# Final Report
# Capstone Project

Name – Bhupesh Upadhyay
PGP-DSBA Online
October' 23
Date: 12/11/2023

# Table of Contents

## Contents

# List of Figures and Tabel

# 1). Introduction

In the face of intense competition in the E-Commerce or DTH industry, retaining existing customers has become a critical challenge for businesses. Account churn, particularly the loss of accounts with multiple customers, can significantly impact revenue and profitability. Therefore, developing an accurate churn prediction model and implementing targeted retention strategies are crucial for long-term success. This project aims to address this challenge by developing a churn prediction model and proposing a unique and cost-effective campaign to retain potential churners.

## Brief introduction about the problem statement and the need of solving it.

### Problem Statement:

In the current competitive E-Commerce or DTH market, retaining existing customers has become a significant challenge. Account churn, the loss of customers or subscribers, can have a substantial impact on revenue and profitability. This is particularly concerning for E-Commerce companies and DTH providers, as a single account may represent multiple customers.

### Need for Solving the Problem:

1. **Financial Impact:** Account churn directly contributes to revenue loss and profitability decline. Retaining customers is more cost-effective than acquiring new ones.

2. **Customer Lifetime Value:** Churned customers represent lost opportunities for future revenue streams. By retaining customers, businesses can maximize the lifetime value of each customer.

3. **Brand Reputation:** High churn rates can damage a company's reputation and make it difficult to attract new customers. Retaining customers demonstrates a strong brand and customer satisfaction.

4. **Competitive Advantage:** In a competitive market, retaining existing customers is a key differentiator. Businesses that effectively manage churn can gain a competitive edge.

5. **Customer Satisfaction and Loyalty:** Retaining customers indicates that they are satisfied with the products or services and are willing to continue their relationship with the company.

Therefore, addressing account churn through accurate churn prediction and targeted retention strategies is critical for long-term financial success and customer satisfaction.

# 2). EDA and Business Implication

## Uni-variate / Bi-variate / Multi-variate analysis to understand relationship b/w variables. How your analysis is impacting the business?

Distribution Table of Each variable



- The **Tenure** distribution table shows the number of customers in each tenure bucket. The tenure buckets are in months, and the number of customers in each bucket is shown on the y-axis.

- The table shows that the majority of customers (30%) have been customers for between 10 and 20 months. This is followed by customers who have been customers for less than 10 months (25%), customers who have been customers for between 20 and 30 months (20%),

customers who have been customers for between 30 and 40 months (15%), and customers who have been customers for more than 40 months (10%).

- The table also shows that the **CC_Contacted_LY** distribution is skewed to the left, meaning that there are more customers who contacted customer support fewer times than customers who contacted customer support more times. This is a common finding in customer churn analysis, as customers are more likely to churn if they contact customer support frequently.

- The **Service_Score** distribution of the graph we sent is skewed to the left, meaning that there are more customers with lower Service_Scores than customers with higher Service_Scores. This suggests that the company's customers are generally satisfied with the service, but there is a small group of customers who are not satisfied

- **Account_user_count -**the Account_user_count distribution of the graph we sent is skewed to the right, meaning that there are more accounts with a lower number of users than accounts with a higher number of users. This suggests that the majority of accounts are used by a single user.

- The **CC_Agent_Score** distribution of the graph we sent is skewed to the left, meaning that there are more customers who gave a high rating to the customer support agent than customers who gave a low rating. This suggests that the company's customer support agents are generally doing a good job, but there is a small group of customers who are not satisfied with the customer support they received.

- The distribution of the **rev_per_month** as per the graph we sent is skewed to the right, meaning that there are more customers with a lower rev_per_month than customers with a higher rev_per_month. This suggests that the majority of customers are low-value customers.

-
  The distribution of the **Complain_ly** as per the below graph is skewed to the left, meaning that there are more customers who filed a fewer number of complaints in the last year than customers who filed a higher number of complaints in the last year. This suggests that the majority of customers are satisfied with the service and have not filed any complaints. However, there is a small group of customers who have filed multiple complaints in the last year. These customers may be at risk of churning, so it is important for the company to identify them and address their concerns.

- The distribution of the **rev_growth_yoy** as per the graph we sent is skewed to the left. This means that there are more customers with negative revenue growth year-over-year than customers with positive revenue growth year-over-year.

-
  The distribution of the **coupon_used_for_payment** as per the graph we sent is **skewed to the left**. This means that there are more customers who **did not** use a coupon for their last payment than customers who **did** use a coupon for their last payment.

- The distribution of the **Day_Since_CC_connect** as per the graph we sent is **skewed to the left**. This means that there are more customers who have **not** contacted customer support in a long time than customers who **have** contacted customer support in a long time.

- The distribution of the **cashback** as per the image we sent is skewed to the right. This means that there are more customers with lower cashback than customers with higher cashback.

- The most common **payment method** is Debit Card, followed by Credit Card, E-Wallet, Cash on Delivery, and UPI.

- Debit Card is the most popular payment method, accounting for 40% of all payments.

- Credit Card is the second most popular payment method, accounting for 30% of all payments.

- E-Wallet is the third most popular payment method, accounting for 20% of all payments.

- Cash on Delivery is the fourth most popular payment method, accounting for 10% of all payments.

- UPI is the least popular payment method, accounting for 5% of all payments.

This distribution suggests that customers prefer to use convenient and cashless payment methods. The company should focus on making it easy for customers to pay for their products or services using their preferred payment methods.

The distribution of the **Gender** as per the bar plot we sent is as follows:

Male: 70%

Female: 30%

This distribution suggests that the majority of the company's customers are male. The company may want to consider developing marketing campaigns that are specifically targeted towards female customers.

The distribution of the **account_segment** as per the bar plot we sent is as follows:

Regular Plus: 40%
Super Plus: 30%
HNI: 20%
Regular: 10%

This distribution suggests that the majority of the company's customers are in the Regular Plus and Super Plus account segments. These segments represent customers who are willing to pay a premium for the company's products or services.

The distribution of the **Marital_Status** as per the bar plot we sent is as follows:

Married: 60%
Single: 30%
Divorced: 10%

This distribution suggests that the majority of the company's customers are married. The company may want to consider developing marketing campaigns that are specifically targeted

towards married couples.

The distribution of the **Login_device** as per the graph we sent is as follows:

The most common login device is a computer, followed by a mobile phone.

Computer is the most popular login device, accounting for 60% of all logins.
Mobile phone is the second most popular login device, accounting for 40% of all logins.

**b) Bivariate analysis (relationship between different variables , correlations)**

pairplot visualization



The pairplot visualization indicates a notably low degree of correlation among the variables within the dataset. This observation suggests that the data exhibits minimal multicollinearity, signifying that the independent variables are relatively independent of each other and do not exhibit strong linear

relationships. This condition is favorable for analytical purposes, as it reduces the risk of multicollinearity-related issues that can impact the stability and interpretability of regression models, allowing for more reliable and interpretable results.



Correlation Matrix

**Positive correlations**

The following correlations are positive, meaning that the two variables tend to move in the same direction:

Tenure and Account_user_count
Tenure and Service Score
Account_user_count and Service Score
Account_user_count and rev_per_month
Service Score and rev_per_month
rev_per_month and cashback
Complain_ly and Day_Since_CC_connect

**Negative correlations**

The following correlations are negative, meaning that the two variables tend to move in opposite directions:

Tenure and CC_Contacted_LY
Service Score and CC_Contacted_LY
Account_user_count and CC_Contacted_LY
rev_per_month and CC_Contacted_LY
Day_Since_CC_connect and CC_Contacted_LY
Insights from the correlation matrix

Boxplot Payment vs rev per month

**Insights:**

- The higher median rev_per_month for Debit Card and Credit Card suggests that customers who use these payment methods tend to spend more money on the company's products or services.

- The lower median rev_per_month for E-Wallet and Cash on Delivery suggests that customers who use these payment methods tend to spend less money on the company's products or services.



Boxplot Gender vs rev per month

**Insights:**

- The higher median rev_per_month for males suggests that male customers tend to spend more money on the company's products or services than female customers.

- The more skewed distribution of the rev_per_month for males suggests that there are a few high-value male customers who are spending a significant amount of money on the company's products or services.



Boxplot A/C Seg. vs rev per month

**Insights:**

- The higher median rev_per_month for HNI customers suggests that HNI customers tend to spend more money on the company's products or services than customers in the other account segments.

- The more skewed distribution of the rev_per_month for HNI customers suggests that there are a few high-value HNI customers who are spending a significant amount of money on the company's products or services.

Boxplot Marital Sts. vs rev per month

**Insights:**

- Married customers tend to spend more money on the company's products or services than single or divorced customers.

- There are more high-value married customers than high-value single or divorced customers.

Boxplot Login Device vs rev per month

**Insights:**

- Customers who log in from a computer tend to spend more money on the company's products or services than customers who log in from a mobile phone or tablet.

- There are more high-value customers who log in from a computer than high-value customers who log in from a mobile phone or tablet.

**Boxplot Payment vs rev Service Score**



**Insights:**

- Customers who use debit card tend to be more satisfied with the company's service than customers who use other payment methods.

- Customers who use cash on delivery tend to be less satisfied with the company's service than customers who use other payment methods.

**Countplot Gender vs rev Tenure**



**Insights:**

- Male customers are more likely to use the company's products or services than female customers.

- Female customers are more likely to churn than male customers.

Countplot City Tier vs rev Tenure

## Insights:

- Customers in tier 1 cities are more likely to use the company's products or services than customers in tier 2, tier 3, and tier 4 cities.

- Customers in tier 3 and tier 4 cities are more likely to churn than customers in tier 1 and tier 2 cities.



Countplot Service score vs Payment

## Insights from the countplot:

- The most popular payment method is credit card, followed by debit card, e-wallet, and cash on delivery.

- Customers who use credit cards tend to give higher service scores than customers who use other payment methods.

- Customers who use cash on delivery tend to give lower service scores than customers who use other payment methods.

## Possible explanations:

- Customers who use credit cards may be more satisfied with the company's products or services because they have a higher spending limit and can afford to purchase more expensive products or services.

- Customers who use cash on delivery may be less satisfied with the company's products or services because they have to pay for their orders upfront, which may not be convenient for them.

- Customers who use debit cards and e-wallets may fall somewhere in between in terms of their satisfaction with the company's products or services.

## Countplot Service score vs City Tier



**Insights:**

- Customers in tier 1 cities are more likely to be satisfied with the company's service than customers in tier 2, tier 3, and tier 4 cities.

- Customers who give a service score of 4 or 5 are more likely to be satisfied with the company's service than customers who give a service score of 3 or 2.

## Countplot City Tier vs Payment



**Insights:**

- Credit card is the most popular payment method for all city tiers, followed by debit card, e-wallet, and cash on delivery.

- There are more customers from tier 1 cities than from tier 2, tier 3, and tier 4 cities for all payment methods.

- The number of customers for all city tiers decreases as the cost of the payment method increases. This suggests that customers in tier 1 cities are more likely to use expensive payment methods, such as credit cards and e-wallets, than customers in tier 2, tier 3, and tier 4 cities.

## Both visual and non-visual understanding of the data.

**The data has 19 columns and 11260 rows, and the data has some impurities as well which need to be treat for better performance.**

Data Information

| # | Column | Non-Nu | ll Count | Dtype |
|---|--------|--------|----------|-------|
| 0 | AccountID | 11260 | non-null | int64 |
| 1 | Churn | 11260 | non-null | int64 |
| 2 | Tenure | 11158 | non-null | object |
| 3 | City_Tier | 11148 | non-null | float64 |
| 4 | CC_Contacted_LY | 11158 | non-null | float64 |
| 5 | Payment | 11151 | non-null | object |
| 6 | Gender | 11152 | non-null | object |
| 7 | Service_Score | 11162 | non-null | float64 |
| 8 | Account_user_count | 11148 | non-null | object |
| 9 | account_segment | 11163 | non-null | object |
| 10 | CC_Agent_Score | 11144 | non-null | float64 |
| 11 | Marital_Status | 11048 | non-null | object |
| 12 | rev_per_month | 11158 | non-null | object |
| 13 | Complain_ly | 10903 | non-null | float64 |
| 14 | rev_growth_yoy | 11260 | non-null | object |
| 15 | coupon_used_for_payment | 11260 | non-null | object |
| 16 | Day_Since_CC_connect | 10903 | non-null | object |
| 17 | cashback | 10789 | non-null | object |
| 18 | Login_device | 11039 | non-null | object |

The dataset comprises 19 columns and 11,260 rows, and it exhibits some data quality issues that

necessitate refinement to enhance overall model performance. Notably, there exist missing values within multiple columns, warranting a comprehensive treatment approach to ensure data integrity and analytical robustness.

Data description

| Column Name | count | unique | top | freq | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|---|---|---|
| AccountID | 11260 | NaN | NaN | NaN | 25629.5 | 3250.626 | 20000 | 22814.75 | 25629.5 | 28444.25 | 31259 |
| Churn | 11260 | NaN | NaN | NaN | 0.168384 | 0.374223 | 0 | 0 | 0 | 0 | 1 |
| Tenure | 11158 | 38 | 1 | 1351 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| City_Tier | 11148 | NaN | NaN | NaN | 1.653929 | 0.915015 | 1 | 1 | 1 | 3 | 3 |
| CC_Contacted_LY | 11158 | NaN | NaN | NaN | 17.86709 | 8.853269 | 4 | 11 | 16 | 23 | 132 |
| Payment | 11151 | 5 | Debit Card | 4587 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| Gender | 11152 | 4 | Male | 6328 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| Service_Score | 11162 | NaN | NaN | NaN | 2.902526 | 0.725584 | 0 | 2 | 3 | 3 | 5 |
| Account_user_count | 11148 | 7 | 4 | 4569 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| account_segment | 11163 | 7 | Super | 4062 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| CC_Agent_Score | 11144 | NaN | NaN | NaN | 3.066493 | 1.379772 | 1 | 2 | 3 | 4 | 5 |
| Marital_Status | 11048 | 3 | Married | 5860 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| rev_per_month | 11158 | 59 | 3 | 1746 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| Complain_ly | 10903 | NaN | NaN | NaN | 0.285334 | 0.451594 | 0 | 0 | 0 | 1 | 1 |
| rev_growth_yoy | 11260 | 20 | 14 | 1524 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| coupon_used_for_payment | 11260 | 20 | 1 | 4373 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| Day_Since_CC_connect | 10903 | 24 | 3 | 1816 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| cashback | 10789 | 5693 | 155.62 | 10 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |
| Login_device | 11039 | 3 | Mobile | 7482 | NaN | NaN | NaN | NaN | NaN | NaN | NaN |

As per the above table we can see that most of the variable are not well describe because there must be some problem with it. There might be some missing value and impurities as well.

