

✓ UNVEILING CUSTOMER CHURN

Problem statement:

An E Commerce company or DTH (you can choose either of these two domains) provider is facing a lot of competition in the current market and it has become a challenge to retain the existing customers in the current situation. Hence, the company wants to develop a model through which they can do churn prediction of the accounts and provide segmented offers to the potential churners. In this company, account churn is a major thing because 1 account can have multiple customers. hence by losing one account the company might be losing more than one customer.

You have been assigned to develop a churn prediction model for this company and provide business recommendations on the campaign. Your campaign suggestion should be unique and be very clear on the campaign offer because your recommendation will go through the revenue assurance team. If they find that you are giving a lot of free (or subsidized) stuff thereby making a loss to the company; they are not going to approve your recommendation.

Hence be very careful while providing campaign recommendation.

Variable	Description
AccountID	account unique identifier
Churn	account churn flag (Target)
Tenure	Tenure of account
City_Tier	Tier of primary customer's city
CC_Contacted_LY	How many times all the customers of the account has contacted customer care in last 12 months
Payment	Preferred Payment mode of the customers in the account
Gender	Gender of the primary customer of the account
Service_Score	Satisfaction score given by customers of the account on service provided by company
Account_user_count	Number of customers tagged with this account

Account_user_count	Number of customers tagged with this account
account_segment	Account segmentation on the basis of spend
CC_Agent_Score	Satisfaction score given by customers of the account on customer care service provided by company
Marital_Status	Marital status of the primary customer of the account
rev_per_month	Monthly average revenue generated by account in last 12 months
Complain_ly	Any complaints has been raised by account in last 12 months
rev_growth_yoy	revenue growth percentage of the account (last 12 months vs last 24 to 13 months)
coupon_used_for_payment	How many times customers have used coupons to do the payment in last 12 months
Day_Since_CC_connect	Number of days since no customers in the account has contacted the customer care
cashback	Monthly average cashback generated by account in last 12 months
Login_device	Preferred login device of the customers in the account

✓ Installing and Importing Dependencies

```
1 !pip install ydata_profiling
```

→ Collecting ydata_profiling

 Downloading ydata_profiling-4.8.3-py2.py3-none-any.whl (359 kB)

 359.5/359.5 kB 2.7 MB/s eta 0

Requirement already satisfied: scipy<1.14,>=1.4.1 in /usr/local/lib/python3

Requirement already satisfied: pandas!=1.4.0,<3,>1.1 in /usr/local/lib/pyth

Requirement already satisfied: matplotlib<3.9,>=3.2 in /usr/local/lib/pytho

Requirement already satisfied: pydantic>=2 in /usr/local/lib/python3.10/dis

```
Requirement already satisfied: PyYAML<6.1,>=5.0.0 in /usr/local/lib/python3
Requirement already satisfied: jinja2<3.2,>=2.11.1 in /usr/local/lib/python
Collecting visions[type_image_path]<0.7.7,>=0.7.5 (from ydata_profiling)
  Downloading visions-0.7.6-py3-none-any.whl (104 kB)
                                             104.8/104.8 kB 7.6 MB/s eta 0
Requirement already satisfied: numpy<2,>=1.16.0 in /usr/local/lib/python3.1
Collecting htmlmin==0.1.12 (from ydata_profiling)
  Downloading htmlmin-0.1.12.tar.gz (19 kB)
  Preparing metadata (setup.py) ... done
Collecting phik<0.13,>=0.11.1 (from ydata_profiling)
  Downloading phik-0.12.4-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x
                                             686.1/686.1 kB 13.5 MB/s eta
Requirement already satisfied: requests<3,>=2.24.0 in /usr/local/lib/python
Requirement already satisfied: tqdm<5,>=4.48.2 in /usr/local/lib/python3.10
Requirement already satisfied: seaborn<0.14,>=0.10.1 in /usr/local/lib/pyth
Collecting multimethod<2,>=1.4 (from ydata_profiling)
  Downloading multimethod-1.11.2-py3-none-any.whl (10 kB)
Requirement already satisfied: statsmodels<1,>=0.13.2 in /usr/local/lib/pyt
Collecting typeguard<5,>=3 (from ydata_profiling)
  Downloading typeguard-4.2.1-py3-none-any.whl (34 kB)
Collecting imagehash==4.3.1 (from ydata_profiling)
  Downloading ImageHash-4.3.1-py2.py3-none-any.whl (296 kB)
                                             296.5/296.5 kB 6.8 MB/s eta 0
Requirement already satisfied: wordcloud>=1.9.1 in /usr/local/lib/python3.1
Collecting dacite>=1.8 (from ydata_profiling)
  Downloading dacite-1.8.1-py3-none-any.whl (14 kB)
Requirement already satisfied: numba<1,>=0.56.0 in /usr/local/lib/python3.1
Requirement already satisfied: PyWavelets in /usr/local/lib/python3.10/dist-
Requirement already satisfied: pillow in /usr/local/lib/python3.10/dist-pac
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.1
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/di
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.1
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/pytho
Requirement already satisfied: llvmlite<0.42,>=0.41.0dev0 in /usr/local/lib
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/di
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/
Requirement already satisfied: joblib>=0.14.1 in /usr/local/lib/python3.10/
Requirement already satisfied: annotated-types>=0.4.0 in /usr/local/lib/pyt
Requirement already satisfied: pydantic-core==2.18.2 in /usr/local/lib/pyth
Requirement already satisfied: typing-extensions>=4.6.1 in /usr/local/lib/p
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/p
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/di
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3
Requirement already satisfied: patsy>=0.5.6 in /usr/local/lib/python3.10/di
Requirement already satisfied: attrs>=19.3.0 in /usr/local/lib/python3.10/d
Requirement already satisfied: networkx>=2.4 in /usr/local/lib/python3.10/d
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packag
```

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 from ydata_profiling import ProfileReport
```

✓ Loading the data into notebook

```
1 df = pd.read_excel('/content/drive/MyDrive/Colab Notebooks/Jain/da
```

✓ Few observations on the dataset

```
1 df.head()
```

→

	AccountID	Churn	Tenure	City_Tier	CC_Contacted_LY	Payment	Gender	Ser
0	20000	1	4	3.0		6.0	Debit Card	Female
1	20001	1	0	1.0		8.0	UPI	Male
2	20002	1	0	1.0		30.0	Debit Card	Male
3	20003	1	0	3.0		15.0	Debit Card	Male
4	20004	1	0	1.0		12.0	Credit Card	Male

```
1 df.columns
```

→

```
Index(['AccountID', 'Churn', 'Tenure', 'City_Tier', 'CC_Contacted_LY',
       'Payment', 'Gender', 'Service_Score', 'Account_user_count',
       'account_segment', 'CC_Agent_Score', 'Marital_Status',
       'rev_per_month',
       'Complain_ly', 'rev_growth_yoy', 'coupon_used_for_payment',
       'Day_Since_CC_connect', 'cashback', 'Login_device'],
      dtype='object')
```

```
1 df.ndim
```

```
↳ 2
```

```
1 df.size
```

```
↳ 213940
```

```
1 len(df)
```

```
↳ 11260
```

```
1 df.shape
```

```
↳ (11260, 19)
```

✓ Checking the data type of each column

```
1 df.info()
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 11260 entries, 0 to 11259
Data columns (total 19 columns):
 #   Column           Non-Null 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  
dtypes: float64(5), int64(2), object(12)
memory usage: 1.6+ MB
```

Insights:

There are 7 numerical columns and 12 categorical columns.

✓ Statistical analysis of numeric columns on the dataset

```
1 df.describe().T
```

	count	mean	std	min	25%	50%
AccountID	11260.0	25629.500000	3250.626350	20000.0	22814.75	25629.5
Churn	11260.0	0.168384	0.374223	0.0	0.00	0.0
City_Tier	11148.0	1.653929	0.915015	1.0	1.00	1.0
CC_Contacted_LY	11158.0	17.867091	8.853269	4.0	11.00	16.0
Service_Score	11162.0	2.902526	0.725584	0.0	2.00	3.0
CC_Agent_Score	11144.0	3.066493	1.379772	1.0	2.00	3.0
Complain_ly	10903.0	0.285334	0.451594	0.0	0.00	0.0

✓ Statistical analysis of categorical columns on the dataset

```
1 df.describe(include='0').T
```

	count	unique	top	freq
Tenure	11158	38	1	1351
Payment	11151	5	Debit Card	4587
Gender	11152	4	Male	6328
Account_user_count	11148	7	4	4569
account_segment	11163	7	Super	4062
Marital_Status	11048	3	Married	5860
rev_per_month	11158	59	3	1746
rev_growth_yoy	11260	20	14	1524
coupon_used_for_payment	11260	20	1	4373
Day_Since_CC_connect	10903	24	3	1816
cashback	10789.0	5693.0	155.62	10.0
Login_device	11039	3	Mobile	7482

▼ Data Cleaning

```
1 for column in df.select_dtypes('object'):
2     print(f'{column}: {df[column].nunique()}')
3     print(f'{df[column].unique()}')
4     print('\n\n')
5
6 [4 0 2 13 11 '#' 9 99 19 20 14 8 26 18 5 30 7 1 23 3 29 6 28 24 25 16 10
7 15 22 nan 27 12 21 17 50 60 31 51 61]
```

Payment: 5
['Debit Card' 'UPI' 'Credit Card' 'Cash on Delivery' 'E wallet' nan]

Gender: 4
['Female' 'Male' 'F' nan 'M']

Account_user_count: 7
[3 4 nan 5 2 '@' 1 6]

account_segment: 7
['Super' 'Regular Plus' 'Regular' 'HNI' 'Regular +' nan 'Super Plus'
'Super +']

Marital_Status: 3
['Single' 'Divorced' 'Married' nan]

rev_per_month: 59
[9 7 6 8 3 2 4 10 1 5 '+' 130 nan 19 139 102 120 138 127 123 124 116 21
126 134 113 114 108 140 133 129 107 118 11 105 20 119 121 137 110 22 101
136 125 14 13 12 115 23 122 117 131 104 15 25 135 111 109 100 103]

rev_growth_yoy: 20
[11 15 14 23 22 16 12 13 17 18 24 19 20 21 25 26 '\$' 4 27 28]

```
coupon_used_for_payment: 20  
[1 0 4 2 9 6 11 7 12 10 5 3 13 15 8 '#' '$' 14 '*' 16]
```

```
Day_Since_CC_connect: 24  
[5 0 3 7 2 1 8 6 4 15 nan 11 10 9 13 12 17 16 14 30 '$' 46 18 31 47]
```

```
cashback: 5693  
[159.93 120.9 nan ... 227.36 226.91 191.42]
```

```
Login_device: 3  
['Mobile' 'Computer' 'Tablet' 'nan']
```

✓ Correcting all the categorical variables with special symbols

✓ Treating 'Tenure' column

- As you can see there is '#' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '#'s are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['Tenure'] = df['Tenure'].replace('#', np.nan)
```

```
2
```

```
3 df['Tenure'].unique()
```

```
→ array([ 4.,  0.,  2., 13., 11., nan,  9., 99., 19., 20., 14.,  8., 26.,  
        18.,  5., 30.,  7.,  1., 23.,  3., 29.,  6., 28., 24., 25., 16.,  
        10., 15., 22., 27., 12., 21., 17., 50., 60., 31., 51., 61.])
```

✓ Treating 'Gender' column

- As we can see 'F' and 'M' are also known as 'Female' and 'Male'.

```
1 df['Gender'] = df['Gender'].replace('F', 'Female')
2 df['Gender'] = df['Gender'].replace('M', 'Male')
3
4 df['Gender'].unique()

→ array(['Female', 'Male', nan], dtype=object)
```

✓ Treating 'Account_user_count' column

- As you can see there is '@' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '@s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['Account_user_count'] = df['Account_user_count'].replace('@', np.nan)
2
3 df['Account_user_count'].unique()

→ array([ 3.,  4., nan,  5.,  2.,  1.,  6.])
```

✓ Treating 'account_segment' column

- As you can see there is 'Regular +' and 'Super +' symbol need to be treated as 'Regular Plus' and 'Super Plus'.
- Since 'Regular +' is nothing but 'Regular Plus' and same in case of 'Super +'.

```
1 df['account_segment'] = df['account_segment'].replace('Regular +', 'Regular')
2 df['account_segment'] = df['account_segment'].replace('Super +', 'Super')
3
4 df['account_segment'].unique()

→ array(['Super', 'Regular Plus', 'Regular', 'HNI', nan, 'Super Plus'],
       dtype=object)
```

✓ Treating 'rev_per_month' column

- As you can see there is '+' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '+'s are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['rev_per_month'] = df['rev_per_month'].replace('+', np.nan)
2
3 df['rev_per_month'].unique()
```

```
→ array([ 9.,  7.,  6.,  8.,  3.,  2.,  4., 10.,  1.,  5.,  nan,
       130., 19., 139., 102., 120., 138., 127., 123., 124., 116., 21.,
       126., 134., 113., 114., 108., 140., 133., 129., 107., 118., 11.,
       105., 20., 119., 121., 137., 110., 22., 101., 136., 125., 14.,
       13., 12., 115., 23., 122., 117., 131., 104., 15., 25., 135.,
       111., 109., 100., 103.])
```

✓ Treating 'rev_growth_yoy' column

- As you can see there is '\$' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '\$s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['rev_growth_yoy'] = df['rev_growth_yoy'].replace('$', np.nan)
2
3 df['rev_growth_yoy'].unique()
```

```
→ array([11., 15., 14., 23., 22., 16., 12., 13., 17., 18., 24., 19., 20.,
       21., 25., 26., nan, 4., 27., 28.])
```

✓ Treating 'coupon_used_for_payment' column

- As you can see there is '#', '\$' and '*' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '#s', '\$s' and '*s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['coupon_used_for_payment'] = df['coupon_used_for_payment'].replace('$', np.nan)
2 df['coupon_used_for_payment'] = df['coupon_used_for_payment'].replace(' ', np.nan)
3 df['coupon_used_for_payment'] = df['coupon_used_for_payment'].replace(',', np.nan)
4
5 df['coupon_used_for_payment'].unique()

→ array([ 1.,  0.,  4.,  2.,  9.,  6., 11.,  7., 12., 10.,  5.,  3., 13.,
       15.,  8., nan, 14., 16.])
```

✓ Treating 'Day_Since_CC_connect' column

- As you can see there is '\$' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '\$s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['Day_Since_CC_connect'] = df['Day_Since_CC_connect'].replace('$', np.nan)
2
3 df['Day_Since_CC_connect'].unique()

→ array([ 5.,  0.,  3.,  7.,  2.,  1.,  8.,  6.,  4., 15., nan, 11., 10.,
       9., 13., 12., 17., 16., 14., 30., 46., 18., 31., 47.])
```

✓ Treating 'cashback' column

- First need to know what is the special character.
- As you can see there is '\$' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '\$s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 for data in df['cashback'].unique():
2     if not isinstance(data, (int, float)):
3         print(data)

→ $
```

```
1 df['cashback'] = df['cashback'].replace('$', np.nan)
2
3 df['cashback'].unique()
```

```
array([159.93, 120.9 ,      nan, ..., 227.36, 226.91, 191.42])
```

✓ Treating 'Login_device' column

- As you can see there is '&&&&' symbol need to be treated as 'np.nan' value and we can perform imputation in future.
- Once the '&&&&s' are converted to 'np.nan', this column will be converted to numerical column.

```
1 df['Login_device'] = df['Login_device'].replace('&&&&', np.nan)
2
3 df['Login_device'].unique()
```

```
array(['Mobile', 'Computer', nan], dtype=object)
```

✓ Exploratory Data Analysis (EDA)

✓ Pandas Profiling

```
1 profile = ProfileReport(df)
2 profile
```

Summarize dataset: 100%	129/129 [00:54<00:00, 1.57it/s, Completed]
Generate report structure: 100%	1/1 [00:14<00:00, 14.08s/it]
Render HTML: 100%	1/1 [00:03<00:00, 3.76s/it]

Overview

Dataset statistics

Number of variables	19
Number of observations	11260
Missing cells	4361
Missing cells (%)	2.0%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	1.6 MiB
Average record size in memory	152.0 B

Variable types

Numeric	10
Categorical	9

Alerts

Tenure has 218 (1.9%) missing values	Missing
Account_user_count has 444 (3.9%) missing values	Missing
CC_Agent_Score has 116 (1.0%) missing values	Missing
Marital_Status has 212 (1.9%) missing values	Missing

✓ Univariate Analysis

✓ 1st Figure

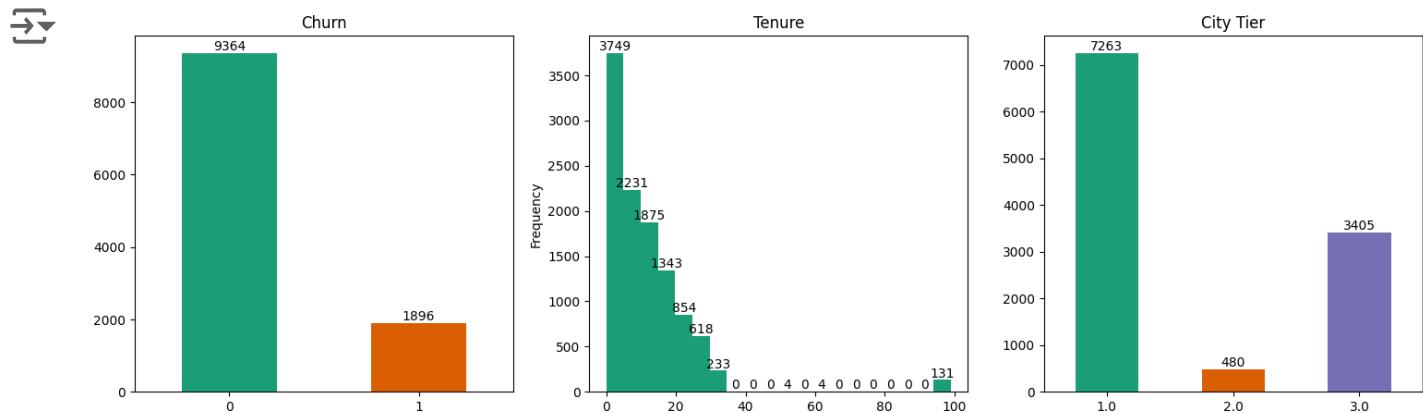
```
1 # @title 1st Figure
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 – Imbalance
5 # @markdown > Plot 2 – More num
6 # @markdown > Plot 3 – More num
7
8
9 fig = plt.figure(figsize=(18, 5
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('Churn').size()
14
15 ax.set_xlabel('')
16 ax.bar_label(ax.containers[0]);
17
18 # 2nd Column
19 plt.subplot(1, 3, 2)
20 ax = df['Tenure'].plot(kind='hi
21
22 ax.bar_label(ax.containers[0]);
23
24 # 3rd Column
25 plt.subplot(1, 3, 3)
26 ax = df.groupby('City_Tier').si
27
28 ax.set_xlabel('')
29 ax.bar_label(ax.containers[0]);
```

Data Insights:

Plot 1 - Imbalance dataset, as Churn samples are less.

Plot 2 - More number of customers are falling under less Tenure .

Plot 3 - More number of customers are from City Tier 1 followed by City Tier 3 and followed by City Tier 2.



▼ 2nd Figure

```

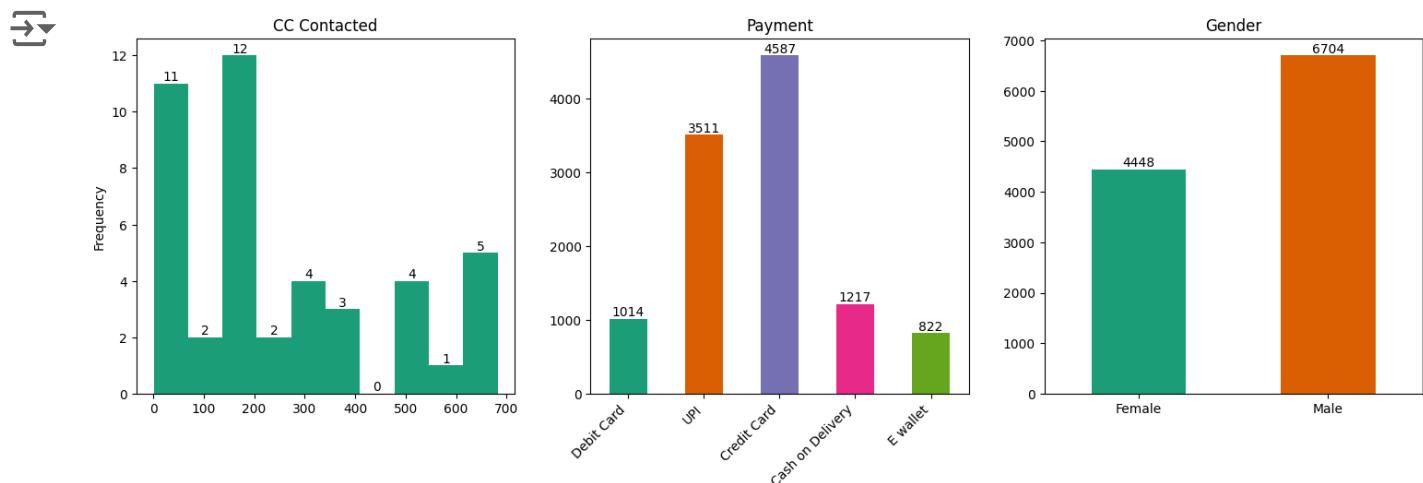
1 # @title 2nd Figure
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 - `CC_Contacted` doesn't follow any distribution.
5 # @markdown > Plot 2 - Most of the customers preferred Payment using credit card followed by UPI.
6 # @markdown > Plot 3 - Most of the customers are Male.
7
8
9 fig = plt.figure(figsize=(18, 5)
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('CC_Contacted_L
14
15
16 ax.set_xlabel('')
17 ax.bar_label(ax.containers[0]);
18
19

```

Data Insights:

Plot 1 - CC_Contacted doesn't follow any distribution.
 Plot 2 - Most of the customers preferred Payment using credit card followed by UPI.
 Plot 3 - Most of the customers are Male.

```
20 # 2nd Column
21 plt.subplot(1, 3, 2)
22 ax = df.groupby('Payment').size()
23
24 ax.set_xlabel('')
25 ax.set_xticklabels(df['Payment'])
26 ax.bar_label(ax.containers[0]);
27
28
29 # 3rd Column
30 plt.subplot(1, 3, 3)
31 ax = df.groupby('Gender').size()
32
33 ax.set_xlabel('')
34 ax.bar_label(ax.containers[0]);
```



▼ 3rd Figure

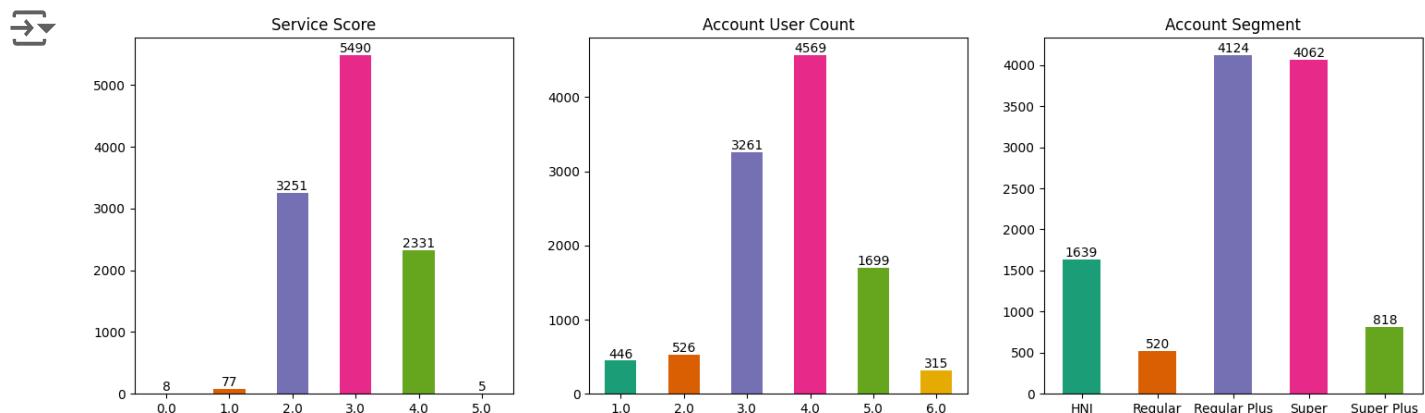
```
1 # @title 3rd Figure
```

Data Insights

```
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 – Most of 
5 # @markdown > Plot 2 – Most of 
6 # @markdown > Plot 3 – Most of 
7
8
9 fig = plt.figure(figsize=(18, 5
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('Service_Score'
14
15
16 ax.set_xlabel('')
17 ax.bar_label(ax.containers[0]);
18
19
20 # 2nd Column
21 plt.subplot(1, 3, 2)
22 ax = df.groupby('Account_user_c
23
24
25 ax.set_xlabel('')
26 ax.bar_label(ax.containers[0]);
27
28
29 # 3rd Column
30 plt.subplot(1, 3, 3)
31 ax = df.groupby('account_segmer
32
33 ax.set_xlabel('')
34 ax.bar_label(ax.containers[0]);
```

Data Insights.

Plot 1 - Most of the customers provided Service Score as 3.
Plot 2 - Most of the customers using 4 accounts.
Plot 3 - Most of the customers are under Regular Plus and Super Account Segment.



4th Figure

```

1 # @title 4th Figure
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 – Most of
5 # @markdown > Plot 2 – Most of
6 # @markdown > Plot 3 – Most of
7
8
9 fig = plt.figure(figsize=(18, 5
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('CC_Agent_Score'
14
15
16 ax.set_xlabel('')
17 ax.bar_label(ax.containers[0]);
18
19

```

Data Insights:

Plot 1 - Most of the customers provided score as 3 for the CC Agent .

