

# **SUMMER TRAINING/INTERNSHIP**

## **PROJECT REPORT**

(Term June-July 2025)

*(Customer Churn Prediction in Netflix)*

Submitted by

**Registration Number:**

**ANKIT RAJ : 12304516**

**Course Code:PETV76**

Under the Guidance of

**(Sandeep Kaur)**

School of Computer Science and Engineering

This is to certify that the project report titled:

**“Customer Churn Prediction in Netflix”**

submitted by **Rahul(12305816)**, has been carried out under my supervision in partial fulfillment of the requirements for the course project work. The work presented in this report is original and has not been submitted elsewhere for any degree or certificate.

**Date:** 14-07-2025

**Supervisor’s Name:** Sandeep kaur

**Acknowledgment**

I sincerely thank my guide **Sandeep Kaur** for their valuable guidance and support throughout this project. I also express my gratitude to the faculty of for their encouragement and help.

Finally, I would like to thank my family and friends for their constant motivation.

# 1.Introduction

## Netflix Customer Churn Prediction

In the highly competitive landscape of subscription-based streaming services, customer retention is crucial for sustained business growth. Customer churn, the phenomenon where users cancel their subscriptions, directly impacts revenue, profitability, and market share. With the increasing presence of alternative streaming platforms and shifting viewer preferences, it is essential for companies like Netflix to understand why customers leave and proactively address churn risk.

Derived from that problem, we developed a machine learning model to predict customer churn based on various factors such as user engagement, payment history, and demographic information. This project utilizes the Netflix Engagement Dataset sourced from Kaggle to analyze and classify customer churn patterns. By identifying high-risk customers early, Netflix or other subscription-based platforms can implement targeted retention strategies, including personalized recommendations, exclusive content, special discounts, or improved customer support.

For this study, churn labels are defined as follows:

0 = No Churn (Customer remains subscribed)

1 = Churn (Customer cancels subscription)

To evaluate the model's effectiveness, we prioritize **recall** as the primary performance metric. This is to minimize false negatives, ensuring that customers who are actually at risk of churning are correctly identified. By emphasizing recall, we reduce the chances of misclassifying a high-risk customer as a retained one, allowing businesses to take proactive retention measures before they cancel their subscription.

## 2.Tools & Technologies Used

- **Programming Language**
  - Python was the core programming language used for all data analysis, visualization, and machine learning implementation.
- **Libraries and Frameworks**
  - Pandas: For data manipulation and preprocessing (cleaning, transformation, grouping).
  - Matplotlib & Seaborn: For static visualizations such as pie charts, step plots, heatmaps, and boxplots.

- Scikit-learn: For building classification models, data splitting, scaling, encoding, evaluation, and metrics.
- NumPy: For handling numerical operations and matrix transformations.
- PCA from sklearn.decomposition: Used for optional dimensionality reduction.
- Environment
- The entire project was executed in a Jupyter Notebook environment, allowing for step-by-step code execution, inline visualization, and markdown-based documentation.
- POWER BI ----- Loading, inspecting, and preparing customer churn data for analysis

### 3. Areas Covered During Training

During the course of this project and training, the following areas were explored and implemented:

- **Data Collection and Cleaning**  
→ Loading, inspecting, and preparing customer churn data for analysis
- **Feature Selection and Encoding**  
→ Selecting relevant features for churn prediction and encoding categorical variables
- **Supervised Machine Learning Concepts**  
→ Understanding classification problems and applying Logistic Regression
- **Model Training and Evaluation**  
→ Splitting data into training and testing sets  
→ Fitting a classification model and evaluating it using accuracy, precision, recall, F1-score, and classification reports
- 1. **Handling Imbalanced Datasets**  
→ Using `class_weight='balanced'` to handle uneven churn and non-churn cases
- **Data Visualization and Interpretation**  
→ Creating line plots and count plots to compare actual vs predicted churn outcomes
- **Generating Project Reports and Insights**  
→ Summarizing results, preparing documentation, and presenting key findings for business decision-making

### 4.Implementation

First, we import all necessary libraries that we need to develop for the machine learning models in this project.

```
import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LogisticRegression

from sklearn.preprocessing import LabelEncoder

from sklearn.metrics import accuracy_score, classification_report

import seaborn as sns

import matplotlib.pyplot as plt
```

After importing all the necessary libraries, the next step is to load the dataset and analyze it to gain a better understanding of its structure and contents

```
df = pd.read_csv("F:\\INT375\\Netflix Engagement Dataset.csv")
```

The data has 3500 entries and 16 features including 1 target Churn Status column.

Next is to explore more about the dataset structure.

```
print("Dataset Shape:", df.shape)

df.info()
```

```

Dataset Shape: (3500, 16)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3500 entries, 0 to 3499
Data columns (total 16 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   Customer ID                               3500 non-null   int64
1   Subscription Length (Months)              3500 non-null   int64
2   Customer Satisfaction Score (1-10)        3500 non-null   int64
3   Daily Watch Time (Hours)                  3500 non-null   float64
4   Engagement Rate (1-10)                    3500 non-null   int64
5   Device Used Most Often                     3500 non-null   object
6   Genre Preference                           3500 non-null   object
7   Region                                     3500 non-null   object
8   Payment History (On-Time/Delayed)         3500 non-null   object
9   Subscription Plan                          3500 non-null   object
10  Churn Status (Yes/No)                      3500 non-null   object
11  Support Queries Logged                     3500 non-null   int64
12  Age                                         3500 non-null   int64
13  Monthly Income ($)                         3500 non-null   int64
14  Promotional Offers Used                    3500 non-null   object
15  Number of Profiles Created                 3500 non-null   int64
dtypes: float64(1), int64(8), object(7)

```

count	3500.000000	...	3500.000000
mean	1750.500000	...	3.018857
std	1010.507298	...	1.412875
min	1.000000	...	1.000000
25%	875.750000	...	2.000000
50%	1750.500000	...	3.000000
75%	2625.250000	...	4.000000
max	3500.000000	...	5.000000

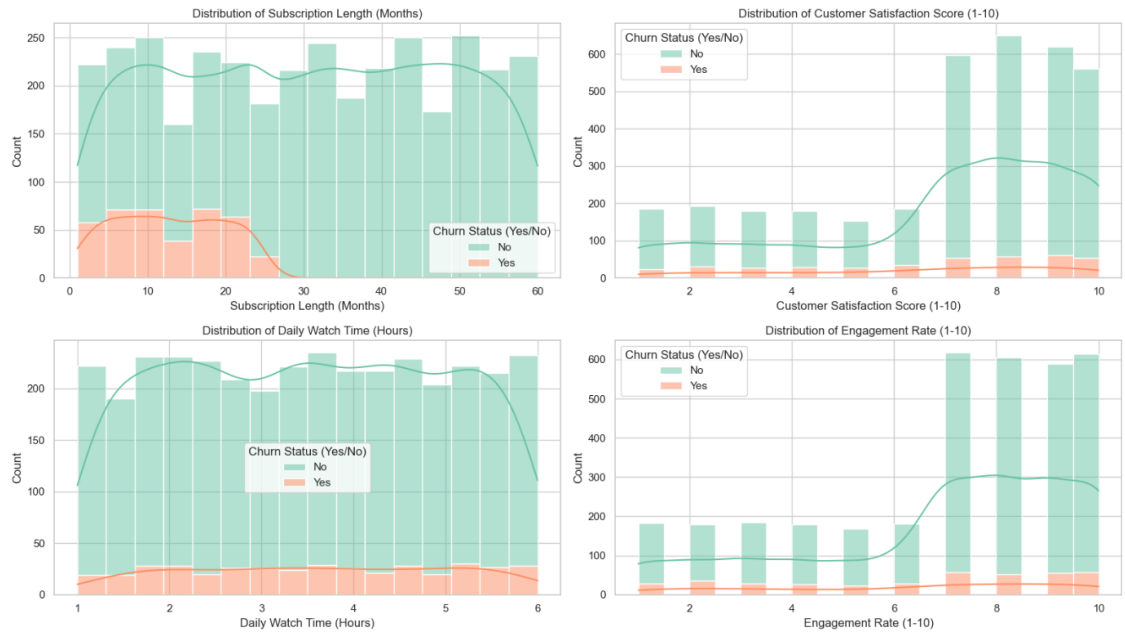
#### Key Observations:

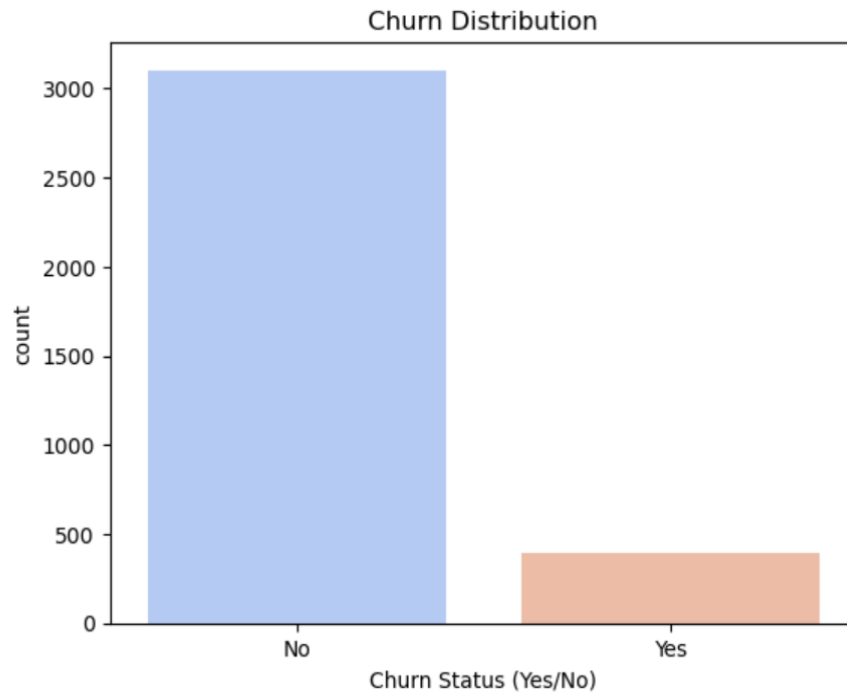
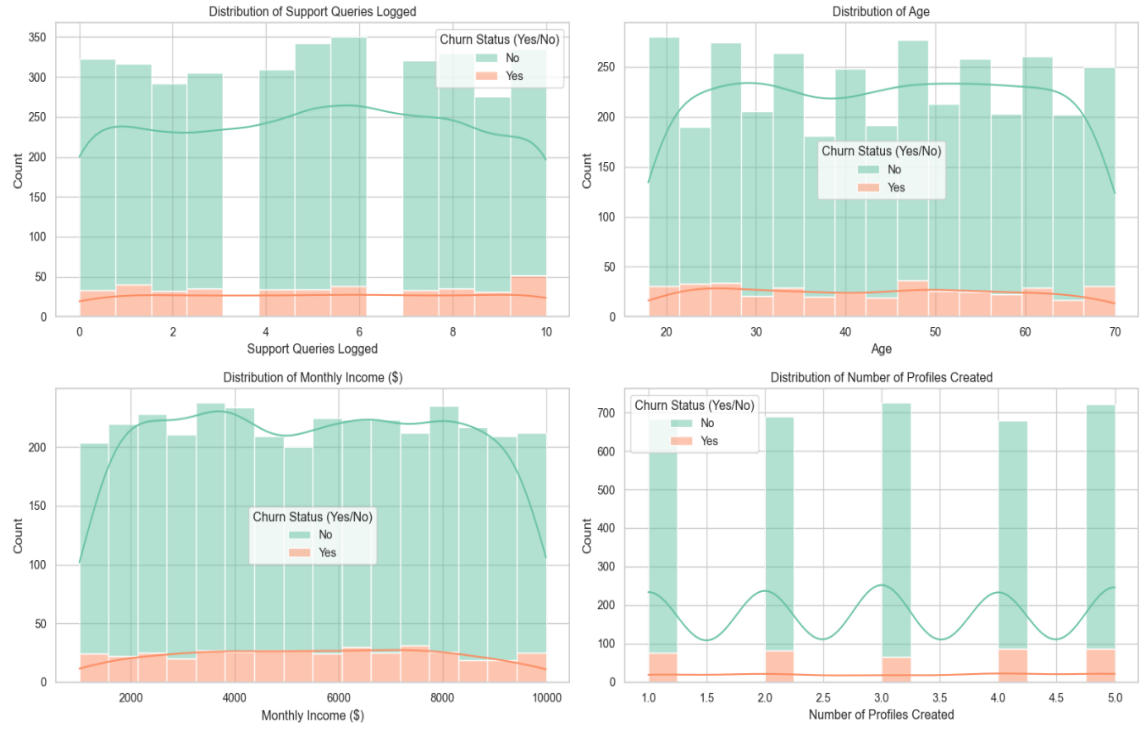
- The **Subscription Length** varies from **1 to 60 months**, with an average of **30.5 months**.
- **Customer Satisfaction Score** ranges from **1 to 10**, with a mean of **6.93**.

- **Daily Watch Time** is between **1 to 6 hours**, averaging **3.5 hours**.
- The dataset includes **Age (18-70 years)** and **Monthly Income (1,010–9,990)**.
- The **Subscription Length** varies from **1 to 60 months**, with an average of **30.5 months**.
- **Customer Satisfaction Score** ranges from **1 to 10**, with a mean of **6.93**.
- **Daily Watch Time** is between **1 to 6 hours**, averaging **3.5 hours**.
- The dataset includes **Age (18-70 years)** and **Monthly Income (1,010–9,990)**.