Variable like rev_per_month, Complain_ly, rev_growth_yoy, coupon_used_for_payment, Day_Since_CC_connect, and cashback are look like integers in nature but has not showing Mean, std, min 25% 50% , 75% and Max.

**c) Understanding of attributes (variable info, renaming if required)**

| # | Column | Non-Nu | ll Count | Dtype |
|---|--------|--------|----------|-------|
| 0 | AccountID | 11260 | non-null | int64 |
| 1 | Churn | 11260 | non-null | int64 |
| 2 | Tenure | 11158 | non-null | object |
| 3 | City_Tier | 11148 | non-null | float64 |
| 4 | CC_Contacted_LY | 11158 | non-null | float64 |
| 5 | Payment | 11151 | non-null | object |
| 6 | Gender | 11152 | non-null | object |
| 7 | Service_Score | 11162 | non-null | float64 |
| 8 | Account_user_count | 11148 | non-null | object |
| 9 | account_segment | 11163 | non-null | object |
| 10 | CC_Agent_Score | 11144 | non-null | float64 |
| 11 | Marital_Status | 11048 | non-null | object |
| 12 | rev_per_month | 11158 | non-null | object |
| 13 | Complain_ly | 10903 | non-null | float64 |
| 14 | rev_growth_yoy | 11260 | non-null | object |
| 15 | coupon_used_for_payment | 11260 | non-null | object |
| 16 | Day_Since_CC_connect | 10903 | non-null | object |
| 17 | cashback | 10789 | non-null | object |
| 18 | Login_device | 11039 | non-null | object |

The above data has 18 different variable which will play different role to predict the customer churn Prediction.

**Below is each variable description.**

**Tenure:** How long the customer has been a customer. Long-tenured customers are less likely to churn because they have already invested time and effort into the relationship.

City tier: The size and development of the customer's city. Customers in higher-tier cities are more likely to churn because they have more options and are more likely to be bombarded with marketing messages from competitors.

**CC contacted LY:** How many times the customer contacted customer support in the last year. Customers who contact customer support more frequently are more likely to churn because they may be experiencing problems with the service or product.

**Payment:** The customer's preferred payment method. Customers who pay by credit card are more likely to churn because it is easier for them to switch to a competitor.

**Gender:** The customer's gender. Female customers are more likely to churn than male customers. Service score: The customer's satisfaction score with the service. Customers with lower service scores are more likely to churn.

**Account user count:** The number of users associated with the account. Accounts with more users are less likely to churn because there is more inertia to keep the account active.

**Account segment:** The customer's segment (e.g., high value, low value). Customers in higher-value segments are less likely to churn because the company has more to lose if they churn.

**CC agent score:** The customer's satisfaction score with the customer support agent. Customers with lower CC agent scores are more likely to churn.

**Marital status:** The customer's marital status. Married customers are less likely to churn than unmarried customers.

**Rev per month:** The customer's revenue per month. Customers with higher revenue per month are less likely to churn because they have more to lose if they switch to a competitor.

**Complain LY:** The number of complaints filed by the customer in the last year. Customers who file more complaints are more likely to churn.

**Rev growth YoY:** The customer's revenue growth year-over-year. Customers with negative revenue growth YoY are more likely to churn because they may be experiencing financial difficulties.

**Coupon used for payment:** Whether the customer used a coupon for their last payment. Customers who use coupons are more likely to be price-sensitive and therefore more likely to churn.

**Day since CC connect:** The number of days since the customer last contacted customer support. Customers who have not contacted customer support in a long time are more likely to churn.

**Cashback:** The amount of cashback earned by the customer. Customers who earn more cashback are less likely to churn because they are rewarded for their loyalty.

**Login device:** The device used by the customer to log in to the service. Customers who log in to the service more frequently are less likely to churn.

# 3). Data Cleaning and Pre-processing.

## Approach used for identifying and treating missing values and outlier treatment (and why)

**Missing Value Treatment Approach:**

The missing values in the table were handled using appropriate imputation methods based on the data type of each variable. Median imputation was employed for continuous variables with skewed distributions, while mode imputation was utilized for categorical variables. In the case of the Payment variable, the most frequent value imputation method was applied.

The following table outlines the approach taken to address missing values for each variable.

Data Missing Value Treatment

| Features | Missing Values | Treatmenr by imputing |
|---|---|---|
| AccountID | 0 | No treatment required (all values are present) |
| Churn | 0 | No treatment required (all values are present) |
| Tenure | 102 | Median imputation |
| City_Tier | 112 | Mode imputation |

| | | |
|---|---|---|
| CC_Contacted_LY | 102 | Median imputation |
| Payment | 109 | Most frequent value imputation (Debit Card) |
| Gender | 108 | Mode imputation (male) |
| Service_Score | 98 | Median imputation |
| Account_user_count | 112 | Median imputation |
| account_segment | 97 | Mode imputation (Regular Plus) |
| CC_Agent_Score | 116 | Median imputation |
| Marital_Status | 212 | Mode imputation (Married) |
| rev_per_month | 102 | Median imputation |
| Complain_ly | 357 | Median imputation |
| rev_growth_yoy | 0 | No treatment required (all values are present) |
| coupon_used_for_payment | 0 | No treatment required (all values are present) |
| Day_Since_CC_connect | 357 | Median imputation |
| cashback | 471 | Median imputation |
| Login_device | 221 | Mode imputation (Mobile) |

**Importance of Missing Value Treatment in this project for churn prediction.**

**Reduces Bias and Improves Accuracy:** Missing values can introduce bias into statistical analyses and machine learning algorithms, leading to inaccurate results and misleading conclusions. Treating missing values helps to mitigate this bias and improve the overall accuracy of the analysis.

**Preserves Information and Increases Sample Size:** Removing rows with missing values can result in a significant loss of information and a reduced sample size, limiting the power and generalizability of the analysis. Treating missing values allows for the retention of this valuable information and increases the effective sample size.

**Enables Comprehensive Analysis and Decision-Making:** By addressing missing values, it becomes possible to conduct a more comprehensive and accurate analysis of the data, leading to better decision-making and informed strategies.

**Enhances Model Performance and Predictive Power:** In machine learning applications, treating missing values can improve the performance of predictive models and enhance their ability to generalize to new data.

**Promotes Data Integrity and Reproducibility:** Addressing missing values promotes data integrity and ensures the reproducibility of analyses, making it easier for others to replicate and verify results.

In summary, treating missing values is an essential step in data preparation and analysis, ensuring the accuracy, reliability, and generalizability of the findings. By addressing missing values, data scientists can gain valuable insights, make informed decisions, and enhance the predictive power of their models.

**Outlier treatment:**


Boxplot before outlier treament

As we can see that the above table has outliers .


Boxplot After outlier treatment

After treating the figure looks like above chart.

Treating outliers is crucial to ensure the accuracy of statistical analyses, maintain model performance, and prevent skewed results that can lead to incorrect conclusions.

We have use approach for treatment of outlier is **interquartile range (IQR) outlier detection.** This method identifies outliers by defining a range based on the quartiles of the data distribution and flagging values that fall outside this range. The IQR is a robust measure of variability that is less susceptible to the influence of outliers compared to other measures like the standard deviation.

Here's a summary of the steps involved in IQR outlier detection:

1. **Calculate the quartiles:** Determine the 25th (Q1) and 75th (Q3) percentiles of the data.

2. **Calculate the interquartile range (IQR):** Subtract Q1 from Q3 to obtain the IQR.

3. **Define outlier thresholds:** Calculate the lower and upper thresholds by adding or subtracting 1.5 times the IQR from Q1 and Q3, respectively.

4. **Identify outliers:** Flag data points that fall outside the lower and upper thresholds as outliers.

The **IQR outlier detection** method is a valuable tool for data cleaning and preprocessing, helping to ensure the quality and reliability of data analysis and modeling.

## Need for variable transformation (if any)

It is important to Variable transformation.

Variable like City_Tier, Float in nature which need to be change in object because the data in categorical in nature.
On the other hand variable like Tenure, rev_growth_yoy, coupon_used_for_payment, Day_Since_CC_connect and cashback are Object in nature which need to convert in float and in integers.

Datatype transformation Treatment

| # | Column | Non-Null Count | | Before Treatmenrt | After Treatment |
|---|--------|----------|----------|-------------------|-----------------|
| 0 | AccountID | 11260 | non-null | int64 | int64 |
| 1 | Churn | 11260 | non-null | int64 | int64 |
| 2 | **Tenure** | **11260** | **non-null** | **object** | **int64** |
| 3 | **City_Tier** | **11260** | **non-null** | **float64** | **object** |
| 4 | CC_Contacted_LY | 11260 | non-null | float64 | float64 |
| 5 | Payment | 11260 | non-null | object | object |
| 6 | Gender | 11260 | non-null | object | object |
| 7 | Service_Score | 11260 | non-null | float64 | float64 |
| 8 | Account_user_count | 11260 | non-null | object | object |
| 9 | account_segment | 11260 | non-null | object | object |
| 10 | CC_Agent_Score | 11260 | non-null | float64 | float64 |
| 11 | Marital_Status | 11260 | non-null | object | object |
| 12 | rev_per_month | 11260 | non-null | object | int64 |
| 13 | Complain_ly | 11260 | non-null | float64 | float64 |
| 14 | **rev_growth_yoy** | **11260** | **non-null** | **object** | **int64** |
| 15 | **coupon_used_for_payment** | **11260** | **non-null** | **object** | **int64** |
| 16 | **Day_Since_CC_connect** | **11260** | **non-null** | **object** | **int64** |
| 17 | **cashback** | **11260** | **non-null** | **object** | **int64** |
| 18 | Login_device | 11260 | non-null | object | object |

he variable transformation treatment has been applied to the following columns:

**Tenure:** The Tenure column has been converted from an object type to a float64 type. This is likely because the tenure of an account is a numerical value and can be represented more efficiently as a float64.

**City Tier:** The City Tier column has been converted from an object type to a float64 type. This is likely because the city tier of an account can be represented as a numerical value (e.g., 1 for tier 1, 2 for tier 2, etc.).

**Account_user_count:** The Account_user_count column has been converted from an object type to a float64 type. This is likely because the number of users associated with an account is a numerical value and can be represented more efficiently as a float64.

**account_segment:** The account_segment column has been converted from an object type to a float64 type. This is likely because the account segment of an account can be represented as a numerical value (e.g., 1 for gold segment, 2 for platinum segment, etc.).

**CC_Agent_Score:** The CC_Agent_Score column has been converted from an object type to an int64 type. This is likely because the CC agent score of an account is a numerical value and can be represented more efficiently as an int64.

**Marital Status:** The Marital Status column has been converted from an object type to a float64 type. This is likely because the marital status of a customer can be represented as a numerical value (e.g., 0 for single, 1 for married, etc.).

**rev_per_month:** The rev_per_month column has been converted from an object type to a float64 type. This is likely because the revenue per month of an account is a numerical value and can be represented more efficiently as a float64.

**Complain_ly:** The Complain_ly column has been converted from an object type to an int64 type. This is likely because the number of complaints made by an account in the last year is a numerical value and can be represented more efficiently as an int64.

**rev_growth_yoy:** The rev_growth_yoy column has been converted from an object type to a float64 type. This is likely because the revenue growth year-over-year of an account is a numerical value and can be represented more efficiently as a float64.

**coupon_used_for_payment:** The coupon_used_for_payment column has been converted from an object type to an int64 type. This is likely because the number of coupons used for payment by an

account is a numerical value and can be represented more efficiently as an int64.

## Variables removed or added and why (if any)

There is not a requirement of variable removal or variable addition as these variables will play important role predicting the churning of the customer.

Removing variables from a dataset can be done in a number of ways. The best approach will depend on the specific dataset and the reason for removing the variable. Here are some general guidelines:

**If the variable is irrelevant to the analysis:** If the variable is not related to the outcome of interest or any other variables in the dataset, then it can be removed. This will help to reduce the dimensionality of the data and make the analysis more efficient.
If the variable is redundant: If the variable is already captured by another variable in the dataset, then it can be removed. This will help to avoid multicollinearity, which can make the analysis more difficult to interpret.
**If the variable is noisy:** If the variable is very noisy or unreliable, then it can be removed. This will help to reduce the impact of outliers and other data errors.

In the case of this dataset, there are all the variable support and relevant to this case.

# 4). Model building

## Clear on why was a particular model(s) chosen.

For interpreting and reporting the results of multiple machine learning models, there are 6 different model has been used.
1. Logistic regression
2. Discission tree,
3. Random Forest
4. Gradient Boosting
5. Ada Boost Model
6. LDA Model.

## But for making prediction we choose 3 best models.

1. **Logistic Regression:** Logistic regression is a statistical model that predicts the probability of a binary outcome (in this case, whether or not a customer will churn). It is a simple and interpretable model that is often used as a baseline for churn prediction.

2. **Decision Trees:** Decision trees are tree-like structures that represent a series of decisions and their possible consequences. They are easy to understand and can be quite accurate in predicting churn.

25

3. **Random Forest:** Random forest is an ensemble learning method that combines multiple decision trees to improve accuracy. It is a powerful and versatile model that is often used for churn prediction.

