

Customer Churn Analysis

Prepared By :
**Shakirah
Benson**

Customer Churn Overview

Executive Summary

This analysis explores the key drivers behind customer churn on an e-commerce platform. Using Python for data cleaning, visualization, segmentation, and predictive modeling, we identified actionable insights that can help improve customer retention. Key findings include the impact of order frequency, satisfaction scores, and complaint status on churn likelihood. We also developed customer segments using clustering and evaluated churn risk using a Random Forest classifier.

Project Goals

1. Clean and prepare the dataset for analysis
2. Explore data to identify patterns and trends
3. Visualize key relationships
4. Predict churn using classification modeling
5. Deliver actionable insights for retention strategies

Key Business Questions

1. Do customers with fewer orders or low engagement tend to churn more often?
2. What is the relationship between satisfaction scores, tenure, and churn?
3. How do other factors like complaints, payment modes, and cashback amounts correlate with churn?

Tools & Methods

- Python (pandas, matplotlib, seaborn, scikit-learn)
- Data Cleaning with Imputation and Type Conversion
- Correlation Analysis and Visualizations
- K-Means Clustering for Segmentation
- Random Forest Classification for Churn Prediction
- Customer Lifetime Value Estimation

Data Cleaning & Preparation

The dataset used in this analysis is sourced from Kaggle, specifically from the Customer Churn Dataset created by Muhammad Shahid Azeem. It contains detailed information about customer behavior on an e-commerce platform, including tenure, order count, satisfaction scores, preferred payment methods, and churn status.

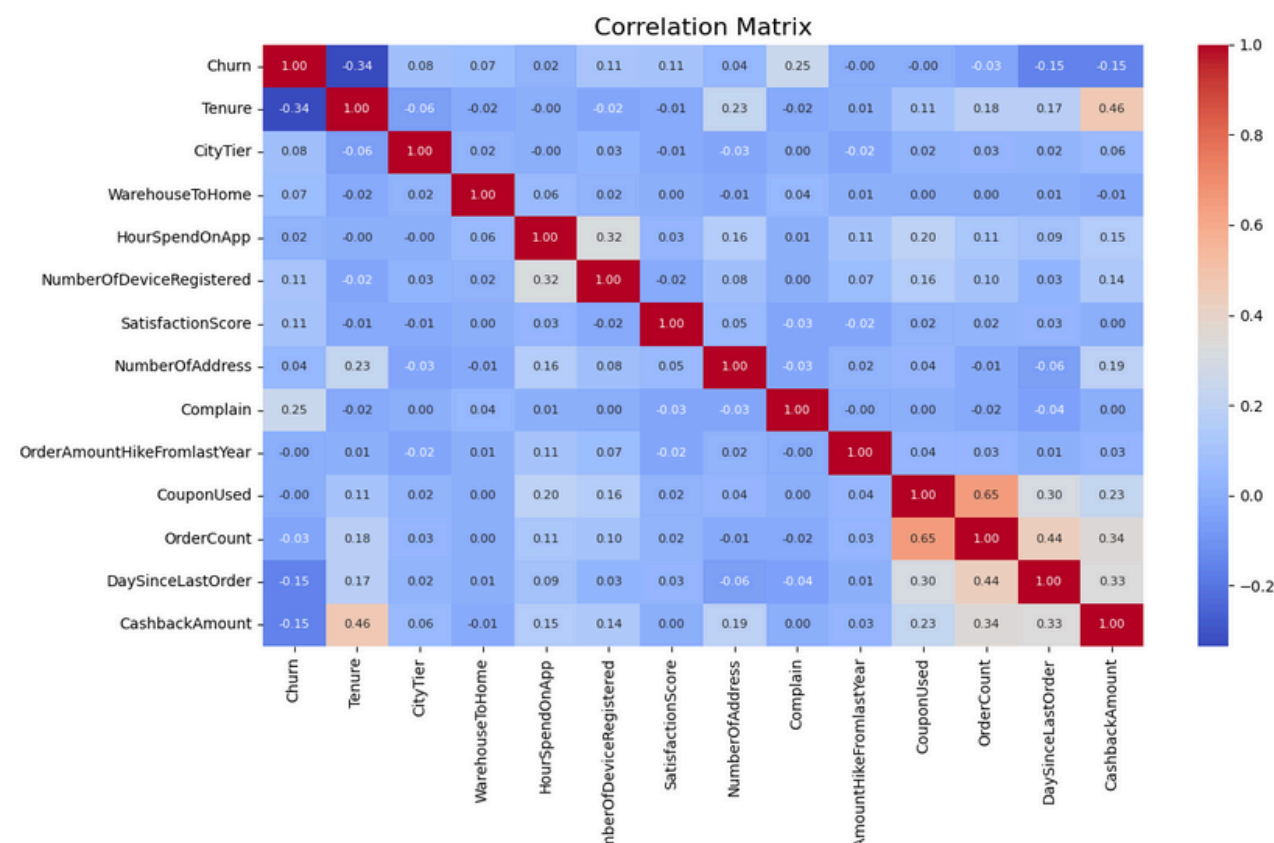
We first load the data, clean any missing or invalid entries, and ensure all numeric columns are correctly formatted. We used mean imputation for numeric columns with missing values. The dataset is loaded into a Pandas DataFrame, and missing or non-numeric values in important columns (e.g., 'Tenure', 'SatisfactionScore', 'OrderCount') are converted to numeric format, filling missing values with the column mean.