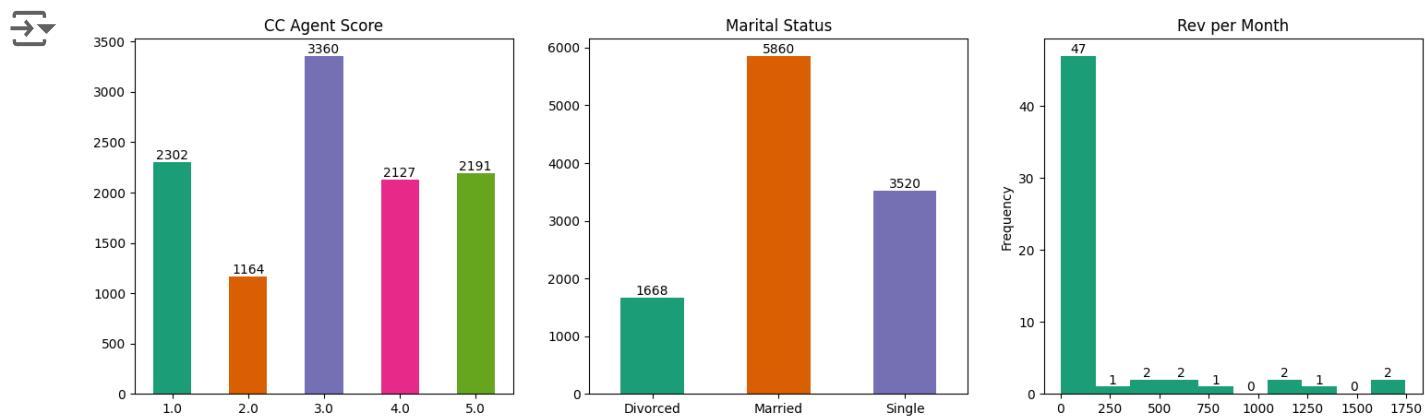
Plot 2 - Most of the customers are under married and followed by Single users.

Plot 3 - Most of the customers are under Regular Plus and Super Account Segment.

```

20 # 2nd Column
21 plt.subplot(1, 3, 2)
22 ax = df.groupby('Marital_Status'
23
24
25 ax.set_xlabel('')
26 ax.bar_label(ax.containers[0]);
27
28
29 # 3rd Column
30 plt.subplot(1, 3, 3)
31 ax = df.groupby('rev_per_month'
32
33 ax.set_xlabel('')
34 ax.bar_label(ax.containers[0]);

```



▼ 5th Figure

```

1 # @title 5th Figure
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 - `Complain_ly

```

Data Insights:

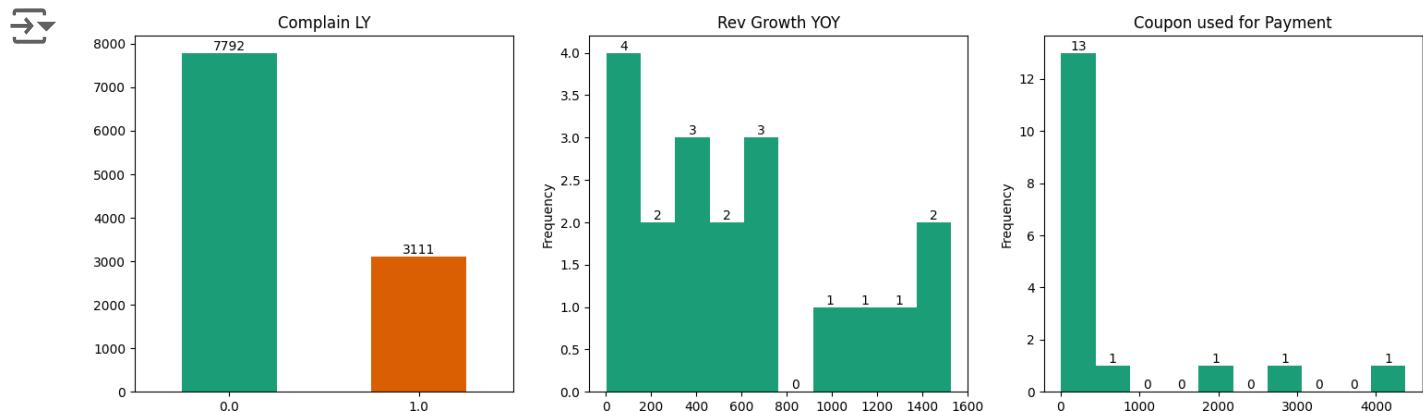
Plot 1 - Complain_ly

```
5 # @markdown > Plot 2 - `Rev_grc
6 # @markdown > Plot 3 - Most of
7
8
9 fig = plt.figure(figsize=(18, 5
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('Complain_ly').
14
15
16 ax.set_xlabel('')
17 ax.bar_label(ax.containers[0]);
18
19
20 # 2nd Column
21 plt.subplot(1, 3, 2)
22 ax = df.groupby('rev_growth_yoy'
23
24
25 ax.set_xlabel('')
26 ax.bar_label(ax.containers[0]);
27
28
29 # 3rd Column
30 plt.subplot(1, 3, 3)
31 ax = df.groupby('coupon_used_fc
32
33
34 ax.set_xlabel('')
35 ax.bar_label(ax.containers[0]);
```

has less complains.

Plot 2 - Rev_growth_yoy
does seems like normal
distribution.

Plot 3 - Most of the
customers used Coupon
for payment between 0
and 500 .



✓ 6th Figure

```

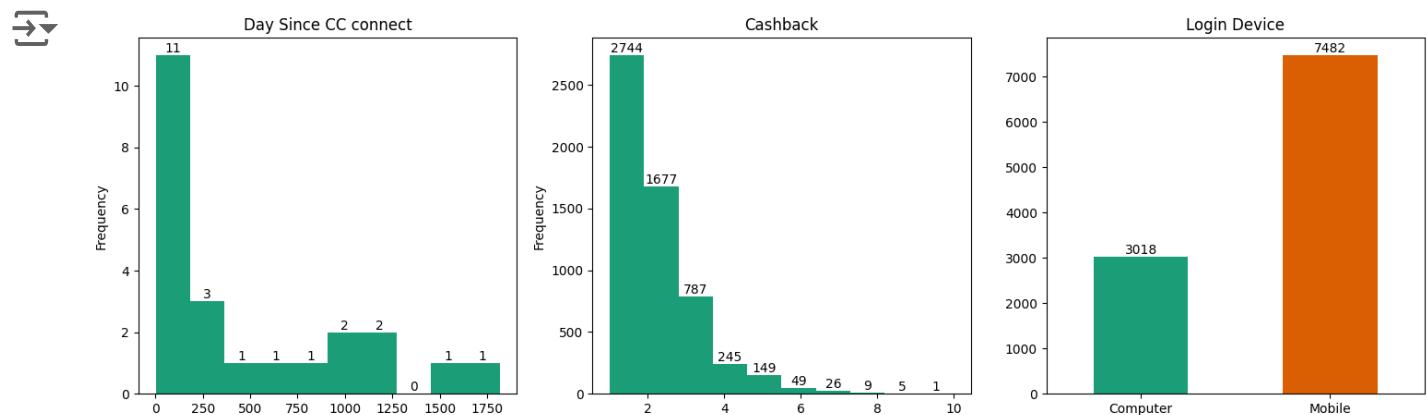
1 # @title 6th Figure
2
3 # @markdown **Data Insights:** 
4 # @markdown > Plot 1 – Most of
5 # @markdown > Plot 2 – Most of
6 # @markdown > Plot 3 – Most of
7
8
9 fig = plt.figure(figsize=(18, 5
10
11 # 1st Column
12 plt.subplot(1, 3, 1)
13 ax = df.groupby('Day_Since_CC_'
14
15
16 ax.set_xlabel('')
17 ax.bar_label(ax.containers[0]);
18
19

```

Data Insights:

Plot 1 - Most of the customers has not contacted CC connect .
 Plot 2 - Most of the customers got Cashback which is less in amount.
 Plot 3 - Most of the customers Login Device is Mobile compared to Computer .

```
20 # 2nd Column
21 plt.subplot(1, 3, 2)
22 ax = df.groupby('cashback').size()
23
24
25 ax.set_xlabel('')
26 ax.bar_label(ax.containers[0]);
27
28
29 # 3rd Column
30 plt.subplot(1, 3, 3)
31 ax = df.groupby('Login_device')
32
33
34 ax.set_xlabel('')
35 ax.bar_label(ax.containers[0]);
```



❖ Bivariate Analysis

✓ Warning

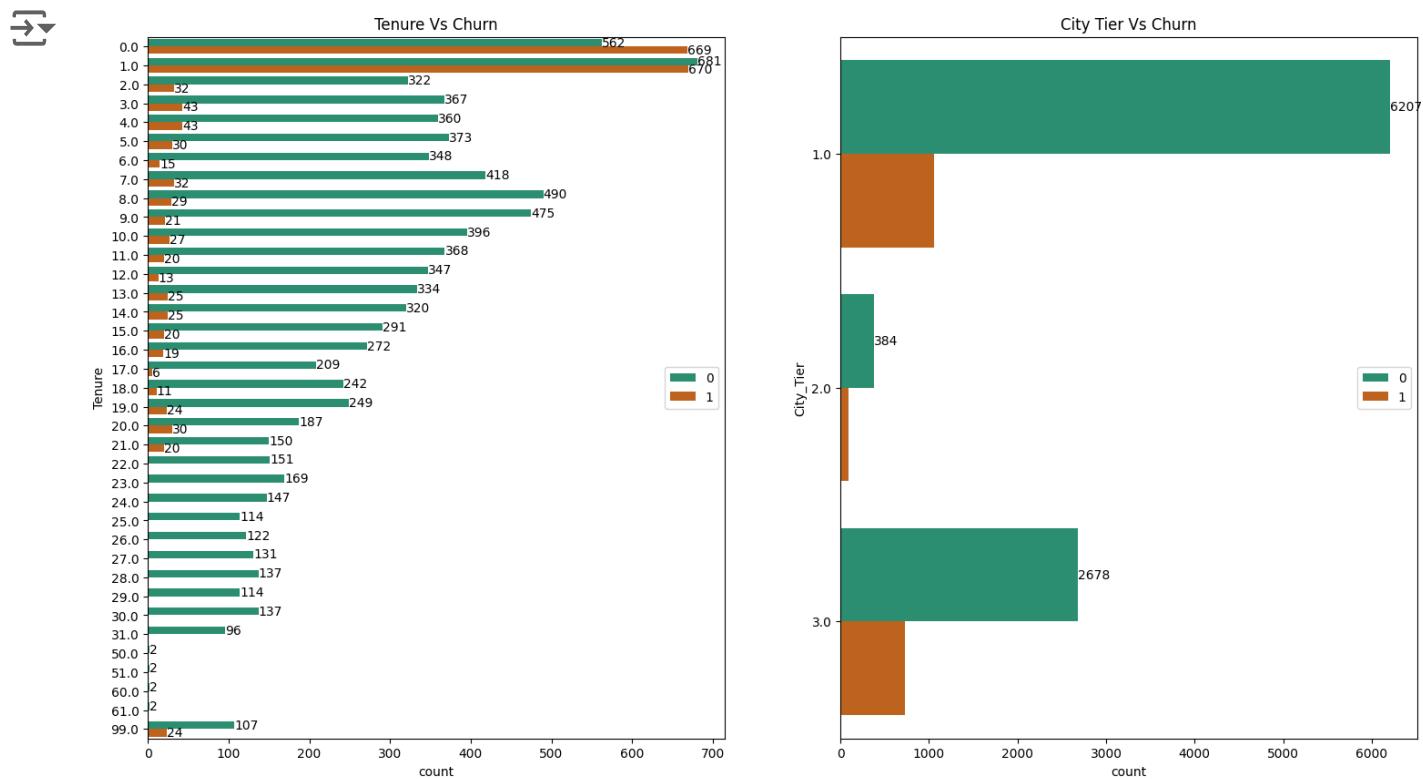
```
1 # @title Warning
2
3 import warnings
4
5 warnings.filterwarnings('ignore')
```

✓ Tenure Vs Churn and City Tier Vs Churn

```
1 # @title `Tenure` Vs `Churn` ar
2 # @markdown **Data Insights:** 
3 # @markdown > Plot 1 – If the `C
4 # @markdown > Plot 2 – `City_Ti
5
6 fig = plt.figure(figsize=(18, 1
7
8 plt.subplot(121)
9 ax = sns.countplot(df, y='Tenur
10 ax.legend(loc='center right')
11 ax.set_title('Tenure Vs Churn')
12 ax.bar_label(ax.containers[0])
13 ax.bar_label(ax.containers[1])
14
15 plt.subplot(122)
16 ax = sns.countplot(df, y='City_
17 ax.legend(loc='center right')
18 ax.set_title('City Tier Vs Chur
19 ax.bar_label(ax.containers[0]);
```

Data Insights:

Plot 1 - If the tenure is less, the Churn is more.
Plot 2 - City_Tier 2 and 3 has comparably more number of Churn rate.



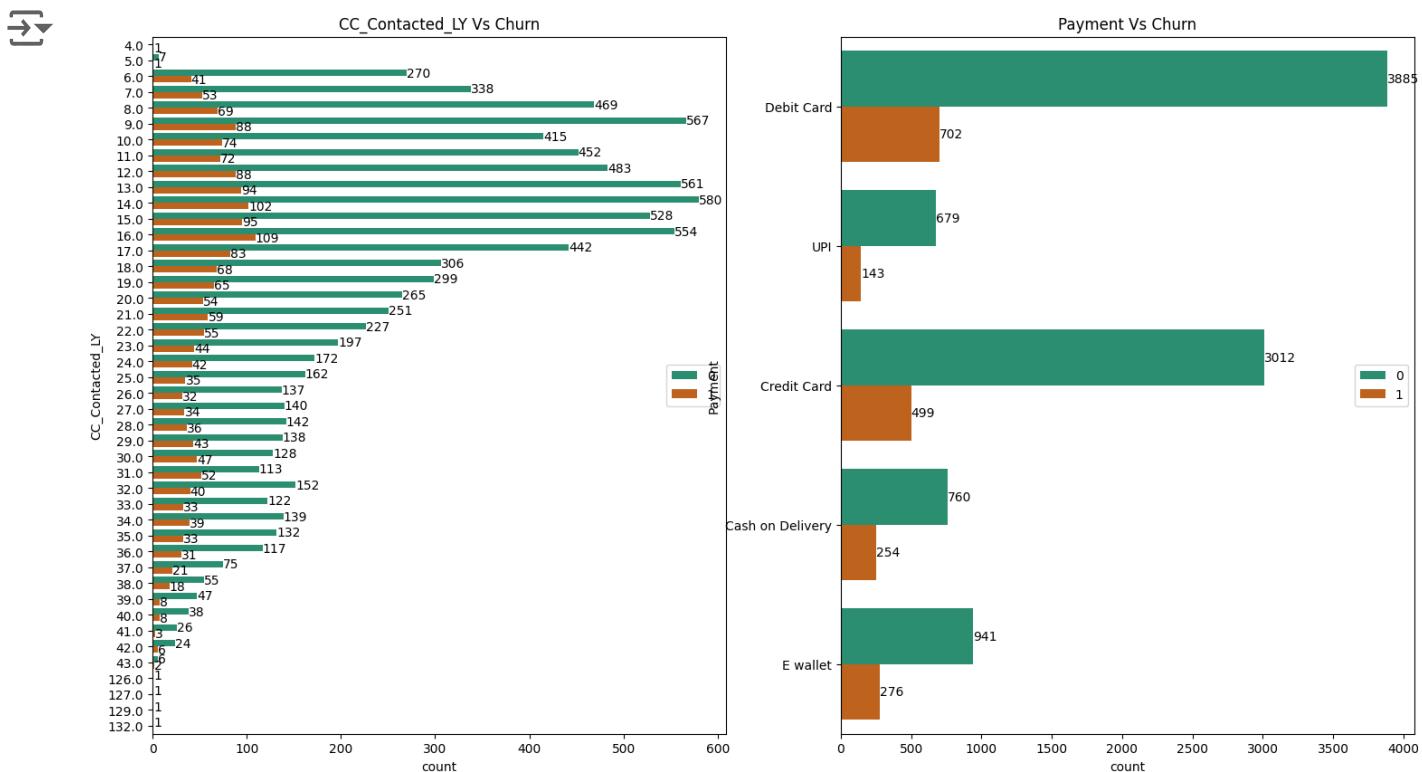
✓ CC_Contacted_LY Vs Churn and Payment Vs Churn

```
1 # @title `CC_Contacted_LY` Vs ` Data Insights:
2 # @markdown **Data Insights:**
```

```
3 # @markdown > Plot 1 - If the `<br>
4 # @markdown > Plot 2 - Except `<br>
5
6 fig = plt.figure(figsize=(18, 1<br>
7
8 plt.subplot(121)<br>
9 ax = sns.countplot(df, y='CC_Cc<br>
10 ax.legend(loc='center right')<br>
11 ax.set_title('CC_Contacted_LY <br>
12 ax.bar_label(ax.containers[0])<br>
13 ax.bar_label(ax.containers[1])<br>
14
15 plt.subplot(122)<br>
16 ax = sns.countplot(df, y='Payme<br>
17 ax.legend(loc='center right')<br>
18 ax.set_title('Payment Vs Churn'<br>
19 ax.bar_label(ax.containers[0])<br>
20 ax.bar_label(ax.containers[1]);
```

Plot 1 - If the CC_Contacted_LY is more, the Churn is also more.

Plot 2 - Except Debit and Credit card all other payment methods are having high Churn rate.

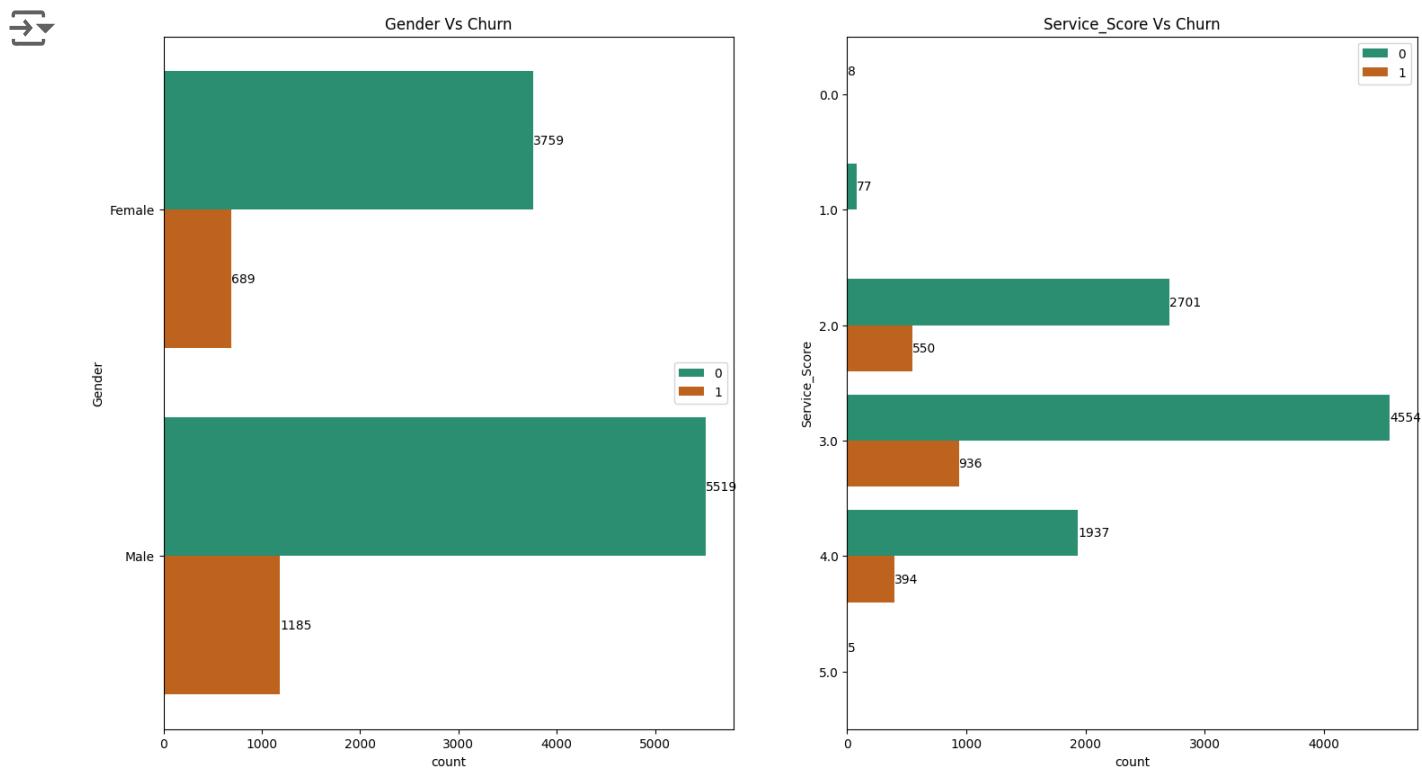


✓ Gender Vs Churn and Service_Score Vs Churn

```
1 # @title `Gender` Vs `Churn` ar Data Insights:
2 # @markdown **Data Insights:**
```

```
3 # @markdown > Plot 1 - `Gender`  
4 # @markdown > Plot 2 - It seems  
5  
6 fig = plt.figure(figsize=(18, 1  
7  
8 plt.subplot(121)  
9 ax = sns.countplot(df, y='Gende  
10 ax.legend(loc='center right')  
11 ax.set_title('Gender Vs Churn')  
12 ax.bar_label(ax.containers[0])  
13 ax.bar_label(ax.containers[1])  
14  
15 plt.subplot(122)  
16 ax = sns.countplot(df, y='Servi  
17 ax.legend(loc='upper right')  
18 ax.set_title('Service_Score Vs  
19 ax.bar_label(ax.containers[0])  
20 ax.bar_label(ax.containers[1]);
```

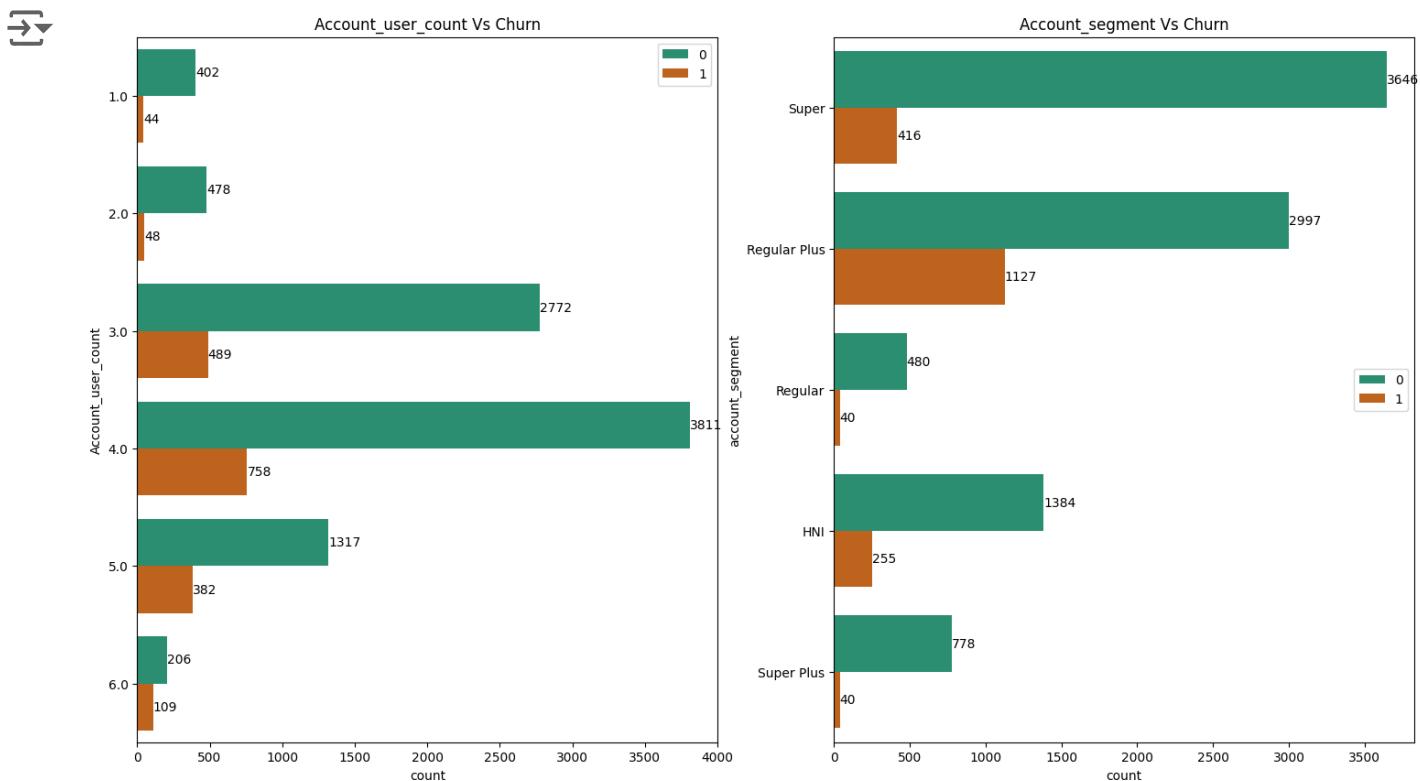
Plot 1 - Gender doesn't show any dependency over Churn.
Plot 2 - It seems like Service_Score has a positive relation with Churn.



