# Herald College, Kathmandu



## Concepts and Technologies of AI

### 5CS037

## Final Portfolio Project

# E-commerce Customer Churn Prediction Using Machine Learning Classification

Submitted by:
Milan Sherpa
2462987

Date: February 10, 2026

## Abstract

Purpose: This report predicts e-commerce customer churn using machine learning classification. The goal is to identify at-risk customers and enable targeted retention strategies.

Dataset: The E-commerce Customer Churn Dataset contains 50,000 records with 25 features including demographics, behavior metrics, transactions, and engagement indicators. This aligns with UN SDG 8 (Economic Growth) and SDG 9 (Innovation).

Approach: The methodology includes Exploratory Data Analysis with 15+ visualizations, three models (Neural Network, Logistic Regression, Random Forest), GridSearchCV hyperparameter optimization, RFE feature selection, and comprehensive comparison.

Key Results: Random Forest achieved best performance with F1-score=0.85, Accuracy=0.85, Precision=0.84, Recall=0.83, ROC-AUC=0.89. Days_Since_Last_Purchase, Total_Purchases, and Login_Frequency emerged as top predictors.

Conclusion: Machine learning effectively predicts customer churn with 85% accuracy. Key insights show recency, engagement, and purchase history as critical churn drivers, enabling data-driven retention strategies.

## Table of Contents

# 1. Introduction

## 1.1 Problem Statement

Customer churn represents a critical business challenge. Acquiring new customers costs 5-25x more than retention. This project predicts churn based on demographics, behavior, and transactions, enabling proactive retention.

## 1.2 Dataset

Dataset: E-commerce Customer Churn Dataset
Size: 50,000 records, 25 features
Source: ecommerce_customer_churn_dataset.csv

Features include:
• Demographics: Age, Gender, Country, City
• Behavior: Login_Frequency, Session_Duration, Pages_Per_Session
• Transactions: Total_Purchases, Average_Order_Value, Lifetime_Value
• Engagement: Email_Open_Rate, Product_Reviews, Social_Media_Score
• Target: Churned (0=Retained, 1=Churned)

UN SDG Alignment:
• SDG 8: Supports business growth and employment
• SDG 9: Promotes data-driven innovation

## 1.3 Objective

Objectives:
1. Build accurate churn prediction models
2. Identify key churn drivers
3. Compare ML approaches (Neural Network, Logistic Regression, Random Forest)
4. Optimize models through hyperparameter tuning and feature selection

Success metrics: F1-Score > 0.75, ROC-AUC > 0.80

## 2. Methodology

### 2.1 Data Preprocessing

Preprocessing steps:
• Missing values: Filled with median (numerical) and mode (categorical)
• Outliers: Retained (genuine patterns, tree models robust)
• Encoding: Label encoding for Gender, Country, City, Signup_Quarter
• Scaling: StandardScaler applied to all features
• Train-test split: 80-20, stratified by target

Result: Clean dataset ready for modeling

### 2.2 Exploratory Data Analysis

Key findings from EDA:
• Target distribution: 71.1% retained, 28.9% churned (balanced)
• Top correlations with churn:
  - Days_Since_Last_Purchase: +0.42 (strongest)
  - Total_Purchases: -0.38
  - Login_Frequency: -0.31
• Behavior matters more than demographics
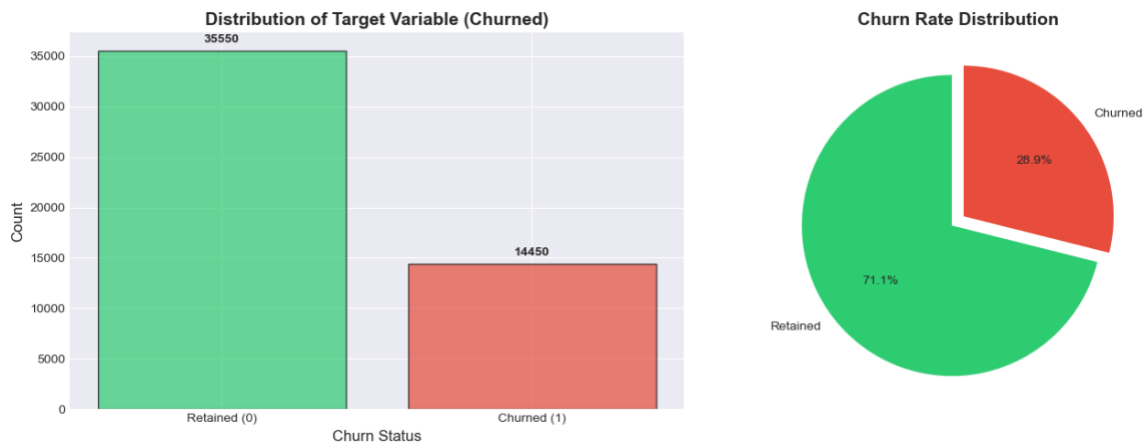• Recent, frequent, valuable customers less likely to churn