Data columns (total 20 columns):				
Dtype	#	Column	Non-Null Count	
-	0	CustomerID	5630	non-null
int64	1	Churn	5630	non-null
int64	2	Tenure	5630	non-null
float64	3	PreferredLoginDevice	5630	non-null
object	4	CityTier	5630	non-null
int64	5	WarehouseToHome	5630	non-null
float64	6	PreferredPaymentMode	5630	non-null
object	7	Gender	5630	non-null
object	8	HourSpendOnApp	5630	non-null
float64	9	NumberOfDeviceRegistered	5630	non-null
int64	10	PreferedOrderCat	5630	non-null
object	11	SatisfactionScore	5630	non-null
int64	12	MaritalStatus	5630	non-null
object	13	NumberOfAddress	5630	non-null
int64	14	Complain	5630	non-null
int64	15	OrderAmountHikeFromLastYear	5630	non-null
float64	16	CouponUsed	5630	non-null
float64	17	OrderCount	5630	non-null
float64	18	DaySinceLastOrder	5630	non-null
float64	19	CashbackAmount	5630	non-null int64

Visualizing the Data

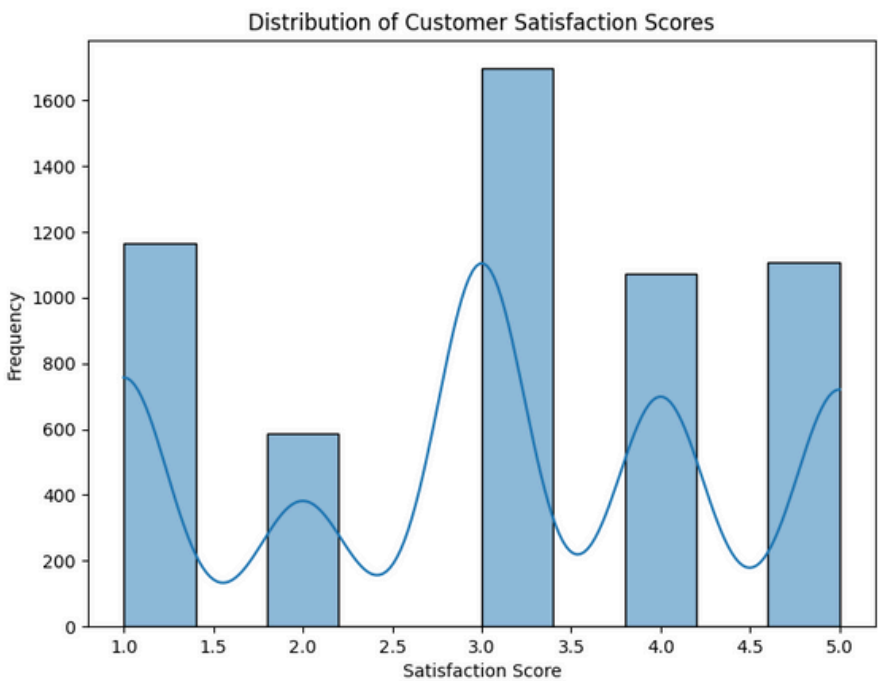
Correlation Matrix