✓ Account_user_count Vs Churn and account_segment Vs Churn

```
1 # @title `Account_user_count` Vs `Churn` and `account_segment` 
2 # @markdown **Data Insights:** 
3 # @markdown > Plot 1 - `Account_user_count` shows positive relat
```

```
4 # @markdown > Plot 2 – It seems like Plot Regular Plus` `Segment` us  
5                                         Account_user_count  
6 fig = plt.figure(figsize=(18, 10))      shows positive relation  
7                                         with Churn rate.  
8 plt.subplot(121)  
9 ax = sns.countplot(df, y='Account_user_count', hue='Churn', palette  
10 ax.legend(loc='upper right')           Regular Plus Segment  
11 ax.set_title('Account_user_count Vs Churn')  
12 ax.bar_label(ax.containers[0])          User are more likely to  
13 ax.bar_label(ax.containers[1])          Churn.  
14  
15 plt.subplot(122)  
16 ax = sns.countplot(df, y='account_segment', hue='Churn', palette  
17 ax.legend(loc='center right')  
18 ax.set_title('Account_segment Vs Churn')  
19 ax.bar_label(ax.containers[0])  
20 ax.bar_label(ax.containers[1]);
```

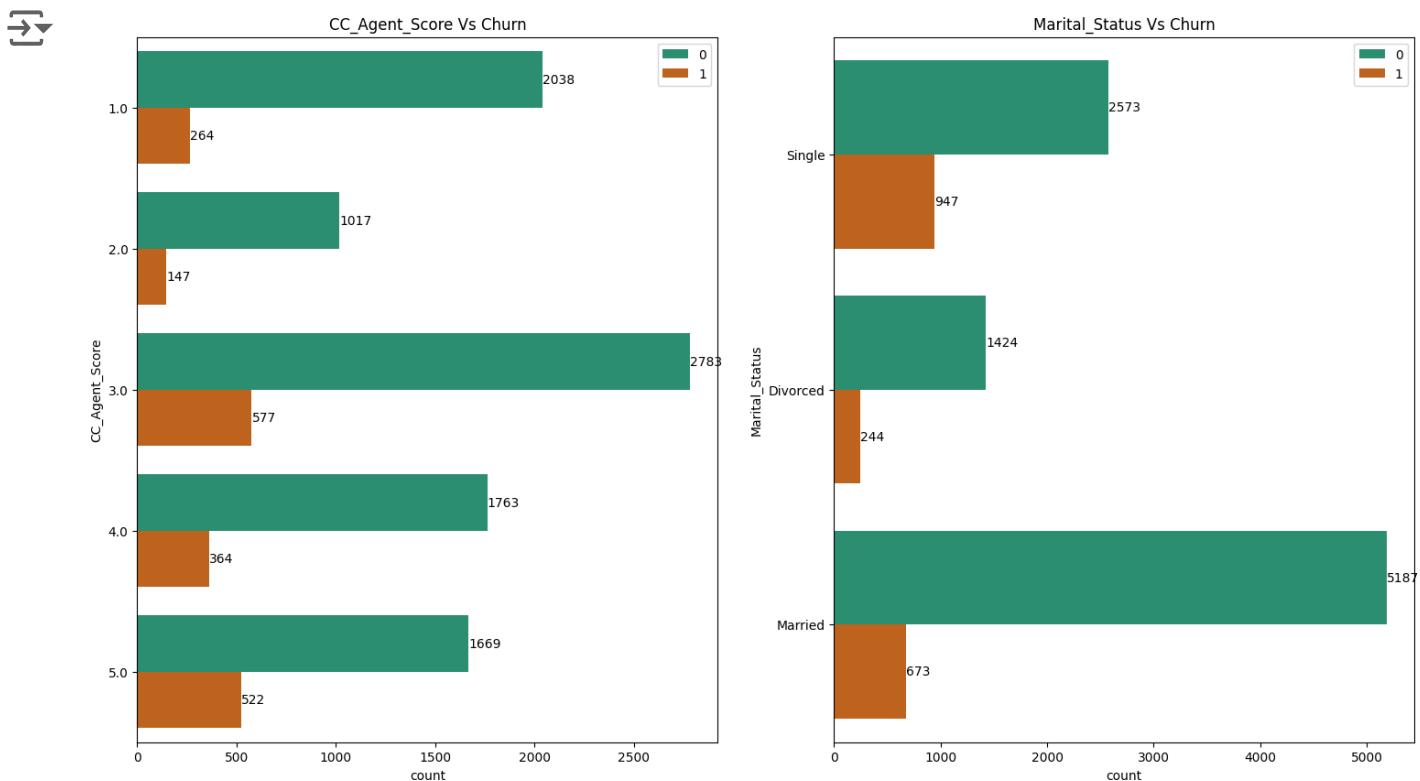


✓ CC_Agent_Score Vs Churn and Marital_Status Vs Churn

```
1 # @title `CC_Agent_Score` Vs `Churn` and `Marital_Status` Vs `Churn`  

2 # @markdown **Data Insights:**
```

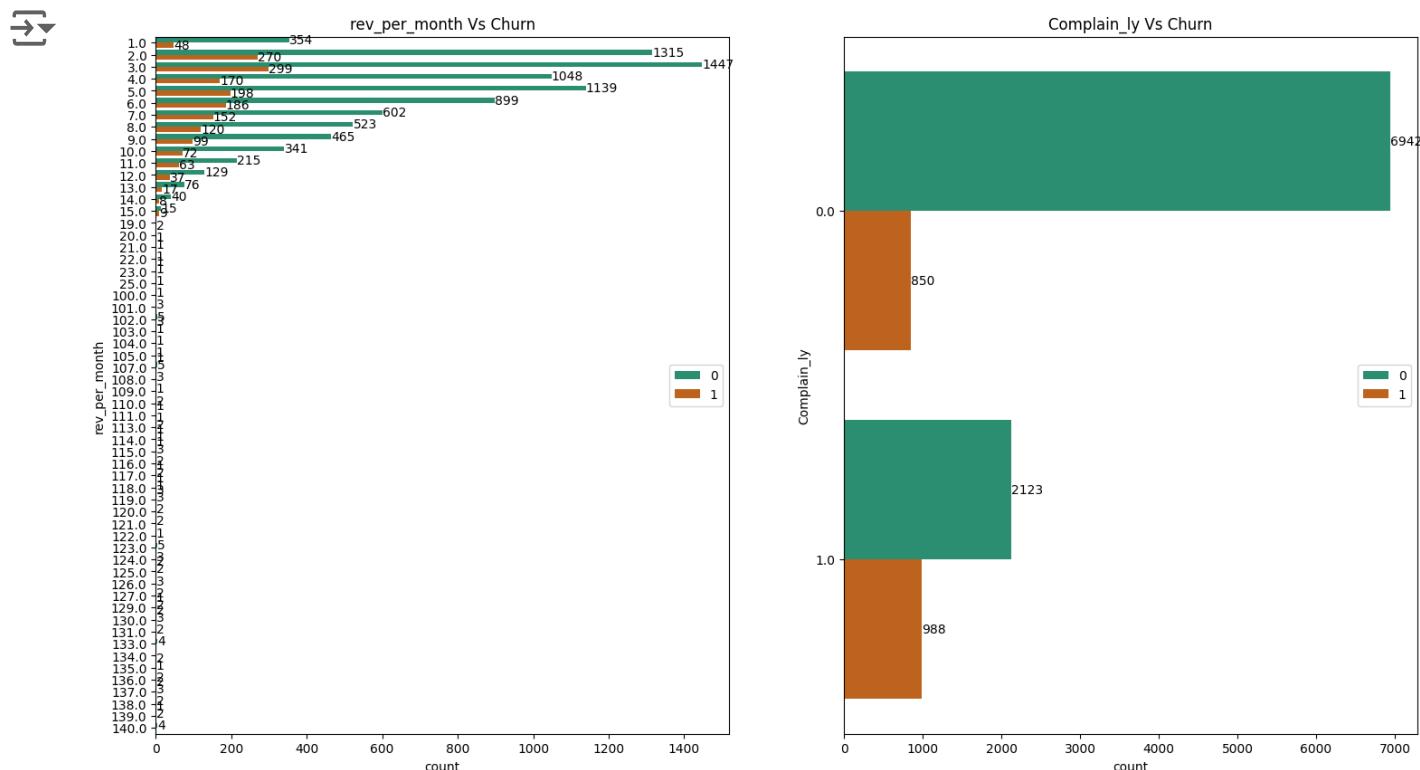
```
3 # @markdown > Plot 1 - `CC_Agent_Score` shows no relation with `Plot 1- CC_Agent_Score`
4 # @markdown > Plot 2 - It seems like `Single` `Marital_status` u shows no relation with
5
6 fig = plt.figure(figsize=(18, 10)) Churn rate.
7
8 plt.subplot(121) Plot 2 - It seems like
9 ax = sns.countplot(df, y='CC_Agent_Score', hue='Churn', palette= Single
10 ax.legend(loc='upper right') Marital_Status user
11 ax.set_title('CC_Agent_Score Vs Churn') are more likely to Churn.
12 ax.bar_label(ax.containers[0])
13 ax.bar_label(ax.containers[1])
14
15 plt.subplot(122)
16 ax = sns.countplot(df, y='Marital_Status', hue='Churn', palette= Marital_Status user
17 ax.legend(loc='upper right')
18 ax.set_title('Marital_Status Vs Churn')
19 ax.bar_label(ax.containers[0])
20 ax.bar_label(ax.containers[1]);
```



✓ rev_per_month Vs Churn and Complain_ly Vs Churn

```
1 # @title `rev_per_month` Vs `Churn` and `Complain_ly` Vs `Churn` Data Insights
2 # @markdown **Data Insights:**
```

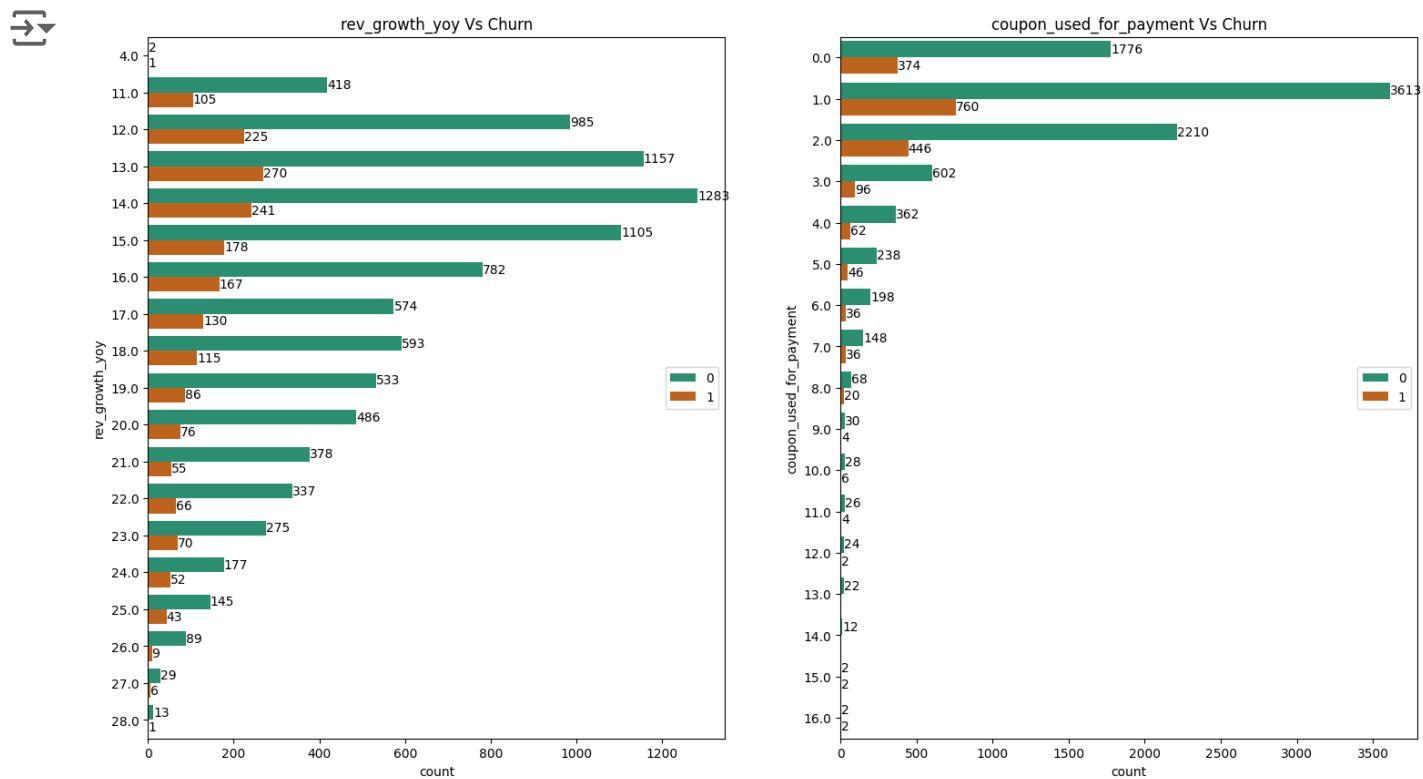
```
3 # @markdown > Plot 1 - `rev_per_month` shows positive relation w  
4 # @markdown > Plot 2 - Positive `Complain_ly` has more likely to  
5  
6 fig = plt.figure(figsize=(18, 10))  
7  
8 plt.subplot(121)  
9 ax = sns.countplot(df, y='rev_per_month', hue='Churn', palette=sns  
10 ax.legend(loc='center right')  
11 ax.set_title('rev_per_month Vs Churn')  
12 ax.bar_label(ax.containers[0])  
13 ax.bar_label(ax.containers[1])  
14  
15 plt.subplot(122)  
16 ax = sns.countplot(df, y='Complain_ly', hue='Churn', palette=sns  
17 ax.legend(loc='center right')  
18 ax.set_title('Complain_ly Vs Churn')  
19 ax.bar_label(ax.containers[0])  
20 ax.bar_label(ax.containers[1]);
```



- ✓ rev_growth_yoy Vs Churn and coupon_used_for_payment Vs Churn

```
1 # @title `rev_growth_yoy` Vs `Churn` and `coupon_used_for_payment`
```

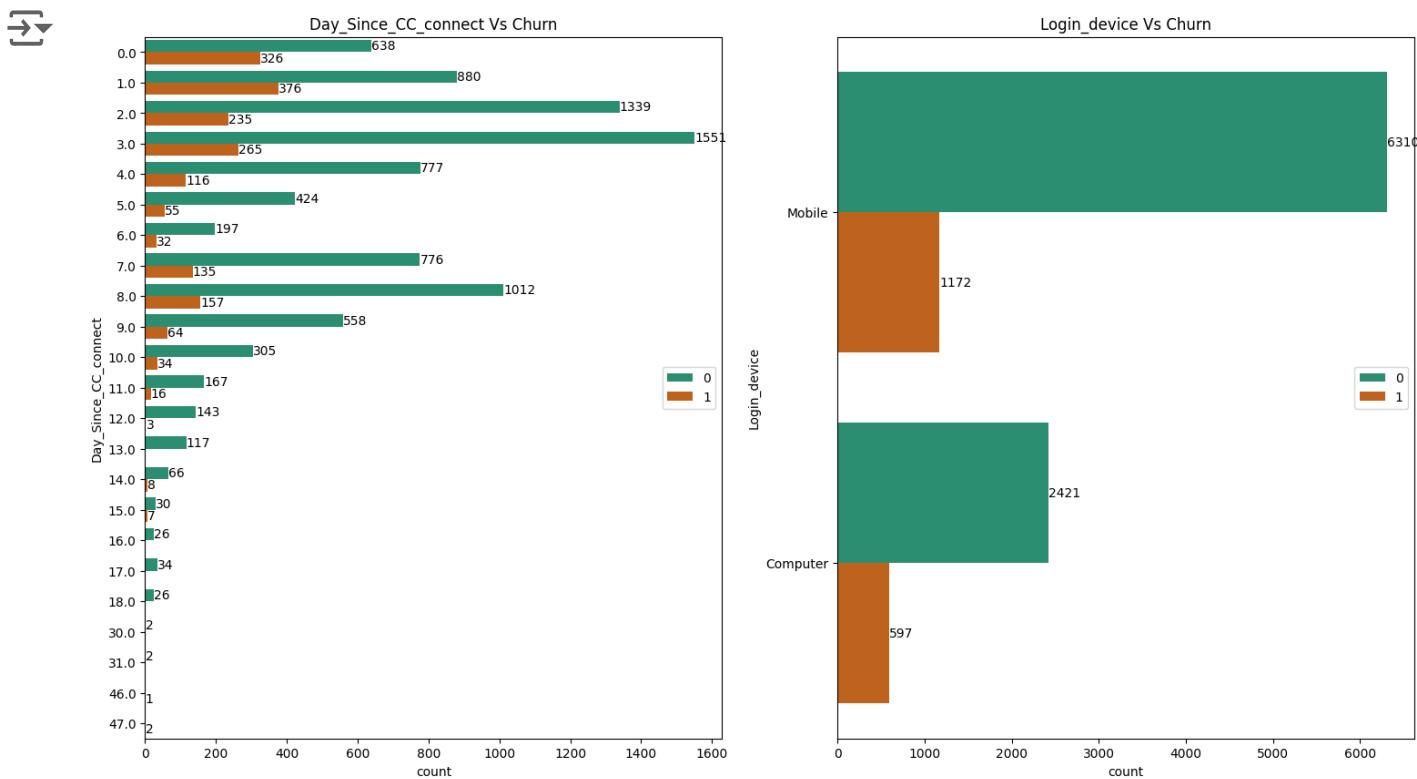
```
2 # @markdown **Data Insights:** Plot 1 - rev_growth_yoy shows positive relation
3 # @markdown > Plot 1 - `rev_growth_yoy` shows positive relation with
4 # @markdown > Plot 2 - `coupon_used_for_payment` shows positive
5 # @markdown > Churn rate.
6 fig = plt.figure(figsize=(18, 10)) Plot 2 -
7
8 plt.subplot(121) coupon_used_for_payment
9 ax = sns.countplot(df, y='rev_growth_yoy', hue='Churn', palette= shows positive relation with
10 ax.legend(loc='center right') Churn rate.
11 ax.set_title('rev_growth_yoy Vs Churn')
12 ax.bar_label(ax.containers[0])
13 ax.bar_label(ax.containers[1])
14
15 plt.subplot(122)
16 ax = sns.countplot(df, y='coupon_used_for_payment', hue='Churn',
17 ax.legend(loc='center right')
18 ax.set_title('coupon_used_for_payment Vs Churn')
19 ax.bar_label(ax.containers[0])
20 ax.bar_label(ax.containers[1]);
```



✓ Day_Since_CC_connect Vs Churn and Login_device Vs Churn

```
1 # @title `Day_Since_CC_connect` Vs `Churn` and `Login_device` Vs Data Insights
2 # @markdown **Data Insights:**
```

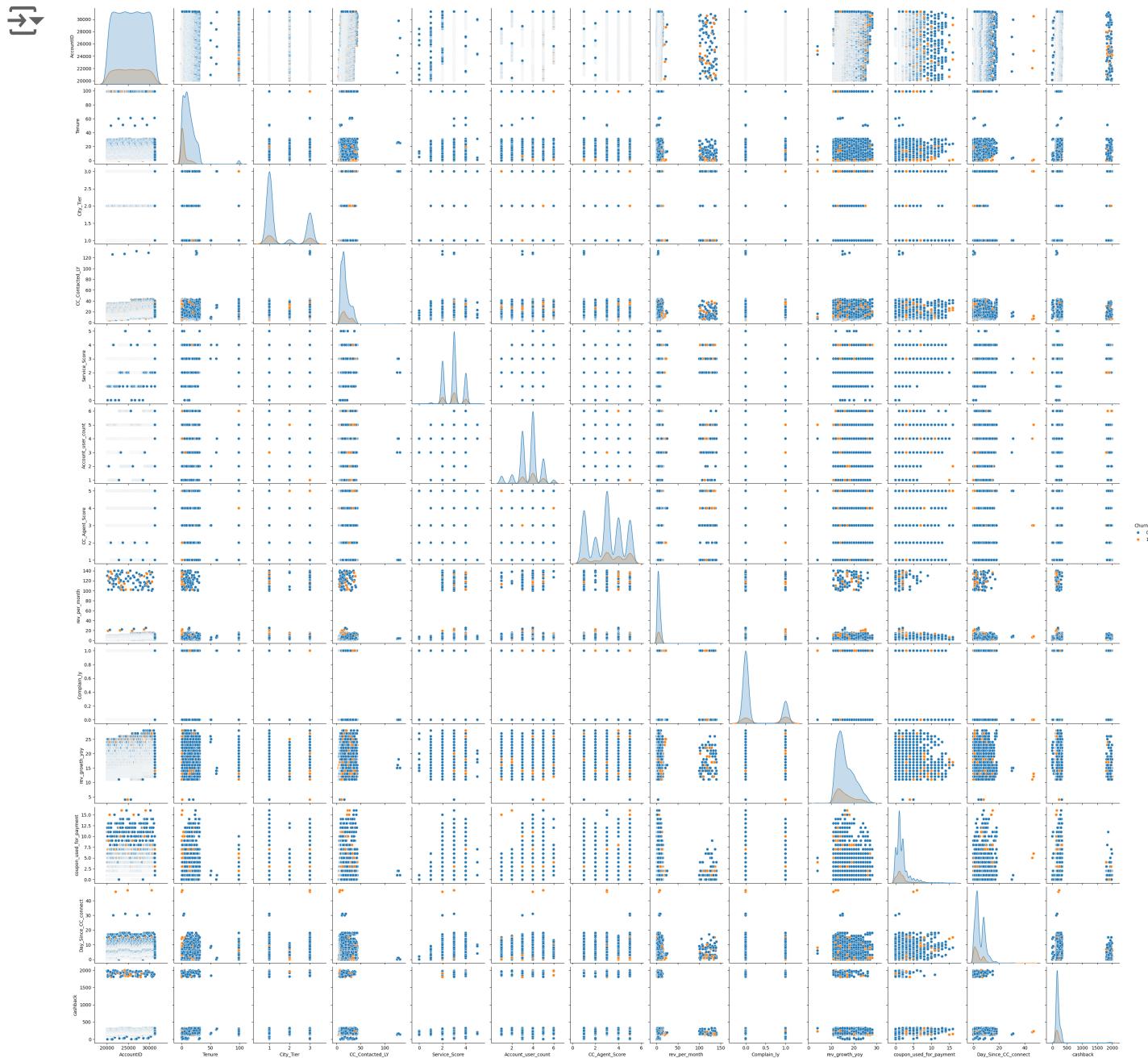
```
3 # @markdown > Plot 1 - `Day_Since_CC_Connect` shows positive relation  
4 # @markdown > Plot 2 - `Login_device` shows not much of the relationship.  
5  
6 fig = plt.figure(figsize=(18, 10))  
7  
8 plt.subplot(121)  
9 ax = sns.countplot(df, y='Day_Since_CC_Connect', hue='Churn', palette=sns  
10 ax.legend(loc='center right')  
11 ax.set_title('Day_Since_CC_Connect Vs Churn')  
12 ax.bar_label(ax.containers[0])  
13 ax.bar_label(ax.containers[1])  
14  
15 plt.subplot(122)  
16 ax = sns.countplot(df, y='Login_Device', hue='Churn', palette=sns  
17 ax.legend(loc='center right')  
18 ax.set_title('Login_Device Vs Churn')  
19 ax.bar_label(ax.containers[0])  
20 ax.bar_label(ax.containers[1]);
```



✓ Multivariate Analysis

✓ Pair plot of the dataset

```
1 sns.pairplot(df, hue='Churn', diag_kind='kde');
```



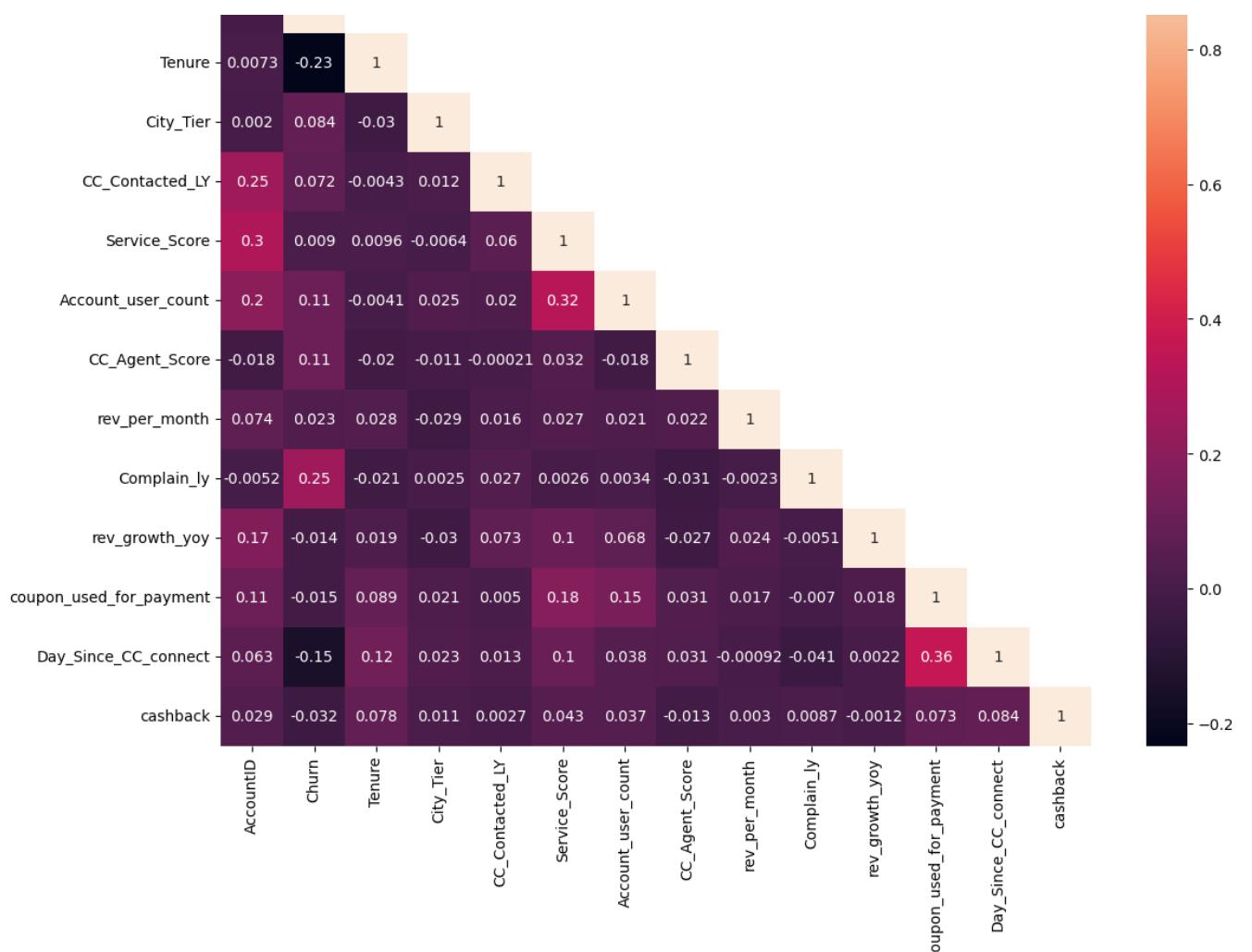
Data Insights:

The pair-plot shown above indicated that the independent variable are week or poor predictors of target variable as we the density of independent variable overlaps with the density of target variable.

