



Churn Prediction

Data-Driven Insights for GTM, and Product

Objective



- Analyze the provided dataset: Customer Churn
- Develop actionable insights and recommendations for CX

Data Exploration



- Load and preprocess the dataset
- Perform exploratory data analysis (EDA) for the dataset
- Visualize key features and relationships

Features:

- Categorical:
 - OrderCat, Gender, PreferredLoginDevice, MaritalStatus, PaymentMode
- Numeric:
 - Churn, Tenure, CityTier, WarehouseToHome, HourSpendOnApp, NumberOfDeviceRegistered, SatisfactionScore, NumberOfAddress, Complain, OrderAmountHikeFromlastQuarter, CouponUsed, OrderCount, DaySinceLastOrder, CashbackAmount

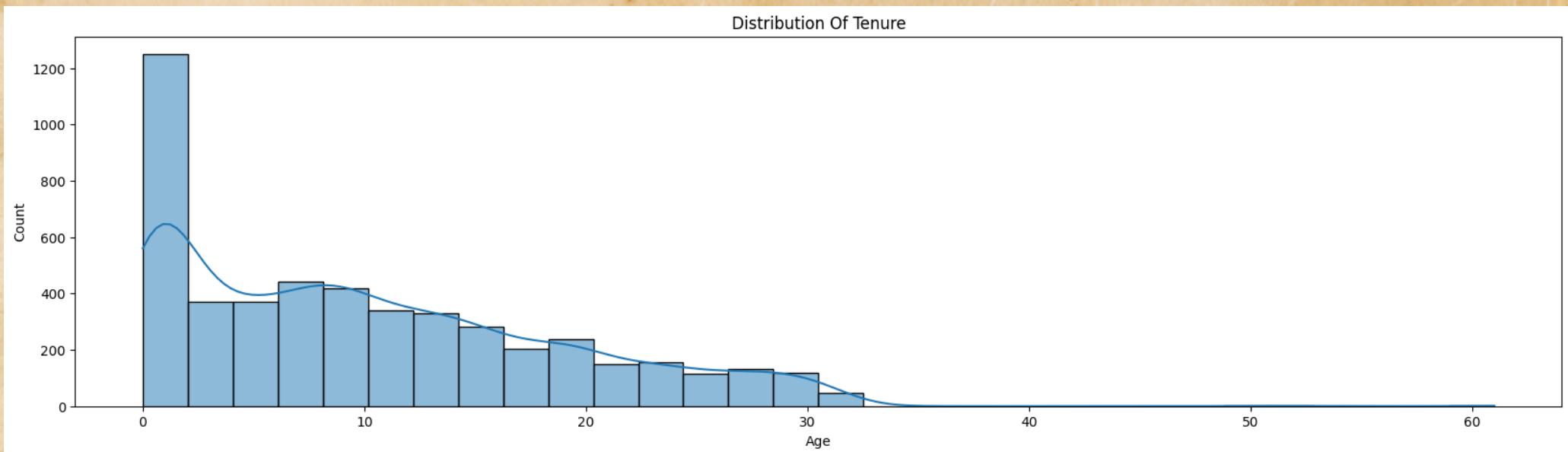
EDA - 1.1 Tenure Distribution vs Churn Rate Analysis



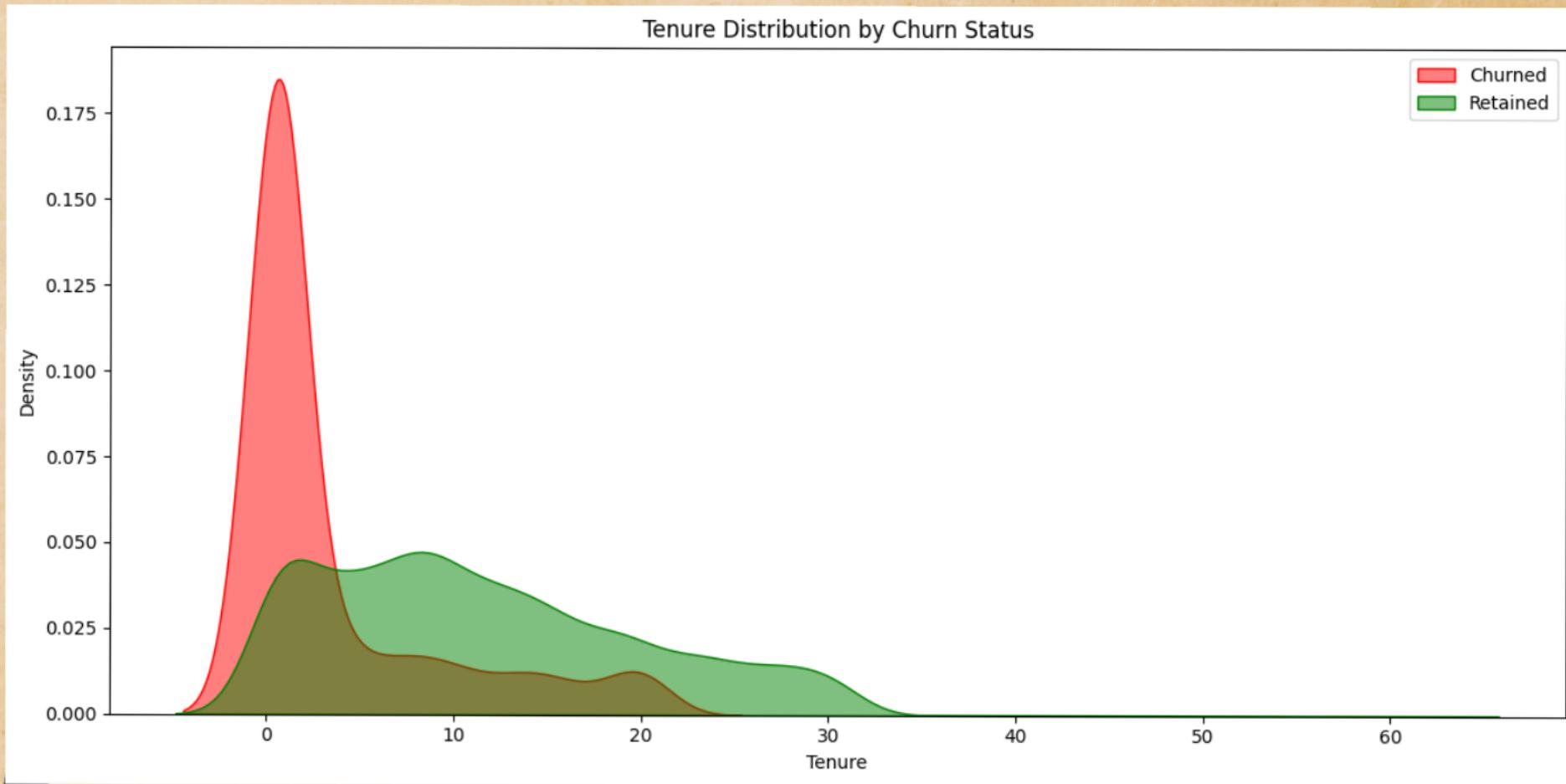
Tenure Statistics

- count 4964.000000
- mean 10.276793
- std 8.630337
- min 0.000000
- 25% 2.000000
- 50% 9.000000
- 75% 16.000000
- max 61.000000
-