To identify relationships between numeric variables, we generated a correlation matrix. This helps highlight key factors that could be contributing to customer churn. A correlation matrix was computed to identify key relationships among numerical variables. The matrix highlights significant correlations, such as between OrderCount and Churn (negative), and between SatisfactionScore and Churn (negative). This indicates that higher satisfaction and order frequency are associated with lower churn. These findings direct attention to customer engagement and satisfaction as key areas for retention.



Customer Satisfaction Scores

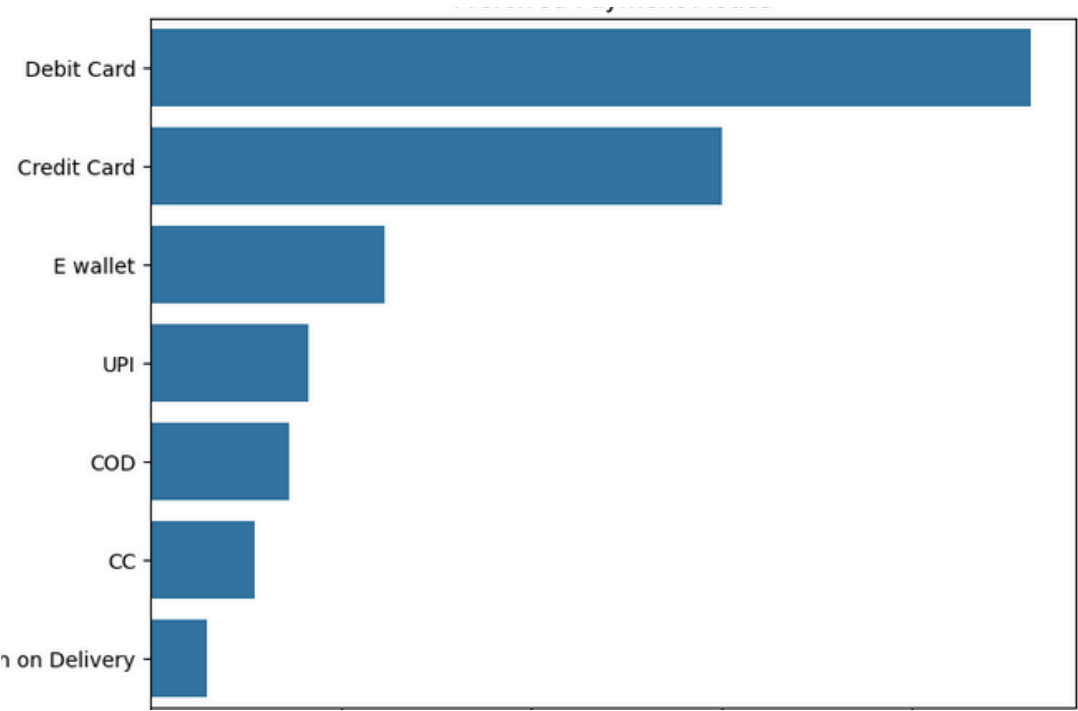
The histogram reveals that the majority of customers have moderate satisfaction scores. However, there is a notable group with low scores, potentially contributing to churn. Addressing these customers' concerns could improve retention.