Correlation Matrix

```
1 cor = df.select_dtypes(exclude='O').corr()  
2 mask = np.triu(cor, +1)  
3 fig,ax= plt.subplots(figsize=(13, 10))  
4 sns.heatmap(cor, mask=mask, annot=True, fmt='.2g');
```





Data Insights:

The Data doesn't shows much of the correlation with each other except few columns (i.e., Day_Since_CC_connect with coupon_used_for_payment, Account_user_count with Service_Score and negative correlation with Tenure and Churn along with Day_Since_CC_connect .

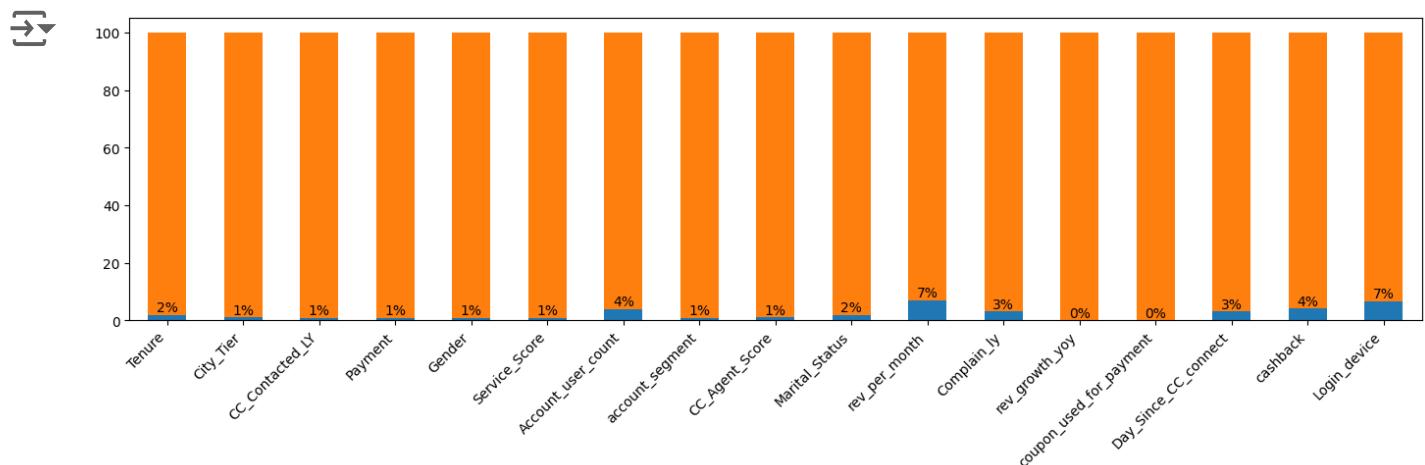
✓ Removing Unwanted column

```
1 df = df.drop('AccountID', axis=1)
2
3 df.head()
```

✓ Handling missing values

➤ Before Treatment

[Show code](#)



▼ After Treatment

```

1 # @title After Treatment
2
3 # @markdown > * Out of 19 variables we have null values in 17 variables
4 # @markdown > * Using "Median" to impute null values where variable is continuous in nature because Median is less prone to outliers
5 # @markdown > * Using "Mode" to impute null values where variable is categorical
6
7
8 for column in df.select_dtypes(exclude='O'):
9     df[column].fillna(df[column].median(), inplace=True)
10
11
12 for column in df.select_dtypes(include='O'):
13     df[column].fillna(df[column].mode()[0], inplace=True)

```

- Out of 19 variables we have null values in 17 variables
- Using "Median" to impute null values where variable is continuous in nature because Median is less prone to outliers
- Using "Mode" to impute null values where variable is categorical

```
15      df['Churn'].mean() > df['Gender'].mode(), replace=True)
14      when compared
15
16 df.isnull().sum()
```

- with mean.
- Using "Mode: to impute null values where variables are categorical in nature. We have treated null values variable by variable as each and every variable is unique in its nature.

```
Churn          0
Tenure         0
City_Tier      0
CC_Contacted_LY 0
Payment        0
Gender         0
Service_Score  0
Account_user_count 0
account_segment 0
CC_Agent_Score  0
Marital_Status  0
rev_per_month   0
Complain_ly    0
rev_growth_yoy 0
coupon_used_for_payment 0
Day_Since_CC_connect 0
cashback       0
Login_device   0
dtype: int64
```

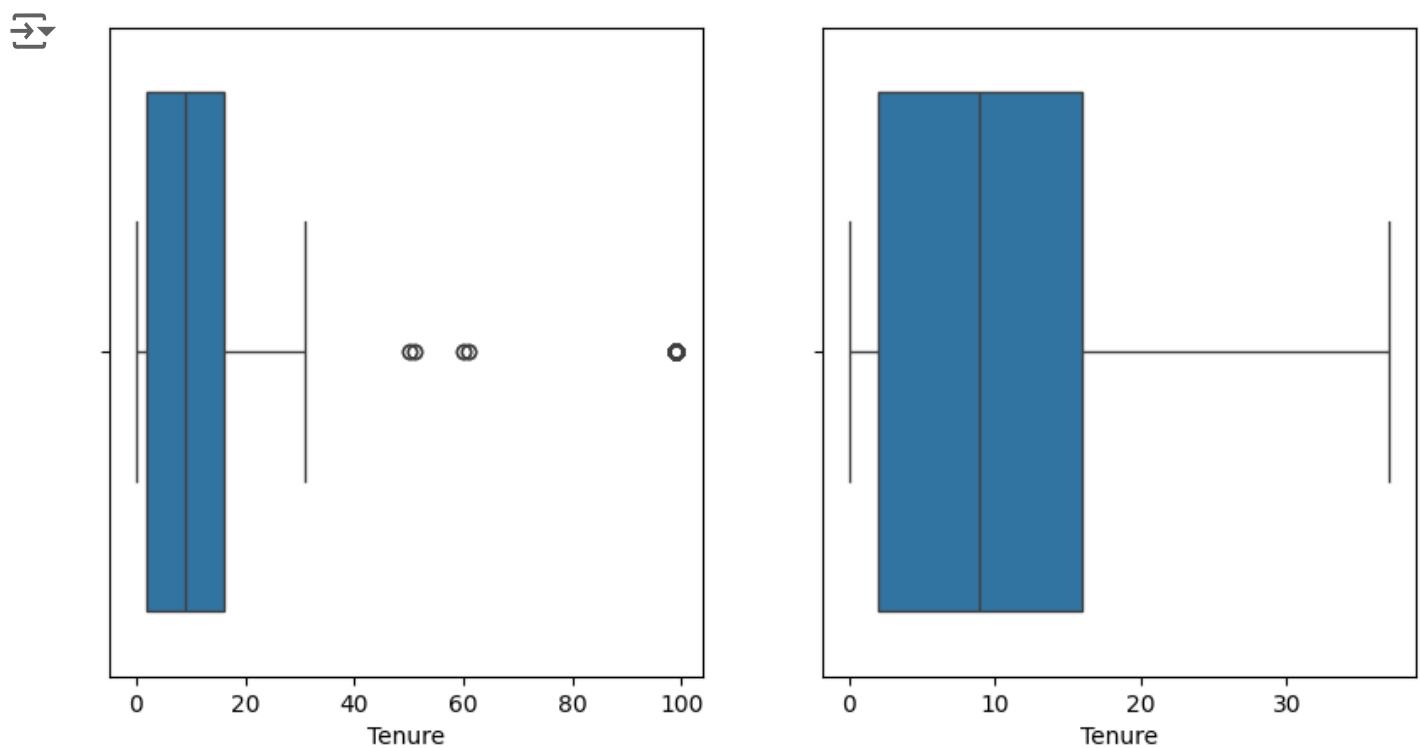
✓ Checking for Outliers on continuous variables

```
1 #treating outlier
2 def remove_outlier(col):
3     sorted(col)
4     Q1,Q3=col.quantile([0.25,0.75])
5     IQR=Q3-Q1
6     lower_range= Q1-(1.5 * IQR)
7     upper_range= Q3+(1.5 * IQR)
8     return lower_range, upper_range
```

▼ Before and After Treatment

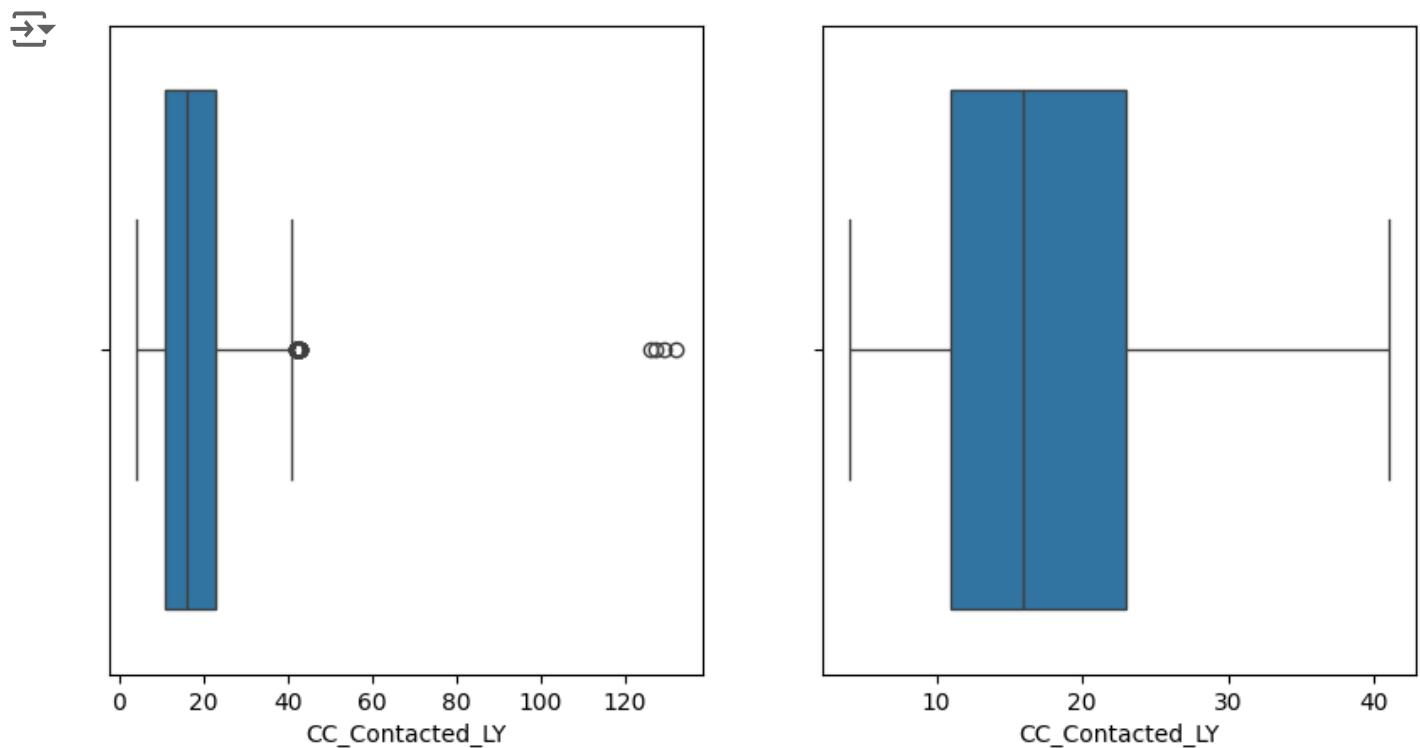
▼ Tenure

```
1 # @title Tenure
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['Tenure'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['Tenure'])
7 df['Tenure']=np.where(df['Tenure']>up,up,df['Tenure'])
8 df['Tenure']=np.where(df['Tenure']<lw,lw,df['Tenure'])
9
10 sns.boxplot(df['Tenure'], orient='h', ax=ax[1]);
```



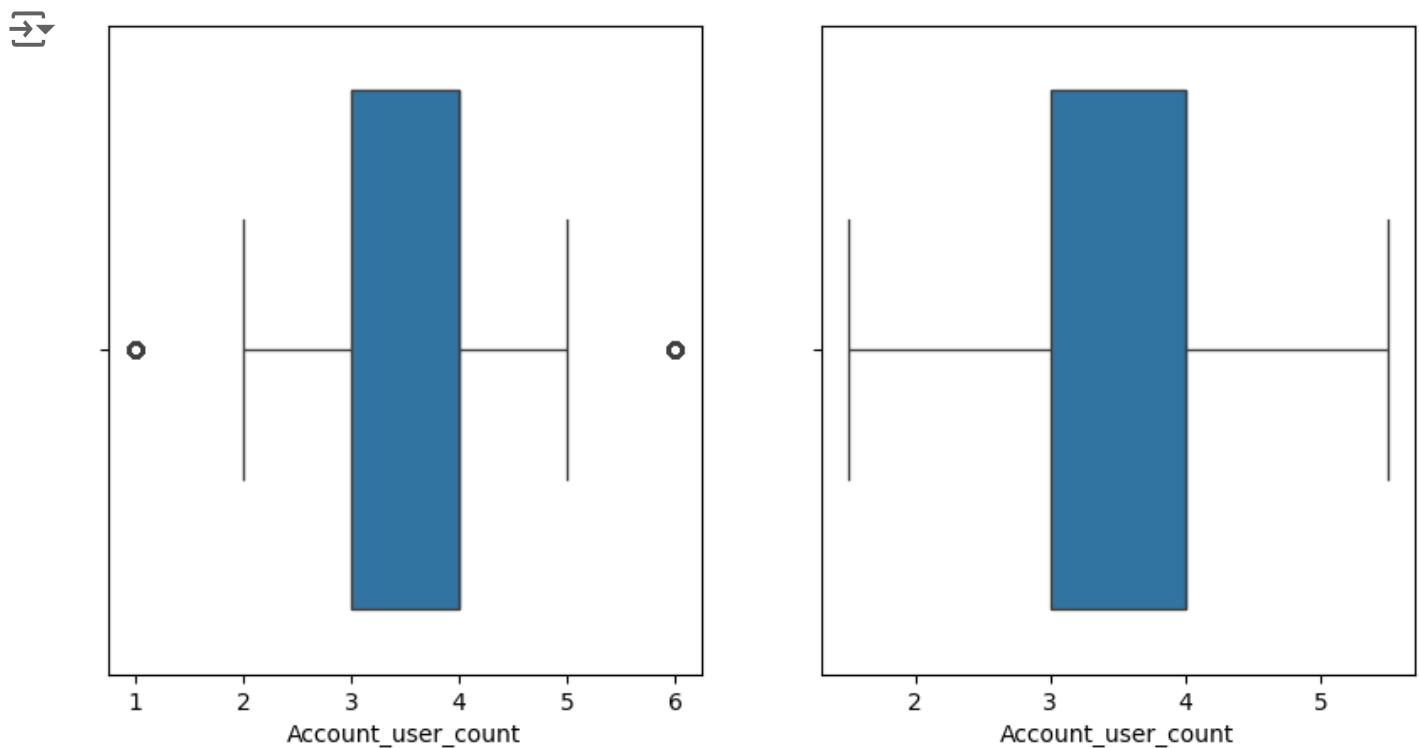
✓ CC_Contacted_LY

```
1 # @title CC_Contacted_LY
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['CC_Contacted_LY'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['CC_Contacted_LY'])
7 df['CC_Contacted_LY']=np.where(df['CC_Contacted_LY']>up,up,df['CC_C
8 df['CC_Contacted_LY']=np.where(df['CC_Contacted_LY']<lw,lw,df['CC_C
9
10 sns.boxplot(df['CC_Contacted_LY'], orient='h', ax=ax[1]);
```



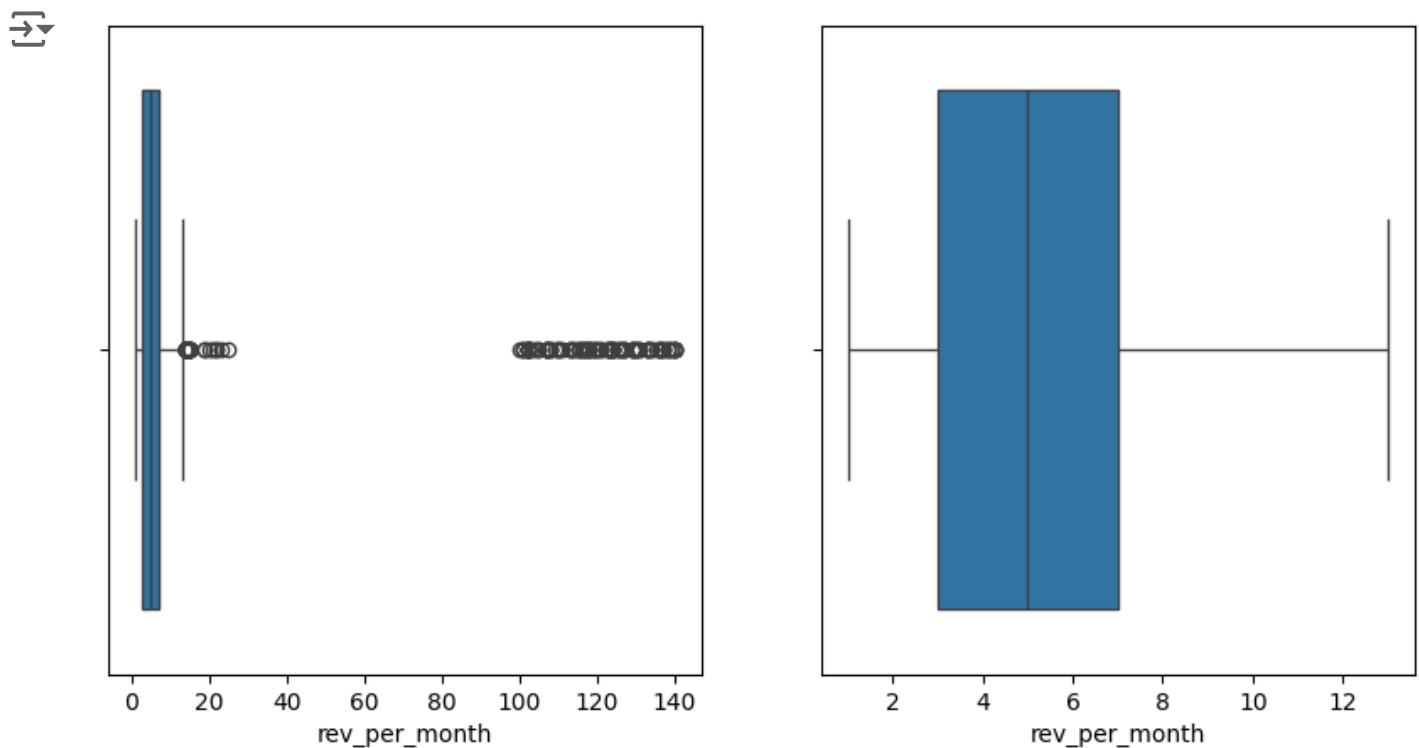
✓ Account_user_count

```
1 # @title Account_user_count
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['Account_user_count'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['Account_user_count'])
7 df['Account_user_count']=np.where(df['Account_user_count']>up,up,d
8 df['Account_user_count']=np.where(df['Account_user_count']<lw,lw,d
9
10 sns.boxplot(df['Account_user_count'], orient='h', ax=ax[1]);
```



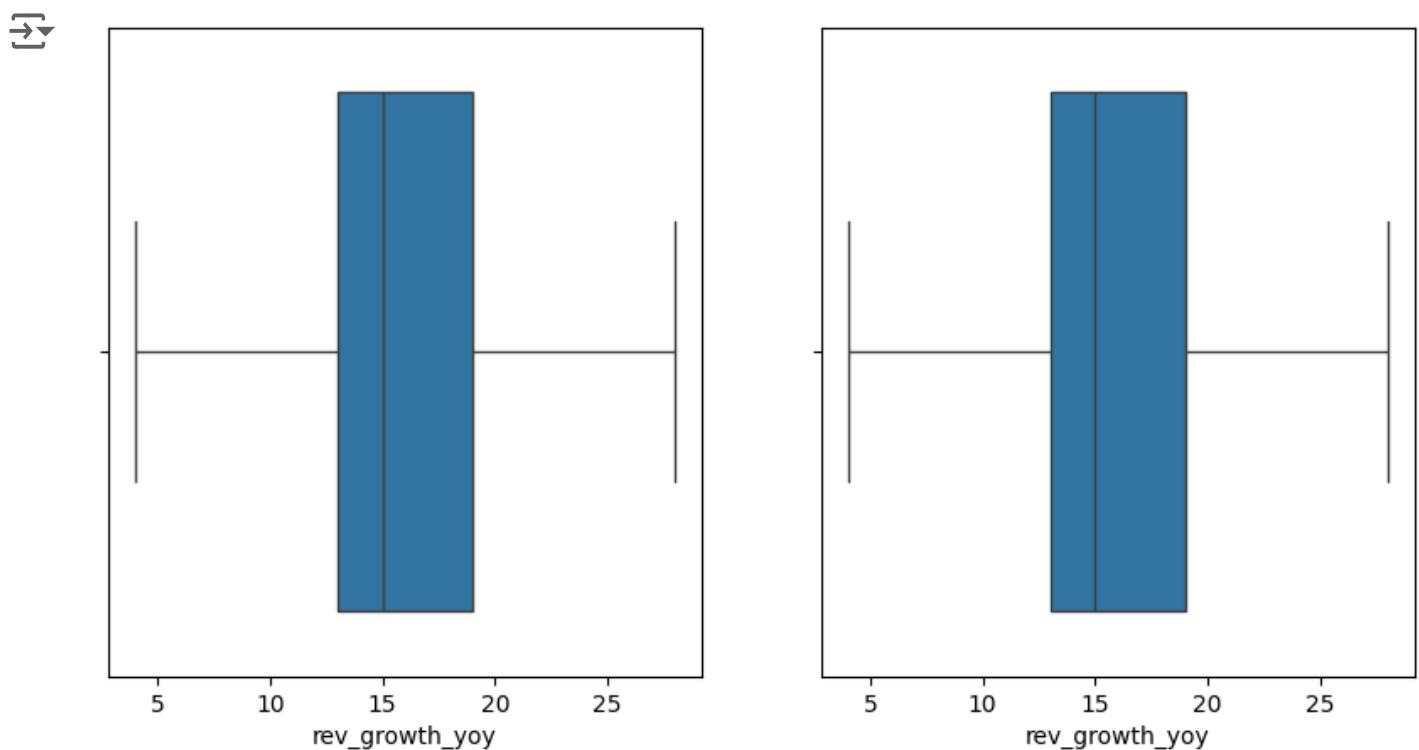
✓ rev_per_month

```
1 # @title rev_per_month
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['rev_per_month'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['rev_per_month'])
7 df['rev_per_month']=np.where(df['rev_per_month']>up,up,df['rev_per_month'])
8 df['rev_per_month']=np.where(df['rev_per_month']<lw,lw,df['rev_per_month'])
9
10 sns.boxplot(df['rev_per_month'], orient='h', ax=ax[1]);
```



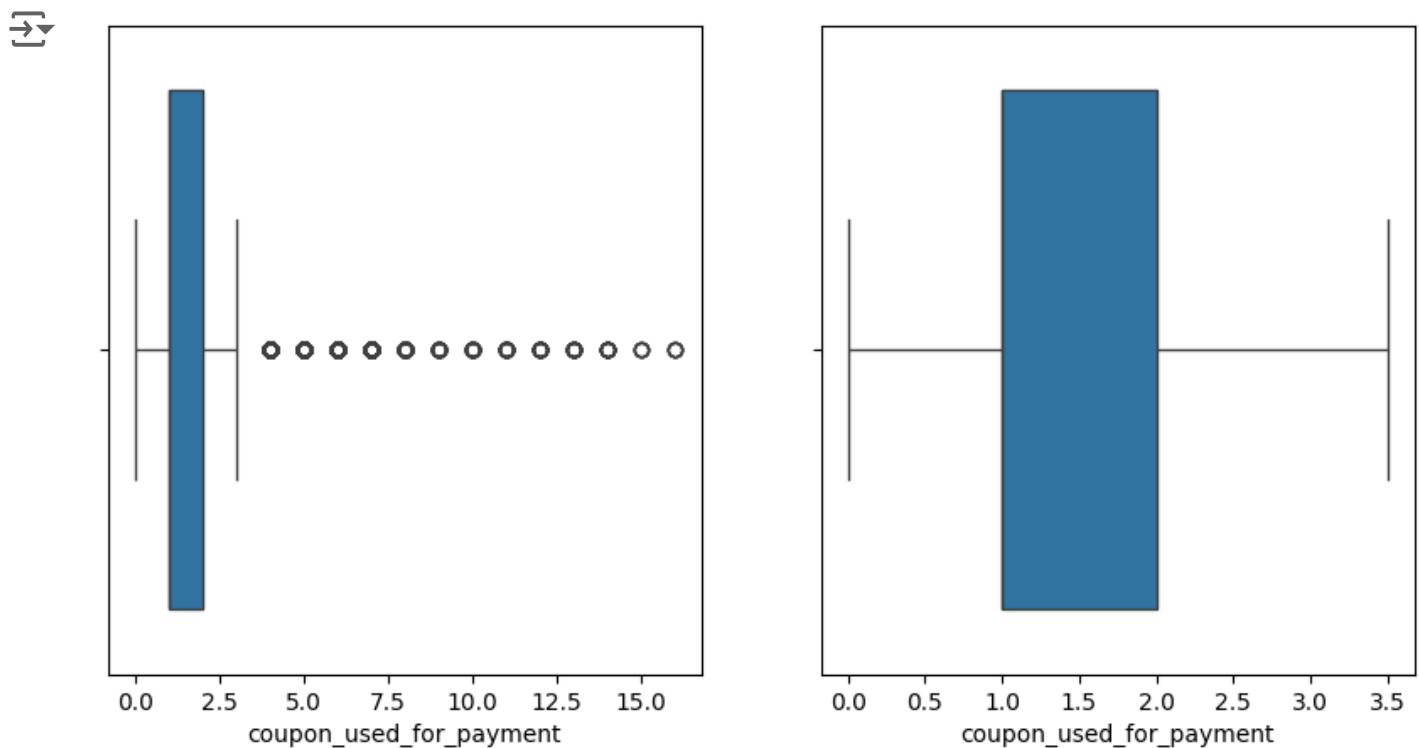
✓ rev_growth_yoy

```
1 # @title rev_growth_yoy
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['rev_growth_yoy'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['rev_growth_yoy'])
7 df['rev_growth_yoy']=np.where(df['rev_growth_yoy']>up,up,df['rev_gr
8 df['rev_growth_yoy']=np.where(df['rev_growth_yoy']<lw,lw,df['rev_gr
9
10 sns.boxplot(df['rev_growth_yoy'], orient='h', ax=ax[1]);
```