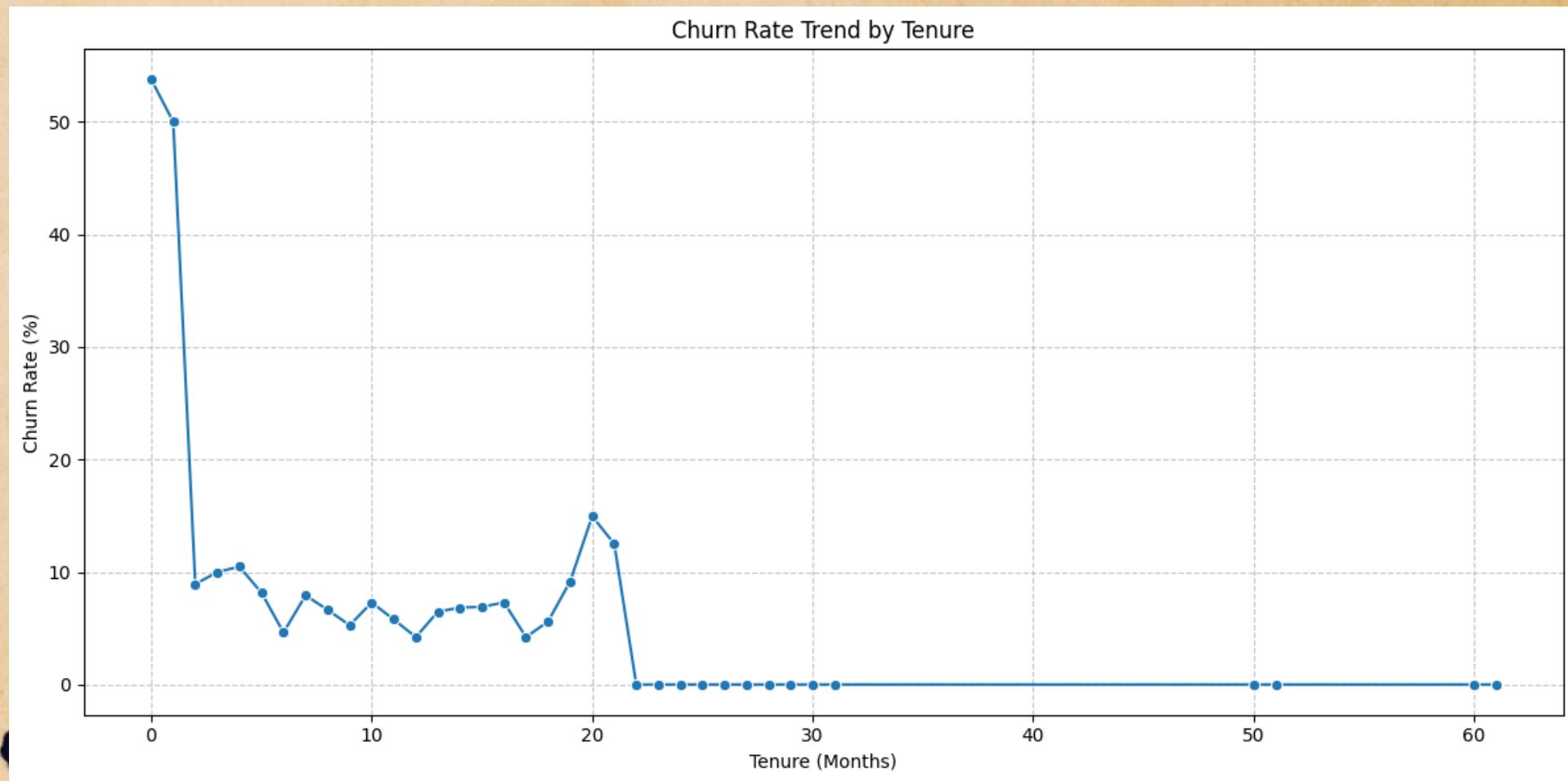
Distribution of Tenure



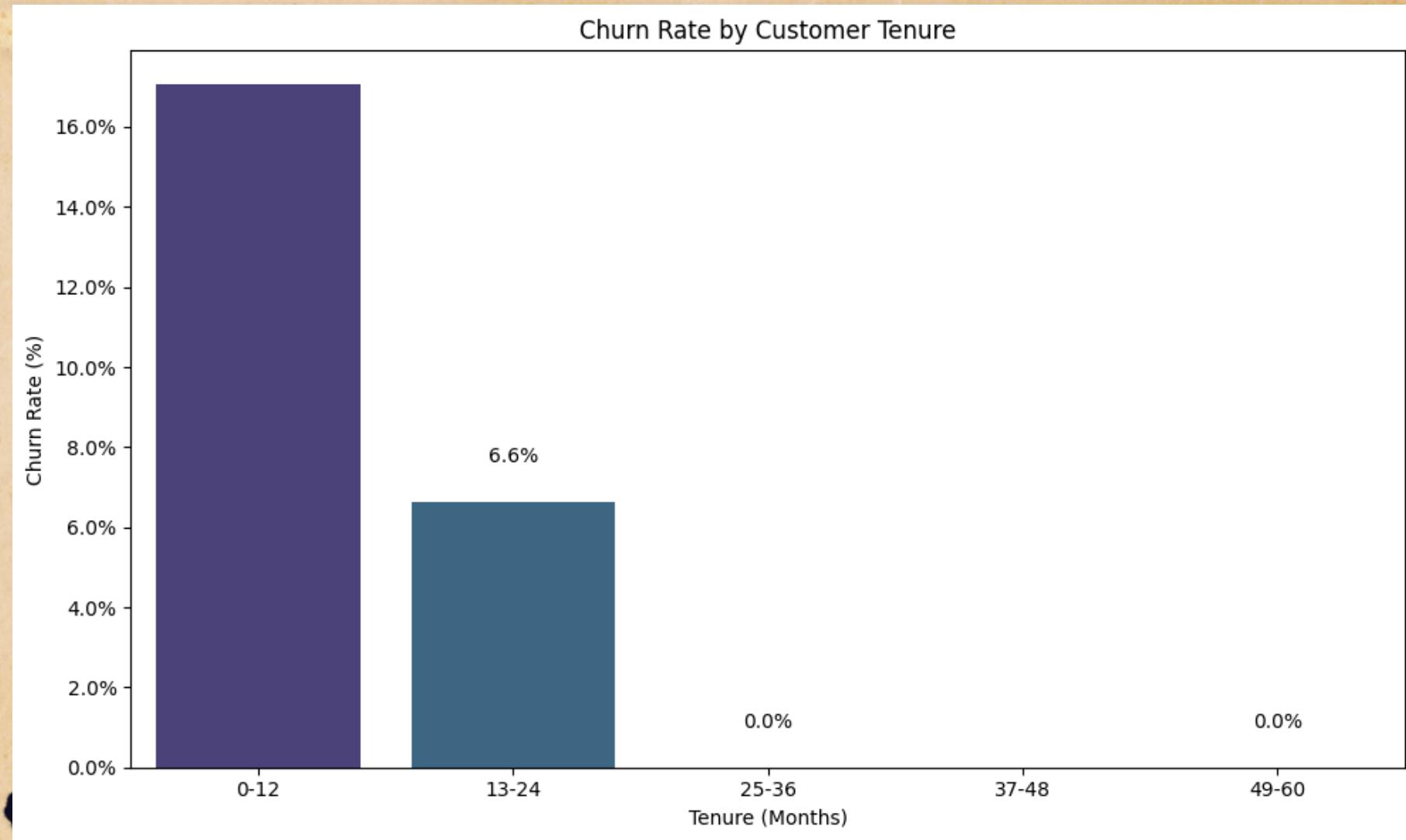
Kernel Density Estimate



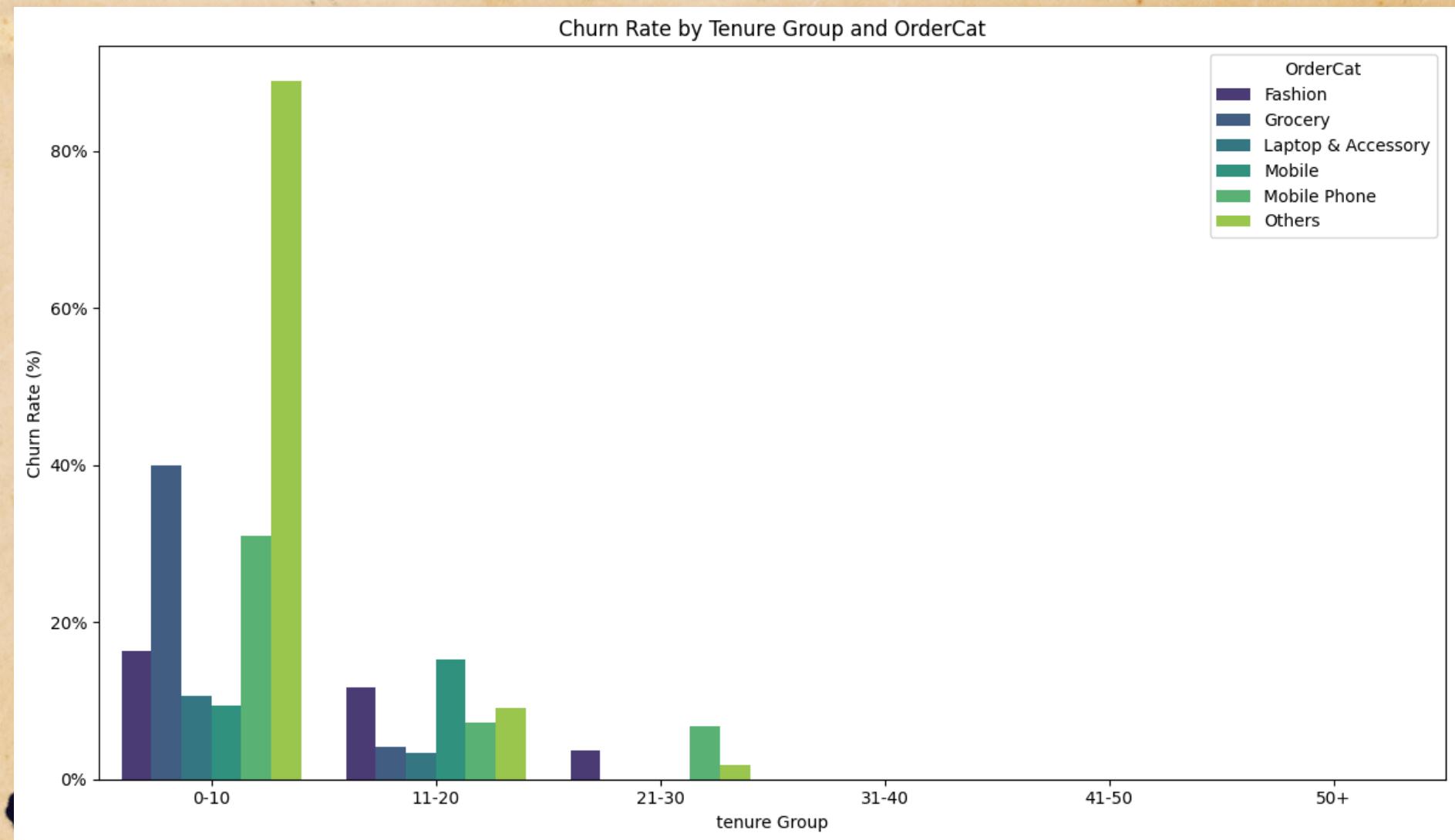
Churn rate for each tenure month



