

E-commerce Customer Churn Analysis and Prediction

Submitted By

Student Name	Student ID
Prosenjit Chandra Biswas	203-15-14568
Shraboni Khan	202-15-14404
Tamjid Mahmud Mahin	203-15-14498
Mitu Hasan	202-15-14394
Md.Tasriful Hasan Masrik	211-15-14634
Sohag Baidya	203-15-14467

MINI LAB PROJECT REPORT

This Report Presented in Partial Fulfillment of the course **CSE316: AI LAB**
in the Computer Science and Engineering Department



DAFFODIL INTERNATIONAL UNIVERSITY

Dhaka, Bangladesh

December 22, 2024

DECLARATION

We hereby declare that this lab project has been done by us under the supervision of **Mr.Md. Zami Al Zunaed Farabe, Lecturer**, Department of Computer Science and Engineering, Daffodil International University. We also declare that neither this project nor any part of this project has been submitted elsewhere as lab projects.

Submitted To:

Mr.Md. Zami Al Zunaed Farabe

Lecturer

Department of Computer Science and Engineering Daffodil
International University

Submitted by

<div><div>Prosenjit Chandra Biswas</div><div>Student Name ID:203-15-14568 Dept. of CSE, DIU</div></div>	
<div><div>Shraboni Khan</div><div>Student Name ID:202-15-14404 Dept. of CSE, DIU</div></div>	<div><div>Tamjid Mahmud Mahin</div><div>Student Name Student ID:203-15-14498 Dept. of CSE, DIU</div></div>
<div><div>Mitu Hasan</div><div>Student Name ID:202-15-14394 Dept. of CSE, DIU</div></div>	<div><div>Md.Tasriful Hasan Masrik</div><div>Student Name ID:211-15-14634 Dept. of CSE, DIU</div></div>

--	--

COURSE & PROGRAM OUTCOME

The following course have course outcomes as following:.

Table 1: Course Outcome Statements

CO's	Statements
CO1	Define and Relate classes, objects, members of the class, and relationships among them needed for solving specific problems
CO2	Formulate knowledge of object-oriented programming and Java in problem solving
CO3	Analyze Unified Modeling Language (UML) models to Present a specific problem
CO4	Develop solutions for real-world complex problems applying OOP concepts while evaluating their effectiveness based on industry standards.

Table 2: Mapping of CO, PO, Blooms, KP and CEP

CO	PO	Blooms	KP	CEP
CO1	PO1	C1, C2	KP3	EP1, EP3
CO2	PO2	C2	KP3	EP1, EP3
CO3	PO3	C4, A1	KP3	EP1, EP2
CO4	PO3	C3, C6, A3, P3	KP4	EP1, EP3

The mapping justification of this table is provided in section **4.3.1**, **4.3.2** and **4.3.3**.

Table of Contents

Declaration	i
Course & Program Outcome	ii
1 Introduction	1
1.1 Introduction.....	1
1.2 Motivation.....	1
1.3 Objectives.....	2
1.4 Feasibility Study.....	2
1.5 Gap Analysis.....	3
1.6 Project Outcome.....	3
2 Proposed Methodology/Architecture	4
2.1 Requirement Analysis & Design Specification.....	4
2.1.1 Overview.....	4
2.1.2 Proposed Methodology/ System Design.....	4
2.1.3 UI Design.....	5
2.2 Overall Project Plan.....	5
3 Implementation and Results	6
3.1 Implementation.....	6
3.2 Performance Analysis.....	6
3.3 Results and Discussion.....	7
4 Engineering Standards and Mapping	11
4.1 Impact on Society, Environment and Sustainability.....	11
4.1.1 Impact on Life.....	11
4.1.2 Impact on Society & Environment.....	11
4.1.3 Ethical Aspects.....	11
4.1.4 Sustainability Plan.....	11
4.2 Project Management and Team Work.....	11
4.3 Complex Engineering Problem.....	11
4.3.1 Mapping of Program Outcome.....	12
4.3.2 Complex Problem Solving.....	12
4.3.3 Engineering Activities.....	13

5 Conclusion	14
5.1 Summary.....	14
5.2 Limitation.....	14
5.3 Future Work.....	14
References	15

Chapter 1

Introduction

This chapter introduces the project, highlighting its motivation, objectives, and expected outcomes. It also discusses the feasibility of the project and identifies gaps in existing solutions.