Visualizing the Data

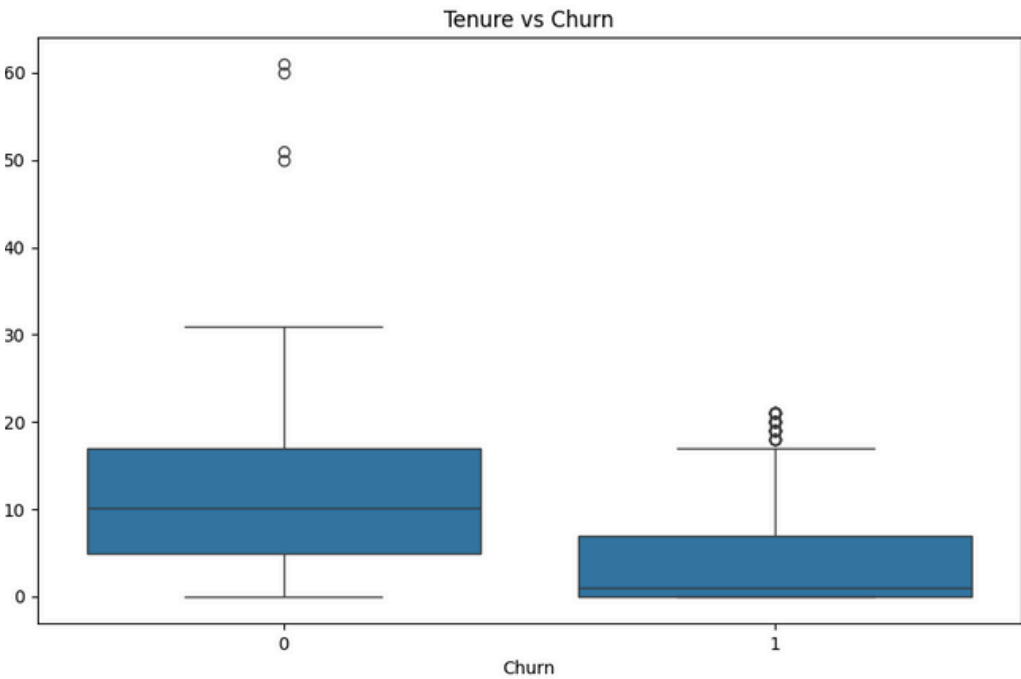
Payment Method Preferences

A count plot was used to visualize the distribution of preferred payment methods among customers. This helped identify the most and least popular options, providing insight into customer preferences that could inform product or UX decisions. Customers using less preferred payment methods may be less satisfied and more prone to churn, warranting exploration into whether certain payment modes affect user experience.



Tenure vs Churn

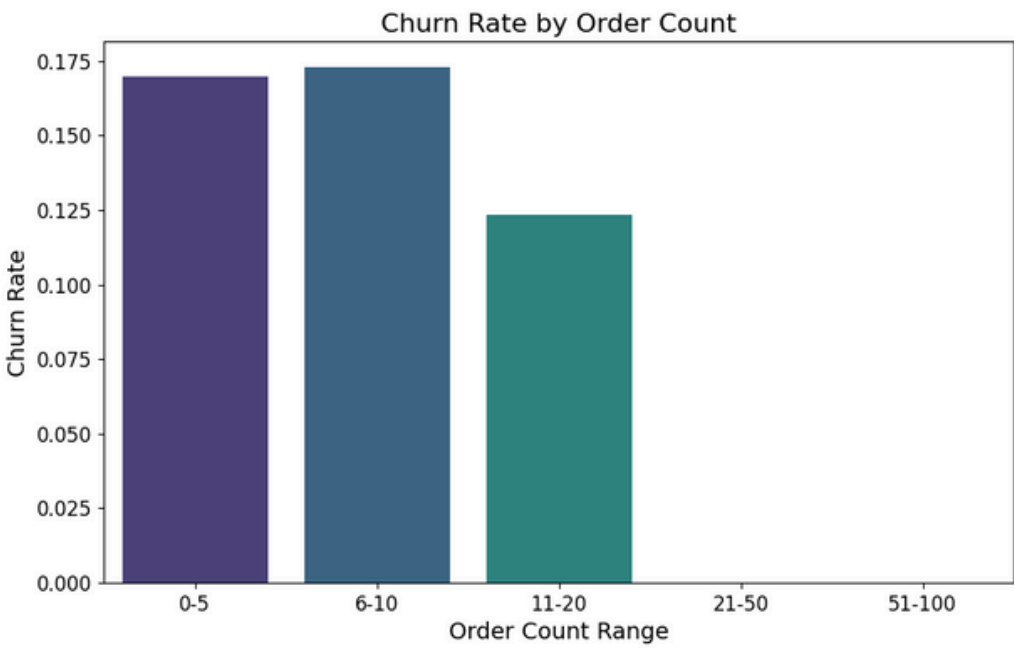
Customers with shorter tenures are significantly more likely to churn, highlighting the importance of strong onboarding and early engagement strategies. This insight was revealed through a boxplot comparing tenure across churned and retained users.



