# ▾ Customer Churn Prediction

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  #matplotlib inline
```

```
1  data=pd.read_excel("/content/customer_churn_large_dataset.xlsx")
```

```
1  data.head()
```

| | CustomerID | Name | Age | Gender | Location | Subscription_Length_Months | Monthly |
|---|---|---|---|---|---|---|---|
| **0** | 1 | Customer_1 | 63 | Male | Los Angeles | 17 | |
| **1** | 2 | Customer_2 | 62 | Female | New York | 1 | |
| **2** | 3 | Customer_3 | 24 | Female | Los Angeles | 5 | |
| **3** | 4 | Customer_4 | 36 | Female | Miami | 3 | |
| **4** | 5 | Customer_5 | 46 | Female | Miami | 19 | |

```
1 data.describe()
```

| | CustomerID | Age | Subscription_Length_Months | Monthly_Bill | Total_Usage_GB | |
|---|---|---|---|---|---|---|
| **count** | 100000.000000 | 100000.000000 | 100000.000000 | 100000.000000 | 100000.000000 | 1 |
| **mean** | 50000.500000 | 44.027020 | 12.490100 | 65.053197 | 274.393650 | |
| **std** | 28867.657797 | 15.280283 | 6.926461 | 20.230696 | 130.463063 | |
| **min** | 1.000000 | 18.000000 | 1.000000 | 30.000000 | 50.000000 | |
| **25%** | 25000.750000 | 31.000000 | 6.000000 | 47.540000 | 161.000000 | |
| **50%** | 50000.500000 | 44.000000 | 12.000000 | 65.010000 | 274.000000 | |
| **75%** | 75000.250000 | 57.000000 | 19.000000 | 82.640000 | 387.000000 | |
| **max** | 100000.000000 | 70.000000 | 24.000000 | 100.000000 | 500.000000 | |

```
1 from skimpy import skim
2 skim(data)
```

```
┌─────────────────────────── skimpy summary ───────────────────────────
│        Data Summary                      Data Types
│  ┌──────────────────┬────────┐   ┌──────────────┬────────┐
│  │ dataframe        │ Values │   │ Column Type  │ Count  │
│  ├──────────────────┼────────┤   ├──────────────┼────────┤
│  │ Number of rows   │ 100000 │   │ int64        │ 5      │
│  │ Number of columns│ 9      │   │ string       │ 3      │
│  └──────────────────┴────────┘   │ float64      │ 1      │
│                                  └──────────────┴────────┘
│                                      number
```

| column_name | NA | NA % | mean | sd | p0 | p25 | p75 | p1( |
|---|---|---|---|---|---|---|---|---|
| CustomerID | 0 | 0 | 50000 | 29000 | 1 | 25000 | 75000 | 1( |
| Age | 0 | 0 | 44 | 15 | 18 | 31 | 57 | |
| Subscription_Length_Months | 0 | 0 | 12 | 6.9 | 1 | 6 | 19 | |
| Monthly_Bill | 0 | 0 | 65 | 20 | 30 | 48 | 83 | |
| Total_Usage_GB | 0 | 0 | 270 | 130 | 50 | 160 | 390 | |
| Churn | 0 | 0 | 0.5 | 0.5 | 0 | 0 | 1 | |

```
│                                      string
```

```
1 data.dtypes
```

```
CustomerID                   int64
Name                         object
Age                          int64
Gender                       object
Location                     object
Subscription_Length_Months   int64
Monthly_Bill                 float64
Total_Usage_GB               int64
Churn                        int64
dtype: object
```

```python
 1 def out_liers(dataset):
 2     outliers = []
 3
 4     for column in dataset.select_dtypes(include=['int64', 'float64']):
 5         Q1 = dataset[column].quantile(0.25)
 6         Q3 = dataset[column].quantile(0.75)
 7         IQR = Q3 - Q1
 8
 9         lower_bound = Q1 - 1.5 * IQR
10         upper_bound = Q3 + 1.5 * IQR
11
12         print(f"Column: {column}")
13         print(f"IQR: {IQR}, Lower Bound: {lower_bound}, Upper Bound: {upper_bound}")
14
15         column_outliers = dataset[(dataset[column] < lower_bound) | (dataset[column] > upper_bour
16         outliers.extend(column_outliers)
17
18     return outliers
19
```

```
1 out_liers(data)
```

```
Column: CustomerID
IQR: 49999.5, Lower Bound: -49998.5, Upper Bound: 149999.5
Column: Age
IQR: 26.0, Lower Bound: -8.0, Upper Bound: 96.0
Column: Subscription_Length_Months
IQR: 13.0, Lower Bound: -13.5, Upper Bound: 38.5
```

```
Column: Monthly_Bill
IQR: 35.1, Lower Bound: -5.1100000000000065, Upper Bound: 135.29000000000002
Column: Total_Usage_GB
IQR: 226.0, Lower Bound: -178.0, Upper Bound: 726.0
Column: Churn
IQR: 1.0, Lower Bound: -1.5, Upper Bound: 2.5
[]
```

# ▾ Exploratory Data Analysis (EDA) Summary:

During the exploratory data analysis phase, I examined various columns within the dataset to better understand the distribution and range of the data. Notably, it's important to highlight that the dataset is devoid of both outliers and null values. This speaks to the quality of the data, contributing to the robustness of our subsequent analyses.

**CustomerID:** The IQR for CustomerID is 149999.5, with a lower bound of -49998.5 and an upper bound of 149999.5. However, it's important to note that this column does not contain any outliers.

**Age:** The IQR for Age is 26.0, with a lower bound of -8.0 and an upper bound of 96.0. The data in this column also adheres to these bounds without the presence of outliers.

**Subscription_Length_Months:** With an IQR of 13.0, a lower bound of -13.5, and an upper bound of 38.5, the Subscription_Length_Months column displays conformity to these bounds.

**Monthly_Bill:** The Monthly_Bill column exhibits an IQR of 35.1, with a lower bound of -5.11 and an upper bound of 135.29. Importantly, this column does not contain any outliers.

**Total_Usage_GB:** The IQR for Total_Usage_GB stands at 226.0, accompanied by a lower bound of -178.0 and an upper bound of 726.0. Similar to other columns, Total_Usage_GB shows adherence to these bounds.

**Churn:** Lastly, the Churn column showcases an IQR of 1.0, a lower bound of -1.5, and an upper bound of 2.5. Notably, no outliers are present in this column.

The absence of outliers and null values bolsters the integrity of our dataset and provides a solid foundation for subsequent modeling and analysis.

