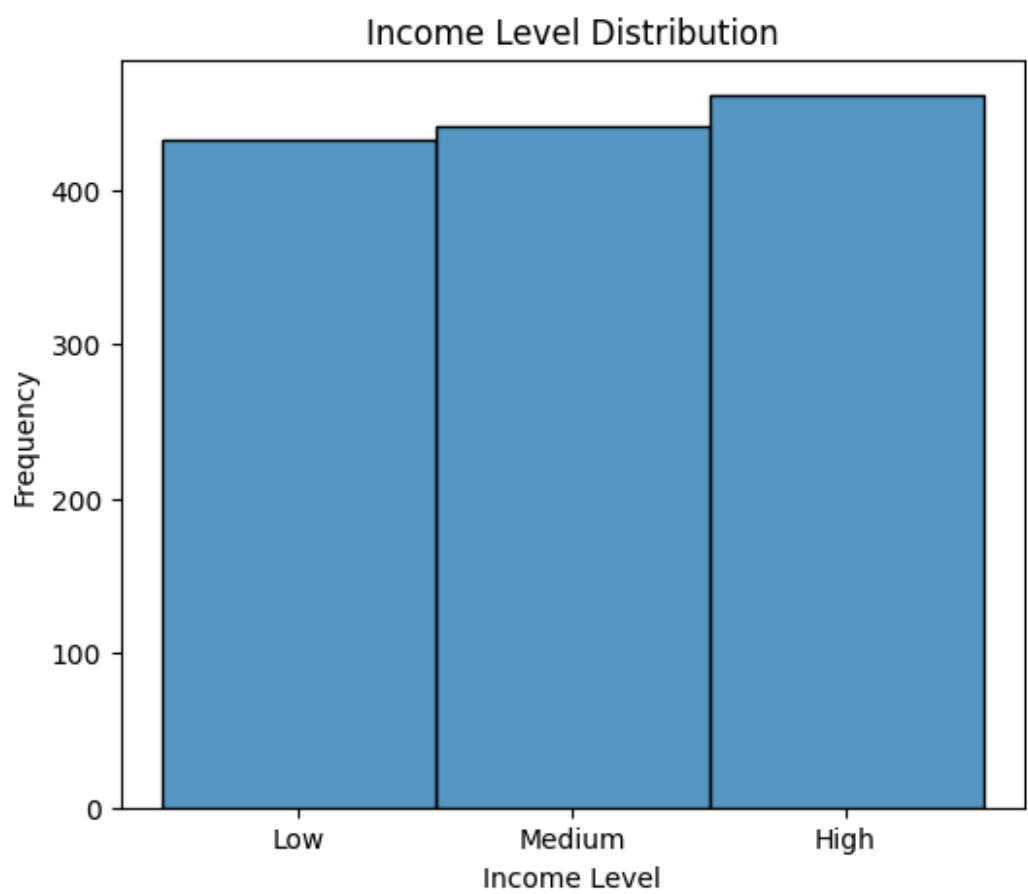


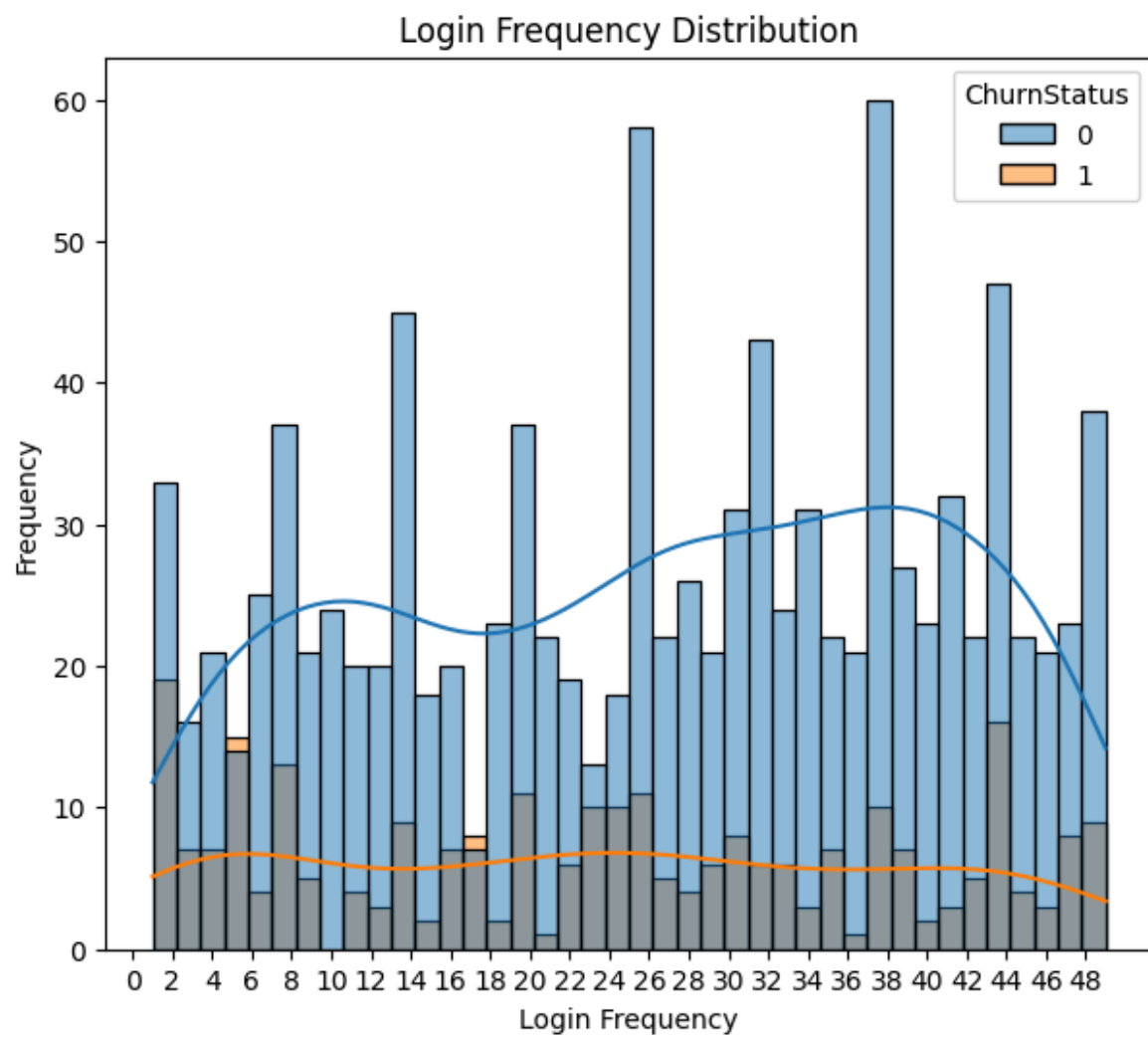
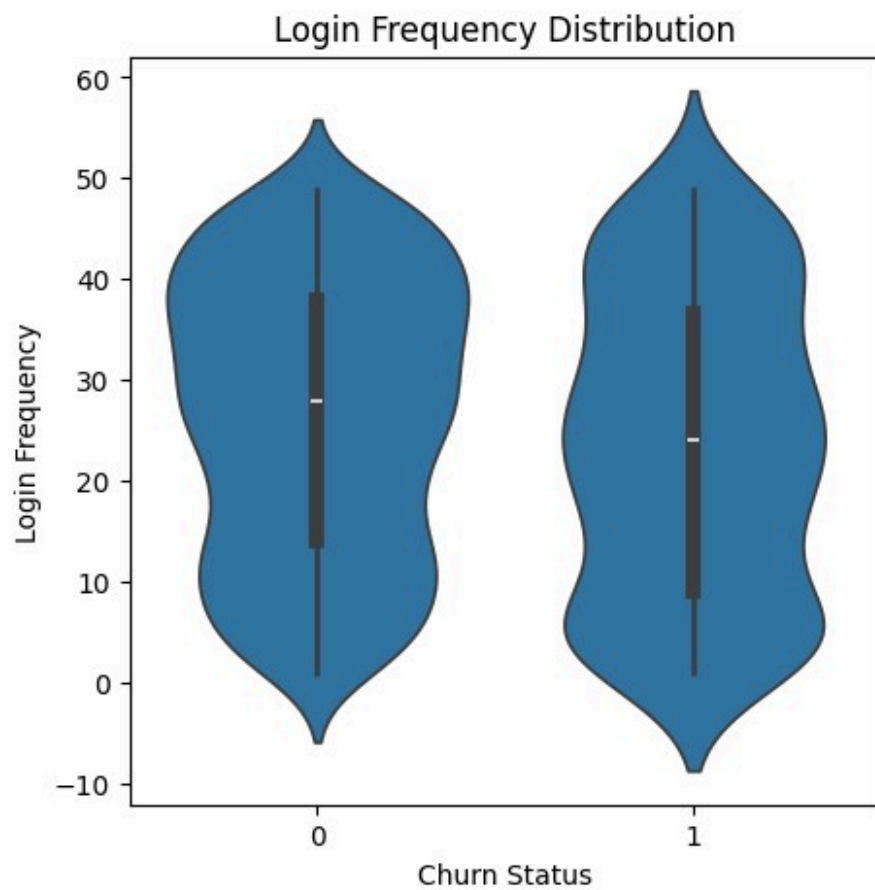
# **AYUSH SAIN DATA ANALYSIS REPORT**