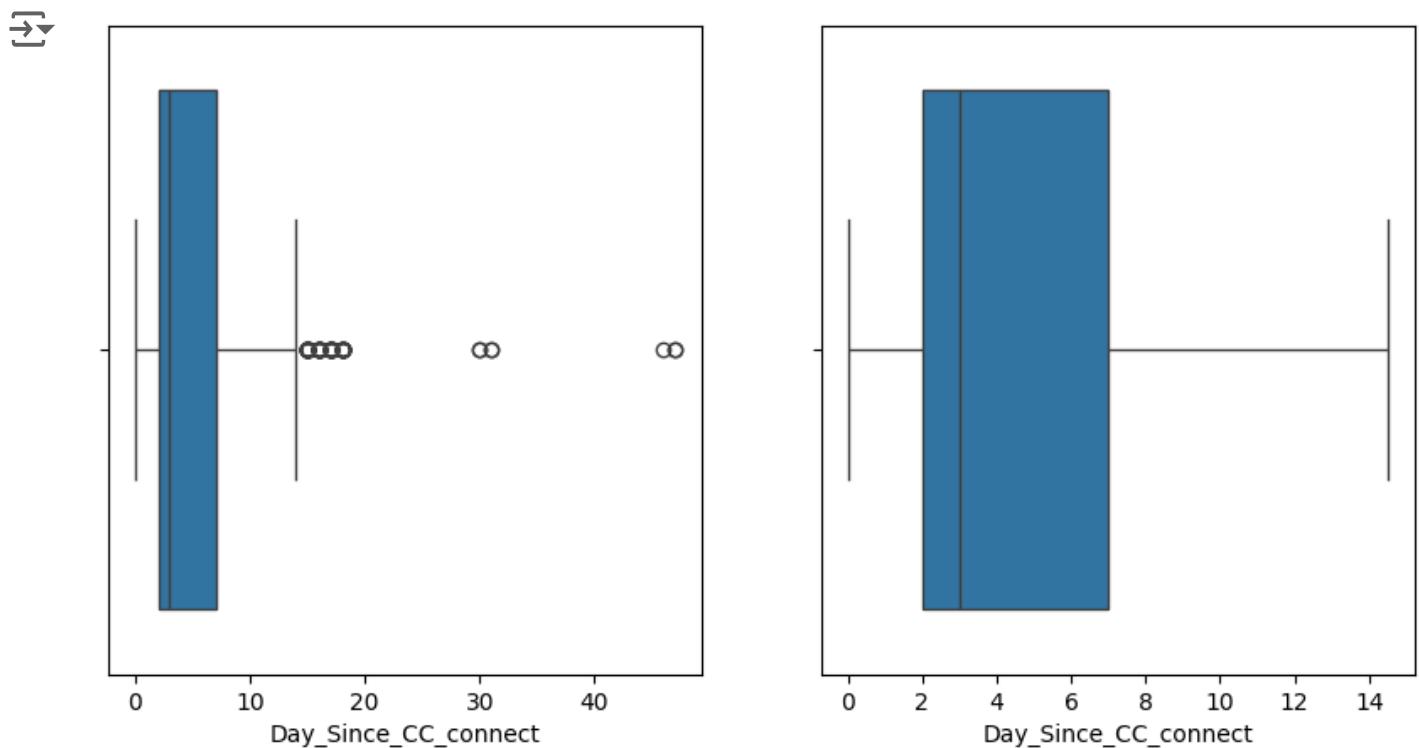
✓ coupon_used_for_payment

```
1 # @title coupon_used_for_payment
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['coupon_used_for_payment'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['coupon_used_for_payment'])
7 df['coupon_used_for_payment']=np.where(df['coupon_used_for_payment'>lw, df['coupon_used_for_payment'], up)
8 df['coupon_used_for_payment']=np.where(df['coupon_used_for_payment'<up, df['coupon_used_for_payment'], df['coupon_used_for_payment']+lw)
9
10 sns.boxplot(df['coupon_used_for_payment'], orient='h', ax=ax[1]);
```



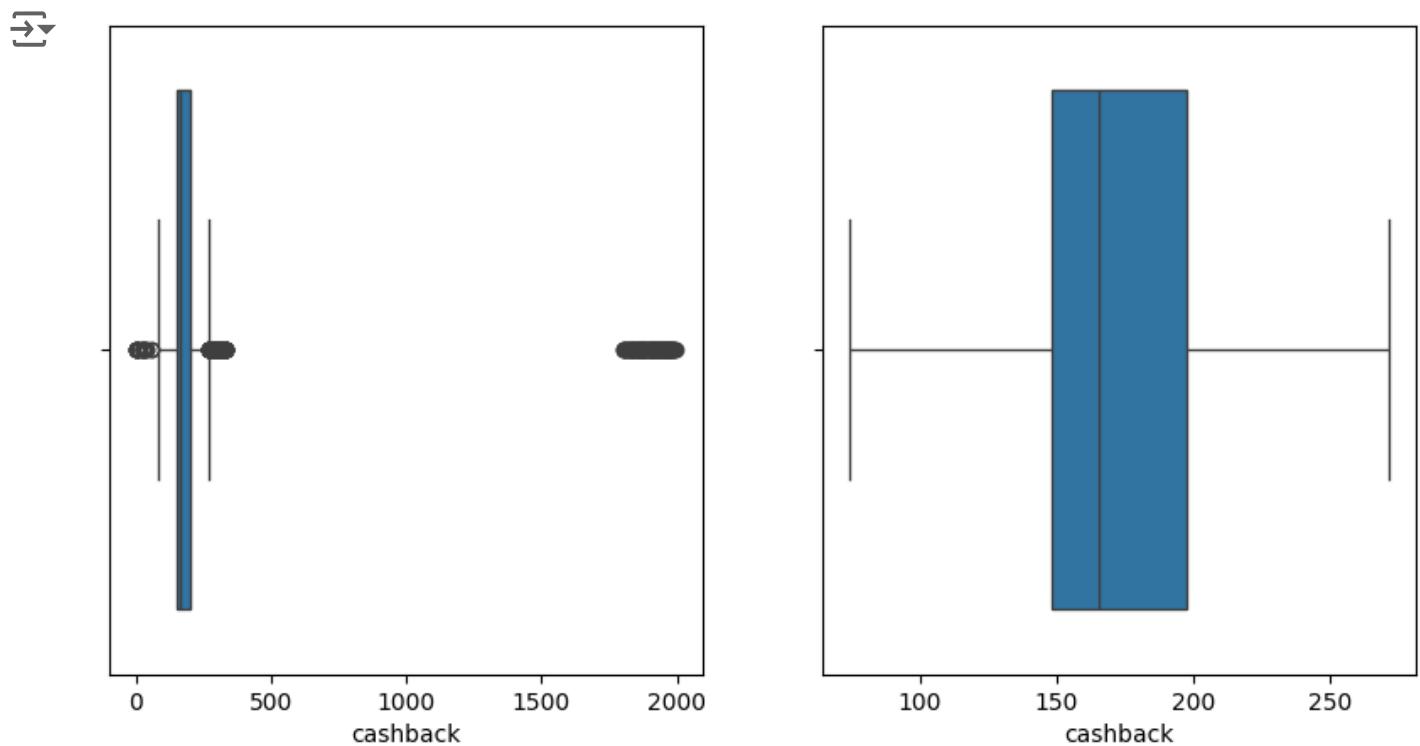
✓ Day_Since_CC_connect

```
1 # @title Day_Since_CC_connect
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['Day_Since_CC_connect'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['Day_Since_CC_connect'])
7 df['Day_Since_CC_connect']=np.where(df['Day_Since_CC_connect']>up,u
8 df['Day_Since_CC_connect']=np.where(df['Day_Since_CC_connect']<lw,l
9
10 sns.boxplot(df['Day_Since_CC_connect'], orient='h', ax=ax[1]);
```



✓ cashback

```
1 # @title cashback
2
3 fig, ax = plt.subplots(1, 2, figsize=(10, 5))
4 sns.boxplot(df['cashback'], orient='h', ax=ax[0])
5
6 lw, up=remove_outlier(df['cashback'])
7 df['cashback']=np.where(df['cashback']>up,up,df['cashback'])
8 df['cashback']=np.where(df['cashback']<lw,lw,df['cashback'])
9
10 sns.boxplot(df['cashback'], orient='h', ax=ax[1]);
```



✓ Data Preprocessing

```
1 from sklearn.model_selection import train_test_split  
2 from sklearn.preprocessing import LabelEncoder  
3 from sklearn.preprocessing import MinMaxScaler  
4 from imblearn.combine import SMOTEENN  
5 from sklearn.cluster import KMeans
```

Splitting the Dataset

```
1 X = df.drop('Churn', axis=1)  
2 y = df[['Churn']]
```

Dividing the dataset into X and y. Where X containing independent variables and y containing dependent variable.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
2 X_train.shape, X_test.shape, y_train.shape, y_test.shape  
3  
4 → ((9008, 17), (2252, 17), (9008, 1), (2252, 1))
```

Splitting the dataset into 80% of train and 20% of test dataset.

- ✓ Encoding the categorical variables using LabelEncoder

```

1 obj_cols = X_train.select_dtypes('object').columns
2
3 for col in X_train.columns:
4     if col in obj_cols:
5         encoder = LabelEncoder()
6         encoder.fit(X_train[col])
7         X_train[col] = encoder.transform(X_train[col])
8         X_test[col] = encoder.transform(X_test[col])
9
10
11 X_train.head()

```

	Tenure	City_Tier	CC_Contacted_LY	Payment	Gender	Service_Score	Acc
6848	2.0	1.0		18.0	2	1	2.0
8131	13.0	1.0		23.0	2	0	3.0
8139	13.0	1.0		29.0	2	0	2.0
9984	19.0	1.0		14.0	1	1	4.0
7701	4.0	1.0		8.0	2	0	3.0

Selecting all the columns which are of object type such as (Payment , Gender , Marital_Status , etc.,) and converting into numerical data using LabelEncoder from scikit-learn library.

✓ Scaling the data using MinMaxScaler

```

1 scaling = MinMaxScaler()
2 scaling.fit(X_train)
3 X_train = pd.DataFrame(scaling.transform(X_train), columns=X_train.columns)
4 X_test = pd.DataFrame(scaling.transform(X_test), columns=X_test.columns)

```

```
1 X_train.head()
```

	Tenure	City_Tier	CC_Contacted_LY	Payment	Gender	Service_Score	Ac
6848	0.054054	0.0	0.378378	0.4	1.0		0.4
8131	0.351351	0.0	0.513514	0.4	0.0		0.6
8139	0.351351	0.0	0.675676	0.4	0.0		0.4
9984	0.513514	0.0	0.270270	0.2	1.0		0.8
7701	0.108108	0.0	0.108108	0.4	0.0		0.6

Data Insights

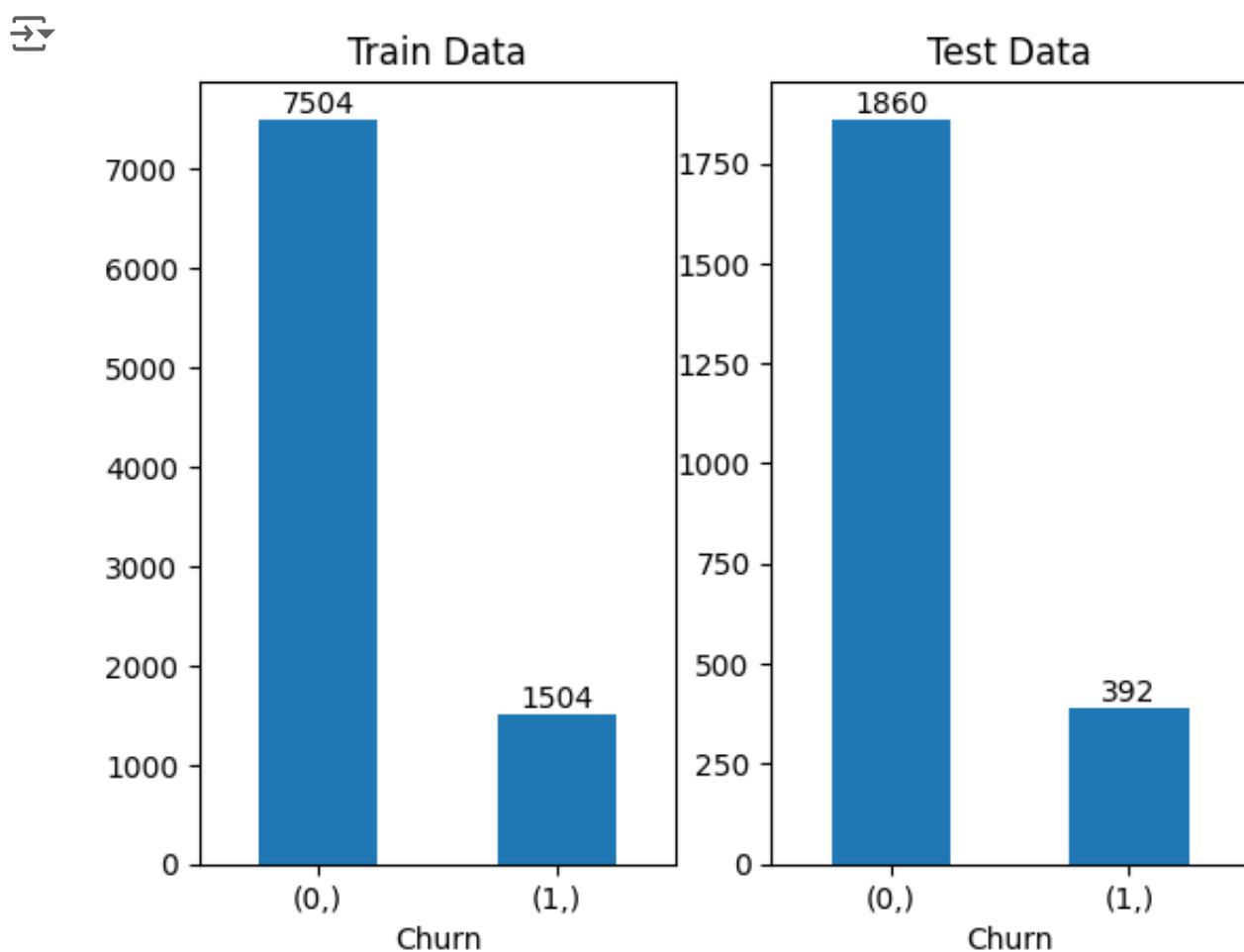
Since the data doesn't follow uniform distribution or bell curve shape, we need to use normalization instead of standardization for better results.

✓ Data balancing using SMOTE

Since, the data is imbalanced we need to try to balance the Churn .

✓ Before sampling

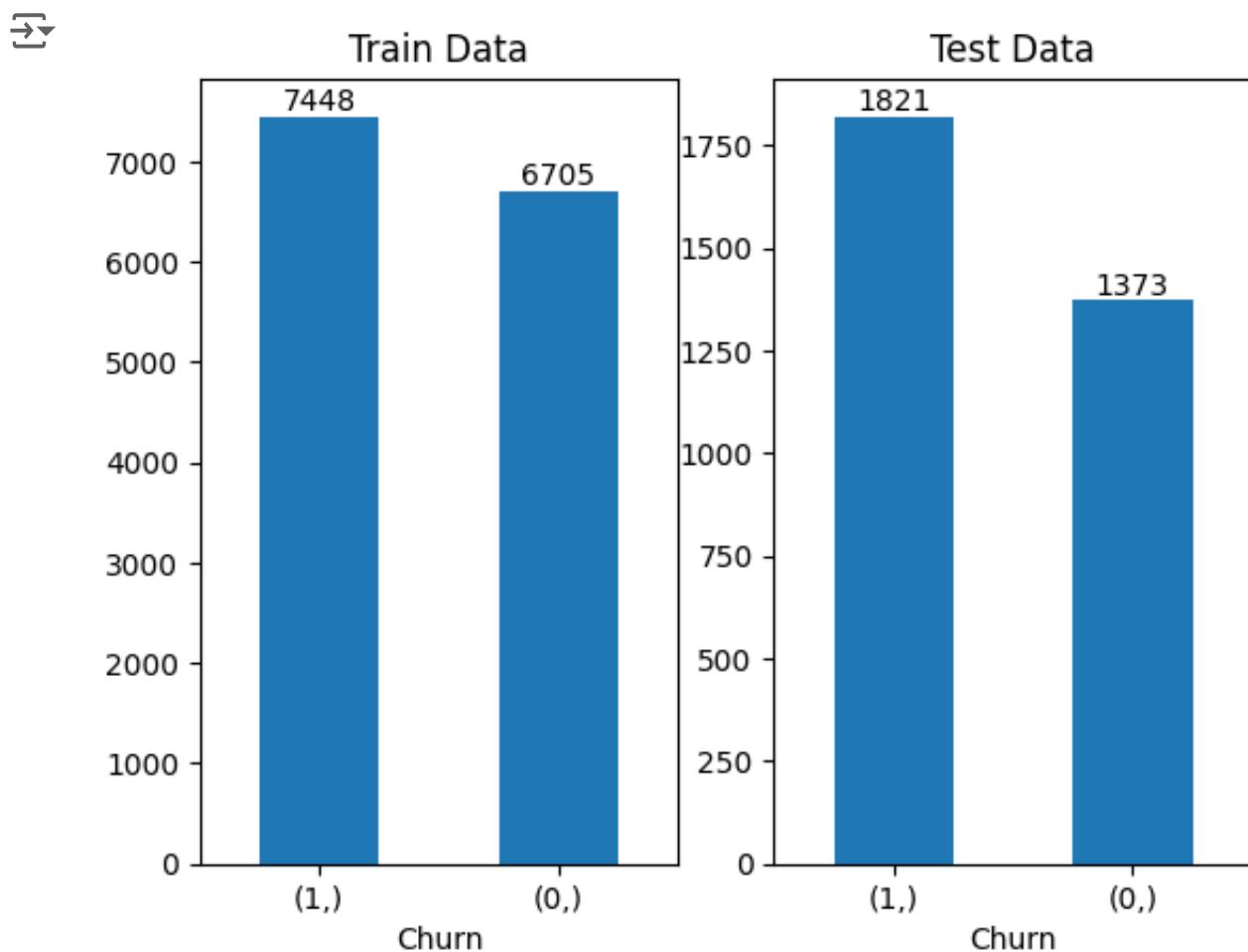
```
1 # @title Before sampling
2
3 plt.subplot(121)
4 ax = y_train.value_counts().plot(kind='bar', rot='horizontal')
5 ax.bar_label(ax.containers[0])
6 ax.set_title('Train Data')
7
8 plt.subplot(122)
9 ax = y_test.value_counts().plot(kind='bar', rot='horizontal')
10 ax.bar_label(ax.containers[0])
11 ax.set_title('Test Data');
```



```
1 sm = SMOTEENN(random_state=8)
2 sm.fit(X_train, y_train)
3 X_train_sm, y_train_sm = sm.fit_resample(X_train, y_train)
4 X_test_sm, y_test_sm = sm.fit_resample(X_test, y_test)
```

✓ After Sampling

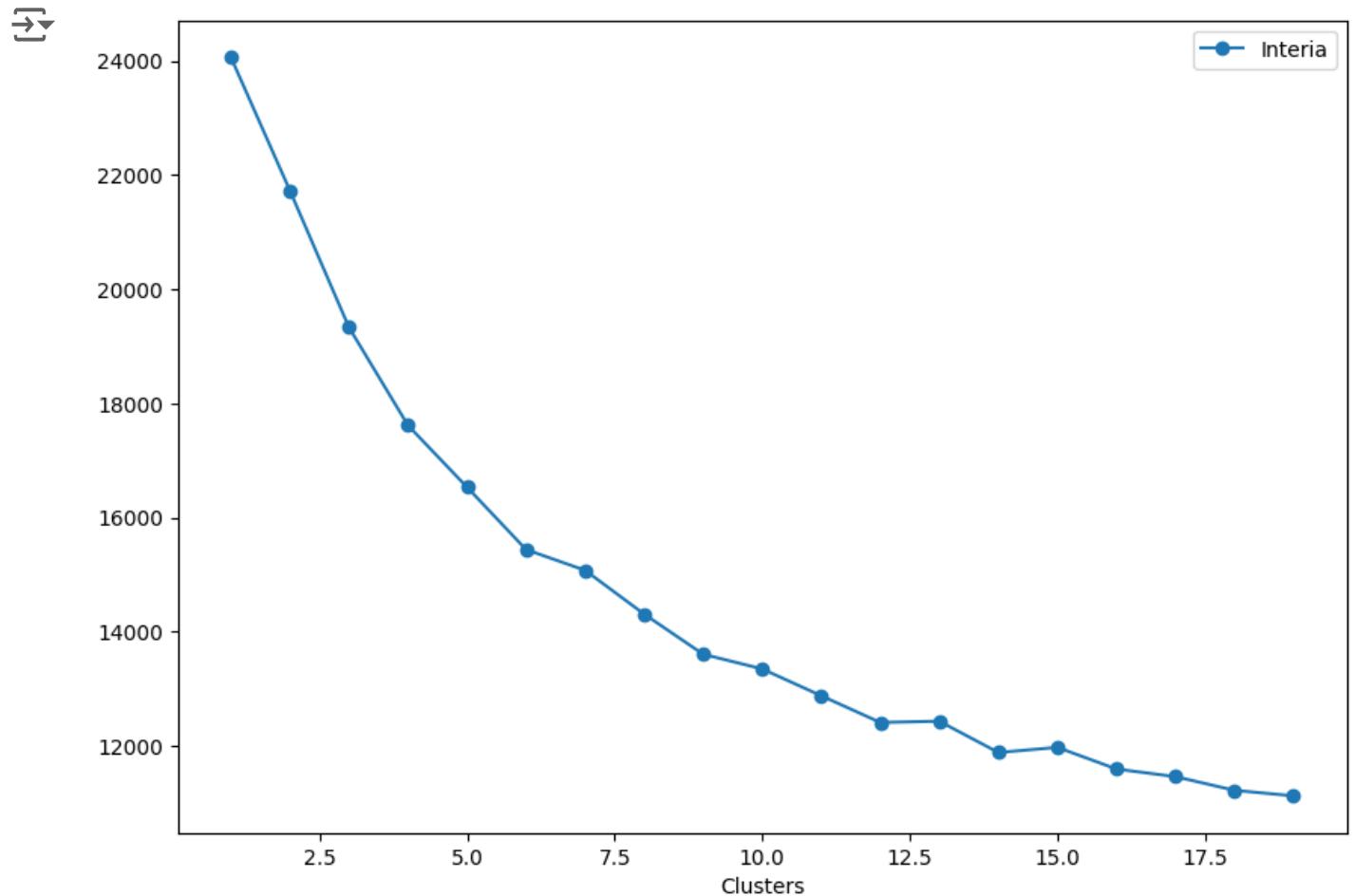
```
1 # @title After Sampling
2
3 plt.subplot(121)
4 ax = y_train_sm.value_counts().plot(kind='bar', rot='horizontal')
5 ax.bar_label(ax.containers[0])
6 ax.set_title('Train Data')
7
8 plt.subplot(122)
9 ax = y_test_sm.value_counts().plot(kind='bar', rot='horizontal')
10 ax.bar_label(ax.containers[0])
11 ax.set_title('Test Data');
```



✓ Clustering using k-means

```
1 interia = []
2
3 for i in range(1, 20):
4     k_means = KMeans(n_clusters=i, n_init='auto')
5     k_means.fit(X_train_sm)
6     inter = k_means.inertia_
7     interia.append(inter)
```

```
1 cluster = pd.DataFrame({  
2     'Clusters': list(range(1, 20)),  
3     'Interia': interia  
4}  
5 })  
6  
7 cluster.plot(x='Clusters', y='Interia', marker='o', figsize=(10, 7)
```



```
1 k_means = KMeans(n_clusters=7, random_state=8, n_init='auto')
2 k_means.fit(X_train_sm)
3 train_labels = k_means.predict(X_train_sm)
4 train_labels

→ array([2, 0, 5, ..., 4, 4, 6], dtype=int32)

1 test_labels = k_means.predict(X_test_sm)
2 test_labels[:10]

→ array([3, 1, 4, 3, 4, 3, 4, 3, 2, 1], dtype=int32)

1 X_train_means = X_train_sm.copy()
2 X_train_means['Clusters'] = train_labels
3
4 X_test_means = X_test_sm.copy()
5 X_test_means['Clusters'] = test_labels

