

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

1. Data Preprocessing:
Load the provided dataset and perform initial data exploration.

In [2]: df = pd.read_excel(r'C:\Users\sutharsan\Downloads\customer_churn_large_dataset.xlsx')

In [3]: df.head()

Out[3]:
CustomerID      Name  Age  Gender  Location  Subscription_Length_Months  Monthly_Bill  Total_Usage_GB  Churn
0      1.0  Customer_1   63.0    Male  Los Angeles          17.0          73.36          236.0      0.0
1      2.0  Customer_2   62.0   Female   New York           1.0          48.76          172.0      0.0
2      3.0  Customer_3   24.0   Female  Los Angeles           5.0          85.47          460.0      1.0
3      4.0  Customer_4   36.0   Female   Miami           3.0          97.94          297.0      0.0
4      5.0  Customer_5   46.0   Female   Miami          19.0          58.14          266.0      0.0

In [4]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 9 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   CustomerID          100000 non-null  float64
 1   Name                100000 non-null  object
 2   Age                 100000 non-null  float64
 3   Gender              100000 non-null  object
 4   Location            100000 non-null  object
 5   Subscription_Length_Months  100000 non-null  float64
 6   Monthly_Bill        100000 non-null  float64
 7   Total_Usage_GB      100000 non-null  float64
 8   Churn               100000 non-null  float64
dtypes: float64(6), object(3)
memory usage: 5.1+ MB

In [5]: # df.shape

Handle missing data and outliers

In [6]: df.isnull().sum() # No Null Values

Out[6]:
CustomerID      0
Name            0
Age             0
Gender          0
Location        0
Subscription_Length_Months  0
Monthly_Bill    0
Total_Usage_GB  0
Churn           0
dtype: int64

In [7]: import seaborn as sns

In [8]: plt.figure(figsize=(10,10))
df = df.drop('CustomerID',axis=1) # no outliers
sns.boxplot(data = df)

Out[8]:
<AxesSubplot:~>

In [9]: df.columns
Index(['Name', 'Age', 'Gender', 'Location', 'Subscription_Length_Months',
       'Monthly_Bill', 'Total_Usage_GB', 'Churn'],
      dtype='object')

In [10]: # as a classification model countplot for balanced data
df.Churn.value_counts().plot(kind="bar") # balanced data

Out[10]:
<AxesSubplot:~>

In [11]: df.dtypes
Name      object
Age      float64
Gender    object
Location  object
Subscription_Length_Months  float64
Monthly_Bill      float64
Total_Usage_GB    float64
Churn             float64
dtype: object

In [12]: sns.countplot(x='Location',hue='Churn',data=df)

Out[12]:
<AxesSubplot:~>

In [13]: sns.countplot(x='Age',hue='Gender',data=df) # In Houston Male customer are higher others Female are higher

Out[13]:
<seaborn.axisgrid.FacetGrid at 0x16d493b8>

In [14]: sns.countplot(x='Churn',hue='Location',data=df) # In NYC and Miami Location--> customer are churn

Out[14]:
<seaborn.axisgrid.FacetGrid at 0x16d49328>

Feature Engineering -- onehotencoding

In [15]: df.Location.unique()

Out[15]:
array(['Los Angeles', 'New York', 'Miami', 'Chicago', 'Houston'],
      dtype=object)

In [16]: df = pd.get_dummies(df,columns=['Location'],drop_first=True)

In [17]: df.replace({'Male':0,'Female':1},inplace=True)

In [18]: df.head(2)

Out[18]:
   Name  Age  Gender  Subscription_Length_Months  Monthly_Bill  Total_Usage_GB  Churn  Location_Houston  Location_Los_Angeles  Location_Miami  Location_New_York
0  Customer_1   63.0      0          17.0          73.36          236.0      0.0              0              1              0              0
1  Customer_2   62.0      1           1.0          48.76          172.0      0.0              0              0              0              1

In [19]: df['Churn'] = df['Churn'].astype(int)

Visualization

In [20]: sns.histplot(x='Age',hue='Churn',data=df) # Here age around 50 is high possible to churn

Out[20]:
<AxesSubplot:~>

In [21]: sns.displot(x='Subscription_Length_Months',hue='Churn',data=df)
# distribution for subscription from this we can say by density is higher in 1,2,6,13,19,25 months are higher subscriber.
# like 2 month, 5 month,10 month,20 month and 25 month subscriber are high compare to other months.
# we can see clearly which month subscriber are leaving.