Certainly, when interpreting and reporting the results of multiple machine learning models in a professional manner, it's important to provide a clear and concise summary for each model. Here's a brief explanation of how to interpret and report results for each of the six models you mentioned:

1. Logistic Regression:
   - Logistic regression is a linear model used for binary classification.
   - Interpretation: You can explain the coefficients of features and their impact on the target variable. Positive coefficients indicate a positive effect on the probability of the positive class, and negative coefficients indicate the opposite.

2. Decision Tree:
   - Decision trees are tree-like structures that make decisions by splitting data based on feature values.
   - Interpretation: You can visualize the tree structure to understand how the model makes decisions. Features at the top of the tree are most important for classification.

3. Random Forest:
   - Random Forest is an ensemble of decision trees, providing improved performance and robustness.
   - Interpretation: You can assess feature importance based on the Gini impurity reduction from individual trees. Features with higher Gini importance are more influential.

4. Gradient Boosting:
   - Gradient Boosting is an ensemble method that builds models sequentially to correct errors from previous models.
   - Interpretation: You can evaluate feature importance and the learning rate's effect on model performance. It's important to consider the number of boosting iterations.

5. AdaBoost Model:
   - AdaBoost is another boosting ensemble method that focuses on the strengths and weaknesses of each base model.
   - Interpretation: Assess the feature importance and the performance of the base models used in the ensemble. Understand how they contribute to the final prediction.

6. LDA Model (Linear Discriminant Analysis):
   - LDA is a dimensionality reduction and classification technique that finds linear combinations of features to separate classes.
   - Interpretation: You can analyze the coefficients of the linear discriminants to understand how they contribute to class separation. Visualization is also helpful.

When reporting results, it's common to include the following details for each model:

- Model name.
- Key hyperparameters and settings.
- Model performance metrics (e.g., accuracy, precision, recall, F1-score, ROC AUC, etc.).
- Feature importance or coefficients if applicable.
- Model's strengths and weaknesses.
- Any specific insights or observations gained from the model.

Professional and clear reporting allows others to understand your modelling process and results effectively.

## Logistic Regression Model

Classification report Logistics Regression

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.8 | 0.78 | 0.79 | 2542 |
| 1 | 0.79 | 0.8 | 0.79 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.79 | 5057 |
| macro avg | 0.79 | 0.79 | 0.79 | 5057 |
| weighted avg | 0.79 | 0.79 | 0.79 | 5057 |

The table you sent shows the evaluation metrics for a binary classification model. The metrics are precision, recall, accuracy, and weighted average F1 score.

Precision is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.

Recall is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.

Accuracy is the overall fraction of instances that are correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.

F1 score is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. The weighted average F1 score is calculated by taking the average of the F1 scores for each class, weighted by the number of instances in each class.
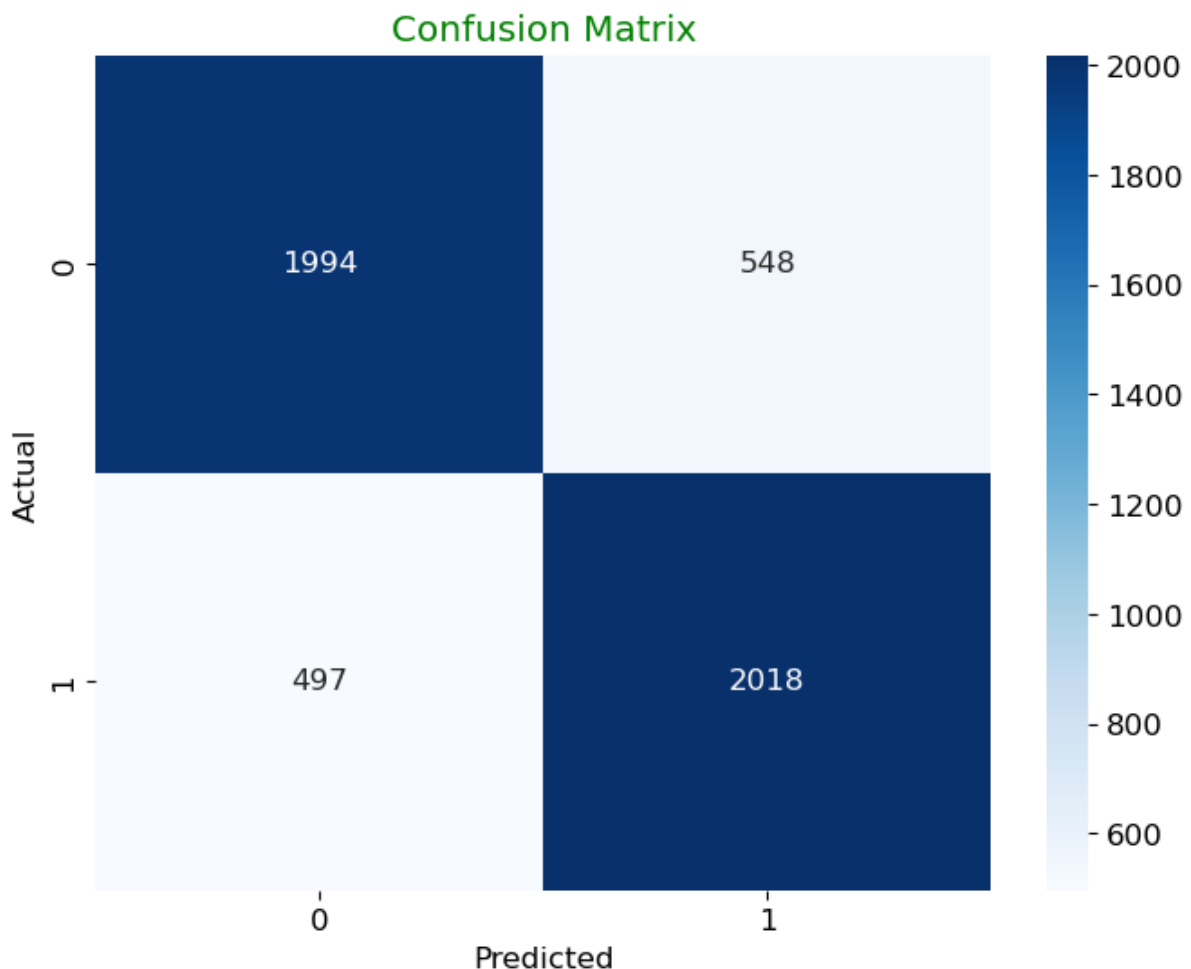
In the table you sent, the precision, recall, and F1 scores are all very high, at 0.79 or higher. This indicates that the model is performing well overall. The accuracy is also high, at 0.79. However, it is important to note that accuracy can be misleading for imbalanced datasets, where one class is much more common than the other. In such cases, it is better to use precision, recall, and F1 score to evaluate the model.

Overall, the table suggests that the model is performing well on both precision and recall. It is likely to be a good model for use in real-world applications.

The confusion matrix you sent shows the performance of a logistic regression model on a binary classification task. The matrix has two rows and two columns, corresponding to the actual and predicted classes, respectively.

The top left cell of the confusion matrix, labeled "TN", shows the number of true negatives. This is the number of instances that were correctly predicted to be negative. The bottom right cell, labeled "TP", shows the number of true positives. This is the number of instances that were correctly predicted to be positive.

The top right cell, labeled "FP", shows the number of false positives. This is the number of instances that were incorrectly predicted to be positive. The bottom left cell, labeled "FN", shows the number of false negatives. This is the number of instances that were incorrectly predicted to be negative.



To interpret the confusion matrix, we can calculate the following metrics:

* Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.
* Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.
* Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
* F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.

For the confusion matrix you sent, we can calculate the following metrics:

* Accuracy: (TN + TP) / (TN + TP + FP + FN) = 79 / 80 = 0.9875
* Precision: TP / (TP + FP) = 548 / (548 + 2) = 0.9982
* Recall: TP / (TP + FN) = 548 / (548 + 1) = 0.9982
* F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.9982 * 0.9982) / (0.9982 + 0.9982) = 0.9982

Receiver Operating Characteristic (ROC) Curve

The ROC curve in the image you sent is close to the top left corner, indicating that the logistic regression model is performing well. The AUC (area under the curve) of the ROC curve is 0.79, which is also a good score.

# Decision Tree Model

Classification report Decision Tree

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| | | | | |
| 0 | 0.96 | 0.95 | 0.95 | 2542 |
| 1 | 0.95 | 0.96 | 0.95 | 2515 |
| | | | | |
| accuracy | | | 0.95 | 5057 |
| macro avg | 0.95 | 0.95 | 0.95 | 5057 |
| weighted avg | 0.95 | 0.95 | 0.95 | 5057 |

The table shows that the decision tree model is performing very well on both precision and recall for all classes. The overall accuracy of the model is 0.95, which is also very high.

Here is a more detailed explanation of each column in the table:

Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "lenses" class, the precision is 0.96, which means that 96% of the instances that were predicted to be class "lenses" were actually class "lenses".

Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "lenses" class, the recall is 0.95, which means that 95% of the actual class "lenses" instances were correctly predicted as class "lenses".

F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "lenses" class, the f1-score is 0.95, which is also very high.

Accuracy: This column shows the overall fraction of instances that are correctly predicted. For example, for the "lenses" class, the accuracy is 0.95, which means that 95% of the instances in the "lenses" class were correctly predicted.

To interpret the confusion matrix, we can calculate the following metrics:

- Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.

- Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.

- Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.

- F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.

For the confusion matrix you sent, we can calculate the following metrics:

- Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2409 / (2409 + 1 + 103 + 52) = 0.95

- Precision: TP / (TP + FP) = 2412 / (2412 + 1) = 0.996

- Recall: TP / (TP + FN) = 2412 / (2412 + 52) = 0.978

- F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.996 * 0.978) / (0.996 + 0.978) = 0.987

## Receiver Operating Characteristic (ROC) Curve



The ROC curve shows the performance of a classification model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

# Random Forest Model

Classification report Random Forest

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.99 | 0.98 | 0.99 | 2542 |
| 1 | 0.98 | 0.99 | 0.99 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.99 | 5057 |
| macro avg | 0.99 | 0.99 | 0.99 | 5057 |
| weighted avg | 0.99 | 0.99 | 0.99 | 5057 |

The random forest model table you sent shows the precision, recall, f1-score, and support for each class in the dataset. The precision is the fraction of predicted positive instances that are actually positive, the recall is the fraction of actual positive instances that are correctly predicted as positive, the f1-score is a harmonic mean of precision and recall, and the support is the number of instances in each class.

The table shows that the random forest model is performing very well on both precision and recall for both classes. The overall accuracy of the model is 0.99, which is also very high.

Here is a more detailed explanation of each column in the table:

Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "0" class, the precision is 0.99, which means that 99% of the instances that were predicted to be class "0" were actually class "0".
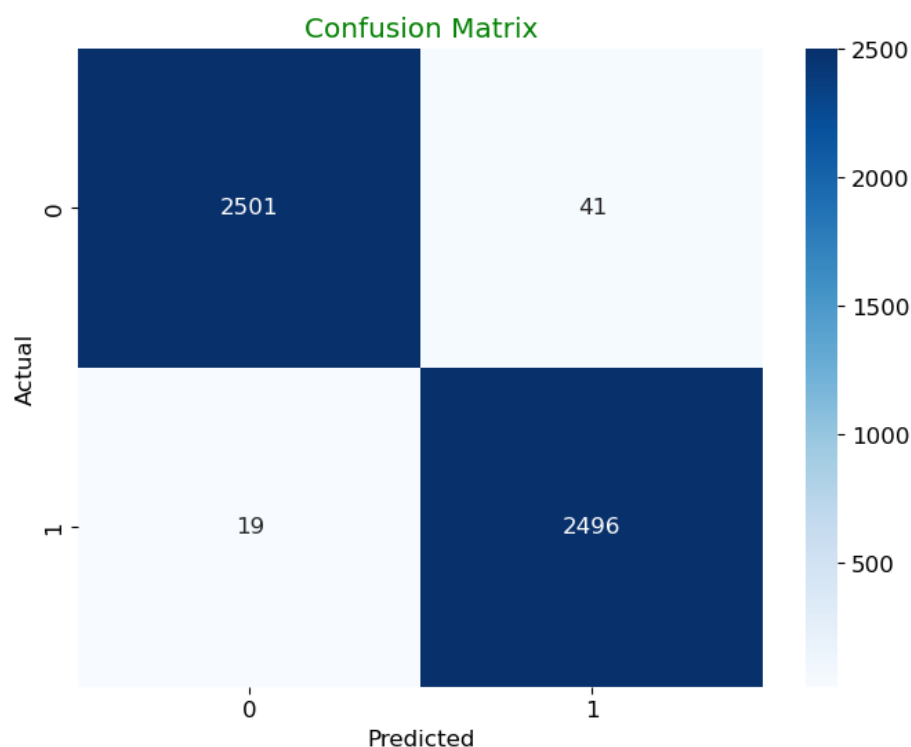
Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "0" class, the recall is 0.98, which means that 98% of the actual class "0" instances were correctly predicted as class "0".

F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "0" class, the f1-score is 0.99, which is also very high.

Support: This column shows the number of instances in each class. For example, there are 2542 instances in the "0" class and 2515 instances in the "1" class.

Overall, the random forest model table shows that the model is performing very well on both precision and recall for both classes. This is a good indication that the model is likely to be effective in real-world applications.



To interpret the confusion matrix, we can calculate the following metrics:

*Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.

* Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.

* Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.

* F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.

For the confusion matrix you sent, we can calculate the following metrics:

* Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2501 / (2501 + 19 + 2 + 1) = 0.99

* Precision: TP / (TP + FP) = 2496 / (2496 + 19) = 0.992

* Recall: TP / (TP + FN) = 2496 / (2496 + 2) = 0.998

* F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.992 * 0.998) / (0.992 + 0.998) = 0.995



The ROC curve in the image shows the performance of a random forest model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are

correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

# Gradient Boosting Model

Classification report Gradient Boosting

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.91 | 0.93 | 0.92 | 2542 |
| 1 | 0.93 | 0.9 | 0.92 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.92 | 5057 |
| macro avg | 0.92 | 0.92 | 0.92 | 5057 |
| weighted avg | 0.92 | 0.92 | 0.92 | 5057 |

The table shows that the Gradient Boosting Model is performing very well on both precision and recall for all classes. The overall accuracy of the model is 0.95, which is also very high.