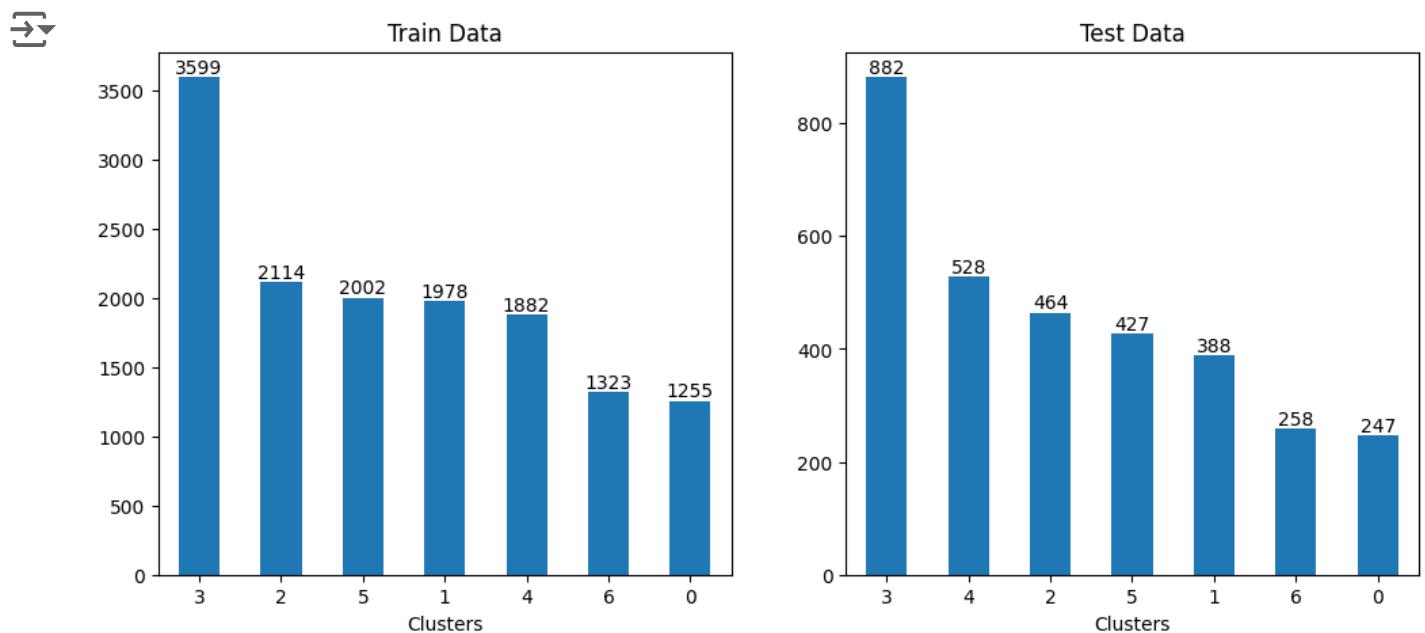
1 X_train_final = X_train_means.copy()
2 X_test_final = X_test_means.copy()
3
4 y_train_final = y_train_sm.copy()
5 y_test_final = y_test_sm.copy()

1 X_train_final.head()
```

→

	Tenure	City_Tier	CC_Contacted_LV	Payment	Gender	Service_Score	Accou
0	0.351351	0.0	0.513514	0.4	0.0		0.6
1	0.351351	0.0	0.675676	0.4	0.0		0.4
2	0.513514	0.0	0.270270	0.2	1.0		0.8
3	0.108108	0.0	0.108108	0.4	0.0		0.6
4	0.216216	0.0	0.243243	0.4	1.0		0.4

```
1 fig = plt.figure(figsize=(12, 5))
2
3 plt.subplot(121)
4 ax = X_train_final['Clusters'].value_counts().plot(kind='bar', rot=
5 ax.bar_label(ax.containers[0])
6 ax.set_title('Train Data')
7
8 plt.subplot(122)
9 ax = X_test_final['Clusters'].value_counts().plot(kind='bar', rot=
10 ax.bar_label(ax.containers[0])
11 ax.set_title('Test Data');
```



✓ Training the model

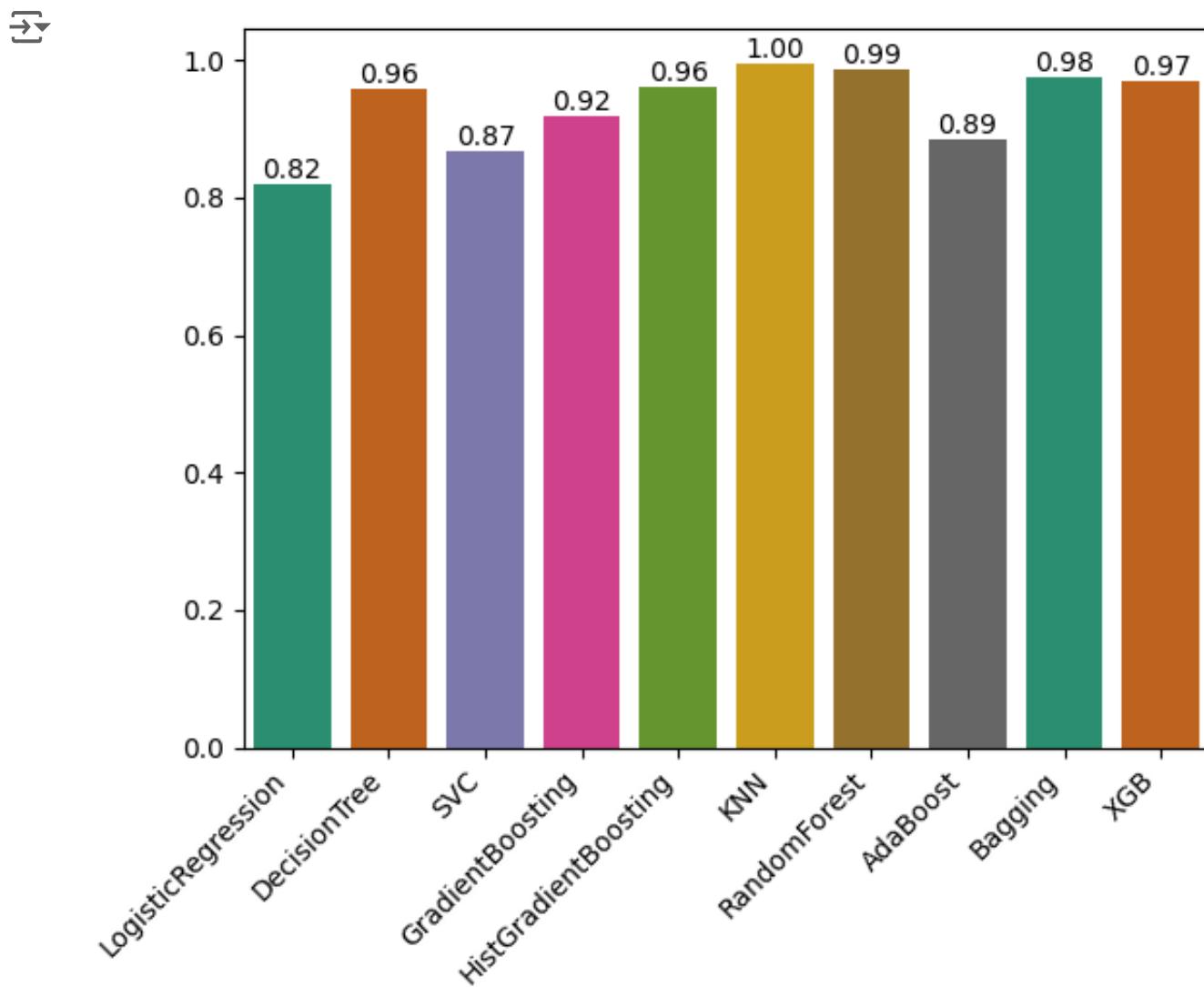
```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.svm import SVC
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.ensemble import (GradientBoostingClassifier,
6                               HistGradientBoostingClassifier,
7                               RandomForestClassifier,
8                               AdaBoostClassifier,
9                               BaggingClassifier,
10                              StackingClassifier)
11 from xgboost import XGBClassifier
12
13
14 from sklearn.model_selection import cross_val_score, GridSearchCV,
15 from sklearn.metrics import roc_curve, RocCurveDisplay, roc_auc_sc
16 from sklearn.metrics import confusion_matrix, ConfusionMatrixDispl
17 from sklearn.metrics import precision_recall_fscore_support, accur

1 models = {
2     'LogisticRegression': LogisticRegression(random_state=8),
3     'DecisionTree': DecisionTreeClassifier(random_state=8),
4     'SVC': SVC(random_state=8),
5     'GradientBoosting': GradientBoostingClassifier(random_state=8)
6     'HistGradientBoosting': HistGradientBoostingClassifier(random_
7     'KNN': KNeighborsClassifier(n_neighbors=2),
8     'RandomForest': RandomForestClassifier(random_state=8),
9     'AdaBoost': AdaBoostClassifier(random_state=8),
10    'Bagging': BaggingClassifier(random_state=8),
11    'XGB': XGBClassifier(),
12 }
```

```
1 # Calculates the accuracy, precision, recall, and f1-score score
2 def calculate_results(name, y_true, y_pred):
3     accuracy = accuracy_score(y_true, y_pred)
4     precision, recall, fscore, _ = precision_recall_fscore_support
5
6
7     return {
8         'Model Name': name,
9         'Accuracy score': accuracy,
10        'Precision score': precision,
11        'Recall score': recall,
12        'F1 score': fscore,
13    }
14
15
16 train_results = []
17 test_results = []
18
19 for name, model in models.items():
20     result = cross_val_score(model, X=X_train_final, y=y_train_final)
21     train_results.append(np.mean(result))
22
23     model.fit(X_train_final, y_train_final)
24     y_pred = model.predict(X_test_final)
25
26 test_results.append(calculate_results(name, y_test_final, y_pred))
```

✓ Train data score

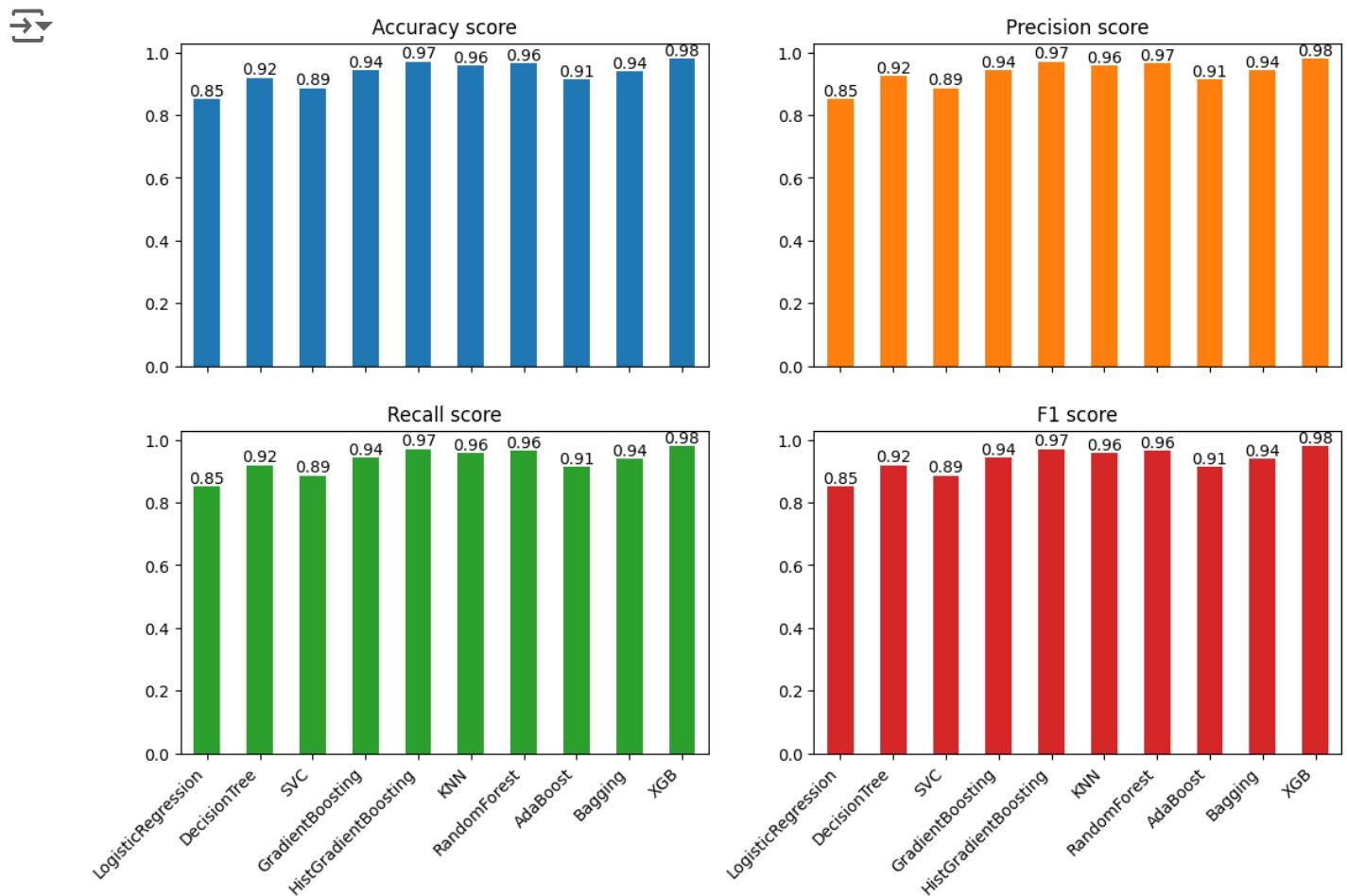
```
1 # @title Train data score
2
3 ax = sns.barplot(x=list(models.keys()), y=train_results, palette=''
4 for container in ax.containers:
5     ax.bar_label(container, fmt='%.2f')
6 ax.set_xticklabels(list(models.keys()), rotation=45, ha='right');
```



✓ Test data score

```
1 # @title Test data score
2
3 axis = pd.DataFrame(test_results).plot(kind='bar', subplots=True,
4 for ax in axis.flatten():
5     ax.get_legend().remove()
```

```
6     ax.set_xticklabels(list(models.keys()), rotation=45, ha='right')
7     ax.bar_label(ax.containers[0], fmt='%.2f');
```



About Data:

Best Models:

- *DecisionTree*
- *Gradient Boosting*
- *HistGradient Boosting*
- *KNN*
- *RandomForest*
- *Bagging*
- *XGB*

```
1 best_models = ['DecisionTree', 'GradientBoosting', 'HistGradientBo
2 best_models_tuned = ['DecisionTreeTuned', 'GradientBoostingTuned',
3
4
5
6
7
8
9
10
11
12
13
14
```

```
1 train_results, test_results = [], []
2
3 def pre_result_calculator(name, model):
4     train_pred = model.predict(X_train_final)
5     test_pred = model.predict(X_test_final)
6
7     train_score = calculate_results(name, y_train_final, train_pred)
8     test_score = calculate_results(name, y_test_final, test_pred)
9
10    train_results.append(train_score)
11    test_results.append(test_score)
12
13    return (y_train_final, train_pred), \
14          (y_test_final, test_pred), \
15          pd.DataFrame(train_score, index=[0]), \
16          pd.DataFrame(test_score, index=[0])
```

```

1 def plot(train, test):
2     fig, ax = plt.subplots(2, 2, figsize=(12, 12))
3     fig.set_tight_layout({'pad': 5.0,
4                           'w_pad': 5.0,
5                           'h_pad': 5.0})
6     ax = ax.flatten()
7
8     ConfusionMatrixDisplay.from_predictions(train[0], train[1], cmap
9     ax[0].set_title('Train Confusion Matrix')
10    RocCurveDisplay.from_predictions(train[0], train[1], ax=ax[1])
11    ax[1].set_title('Train Roc Curve')
12
13    ConfusionMatrixDisplay.from_predictions(test[0], test[1], cmap
14    ax[2].set_title('Test Confusion Matrix')
15    RocCurveDisplay.from_predictions(test[0], test[1], ax=ax[3])
16    ax[3].set_title('Test Roc Curve')

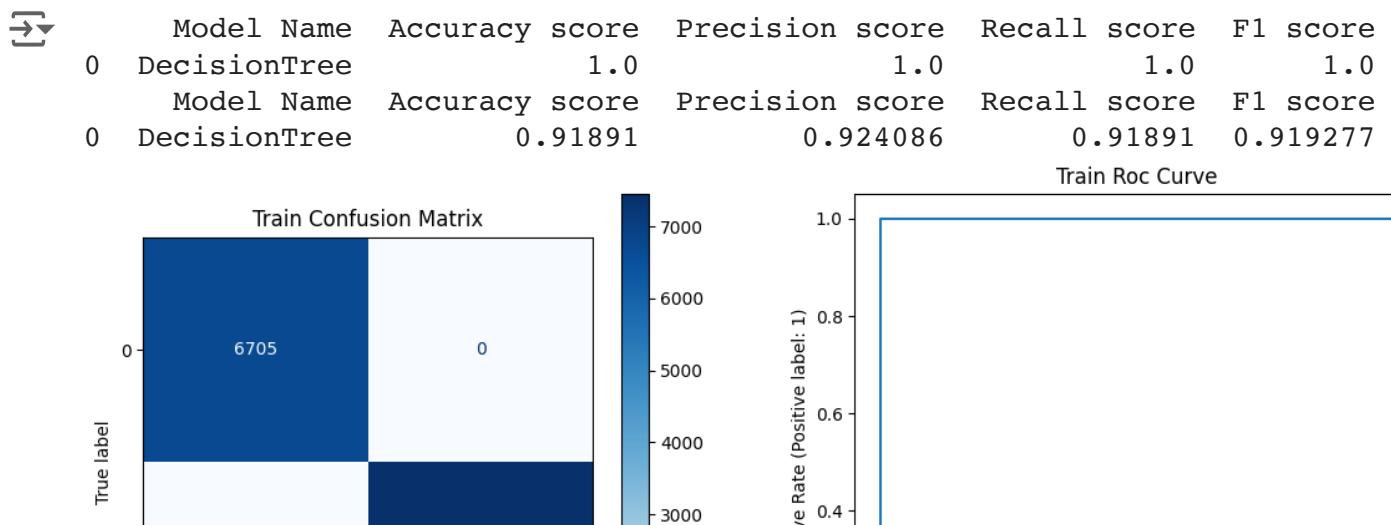
```

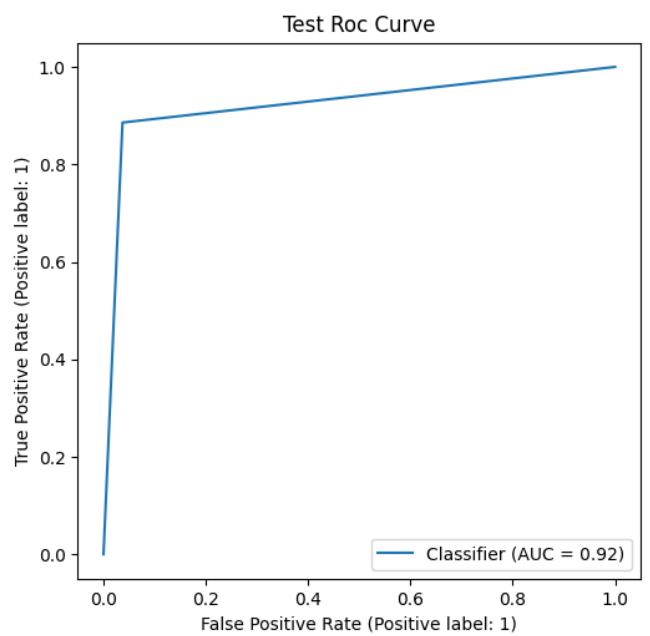
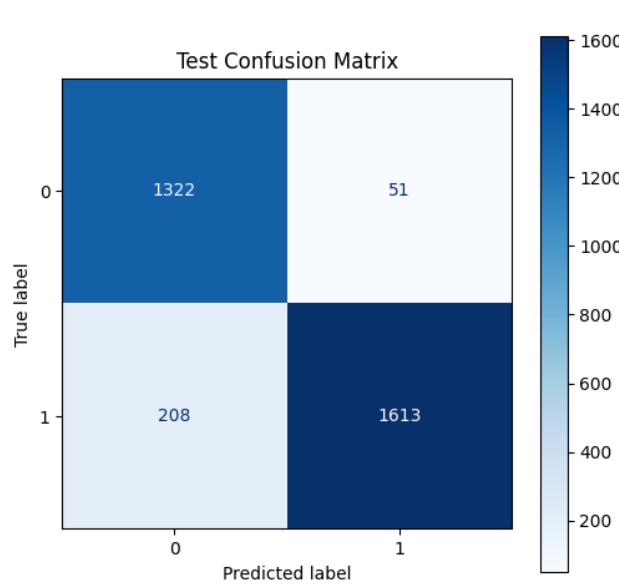
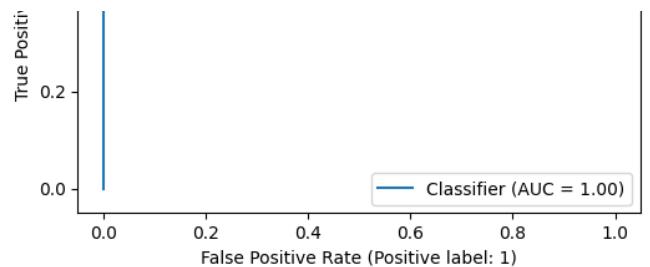
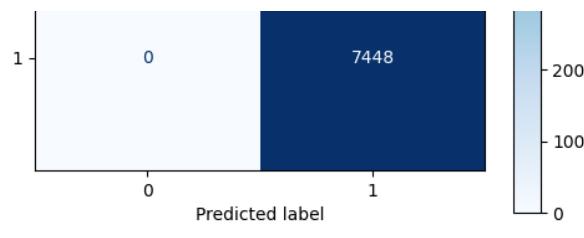
▼ Model 0 - DecisionTree

```

1 dt = DecisionTreeClassifier(random_state=8)
2 dt.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[0], dt)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)

```





✓ Model 1 - GradientBoosting

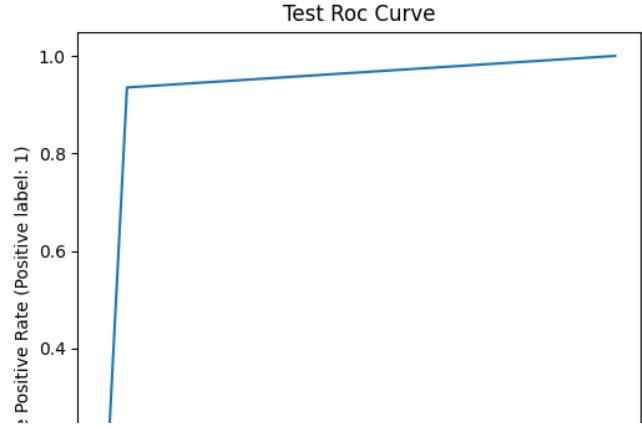
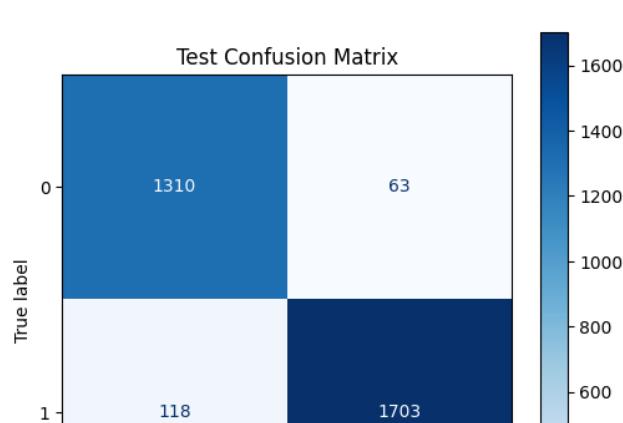
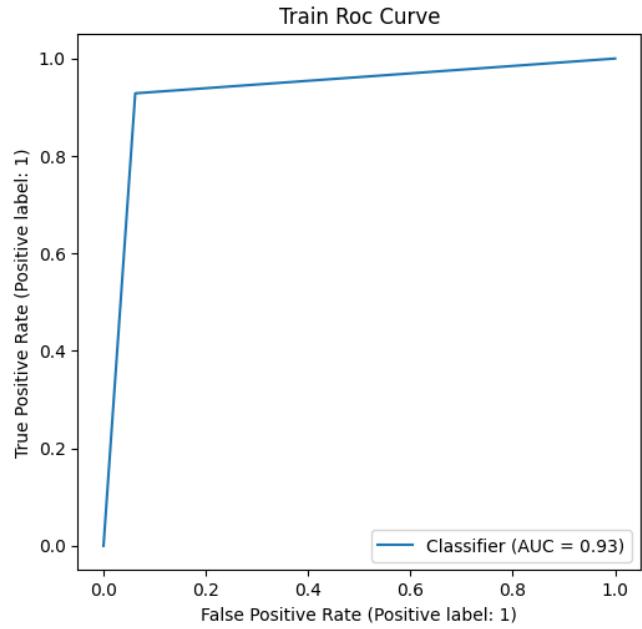
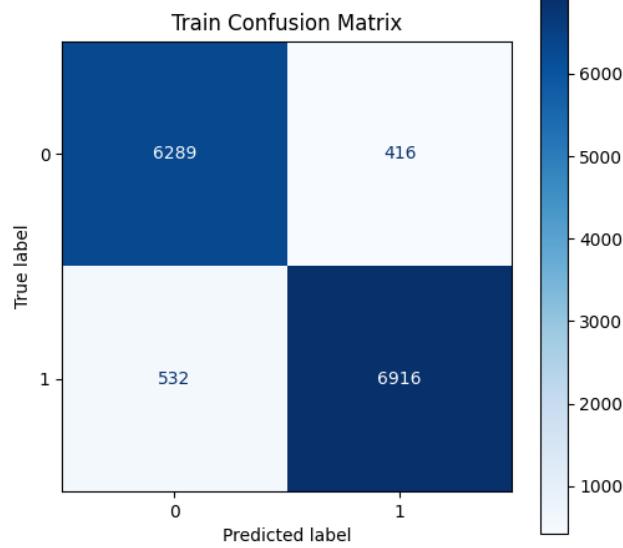
```