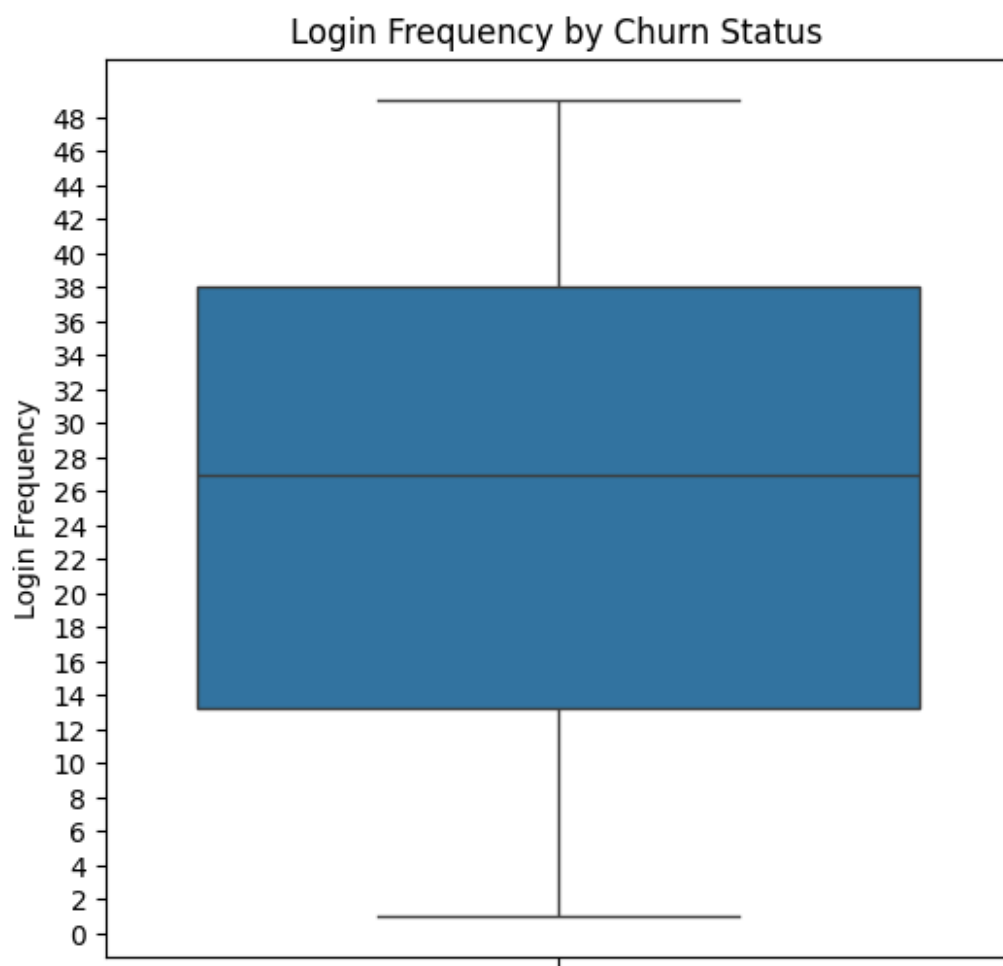
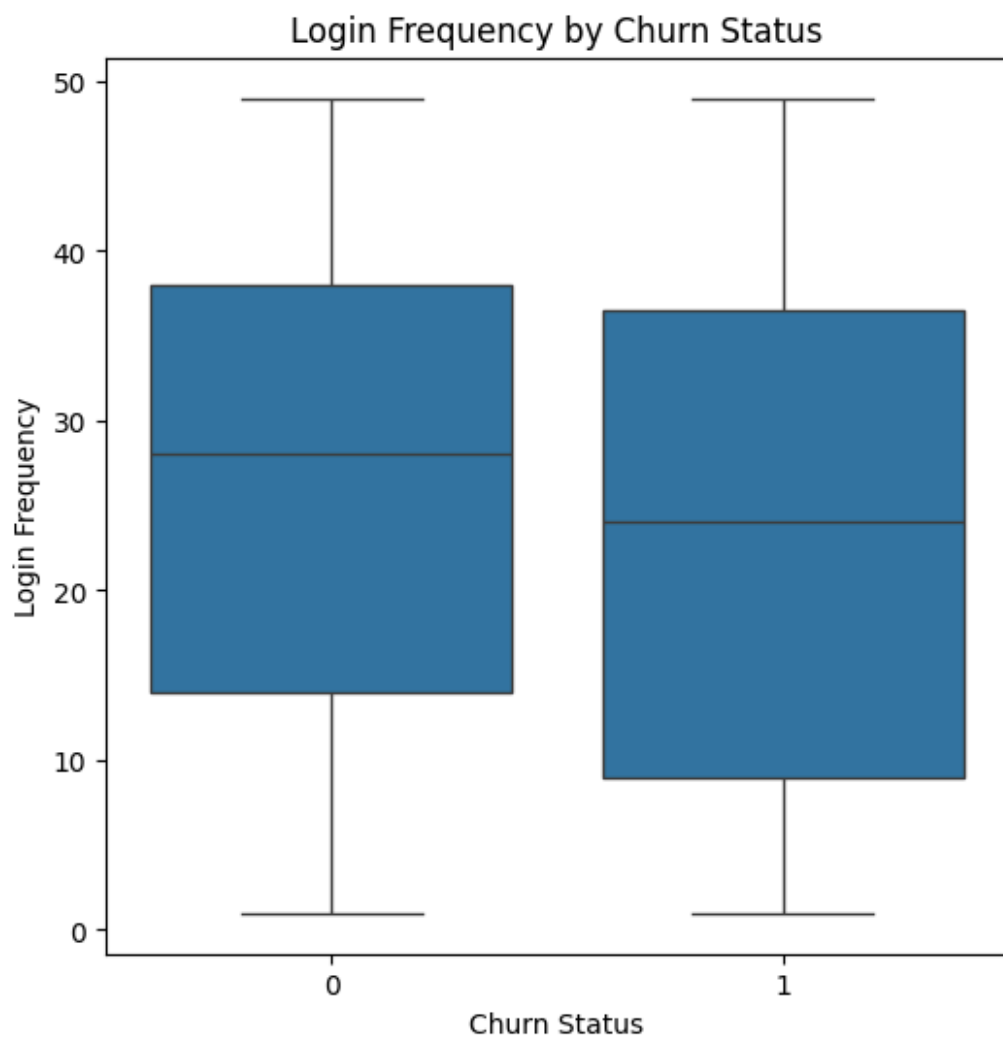
# Data Preprocessing and Cleaning Report

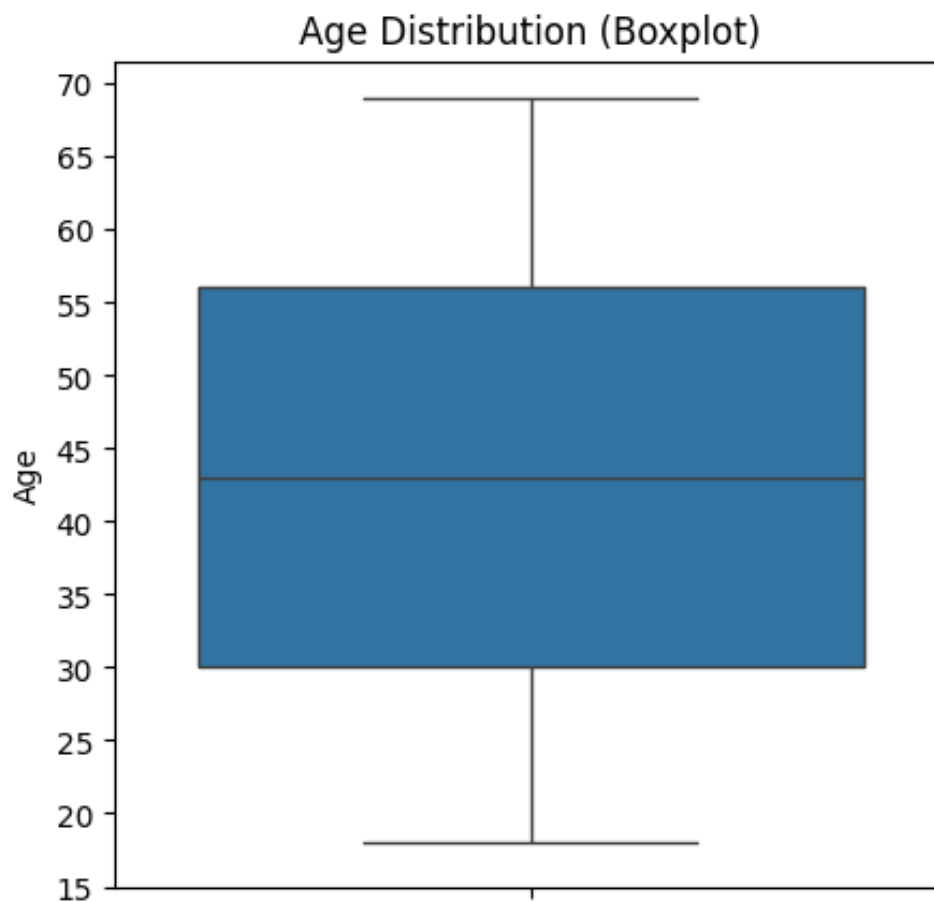
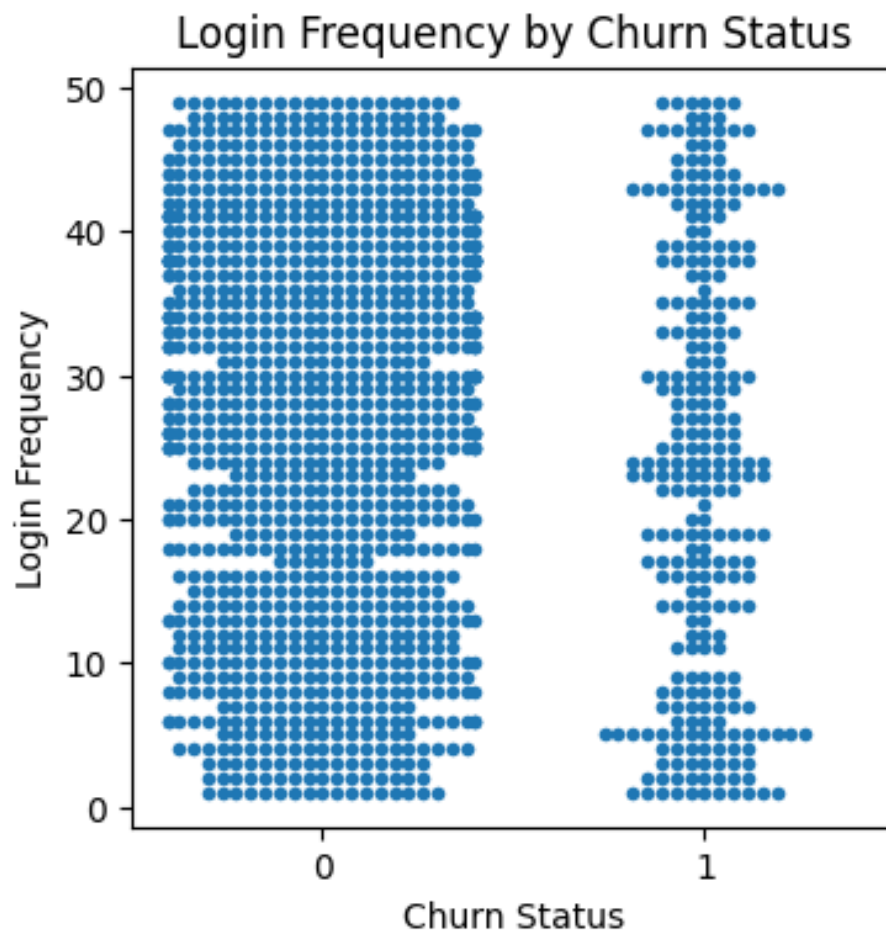
## 1. Introduction

This section of the report details the systematic steps taken to prepare the customer churn dataset for analysis and model building. The dataset comprises multiple features related to customer demographics, behavior, and service usage. The primary goal of this preprocessing phase was to ensure the data is clean, consistent, and in a format suitable for machine learning tasks, particularly churn prediction.

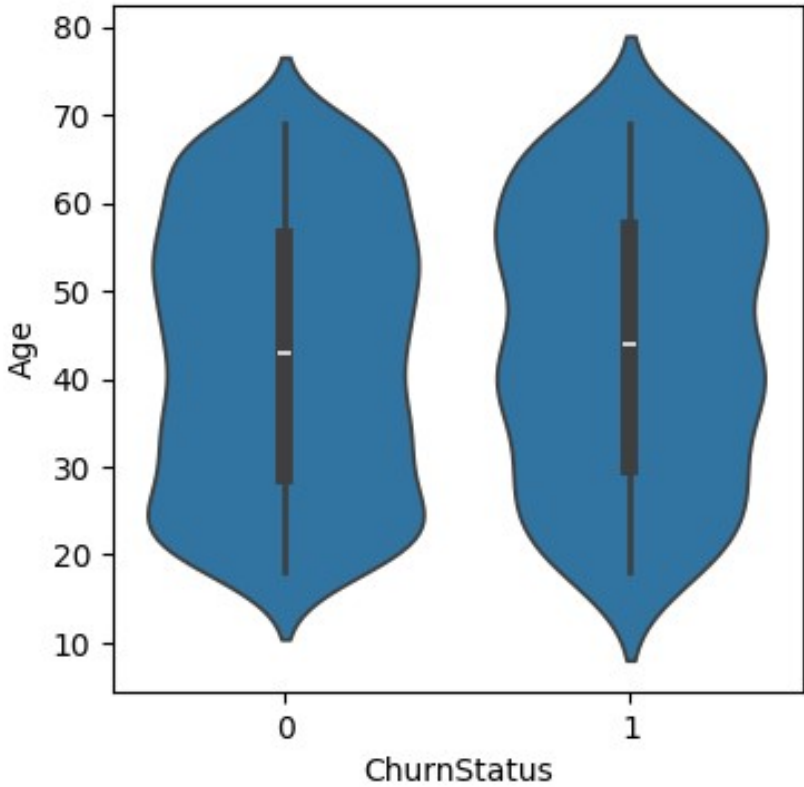
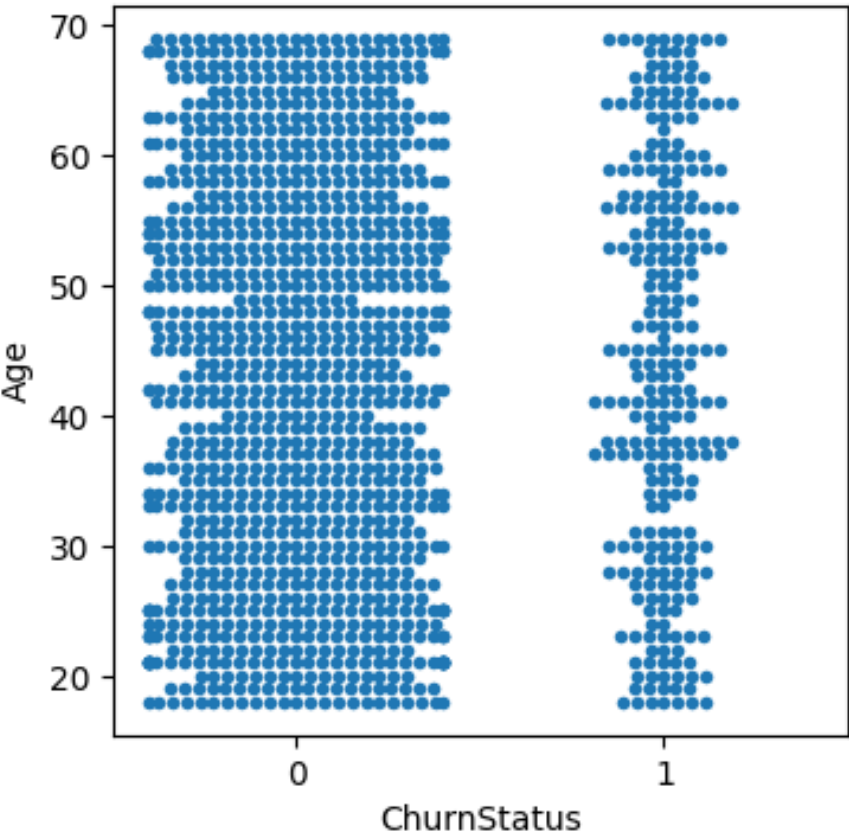