Out[21]:
<seaborn.axisgrid.FacetGrid at 0x16d49580>

In [22]: df.Gender.value_counts().plot(kind="bar") # we have equal no of gender subscriber.

Out[22]:
<AxesSubplot:~>

In [23]: plt.figure(figsize=(10,10))
sns.heatmap(data = df.corr(),annot = True)

Out[23]:
<AxesSubplot:~>

In [24]: y=df['Churn']
x=df.drop('Churn',axis=1,inplace=True)

In [25]: x=df.drop('Name',axis=1,inplace=True)

In [26]: y.isnull().sum()

Out[26]:
0

In [49]: x = df[['Age','Subscription_Length_Months']]

In [50]: df.columns
Index(['Age', 'Gender', 'Subscription_Length_Months', 'Monthly_Bill',
       'Total_Usage_GB', 'Location_Houston', 'Location_Los_Angeles',
       'Location_Miami', 'Location_New York'],
      dtype='object')

In [51]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)

In [52]: # Step 4 - Feature scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
x_train = sc.fit_transform(x_train)
x_test = sc.transform(x_test)

In [53]: # Step 5 - Logistic regression classifier
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier()
rf.fit(x_train, y_train)

Out[53]:
RandomForestClassifier()

In [54]: # Step 6 - Predicting logistic regression model on x_test
lr_pred = rf.predict(x_test)

In [55]: # Step 7 - Confusion matrix
from sklearn import metrics
cm = metrics.confusion_matrix(y_test, y_pred)
print(cm)
accuracy = metrics.accuracy_score(y_test, y_pred)
print("Accuracy score:",accuracy)
precision = metrics.precision_score(y_test, y_pred)
print("Precision score:",precision)
recall = metrics.recall_score(y_test, y_pred)
print("Recall score:",recall)

[[7289 7285]
 [1694 7284]]
Accuracy score: 0.5034666666666666
Precision score: 0.502726266873766
Recall score: 0.48641068447412356

In [56]: # Step 5 - Logistic regression classifier
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state=0, solver="liblinear")
classifier.fit(x_train, y_train)

Out[56]:
LogisticRegression(random_state=0, solver='liblinear')

In [57]: # Step 6 - Predicting logistic regression model on x_test
lr_pred = classifier.predict(x_test)

In [58]: # Step 7 - Confusion matrix
from sklearn import metrics
cm = metrics.confusion_matrix(y_test, lr_pred)
print(cm)
accuracy = metrics.accuracy_score(y_test, lr_pred)
print("Accuracy score:",accuracy)
precision = metrics.precision_score(y_test, lr_pred)
print("Precision score:",precision)
recall = metrics.recall_score(y_test, lr_pred)
print("Recall score:",recall)

[[15025  0]
 [14975  0]]
Accuracy score: 0.5008333333333334
Precision score: 0.0
Recall score: 0.0
C:\Users\sutharsan\anaconda3\lib\site-packages\sklearn\metrics\classification.py:1318: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use 'zero_division' parameter to control this behavior.
  warnings.warn(msg_start, msg_end)

In [59]: # Accuracy is low so HyperParameter tuning
from sklearn.model_selection import GridSearchCV
param_grid = [
    {'penalty': ['l1', 'l2', 'elasticnet', 'none'],
     'C': [0.1, 0.01, 0.001, 0.0001, 0.00001],
     'solver': ['lbfgs', 'newton-cg', 'liblinear', 'sag']}
]

In [60]: clf = GridSearchCV(classifier, param_grid, cv = 3, verbose=True, n_jobs=-1)

In [61]: best_clf = clf.fit(x,y)

In [62]: clf.best_estimator_
LogisticRegression(C=0.001, penalty='l1', random_state=0, solver='liblinear')

Out[62]:
LogisticRegression(C=0.001, penalty='l1', random_state=0, solver='liblinear')

In [63]: print(f'Accuracy - : {best_clf.score(x,y):.3f}')

Accuracy - : 0.502

In [ ]:
```