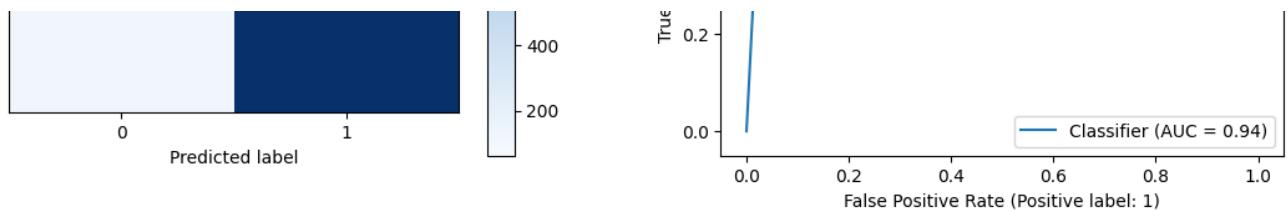
1 gb = GradientBoostingClassifier(random_state=8)
2 gb.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[1], gb)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)

```

Model Name Accuracy score Precision score Recall score F1 scor

	Model Name	Accuracy score	Precision score	Recall score	F1 scor
0	GradientBoosting	0.933018	0.933192	0.933018	0.93304
0	GradientBoosting	0.943331	0.94414	0.943331	0.94345





▼ Model 2 - HistGradientBoosting

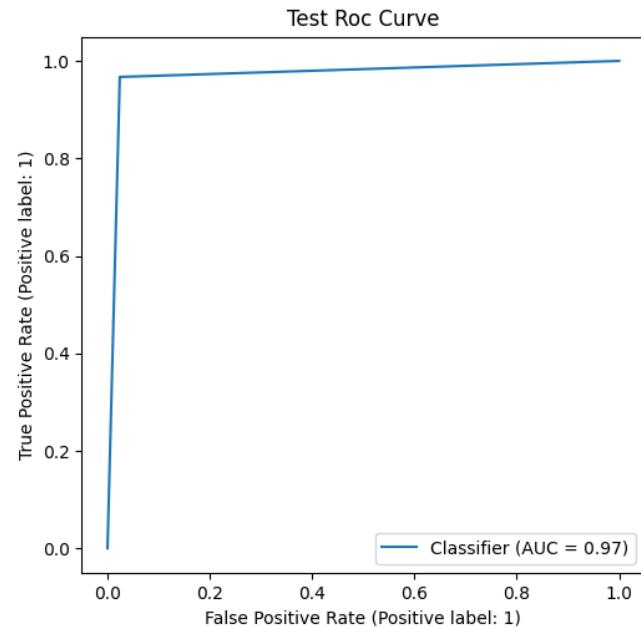
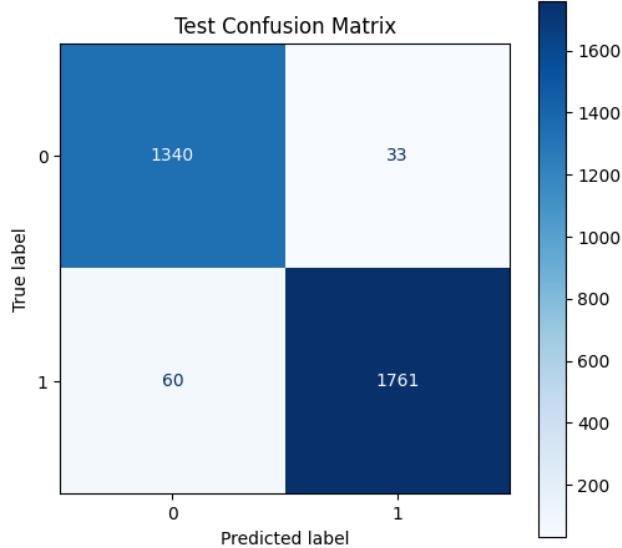
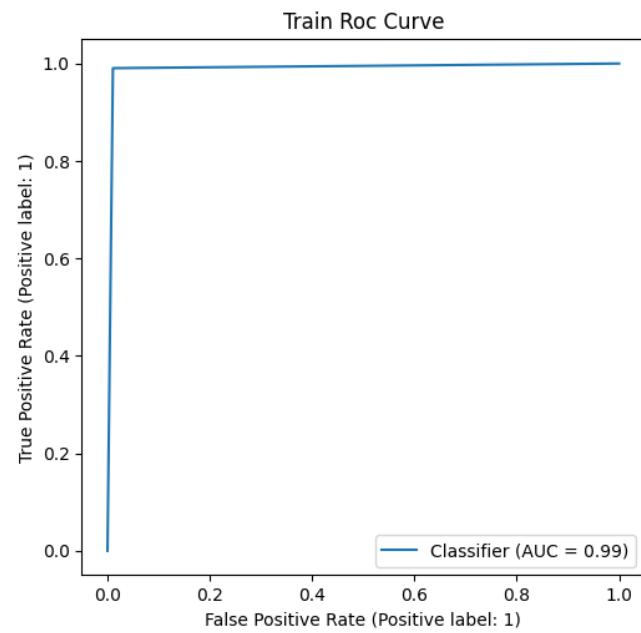
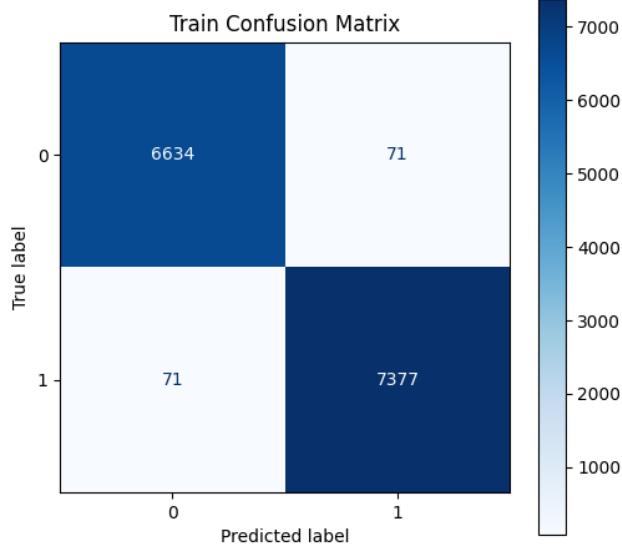
```
1 hgb = HistGradientBoostingClassifier(random_state=8)
2 hgb.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[2], hgb)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)
```

	Model Name	Accuracy score	Precision score	Recall score	\
0	HistGradientBoosting	0.989967	0.989967	0.989967	

F1 score
0 0.989967

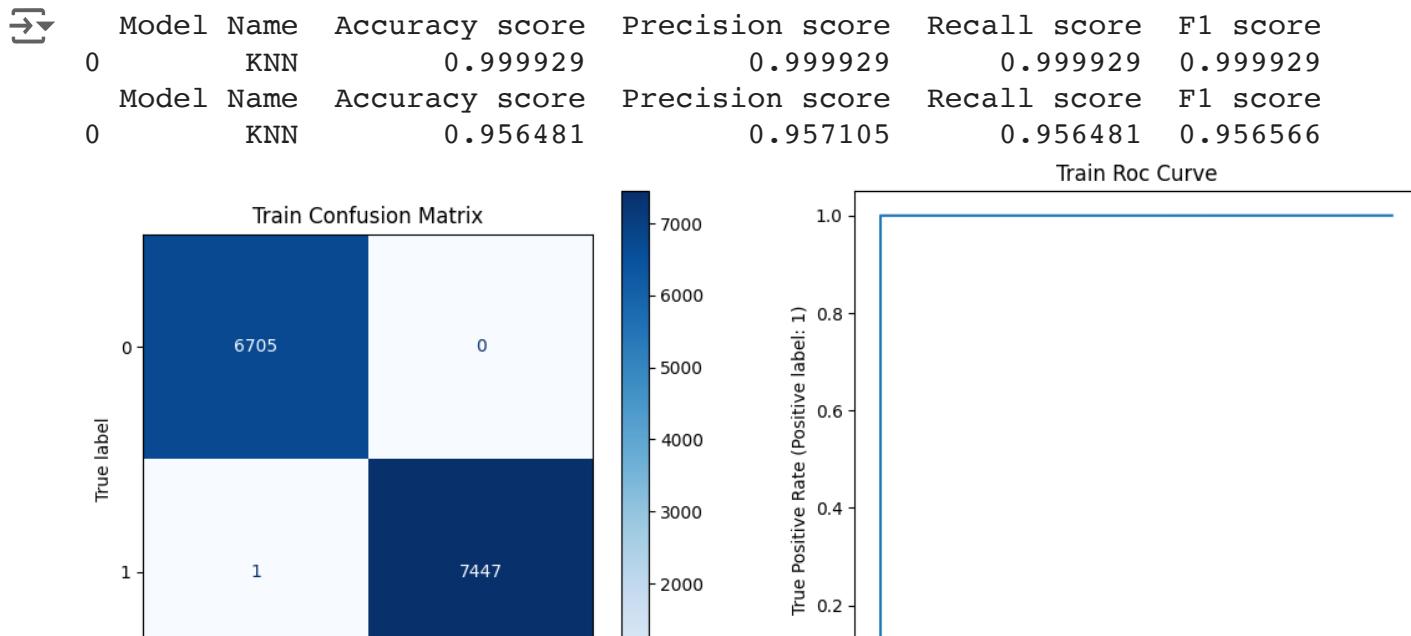
	Model Name	Accuracy score	Precision score	Recall score	\
0	HistGradientBoosting	0.970883	0.97109	0.970883	

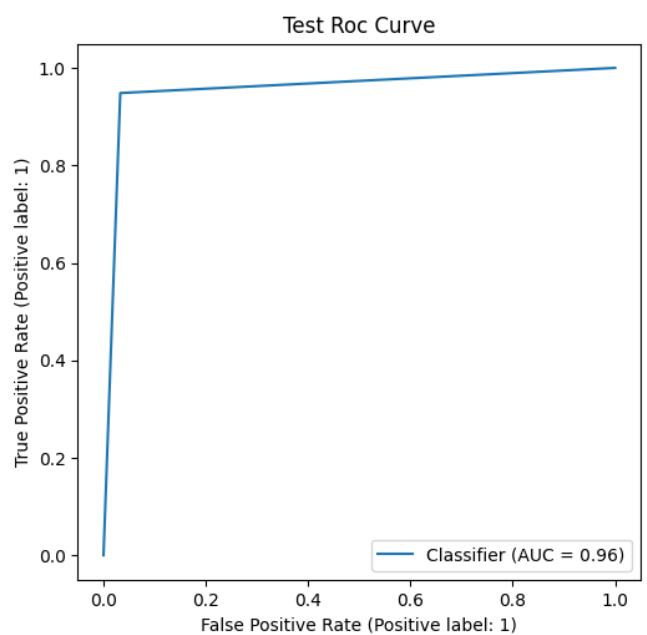
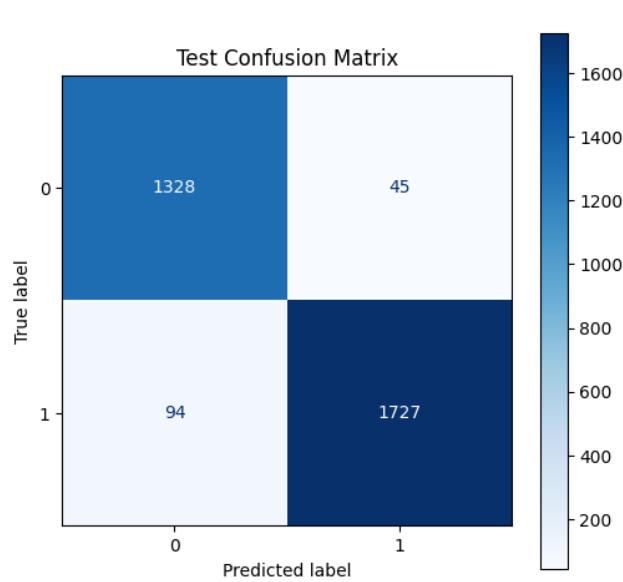
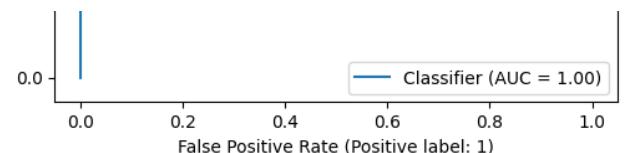
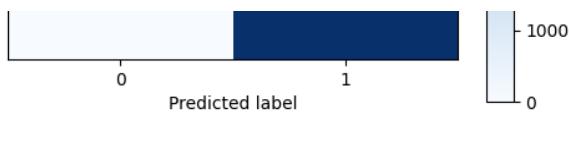
F1 score
0 0.970916



▼ Model 3 - KNN

```
1 knn = KNeighborsClassifier(n_neighbors=2)
2 knn.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[3], knn)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)
```





Model 4 - RandomForest

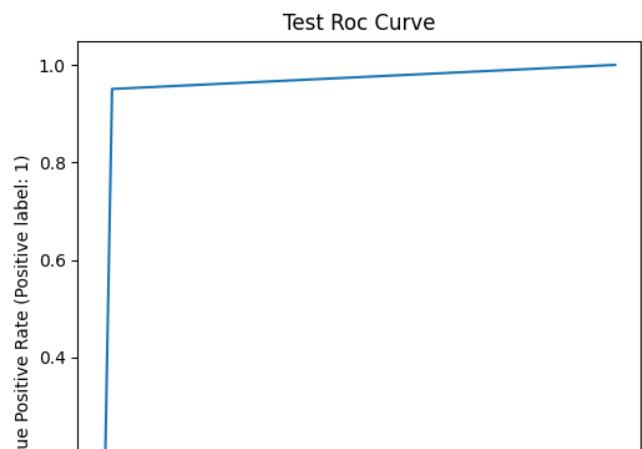
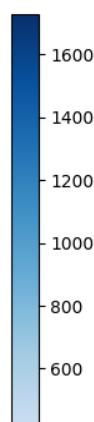
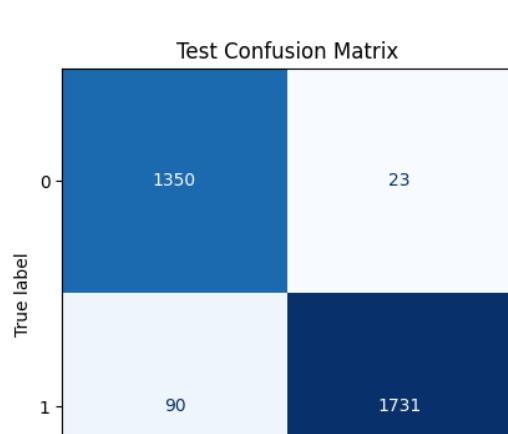
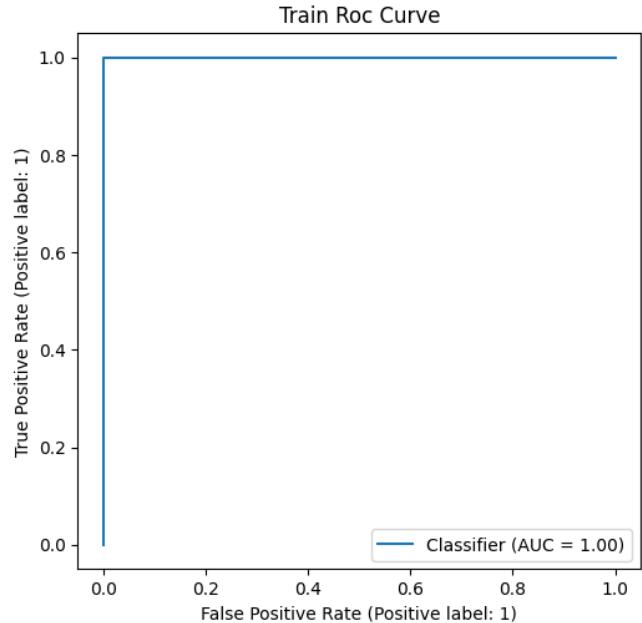
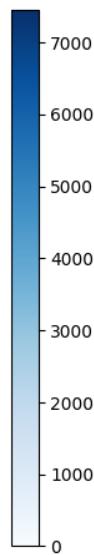
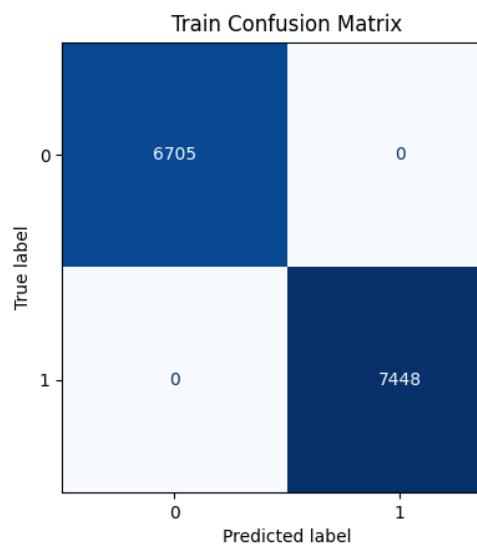
```

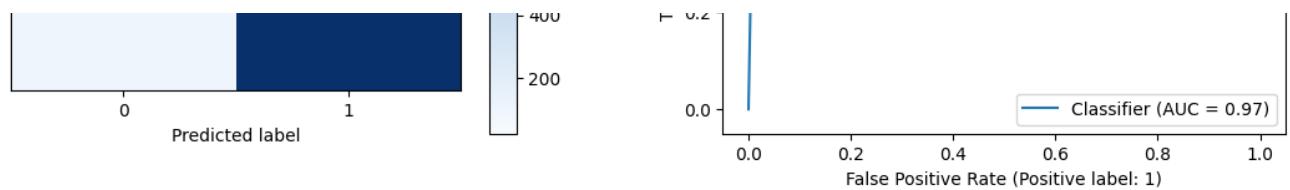
1 rf = RandomForestClassifier(random_state=8)
2 rf.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[4], rf)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)

```

Model Name Accuracy score Precision score Recall score F1 score

0	RandomForest	1.0	1.0	1.0	1.0
0	RandomForest	0.964621	0.965657	0.964621	0.964711

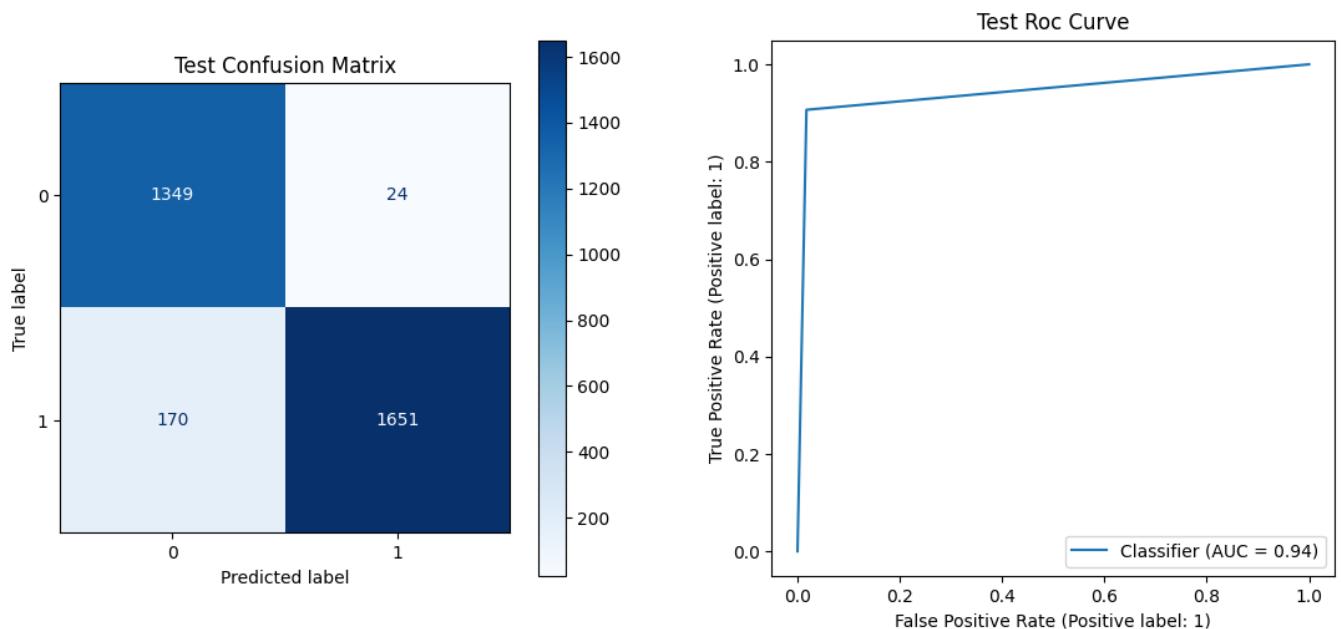
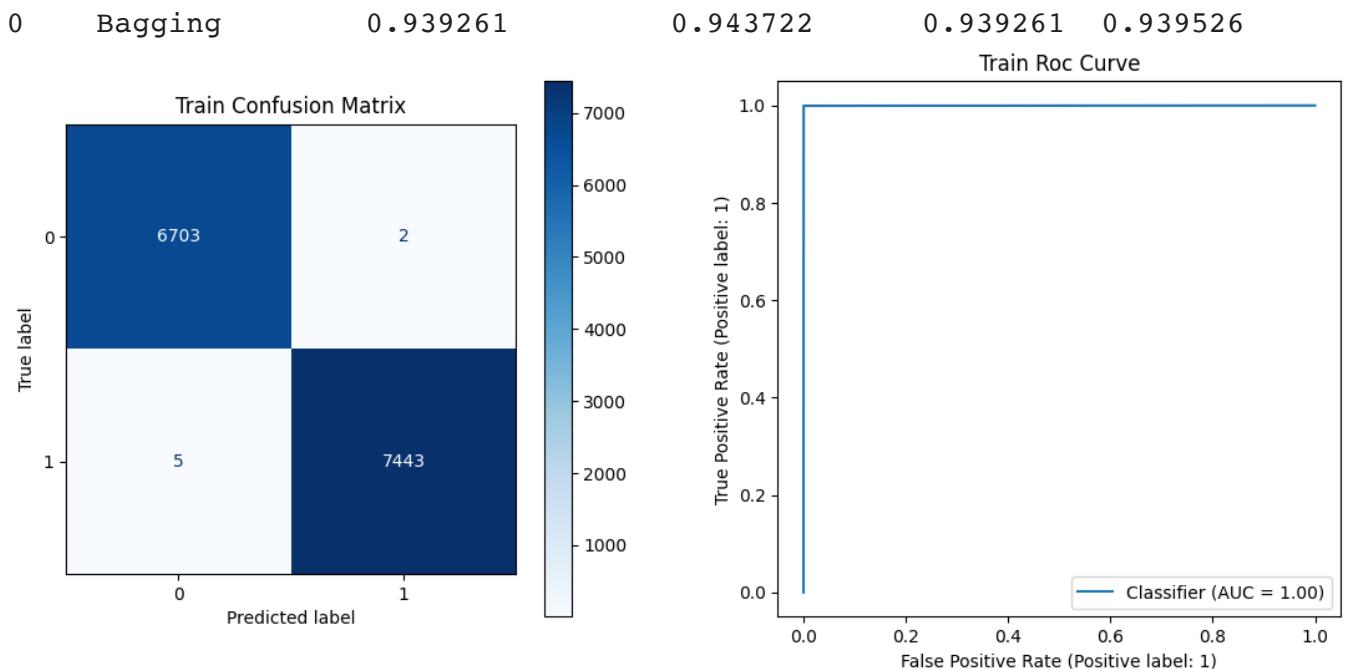




▼ Model 5 - Bagging

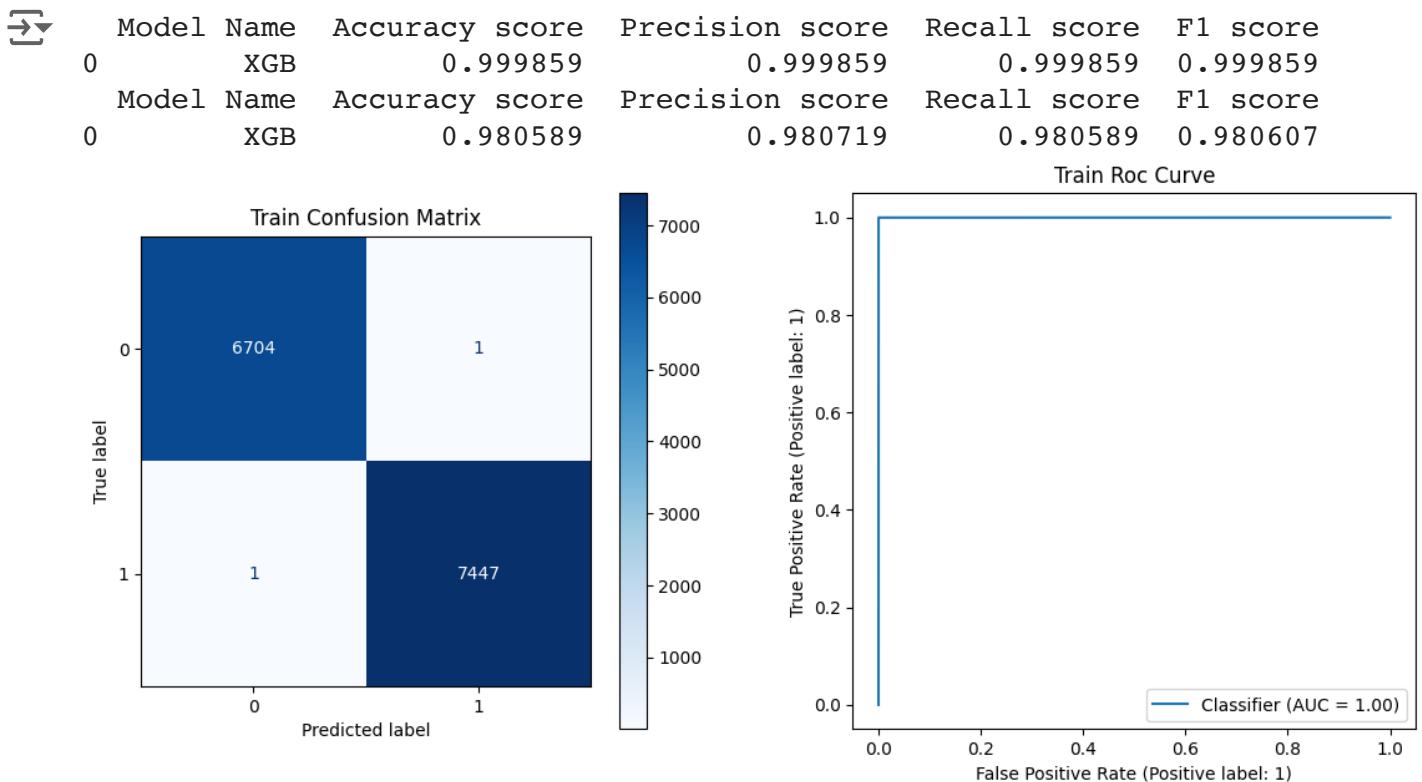
```
1 bg = BaggingClassifier(random_state=8)
2 bg.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[5], bg)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)
```

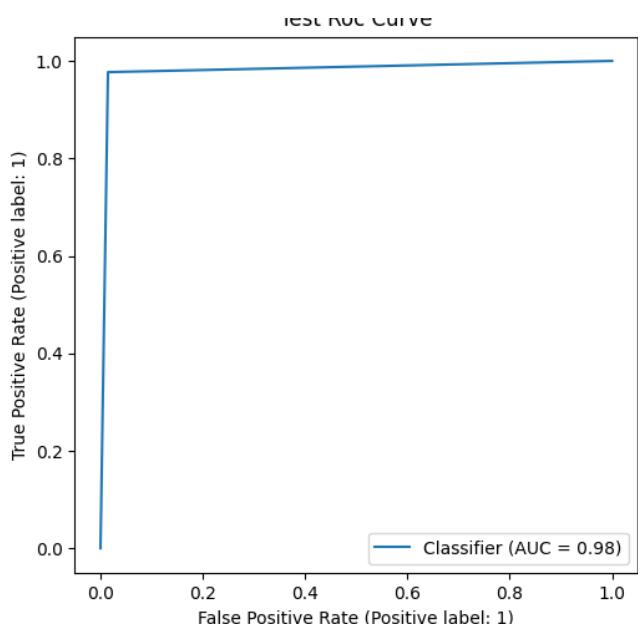
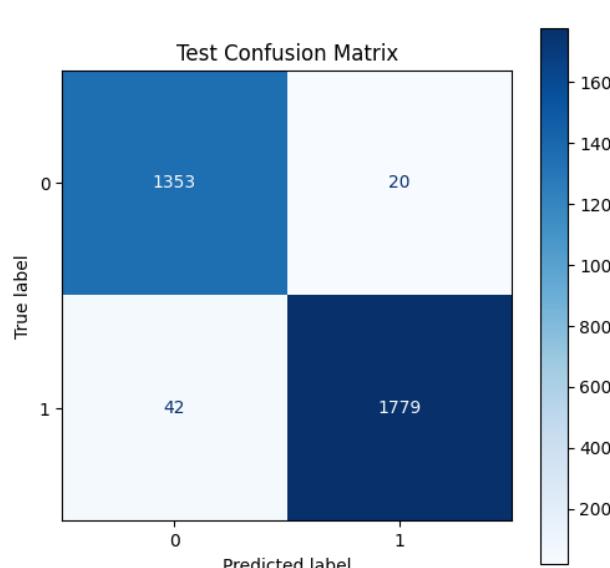
	Model Name	Accuracy score	Precision score	Recall score	F1 score
0	Bagging	0.999505	0.999506	0.999505	0.999505
	Model Name	Accuracy score	Precision score	Recall score	F1 score



✓ Model 6 - XGB