Tenure Group: 0-12, 13-24, 25-36, 37-48, 49-60

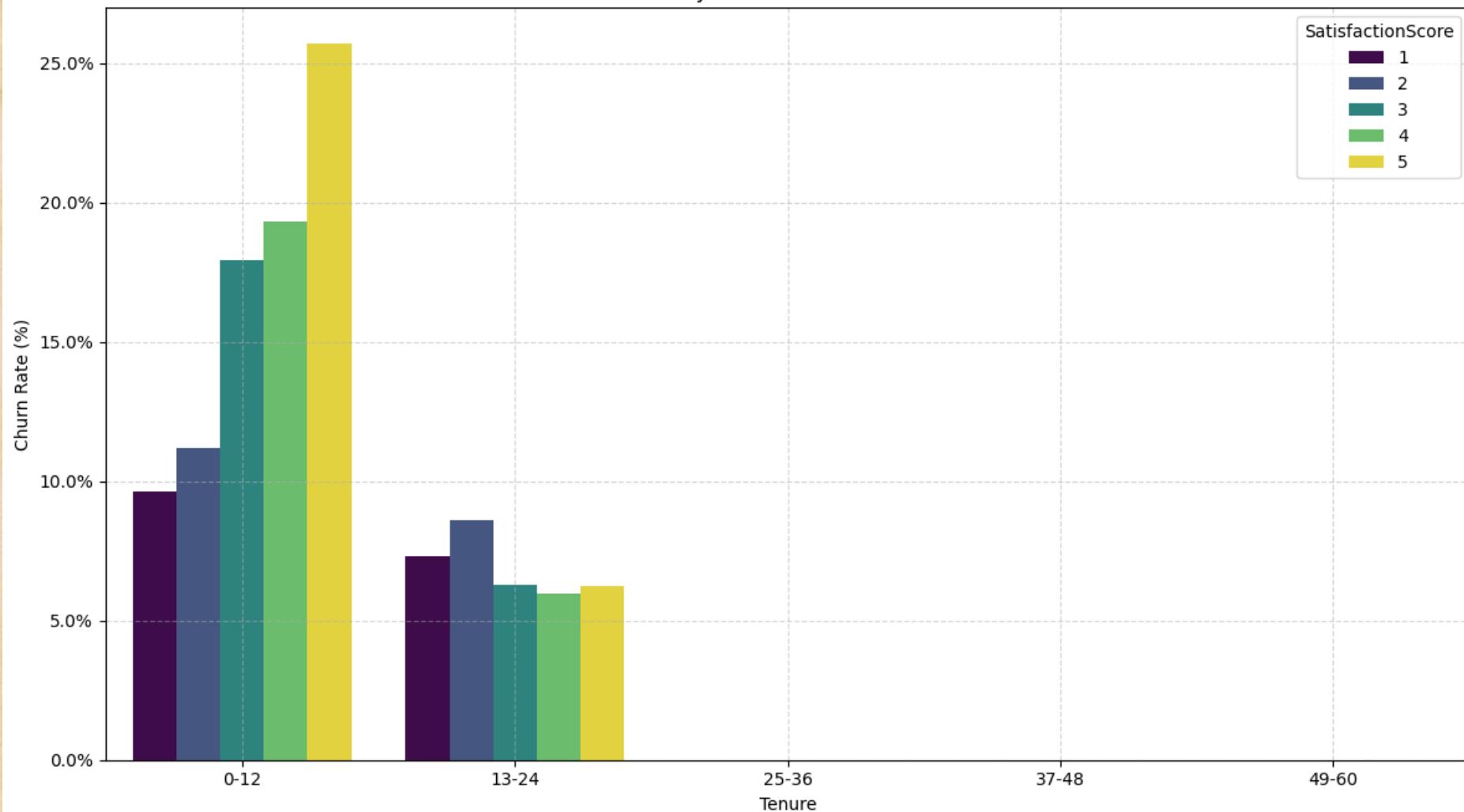


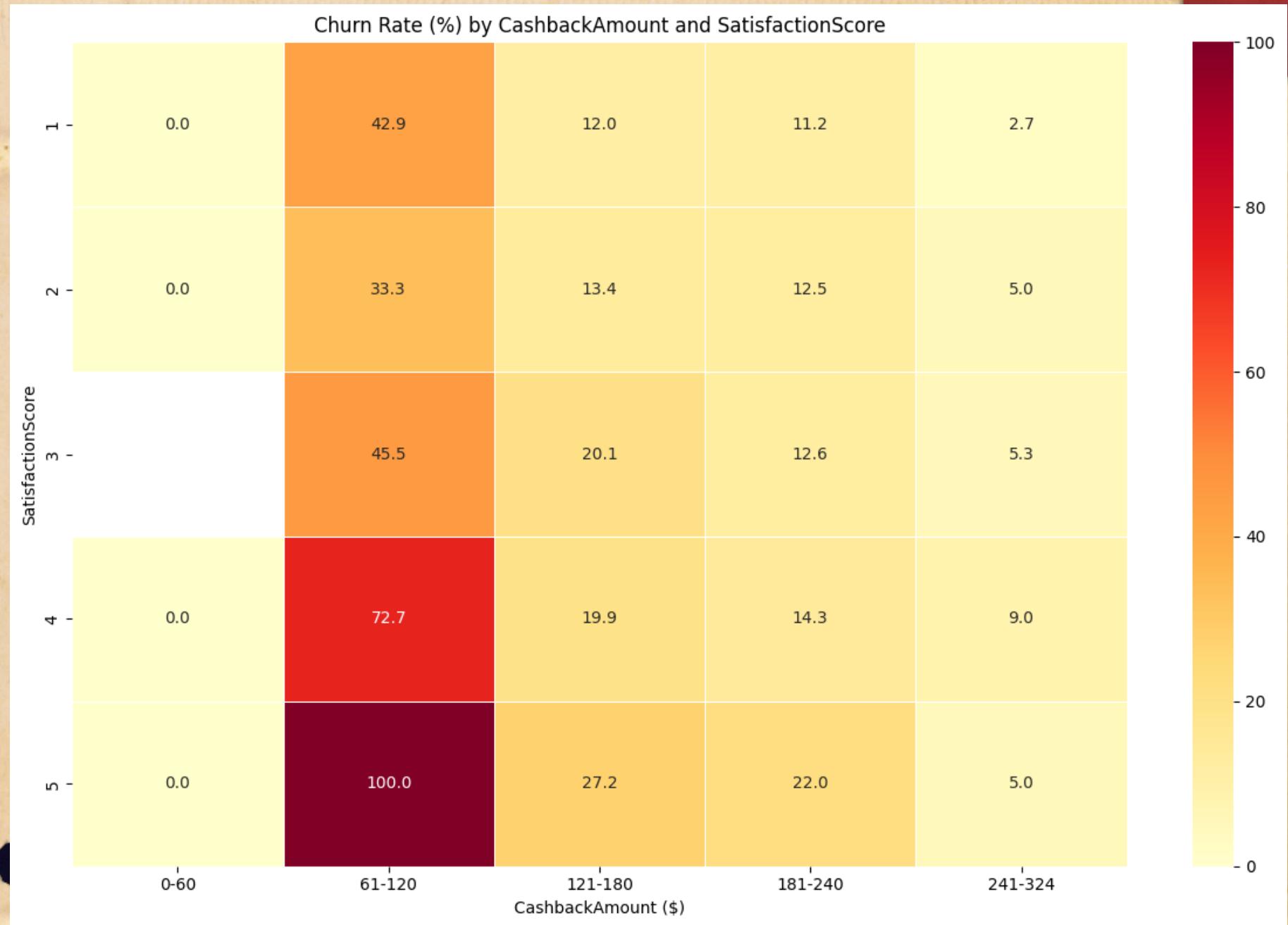
Churn Rate by multiple features





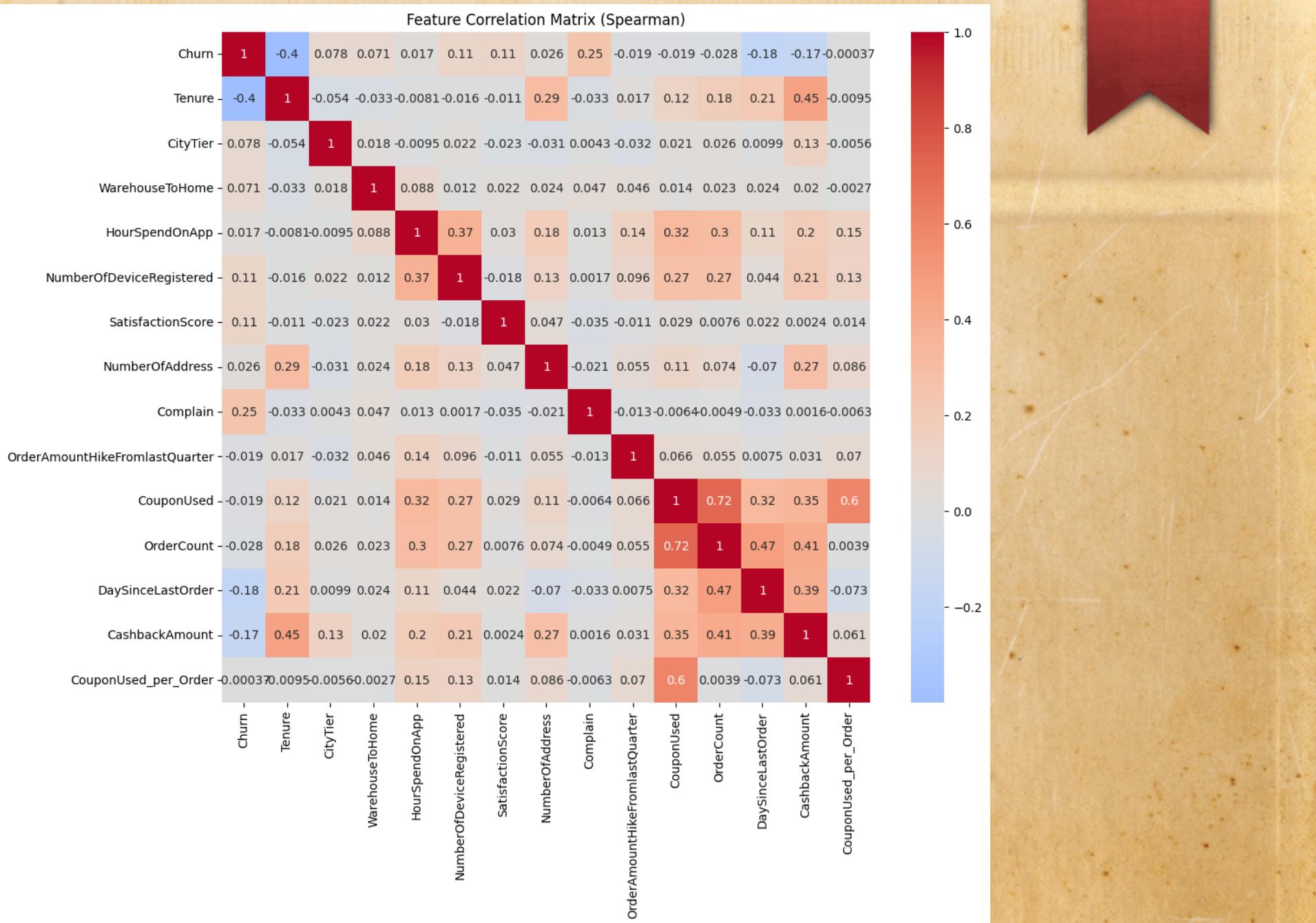
Churn Rate by Tenure and SatisfactionScore