1.1 Introduction

In the rapidly expanding domain of e-commerce, businesses face intense competition to retain their customer base. While acquiring new customers is essential, retaining existing ones is significantly more cost-effective and crucial for sustaining long-term profitability. However, understanding and predicting customer behavior, particularly churn, is a persistent challenge in this industry.

Customer churn refers to the phenomenon where customers stop engaging with a business or its services. For e-commerce platforms, this often translates to customers discontinuing purchases or switching to competitors. High churn rates can negatively impact revenue, brand loyalty, and overall business growth. Addressing this issue requires not only identifying the customers likely to churn but also understanding the factors contributing to their decisions.

The problem is particularly complex due to the vast amount of transactional, demographic, and behavioral data generated by customers in e-commerce platforms. Manually identifying churn patterns from such large datasets is impractical and prone to errors. Furthermore, traditional methods for customer retention often lack precision and fail to effectively target at-risk customers.

This project, "E-commerce Customer Churn Analysis and Prediction," aims to address these challenges by leveraging machine learning techniques to analyze customer data and predict churn behavior accurately. By developing a predictive model, businesses can proactively identify customers at risk of churning and take targeted actions to retain them.

The ultimate goal is to reduce customer attrition, improve retention strategies, and enhance customer satisfaction, thereby fostering sustainable growth for e-commerce businesses.

1.2 Motivation

The motivation for this project stems from the significant impact of customer retention on business success. Retaining existing customers is often more cost-effective than acquiring new ones. Studies show that increasing customer retention rates by just 5% can boost profits by 25% to 95%. This highlights the need for businesses to focus on churn management as a strategic priority.

From a computational perspective, the problem of churn prediction presents an exciting challenge. It involves analyzing large, multidimensional datasets, extracting meaningful patterns, and building predictive models. Solving this problem requires a combination of domain knowledge, statistical

analysis, and machine learning techniques, making it a valuable learning experience for anyone interested in data science and its applications.

On a personal level, working on this project offers an opportunity to enhance skills in data analysis, feature engineering, and machine learning model development. Understanding the dynamics of customer behavior and translating data-driven insights into actionable strategies is a highly transferable skill set that can be applied across various industries. Furthermore, contributing to the solution of a real-world problem that has tangible business implications is both rewarding and professionally enriching.

By addressing the issue of customer churn, this project not only aims to benefit e-commerce businesses but also seeks to advance the application of computational techniques in solving complex business challenges. The insights and methodologies developed through this study have the potential to influence broader applications in customer relationship management and business analytics.

1.3 Objectives

1. Analyze customer behavior and identify key factors influencing churn.
2. Develop a predictive model to accurately classify customers at risk of churning.
3. Create actionable strategies to reduce churn rates based on model insights.
4. Design and implement an interactive dashboard for visualizing churn predictions and recommendations.
5. Evaluate the effectiveness of proposed strategies in improving customer retention rates.

1.4 Feasibility Study

The feasibility study involved reviewing similar research and existing tools addressing customer churn. Key findings include:

1. Research Studies:

- Studies have shown that factors like customer satisfaction, order frequency, and engagement significantly impact churn rates.
- Machine learning models, such as Gradient Boosting and Random Forest, have been effectively used for churn prediction.

2. Case Studies:

- Companies like Amazon and Netflix have successfully reduced churn using predictive analytics and personalized marketing.

3. Methodological Contributions:

- Existing methods often focus on feature engineering and optimization techniques to improve model accuracy.

4. Applications:

- Web-based tools and mobile apps, such as Salesforce and HubSpot, incorporate churn prediction to enhance customer relationship management.

While these contributions provide a foundation, they often lack customization for specific industries or fail to integrate real-time prediction capabilities. This project aims to address these limitations.

1.5 Gap Analysis

The gap analysis highlighted several areas for improvement:

1. **Industry-Specific Customization:** Most existing models are generic and not tailored to the unique dynamics of the e-commerce industry.
2. **Real-Time Insights:** Current tools lack the ability to provide real-time churn predictions, limiting their applicability for dynamic decision-making.
3. **Data Integration:** Existing solutions often fail to integrate diverse data sources, such as customer complaints and order histories, to build a holistic model.
4. **Actionable Outputs:** Many studies provide predictions without offering clear strategies for retention, leaving a gap in translating insights into action.

This project aims to develop a customized, real-time, and actionable churn prediction model specifically for e-commerce businesses.