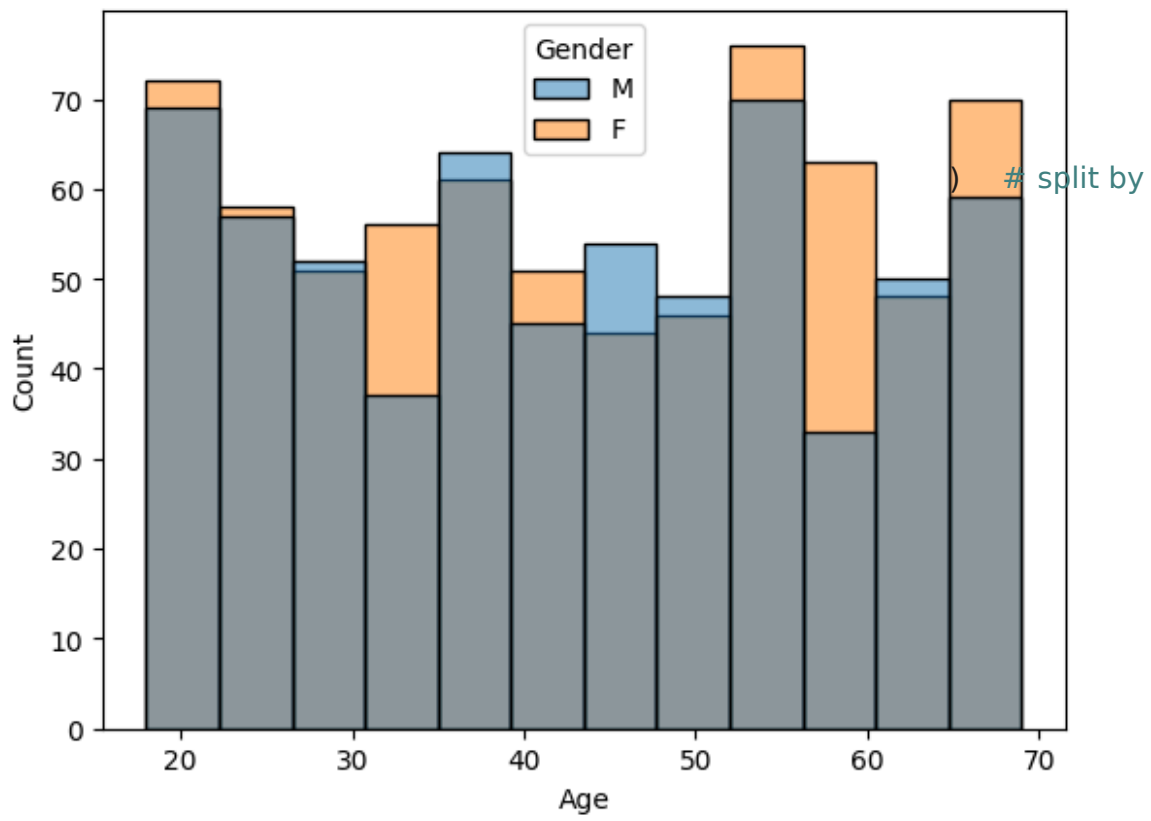




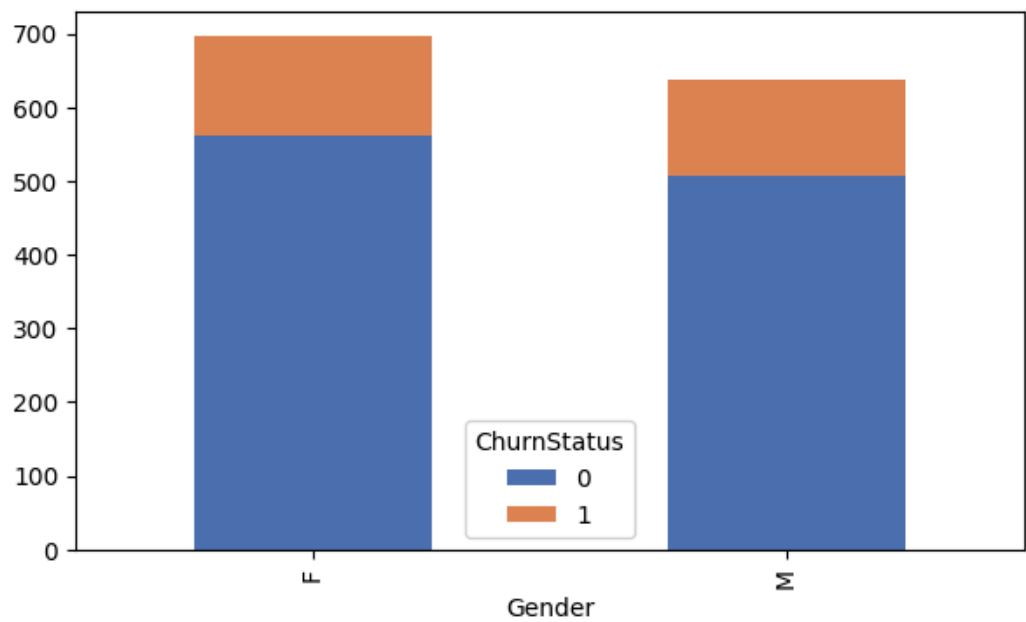
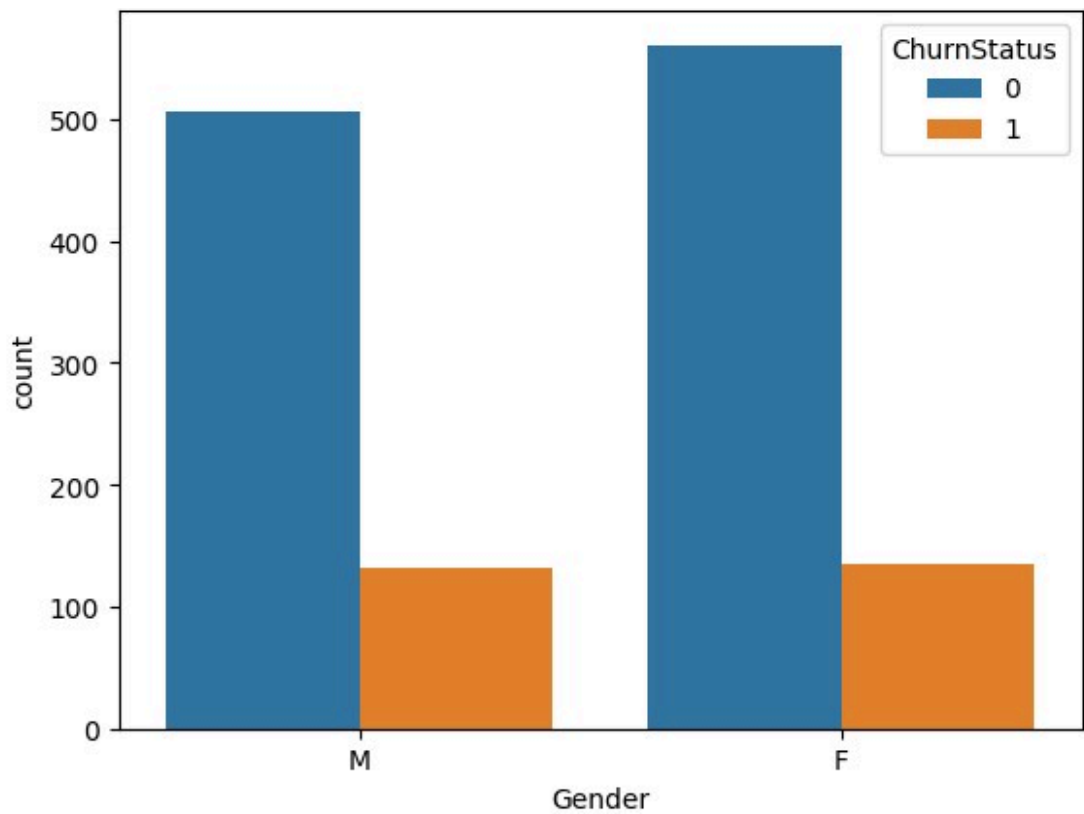


bi variate analysis using the target variable

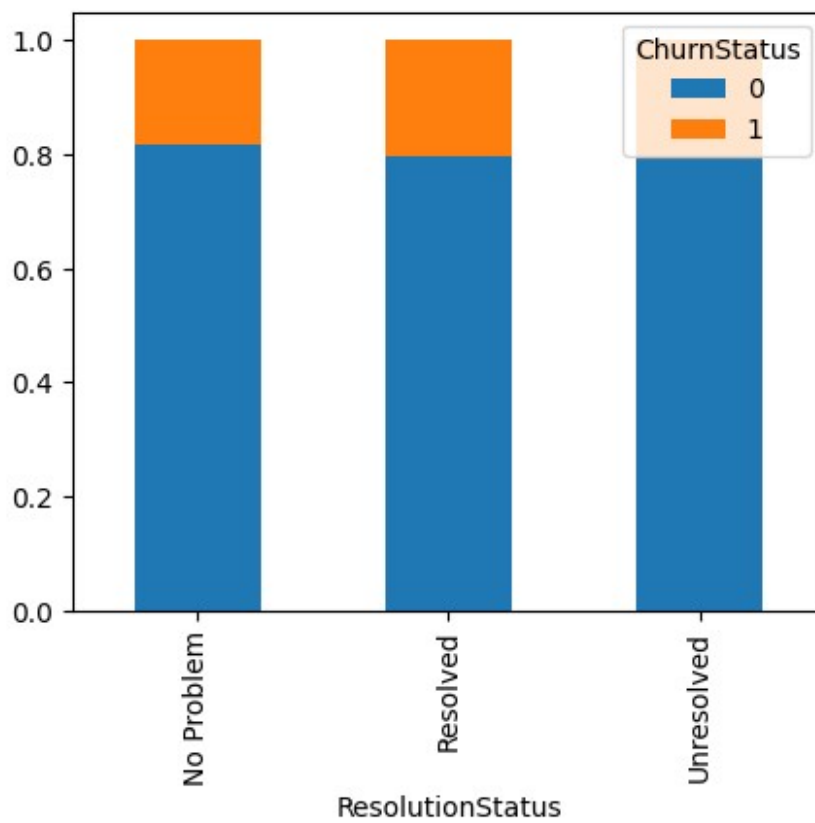
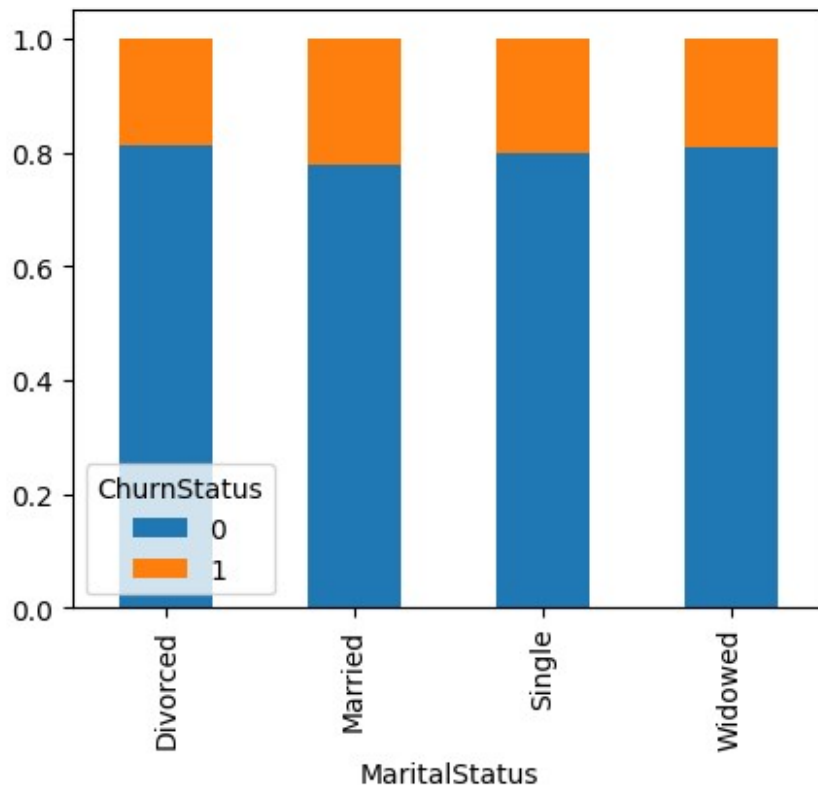