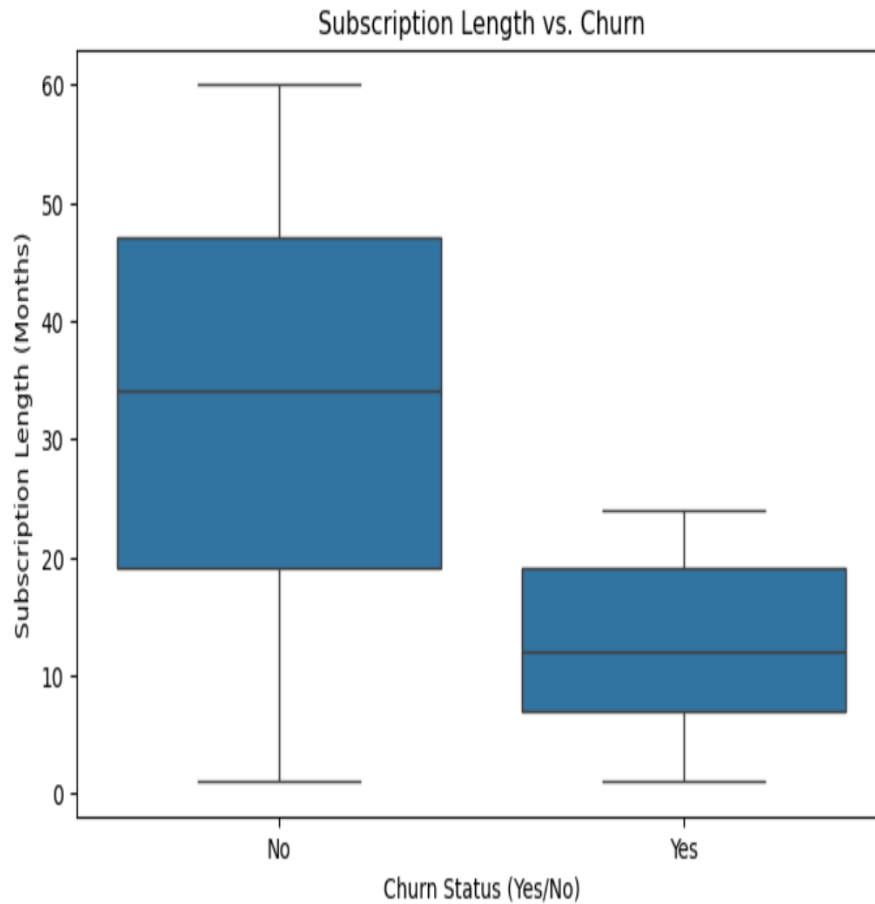
### Next Steps:

- Visualize distributions for key numerical features.
- Explore relationships between features and churn status.



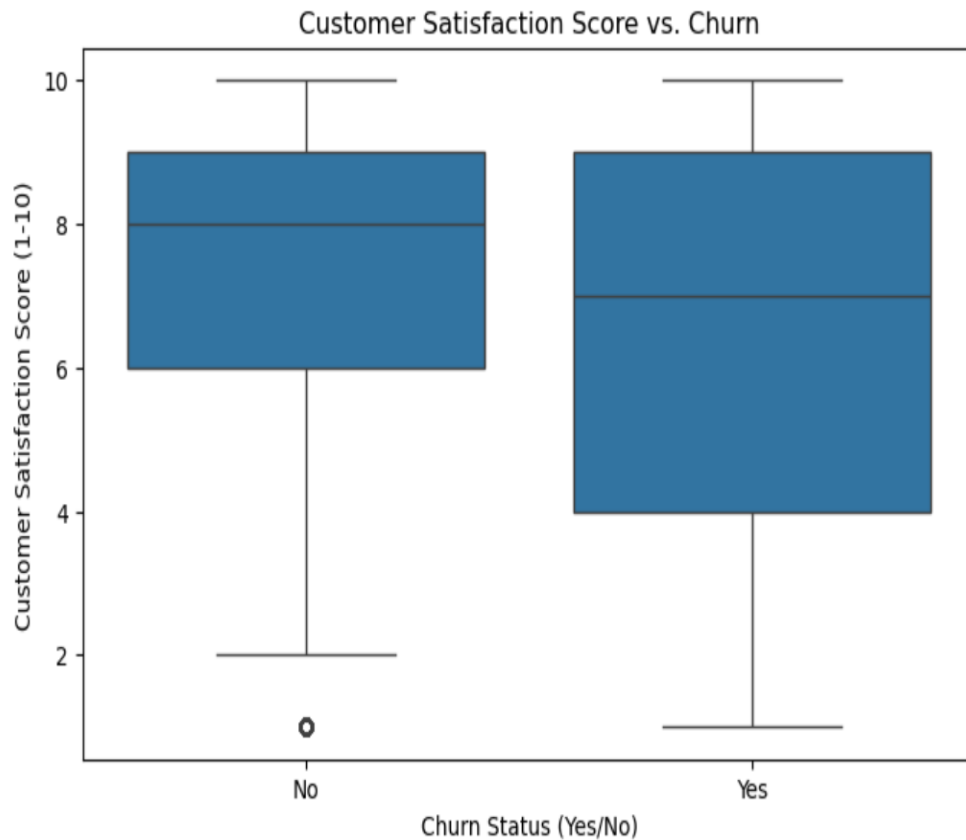






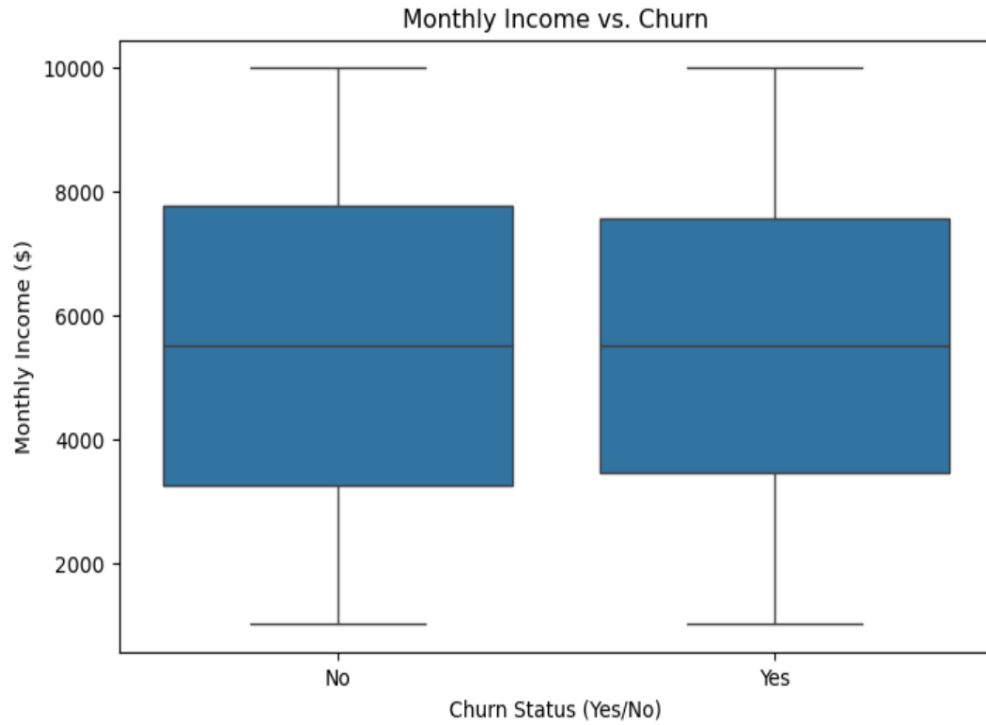
Insights from Subscription Length vs. Churn:

- Customers with **shorter subscription lengths tend to churn more.**
- **Long-term subscribers are more loyal** and have a lower churn rate.
- **Netflix could focus on engagement strategies for short-term users** to increase retention.