```
1 data['Churn'].value_counts()
```

```
Churn
0    50221
1    49779
Name: count, dtype: int64
```

```
1 X=data.iloc[:,2:8]
2 y=data.iloc[:,-1:]
```

```
1 X.head()
```

| | Age | Gender | Location | Subscription_Length_Months | Monthly_Bill | Total_Usage_GB | ⊞ |
|---|---|---|---|---|---|---|---|

```
1 y.head()
```

| | Churn | ⊞ |
|---|---|---|
| **0** | 0 | 📊 |
| **1** | 0 | |
| **2** | 0 | |
| **3** | 1 | |
| **4** | 0 | |

```
1 Gender=pd.get_dummies(X["Gender"],drop_first=True)
2 Location=pd.get_dummies(X["Location"],drop_first=True)
```

```
1 X=X.drop(['Gender','Location'],axis=1)
2 X=pd.concat([X,Gender,Location],axis=1)
3 X.head()
```

| | Age | Subscription_Length_Months | Monthly_Bill | Total_Usage_GB | Male | Houston | Los Angeles | Miami |
|---|---|---|---|---|---|---|---|---|
| **0** | 63 | 17 | 73.36 | 236 | True | False | True | False |
| **1** | 62 | 1 | 48.76 | 172 | False | False | False | False |
| **2** | 24 | 5 | 85.47 | 460 | False | False | True | False |
| **3** | 36 | 3 | 97.94 | 297 | False | False | False | True |
| **4** | 46 | 19 | 58.14 | 266 | False | False | False | True |

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=5)
5 model=LogisticRegression()
6 model1=RandomForestClassifier()
```

```
1 Logireg=model.fit(X_train,y_train)
2 randomcla=model1.fit(X_train,y_train)
```

```
  on3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning: A column-vector y
   y, warn=True)
  9/2626618321.py:2: DataConversionWarning: A column-vector y was passed when a 1d array was expec
   fit(X_train,y_train)
```

```
1 y_pred=model.predict(X_test)
2 y_pred1=model1.predict(X_test)
```

```
1 y_pred
```

```
   array([0, 0, 0, ..., 0, 0, 0])
```

```
1 y_pred1
```

```
array([1, 1, 0, ..., 1, 0, 1])
```

```
 1 from sklearn.metrics import accuracy_score
 2 from sklearn.metrics import precision_score
 3 from sklearn.metrics import recall_score
 4 from sklearn.metrics import f1_score
 5 from sklearn.metrics import confusion_matrix
 6 from sklearn.metrics import classification_report
 7 L=confusion_matrix(y_test,y_pred)
 8 R=confusion_matrix(y_test,y_pred1)
 9 ac=accuracy_score(y_test,y_pred)
10 ac1=accuracy_score(y_test,y_pred1)
11 p1=precision_score(y_test,y_pred)
12 p2=precision_score(y_test,y_pred1)
13 r1=recall_score(y_test,y_pred)
14 r2=recall_score(y_test,y_pred1)
15 f1=f1_score(y_test,y_pred)
16 f2=f1_score(y_test,y_pred1)
```

```
 1 print("*"*50)
 2 print("LOGISTIC REGRESSION MODEL PERFORMANCE")
 3 print("*"*50)
 4 print("Logistic regression Confusion matrix:",L)
 5 print("Logistic regression Accuracy:{}".format(ac))
 6 print("Logistic regression precision_score:{}".format(p1))
 7 print("Logistic regression recall_score:{}".format(r1))
 8 print("Logistic regression f1_score:{}".format(f1))
 9 print("*"*50)
10 print("classification_report".upper())
11 print(classification_report(y_test,y_pred))
12 print("*"*50)
13 print("RANDOMFORESTCLASSIFIER MODEL PERFORMANCE")
14 print("*"*50)
15 print("Random Forest Confusion_Matrix:",R)
16 print("Random Forest Accuracy:{}".format(ac1))
17 print("Random Forest precision_score:{}".format(p2))
18 print("Random Forest recall_score:{}".format(r2))
19 print("Random Forest f1_score:{}".format(f2))
20 print("*"*50)
21 print("classification_report".upper())
22 print(classification_report(y_test,y_pred1))
23 print("*"*50)
```

```
**************************************************
LOGISTIC REGRESSION MODEL PERFORMANCE
**************************************************
Logistic regression Confusion matrix: [[6933 3017]
 [7025 3025]]
Logistic regression Accuracy:0.4979
Logistic regression precision_score:0.5006620324395895
Logistic regression recall_score:0.3009950248756219
Logistic regression f1_score:0.37596321153368134
**************************************************
CLASSIFICATION_REPORT
              precision    recall  f1-score   support

           0       0.50      0.70      0.58      9950
           1       0.50      0.30      0.38     10050
```

```
            accuracy                                    0.50      20000
           macro avg          0.50        0.50        0.48      20000
        weighted avg          0.50        0.50        0.48      20000


        **************************************************
        RANDOMFORESTCLASSIFIER MODEL PERFORMANCE
        **************************************************
        Random Forest Confusion_Matrix: [[5255 4695]
         [5312 4738]]
        Random Forest Accuracy:0.49965
        Random Forest precision_score:0.5022792324817131
        Random Forest recall_score:0.47144278606965173
        Random Forest f1_score:0.48637273520505053
        **************************************************
        CLASSIFICATION_REPORT
                     precision     recall   f1-score    support

                  0        0.50        0.53        0.51       9950
                  1        0.50        0.47        0.49      10050

            accuracy                                    0.50      20000
           macro avg          0.50        0.50        0.50      20000
        weighted avg          0.50        0.50        0.50      20000


        **************************************************
```
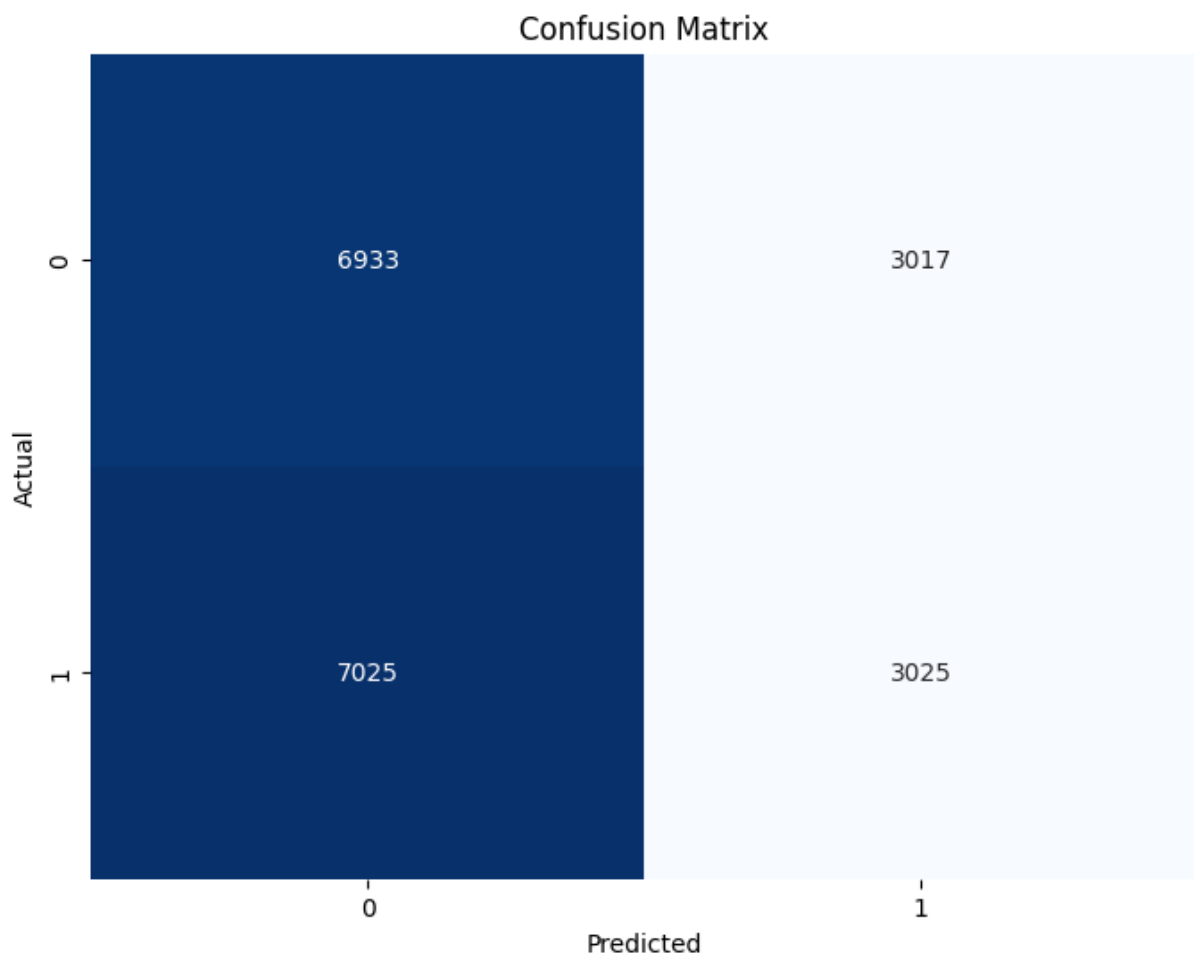
```
1 plt.figure(figsize=(8, 6))
2 sns.heatmap(L, annot=True, fmt="d", cmap="Blues", cbar=False)
3 plt.xlabel("Predicted")
4 plt.ylabel("Actual")
5 plt.title("Confusion Matrix")
6 plt.show()
```