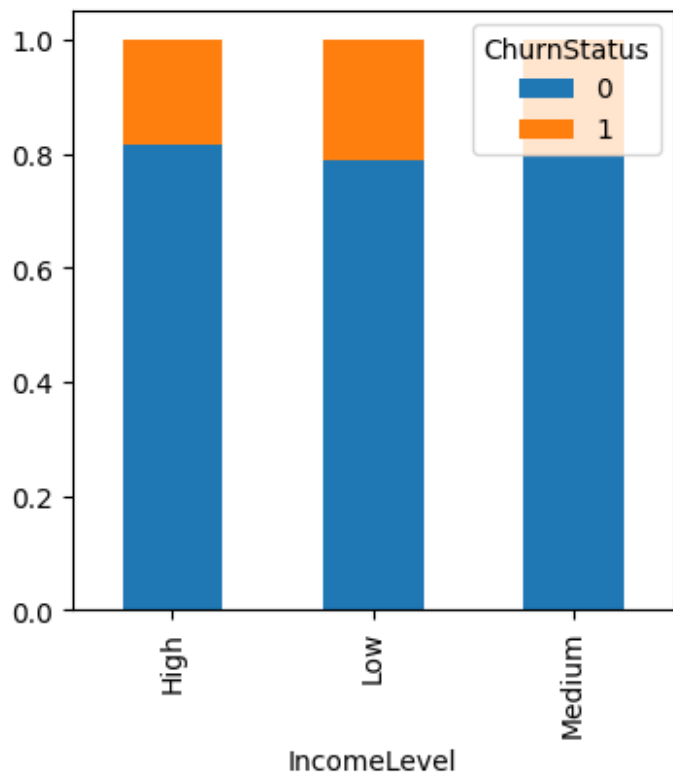


categorical vs target variable

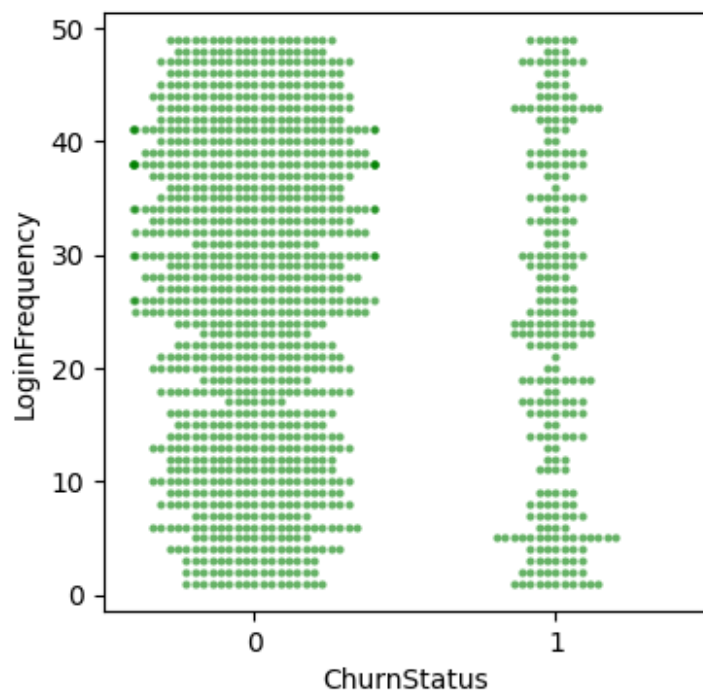


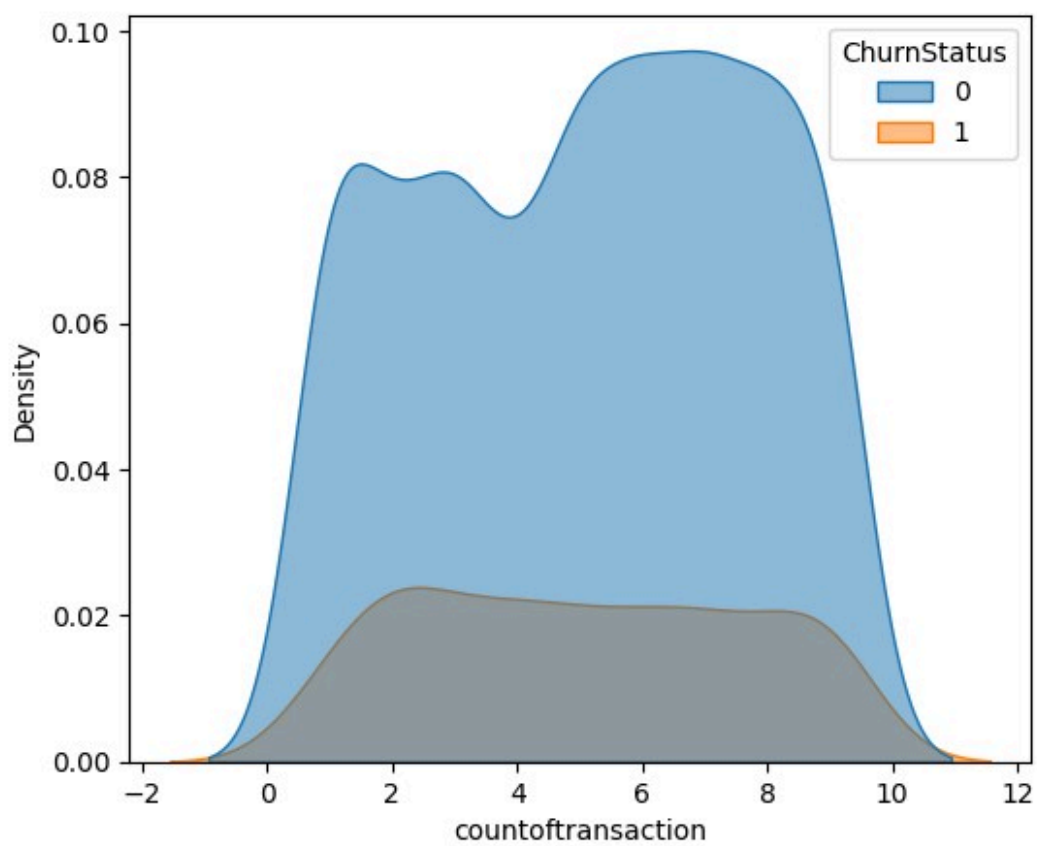
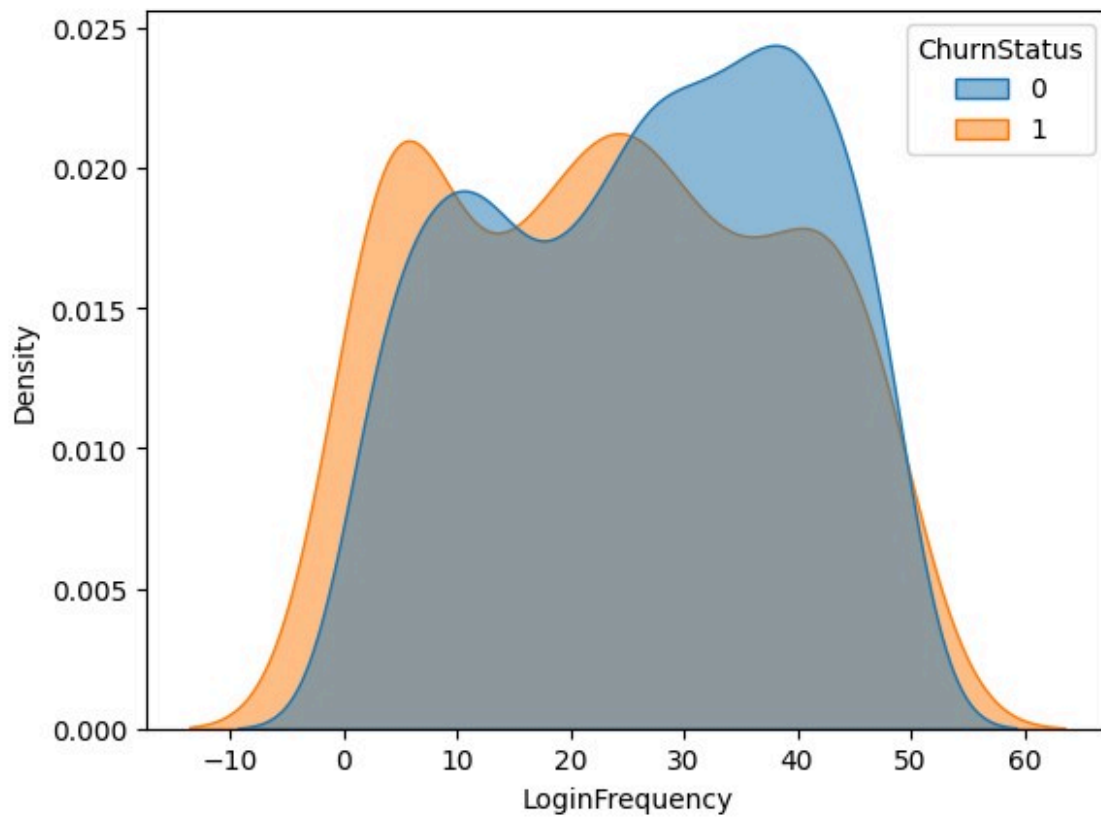