EDA Statistical Correlation Analysis

- Numeric Features

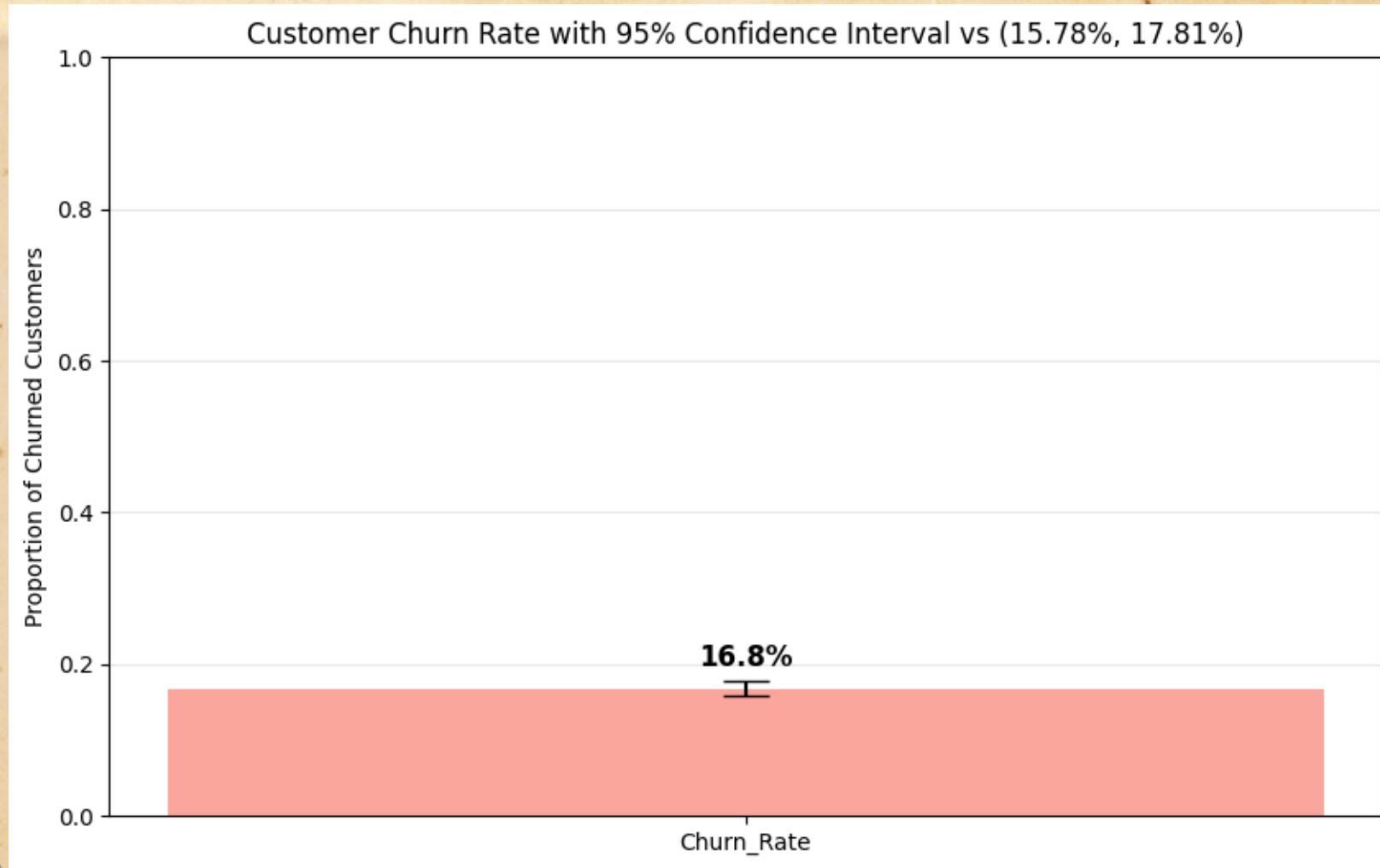


EDA Statistical Correlation Analysis



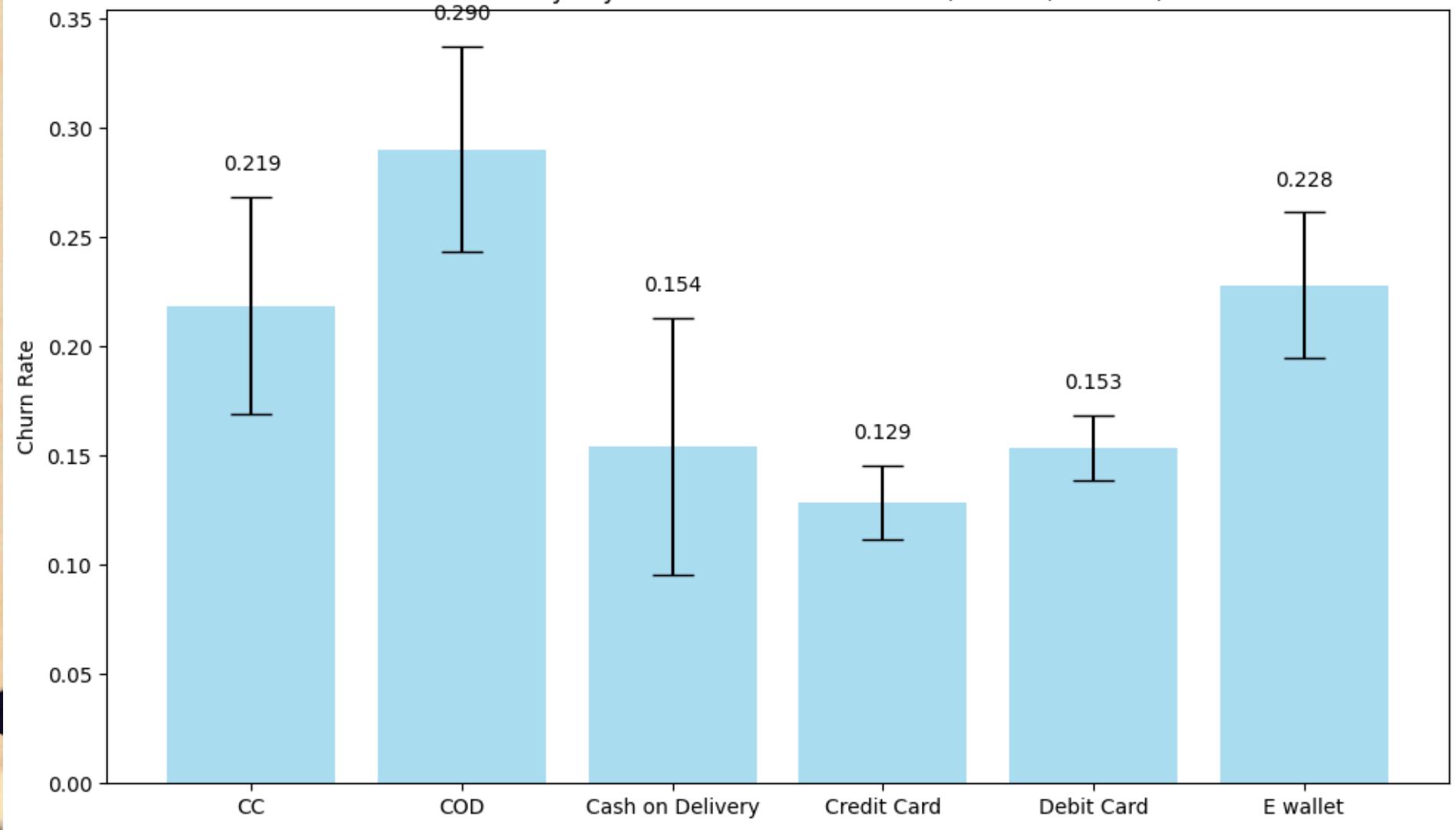
- Categorical Features

Confidence Interval churn rate



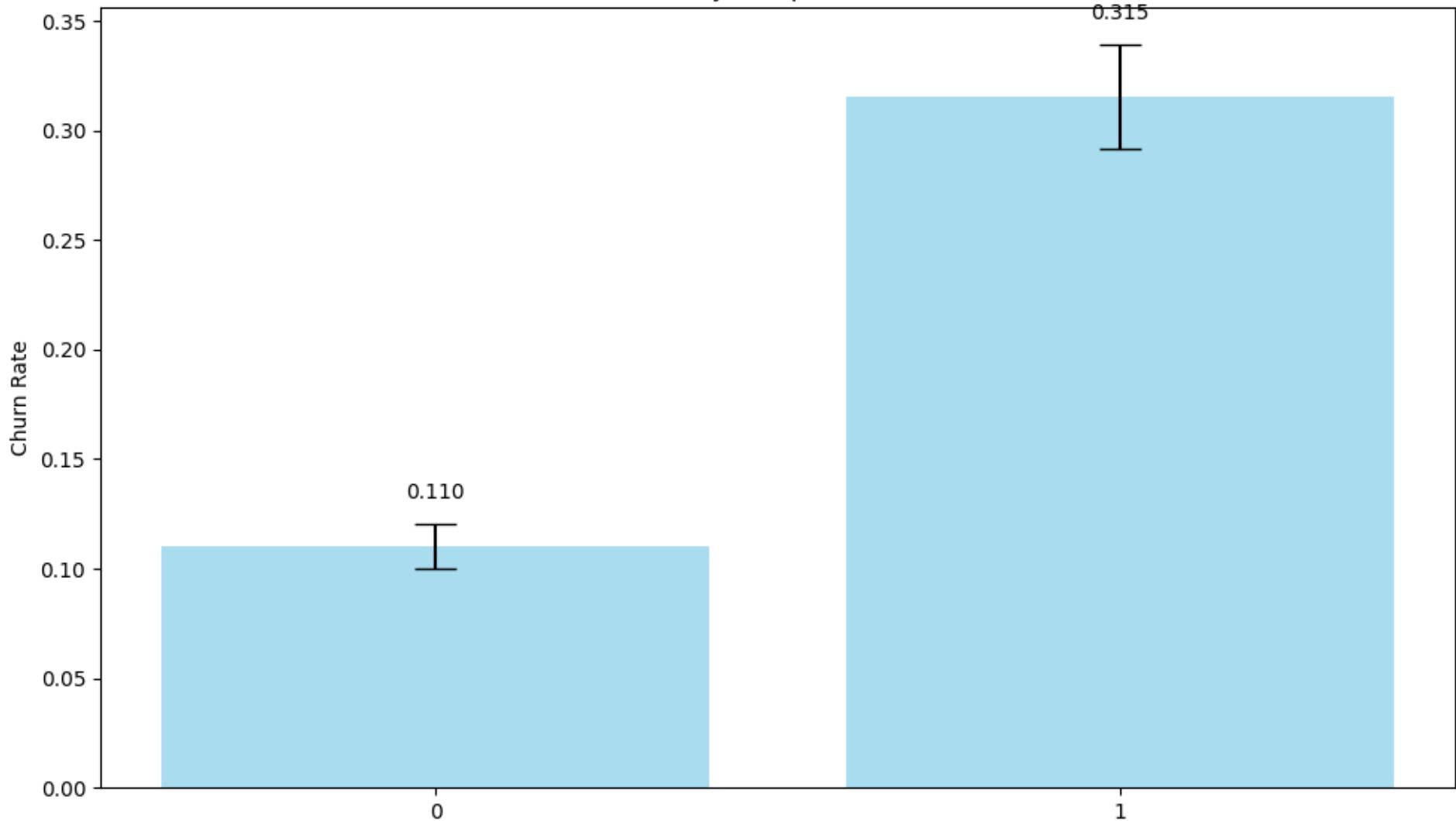


Churn Rate by PaymentMode with 95% CI vs (15.78%, 17.81%)