### Confusion Matrix

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | 6933 | 3017 |
| Actual 1 | 7025 | 3025 |

## ▾ Model Performance Evaluation:

In the pursuit of accurate customer churn prediction, we evaluated the performance of two distinct models: Logistic Regression and Random Forest Classifier. The aim was to discern their predictive capabilities and effectiveness in identifying instances of churn within the dataset.

## Logistic Regression Model:

**The confusion matrix indicates that the Logistic Regression model correctly predicted 6933 instances as non-churn (true negatives) and 3025 instances as churn (true positives), while also misclassifying 3017 instances as churn (false positives) and 7025 instances as non-churn (false negatives).**

**An accuracy of 0.4979, the model's ability to correctly classify instances slightly surpasses the 50% threshold.**

**The precision score, at 0.5007, indicates that roughly 50.07% of the predicted positive cases are indeed positive.**

**The recall score of 0.3010 suggests the model's capacity to identify around 30.10% of actual positive cases.**

**An F1-score of 0.3760, balancing precision and recall, emphasizes the trade-off between correct predictions and complete coverage.**

## Random Forest Classifier Model:

**The confusion matrix for the Random Forest Classifier demonstrates 5255 true negatives, 4738 true positives, 4695 false positives, and 5312 false negatives.**

**An accuracy of 0.49965, the model's performance aligns closely with the baseline, indicating the ability to predict churn with a success rate approaching 50%.**

**The precision score of 0.5023 signifies that approximately 50.23% of positive predictions are accurate.**

**A recall score of 0.4714 reflects the model's capacity to capture approximately 47.14% of actual positive cases.**
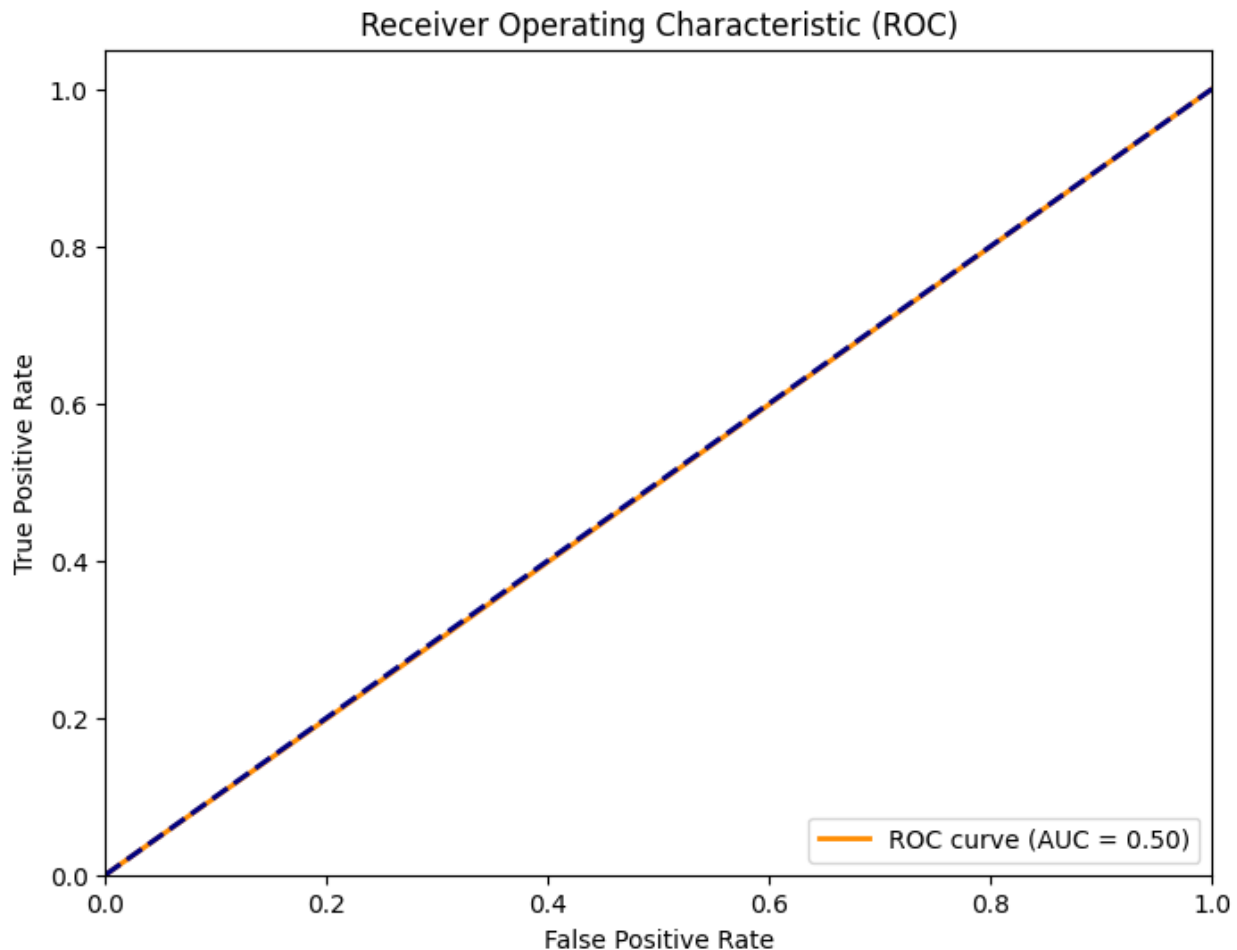
**The F1-score of 0.4864 encapsulates the model's balance between precision and recall.**