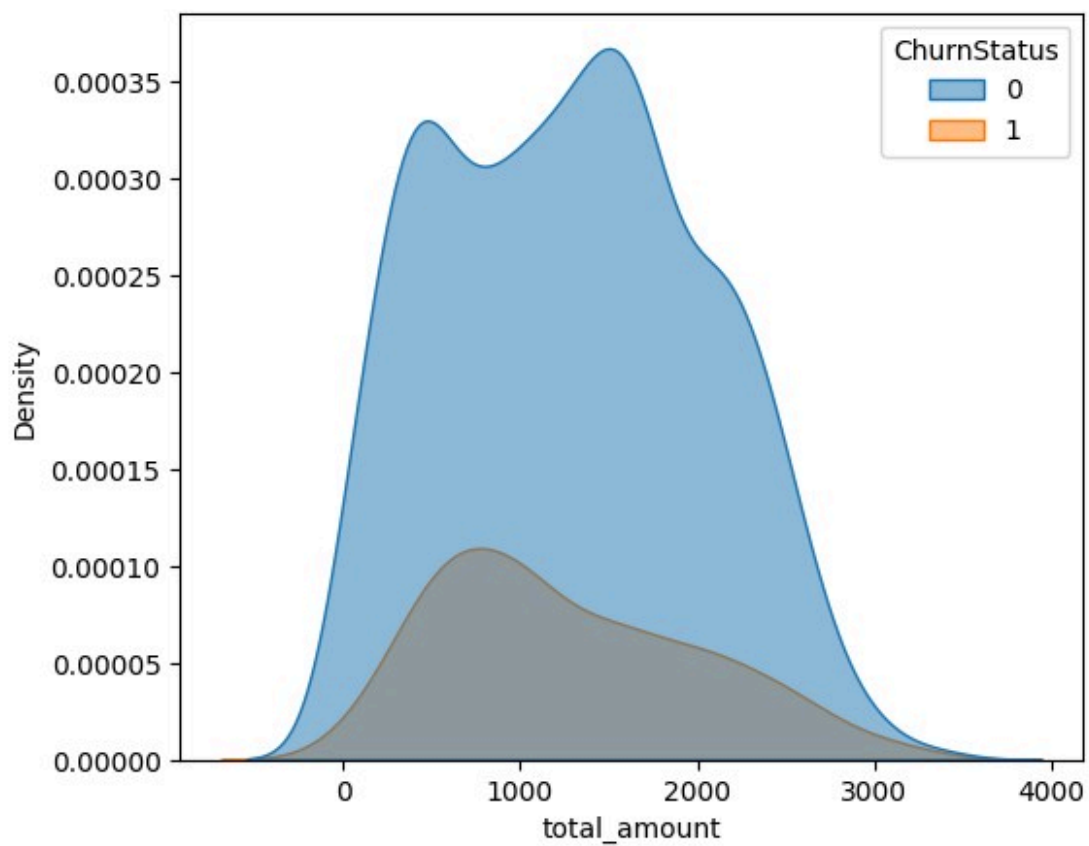
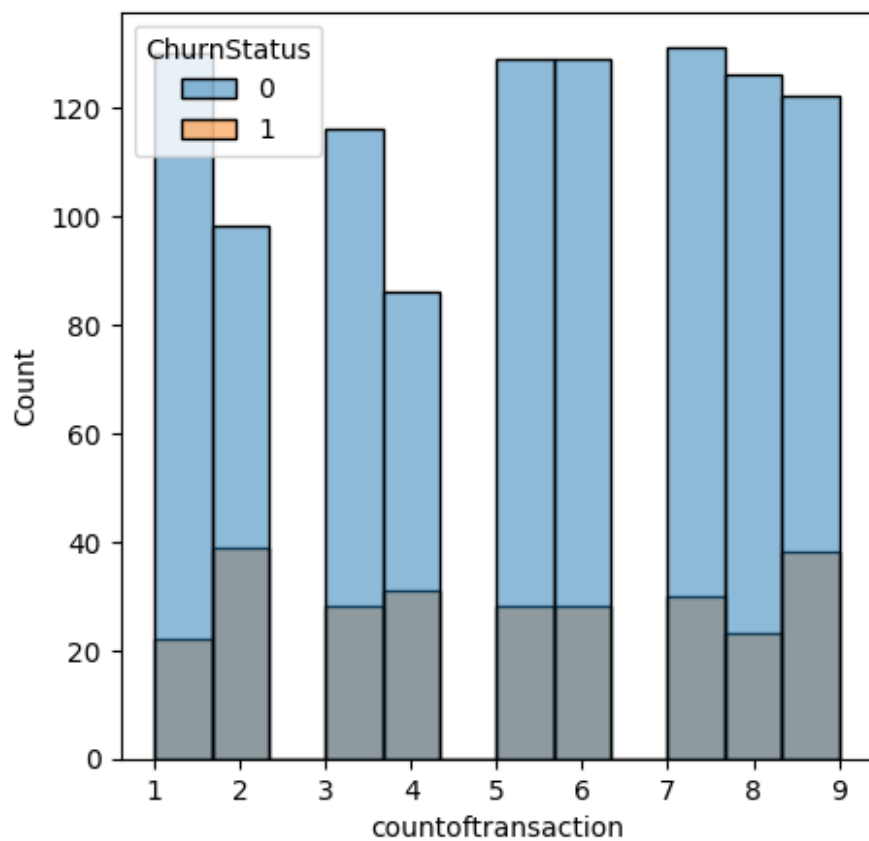


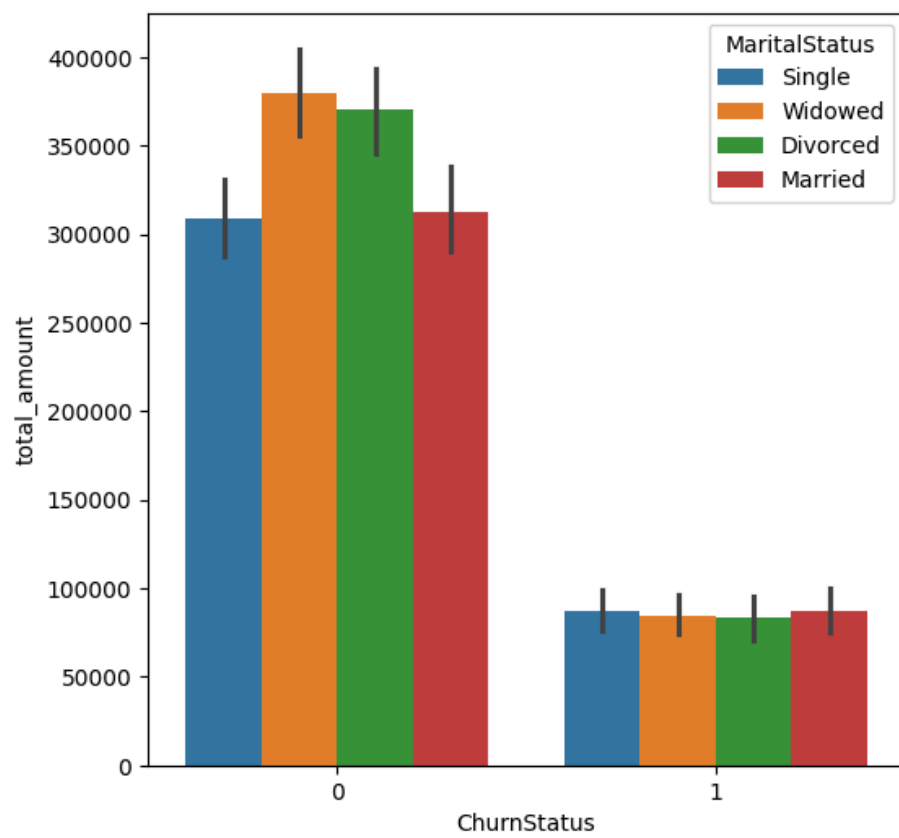
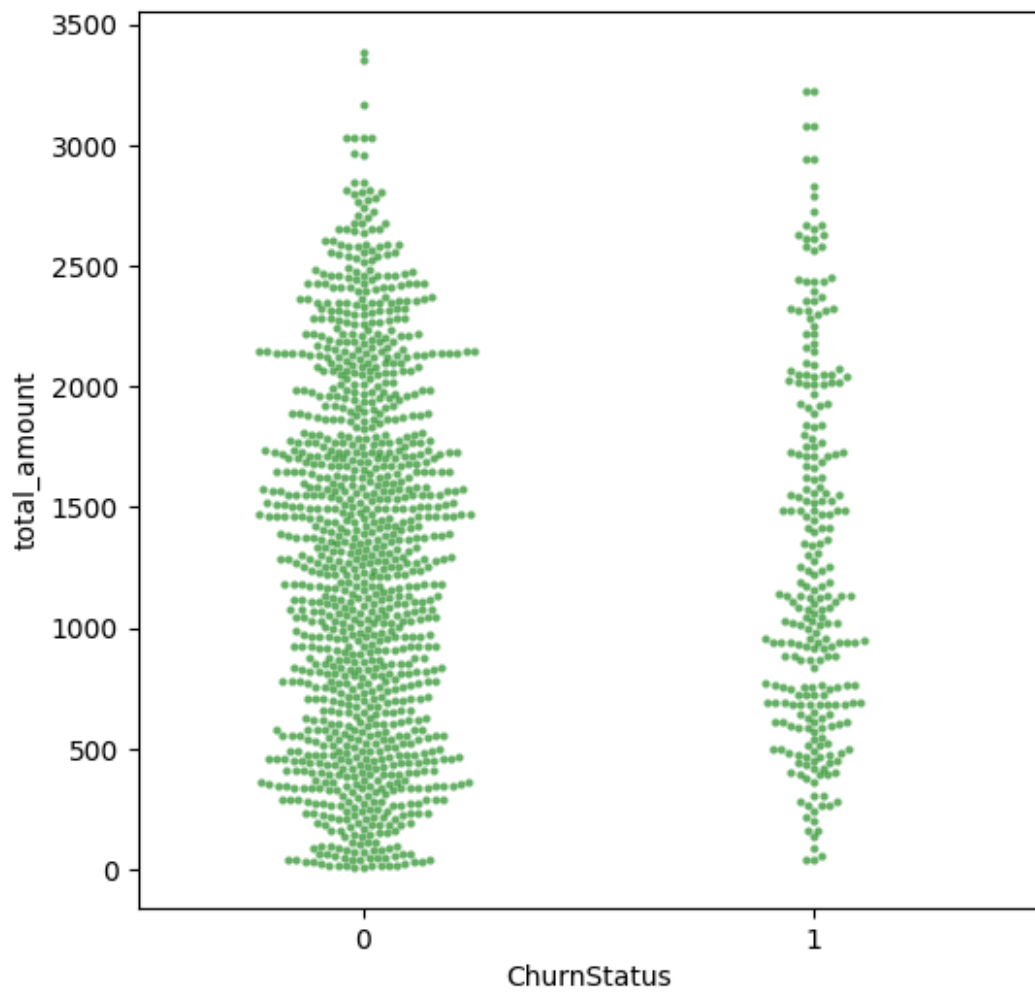