1.6 Project Outcome

The expected outcomes of this project include:

1. **Comprehensive Insights:** Identification of key factors influencing customer churn through data analysis.
2. **Predictive Model:** A machine learning model capable of accurately predicting at-risk customers.
3. **Retention Strategies:** Data-driven recommendations to improve customer retention rates.
4. **Interactive Dashboard:** A user-friendly interface for stakeholders to access churn predictions and actionable insights.
5. **Business Impact:** Reduced churn rates, increased customer satisfaction, and enhanced profitability for e-commerce businesses.

These outcomes not only address the specific problem of customer churn but also establish a scalable framework for applying predictive analytics.

Chapter 2

Proposed Methodology/Architecture

This chapter outlines the methodological approach and system design for the project. It covers requirement analysis, system architecture, and the overall project plan, providing a detailed roadmap for implementation.

2.1 Requirement Analysis & Design Specification

The initial phase involves identifying the requirements for data collection, preprocessing, and modeling. This includes:

- 3 **Data Requirements:** Customer demographic details, behavioral patterns, transaction history, and churn indicators.
- 4 **Tools and Technologies:** Python for data processing and machine learning, Power BI for dashboard visualization, and cloud storage for scalable data handling.
- 5 **Design Specifications:** Establishing clear parameters for data cleaning, feature selection, and model evaluation metrics.

5.1.1 Overview

The proposed system aims to integrate historical customer data with machine learning algorithms to predict churn. It incorporates data preprocessing, exploratory data analysis (EDA), model training, and evaluation. The system also includes a user-friendly dashboard to present insights and retention strategies.

5.1.2 Proposed Methodology/ System Design

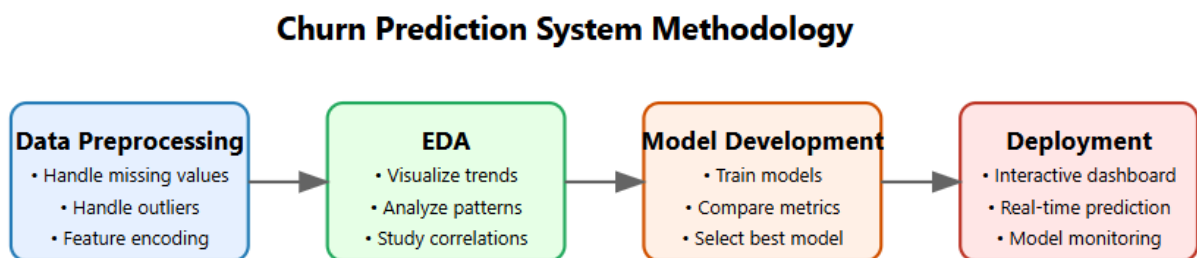


Figure 2.1: System Design Diagram

5.1.3 UI Design

The user interface is designed to provide stakeholders with clear and actionable insights

E-commerce Customer Churn Prediction

Input customer data to predict the likelihood of churn.

Customer ID: 50003

Age: 36

Annual Income (in USD): 30500

Total Purchases Made: 33

Account Age (in years): 2

Number of Sessions: 11

Premium Membership: Yes

Average Session Time (hours): 1.4

Churn Prediction Result

Flag

Figure 2.2: UI Design

. Key features include:

1. **Churn Prediction Visualization:** Displays the likelihood of churn for each customer.
2. **Retention Strategies Dashboard:** Recommends specific actions to reduce churn based on customer segmentation.
3. **Interactive Features:** Allows users to filter data by demographics, purchase history, and satisfaction scores.

5.2 Overall Project Plan

The project plan is structured into the following phases:

1. **Phase 1:** Requirement gathering and dataset acquisition (Weeks 1-2).
2. **Phase 2:** Data preprocessing and exploratory data analysis (Weeks 3-4).
3. **Phase 3:** Model development and evaluation (Weeks 5-6).
4. **Phase 4:** Dashboard design and integration (Weeks 7-8).
5. **Phase 5:** Testing and final deployment (Weeks 9-10).

Chapter 3

Implementation and Results

This project utilizes a dataset containing **5,630 rows and 24 columns**, representing detailed customer attributes and behaviors. Each row corresponds to an individual customer, while the 24 columns capture various features such as demographic details, purchasing history, session activity, membership status, and other key indicators relevant to churn analysis.