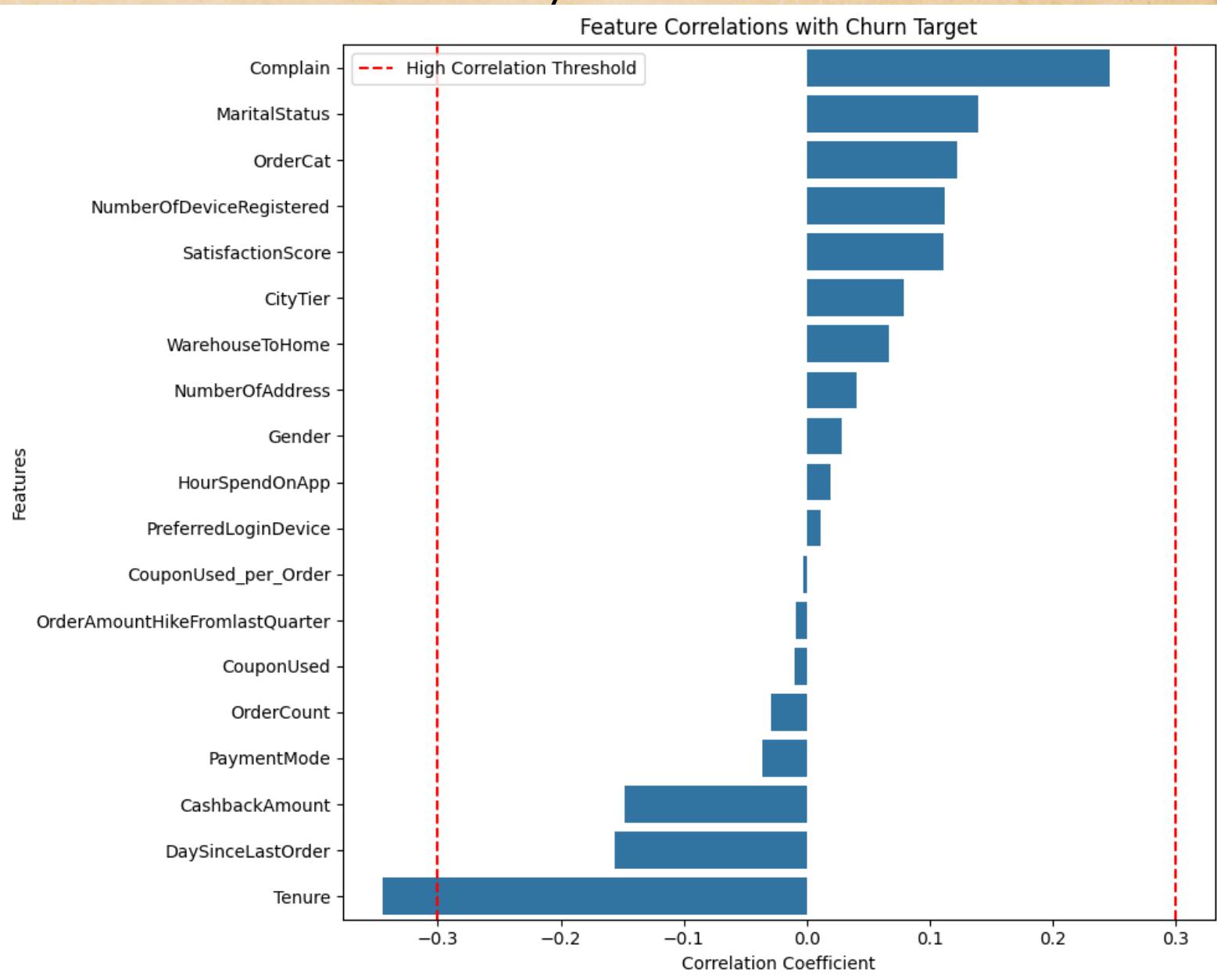
Churn Rate by Complain with 95% CI



Statistical-test for Complain Distribution by Churn Status

- T-test Results:
- T-statistic: 16.45
- P-value: 0.0000
-
- Chi-square Results:
- Chi2-stat: 315.88
- P-value: 0.0000

Correlation Analysis



Feature Selection



- Select features with $\text{corr} \geq 0.05$ and $\text{corr} \leq 0.7$

- Complain
- MaritalStatus
- OrderCat
- NumberOfDeviceRegistered
- SatisfactionScore
- CityTier
- WarehouseToHome
- CashbackAmount
- DaySinceLastOrder
- Tenure

Model Selection



- Logistic Regression

pro: simple, fast, linear boundary, lower variance

cons: pron to underfitting, lower accuracy

- Random Forest

pro: high capacity, high accuracy, robust to overfitting, model explainability, no need for feature scaling, low bias and low variance

cons: training slow, not as good on unbalanced data as XGBoost.

- XGBoost

Pros: usually with highest accuracy of the 3, very good at unbalanced data, distributed computation, low bias and low variance