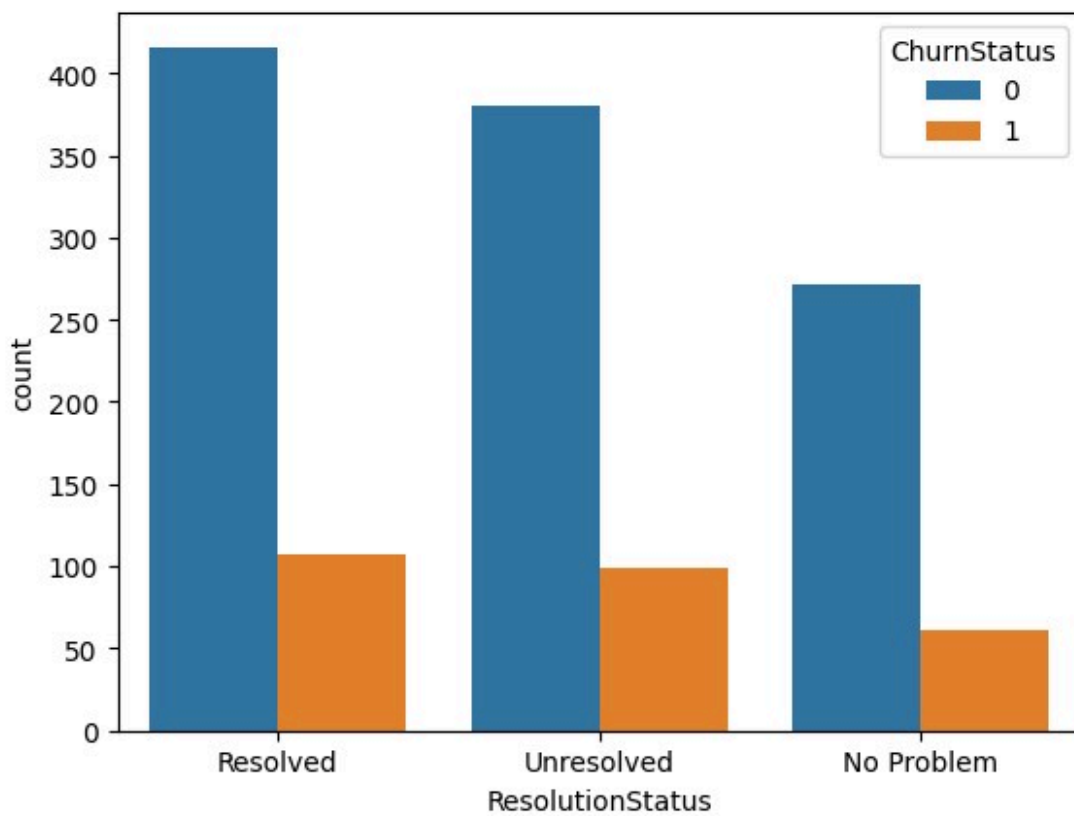
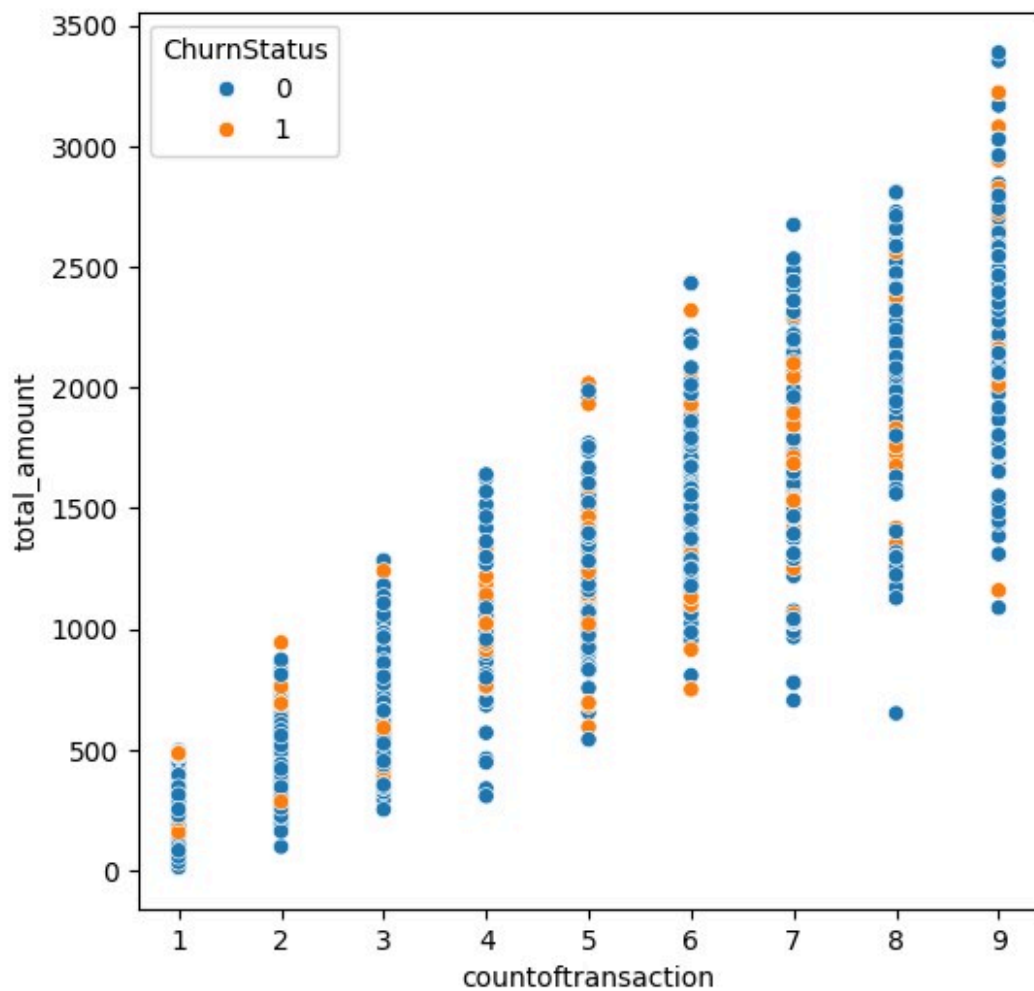
numerical vs target