Visualizing the Data

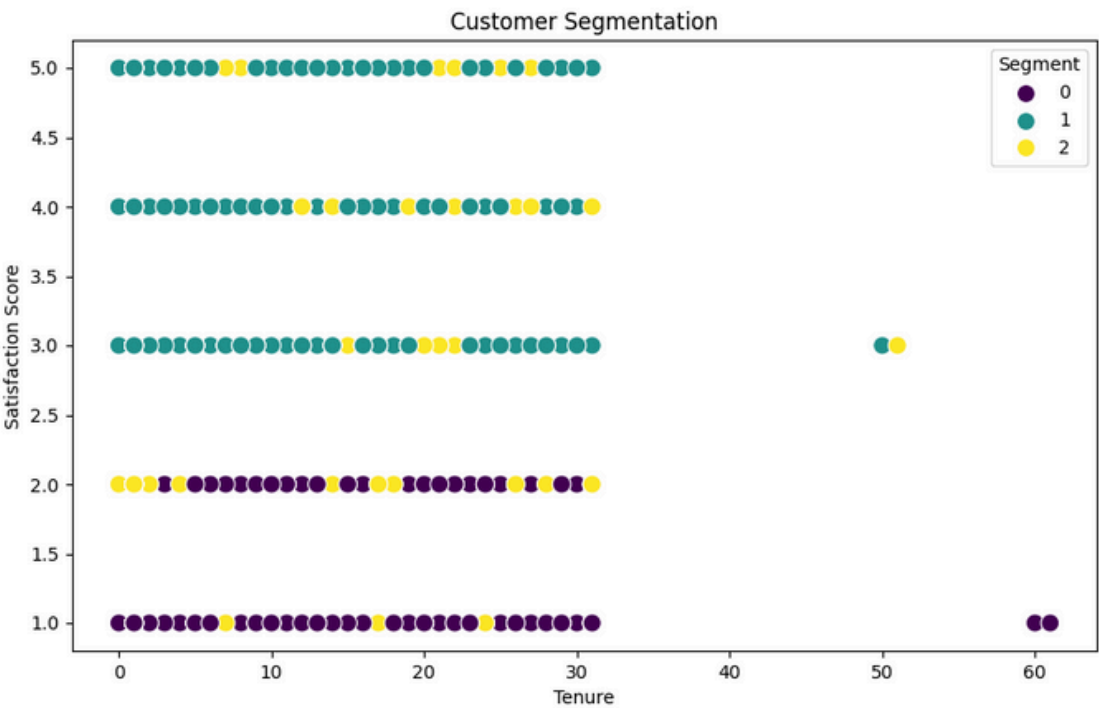
Churn Analysis by Order Count

A bar plot was used to examine how churn rates vary across different ranges of order counts. The results showed that customers with lower order counts (particularly in the 0–5 range) exhibited the highest churn rates. This insight highlights the importance of encouraging early repeat purchases to improve retention.