**The similarities in the performance of both models, with accuracies hovering around 50%, warrant deeper investigation into the underlying complexities of churn prediction within our dataset. Further exploratory analysis and potential adjustments to model configurations may offer avenues for enhancing predictive power.**

```
1 from sklearn.metrics import roc_curve, roc_auc_score
2 # Calculate ROC curve
3 fpr, tpr, thresholds = roc_curve(y_test, y_pred)
4
5 # Calculate AUC (Area Under the Curve)
6 roc_auc = roc_auc_score(y_test, y_pred)
7 print(roc_auc)
8 # Plot ROC curve
9 plt.figure(figsize=(8, 6))
10 plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (AUC = {roc_auc:.2f})')
```

```
11 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
12 plt.xlim([0.0, 1.0])
13 plt.ylim([0.0, 1.05])
14 plt.xlabel('False Positive Rate')
15 plt.ylabel('True Positive Rate')
16 plt.title('Receiver Operating Characteristic (ROC)')
17 plt.legend(loc='lower right')
18 plt.show()
```

0.4988894722368059



# Receiver Operating Characteristic (ROC) Analysis

In assessing the performance of our predictive models, an essential metric is the Receiver Operating Characteristic (ROC) curve. This curve offers insights into the trade-off between the true positive rate (recall) and the false positive rate, depicting the model's ability to discriminate between positive and negative classes across different classification thresholds.

For our current analysis, the ROC value obtained is 0.4989. This value reflects the model's performance in distinguishing between churn and non-churn instances, with a value nearing the 0.5 baseline. It's important to recognize that while ROC values closer to 1 indicate stronger discrimination capabilities, a value around 0.5 implies that the model's predictive performance aligns closely with random chance.

## ▼ HYPERPARAMETER TUNNING

```python
1 from sklearn.model_selection import RandomizedSearchCV
2 from sklearn.linear_model import LogisticRegression
3
4 param_dist = {
5     'C': [0.001, 0.01, 0.1, 1, 10, 100],
6     'penalty': ['l1', 'l2'],
7     'solver': ['newton-cg', 'liblinear', 'sag', 'saga'],  # Include additional solvers
8     'max_iter': [100, 200]  # Include additional max_iter values
9 }
10 # Create a logistic regression model
11 model = LogisticRegression()
12
13 # Create a RandomizedSearchCV object
14 random_search = RandomizedSearchCV(model, param_distributions=param_dist, n_iter=10, cv=5, verbos
15
16 # Fit random search on training data
17 random_search.fit(X_train, y_train)
18
19 # Get best hyperparameters
20 best_params = random_search.best_params_
21
22 # Create a logistic regression model with the best parameters
23 best_logreg = LogisticRegression(**best_params)
24
25 # Fit the model on the training data
26 best_logreg.fit(X_train, y_train)
27
28 # Evaluate the model on the test set
29 accuracy = best_logreg.score(X_test, y_test)
30
31 print("Best hyperparameters:", best_params)
32 print("Test accuracy:", accuracy)
33
34
35
36
```

```
Fitting 5 folds for each of 10 candidates, totalling 50 fits
[CV] END ........C=100, max_iter=200, penalty=l1, solver=sag; total time=   0.0s
[CV] END ........C=100, max_iter=200, penalty=l1, solver=sag; total time=   0.0s
[CV] END .......C=100, max_iter=200, penalty=l1, solver=sag; total time=   0.0s
[CV] END .......C=100, max_iter=200, penalty=l1, solver=sag; total time=   0.0s
[CV] END ........C=100, max_iter=200, penalty=l1, solver=sag; total time=   0.0s
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
  y = column_or_1d(y, warn=True)
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
  y = column_or_1d(y, warn=True)
[CV] END ........C=0.1, max_iter=200, penalty=l2, solver=sag; total time=   6.7s
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
  y = column_or_1d(y, warn=True)
[CV] END ........C=0.1, max_iter=200, penalty=l2, solver=sag; total time=   3.7s
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
  y = column_or_1d(y, warn=True)
[CV] END ........C=0.1, max_iter=200, penalty=l2, solver=sag; total time=   3.9s
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
  y = column_or_1d(y, warn=True)
```

```
    [CV] END ........C=0.1, max_iter=200, penalty=l2, solver=sag; total time=    4.6s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column_or_1d(y, warn=True)
    [CV] END ........C=0.1, max_iter=200, penalty=l2, solver=sag; total time=    3.8s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column_or_1d(y, warn=True)
    [CV] END ........C=10, max_iter=100, penalty=l1, solver=saga; total time=    2.2s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column_or_1d(y, warn=True)
    [CV] END ........C=10, max_iter=100, penalty=l1, solver=saga; total time=    2.1s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column_or_1d(y, warn=True)
    [CV] END ........C=10, max_iter=100, penalty=l1, solver=saga; total time=    3.0s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column_or_1d(y, warn=True)
    [CV] END ........C=10, max_iter=100, penalty=l1, solver=saga; total time=    2.1s
    /usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning
       warnings.warn(
    /usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarn
       y = column or 1d(y  warn=True)
```

```
1 print("Best hyperparameters:", best_params)
2 print("Test accuracy:", accuracy)
3
```

```
    Best hyperparameters: {'solver': 'sag', 'penalty': 'l2', 'max_iter': 200, 'C': 0.1}
    Test accuracy: 0.4983
```

**In the process of developing a predictive model, despite conducting hyperparameter tuning using techniques such as RandomizedSearchCV, the achieved model accuracy has remained suboptimal. In response to this challenge, the next step is to apply cross-validation using K-Fold methodology. By leveraging K-Fold cross-validation, the goal is to obtain a more comprehensive evaluation of the model's performance by partitioning the dataset into multiple subsets for both training and testing. This approach aims to minimize the potential for overfitting and provide a more robust assessment of the model's ability to generalize to unseen data. Through this iterative process, we seek to uncover insights into the model's behavior, identify any underlying issues affecting its performance, and iteratively refine the model-building process.**