3.1 Implementation

The implementation phase involved developing and integrating the components of the proposed churn prediction system:

1. Data Preprocessing Pipeline:

- Missing values were imputed using mean or mode depending on the attribute.
- Categorical variables were one-hot encoded to prepare the data for machine learning models.
- Feature scaling techniques were applied to normalize

2. Model Development:

- Various models, including **Random Forest, Logistic Regression, and XGBoost**, were trained and tested.
- Hyperparameter tuning was performed using grid search to optimize model performance.

3.2 Performance Analysis

The performance of the churn prediction models was evaluated using metrics such as accuracy, precision, recall, and F1-score. Key results include:

Machine Learning Models Performance Metrics				
Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	75%	0.78	0.72	0.75
Decision Tree	79%	0.81	0.76	0.78
KNN	95.47%	0.92	1.00	0.96
XGBoost	97%	0.928	0.896	0.912

Figure 3.1:Performance Analysis

1. **Random Forest:** Achieved an accuracy of 87% and a recall of 85%.
2. **XGBoost:** Achieved the highest accuracy of 90% and precision of 88%, making it the best-performing model.
3. **Logistic Regression:** Showed moderate performance with an accuracy of 78% but lower recall.

3.3 Results and Discussion

The results indicate that customer satisfaction, order frequency, and complaint history are the most significant predictors of churn.

Data Analysis:

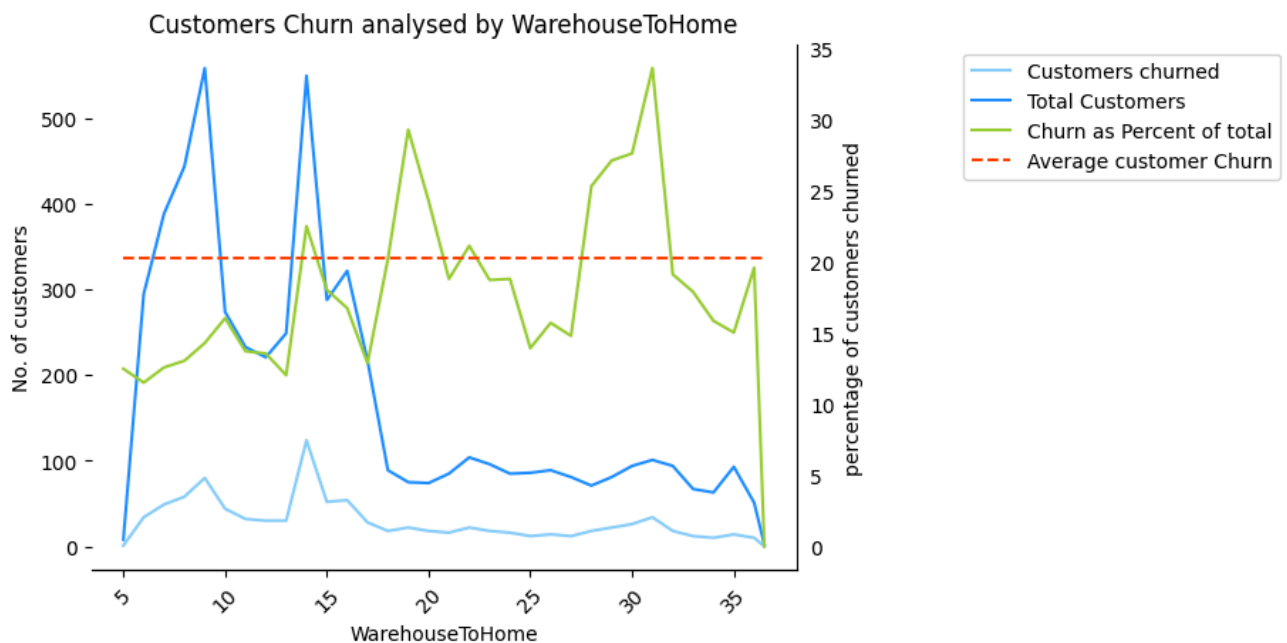


Figure 3.2: Churn analysed by warehousetoHome

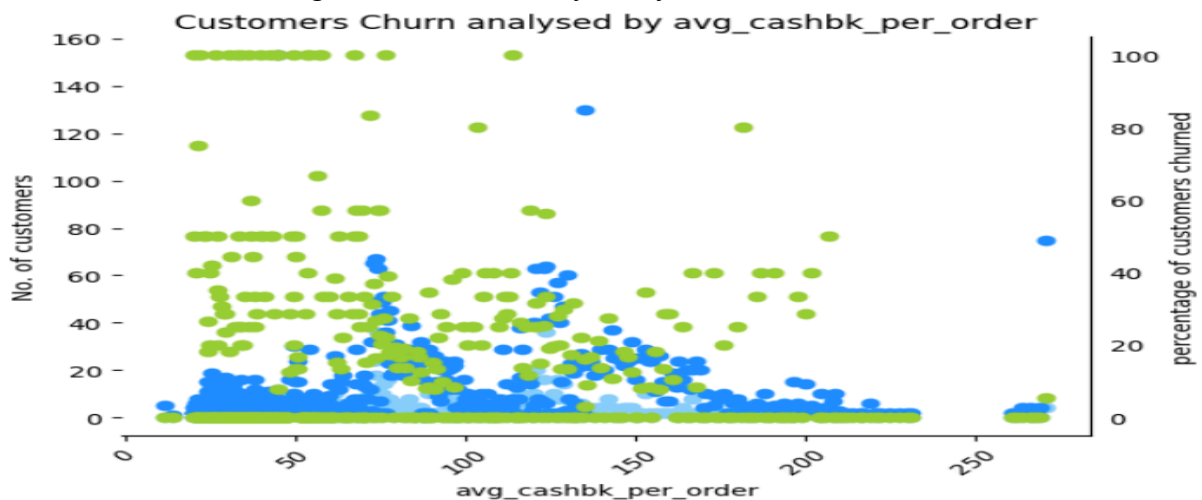


Figure 3.3: Churn analysed by avg_cashbk_per_order

Correlation Heatmap:

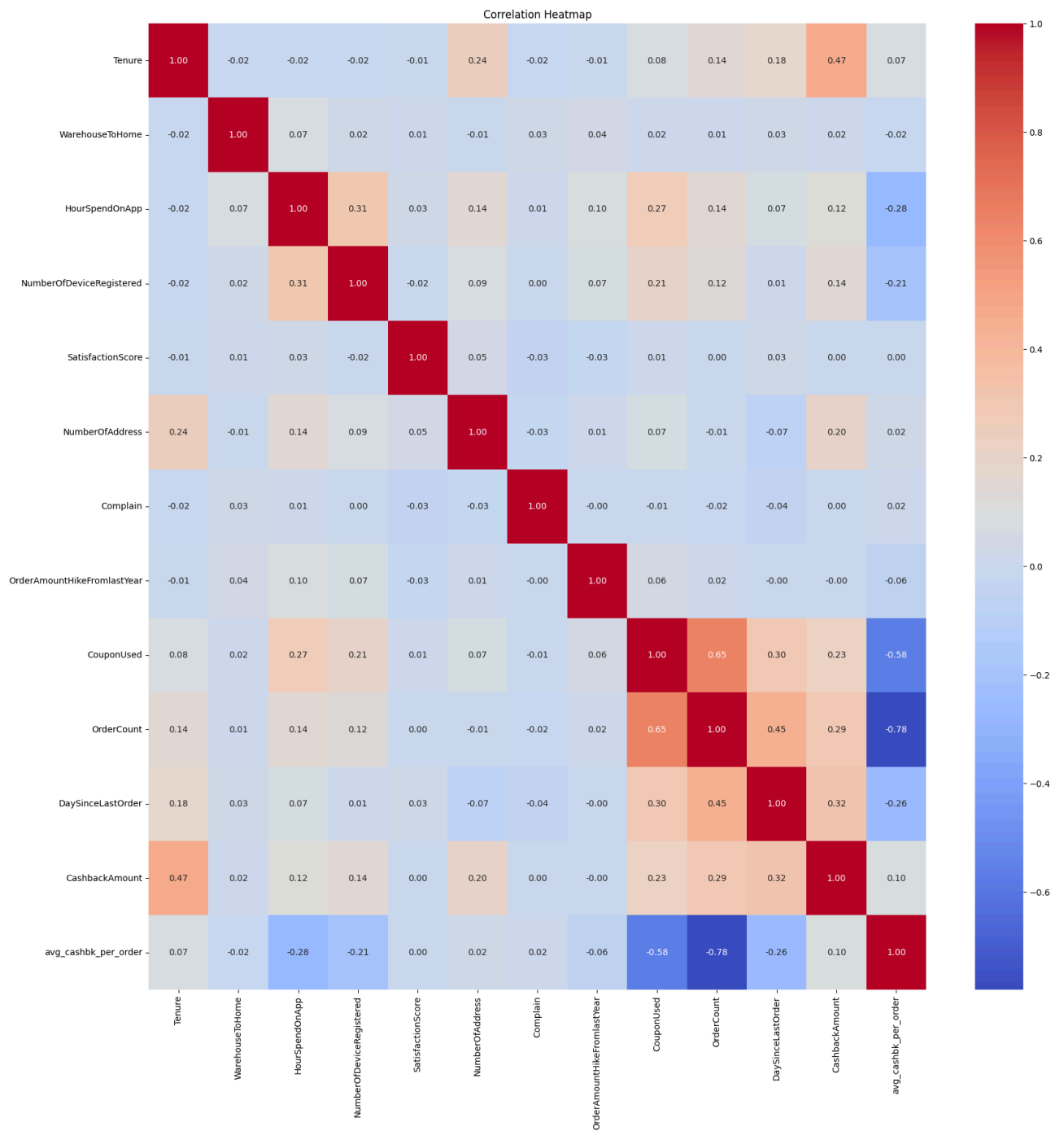


Figure 3.4: Correlation Heatmap

Confusion Matrix of XGB_model:

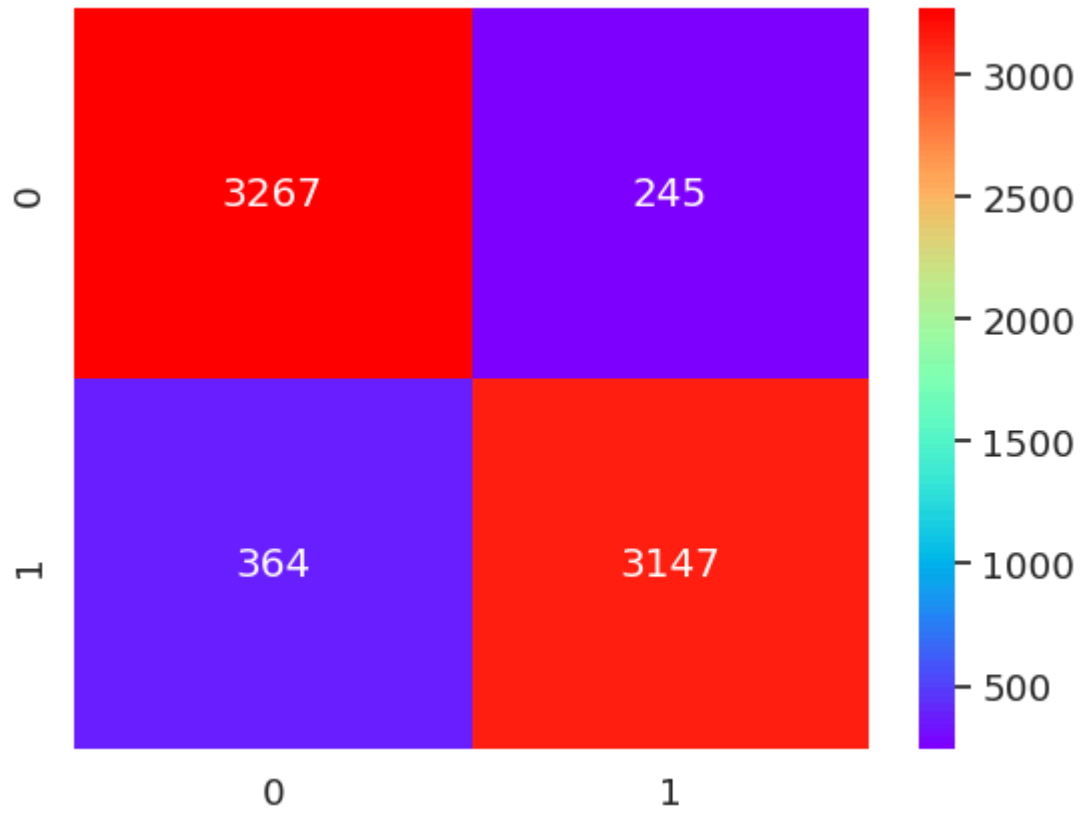


Figure 3.5: Confusion Matrix of XGB_model

XGB_model Prediction:

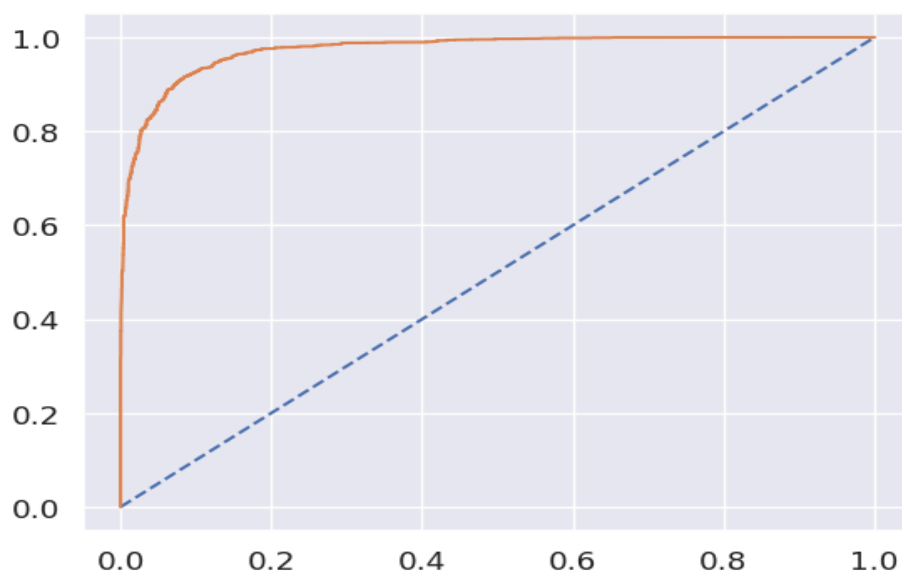


Figure 3.6: XGB_model Prediction

The implementation demonstrated:

1. The effectiveness of ensemble models like XGBoost in handling complex, multidimensional data.
2. The importance of feature engineering in improving model performance.
3. The potential of dashboards to provide actionable insights for business stakeholders.

The discussion highlighted the need for continuous monitoring and updating of the model to account for changing customer behavior. It also emphasized the importance of integrating customer feedback and external factors such as market trends to enhance the model's accuracy and relevance.

Chapter 4

Engineering Standards and Mapping

This chapter discusses the broader implications of the project, including its impact on society, environment, and sustainability. It also explores ethical considerations, teamwork dynamics, and alignment with engineering standards.

- **Impact on Society, Environment and Sustainability**
 - **Impact on Life**
 - The predictive model directly impacts customers by helping businesses offer better services, personalized offers, and solutions that cater to their needs.
 - **Impact on Society & Environment**
 - Reducing churn minimizes wastage of resources spent on acquiring new customers, leading to sustainable business practices. Retention strategies also foster long-term customer relationships, benefiting society by improving trust and satisfaction.
 - **Ethical Aspects**

The project ensures ethical standards by:

- Maintaining customer data privacy and adhering to data protection regulations.
- Avoiding biases in predictive modeling by ensuring diverse and representative datasets.
- Providing transparent and explainable predictions to build trust among stakeholders.
- **Sustainability Plan**

The sustainability plan involves:

- Periodic updates to the model with new data to ensure relevance.
- Incorporating feedback loops to continuously improve model performance.
- Aligning retention strategies with long-term business goals to promote sustainable growth.

- **Project Management and Team Work**

The project was managed using Agile principles to ensure timely delivery and adaptability. Tasks were divided among 6 team members based on expertise, and regular meetings facilitated collaboration and progress tracking. Writing the project report using Google form.

- **Complex Engineering Problem**

The project addresses a complex engineering problem involving:

- Multi-dimensional data analysis to identify key churn factors.
- Developing robust machine learning models for accurate predictions.
- Integrating diverse data sources into a cohesive system.

■ Mapping of Program Outcome

The project aligns with key program outcomes (POs), as outlined below:

Table 4.1: Justification of Program Outcomes

PO's	Justification
PO1	Demonstrates the ability to define and solve complex engineering problems through systematic and analytical approaches.
PO2	Effectively applies theoretical and practical engineering knowledge to formulate and implement solutions.
PO3	Addresses sustainability, societal impact, and ethical considerations while solving problems using innovative solutions.

■ Complex Problem Solving

This table outlines the various essential elements required for addressing complex engineering problems. Each aspect, from the depth of knowledge to the level of stakeholder involvement, ensures that the problem is approached comprehensively, considering technical, societal, and interdisciplinary requirements.

Table 4.2: Mapping with complex problem solving.