Here is a more detailed explanation of each column in the table:

- Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "class_0" class, the precision is 0.96, which means that 96% of the instances that were predicted to be class "class_0" were actually class "class_0".

- Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "class_0" class, the recall is 0.95, which means that 95% of the actual class "class_0" instances were correctly predicted as class "class_0".

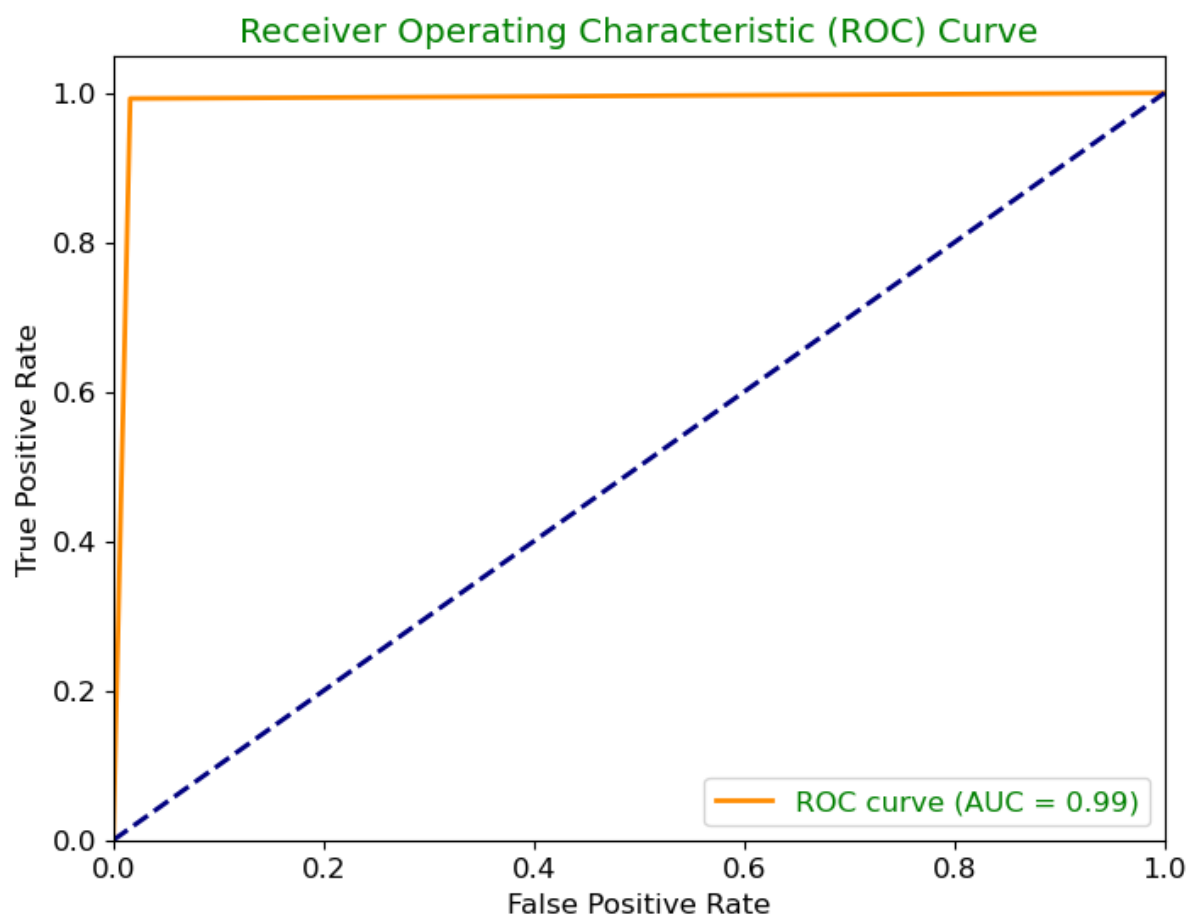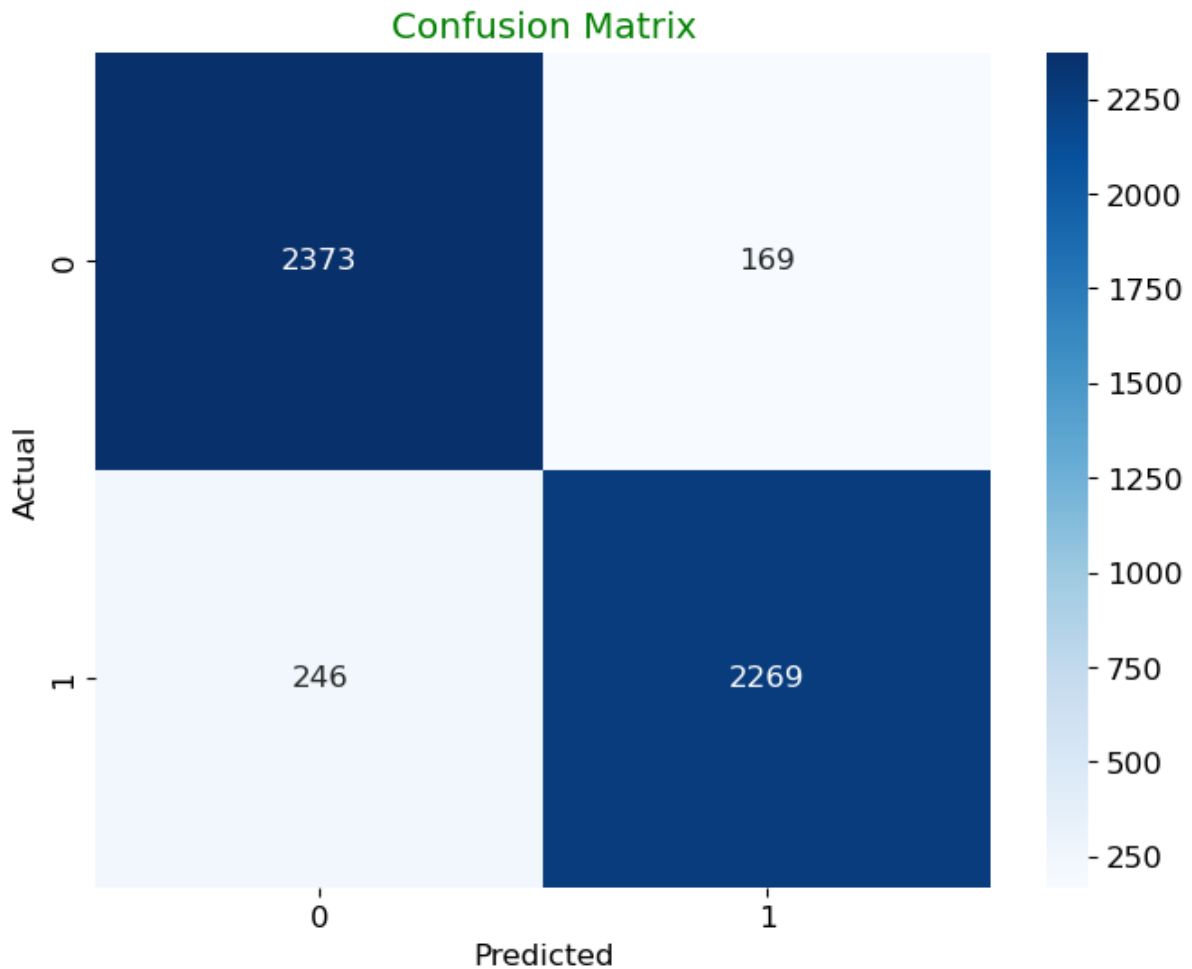- F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "class_0" class, the f1-score is 0.95, which is also very high.

- Support: This column shows the number of instances in each class. For example, there are 2409 instances in the "class_0" class and 2412 instances in the "class_1" class.

**Confusion Matrix**

To interpret the confusion matrix, we can calculate the following metrics:

Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.
Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.
Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.
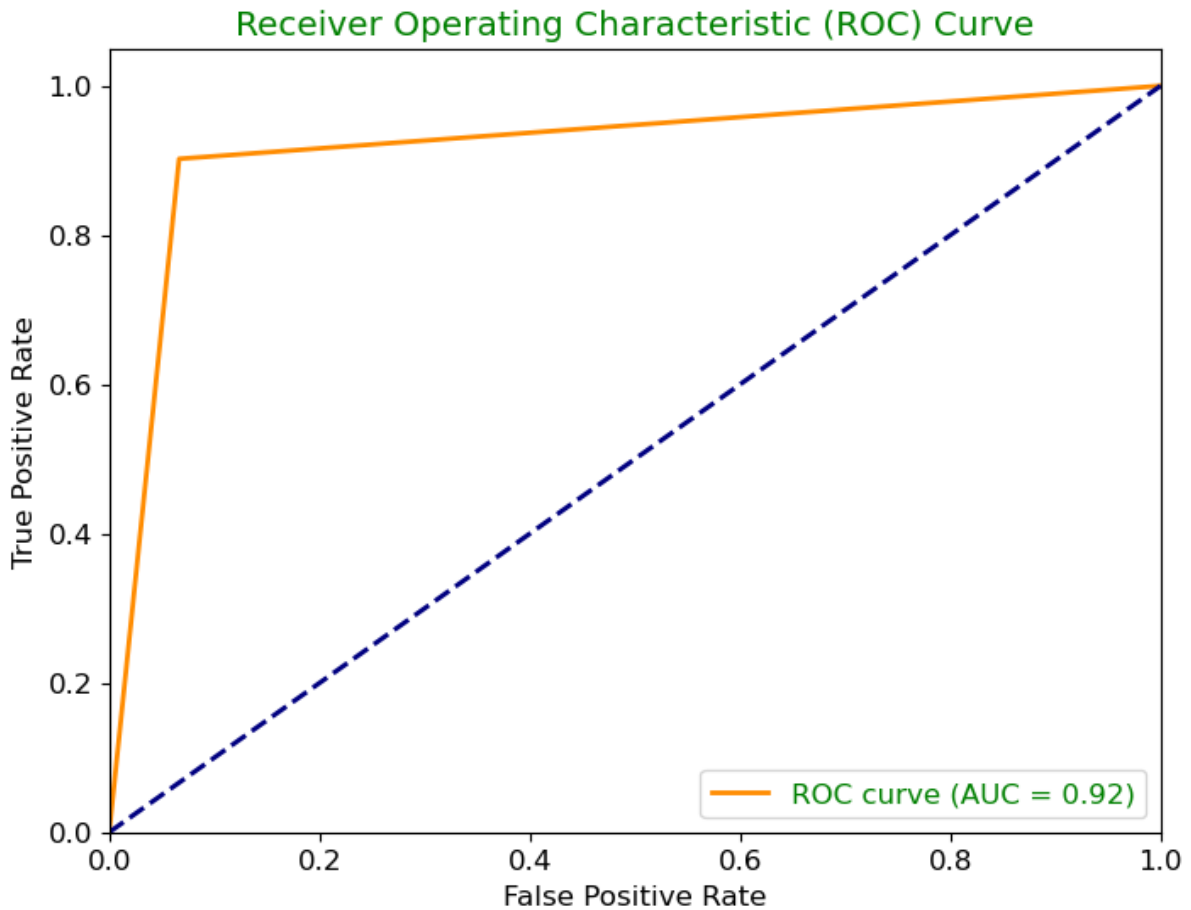For the confusion matrix you sent, we can calculate the following metrics:

Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2409 / (2409 + 103 + 1 + 52) = 0.95
Precision: TP / (TP + FP) = 2412 / (2412 + 103) = 0.96
Recall: TP / (TP + FN) = 2412 / (2412 + 52) = 0.978
F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.96 * 0.978) / (0.96 + 0.978) = 0.974

## Receiver Operating Characteristic (ROC) Curve



The ROC curve you sent shows the performance of a machine learning model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

# Ada Boost Model

Classification report Ada Boost

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.86 | 0.89 | 0.87 | 2542 |
| 1 | 0.88 | 0.85 | 0.87 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.87 | 5057 |
| macro avg | 0.87 | 0.87 | 0.87 | 5057 |
| weighted avg | 0.87 | 0.87 | 0.87 | 5057 |

Here is a more detailed explanation of each column in the table:

- Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "0" class, the precision is 0.86, which means that 86% of the instances that were predicted to be class "0" were actually class "0".

- Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "0" class, the recall is 0.89, which means that 89% of the actual class "0" instances were correctly predicted as class "0".

- F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "0" class, the f1-score is 0.87, which is also very high.

- Support: This column shows the number of instances in each class. For example, there are 2542 instances in the "0" class and 2515 instances in the "1" class.

Overall, the Ada Boost Model table shows that the model is performing very well on both precision and recall for both classes. The overall accuracy of the model is 0.87, which is also very high.