Customer Segmentation

K-Means clustering was applied to tenure, satisfaction scores, and order volume, resulting in three distinct customer segments. Segment 2 included high-tenure, high-frequency customers, our most valuable group for retention. Segment 0 showed low satisfaction and low order volume, signaling a churn-risk group while 1 showed satisfaction, but low activity. This segmentation supported data-driven targeting strategies for loyalty programs and re-engagement.



Visualizing the Data

Customer Lifetime Value (CLV)

Customer Lifetime Value is calculated by combining cashback earnings with churn risk estimates. This allowed us to identify high-value customers who were likely to stay, helping prioritize retention and upsell efforts to maximize ROI. Customers who receive higher cashback and have a lower likelihood of churning represent greater long-term value. By calculating Customer Lifetime Value (CLV), the top 10 most valuable customers were identified and prioritized for targeted retention efforts.

C U S T O M E R I D		C L V
3 5 4 1	5 3 5 4 2	3 2 2 . 4 4 0 1 3 6
5 0 1 1	5 5 0 1 2	3 2 2 . 2 3 1 4 4 2
2 8 8 0	5 2 8 8 1	3 2 2 . 0 5 4 8 9 5
4 3 5 0	5 4 3 5 1	3 2 1 . 4 4 7 7 6 2
5 1 6 9	5 5 1 7 0	3 2 0 . 8 2 9 1 0 7
3 6 9 9	5 3 7 0 0	3 2 0 . 8 2 9 1 0 7
3 6 4 6	5 3 6 4 7	3 2 0 . 5 2 7 8 1 2
5 1 1 6	5 5 1 1 7	3 2 0 . 4 8 0 8 4 2
5 2 7 6	5 5 2 7 7	3 2 0 . 2 7 8 6 0 0
4 3 2 5	5 4 3 2 6	3 2 0 . 2 7 2 9 9 7

Customer Churn Predictor

A Random Forest Classifier was used to predict customer churn based on features such as tenure, satisfaction score, order frequency, and coupon usage. The model achieved an overall accuracy of 86%, performing especially well in identifying customers likely to remain active. Although recall for churned customers was lower, the model provides a strong foundation for an early warning system. With further tuning or class rebalancing, it could be deployed to trigger automated retention efforts for at-risk users.

P R E C I S I O N		R E C A L L
F 1 - S C O R E	S U P P O R T	
0	0 . 8 9	0 . 9 4
0 . 9 2	9 4 1	
1	0 . 5 8	0 . 4 2
0 . 4 9	1 8 5	
A C C U R A C Y		
0 . 8 6	1 1 2 6	
M A C R O A V G	0 . 7 4	0 . 6 8
0 . 7 0	1 1 2 6	
W E I G H T E D A V G	0 . 8 4	0 . 8 6
0 . 8 5	1 1 2 6	

Insights & Recommendations

1. High-risk segments:

- a. Target customers with 0-5 orders with retention campaigns
- b. Implement immediate response protocols for customer complaints, as complaint-filing

customers show higher churn rates

2. Customer Lifetime Value (CLV):

- a. Prioritize retention efforts on Segment 2 customers (high-tenure, high-satisfaction, frequent orders)
- b. Develop programs targeted towards customers with high CLV scores above 320

3. Measurement and Monitoring:

- a. Track customer satisfaction scores monthly
- b. Monitor complaint resolution time and impact
- c. Measure the effectiveness of retention campaigns using the developed metrics
- d. Develop an early warning system using the Random Forest model (86% accuracy) to identify potential churners

Next Steps

Future improvements could include tuning the Random Forest model for better churn recall, incorporating additional behavioral features (e.g., browsing activity or customer support interactions), and testing A/B retention strategies on high-risk segments. Deploying a real-time dashboard to monitor churn indicators and CLV could further support proactive decision-making.

Github Script Link

<https://github.com/Skirah/customer-churn-analysis>