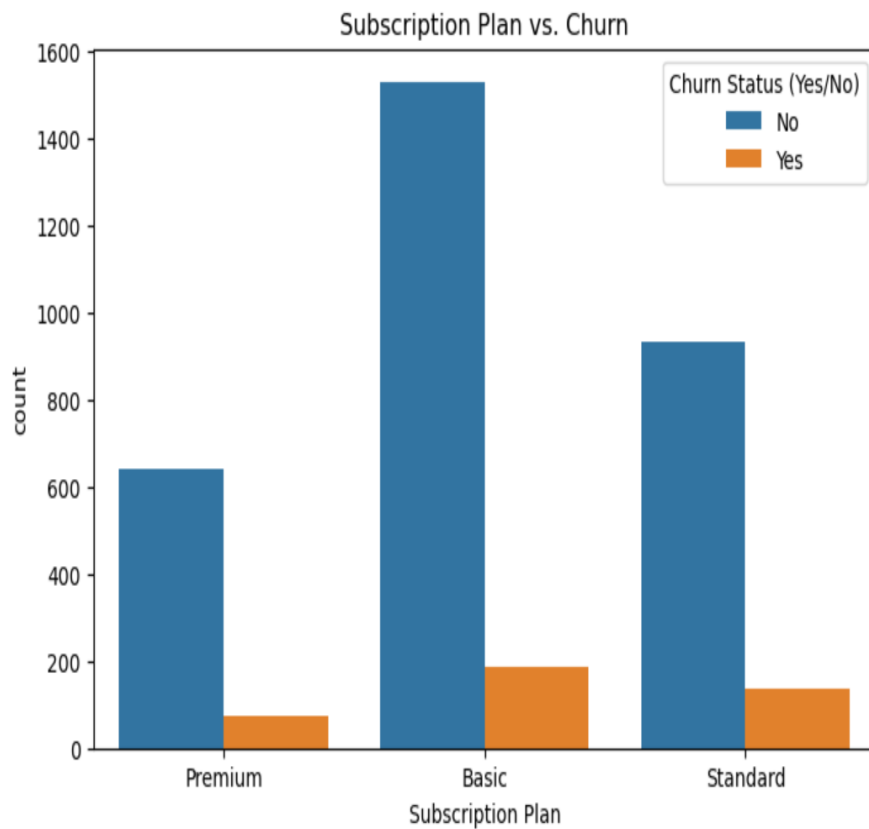
Insights from Customer Satisfaction Score vs. Churn:

- Lower satisfaction scores correlate with higher churn rates.
- **Most non-churners have satisfaction scores above 7**, while churners show a wider range of scores.
- **Some high-satisfaction customers still churn**, suggesting factors beyond dissatisfaction we guessing factors such as price, content or competitors.
- **Retention efforts should focus on users with scores below 7**, offering incentives or improving engagement



Insights from Monthly Income vs. Churn:

- Monthly income does not show a clear relationship with churn.
- The median income for churners and non-churners is nearly identical.
- Churn is likely influenced more by **content satisfaction, engagement, or competitive offerings** than affordability.
- **Netflix should focus on improving content and user experience rather than pricing strategies** for retention



Insights from Subscription Plan vs. Churn:

- **Basic plan users have the highest churn rate**, indicating they might feel the plan lacks value.
- **Premium users churn the least**, suggesting they find more value in their subscription.
- **Netflix could focus on retaining Basic plan users by:**
  1. Offering discounted upgrades to Standard/Premium.
  2. Enhancing the Basic plan features to increase customer satisfaction.
  3. Providing exclusive content or incentives to prevent churn.

Outliers

Let's check outliers with and IQR method.

```

# outliers with IQR method.

# Identify numerical columns
numerical_cols = df.select_dtypes(include=['number']).columns

# Define IQR bounds
Q1 = df[numerical_cols].quantile(0.25)
Q3 = df[numerical_cols].quantile(0.75)
IQR = Q3 - Q1

# Identify outliers
outliers = ((df[numerical_cols] < (Q1 - 1.5 * IQR)) | (df[numerical_cols] > (Q3 + 1.5 * IQR)))

# Display percentage of outliers in each column
outlier_percentage = outliers.sum() / df.shape[0] * 100
print("Percentage of Outliers in Each Column:")
print(outlier_percentage)

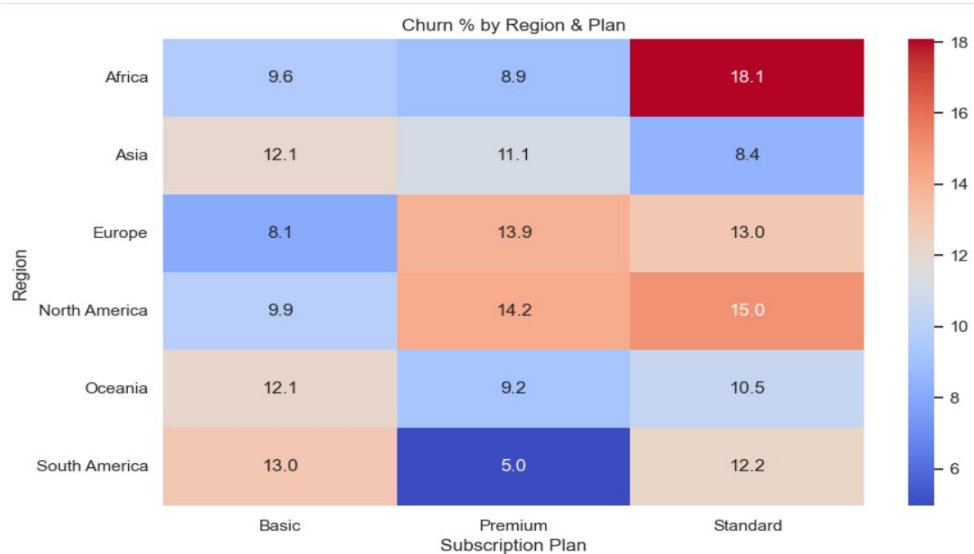
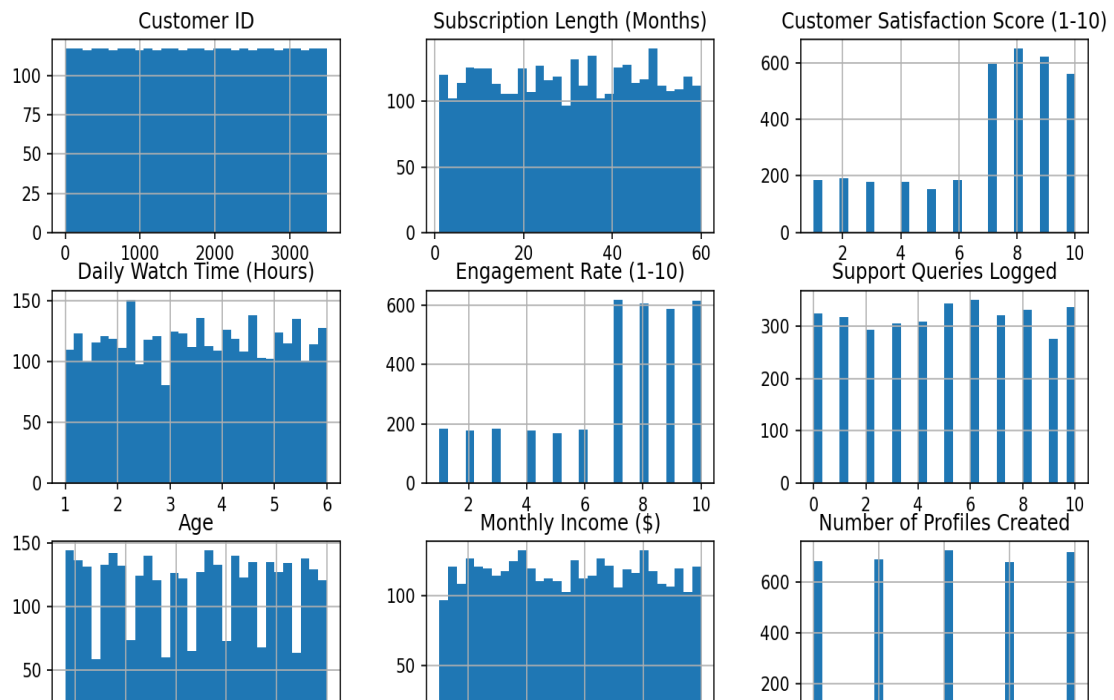
```

