

Combined all classifier model

In [6]:

```
1 #importing all library
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.compose import ColumnTransformer
7 from sklearn.preprocessing import OneHotEncoder
8 from sklearn.metrics import confusion_matrix
9 from sklearn.metrics import accuracy_score
10 from sklearn.metrics import roc_auc_score, roc_curve
```

In [7]:

```
1 df = pd.read_csv("/Users/myyntiimac/Desktop/Churn_Modelling.csv")
2 df.head()
```

Out[7]:

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0
2	3	15619304	Onio	502	France	Female	42	8	159660.80	3	1
3	4	15701354	Boni	699	France	Female	39	1	0.00	2	0
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	1	1

In [66]:

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  -
0   RowNumber             10000 non-null  int64
1   CustomerId            10000 non-null  int64
2   Surname               10000 non-null  object
3   CreditScore           10000 non-null  int64
4   Geography             10000 non-null  object
5   Gender                10000 non-null  object
6   Age                  10000 non-null  int64
7   Tenure               10000 non-null  int64
8   Balance              10000 non-null  float64
9   NumOfProducts        10000 non-null  int64
10  HasCrCard            10000 non-null  int64
11  IsActiveMember       10000 non-null  int64
12  EstimatedSalary      10000 non-null  float64
13  Exited               10000 non-null  int64
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

In [8]:

```
1 X = df.iloc[:, 3:-1].values
2 X
```

Out[8]:

```
array([[619, 'France', 'Female', ..., 1, 1, 101348.88],
       [608, 'Spain', 'Female', ..., 0, 1, 112542.58],
       [502, 'France', 'Female', ..., 1, 0, 113931.57],
       ...,
       [709, 'France', 'Female', ..., 0, 1, 42085.58],
       [772, 'Germany', 'Male', ..., 1, 0, 92888.52],
       [792, 'France', 'Female', ..., 1, 0, 38190.78]], dtype=object)
```

In [9]:

```
1 y = df.iloc[:, -1].values
2 y
```

Out[9]:

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

In [10]:

```
1 #converting gender column into numerical
2 le = LabelEncoder()
3 X[:, 2] = le.fit_transform(X[:, 2])
```

In [11]:

```
1 # Assuming 'X' is your input array and 'column_index' is the index of the column you want to one-hot
2 column_index = 1 # Example column index
3
4 # Convert the array to a DataFrame
5 df = pd.DataFrame(X)
6
7 # Perform one-hot encoding using get_dummies
8 encoded_df = pd.get_dummies(df, columns=[column_index], drop_first=True)
9
10 # Extract the values from the encoded DataFrame
11 X_encoded = encoded_df.values
```

In [12]:

```
1 X_encoded
```

Out[12]:

```
array([[619, 0, 42, ..., 101348.88, 0, 0],
       [608, 0, 41, ..., 112542.58, 0, 1],
       [502, 0, 42, ..., 113931.57, 0, 0],
       ...,
       [709, 0, 36, ..., 42085.58, 0, 0],
       [772, 1, 42, ..., 92888.52, 1, 0],
       [792, 0, 28, ..., 38190.78, 0, 0]], dtype=object)
```

In [13]:

```
1 #Split the dataset
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X_encoded, y, test_size = 0.2, random_state = 0)
```

In [14]:

```
1
2 from sklearn.ensemble import AdaBoostClassifier
3
```

In [15]:

```
1 AdaBoostClassifier= AdaBoostClassifier()
```

In [16]:

```
1 AdaBoostClassifier.fit(X_train, y_train)
```

Out[16]:

```
▼ AdaBoostClassifier
AdaBoostClassifier()
```



In [25]:

```
1 # Make predictions on the test data
2 y_pred = AdaBoostClassifier.predict(X_test)
3 y_pred
```

Out[25]:

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

In [21]:

```
1 # Calculate the accuracy score
2 accuracy = accuracy_score(y_test, y_pred)
3 accuracy
```

Out[21]:

```
0.8655
```

In [38]:

```
1 #ROC AND AUC
2 from sklearn.metrics import roc_curve, roc_auc_score
3
```

In [22]:

```
1 # Compute the False Positive Rate (FPR), True Positive Rate (TPR), and thresholds
2 fpr, tpr, thresholds = roc_curve(y_test, y_pred)
```