Evaluation metrics

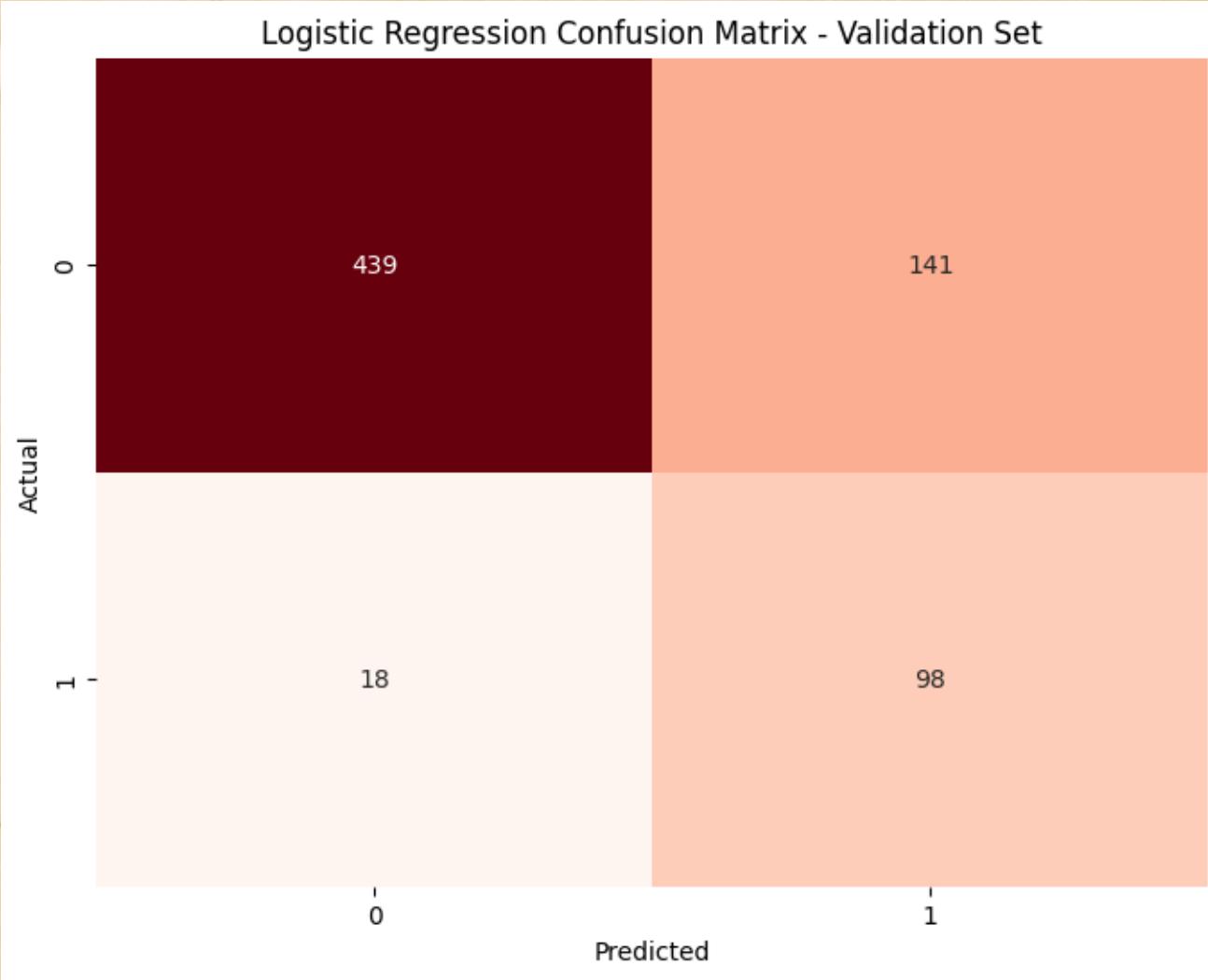


- Accuracy
- Precision
- Recall
- F1
- ROC AUC Score
- Confusion Matrix
-

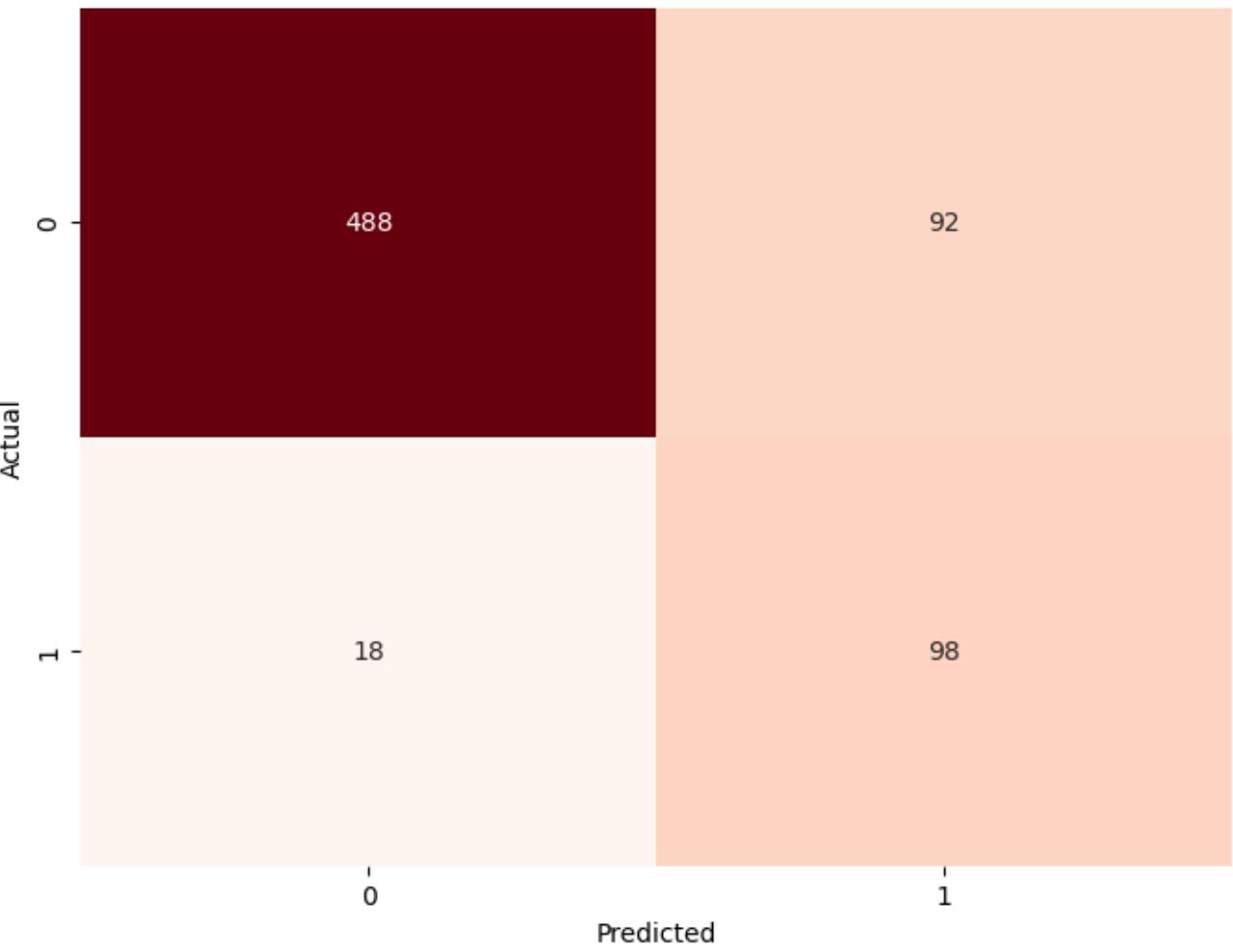
Evaluation Matrics

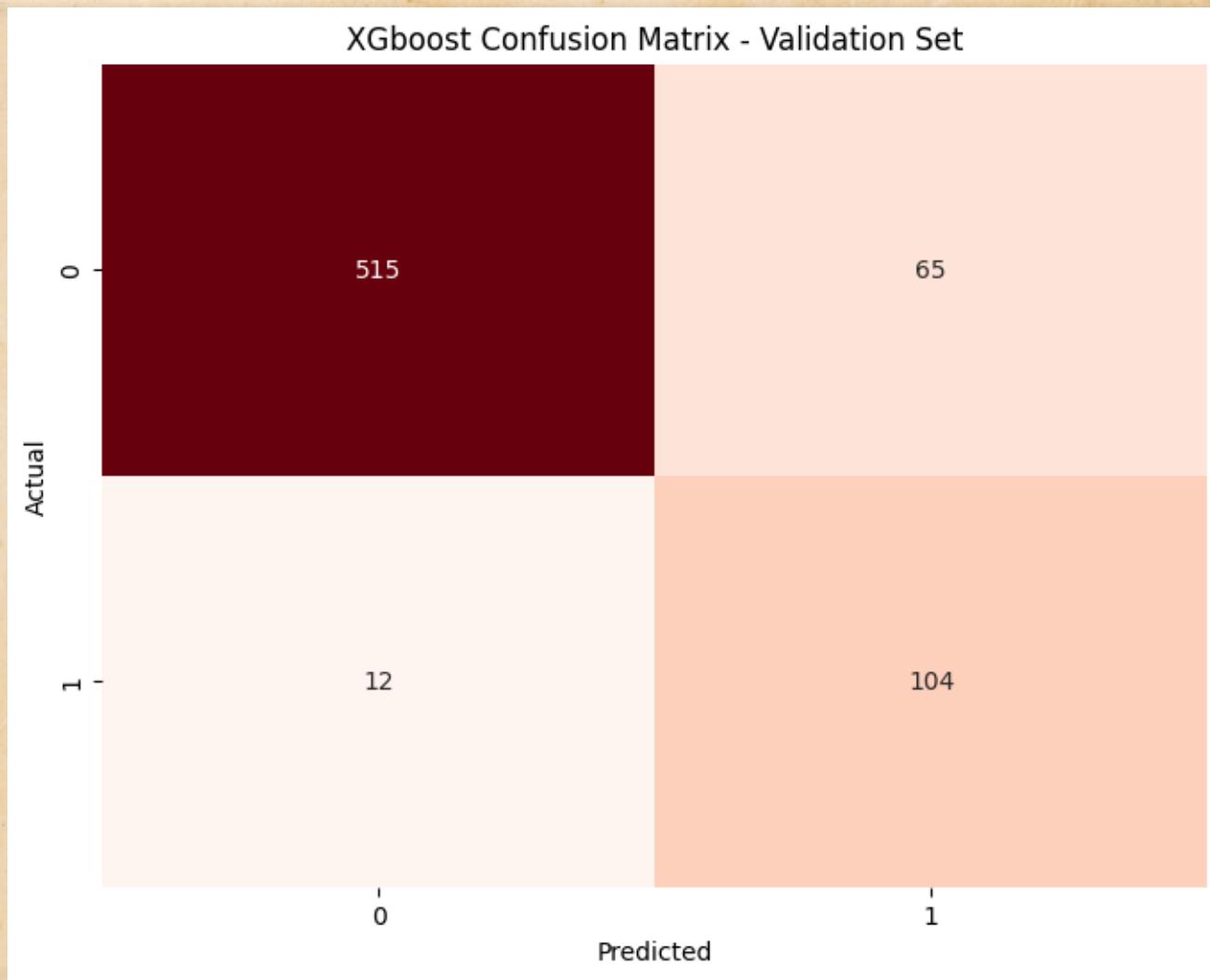
Evaluation set metrics	Logistic Regression	Random Forest	XGBoost
Accuracy	0.7716	0.8420	0.8894
ROC AUC Score	0.8735	0.9146	0.9536
Precision	0.41	0.52	0.62
Recall	0.84	0.84	0.90
F1-score	0.55	0.64	0.73

Logistic Regression Confusion Matrix - Validation Set

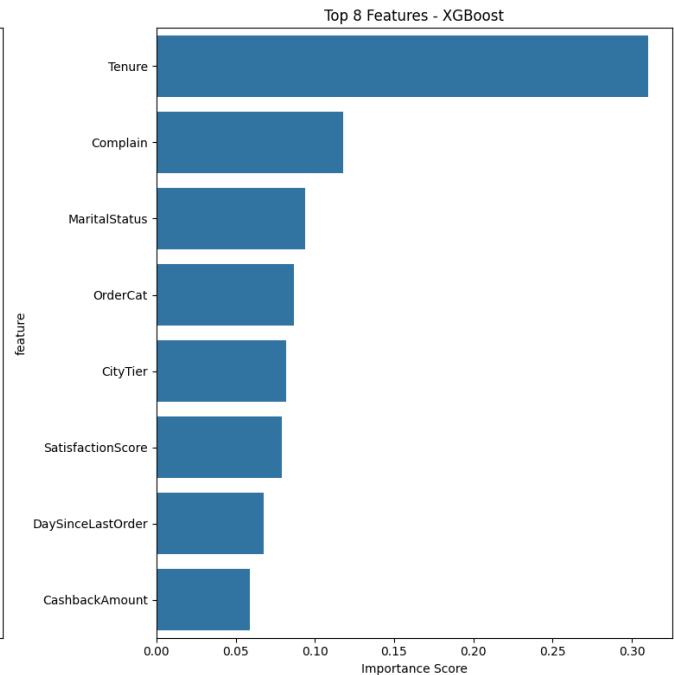
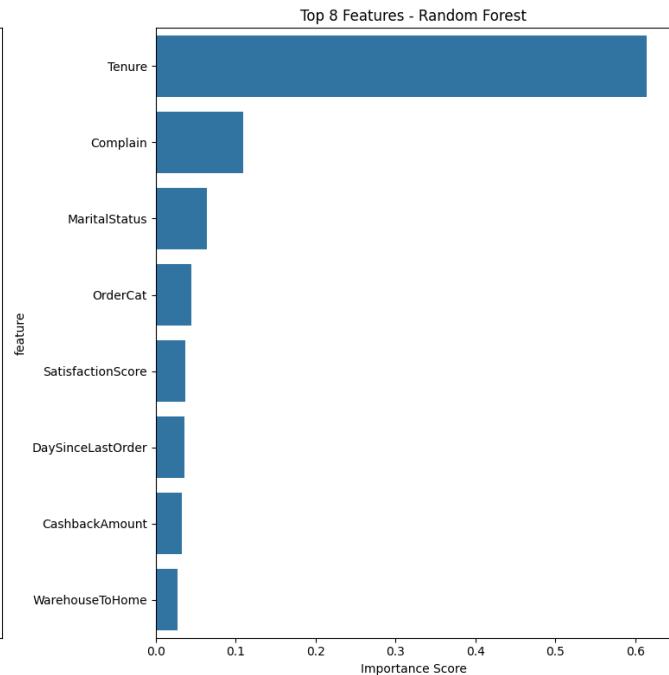
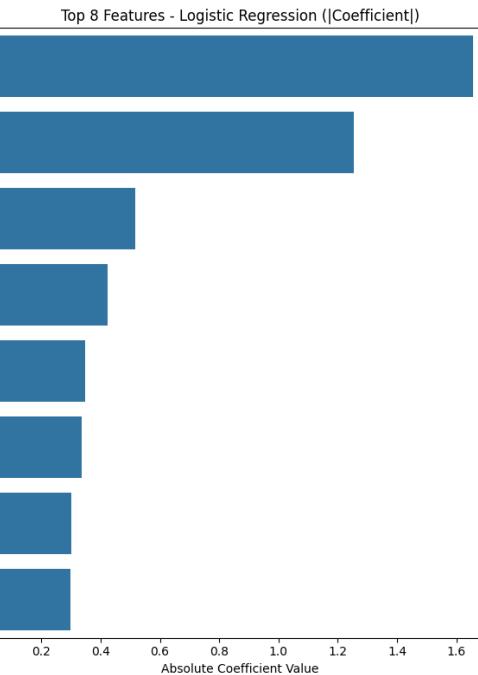