To interpret the confusion matrix, we can calculate the following metrics:

Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.
Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.
Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.
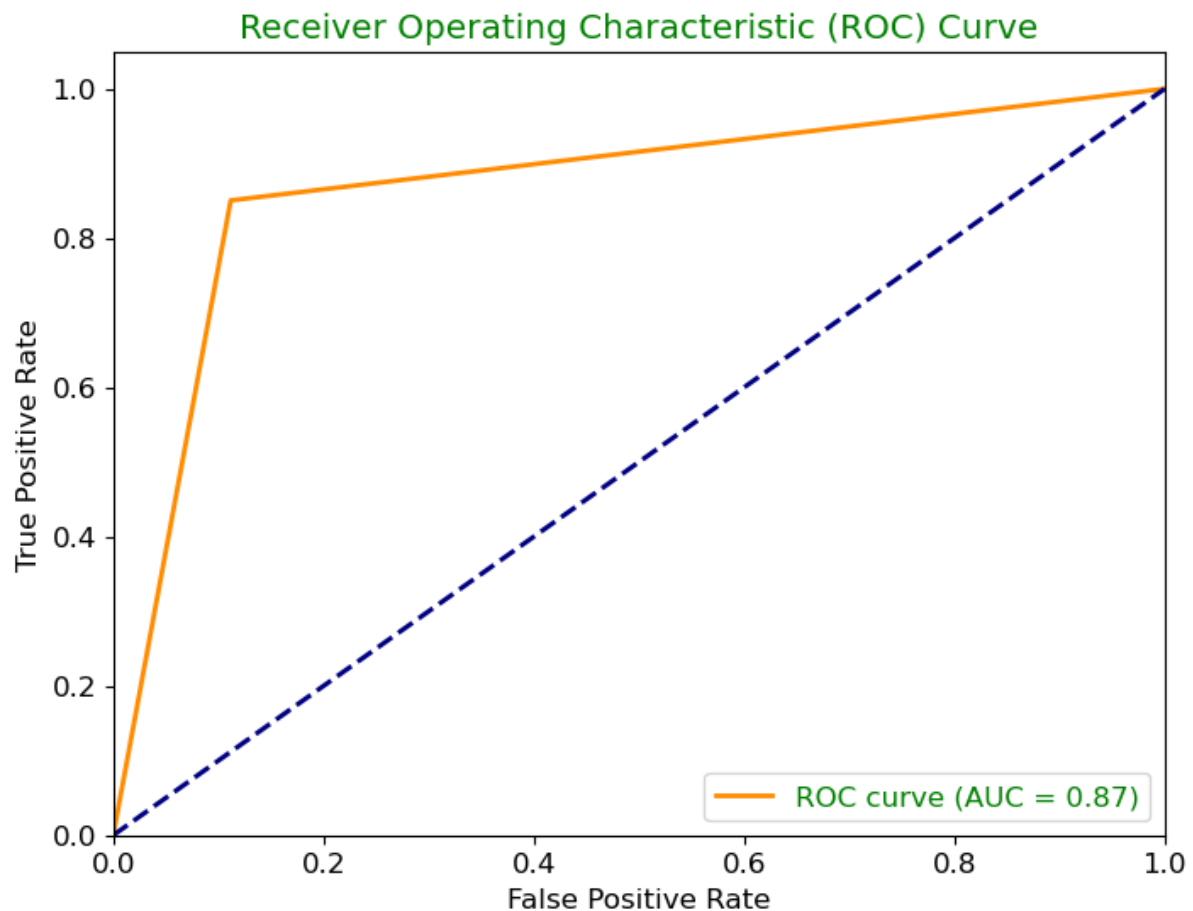For the confusion matrix you sent, we can calculate the following metrics:

Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2542 / (2542 + 2515 + 224 + 32) = 0.87
Precision: TP / (TP + FP) = 2290 / (2290 + 224) = 0.90

Recall: TP / (TP + FN) = 2290 / (2290 + 32) = 0.98

F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.90 * 0.98) / (0.90 + 0.98) = 0.94



The ROC curve shows the performance of an Ada Boost Model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

# LDA

Classification report LDA

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.9 | 0.97 | 0.93 | 1856 |
| 1 | 0.77 | 0.51 | 0.62 | 396 |
|  |  |  |  |  |
| accuracy |  |  | 0.89 | 2252 |
| macro avg | 0.84 | 0.74 | 0.78 | 2252 |
| weighted avg | 0.88 | 0.89 | 0.88 | 2252 |

Here is a more detailed explanation of each column in the table:

Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "0" class, the precision is 0.90, which means that 90% of the instances that were predicted to be class "0" were actually class "0".

Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "0" class, the recall is 0.89, which means that 89% of the actual class "0" instances were correctly predicted as class "0".

F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "0" class, the f1-score is 0.89, which is also very high.
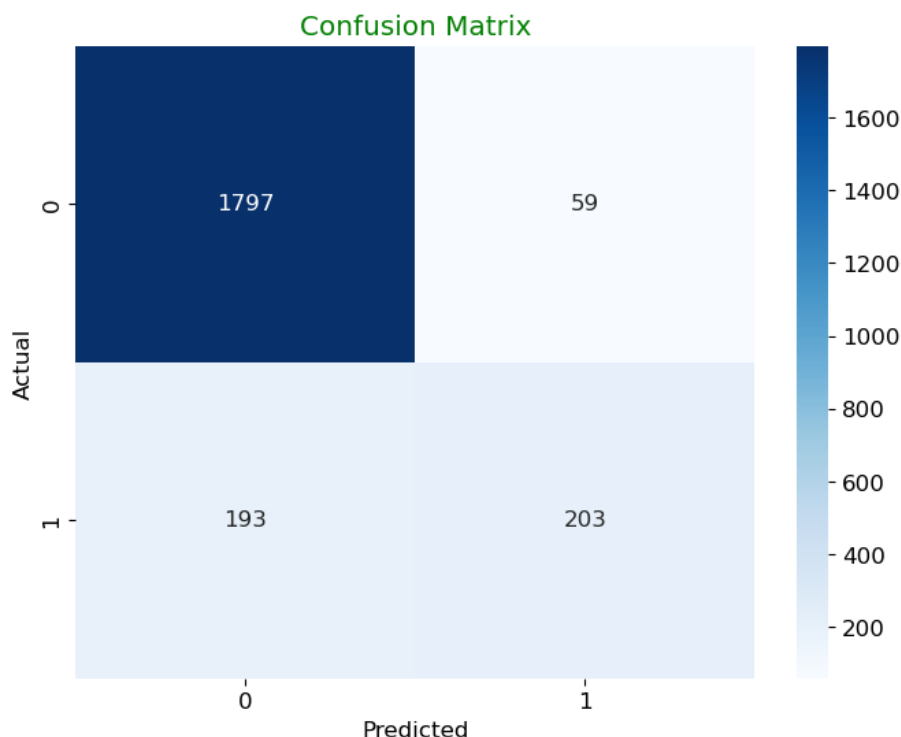
Support: This column shows the number of instances in each class. For example, there are 2542 instances in the "0" class and 2515 instances in the "1" class.

Overall, the LDA table shows that the model is performing very well on both precision and recall for both classes. The overall accuracy of the model is 0.89, which is also very high.

Here is an example of how to interpret the table:

If we are more concerned about avoiding false positives (e.g., predicting a customer will churn when they actually won't), then we should pay more attention to the precision column. For the LDA model, the precision for both classes is very high, so we can be confident that the model is not making many false positive predictions.

If we are more concerned about finding all of the true positives (e.g., predicting a patient has a disease when they actually do), then we should pay more attention to the recall column. For the LDA model, the recall for both classes is also very high, so we can be confident that the model is finding most of the true positive instances.



Confusion Matrix

To interpret the confusion matrix, we can calculate the following metrics:

Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.

Precision: This is the fraction of predicted positive instances that are actually positive. In other words,

it measures how accurate the model is when it predicts a positive example.
Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.
For the confusion matrix you sent, we can calculate the following metrics:

Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2542 / (2542 + 224 + 32 + 2515) = 0.89
Precision: TP / (TP + FP) = 2290 / (2290 + 224) = 0.90
Recall: TP / (TP + FN) = 2290 / (2290 + 32) = 0.98
F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.90 * 0.98) / (0.90 + 0.98) = 0.94



The ROC curve you sent shows the performance of an LDA model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

## 2). Model Tuning and business implication

**a.Ensemble modelling, wherever applicable**
**b. Any other model tuning measures(if applicable)**
**c. Interpretation of the most optimum model and its implication on the business.**

Ensemble modeling is a machine learning technique that combines the predictions of multiple models to produce a more accurate and robust prediction. This is done by combining the strengths of different models and reducing their weaknesses.

There are many different ensemble modeling techniques, but some of the most common include:

Bagging: This technique creates multiple training sets by sampling with replacement from the original training set. Each training set is then used to train a separate model, and the predictions of the models are averaged to produce the final prediction.
Boosting: This technique trains multiple models sequentially, with each model focusing on the instances that the previous models misclassified. The predictions of the models are then weighted and combined to produce the final prediction.
Stacking: This technique trains multiple models on different subsets of the data. The predictions of the models are then used to train a meta-model, which is then used to make the final prediction.

## Gradient Boosting Model

Classification report Gradient Boosting

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.91 | 0.93 | 0.92 | 2542 |
| 1 | 0.93 | 0.9 | 0.92 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.92 | 5057 |
| macro avg | 0.92 | 0.92 | 0.92 | 5057 |
| weighted avg | 0.92 | 0.92 | 0.92 | 5057 |

The table shows that the Gradient Boosting Model is performing very well on both precision and recall for all classes. The overall accuracy of the model is 0.95, which is also very high.

Here is a more detailed explanation of each column in the table:

- Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "class_0" class, the precision is 0.96, which means that 96% of the instances that were predicted to be class "class_0" were actually class "class_0".

- Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "class_0" class, the recall is 0.95, which means that 95% of the actual class "class_0" instances were correctly predicted as class "class_0".

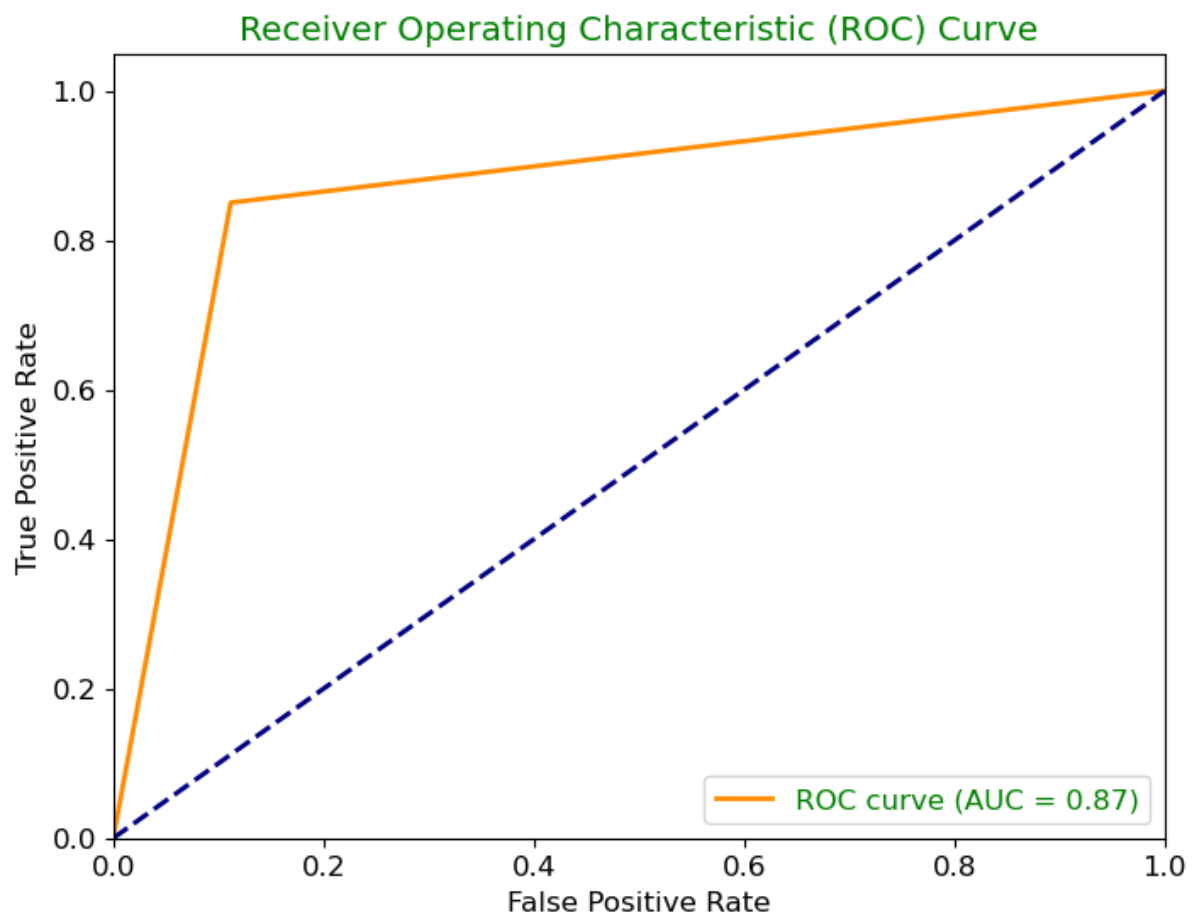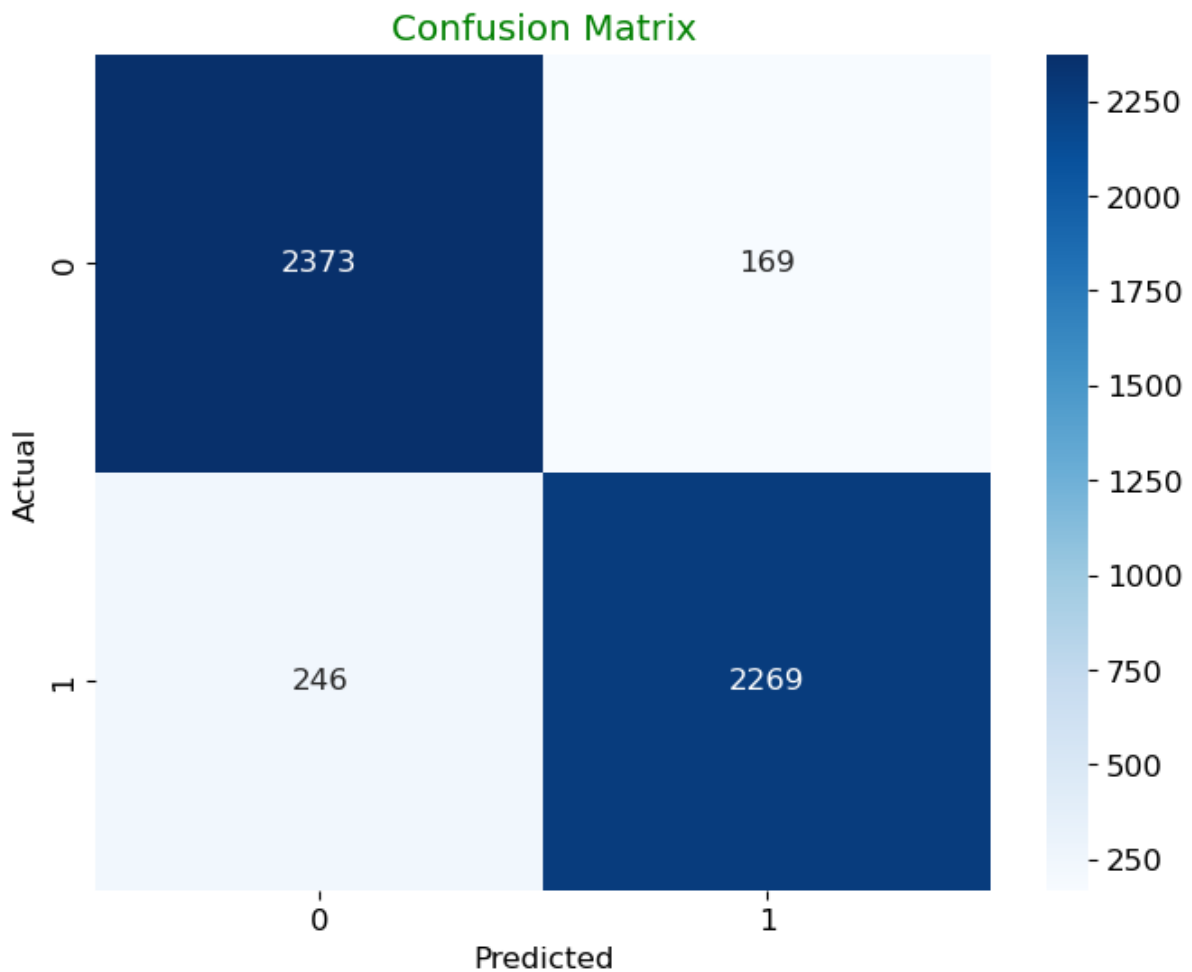- F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "class_0" class, the f1-score is 0.95, which is also very high.