In [23]:

```
1 auc = roc_auc_score(y_test, y_pred)
2 auc
```

Out[23]:

```
0.7489570029799915
```

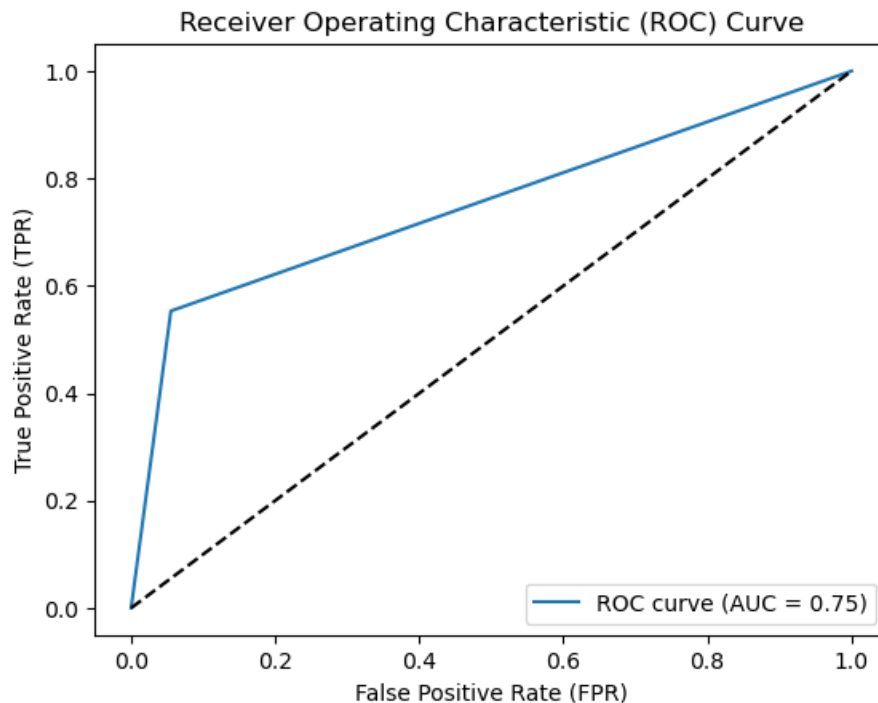


In [24]:

```

1 # Plotting the ROC curve
2 plt.plot(fpr, tpr, label='ROC curve (AUC = {:.2f})'.format(auc))
3 plt.plot([0, 1], [0, 1], 'k--') # Random guess line
4 plt.xlabel('False Positive Rate (FPR)')
5 plt.ylabel('True Positive Rate (TPR)')
6 plt.title('Receiver Operating Characteristic (ROC) Curve')
7 plt.legend(loc='lower right')
8 plt.show()

```



1 Insight: In the case of AUC = 0.74, the model demonstrates reasonable discriminative ability, but there is still room for improvement. It correctly ranks 74% of the positive samples higher than the negative samples, on average, across different classification thresholds. However, it might misclassify some instances, leading to false positives or false negatives.

In [26]:

```

1
2 from sklearn.model_selection import cross_val_score

```

In [27]:

```

1 scores = cross_val_score(AdaBoostClassifier, X_train, y_train, cv=5)
2 scores
3

```

Out[27]:

```
array([0.85625 , 0.865   , 0.84    , 0.85375 , 0.849375])
```

In [28]:

```

1 # Calculate the mean cross-validation score
2 mean_score = np.mean(scores)
3 mean_score

```

Out[28]:

```
0.852875
```

In [29]:

```

1 from sklearn.model_selection import GridSearchCV

```

In [30]:

```
1 # Define the parameter grid
2 param_grid = {
3     'n_estimators': [50, 100, 200],
4     'learning_rate': [0.1, 0.5, 1.0]
5 }
```

In [32]:

```
1 # Create an object of GridSearchCV
2 grid_search = GridSearchCV(AdaBoostClassifier, param_grid, cv=5)
```

In [33]:

```
1 # Fit the GridSearchCV on the training data
2 grid_search.fit(X_train, y_train)
3
```

Out[33]:

```
GridSearchCV
├── estimator: AdaBoostClassifier
│   └── AdaBoostClassifier
```

In [34]:

```
1 # Get the best parameters and best score
2 best_params = grid_search.best_params_
3 best_score = grid_search.best_score_
4 # Print the best parameters and best score
5 print("Best Parameters:", best_params)
6 print("Best Score:", best_score)
```

Best Parameters: {'learning_rate': 0.1, 'n_estimators': 200}
Best Score: 0.8551249999999999

In [35]:

```
1 #Apply random search CV
2 from sklearn.model_selection import RandomizedSearchCV
3 from scipy.stats import uniform#learning rate need to be uniform
```

In [36]:

```
1 # Define the parameter distribution
2 param_dist = {
3     'n_estimators': [50, 100, 200],
4     'learning_rate': uniform(0.1, 1.0)
5 }
```

In [37]:

```
1 # Create an instance of RandomizedSearchCV
2 random_search = RandomizedSearchCV(AdaBoostClassifier, param_distributions=param_dist, n_iter=5, cv=5)
3
```

In [38]:

```
1 # Fit the RandomizedSearchCV on the training data
2 random_search.fit(X_train, y_train)
```

Out[38]:

```
RandomizedSearchCV
├── estimator: AdaBoostClassifier
│   └── AdaBoostClassifier
```

In [39]:

```
1 # Get the best parameters and best score
2 best_params = random_search.best_params_
3 best_score = random_search.best_score_
4
5 # Print the best parameters and best score
6 print("Best Parameters:", best_params)
7 print("Best Score:", best_score)
```

Best Parameters: {'learning_rate': 0.7184721402879922, 'n_estimators': 200}
Best Score: 0.851875

```
1 Insight: we got 4 accuracy score
2 Adaboost classifier=0.8655
3 Kfold validation=0.852875
4 Gridsearch=0.8551249999999999
5 Randomsearch=0.851875
6 best learning rate .1 and .7 and n_estimator is 200
```