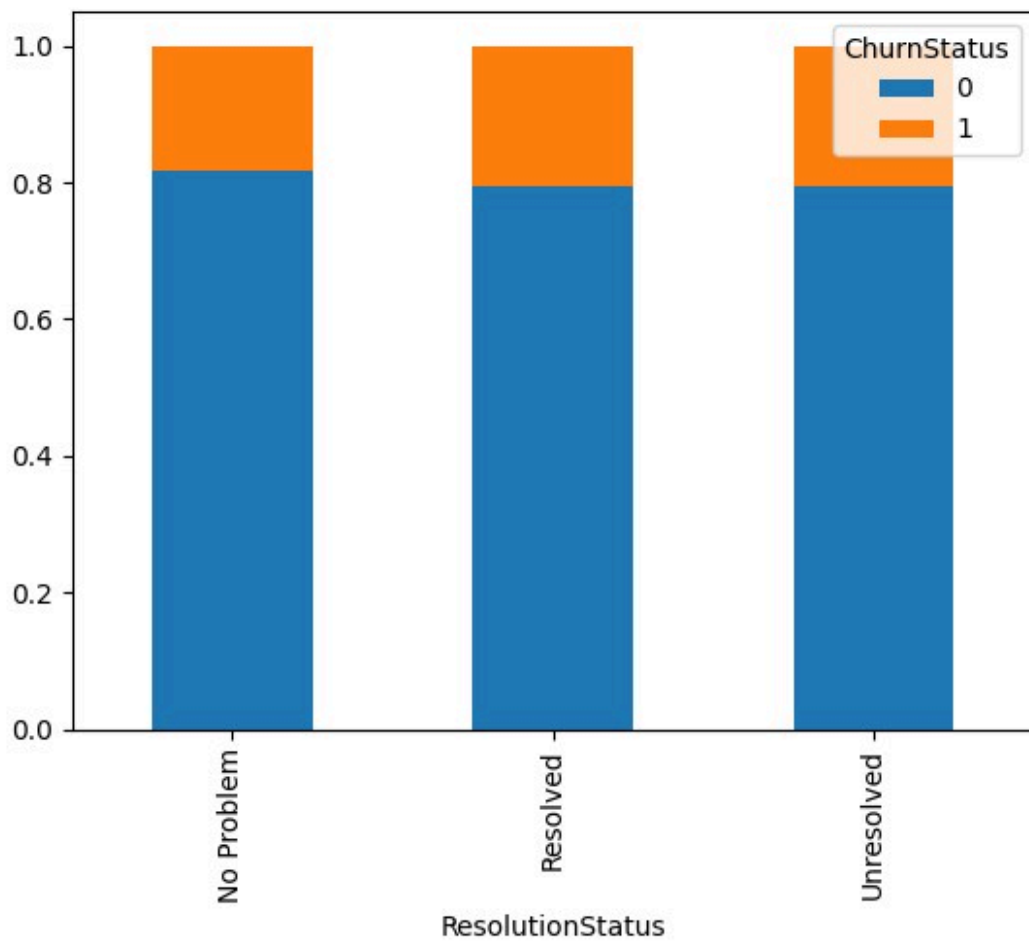


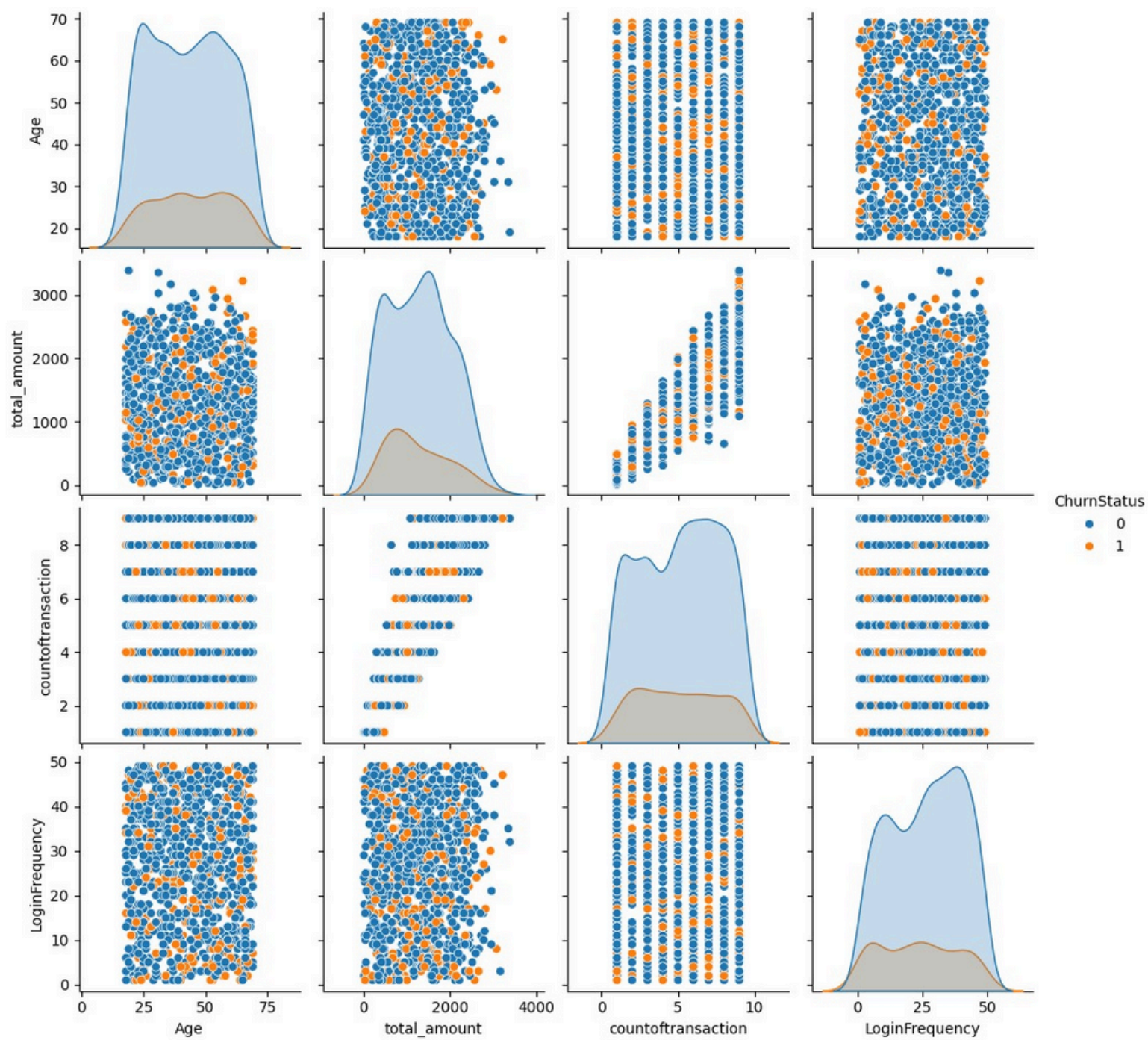




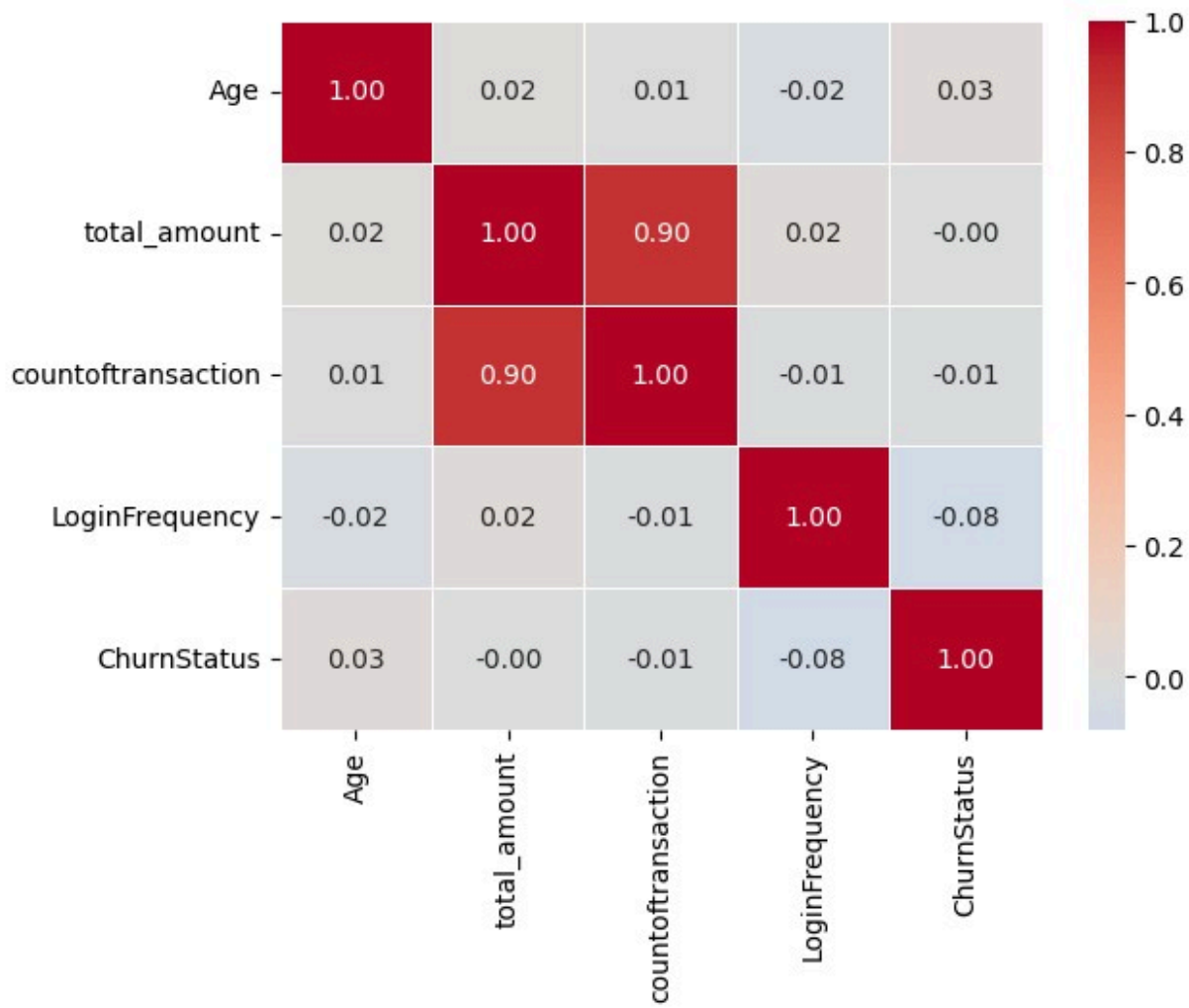
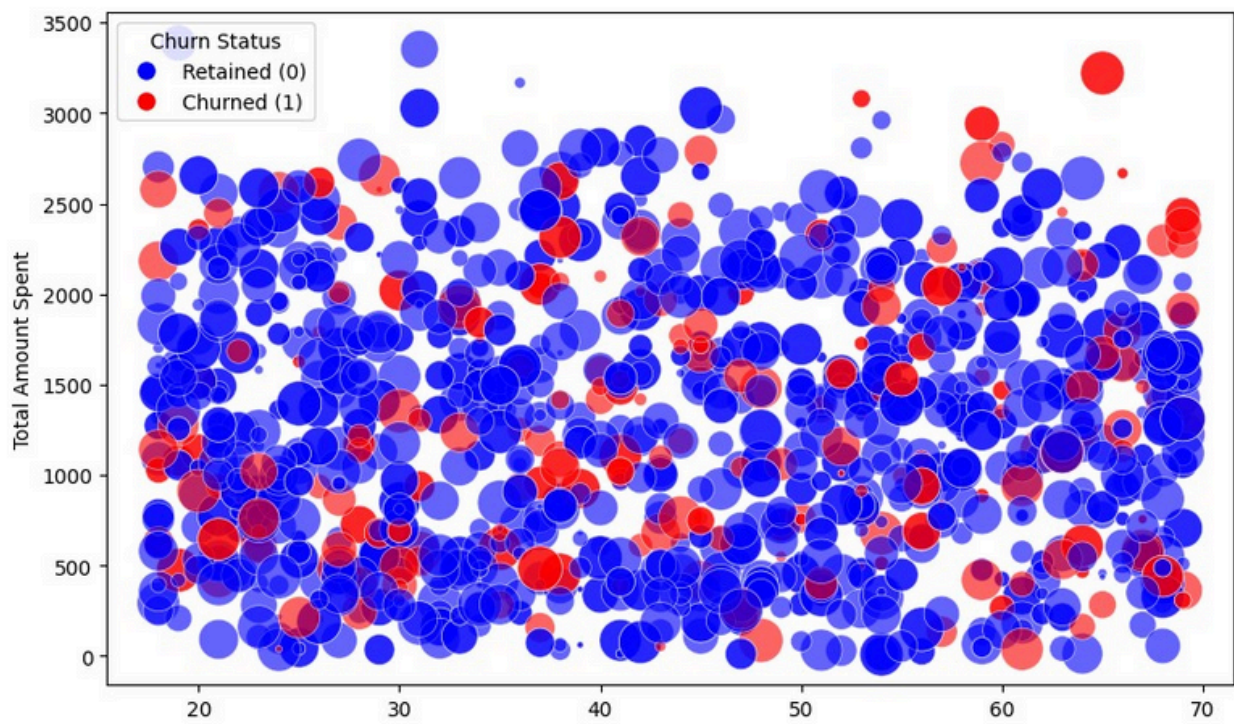












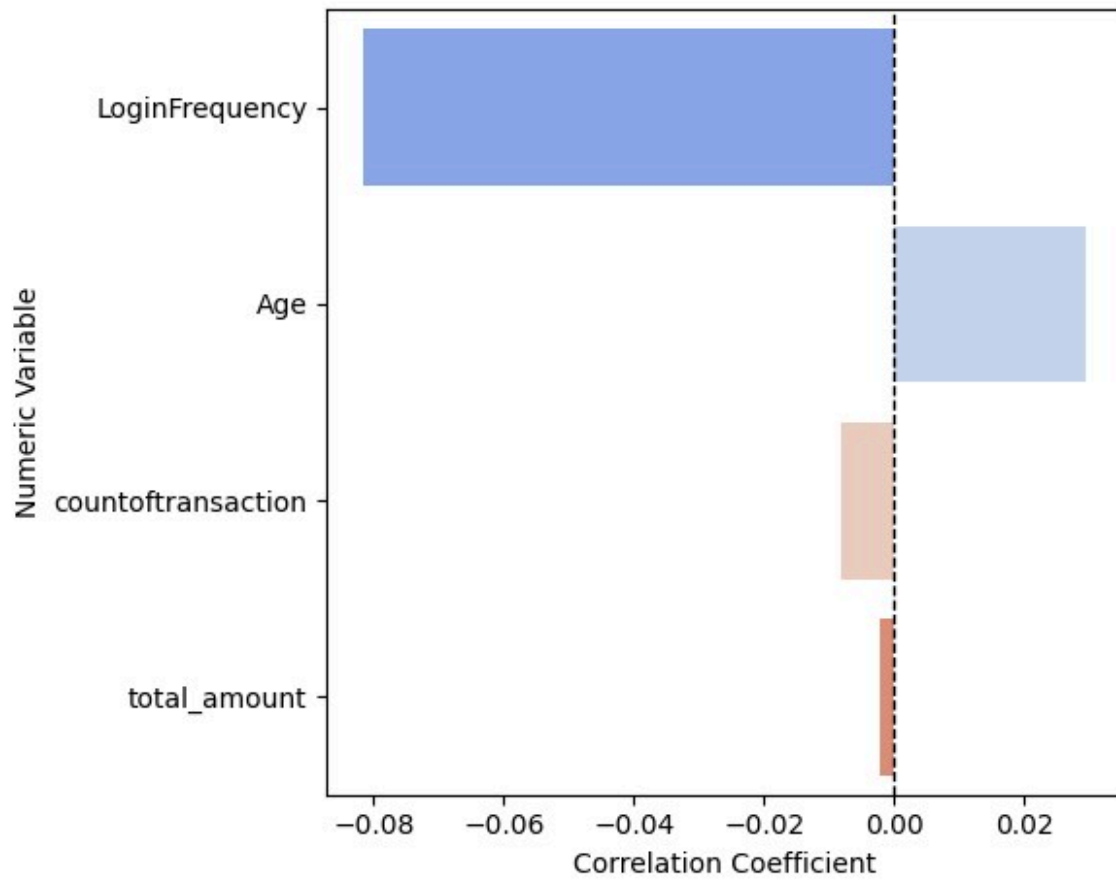
## **2. Handling Missing Values**

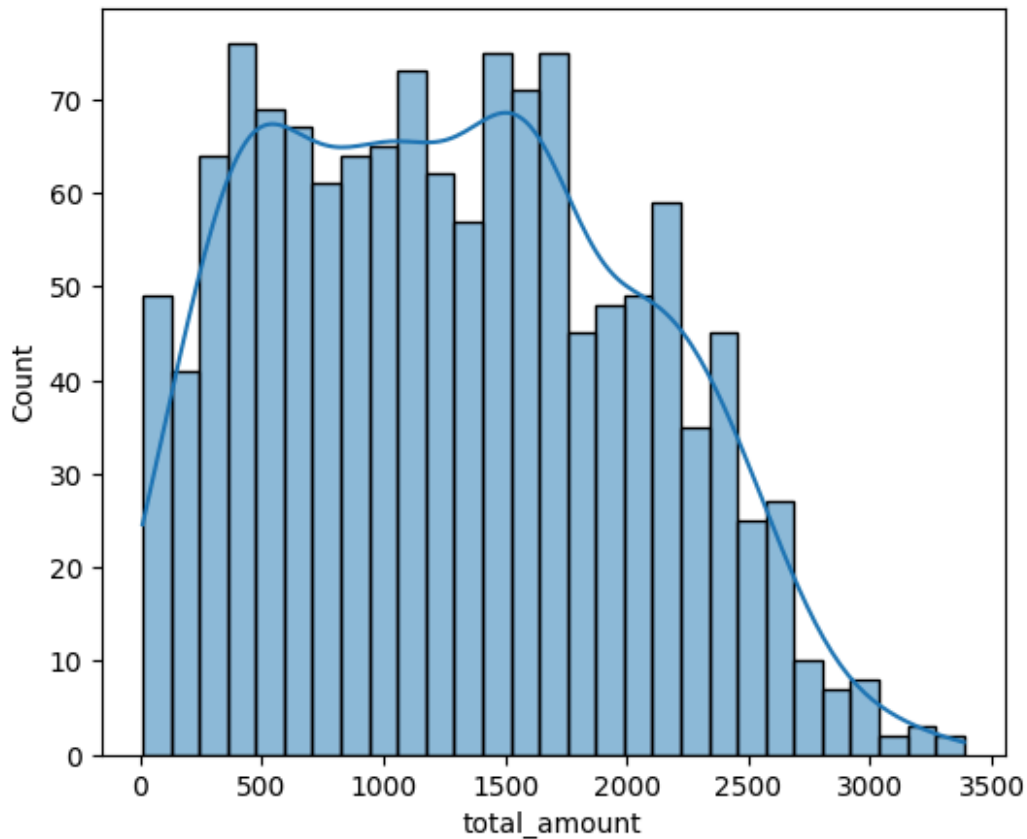
Missing values were identified across both numerical and categorical features. Numerical columns such as Age, IncomeLevel, and LoginFrequency had minimal missing data, which were imputed using the median value to preserve distribution and avoid skewness. For categorical columns like Gender and ResolutionStatus, missing values were imputed using a new category labelled "Unknown" to avoid dropping informative instances. A binary flag indicating the presence of missing values was also created where necessary, to allow the model to capture any potential predictive signal associated with missingness. Records with missing target labels (ChurnStatus) were dropped, as they cannot contribute to supervised learning.

## **3. Outlier Detection and Treatment**

Outliers were assessed using the interquartile range (IQR) method and visual inspection via boxplots. Key numeric features such as total\_amount, countoftransaction, and IncomeLevel exhibited significant right-skewed distributions, indicating the presence of high-end outliers. Rather than removing these values outright—which could eliminate high-value customers—extreme values were capped at the 99th percentile. Additionally, log transformation was applied to reduce skewness and compress scale differences for better model performance. Erroneous values in Age (e.g., below 0 or above 120) were considered data entry errors and removed.