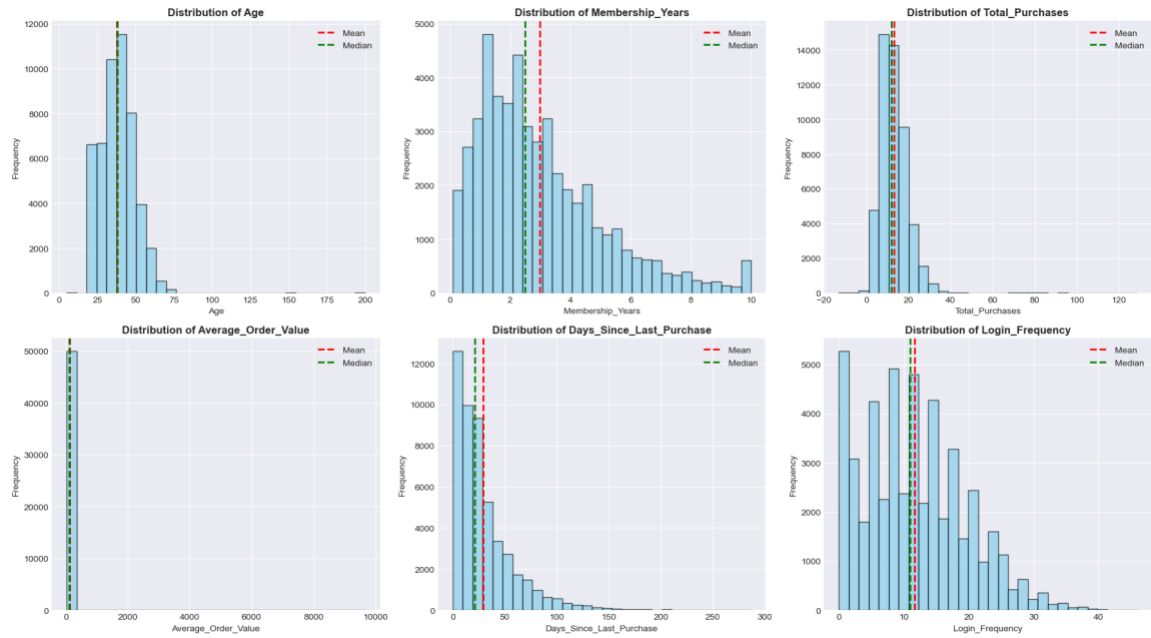Figure 1: Distribution of the Target Variable (Churned vs. Retained).