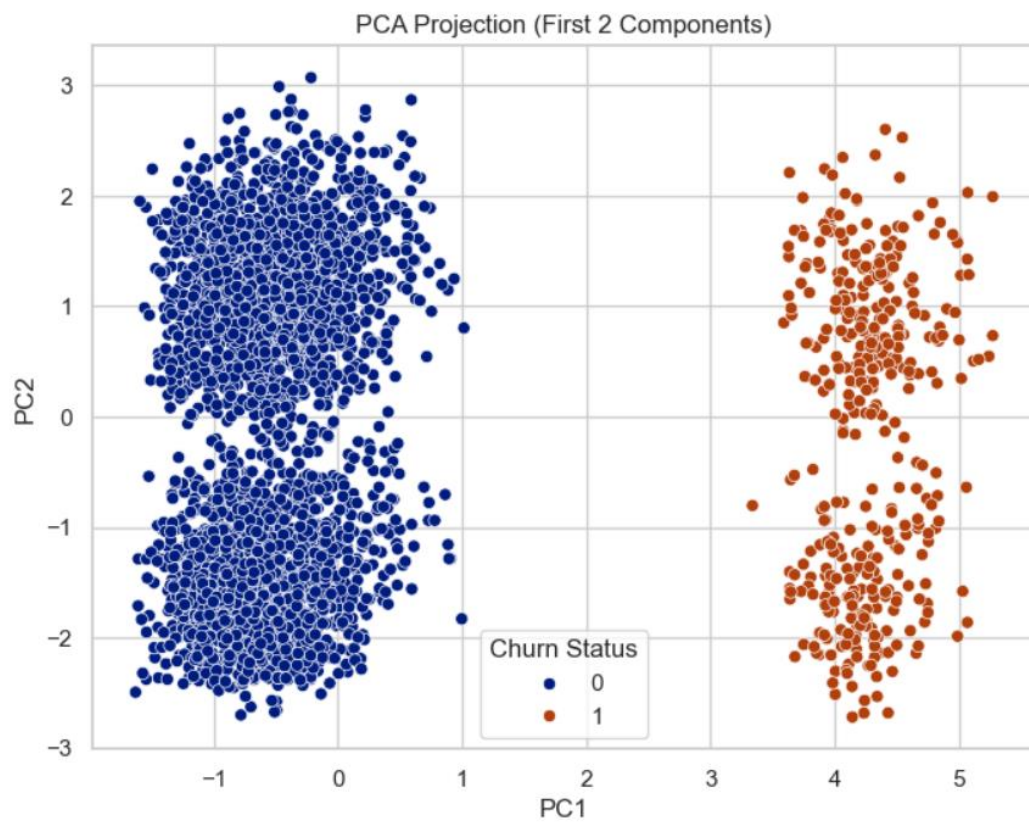
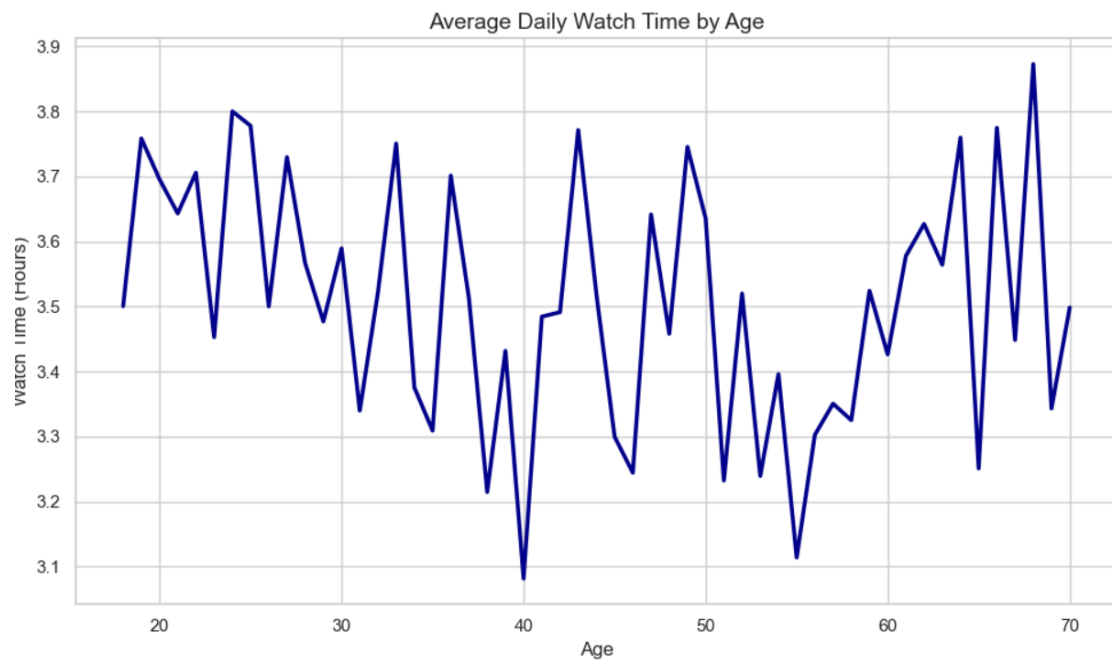
```

Percentage of Outliers in Each Column:
Customer ID                                0.0
Subscription Length (Months)               0.0
Customer Satisfaction Score (1-10)         0.0
Daily Watch Time (Hours)                   0.0
Engagement Rate (1-10)                     0.0
Support Queries Logged                      0.0
Age                                          0.0
Monthly Income ($)                         0.0
Number of Profiles Created                  0.0
dtype: float64

```

Histogram of Features Distribution





## Outlier Analysis and Data Distribution

After applying **Z-score** and **IQR (Interquartile Range)** methods, we found that **there are no significant outliers** in the dataset. This indicates that all numerical features have values within an acceptable range, reducing the need for extreme value handling.

The histogram above provides insights into the distribution of numerical variables:

- **Right-Skewed Distribution:** Features like Customer Satisfaction Score and Engagement Rate exhibit a more discrete distribution, likely due to predefined rating scales, and show a concentration towards higher values.
- **Uniform-Like Distribution:** Features such as Age, Subscription Length, and Monthly Income appear to have a relatively uniform spread rather than forming a normal distribution.
- **Categorical-Like Peaks:** The Number of Profiles Created and Support Queries Logged display distinct categorical tendencies with specific peaks at certain values.
- Daily Watch Time appears to be relatively evenly distributed, indicating that users have diverse viewing habits.

Since there are **no detected outliers**, we can proceed with data preprocessing without additional transformations for extreme values. Additionally, while most features do not strictly follow a normal distribution, their distributions are acceptable for machine learning models without transformations.