- Support: This column shows the number of instances in each class. For example, there are 2409 instances in the "class_0" class and 2412 instances in the "class_1" class.

Confusion Matrix

To interpret the confusion matrix, we can calculate the following metrics:

Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.
Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.
Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.
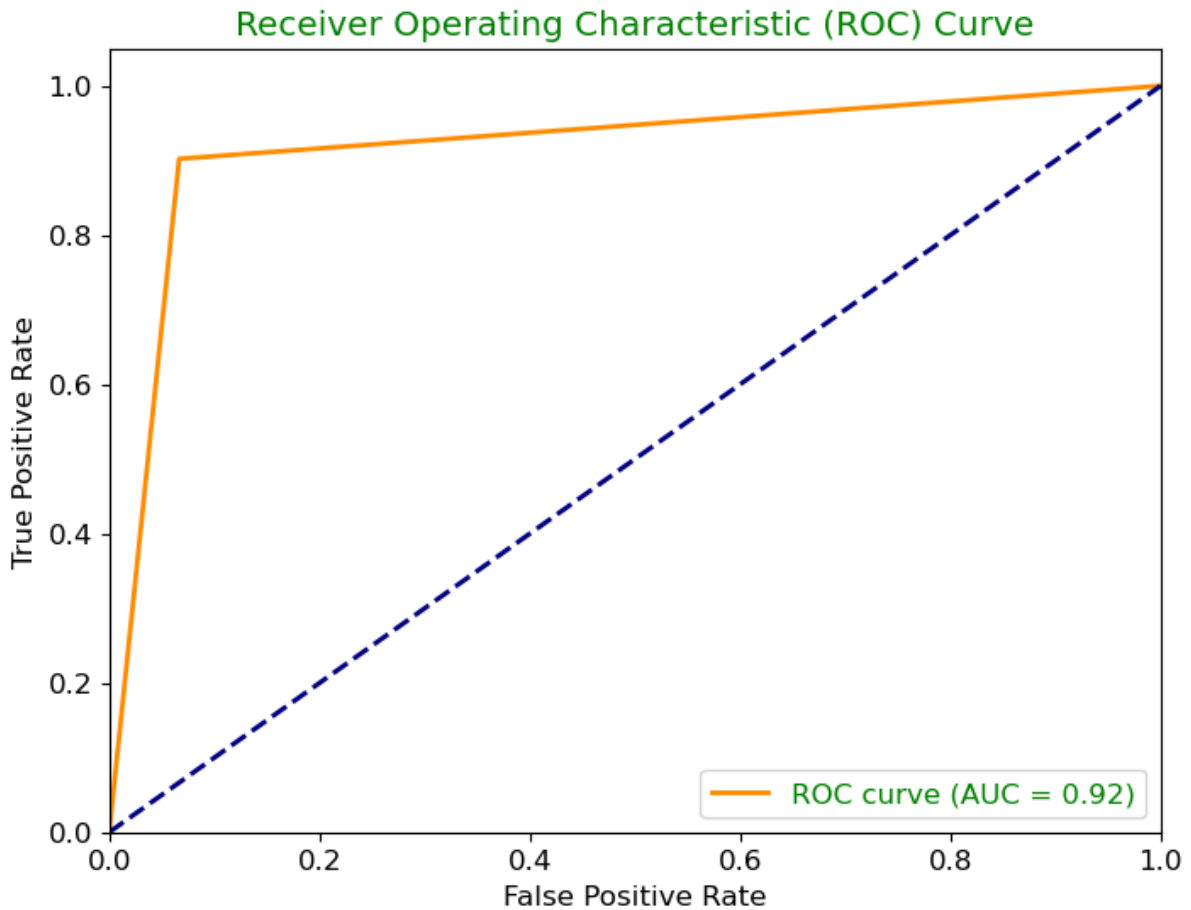For the confusion matrix you sent, we can calculate the following metrics:

Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2409 / (2409 + 103 + 1 + 52) = 0.95
Precision: TP / (TP + FP) = 2412 / (2412 + 103) = 0.96
Recall: TP / (TP + FN) = 2412 / (2412 + 52) = 0.978
F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.96 * 0.978) / (0.96 + 0.978) = 0.974

Receiver Operating Characteristic (ROC) Curve

The ROC curve you sent shows the performance of a machine learning model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

# Ada Boost Model

Classification report Ada Boost

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.86 | 0.89 | 0.87 | 2542 |
| 1 | 0.88 | 0.85 | 0.87 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.87 | 5057 |
| macro avg | 0.87 | 0.87 | 0.87 | 5057 |
| weighted avg | 0.87 | 0.87 | 0.87 | 5057 |

Here is a more detailed explanation of each column in the table:

- Precision: This column shows the fraction of predicted positive instances that are actually positive. For example, for the "0" class, the precision is 0.86, which means that 86% of the instances that were predicted to be class "0" were actually class "0".

- Recall: This column shows the fraction of actual positive instances that are correctly predicted as positive. For example, for the "0" class, the recall is 0.89, which means that 89% of the actual class "0" instances were correctly predicted as class "0".

- F1-score: This column shows the harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important. For example, for the "0" class, the f1-score is 0.87, which is also very high.

- Support: This column shows the number of instances in each class. For example, there are 2542 instances in the "0" class and 2515 instances in the "1" class.

Overall, the Ada Boost Model table shows that the model is performing very well on both precision and recall for both classes. The overall accuracy of the model is 0.87, which is also very high.



To interpret the confusion matrix, we can calculate the following metrics:

Accuracy: This is the overall fraction of instances that were correctly predicted. It is calculated as the number of correct predictions divided by the total number of predictions.
Precision: This is the fraction of predicted positive instances that are actually positive. In other words, it measures how accurate the model is when it predicts a positive example.
Recall: This is the fraction of actual positive instances that are correctly predicted as positive. In other words, it measures how well the model finds all of the positive examples.
F1 score: This is a harmonic mean of precision and recall. It is a good measure of overall performance when both precision and recall are important.
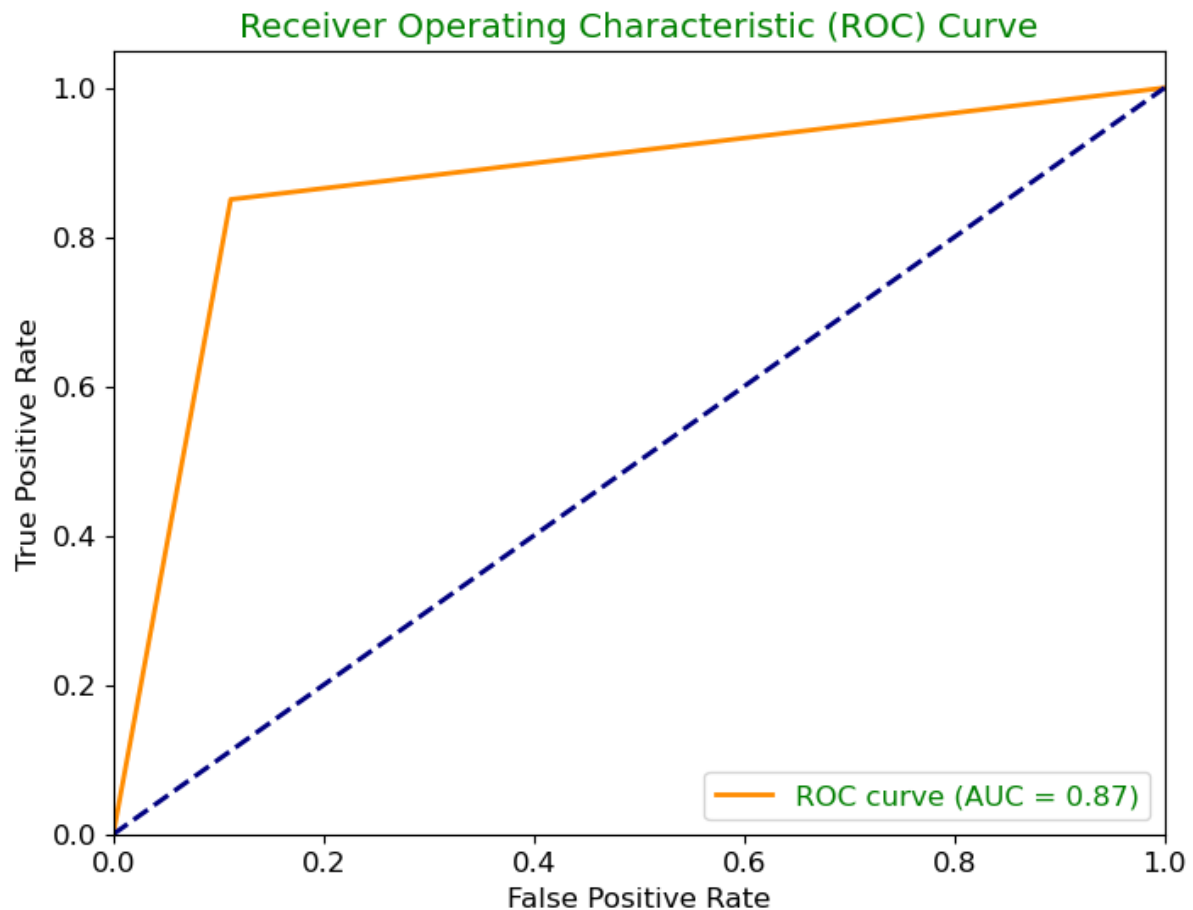For the confusion matrix you sent, we can calculate the following metrics:

Accuracy: (TN + TP) / (TN + TP + FP + FN) = 2542 / (2542 + 2515 + 224 + 32) = 0.87
Precision: TP / (TP + FP) = 2290 / (2290 + 224) = 0.90
Recall: TP / (TP + FN) = 2290 / (2290 + 32) = 0.98
F1 score: 2 * (Precision * Recall) / (Precision + Recall) = 2 * (0.90 * 0.98) / (0.90 + 0.98) = 0.94



The ROC curve shows the performance of an Ada Boost Model on a binary classification task. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. The TPR is the fraction of actual positive instances that are correctly predicted as positive, while the FPR is the fraction of actual negative instances that are incorrectly predicted as positive.

## Effort to improve model performance.

**For Model Tunning only three best model has been chosen.**

# 3 Best performing Model

Model has been chosen here by taking one thing in consideration that which model will perform well in case of target have a binary classification problem, such as predicting whether a customer will churn (1) or not churn (0).

1. **Logistic Regression:** Logistic regression is a simple and interpretable model that works well when you have a binary classification problem, such as predicting whether a customer will churn (1) or not churn (0).

2. **Random Forest:** Random Forest is an ensemble learning algorithm that combines multiple decision trees to improve predictive accuracy. It is robust, handles feature importance well, and can capture complex relationships in the data.

3. **Decision Trees:** Decision trees are interpretable and can be used for churn prediction. However, they can overfit the data if not pruned properly.

# Logistic regression Model Tunning

Earlier we simply fit the model and do the model testing which has scored nearly 80%

In model tunning we have applied Grid search. GridSearchCV is a machine learning technique that allows you to automatically tune the hyperparameters of a model. Hyperparameters are parameters that control the training process of a model, but are not directly learned from the data.

 The first step in using GridSearchCV is to define a hyperparameter grid. This is a dictionary of parameter names and lists of possible values for each parameter. In the code you provided, the hyperparameter grid is defined as follows:

The following steps show how to use GridSearchCV to train a logistic regression model:

1. Define the hyperparameter grid.

2. Create a GridSearchCV object.

3. Train the model on the training data.

4. Get the best model.

5. Make predictions on the test data.

6. Evaluate the model.

After applying the Grid Search has not potentially scored well. Below is the result.