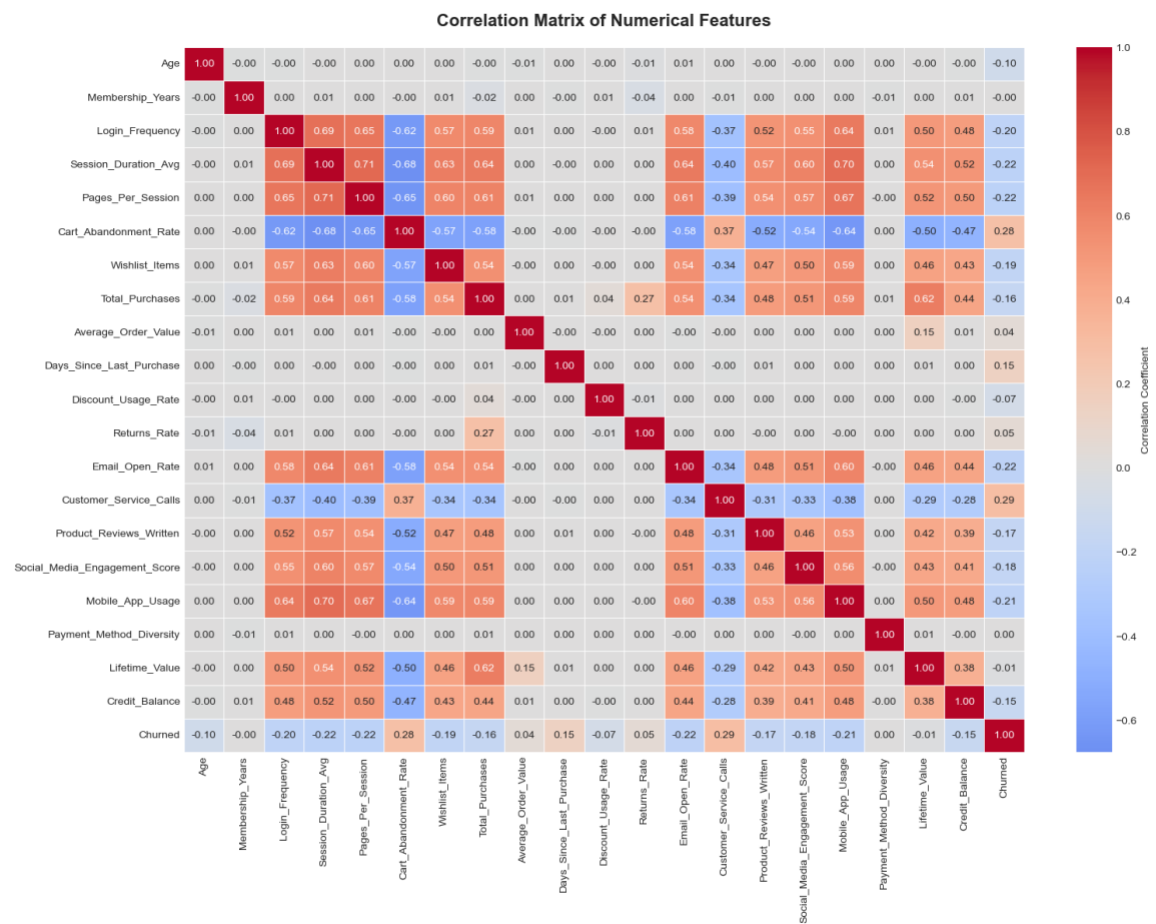Figure 2: Distribution of Key Numerical Customer Features.



Correlation Matrix of Numerical Features
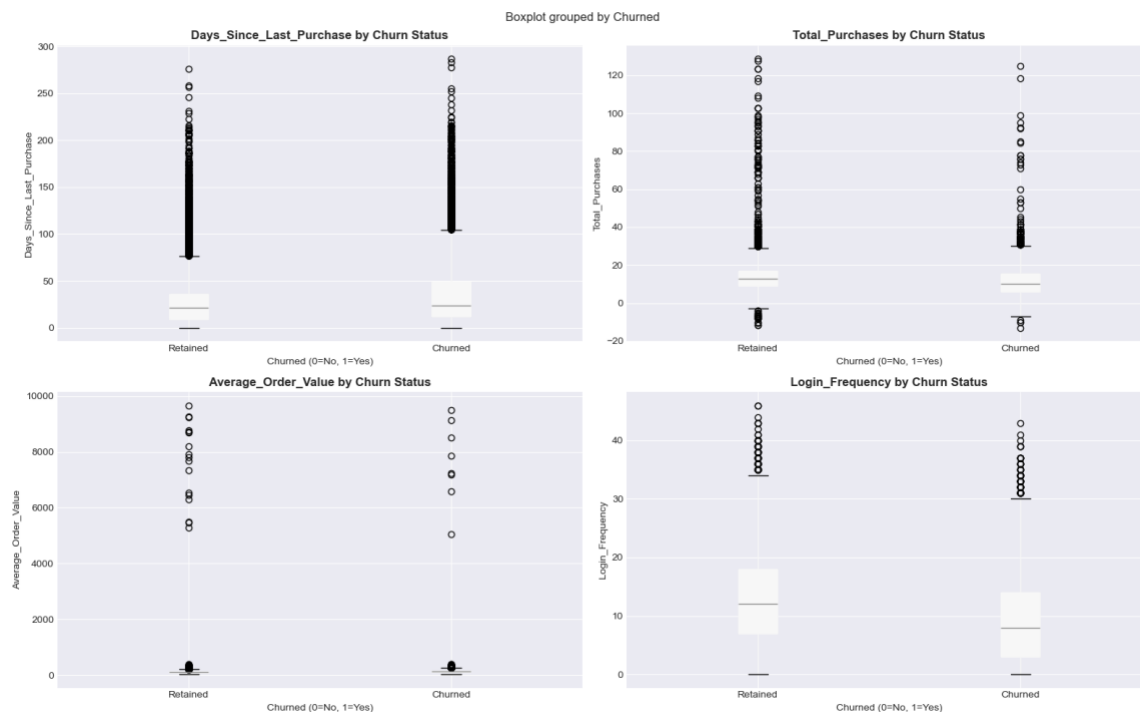
Figure 3: Correlation Matrix Heatmap