Random Forest Confusion Matrix - Validation Set





Feature Importance Analysis



Evaluation – 5-Fold Cross-Validation

- Logistic Regression

- Accuracy: 0.7608 (+/- 0.0119)
- ROC AUC: 0.8491 (+/- 0.0253)
- Average Precision: 0.6492 (+/- 0.0452)

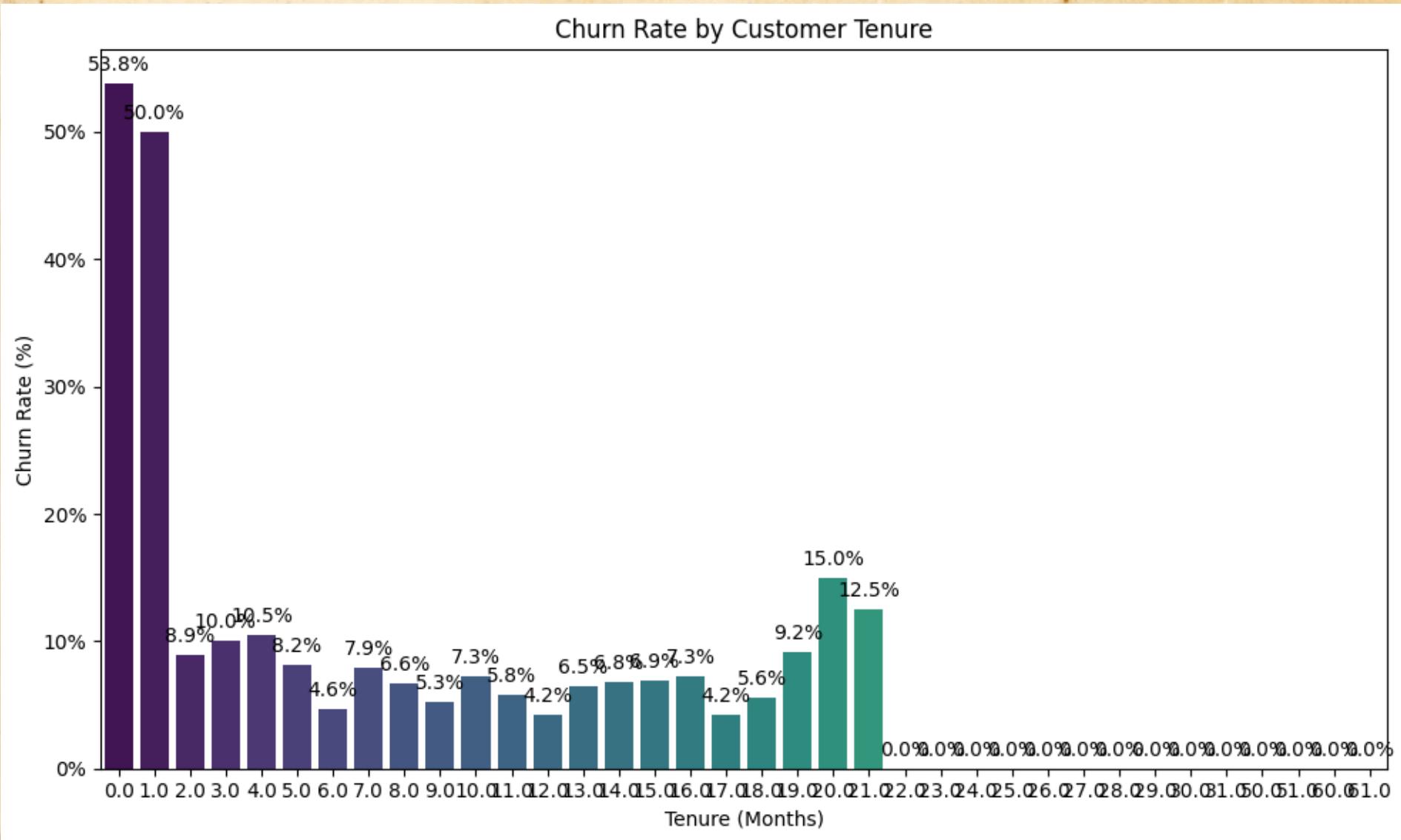
- Random Forest

- Accuracy: 0.8344 (+/- 0.0142)
- ROC AUC: 0.8969 (+/- 0.0255)
- Average Precision: 0.7028 (+/- 0.0425)

- XGBoost

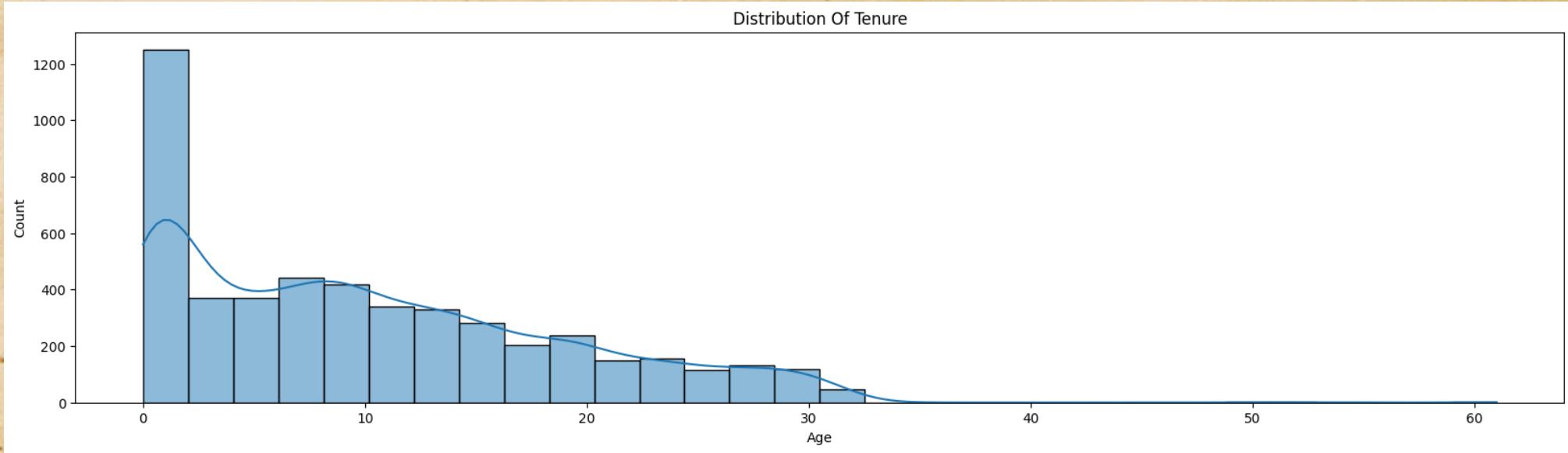
- Accuracy: 0.9080 (+/- 0.0097)
- ROC AUC: 0.9498 (+/- 0.0173)
- Average Precision: 0.8137 (+/- 0.0383)