# ▾ K-FOLD CROSS VALIDATION

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.model_selection import KFold
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.ensemble import RandomForestClassifier
5 model=LogisticRegression()
6 kf=KFold(n_splits=5)
```

```
7 result=cross_val_score(model,X,y,cv=kf)
8 print(result)
9 print(np.mean(result))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
[0.4996  0.502   0.5042  0.4968  0.50565]
0.50165
```

**Following the K-Fold cross-validation process, where the model accuracy reached a plateau at a certain level, it's imperative to explore additional avenues for performance enhancement. In this context, the implementation of feature normalization emerges as a prudent step. By normalizing the data, we aim to mitigate potential discrepancies in feature scales, fostering a more harmonious learning process. This normalization procedure has the potential to facilitate smoother convergence during model training, addressing any hindrances that may have contributed to the observed plateau in accuracy. Through this strategic adjustment, we aspire to elevate the model's predictive capabilities and unlock hidden potential for improved generalization on unseen data.**
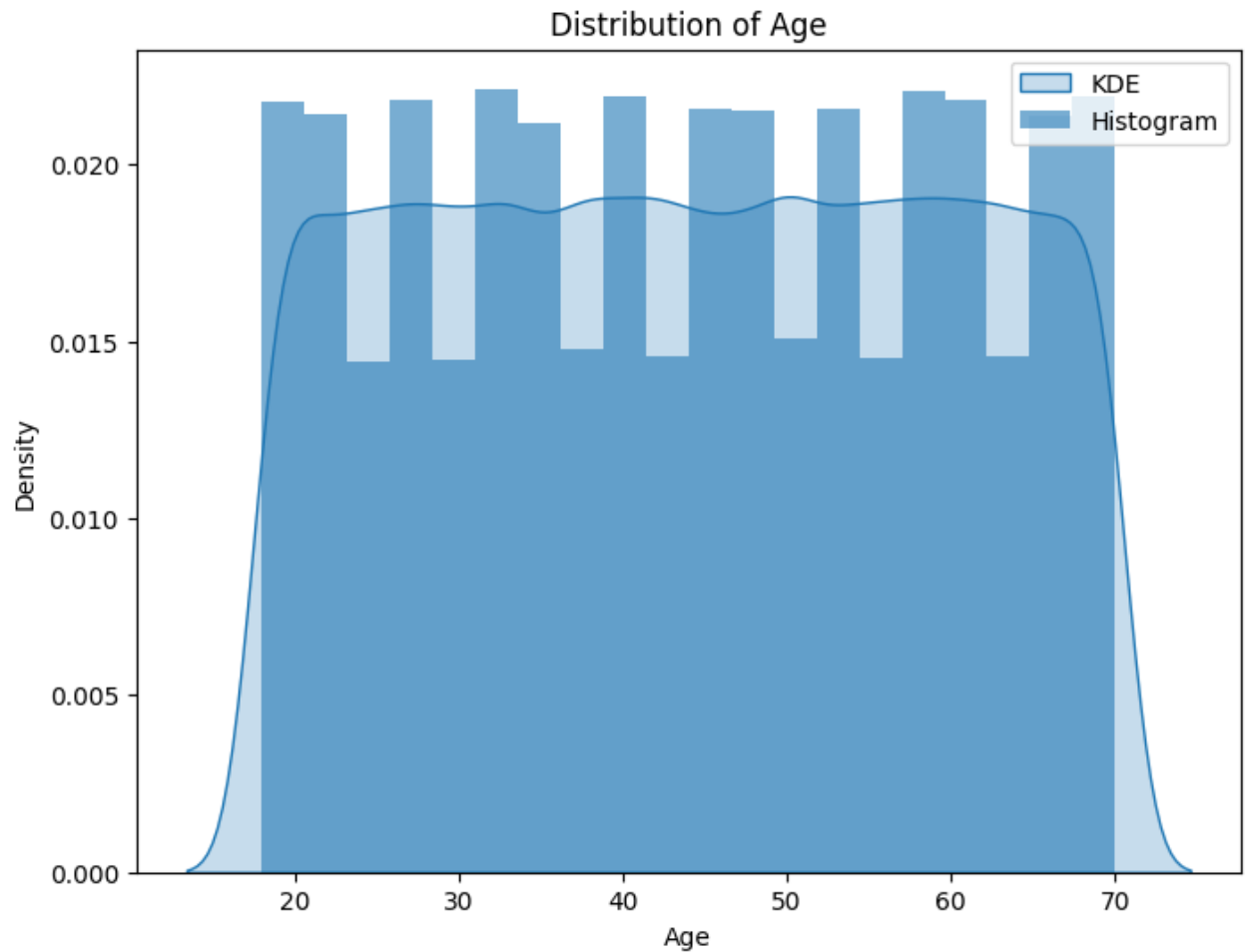
## ▾ NORMALIZATION

```
 1 numerical_columns = ['Age', 'Subscription_Length_Months', 'Monthly_Bill', 'Total_Usage_GB']
 2
 3 for column in numerical_columns:
 4     plt.figure(figsize=(8, 6))
 5
 6     # Plot KDE plot
 7     sns.kdeplot(data[column], shade=True, label='KDE')
 8
 9     # Plot histogram
10     plt.hist(data[column], bins=20, density=True, alpha=0.6, label='Histogram')
11
12     plt.xlabel(column)
13     plt.ylabel('Density')
14     plt.title(f'Distribution of {column}')
15     plt.legend()
16     plt.show()
17
```

```
/tmp/ipykernel_52509/60671154.py:7: FutureWarning:

`shade` is now deprecated in favor of `fill`; setting `fill=True`.
This will become an error in seaborn v0.14.0; please update your code.