Point-Biserial Correlation: Numeric Variables vs. Churn

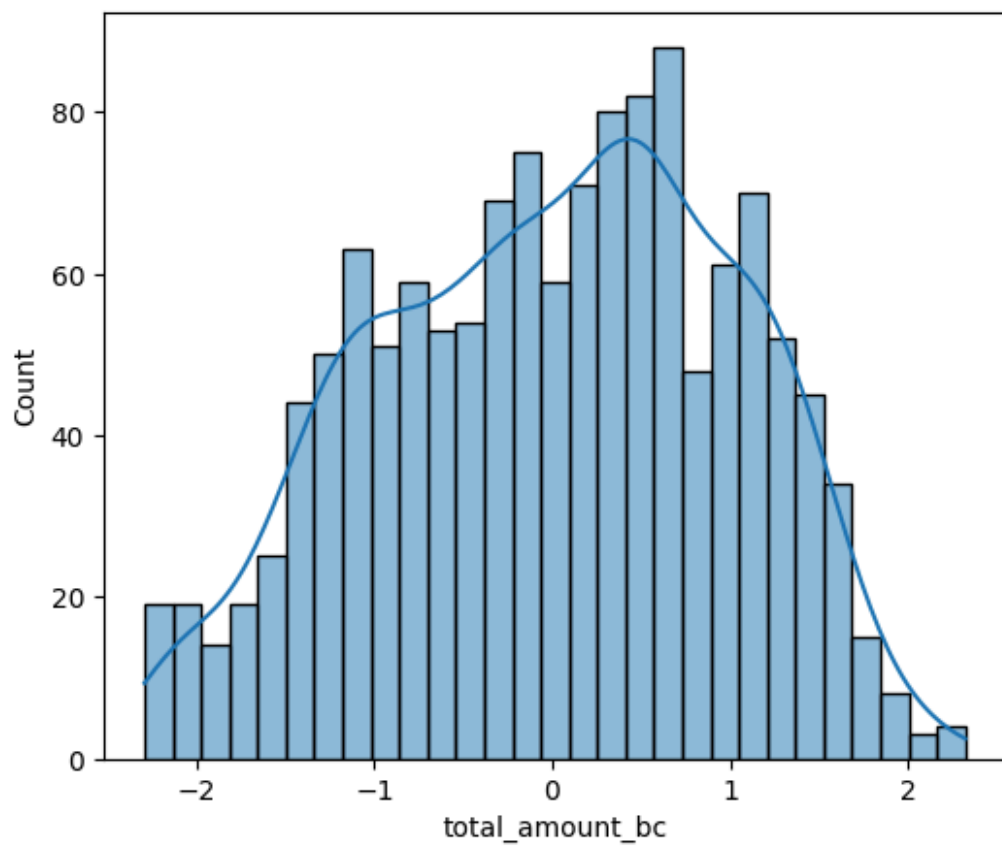


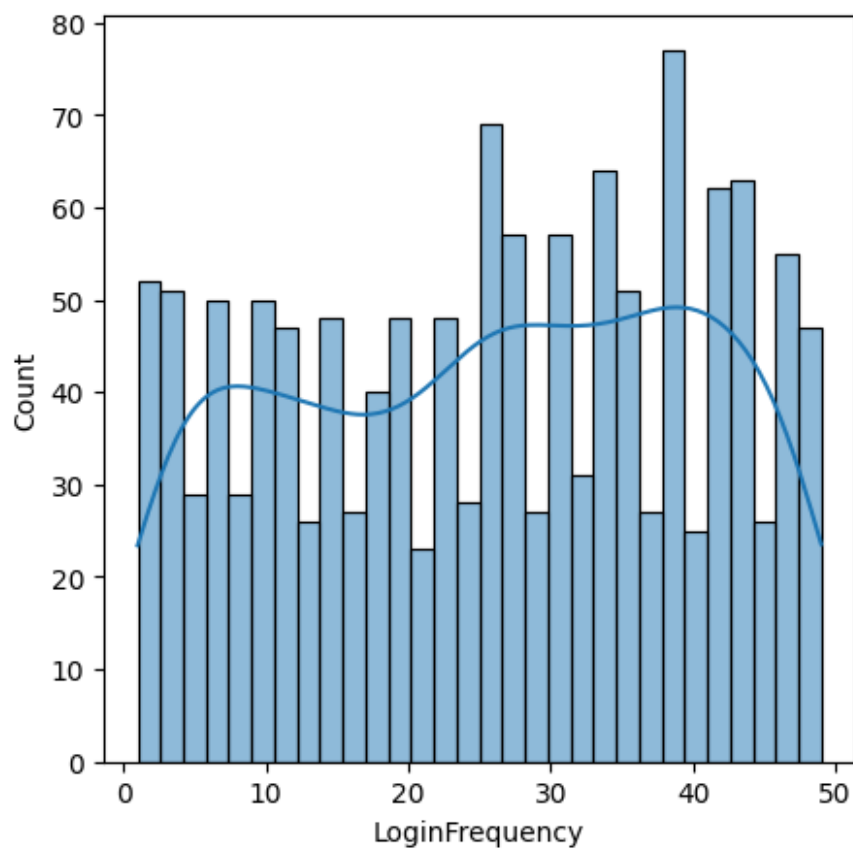
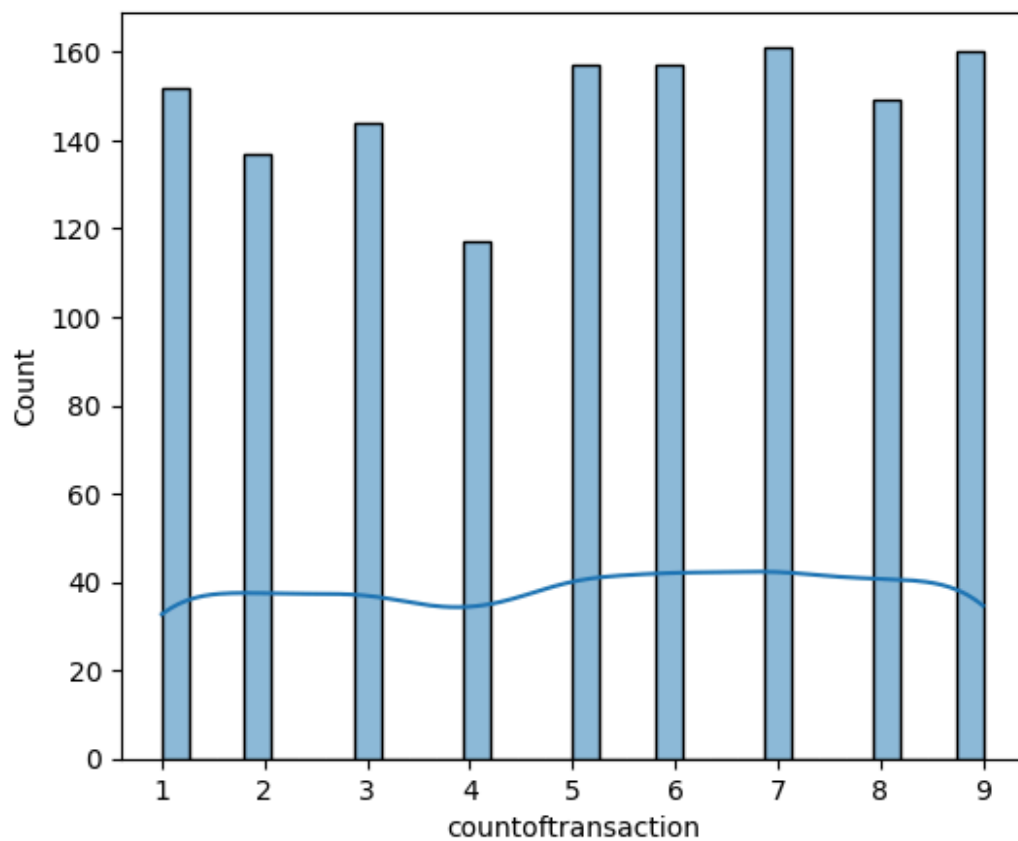


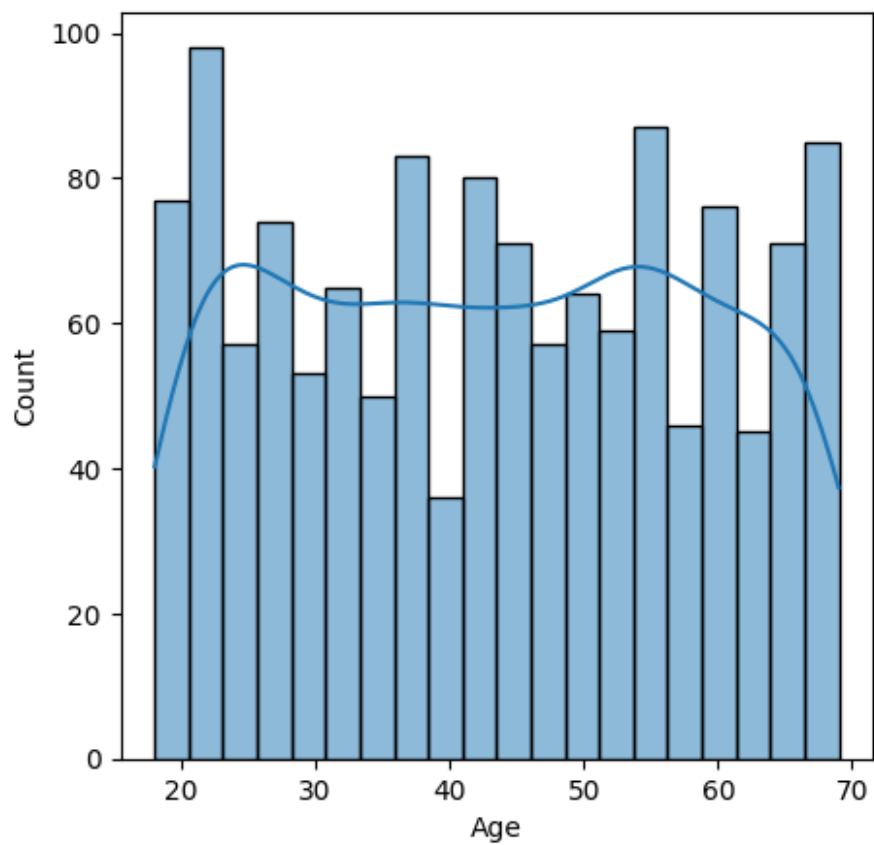
```
from sklearn.preprocessing import PowerTransformer

# For strictly positive data (Box-Cox)
pt = PowerTransformer(method='box-cox')
df_final['total_amount_bc'] =
pt.fit_transform(df_final[['total_amount']])
plt.figure(figsize=(6, 5))
sns.histplot( data = df_final                , x='total_amount_bc' , bins = 29          ,
kde= True )

<Axes: xlabel='total_amount_bc', ylabel='Count'>
```







## **4. Feature Scaling**

To ensure consistent scales across features, especially for distance-based and gradient-based models, all numeric variables were standardized using Z-score normalization. This method centers the features around a mean of zero and standard deviation of one, helping algorithms like logistic regression and SVM converge faster and perform more accurately. Log-transformed features were scaled after transformation to maintain consistency in range and distribution.

## **5. Encoding Categorical Variables**

Categorical variables including Gender, MaritalStatus, and ResolutionStatus were transformed using one-hot encoding. This approach preserves the non-ordinal nature of these categories and prevents the model from inferring an artificial ranking. To avoid dimensionality explosion from sparse categories, rare levels were grouped under an "Other" category where necessary. The target variable ChurnStatus was already binary (0 for retained customers, 1 for churned customers) and required no further encoding.



## 6. CONCLUSION

Through careful handling of missing data, strategic treatment of outliers, consistent feature scaling, and appropriate encoding of categorical variables, the dataset is now clean and structured for effective modeling. These preprocessing steps lay a robust foundation for the next phase of the project: exploratory data analysis and predictive modeling. The final dataset balances statistical integrity with business relevance, ensuring high-quality insights into customer churn behavior.