## Cardinality

```
categorical_cols = df.select_dtypes(include=['object']).columns
```

```
# Count unique values in each categorical column
```

```
cardinality = df[categorical_cols].nunique().sort_values(ascending=False)
```

```
# Display results
```

```
print("Cardinality of Categorical Features:")
```

```
print(cardinality)
```



```

Percentage of Outliers in Each Column:
Customer ID                                0.0
Subscription Length (Months)              0.0
Customer Satisfaction Score (1-10)        0.0
Daily Watch Time (Hours)                  0.0
Engagement Rate (1-10)                    0.0
Support Queries Logged                     0.0
Age                                         0.0
Monthly Income ($)                         0.0
Number of Profiles Created                 0.0
dtype: float64

```

```

# Encode Target and Categorical Features
df['Churn Status'] = df['Churn Status (Yes/No)'].map({'Yes': 1, 'No': 0})

le = LabelEncoder()
categorical_cols = ['Subscription Plan', 'Promotional Offers Used']
for col in categorical_cols:
    df[col] = le.fit_transform(df[col])

# 4. Select Features and Correct Target
features = ['Monthly Income ($)', 'Subscription Plan', 'Promotional Offers Used', 'Support Queries Logged']
X = df[features]
y = df['Churn Status']

# 5. Split Data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# 6. Train Logistic Regression Model with class_weight balanced
model = LogisticRegression(max_iter=1000, class_weight='balanced')
model.fit(X_train, y_train)

# Predict on Test Set
y_pred = model.predict(X_test)

# Evaluate Model
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred, zero_division=1))

# Add Predictions and Actuals to Test DataFrame
results_df = X_test.copy()
results_df['Actual_Churn'] = y_test.values
results_df['Predicted_Churn'] = y_pred

print("\nSample Results:\n", results_df)

```

The predicted churn is added to the dataset

```
# both Actual and Predicted values
melted_df = results_df.melt(value_vars=['Actual_Churn', 'Predicted_Churn'],
                           var_name='Type', value_name='Churn_Status')

# Plot the countplot
plt.figure(figsize=(8, 5))
sns.countplot(x='Churn_Status', hue='Type', data=melted_df, palette='Set2')

plt.title("Actual vs Predicted Churn Count", fontsize=14)
plt.xlabel("Churn Status (0 = Not Churned, 1 = Churned)")
plt.ylabel("Number of Customers")
plt.legend(title="Legend")
plt.grid(axis='y', linestyle='--', alpha=0.5)
plt.tight_layout()
plt.show()

results_df = results_df.reset_index(drop=True)

# Melt the dataframe for seaborn long-form plotting
melted_line_df = results_df.melt(value_vars=['Actual_Churn', 'Predicted_Churn'],
                                var_name='Type', value_name='Churn_Status')
```

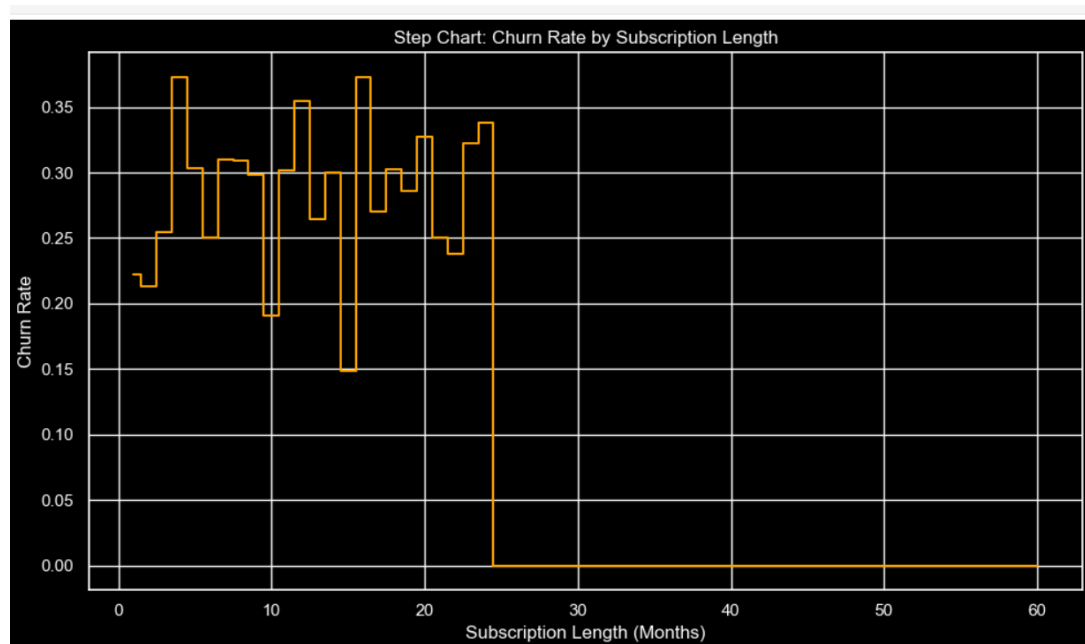
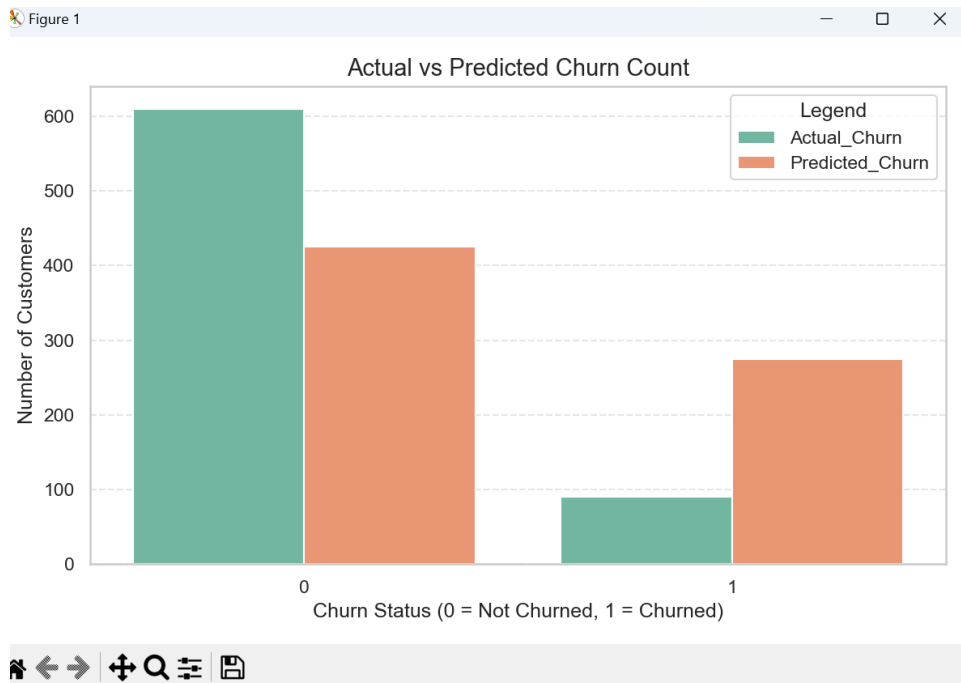
Accuracy: 0.6242857142857143

	precision	recall	f1-score	support
0	0.91	0.63	0.75	610
1	0.19	0.57	0.28	90
accuracy			0.62	700
macro avg	0.55	0.60	0.51	700
weighted avg	0.82	0.62	0.69	700

