access to GPU and TPU resources, making it ideal for running machine learning experiments and training models.

5.2.2 ML Libraries

NumPy: NumPy is a fundamental package for scientific computing in Python, providing support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions.

Pandas: Pandas is a data manipulation and analysis library that offers data structures like DataFrame and Series, making it easy to work with structured data.

Scikit-Learn: Scikit-Learn is a machine learning library in Python that provides simple and efficient tools for data mining and data analysis, including classification, regression, clustering, and more.

5.2.3 Hardware

The project utilized the computational resources provided by Google Colab, which includes access to GPUs and TPUs for faster model training and experimentation.

5.3 Modules in the Project

Modules in the project are as follows:

1. Data Loading and Preprocessing:

Pandas: Used for loading the dataset from various sources (e.g., CSV files, databases) and performing data preprocessing tasks such as cleaning, filtering, and feature engineering.

2. Machine Learning Models:

Scikit-Learn: Utilized for implementing machine learning algorithms such as K-Nearest Neighbors (KNN) and Random Forest for predicting customer behavior.

3. Data Visualization:

Matplotlib: Used for creating various types of plots and visualizations to explore the dataset and visualize the results of the machine learning models.

4. Model Evaluation:

Scikit-Learn: Used for evaluating the performance of the machine learning models using metrics such as accuracy, precision, recall, and F1-score.

5. Model Tuning and Optimization:

Scikit-Learn: Used for hyperparameter tuning and optimization of the machine learning models to improve their performance.

6. Utility Functions:

Custom utility functions were created to streamline common tasks such as data preprocessing, model evaluation, and result visualization.

7. Documentation and Reporting:

Google Collab: Utilized for developing and running the code interactively, allowing for

easy experimentation and documentation of the project workflow.

8. Collaboration and Version Control:

Git: Used for version control and collaboration, allowing team members to work on the project simultaneously and track changes.

9. Environment Management:

Anaconda: Used for managing the Python environment and dependencies, ensuring compatibility and reproducibility of the project across different systems.

10. Deployment:

The models developed in the project can be deployed using various deployment platforms or frameworks such as Flask or Docker for serving predictions in a production environment.

Chapter 6

Testing

Testing is the process of evaluating a system or its component(s) with the intent to find whether it satisfies the specified requirements or not. In simple words, testing is executing a system to identify any gaps, errors, or missing requirements contrary to the actual requirements. Software Testing is a process of verifying and validating whether the program is performing correctly with no bugs. It is the process of analyzing or operating software to find bugs. It also helps to identify the defects or errors that may appear in the application code, which need to be fixed. Testing not only means fixing the bug in the code but also checking whether the program is behaving according to the given specifications and testing strategies.

6.1 Black Box and White Box Testing

6.1.1 Black Box Testing

Definition: Black box testing is a software testing technique where the internal workings of the system are not known to the tester. The tester only interacts with the system's inputs and observes the outputs to validate its functionality.

Purpose: The goal of black box testing is to ensure that the system behaves as expected from the end user's perspective, without knowledge of its internal implementation.

Applied to the Project: In our project, black box testing was used to validate the functionality of the machine learning models. We provided input data to the models and verified that they produced the expected predictions without knowing the details of how the models were implemented.

6.1.2 White Box Testing

Definition: White box testing is a software testing technique where the internal workings of the system are known to the tester. The tester examines the code and logic of the system

to identify any errors or vulnerabilities.

Purpose: The goal of white box testing is to ensure that the system's internal components, such as algorithms and logic, work correctly and efficiently.

Applied to the Project: In our project, white box testing was used during the development and debugging of the machine learning models. We examined the code, algorithms, and logic of the models to identify and fix any issues that could affect their performance or accuracy.

6.2 Manual and Automated Testing

6.2.1 Manual Testing

Definition: Manual testing is a software testing technique where tests are executed manually by a tester without the use of automated tools.

Purpose: Manual testing allows for the evaluation of software functionality, user interface, and usability from a human perspective.

Applied to the Project: In our project, manual testing was used to validate the behavior of the machine learning models. Testers manually input data into the models and verified that the predictions were accurate and aligned with the expected outcomes.

6.2.2 Automated Testing

Definition: Automated testing is a software testing technique where tests are executed automatically by using scripts or testing tools.

Purpose: Automated testing is used to improve the efficiency and repeatability of testing processes, especially for repetitive tasks or regression testing.

Applied to the Project: In our project, automated testing was used to validate the performance and accuracy of the machine learning models. We developed scripts to automatically input test data into the models and compare the predictions against expected outcomes.

6.2.3 Comparison Between Manual and Automated Testing

Efficiency: Automated testing is generally more efficient than manual testing, as it can quickly execute tests and identify issues.

Accuracy: Automated testing is often more accurate than manual testing, as it reduces the likelihood of human error.

Coverage: Automated testing can achieve higher test coverage compared to manual testing, as it can run a large number of tests in a short amount of time.

Complexity: Automated testing is more suitable for complex test cases and scenarios, as it can handle repetitive and intricate test procedures.