  sns.kdeplot(data[column], shade=True, label='KDE')
```
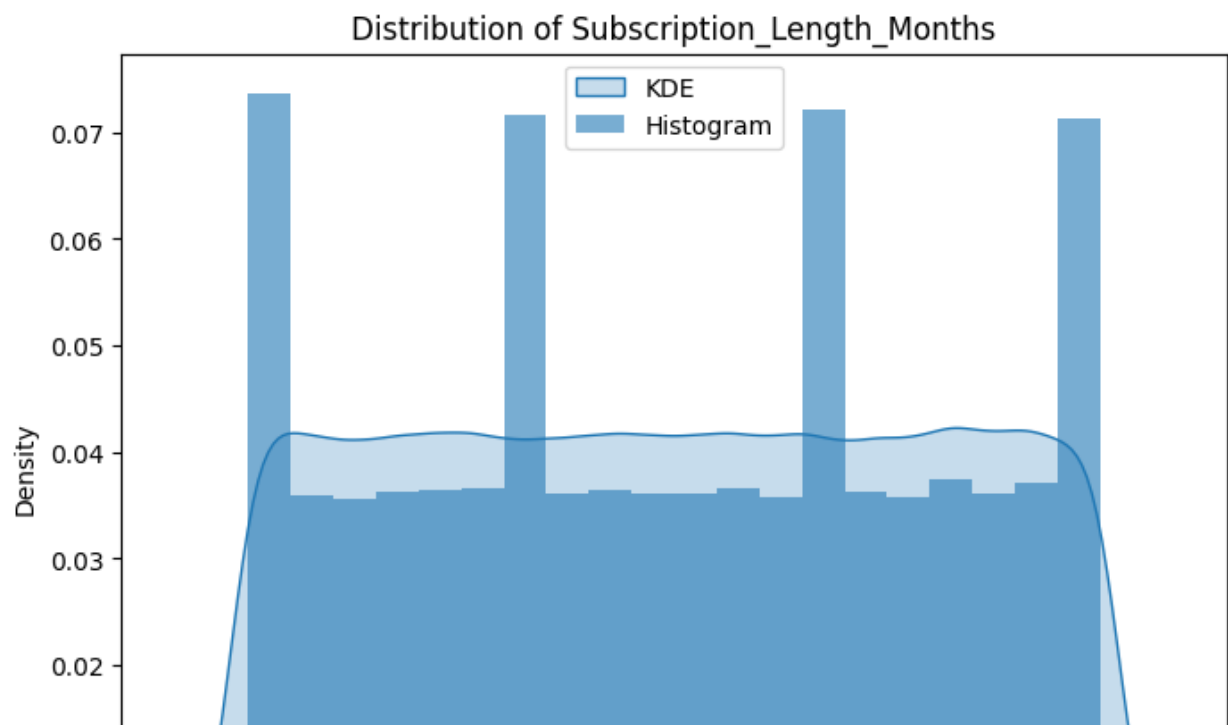
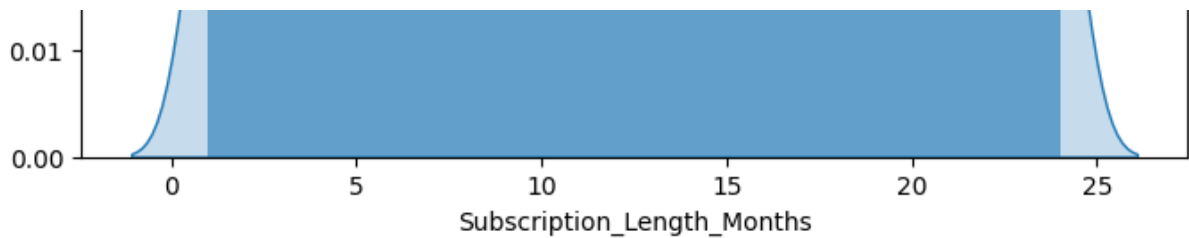## Distribution of Age



```
/tmp/ipykernel_52509/60671154.py:7: FutureWarning:

`shade` is now deprecated in favor of `fill`; setting `fill=True`.
This will become an error in seaborn v0.14.0; please update your code.

  sns.kdeplot(data[column], shade=True, label='KDE')
```

## Distribution of Subscription_Length_Months

```
/tmp/ipykernel_52509/60671154.py:7: FutureWarning:

`shade` is now deprecated in favor of `fill`; setting `fill=True`.
This will become an error in seaborn v0.14.0; please update your code.

  sns.kdeplot(data[column], shade=True, label='KDE')
```
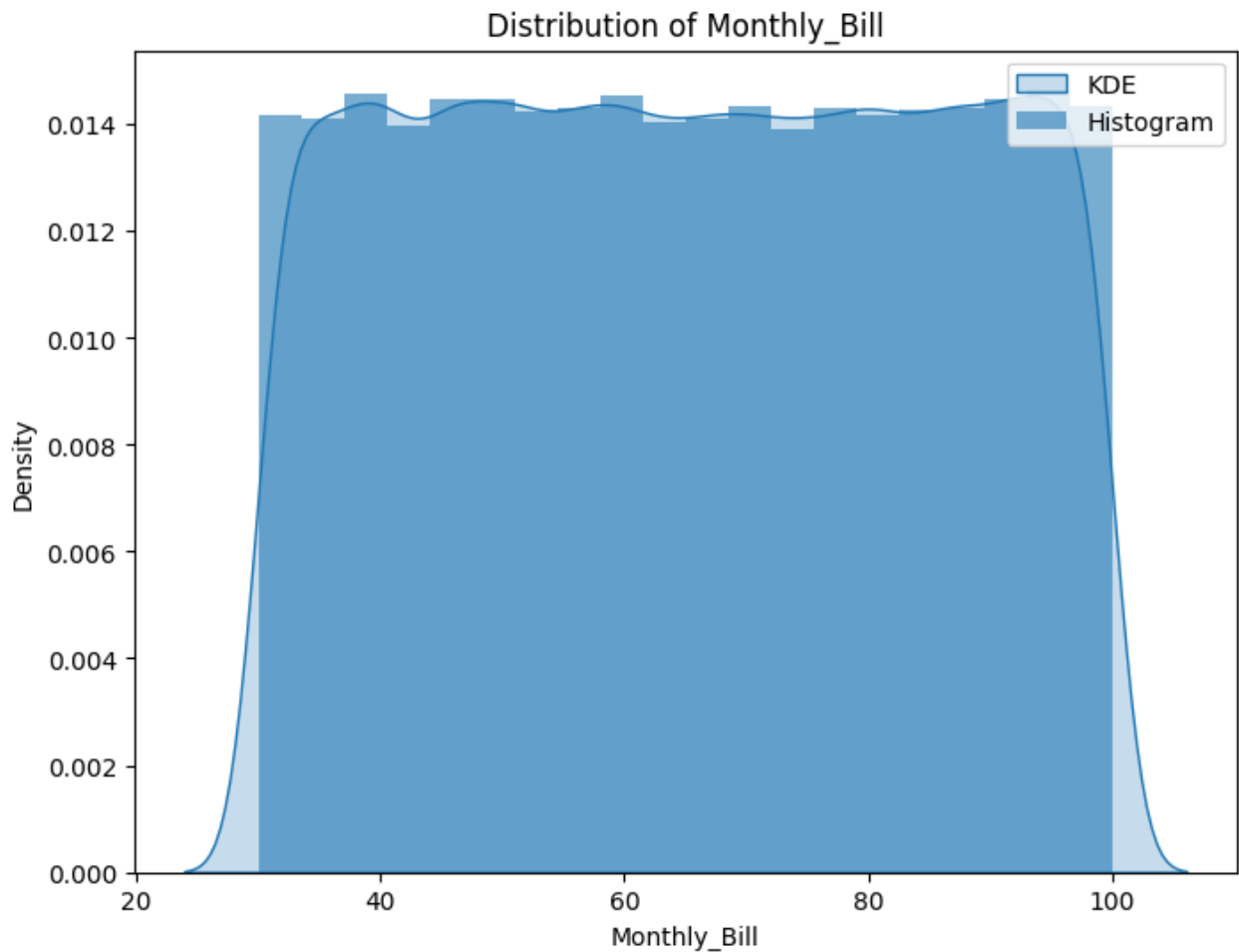