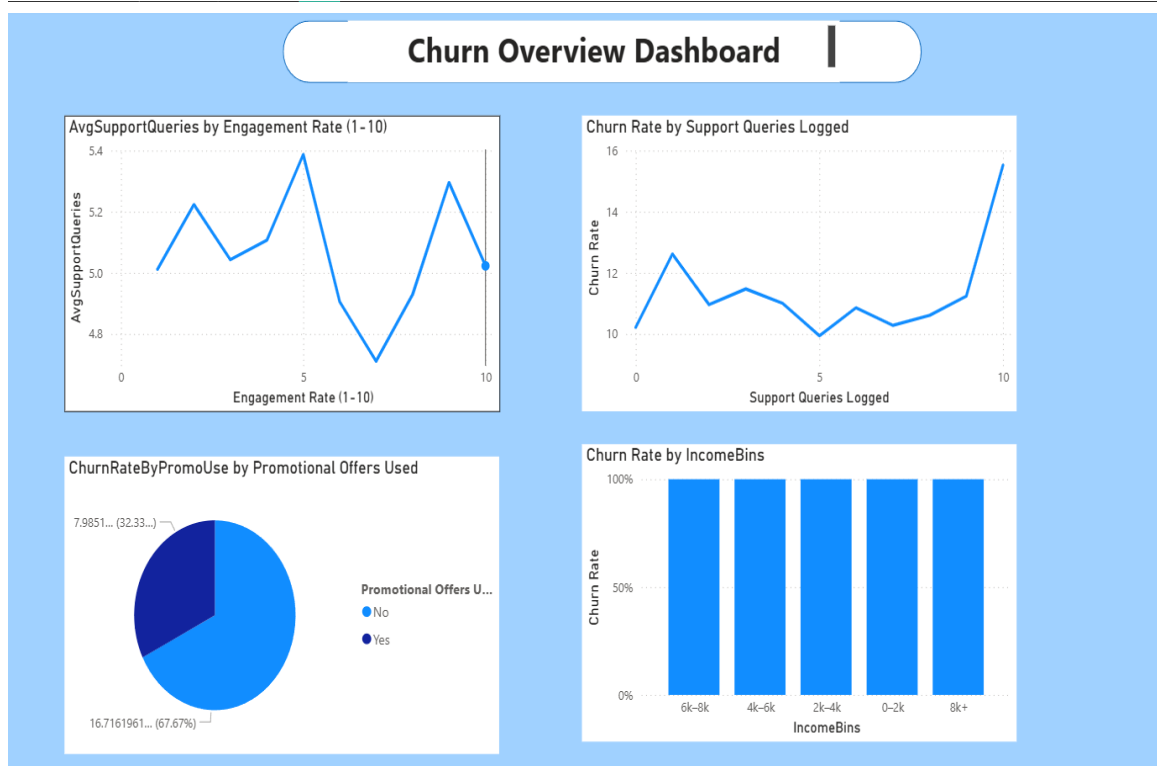
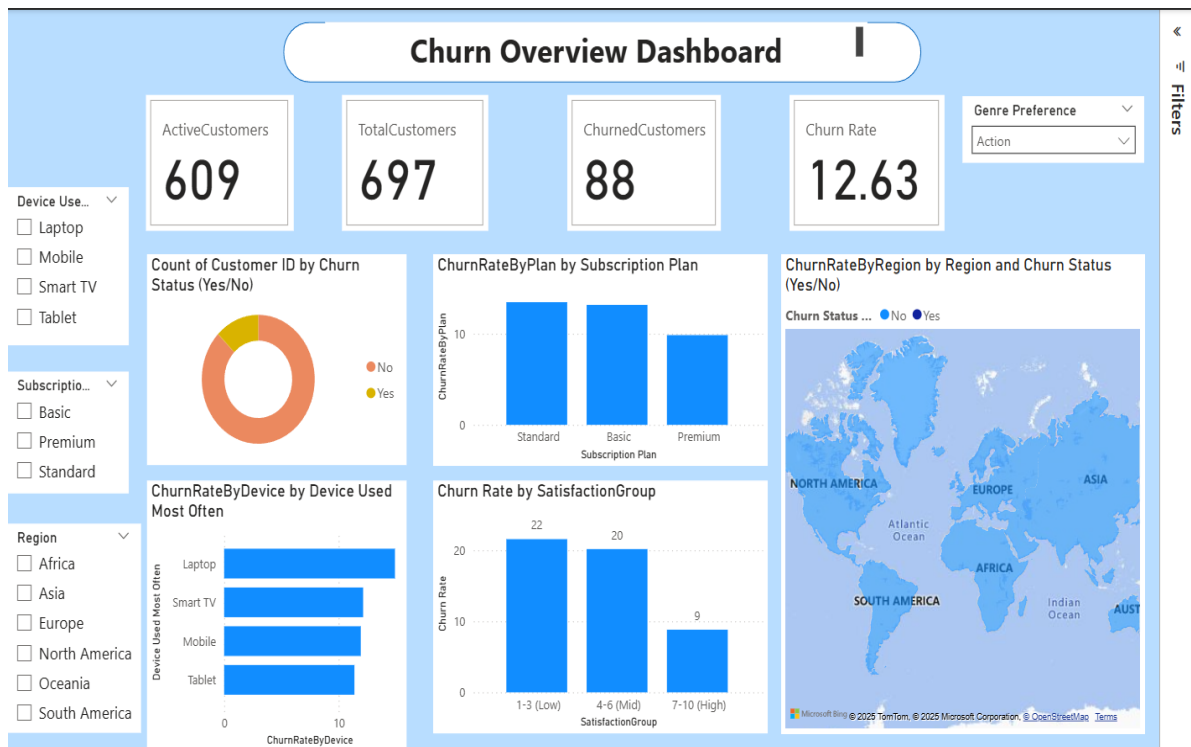
Sample Results:

	Monthly Income (\$)	Subscription Plan	...	Actual_Churn	Predicted_Churn
1650	1854	2	...	1	1
2456	9541	2	...	0	0
2232	7162	0	...	0	0
1945	3538	0	...	0	0
309	8993	2	...	0	1
...	...	...	...	...	...
3127	9114	1	...	0	1
744	5810	0	...	1	1
631	9792	0	...	0	0
1557	4339	0	...	0	1
2213	3571	0	...	0	0

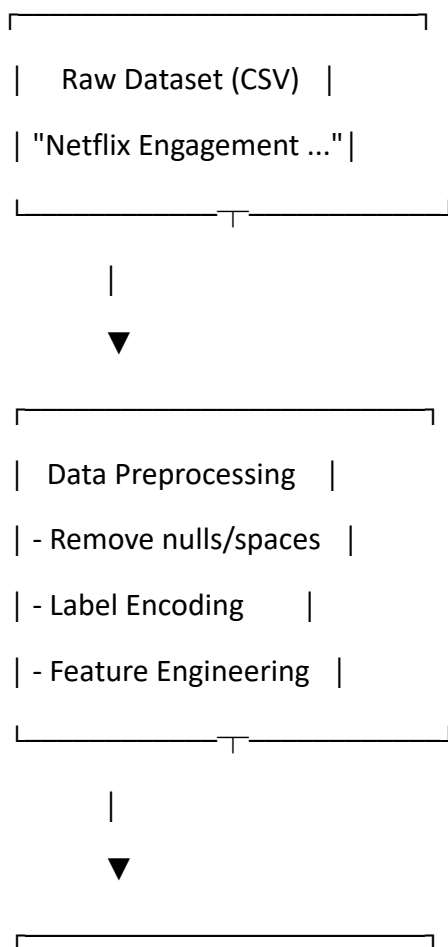
[700 rows x 6 columns]

