

Gradient boosting

In [3]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
```

In [4]:

```
1 dataset = pd.read_csv('/Users/myyntiimac/Desktop/Social_Network_Ads.csv')
2 X = dataset.iloc[:, [2, 3]].values
3 y = dataset.iloc[:, -1].values
```

In [5]:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
```

In [7]:

```
1 from sklearn.ensemble import GradientBoostingClassifier
```

In [9]:

```
1 classifier=GradientBoostingClassifier()
```

In [10]:

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

Out[10]:

```
▼ GradientBoostingClassifier
GradientBoostingClassifier()
```

In [11]:

```
1 y_pred = classifier.predict(X_test)
```

In [12]:

```
1 from sklearn.metrics import confusion_matrix
2 cm = confusion_matrix(y_test, y_pred)
3 print(cm)
```

```
[[64  4]
 [ 5 27]]
```

In [14]:

```
1 from sklearn.metrics import accuracy_score
2 # Calculate the accuracy score
3 accuracy = accuracy_score(y_test, y_pred)
4 accuracy
```

Out[14]:

0.91

In [15]:

```
1 bias = classifier.score(X_train, y_train)
2 bias
```

Out[15]:

0.9733333333333334

In [16]:

```
1 variance = classifier.score(X_test, y_test)
2 variance
```

Out[16]:

0.91

In [17]:

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

In [18]:

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

In [19]:

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

Out[19]:

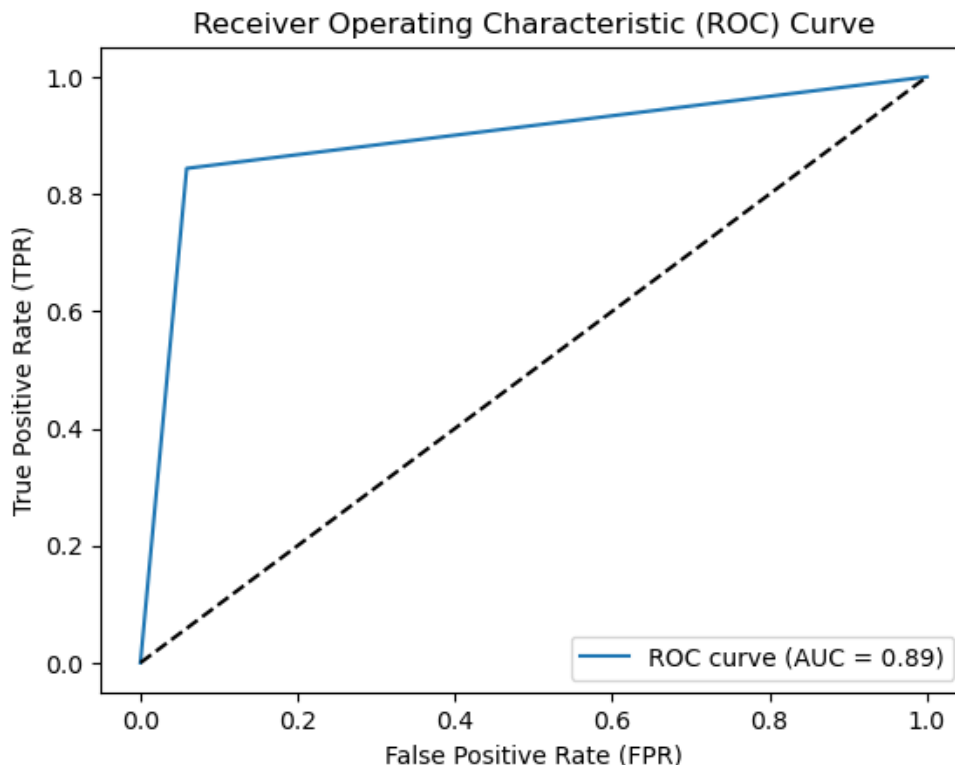
0.8924632352941176

In [20]:

```

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 = 0.89, the model demonstrates reasonable discriminative ability, but there is still room for improvement. It correctly ranks 89% 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 [21]:

```

1 from sklearn.model_selection import cross_val_score
2 accuracies = cross_val_score(estimator = classifier, X = X_train, y = y_train, cv = 10)
3 print("Accuracy: {:.2f} %".format(accuracies.mean()*100))
4 print("Standard Deviation: {:.2f} %".format(accuracies.std()*100))

```

Accuracy: 88.33 %

Standard Deviation: 5.43 %

In [22]:

```

1 from sklearn.model_selection import GridSearchCV

```

In [30]:

```
1 # Define the parameter grid
2 param_grid = {
3     'n_estimators': [100, 200, 300], # Number of trees
4     'learning_rate': [0.1, 0.01, 0.001], # Learning rate
5     'max_depth': [3, 5, 7], # Maximum depth of trees
6     'min_samples_split': [2, 5, 10], # Minimum number of samples required to split a node
7     'min_samples_leaf': [1, 2, 4], # Minimum number of samples required at a leaf node
8     'subsample': [0.8, 1.0], # Fraction of samples to be used for fitting each tree
9     'max_features': [1.0, 'sqrt'] # Number of features to consider for the best split
10 }
```

In [32]:

```
1 # Create the GridSearchCV object
2 grid_search = GridSearchCV(estimator=classifier, param_grid=param_grid, cv=10)
3 grid_search = GridSearchCV(classifier, param_grid, cv=10)
4 grid_search = grid_search.fit(X_train, y_train)
5 best_accuracy = grid_search.best_score_
6 best_parameters = grid_search.best_params_
7 print("Best Accuracy: {:.2f} %".format(best_accuracy*100))
8 print("Best Parameters:", best_parameters)
9
```

Best Accuracy: 90.67 %

Best Parameters: {'learning_rate': 0.01, 'max_depth': 3, 'max_features': 'sqrt', 'min_samples_leaf': 4, 'min_samples_split': 2, 'n_estimators': 100, 'subsample': 0.8}

In []:

1