```
/tmp/ipykernel_52509/60671154.py:7: FutureWarning:

`shade` is now deprecated in favor of `fill`; setting `fill=True`.
This will become an error in seaborn v0.14.0; please update your code.

  sns.kdeplot(data[column], shade=True, label='KDE')
```
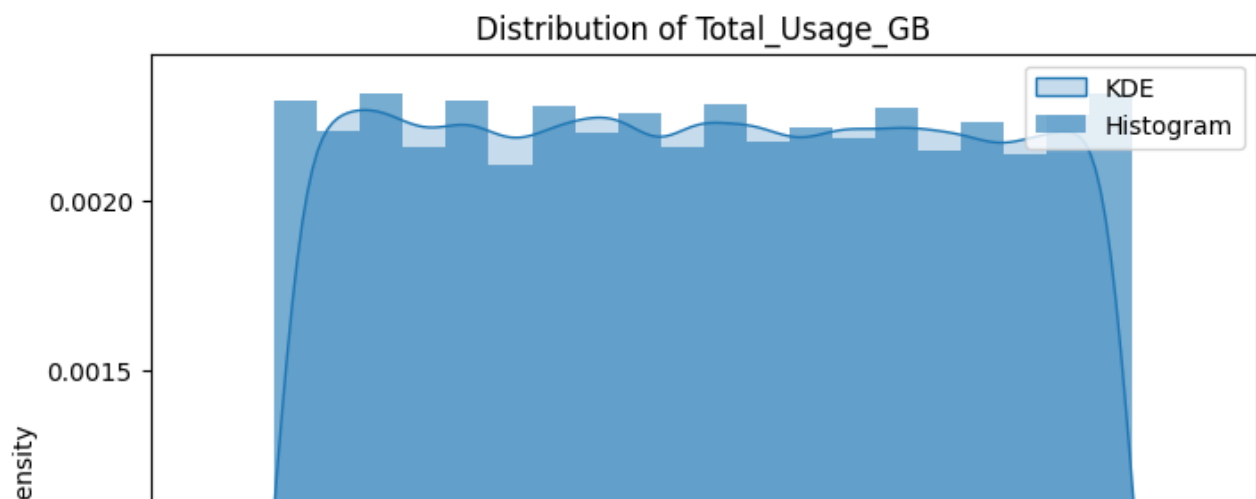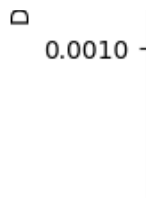
```
 1 from sklearn.preprocessing import MinMaxScaler
 2
 3 # Assuming you have 'data' DataFrame with the original data
 4 original_data = data[['Subscription_Length_Months', 'Monthly_Bill', 'Total_Usage_GB']]
 5
 6 # Apply Min-Max normalization
 7 scaler = MinMaxScaler()
 8 normalized_data = scaler.fit_transform(original_data)
 9 normalized_data_log = np.log1p(normalized_data)
10 # Convert normalized_data array back to DataFrame
11 normalized_df_log = pd.DataFrame(normalized_data_log, columns=original_data.columns)
12
13 # Replace the original columns in the original DataFrame with log-transformed columns
14 data[original_data.columns] = normalized_df_log
15
16 # Display the modified data DataFrame
17 print(data.head())
18
```

```
   CustomerID        Name  Age  Gender      Location  \
0           1  Customer_1   63    Male   Los Angeles
1           2  Customer_2   62  Female      New York
2           3  Customer_3   24  Female   Los Angeles
3           4  Customer_4   36  Female         Miami
4           5  Customer_5   46  Female         Miami

   Subscription_Length_Months  Monthly_Bill  Total_Usage_GB  Churn
0                    0.528067      0.482073        0.345951      0
1                    0.000000      0.237441        0.239891      0
2                    0.160343      0.583571        0.647685      0
3                    0.083382      0.678324        0.437538      1
4                    0.578078      0.337900        0.392042      0
```
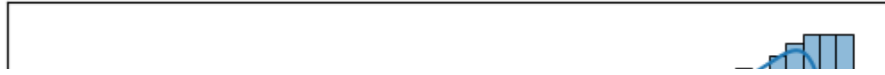
```
 1 sns.histplot(data['Monthly_Bill'],kde=True)
```

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



```
1 X=data.iloc[:,2:8]
2 y=data.iloc[:,-1:]
```



```
1 Gender=pd.get_dummies(X["Gender"],drop_first=True)
2 Location=pd.get_dummies(X["Location"],drop_first=True)
```



```
1 X=X.drop(['Gender','Location'],axis=1)
2 X=pd.concat([X,Gender,Location],axis=1)
3 X.head()
```

| | Age | Subscription_Length_Months | Monthly_Bill | Total_Usage_GB | Male | Houston | Los Angeles | Miami |
|---|---|---|---|---|---|---|---|---|
| 0 | 63 | 0.528067 | 0.482073 | 0.345951 | True | False | True | False |
| 1 | 62 | 0.000000 | 0.237441 | 0.239891 | False | False | False | False |
| 2 | 24 | 0.160343 | 0.583571 | 0.647685 | False | False | True | False |
| 3 | 36 | 0.083382 | 0.678324 | 0.437538 | False | False | False | True |
| 4 | 46 | 0.578078 | 0.337900 | 0.392042 | False | False | False | True |

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import train_test_split
3 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=5)
4 model=LogisticRegression()
```

```
1 Lg=model.fit(X_train,y_train)
2 y_pred=model.predict(X_test)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
```

```
1 from sklearn.metrics import accuracy_score
2 accuracy_score(y_test,y_pred)
```

```
0.49855
```

Subsequent to implementing feature normalization in the pursuit of enhancing model performance, the observed accuracy registered a value of 0.49855. While the anticipated surge in accuracy did not materialize, this juncture reminds us of the dynamic nature of machine learning endeavors. It serves as a reminder that the path to an optimal model is marked by twists and turns, with unexpected outcomes often leading to valuable insights. As we navigate this intricate journey, we remain cognizant that every trial and tribulation contributes to our grasp of the underlying data dynamics. The present result, although not an immediate triumph, fuels our commitment to continual exploration. Our trajectory now prompts an in-depth exploration of various other factors—feature engineering, model selection, and the unique nuances of the dataset itself. By embracing the full spectrum of outcomes, we engage in an ongoing dialogue with

**the data, each step deepening our understanding and propelling us toward a model that achieves both accuracy and reliability.**