1. More Tenure Analysis & Related Predict Analysis





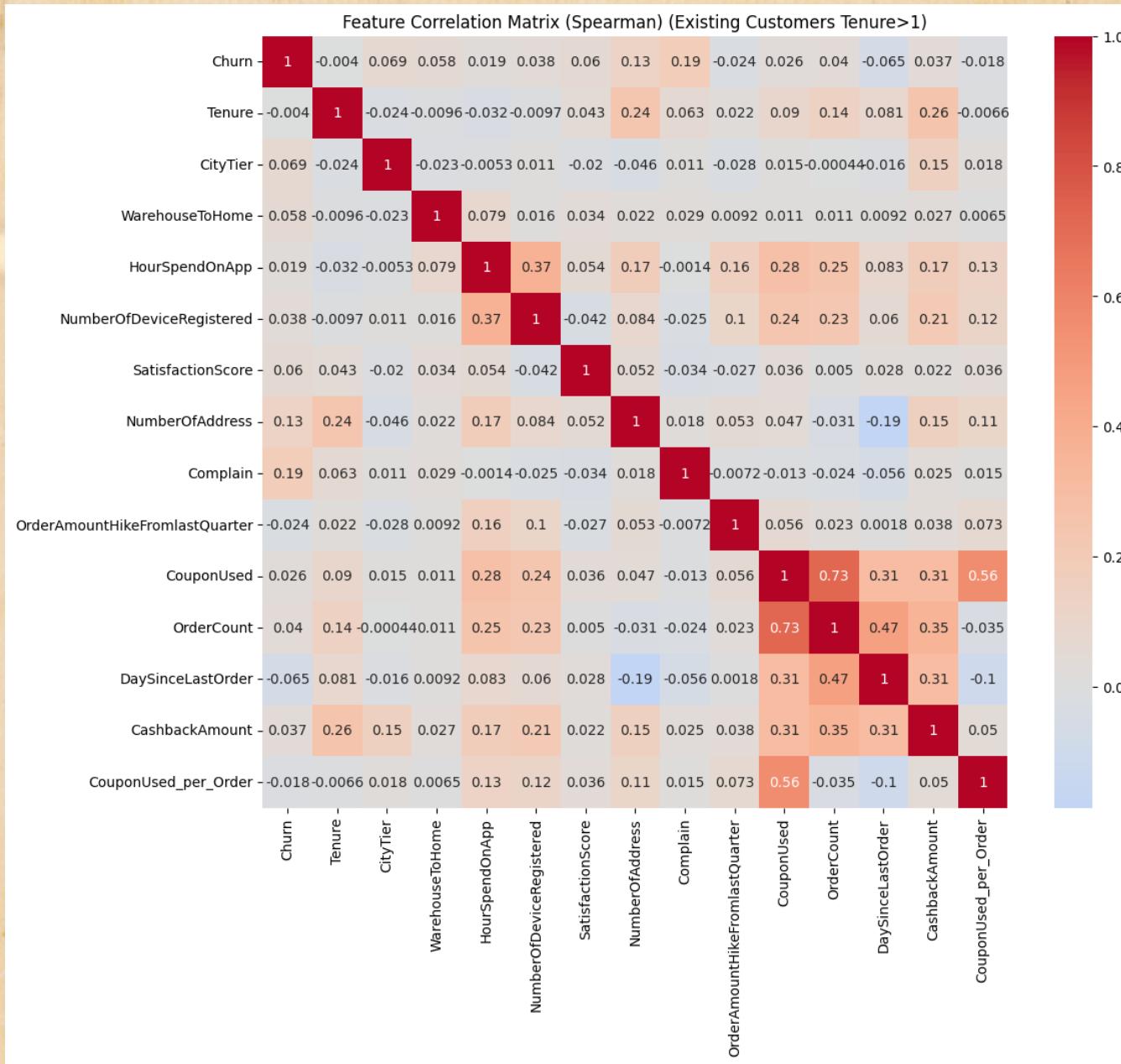
Distribution Of Tenure

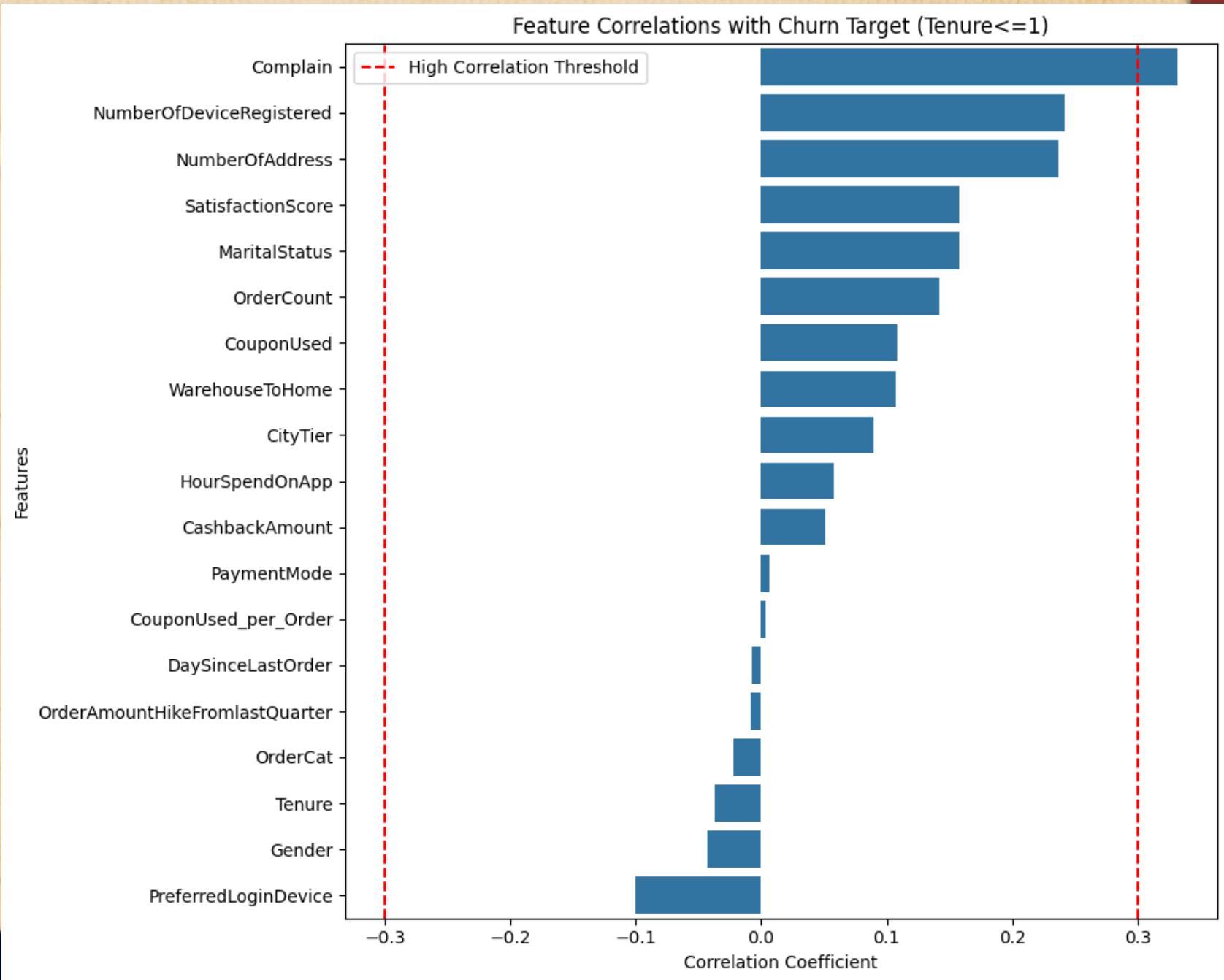


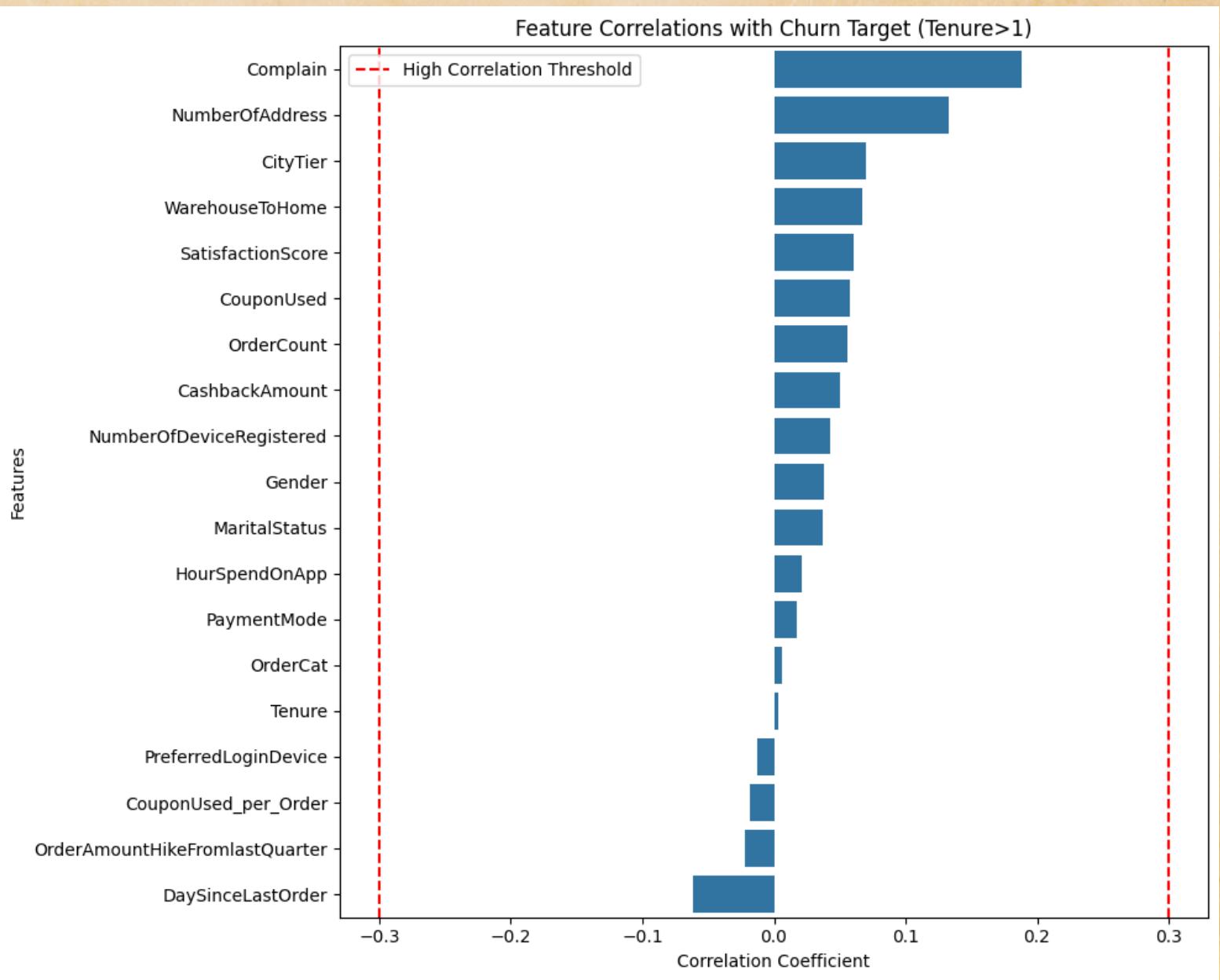
2 models

- 2 dataset
- Tenure ≤ 1
- Tenure > 1









2. Actionable Inside and recommendations

- Churn prediction with probability for customer retention
- Different customer retention strategy for new customers($\text{tenure} \leq 1$) and existing ones($\text{tenure} > 1$)
- Consider remove sensitive features like Sex and MaritalStatus for AI data regulation related DEI
- Collect data with actionable features like sales/marketing related ones, promotion by email, samples sent, visual/phone communication logs... which can be experimented to lower the churn

Questions?



Thanks





This work is licensed under
a Creative Commons Attribution-ShareAlike 3.0 Unported License.
It makes use of the works of
Kelly Loves Whales and Nick Merritt.