```
1 xgb = XGBClassifier(random_state=8)
2 xgb.fit(X_train_final, y_train_final)
3
4 results = pre_result_calculator(best_models[6], xgb)
5 train, test = results[0], results[1]
6 train_score, test_score = results[2], results[3]
7
8 print(train_score)
9 print(test_score)
10 plot(train, test)
```



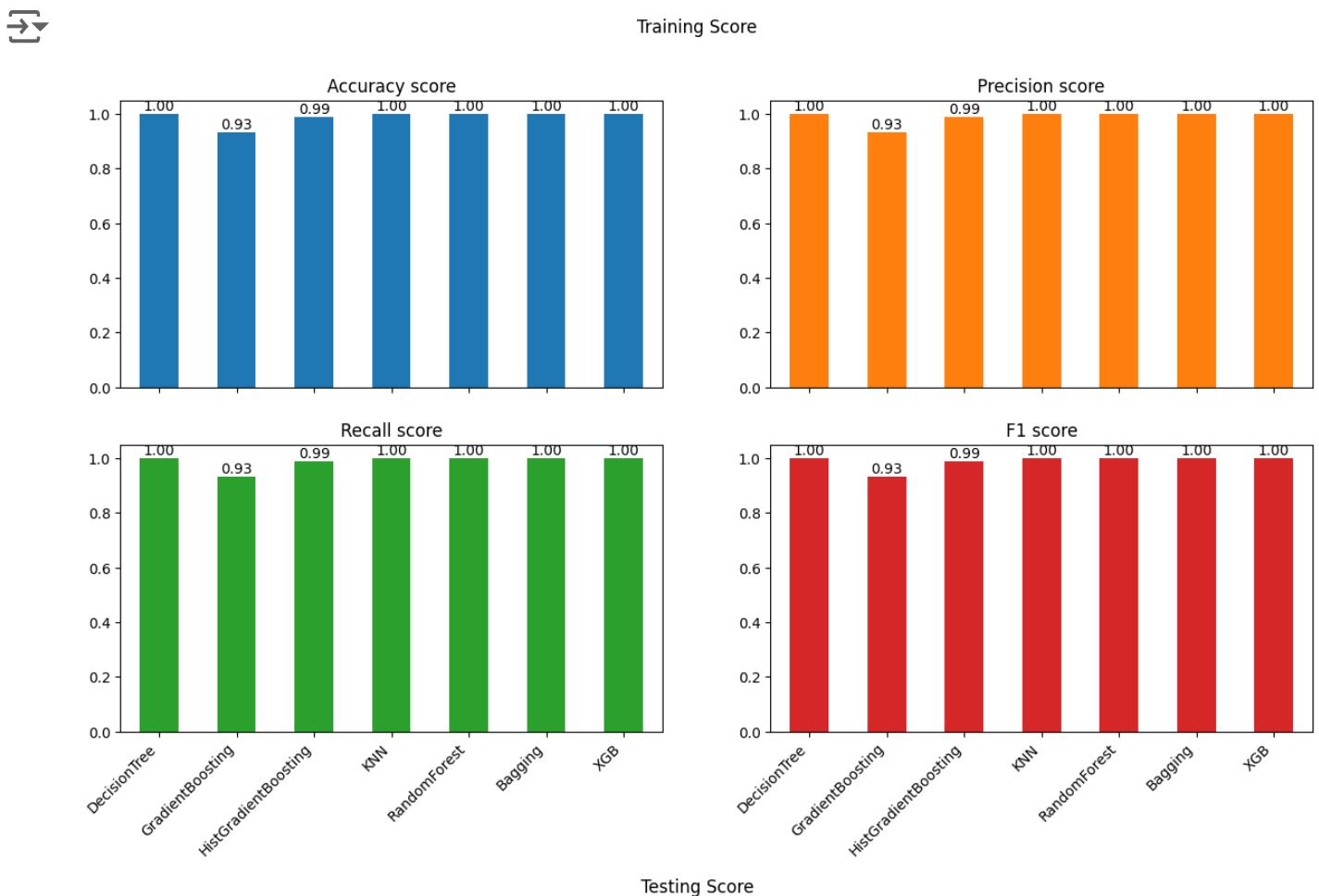


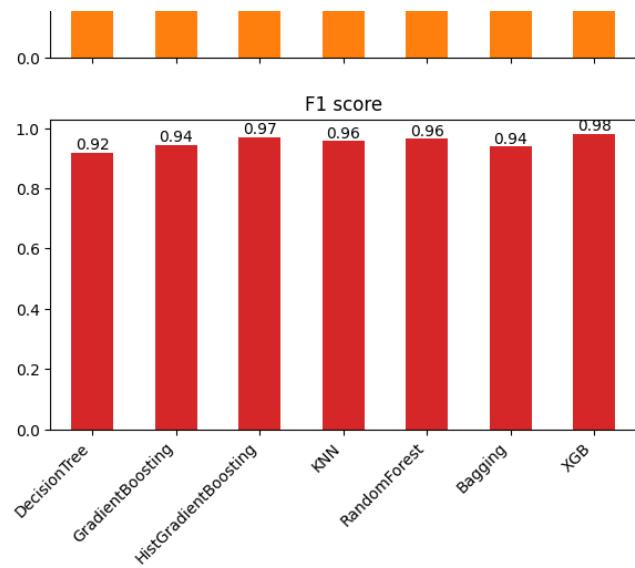
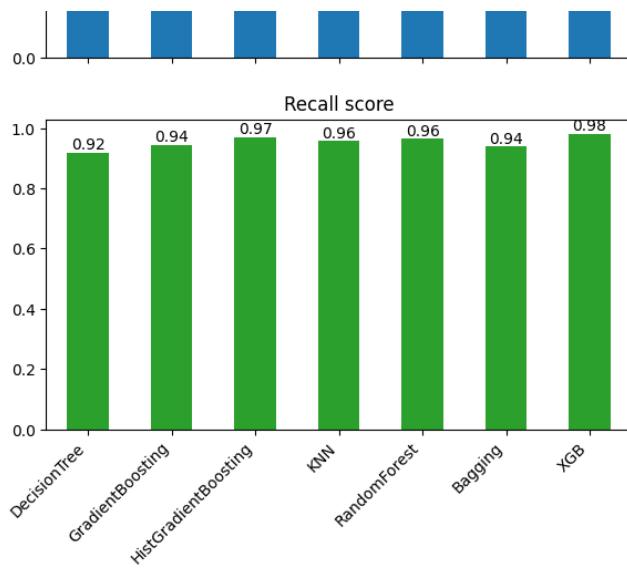
✓ Results

✓ Graph

```
1 # @title Graph
```

```
2
3 ax = pd.DataFrame(train_results).plot(kind='bar', subplots=True, la
4 ax = ax.flatten()
5 for axes in ax:
6     axes.get_legend().remove()
7     axes.bar_label(axes.containers[0], fmt='%.2f')
8     axes.set_xticklabels(best_models, rotation=45, ha='right')
9
10 ax = pd.DataFrame(test_results).plot(kind='bar', subplots=True, lay
11 ax = ax.flatten()
12 for axes in ax:
13     axes.get_legend().remove()
14     axes.bar_label(axes.containers[0], fmt='%.2f')
15     axes.set_xticklabels(best_models, rotation=45, ha='right');
```





✓ Tuning the model

```
1 train_tuned_results, test_tuned_results = [], []
2
3 def pre_tuned_result_calculator(name, model):
4     train_pred = model.predict(X_train_final)
5     test_pred = model.predict(X_test_final)
6
7     train_score = calculate_results(name, y_train_final, train_pred)
8     test_score = calculate_results(name, y_test_final, test_pred)
9
10    train_tuned_results.append(train_score)
11    test_tuned_results.append(test_score)
12
13    return (y_train_final, train_pred), \
14          (y_test_final, test_pred), \
15          pd.DataFrame(train_score, index=[0]), \
16          pd.DataFrame(test_score, index=[0])
```

✓ GradientBoosting Tuning

```
1 gb_tuned = GradientBoostingClassifier(random_state=8)
2
3 weights = np.linspace(0.0, 0.50, 100)
4
5 params = {
6     # 'loss': ['log_loss', 'exponential'],
7     # 'criterion': ['friedman_mse', 'squared_error'],
8     # 'learning_rate': [0.1, 0.01, 0.001],
9     'max_depth': [3, 4, 5],
10    'max_features': [None, 'auto', 'sqrt'],
11    'n_estimators': [100, 120, 140, 160],
12 }
13
14 grid = GridSearchCV(gb_tuned, params, scoring='f1', cv=StratifiedKFold(5))
15 grid.fit(X_train_final, y_train_final)
16 gb_tuned = grid.best_estimator_
17 gb_tuned
```

```
GradientBoostingClassifier
GradientBoostingClassifier(max_depth=5, max_features='sqrt', n_estimators=100,
                           random_state=8)
```

```
1 GradientBoostingClassifier().get_params()
```

```
GradientBoostingClassifier()
{'ccp_alpha': 0.0,
 'criterion': 'friedman_mse',
 'init': None,
 'learning_rate': 0.1,
 'loss': 'log_loss',
 'max_depth': 3,
 'max_features': None,
 'max_leaf_nodes': None,
 'min_impurity_decrease': 0.0,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'min_weight_fraction_leaf': 0.0,
 'n_estimators': 100,
 'n_iter_no_change': None,
 'random_state': None,
 'subsample': 1.0,
 'tol': 0.0001,
 'validation_fraction': 0.1,
 'verbose': 0,
 'warm_start': False}
```

```
1 print(GradientBoostingClassifier.__doc__)
```

```
-> ...
```

Gradient Boosting for classification.

This algorithm builds an additive model in a forward stage-wise fashion allows for the optimization of arbitrary differentiable loss functions. each stage ``n_classes_`` regression trees are fit on the negative grad of the loss function, e.g. binary or multiclass log loss. Binary classification is a special case where only a single regression tree is induced.

:class:`sklearn.ensemble.HistGradientBoostingClassifier` is a much fast variant of this algorithm for intermediate datasets (`n_samples >= 10_0

Read more in the :ref:`User Guide <gradient_boosting>`.

Parameters

loss : {'log_loss', 'deviance', 'exponential'}, default='log_loss'
The loss function to be optimized. 'log_loss' refers to binomial and multinomial deviance, the same as used in logistic regression.
It is a good choice for classification with probabilistic outputs.
For loss 'exponential', gradient boosting recovers the AdaBoost alg
.. deprecated:: 1.1
The loss 'deviance' was deprecated in v1.1 and will be removed version 1.3. Use `loss='log_loss'` which is equivalent.

learning_rate : float, default=0.1

Learning rate shrinks the contribution of each tree by `learning_rate`. There is a trade-off between learning_rate and n_estimators.
Values must be in the range `[0.0, inf)`.

n_estimators : int, default=100

The number of boosting stages to perform. Gradient boosting is fairly robust to over-fitting so a large number usually results in better performance.
Values must be in the range `[1, inf)`.

subsample : float, default=1.0

The fraction of samples to be used for fitting the individual base learners. If smaller than 1.0 this results in Stochastic Gradient Boosting. `subsample` interacts with the parameter `n_estimators`. Choosing `subsample < 1.0` leads to a reduction of variance and an increase in bias.
Values must be in the range `(0.0, 1.0)`.

criterion : {'friedman_mse', 'squared_error'}, default='friedman_mse'

The function to measure the quality of a split. Supported criteria 'friedman_mse' for the mean squared error with improvement score by Friedman, 'squared_error' for mean squared error. The default value 'friedman_mse' is generally the best as it can provide a better approximation in some cases.

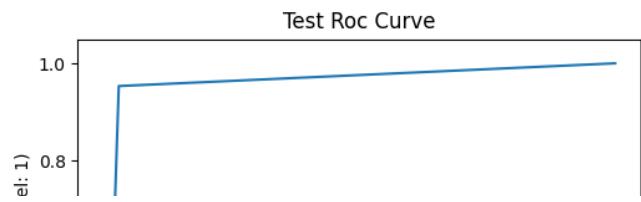
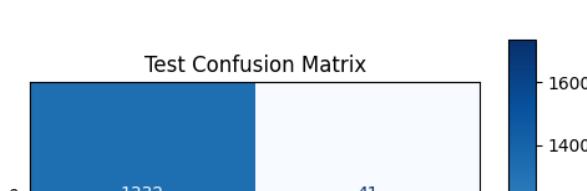
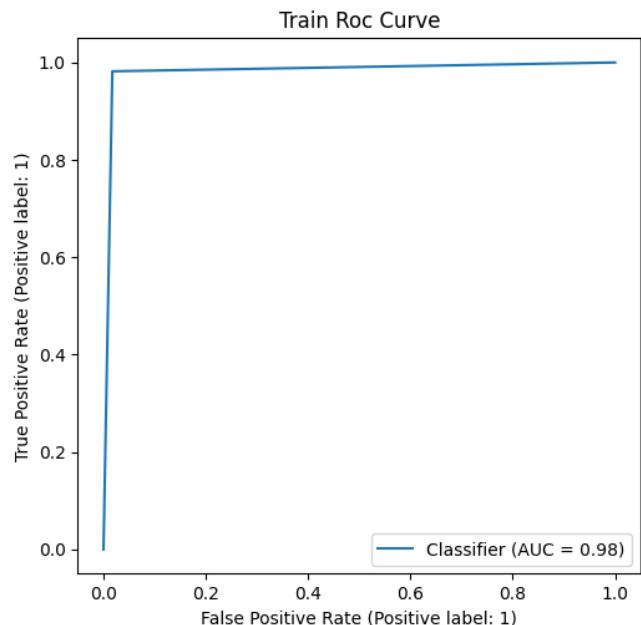
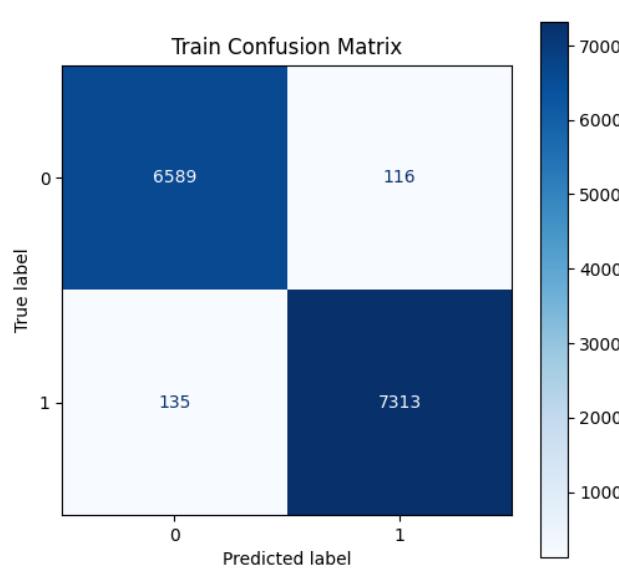
```
.. versionadded:: 0.18
```

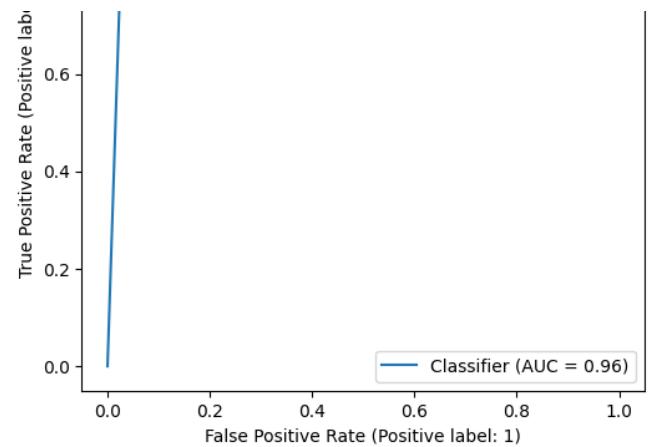
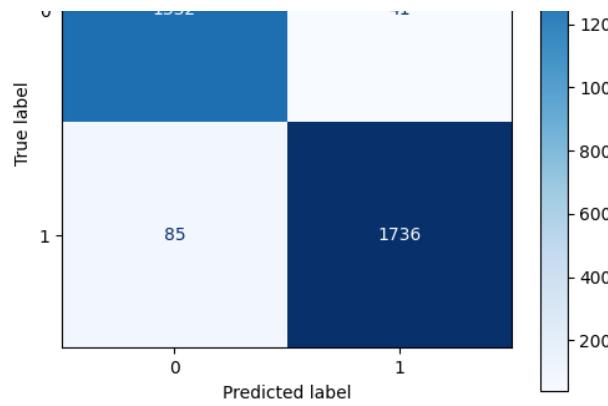
```
min_samples_split : int or float, default=2
The minimum number of samples required to split an internal node:
- If int, values must be in the range `[2, inf)`.
- If float, values must be in the range `(0.0, 1.0)` and `min_sampl
```

```
1 gb_tuned.fit(X_train_final, y_train_final)
2 results = pre_tuned_result_calculator(best_models[1] + 'Tuned', gb_t
3 train, val = results[0], results[1]
4 train_score, val_score = results[2], results[3]
5
6 print(train_score)
7 print(val_score)
8 plot(train, val)
```

	Model Name	Accuracy score	Precision score	Recall score	\
0	GradientBoostingTuned	0.982265	0.982271	0.982265	
	F1 score				
0	0.982266				

	Model Name	Accuracy score	Precision score	Recall score	\
0	GradientBoostingTuned	0.960551	0.96106	0.960551	
	F1 score				
0	0.960621				





✓ Final Results

```
1 pd.DataFrame(test_results).sort_values('Recall score', ascending=F
```

	Model Name	Accuracy score	Precision score	Recall score	F1 score
6	XGB	0.980589	0.980719	0.980589	0.980607
2	HistGradientBoosting	0.970883	0.971090	0.970883	0.970916
4	RandomForest	0.964621	0.965657	0.964621	0.964711
3	KNN	0.956481	0.957105	0.956481	0.956566
1	GradientBoosting	0.943331	0.944140	0.943331	0.943453
5	Bagging	0.939261	0.943722	0.939261	0.939526
0	DecisionTree	0.918910	0.924086	0.918910	0.919277

Data Insights (Best Model)

It's found the best model so far is XGB . We got 100% results on training dataset on all the metrics (Accuracy, Precision, Recall, F1-score) and around 98% on testing dataset on all the metrics.

✓ Saving the Best Model

```
1 from joblib import load, dump
```

```
1 path = '/content/drive/MyDrive/Colab Notebooks/Jain/best_model.job
2
3 best_model = dump(xgb, path)
```

✓ Loading the Best Model

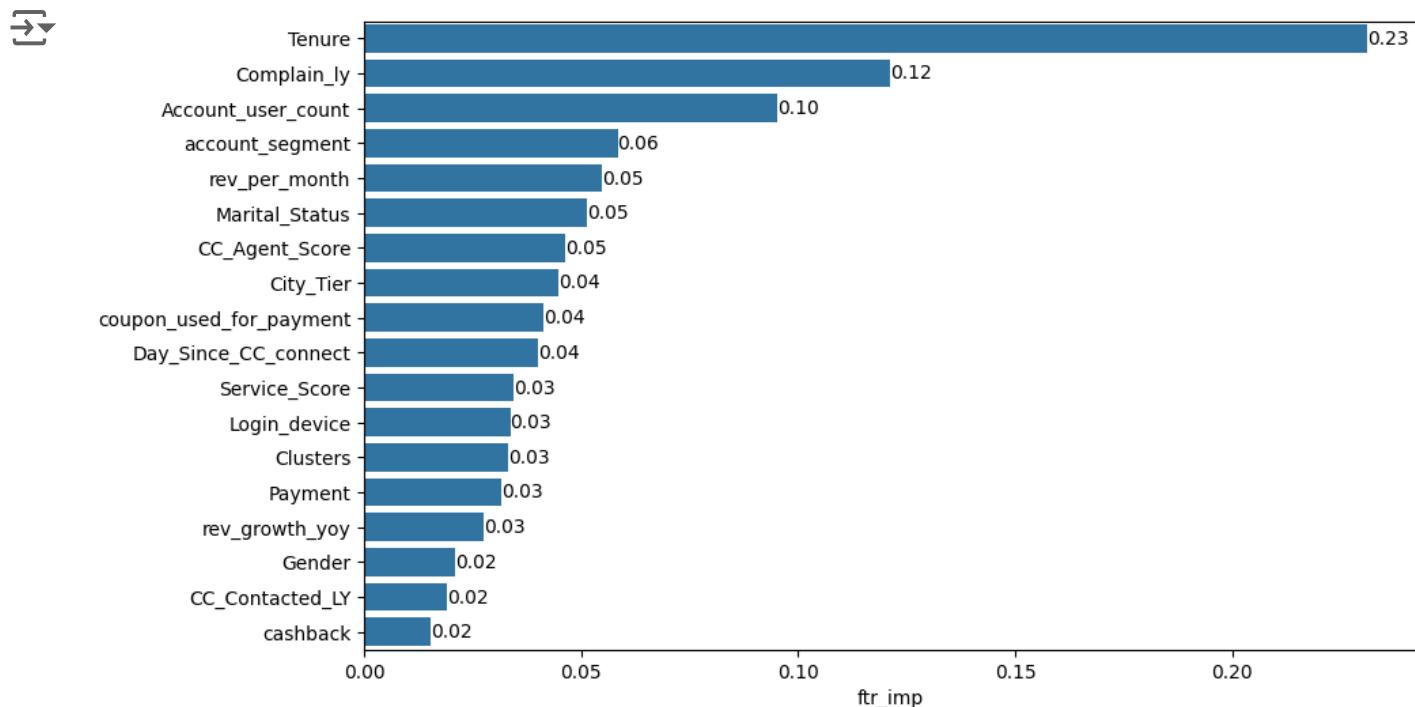
```
1 saved_model = load(path)
```

```
1 saved_model.score(X_test_final, y_test_final)
```

↳ 0.9805886036318097

▼ Feature Importance

```
1 ftr_data = pd.DataFrame({'ftr_names': saved_model.feature_names_in_
2                               'ftr_imp': saved_model.feature_importances_
3
4 ftr_data = ftr_data.sort_values('ftr_imp', ascending=False).reset_
5
6 fig = plt.figure(figsize=(10, 6))
7 ax = sns.barplot(ftr_data.ftr_imp, orient='h')
8 ax.set_yticklabels(ftr_data.ftr_names)
9 ax.bar_label(ax.containers[0], fmt='%.2f');
```



❖ Data Insights

It is found that the top most important features of the dataset which is contributing most to the model are Tenure , Complain_ly and Account_user_count .