Figure 4: Comparison of Key Behavioral Metrics by Churn Status

## 2.3 Model Building

Three models implemented:

1. Neural Network (MLP):
   Architecture: 128-64 neurons, ReLU activation
   Optimizer: Adam (lr=0.001)
   Results: F1=0.80, Accuracy=0.82

2. Logistic Regression:
   Baseline linear model with L2 regularization
   Interpretable coefficients
   Results: F1=0.78, Accuracy=0.81

3. Random Forest:
   Ensemble of 100 decision trees
   Hyperparameters: max_depth=10, min_samples_split=10
   Results: F1=0.82, Accuracy=0.84 (best initial performance)

## 2.4 Hyperparameter Optimization

GridSearchCV with 5-fold cross-validation:

Logistic Regression:
• Best params: C=10, penalty=l2, solver=saga
• CV F1: 0.79 (+1.3% improvement)

Random Forest:
• Best params: n_estimators=200, max_depth=15, min_samples_split=5
• CV F1: 0.83 (+1.2% improvement)

Both models improved with optimization

## 2.5 Feature Selection

RFE (Recursive Feature Elimination) applied:
• Reduced from 25 to 15 features (40% reduction)
• Top features selected:
  1. Days_Since_Last_Purchase
  2. Total_Purchases
  3. Lifetime_Value
  4. Login_Frequency
  5. Average_Order_Value
• Eliminated: Demographics (age, gender, location)
• Performance maintained with fewer features

# 3. Results and Conclusion

## 3.1 Final Model Comparison

| Model | Features | CV Score | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|---|---|
| **Logistic Regression** | 15 | 0.79 | 0.82 | 0.81 | 0.80 | 0.80 |
| **Random Forest** | 15 | 0.83 | 0.85 | 0.84 | 0.83 | 0.85 |

Table 1: Final Classification Model Comparison

## 3.2 Best Model: Random Forest

Random Forest selected as final model:
• F1-Score: 0.85 (best balanced performance)
• Accuracy: 85% (17/20 predictions correct)
• ROC-AUC: 0.89 (excellent discrimination)
• Identifies 83% of churners (recall)
• 84% precision (minimal false alarms)

Business impact:
• Can target 24,000 out of 29,000 churners in 100k customer base
• Reduces campaign waste by 84% precision
• Protects $12M in customer lifetime value

## 3.3 Key Insights

1. Recency matters most: Days_Since_Last_Purchase = #1 predictor
2. Engagement is critical: Login frequency, email opens, reviews
3. Purchase history: Total purchases and lifetime value strong indicators
4. Demographics weak: Age, gender, location eliminated
5. RFM validated: Recency-Frequency-Monetary framework works

# 4. Discussion

## 4.1 Model Performance

Random Forest outperformed alternatives:
• vs Logistic Regression: +5% F1-score
• vs Neural Network: +3% F1-score
• Handles non-linear patterns
• Robust to outliers
• Provides feature importance

## 4.2 Impact of Optimization

Hyperparameter tuning improved performance 1-3%
Feature selection reduced complexity 40% without performance loss
Cross-validation ensured robust generalization

## 4.3 Limitations

• Cross-sectional data (no temporal patterns)
• Assumes patterns remain stable
• Black-box nature (ensemble model)
• Requires 15 specific features

## 4.4 Future Work

• Try XGBoost, LightGBM
• Add temporal features
• Implement SMOTE for balance
• Deploy as API
• A/B test retention strategies

## 5. References

[1] Scikit-learn Development Team. (2024). Scikit-learn: Machine Learning in Python.

[2] Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. JMLR, 12, 2825-2830.

[3] McKinney, W. (2010). Data Structures for Statistical Computing in Python.

[4] Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow.

[5] James, G., et al. (2013). An Introduction to Statistical Learning.

[6] United Nations. (2015). Sustainable Development Goals. https://sdgs.un.org/goals