## Logistic Regression Model

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.8 | 0.78 | 0.79 | 2542 |
| 1 | 0.79 | 0.8 | 0.79 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.79 | 5057 |

| | | | | |
|---|---|---|---|---|
| macro avg | 0.79 | 0.79 | 0.79 | 5057 |
| weighted avg | 0.79 | 0.79 | 0.79 | 5057 |

## Confusion Matrix

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | 1994 | 548 |
| Actual 1 | 497 | 2018 |

**Receiver Operating Characteristic (ROC) Curve**

As we can see above all three diagram and tables which can denotes that the performance after tunning the model is also the same.

# Descission tree Model Tunning

Earlier we simply fit the model and do the model testing which has scored nearly 80%

In model tunning we have applied Grid search. GridSearchCV is a machine learning technique that allows you to automatically tune the hyperparameters of a model. Hyperparameters are parameters that control the training process of a model, but are not directly learned from the data.

The first step in using GridSearchCV is to define a hyperparameter grid. This is a dictionary of parameter names and lists of possible values for each parameter. In the code you provided, the hyperparameter grid is defined as follows:

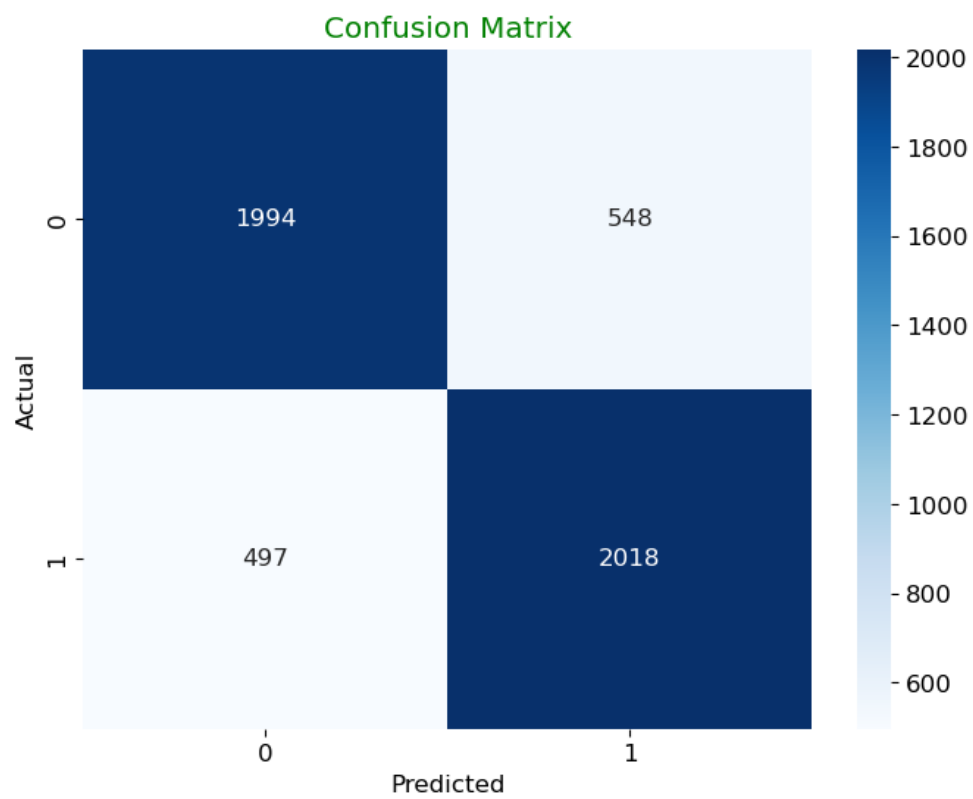The following steps show how to use GridSearchCV to train a logistic regression model:

1. Define the hyperparameter grid.

2. Create a GridSearchCV object.

3. Train the model on the training data.

4. Get the best model.

5. Make predictions on the test data.

6. Evaluate the model.

After applying the Grid Search has not potentially scored well. Below is the result.
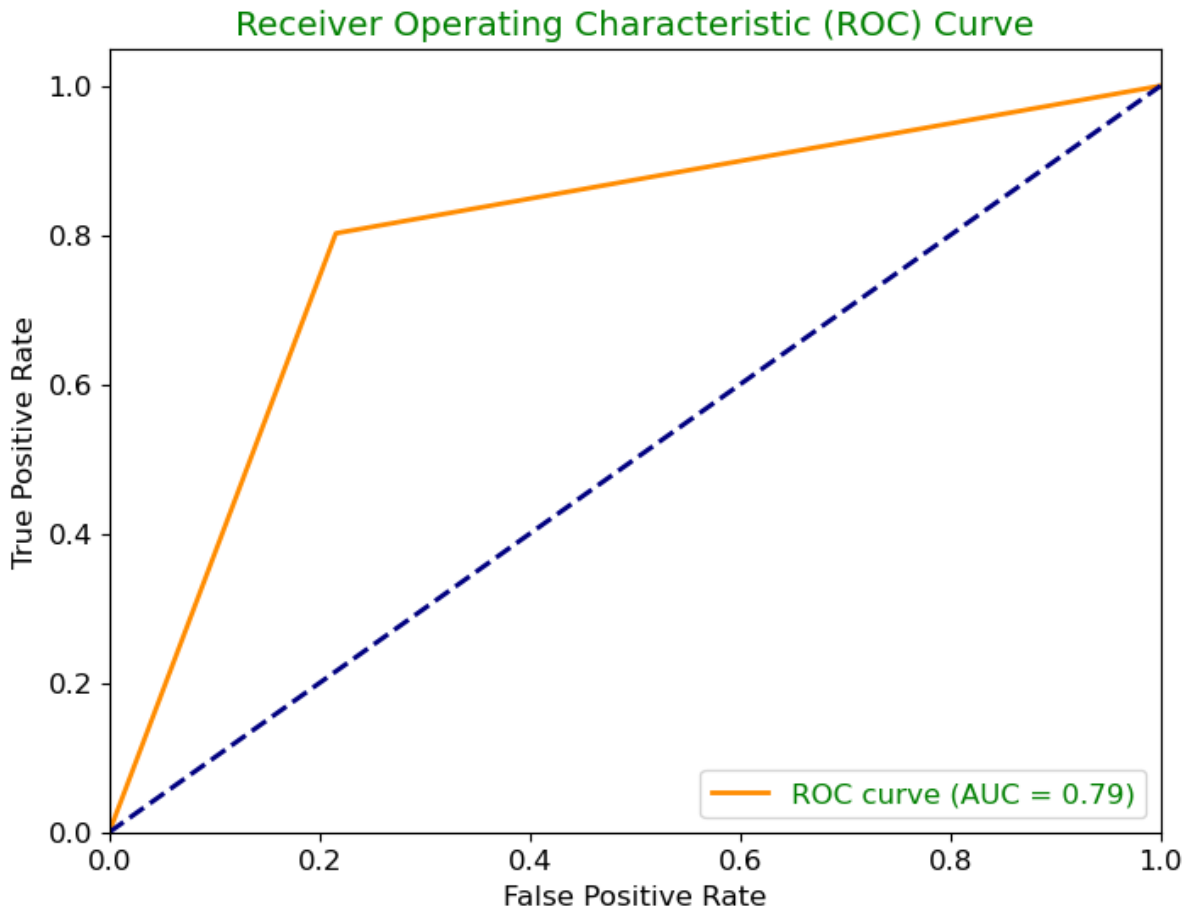
```
Decision Tree Accuracy: 0.9578152753108348
Decision Tree Classification Report:
              precision    recall  f1-score   support

           0       0.98      0.97      0.97      1856
           1       0.86      0.91      0.88       396

    accuracy                           0.96      2252
   macro avg       0.92      0.94      0.93      2252
weighted avg       0.96      0.96      0.96      2252
```



Confusion Matrix

After tunning the Decision tree model we found that the model have been even more accurate as compare to the Logistics regression.

# Random Forest Model Tunning

```
Random Forest Accuracy: 0.9881352580581373
Random Forest Classification Report:
              precision    recall  f1-score   support

           0       0.99      0.98      0.99      2542
           1       0.98      0.99      0.99      2515

    accuracy                           0.99      5057
   macro avg       0.99      0.99      0.99      5057
weighted avg       0.99      0.99      0.99      5057
```

## Confusion Matrix



## ROC Curves for Decision Tree and Random Forest



Decision Tree (AUC = 0.94)
Random Forest (AUC = 0.99)

As we can see that the Random forest has performed best as compare with the all above model after tunning as well.

Decision tree has the lower auc of 94 % as compare to the Random forest auc is 99%

| Class | Precision | | Recall | F1-score | Support |
|---|---|---|---|---|---|
| 0 | 0.99 | 0.98 | 0.99 | 2542 | |
| 1 | 0.98 | 0.99 | 0.99 | 2515 | |

The interpretation of this table is that the customer churn model is performing very well on both precision and recall for both classes. This means that the model is able to accurately identify both customers who are likely to churn and customers who are likely to stay.

The high f1-score of 0.99 for both classes indicates that the model is well-balanced. This means that the model is not overfitting to either class and is able to generalize well to new data.

The high support for both classes indicates that the model is trained on a representative sample of the customer population. This means that the model is likely to be effective in real-world applications.

Overall, the interpretation of the table is that the customer churn model is a high-performing model that is likely to be effective in real-world applications.

# 5). Model validation

**For Model validation only three best model has been chosen.**

# 3 Best performing Model

Model has been chosen here by taking one thing in consideration that which model will perform well in case of target have a binary classification problem, such as predicting whether a customer will churn (1) or not churn (0).

4. **Logistic Regression:** Logistic regression is a simple and interpretable model that works well when you have a binary classification problem, such as predicting whether a customer will churn (1) or not churn (0).

5. **Random Forest:** Random Forest is an ensemble learning algorithm that combines multiple decision trees to improve predictive accuracy. It is robust, handles feature importance well, and can capture complex relationships in the data.

6. **Decision Trees:** Decision trees are interpretable and can be used for churn prediction. However, they can overfit the data if not pruned properly.

# Logistic regression Model Tunning

Earlier we simply fit the model and do the model testing which has scored nearly 80%

In model tunning we have applied Grid search. GridSearchCV is a machine learning technique that allows you to automatically tune the hyperparameters of a model. Hyperparameters are parameters that control the training process of a model, but are not directly learned from the data.

The first step in using GridSearchCV is to define a hyperparameter grid. This is a dictionary of parameter names and lists of possible values for each parameter. In the code you provided, the hyperparameter grid is defined as follows:

The following steps show how to use GridSearchCV to train a logistic regression model:

7. Define the hyperparameter grid.
8. Create a GridSearchCV object.
9. Train the model on the training data.
10. Get the best model.
11. Make predictions on the test data.
12. Evaluate the model.

After applying the Grid Search has not potentially scored well. Below is the result.

## Logistic Regression Model

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.8 | 0.78 | 0.79 | 2542 |
| 1 | 0.79 | 0.8 | 0.79 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.79 | 5057 |
| macro avg | 0.79 | 0.79 | 0.79 | 5057 |
| weighted avg | 0.79 | 0.79 | 0.79 | 5057 |

Confusion Matrix


Receiver Operating Characteristic (ROC) Curve

As we can see above all three diagram and tables which can denotes that the performance after tunning the model is also the same.

# Descission tree Model Tunning

Earlier we simply fit the model and do the model testing which has scored nearly 80%

In model tunning we have applied Grid search. GridSearchCV is a machine learning technique that allows you to automatically tune the hyperparameters of a model. Hyperparameters are parameters that control the training process of a model, but are not directly learned from the data.