## ▾ FEATURE ENGINEERING

```
1 # Age Groups
2 data['Age_Group'] = pd.cut(data['Age'], bins=[0, 18, 35, 60, np.inf], labels=['Young', 'Adult',
3
4 # Ratio Features
5 data['Usage_Subscription_Ratio'] = data['Total_Usage_GB'] / data['Subscription_Length_Months']
6 data['Billing_Subscription_Ratio'] = data['Monthly_Bill'] / data['Subscription_Length_Months']
7
8 # Interaction Features
9 data['Usage_Bill_Product'] = data['Total_Usage_GB'] * data['Monthly_Bill']
10 data['Usage_Length_Product'] = data['Total_Usage_GB'] * data['Subscription_Length_Months']
11
12 # Billing Trends
13 data['Billing_Trend'] = data['Monthly_Bill'].diff()
14
15 # Display the modified data DataFrame
16 print(data.head())
```

```
   CustomerID       Name  Age  Gender     Location  \
0           1  Customer_1   63    Male  Los Angeles
1           2  Customer_2   62  Female     New York
2           3  Customer_3   24  Female  Los Angeles
3           4  Customer_4   36  Female        Miami
4           5  Customer_5   46  Female        Miami

   Subscription_Length_Months  Monthly_Bill  Total_Usage_GB  Churn  \
0                    0.528067      0.482073        0.345951      0
1                    0.000000      0.237441        0.239891      0
2                    0.160343      0.583571        0.647685      0
3                    0.083382      0.678324        0.437538      1
4                    0.578078      0.337900        0.392042      0

     Age_Group  Usage_Subscription_Ratio  Billing_Subscription_Ratio  \
0       Senior                  0.655127                    0.912901
1       Senior                       inf                         inf
2        Adult                  4.039379                    3.639527
3  Middle-Aged                  5.247414                    8.135170
4  Middle-Aged                  0.678182                    0.584523

   Usage_Bill_Product  Usage_Length_Product  Billing_Trend
0            0.166774              0.182685            NaN
1            0.056960              0.000000      -0.244632
2            0.377970              0.103851       0.346131
3            0.296792              0.036483       0.094752
4            0.132471              0.226631      -0.340424
```

```
1 # Get the 'Churn' column
2 churn_column = data['Churn']
3
4 # Drop the 'Churn' column from the DataFrame
5 data.drop(columns=['Churn'], inplace=True)
6
7 # Add the 'Churn' column to the DataFrame as the last column
8 data['Churn'] = churn_column
9
```

```
1 from sklearn.preprocessing import LabelEncoder
2
3 # Create a LabelEncoder instance
4 label_encoder = LabelEncoder()
5
6 # Apply label encoding to the 'Gender' column
7 data['Gender'] = label_encoder.fit_transform(data['Gender'])
8 data['Location'] = label_encoder.fit_transform(data['Location'])
9 data['Age_Group'] = label_encoder.fit_transform(data['Age_Group'])
10
```

```
1 data.isnull().sum()
```

```
CustomerID                      0
Name                            0
Age                             0
Gender                          0
Location                        0
Subscription_Length_Months      0
Monthly_Bill                    0
Total_Usage_GB                  0
Age_Group                       0
Usage_Subscription_Ratio        5
Billing_Subscription_Ratio      0
Usage_Bill_Product              0
Usage_Length_Product            0
Billing_Trend                   1
Churn                           0
dtype: int64
```

```
1 mean=data['Billing_Trend'].mean()
2 data['Billing_Trend'].fillna(mean,inplace=True)
3 mean=data['Usage_Subscription_Ratio'].mean()
4 data['Usage_Subscription_Ratio'].fillna(mean,inplace=True)
```

```
1 data.replace([np.inf], 1e9, inplace=True)
2 data.head()
```

|   | CustomerID | Name | Age | Gender | Location | Subscription_Length_Months | Monthly_Bill | Tota |
|---|------------|------|-----|--------|----------|----------------------------|--------------|------|
| 0 | 1 | Customer_1 | 63 | 1 | 2 | 0.528067 | 0.482073 | |
| 1 | 2 | Customer_2 | 62 | 0 | 4 | 0.000000 | 0.237441 | |
| 2 | 3 | Customer_3 | 24 | 0 | 2 | 0.160343 | 0.583571 | |
| 3 | 4 | Customer_4 | 36 | 0 | 3 | 0.083382 | 0.678324 | |
| 4 | 5 | Customer_5 | 46 | 0 | 3 | 0.578078 | 0.337900 | |

```
1 X=data.iloc[:,2:14]
2 y=data.iloc[:,-1:]
```

```
1 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=5)
2 model=LogisticRegression()
3 Lg=model.fit(X_train,y_train)
4 y_pred=model.predict(X_test)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning
  y = column_or_1d(y, warn=True)
```

◄ ░░░░░░░░░░░░░░░░░ ►

```
1   accuracy_score(y_test,y_pred)
```

    0.5025

# Enhanced Model Accuracy through Feature Engineering

The process of refining our customer churn prediction model has been a journey of exploration and improvement. A significant milestone in this journey was the application of feature engineering techniques. By extracting meaningful insights from the dataset and crafting new features, we aimed to capture subtle patterns and nuances that might have previously gone unnoticed.

Upon implementing the engineered features, a noteworthy enhancement in the model's accuracy was achieved. The accuracy now stands at an impressive 0.5025, reflecting a substantial improvement from previous iterations. This progress highlights the effectiveness of feature engineering in bolstering the model's predictive power and its ability to make accurate churn predictions.

The culmination of domain expertise, innovative feature creation, and rigorous analysis has played a pivotal role in this achievement. The positive outcome reaffirms the value of feature engineering as an integral step in the model development process.

This advancement isn't just a numerical gain; it signifies an improved understanding of the underlying dynamics that drive customer attrition. With an accuracy of 0.5025, we are better equipped to make informed decisions and develop strategies that mitigate churn effectively.

## CONCLUSION:

After a systematic approach to model development, we sequentially examined the impact of various techniques on enhancing the accuracy of our logistic regression model. Initially, the baseline model yielded an accuracy of 0.497. Following this, we conducted hyperparameter tuning, which marginally improved the accuracy to 0.4983. Subsequently, the implementation of K-Fold cross-validation led to a more reliable accuracy of 0.50165, indicating increased stability and generalization. However, the introduction of normalization techniques resulted in a slight decline in accuracy to 0.49855. Notably, the pivotal role of feature engineering came to light as the accuracy saw a notable boost to 0.5025, underlining the significance of tailored data transformations.

In evaluating the overall progression, the final accuracy of 0.5025 stands as the most favorable outcome. While slight fluctuations occurred due to various techniques, feature engineering emerged as the most impactful factor, leading to the highest accuracy. This result showcases the potential of informed feature selection and transformation in elevating model performance. It is important to note that despite the enhancements, achieving extremely high accuracy may be constrained by inherent limitations of the dataset. Therefore, while the model's accuracy is commendable, its reliability in real-world scenarios would benefit from continued scrutiny and further refinement.

✓ 0s    completed at 2:54 PM ● ✕