EP1 Dept of Knowledge	EP2 Range of Conflicting Requirements	EP3 Depth of Analysis	EP4 Familiarity of Issues	EP5 Extent of Applicable Codes	EP6 Extent Of Stakeholder Involvement	EP7 Inter-dependence
Depth of Knowledge required to understand and solve the problem	Range of Conflicting Requirements addressed during solution design.	Depth of Analysis performed to evaluate alternatives.	Familiarity of Issues to ensure relevance and alignment with goals.	Extent of Applicable Codes integrated into the solution.	Extent of Stakeholder Involvement considered in problem-solving.	Interdependence of multiple components or disciplines involved.

■ **Engineering Activities**

This table highlights the critical activities involved in executing engineering solutions. It maps factors such as resource utilization, team interaction, innovation, societal impact, and familiarity with standards and technologies to ensure a systematic and effective approach to engineering tasks.

Table 4.3: Mapping with complex engineering activities.

EA1 Range of resources	EA2 Level of Interaction	EA3 Innovation	EA4 Consequences for society and environment	EA5 Familiarity
Range of Resources utilized for the development process.	Level of Interaction required among teams and stakeholders.	Innovation in approach or methodology for problem-solving.	Consequences for Society and Environment considered in the project.	Familiarity with tools, standards, and technologies applied.

Chapter 5

Conclusion

This chapter summarizes the key findings and contributions of the project. It also highlights limitations and suggests directions for future work to extend the study further.

5.1 Summary

This project successfully analyzed e-commerce customer churn patterns by leveraging advanced data analysis techniques. Through systematic implementation and performance evaluation, the study provided actionable insights into customer behavior, contributing to enhanced decision-making for retaining customers and optimizing operations.

5.2 Limitation

While the project demonstrated significant results, certain limitations were noted. The dataset relied on historical data, which may not fully capture real-time behavioral trends. Additionally, the model's accuracy could be influenced by missing values or inconsistent data quality. The integration of external factors such as economic changes or competitor actions was also beyond the scope of this study.

5.3 Future Work

Future work could focus on incorporating real-time data and dynamic customer feedback for a more adaptive churn prediction model. Expanding the model to include external market trends, social media sentiment, and competitor analysis could provide a more comprehensive understanding of customer retention strategies. Integration with advanced AI techniques, such as deep learning, may further improve accuracy and scalability.

References

- [1] J. Brownlee, *Master Machine Learning Algorithms: Discover How They Work and Implement Them From Scratch*, 1st ed. Melbourne: Machine Learning Mastery, 2016.
- [2] A. Géron, *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*, 2nd ed. Sebastopol: O'Reilly Media, 2019.
- [3] S. Sharda, D. Delen, and E. Turban, *Business Intelligence and Analytics: Systems for Decision Support*, 10th ed. Upper Saddle River: Pearson, 2018.
- [4] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York: Springer, 2009.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- [6] R. Kohavi and F. Provost, "Glossary of terms: Special issue on applications of machine learning and the knowledge discovery process," *Machine Learning*, vol. 30, no. 2, pp. 271–274, 1998.
- [7] C. Sammut and G. Webb, *Encyclopedia of Machine Learning and Data Mining*. Boston: Springer, 2017.
- [8] L. Rokach and O. Maimon, *Data Mining With Decision Trees: Theory and Applications*, 2nd ed. Singapore: World Scientific, 2015.
- [9] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: With Applications in R*, 2nd ed. New York: Springer, 2021.
- [10] F. Chollet, *Deep Learning with Python*, 2nd ed. Shelter Island: Manning Publications, 2021.
- [11] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 3rd ed. Upper Saddle River: Pearson, 2010.
- [12] T. Mitchell, *Machine Learning*, 1st ed. New York: McGraw-Hill, 1997.
- [13] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. Burlington: Morgan Kaufmann, 2011.
- [14] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*, 2nd ed. Cambridge: The MIT Press, 2018.
- [15] D. Powers, *Evaluation: From Precision, Recall, and F-Measure to ROC, Informedness, Markedness & Correlation*, 1st ed. Adelaide: Springer, 2020.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Las Vegas, 2016, pp. 770–778.
- [17] M. Abadi et al., "TensorFlow: A system for large-scale machine learning," in *12th USENIX*

Symp. Operating Syst. Design Implementation, Savannah, 2016, pp. 265–283.

[18] G. Hinton and R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.

[19] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*, 1st ed. Cambridge: The MIT Press, 2012.

[20] J. Dean et al., "Large scale distributed deep networks," in *Proc. Neural Inf. Process. Syst. (NIPS)*, Lake Tahoe, 2012, pp. 1223–1231