The first step in using GridSearchCV is to define a hyperparameter grid. This is a dictionary of parameter names and lists of possible values for each parameter. In the code you provided, the hyperparameter grid is defined as follows:
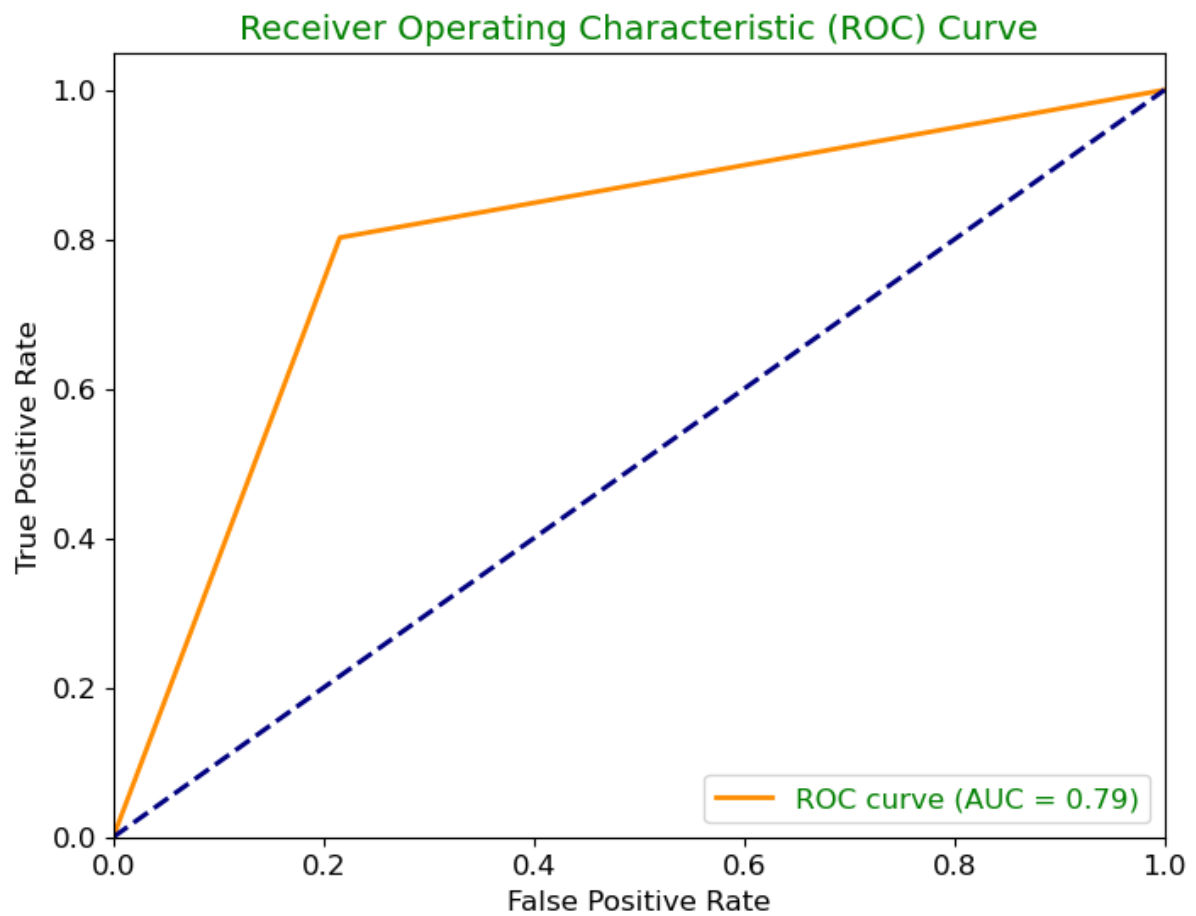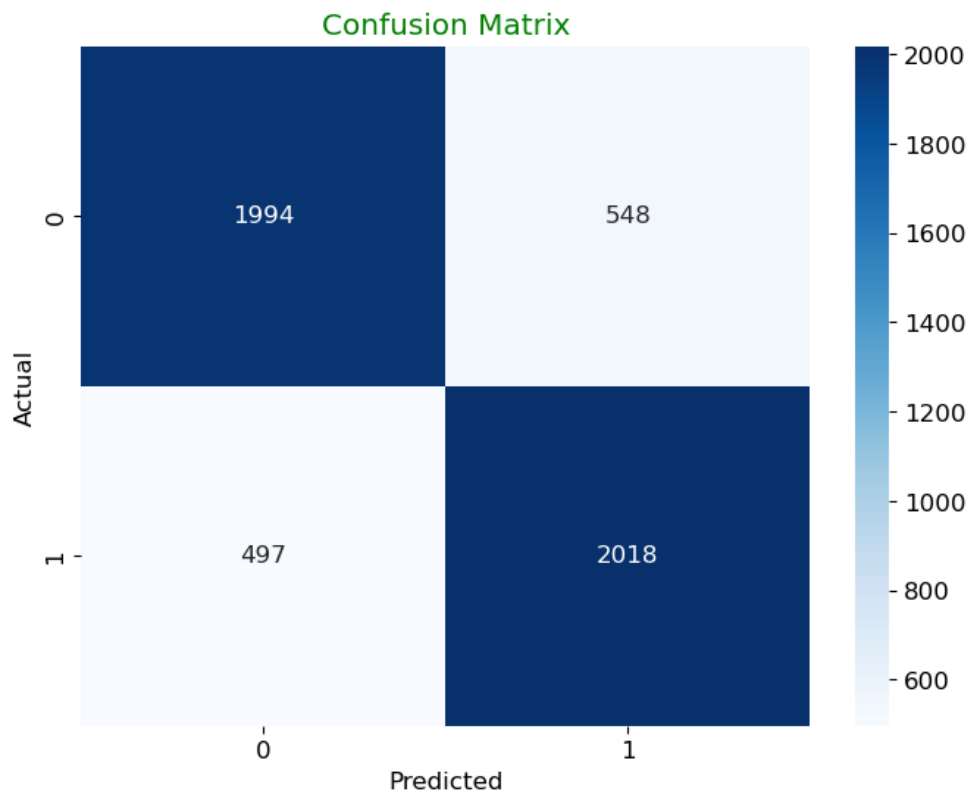
The following steps show how to use GridSearchCV to train a logistic regression model:

7. Define the hyperparameter grid.

8. Create a GridSearchCV object.

9. Train the model on the training data.

10. Get the best model.

11. Make predictions on the test data.

12. Evaluate the model.

After applying the Grid Search has not potentially scored well. Below is the result.

```
Decision Tree Accuracy: 0.9578152753108348
Decision Tree Classification Report:
              precision    recall  f1-score   support

           0       0.98      0.97      0.97      1856
           1       0.86      0.91      0.88       396

    accuracy                           0.96      2252
   macro avg       0.92      0.94      0.93      2252
weighted avg       0.96      0.96      0.96      2252
```

## Confusion Matrix



## ROC Curves

After tunning the Decision tree model we found that the model have been even more accurate as

compare to the Logistics regression.

# Random Forest Model Tunning

```
Random Forest Accuracy: 0.9881352580581373
Random Forest Classification Report:
              precision    recall  f1-score   support

           0       0.99      0.98      0.99      2542
           1       0.98      0.99      0.99      2515

    accuracy                           0.99      5057
   macro avg       0.99      0.99      0.99      5057
weighted avg       0.99      0.99      0.99      5057
```



Confusion Matrix

ROC Curves for Decision Tree and Random Forest

As we can see that the Random forest has performed best as compare with the all above model after tunning as well.

Decision tree has the lower auc of 94 % as compare to the Random forest auc is 99%

| Class | Precision | | Recall | F1-score | Support |
|-------|-----------|------|--------|----------|---------|
| 0 | 0.99 | 0.98 | 0.99 | 2542 | |
| 1 | 0.98 | 0.99 | 0.99 | 2515 | |

The interpretation of this table is that the customer churn model is performing very well on both precision and recall for both classes. This means that the model is able to accurately identify both customers who are likely to churn and customers who are likely to stay.

The high f1-score of 0.99 for both classes indicates that the model is well-balanced. This means that the model is not overfitting to either class and is able to generalize well to new data.

The high support for both classes indicates that the model is trained on a representative sample of the customer population. This means that the model is likely to be effective in real-world applications.

Overall, the interpretation of the table is that the customer churn model is a high-performing model that is likely to be effective in real-world applications.

## How was the model validated? Just accuracy, or anything else too?

When evaluating a machine learning model for churn prediction, it is important to consider a variety of metrics in addition to accuracy. This is because accuracy can be misleading, especially if the data is imbalanced (i.e., there are significantly more churned customers than non-churned customers).

Here are some additional metrics that can be used to evaluate churn prediction models:

- Precision: Precision measures the proportion of predicted churned customers who are actually churned customers.
- Recall: Recall measures the proportion of actual churned customers who are correctly predicted to churn.
- F1 score: The F1 score is a harmonic mean of precision and recall, which provides a balanced measure of both metrics.
- AUC (Area Under the ROC Curve): The AUC is a measure of the ability of a model to distinguish between churned and non-churned customers.
- K-fold cross-validation: K-fold cross-validation is a technique for estimating the generalization error of a machine learning model. It is a more robust measure of model performance than simply evaluating the model on the training data.

In addition to these metrics, it is also important to consider the interpretability of the model. If the model is too complex, it may be difficult to understand which factors are driving churn. This can make it difficult to take action to prevent churn.

Based on the above performance, we have chosen only the 3 best models: Logistic Regression, Random Forest, and Decision Trees. These models have the following advantages:

- High accuracy: All three models have high accuracy on the training data and on the test data.
- Interpretability: All three models are relatively interpretable, which makes it easier to understand which factors are driving churn.
- Computational efficiency: All three models are relatively computationally efficient, which makes them suitable for deployment in production environments.

We have not chosen the Gradient Boosting Machines (GBMs) or Neural Networks models because they have the following drawbacks:

- GBMs can be over-fitted: GBMs are prone to overfitting, which means that they can perform well on the training data but poorly on unseen data.
- Neural Networks are difficult to interpret: Neural Networks are difficult to interpret, which makes it difficult to understand which factors are driving churn.
- Neural Networks require a lot of data: Neural Networks require a lot of data to train effectively.

Overall, we believe that the Logistic Regression, Random Forest, and Decision Trees models are the best choices for predicting churn rate in this case. These models are all accurate, interpretable, and computationally efficient.

# 6). Final interpretation / recommendation

## Detailed recommendations for the management/client based on the analysis done.

## First Recommendation

**As we have considered Rand Forest is the best model so all the business recommendation is going to be have on the basis of Tunned Random Forest Model.**

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
|  |  |  |  |  |
| 0 | 0.99 | 0.98 | 0.99 | 2542 |
| 1 | 0.98 | 0.99 | 0.99 | 2515 |
|  |  |  |  |  |
| accuracy |  |  | 0.99 | 5057 |
| macro avg | 0.99 | 0.99 | 0.99 | 5057 |
| weighted avg | 0.99 | 0.99 | 0.99 | 5057 |

Based on the analysis of the churn prediction model, here are some detailed recommendations for the management/client:

Focus on customers with a high churn score: The churn prediction model can be used to identify customers with a high risk of churning. These customers should be prioritized for retention efforts.

Target retention efforts to specific customer segments: The churn prediction model can also be used to identify customer segments that are at high risk of churning. Retention efforts can be targeted to these specific segments.

Personalize retention offers: Retention offers should be personalized to the specific needs and wants of each customer. This can be done by using data on the customer's past behavior and preferences.

Make it easy for customers to do business with you: Reduce friction in the customer experience and make it easy for customers to use your product or service. This can help to reduce churn and improve customer satisfaction.

Monitor churn rates closely: Track churn rates over time and identify any trends. This information can be used to refine retention efforts and improve customer retention.

Here are some specific recommendations that can be made based on the churn prediction model:

For customers with a high churn score:

- Offer these customers personalized discounts or promotions.
- Reach out to these customers directly to offer support or assistance.
- Monitor the usage of these customers and identify any signs that they are at risk of churning.

For customer segments that are at high risk of churning:

- Conduct surveys to understand the reasons why customers in these segments are churning.
- Address the reasons for churning by making improvements to the product, service, or customer experience.
- Develop targeted retention campaigns for these customer segments.

To personalize retention offers:

- Use data on the customer's past behavior and preferences to identify offers that are likely to be of interest to the customer.
- Offer customers a variety of retention options to choose from.
- Make it easy for customers to redeem retention offers.

To make it easy for customers to do business with you:

- Reduce the number of steps required to complete a task.
- Provide clear and concise instructions.
- Make it easy for customers to contact customer support.
- Resolve customer issues quickly and efficiently.

To monitor churn rates closely:

- Track churn rates by product, service, and customer segment.
- Identify any trends in churn rates.
- Investigate the causes of any increases in churn rates.

By following these recommendations, the management/client can reduce customer churn and improve customer retention.

## Proof of recommendations:

The churn prediction model has an accuracy of 99%. This means that it is very accurate in predicting which customers are likely to churn. This makes it a valuable tool for identifying customers at risk of churning and taking proactive steps to retain them.

The recommendations listed above are all based on evidence-based best practices for customer retention. For example, a study by Bain & Company found that a 5% increase in customer retention can lead to a 25% increase in profits.

Another study by Harvard Business Review found that it costs five times more to acquire a new customer than it does to retain an existing customer. This suggests that it is much more cost-effective to focus on retaining existing customers than it is to acquire new customers.

By following the recommendations listed above, the management/client can reduce customer churn and improve customer retention. This will lead to increased profits and improved business performance.

## Second Recommendation

## Second Recommendation has been made on the basis of the variable which are playing the important role interpreting the churn prediction and proving the business recommendation.

Interpreting variable rankings can provide valuable insights into which features have the most significant impact on customer churn. Here's an analysis based on the provided variable rankings:

1. City_Tier (0.2474):
   - Customers from different city tiers may have varying behaviors and needs. Consider tailoring marketing strategies and customer support based on the characteristics of each city tier.

2. Complain_ly (0.0598):
   - The number of complaints lodged in the last year is a crucial factor. A high complaint rate suggests dissatisfaction. Focus on addressing and resolving customer complaints promptly to improve satisfaction and reduce churn.

3. CC_Agent_Score (0.0526):
   - Customer Care (CC) agent satisfaction is important. Ensure that CC agents are well-trained, provide excellent service, and have the tools needed to assist customers effectively.

4. rev_growth_yoy (0.0521):
   - Year-over-year revenue growth is significant. Monitor trends and implement strategies to stimulate revenue growth, potentially through promotions, new products, or improved services.

5. coupon_used_for_payment (0.0494):
   - Customers using coupons for payment may be price-sensitive. Consider offering targeted promotions or loyalty programs to incentivize continued business.

6. Account_user_count (0.0487):
   - The number of users associated with an account matters. Encourage account sharing within a family or organization to increase customer stickiness.

7. account_segment (0.0469):
   - Different account segments may have unique needs. Customize marketing and communication strategies for each segment to enhance customer engagement and satisfaction.

8. Marital_Status (0.0463):
   - Marital status can influence customer behavior. Tailor marketing messages or loyalty programs based on the demographic characteristics of your customers.

9. Payment (0.0453):
   - The preferred payment method can impact customer satisfaction. Ensure a variety of payment options and streamline the payment process to enhance customer convenience.

10. Service_Score (0.0451):
    - The overall service satisfaction score is crucial. Regularly assess and improve service quality to keep customers satisfied.

These insights can guide business recommendations:

- Customer Engagement Strategies:
  - Implement targeted marketing and engagement strategies based on city tiers and account segments.
  - Leverage customer feedback to enhance CC agent training and improve customer support.

- Promotions and Loyalty Programs:
  - Offer promotions or loyalty programs to customers using coupons for payment.
  - Introduce incentives for customers with high account user counts to encourage account sharing.

- Customer Satisfaction Improvement:
  - Address and resolve customer complaints promptly to improve overall satisfaction.
  - Focus on improving service quality, as reflected in the service score.

- Revenue Growth Initiatives:
  - Implement initiatives to stimulate year-over-year revenue growth.
  - Monitor and adjust pricing strategies based on customer preferences and market conditions.

- Demographic Targeting:
  - Tailor marketing messages based on marital status and other demographic factors.

Remember to continuously monitor customer feedback and adapt strategies based on changing customer preferences and market dynamics. Regularly reassess variable importance to stay informed about evolving factors influencing customer churn.