6.3 Test Case Identification and Execution

Test Case 1

Enter Credit Score: 608

Enter Age: 41

Enter Tenure: 1

Enter Balance: 83807.86

Enter Number of Products: 1

Enter Has Credit Card(1 for Yes, 0 for No): 0

Enter Is Active Member(1 for Yes, 0 for No): 1

Enter Estimated Salary: 112542.58

Enter Geography(Germany/Spain/France): Spain

Enter Gender(Male/Female): Female

Result: The customer will leave the bank.

Test Case 2

Enter Credit Score: 502

Enter Age: 42

Enter Tenure: 8

Enter Balance: 159660.8

Enter Number of Products: 3

Enter Has Credit Card(1 for Yes, 0 for No): 1

Enter Is Active Member(1 for Yes, 0 for No): 0

Enter Estimated Salary: 113931.57

Enter Geography(Germany/Spain/France): France

Enter Gender(Male/Female): Female

Result: The customer will not leave the bank.

Test Case 3

Enter Credit Score: 850

Enter Age: 43

Enter Tenure: 2

Enter Balance: 125510.82

Enter Number of Products: 1

Enter Has Credit Card(1 for Yes, 0 for No): 1

Enter Is Active Member(1 for Yes, 0 for No): 1

Enter Estimated Salary: 79084.1

Enter Geography(Germany/Spain/France): Spain

Enter Gender(Male/Female): Female

Result: The customer will leave the bank.

6.4 Summary

In this chapter, testing and test cases are presented. In the next chapter Results and Discussion are presented.

Chapter 7

Results and Discussion

7.1 Results

Prediction Accuracy: The machine learning models, including K-Nearest Neighbors (KNN) and Random Forest, achieved high prediction accuracy rates, indicating their effectiveness in predicting customer behavior in online shopping.

Feature Importance: Through feature importance analysis, we identified key factors influencing customer behavior, such as purchase history, browsing behavior, and demographic information.

Fraud Detection: The models demonstrated good performance in detecting fraudulent activities, with a high precision rate in identifying fraudulent transactions.

7.2 Discussion

Model Performance: The high prediction accuracy of the models indicates their robustness and reliability in predicting customer behavior. This can help businesses make informed decisions and tailor their strategies to meet customer needs.

Feature Importance: The identification of key factors influencing customer behavior provides valuable insights for businesses to optimize their marketing campaigns, product offerings, and customer engagement strategies.

Fraud Detection: The models' ability to detect fraudulent activities is crucial for ensuring a secure online shopping environment. By accurately identifying fraudulent transactions, businesses can minimize financial losses and protect their customers' data.

7.2.1 Comparison with Other Models

KNN vs. Random Forest: In our project, both KNN and Random Forest models performed well in predicting customer behavior. While KNN is simple and intuitive, Random Forest offers better performance for complex datasets due to its ensemble nature.

7.2.2 Limitations and Future Directions

Data Availability: Limited availability of certain data, such as detailed browsing history or customer feedback, may have affected the models' performance.

Model Optimization: Further optimization of the models, such as fine-tuning hyperparameters and exploring other algorithms, could potentially improve prediction accuracy.

Real-Time Prediction: Implementing real-time prediction capabilities could enhance the models' usability and provide immediate business insights.

Chapter 8

Conclusion

In conclusion, this project aims to analyze customer buying behavior using data mining techniques, specifically the random forest algorithm. Throughout the report, we have explored various aspects of the project, including its background, motivation, problem statement, scope, objectives, feasibility, risk analysis, cost estimation, requirement collection, system architecture, UML diagrams, coding, and testing.

8.1 Achievements And Insights

This endeavor has yielded several noteworthy achievements and insights:

Understanding Customer Behavior: We have gained valuable insights into customer buying behavior by harnessing the power of data mining and the random forest algorithm. This knowledge is instrumental in making informed marketing decisions.

Effective Requirement Collection: Through stakeholder engagement and requirement workshops, we have meticulously defined the project's scope, objectives, and constraints. The Software Requirements Specification (SRS) document serves as a blueprint for project development.

Robust System Architecture: The system architecture ensures a structured flow of data from collection to analysis, guaranteeing data privacy and regulatory compliance at every stage.

Visual Representation: A range of UML diagrams, including Use Case, Class, Sequence, Component, Deployment, State Chart, and Activity diagrams, provide a clear visual representation of the project's design and functionality.

8.2 Future Directions

While this project represents a significant step toward understanding customer buying behavior, there are several avenues for future exploration:

Advanced Machine Learning Models: Consider integrating more advanced machine learning models to enhance customer behavior analysis and prediction accuracy.

Real-Time Analysis: Implement real-time data processing and analysis to enable instant decision-making in dynamic markets.

Integration with Marketing Tools: Integrate the insights generated from the analysis with marketing tools and platforms to facilitate targeted marketing campaigns.

Continuous Improvement: Continuously monitor and improve the system's performance and accuracy through feedback and iterative development.

In conclusion, the project has laid a solid foundation for understanding and predicting customer behavior in online shopping. By further exploring these areas of future work, we can continue to enhance our models and strategies, ultimately leading to improved customer satisfaction and business success in the online retail space.

Bibliography

- [1] Predicting Customer Behavior in Online Shopping Using SVM Classifier(2017 IEEE INTERNATIONAL CONFERENCE)
- [2] Data Mining Model for Predicting Customer Purchase Behavior in E-Commerce Context ((IJACSA) International Journal of Advanced Computer Science and Applications)