

FPCP3

Robin Tran

Introduction

This project explores customer behavior using a dataset with demographic, transactional, and engagement features. There are two main sections. In the first section, we implemented linear regression models to predict each customer's average transaction value, comparing different model specifications using 10-fold cross-validation. In the second section, we used logistic regression to classify whether a customer is at high churn risk based on selected predictors. The goal is to understand what factors are associated with customer spending and retention, and to evaluate model performance using appropriate validation techniques.

Variables of Interest

I will be considering 7 numerical variables and 12 categorical variables. They are listed below along.

Table 1: Numerical Variables

Variable	Description
age	The age of the user, measured in years.
days_since_last_login	The number of days since the user last logged into the website, measured in days.
avg_time_spent	The average amount of time the user spends per visit on the website, measured in minutes.
avg_transaction_value	The average value of transactions made by the user, measured in monetary units (currency).
avg_frequency_login_days	The average number of days between the user's consecutive logins, measured in days.
points_in_wallet	The amount of loyalty or reward points currently available in the user's wallet, measured in points.
churn_risk_score	The churn risk score assigned to the user, ranging from 1 to 5, indicating the likelihood of the user leaving the service (higher scores indicate higher risk).

Table 2: Categorical Variables

Variable	Categories
gender	F, M, Unknown
region_category	Village, City, Town
membership_category	Platinum Membership, Premium Membership, No Membership, Gold Membership, Silver Membership, Basic Membership
joined_through_referral	Yes, No
preferred_offer_types	Gift Vouchers/Coupons, Credit/Debit Card Offers, Without Offers
medium_of_operation	Desktop, Smartphone, Both
internet_option	Wi-Fi, Mobile_Data, Fiber_Optic
used_special_discount	Yes, No
offer_application_preference	Yes, No
past_complaint	Yes, No
complaint_status	Not Applicable, Solved, Solved in Follow-up, Unsolved, No Information Available
feedback	Products always in Stock, Quality Customer Care, Poor Website, No reason specified, Poor Product Quality, Poor Customer Service, Too many ads, User Friendly Website, Reasonable Price

Observational Unit

Each row in the data set represents one individual customer who has engaged with the platform and has at least one recorded purchase.

We also verified the categorical variables' types.

Exploratory Data Analysis

Summaries

We included statistical summary for our numerical variables as below.

We also included statistical summaries for some of our categorical variables as below. The summaries are quite insightful. We observe balanced classes in **gender**, but there might be some imbalances in **membership_category** and **region_category**.

Visualizations

To better understand the relationship between user behavior and churn risk, we included a variety of plots.

Table 3: Summary Statistics for Numerical Variables

Variable	mean	median	sd	IQR	min	max	n
age	37.08	37.00	15.91	28.00	10.00	64.00	24853
days_since_last_login	-42.61	12.00	230.18	8.00	-999.00	26.00	24853
avg_time_spent	244.14	162.37	398.10	295.65	-2281.24	3040.41	24853
avg_transaction_value	29321.54	27534.68	19499.51	26605.22	800.46	99914.05	24853
avg_frequency_login_days	15.97	16.00	9.23	14.00	-43.65	73.06	24853
points_in_wallet	688.30	698.66	195.79	148.52	-760.66	2069.07	24853
churn_risk_score	3.61	4.00	1.18	2.00	1.00	5.00	24853

Table 4: Distribution of Gender

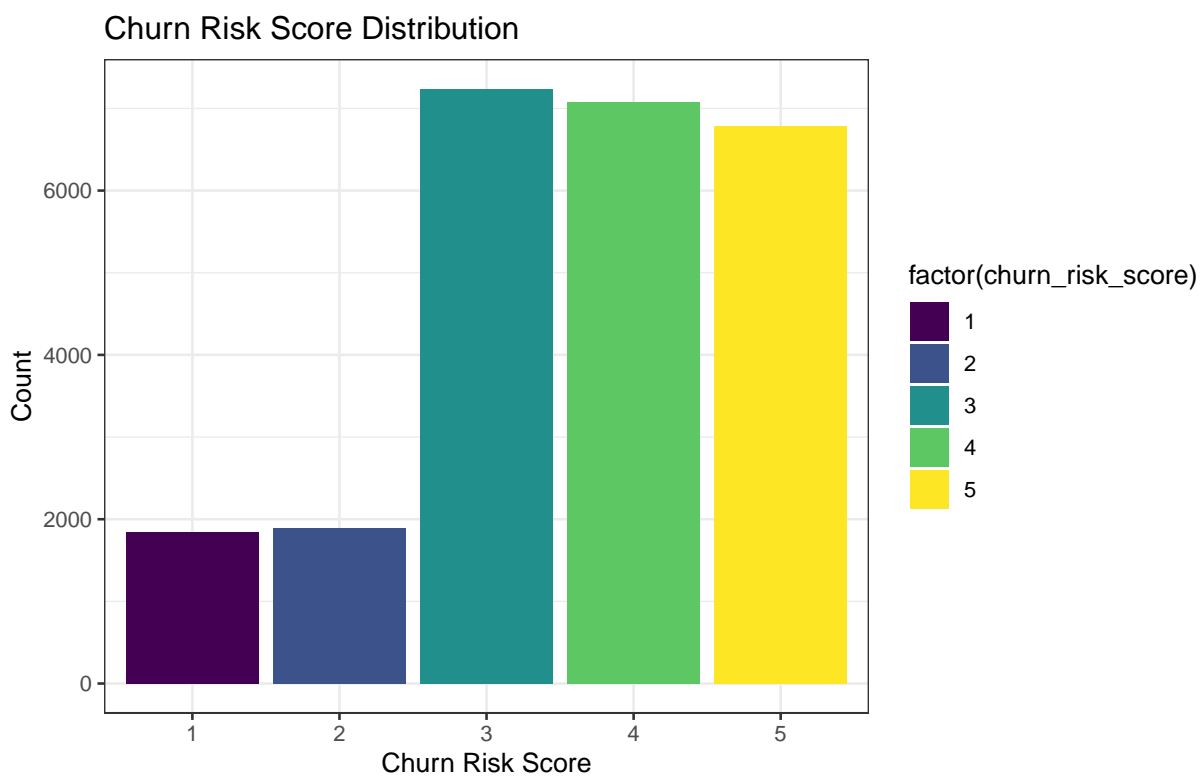
gender	n	percentage
F	12463	50.15
M	12351	49.70
Unknown	39	0.16

Table 5: Distribution of Membership

membership_category	n	percentage
Basic Membership	5131	20.65
Gold Membership	4535	18.25
No Membership	5198	20.91
Platinum Membership	2912	11.72
Premium Membership	2982	12.00
Silver Membership	4095	16.48

Table 6: Distribution of Region

region_category	n	percentage
City	10020	40.32
Town	11090	44.62
Village	3743	15.06



Membership Category Distribution