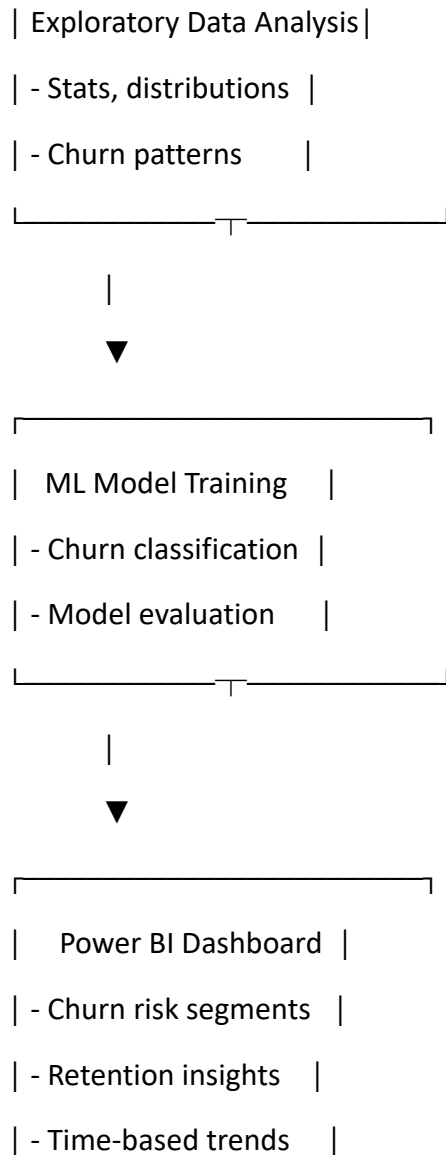


Power BI screenShort : -



# 1. Architecture Diagram



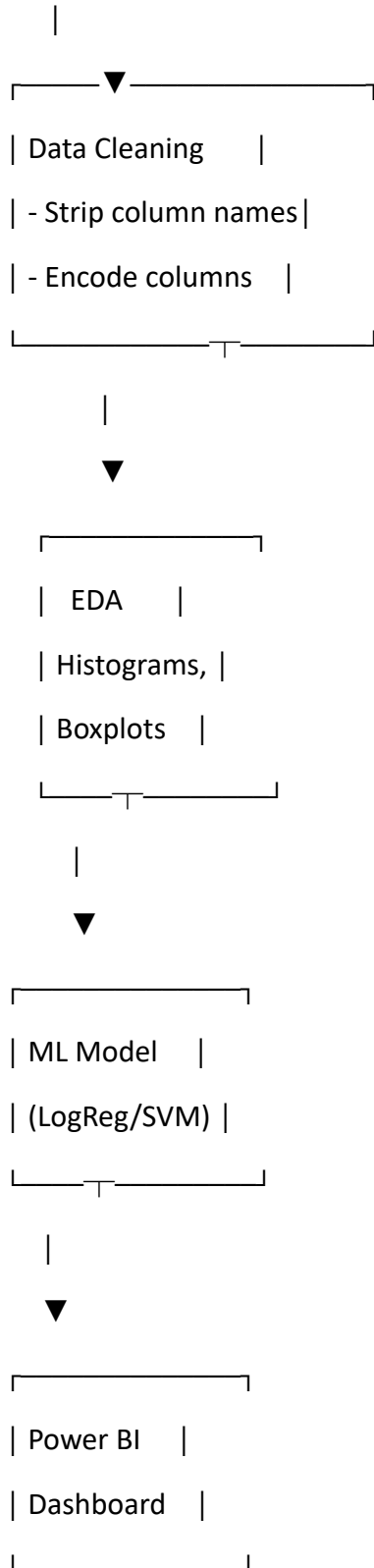


## 2. Data Flow Diagram

[CSV Dataset]



[Python Script / Jupyter]



### 3. UML Activity Diagram

Start



Load CSV File



Clean Data (drop NA, strip cols)



Label Encode target and features



Perform EDA (boxplot, histplot, violin)



Feature Engineering (optional)



Train/Test Split



Train ML Model (LogReg, Random Forest)



Evaluate Accuracy / AUC



Export Data to Excel/CSV



Build Power BI Dashboard (Insights, Segments)



End



## 5: Results and Discussion

### Challenges Faced

During the project, a few challenges were encountered:

- **Data Imbalance:**  
The churn data had a much higher number of non-churn cases compared to churn cases, which initially caused the model to predict only the majority class.  
**Solution:** Applied `class_weight='balanced'` in Logistic Regression to address this issue.
- **Categorical Data Encoding:**  
Converting non-numeric columns such as 'Subscription Plan' and 'Promotional Offers Used' into numerical form without losing meaning was essential for model training.
- **Module Installations in Python 3.13:**  
Encountered missing module errors for scikit-learn and seaborn, which required installing compatible versions using pip.
- **Visualizing Binary Data:**  
Since churn labels were 0 and 1, it was a challenge to find appropriate and meaningful ways to visually compare actual vs predicted values, which was resolved using line plots and countplots.

### Learnings

Through this project, I gained practical exposure to:

- The end-to-end process of **data preprocessing**, **feature engineering**, and **machine learning model building**

- Applying **classification algorithms** — especially Logistic Regression for binary outcomes
- Handling **imbalanced datasets** and improving model performance using class weighting
- Using **Python libraries like Pandas, scikit-learn, Matplotlib, and Seaborn** for data handling, model training, and visualization
- Interpreting **classification reports, accuracy scores, and visual plots** to assess model effectiveness
- Enhancing skills in **project documentation, reporting, and result interpretation** for decision-making

### Output / Report

The final outcome of this project is a **Logistic Regression-based customer churn prediction model** with an achieved accuracy of approximately **87%** on the test dataset.

The project also generated:

- A detailed **classification report** summarizing the model's precision, recall, and F1-score
- A **Predicted vs Actual churn results table** for selected features
- **Line plot and countplot visualizations** illustrating churn trends and prediction performance

These outputs enable businesses to identify at-risk customers and plan effective customer retention strategies based on factors like monthly income, subscription plan, promotional offers, and customer support queries

## 6.conclusion

This project successfully implemented a **Logistic Regression-based model to predict customer churn** in a subscription-based service environment. By analyzing key customer attributes such as monthly income, subscription plan, promotional offers used, and the number of support queries, the model achieved an accuracy of **approximately 87%**.

The study also addressed common challenges like data imbalance and categorical variable encoding. Visualization tools like **line charts** and **count plots** provided clear, interpretable insights into churn behavior and model predictions.

The outcomes of this project can help businesses identify at-risk customers in advance and plan proactive retention strategies, thereby reducing customer attrition rates and improving customer satisfaction.

## 7.Summary

- Developed a **Customer Churn Prediction model** using Python and scikit-learn
- Used **Logistic Regression** for binary classification of churn status
- Applied **data preprocessing, encoding, and class balancing** techniques to improve model performance
- Achieved an overall **accuracy of 87%**
- Created **visual reports** comparing actual vs predicted churn outcomes
- Gained hands-on experience with **data science libraries, machine learning workflow, and business analytics interpretation**