

Chetna Mundra - 1BM21AI036 - ACY AAT

Importing libraries

```
In [ ]: 1 from imblearn.over_sampling import RandomOverSampler, ADASYN
2 from collections import Counter
3
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import pandas as pd
7 import numpy as np
8
9 from sklearn.model_selection import train_test_split, KFold,
GridSearchCV, ParameterGrid, cross_val_score
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.tree import DecisionTreeClassifier
12 #from sklearn.ensemble import RandomForestClassifier,
AdaBoostClassifier
13 #from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.preprocessing import MinMaxScaler, PowerTransformer,
StandardScaler
15 from sklearn.metrics import roc_auc_score
16 from sklearn.svm import SVC
17 from sklearn import metrics
18
19
20
21
22 import warnings
23 warnings.filterwarnings("ignore")
```

Reading Data

```
In [ ]: 1 df = pd.read_csv("creditcard.csv")
```

```
In [ ]: 1 classes=df['Class'].value_counts()
2 normal_share=classes[0]/df['Class'].count()*100
3 fraud_share=classes[1]/df['Class'].count()*100
```

```
In [ ]: 1 classes
```

Out[11]:

	count
0	284315
1	492

Class	
0	284315
1	492

dtype: int64

```
In [ ]: 1 df['TimeMin'] = df['Time'] / 60
        2 df['TimeHour'] = df['Time'] / 60**2
```

```
In [ ]: 1 df.drop(['Time', 'TimeMin', 'TimeHour'], axis=1, inplace=True)
```

Train test split

```
In [ ]: 1 # Train test split
        2 y = df.pop("Class")
        3 X = df
        4
        5 # Using stratify=y for splitting data into stratified fashion
        6 X_train, X_test, y_train, y_test = train_test_split(X, y,
        train_size=0.75, stratify=y, random_state=42)
```

Results

```
In [ ]: 1 # creating a results dataframe for later to evaluate the models
        2 results = pd.DataFrame(columns = ['model_name', 'threshold',
        'recall', 'roc_auc_score'])
```

Model Evaluation Function

```

In [ ]: 1 def model_evaluation(y_pred_proba_test):
2 #     Report for different thresholds
3     thresholds = [i * 0.1 for i in range(0, 10)]
4     print("-----Results-----")
5     #     best ROC score initialisation
6     best_roc_score = 0
7
8     #     Iterating through every threshold from 0.1 to 0.9
9     for threshold in thresholds:
10         y_pred = np.where(y_pred_proba_test[:, 1] > threshold, 1, 0)
11     #     Calculating different metrics
12         accuracy = str(round(metrics.accuracy_score(y_test, y_pred),
13             3))
14         precision = str(round(metrics.precision_score(y_test,
15             y_pred), 3))
16         recall = str(round(metrics.recall_score(y_test, y_pred), 3))
17         roc_auc = str(round(metrics.roc_auc_score(y_test, y_pred),
18             3))
19     #     Setting the best roc score, threshold, recall scores.
20     if float(roc_auc) > best_roc_score:
21         best_roc_score = float(roc_auc)
22         best_threshold = threshold
23         best_recall_score = recall
24     #     printing the results for every threshold
25     print("-----for Test with threshold", round(threshold,
26         2), "-----")
27     print("accuracy\tprecision\trecall\t\troc_auc")
28     print("\t\t".join([accuracy, precision, recall, roc_auc]))
29     print("\n")
30     #     Confusion Matrix
31     print("\t\tCONFUSION MATRIX")
32     confusion_matrix =
33     pd.DataFrame(metrics.confusion_matrix(y_test, y_pred),
34         columns=['Predicted
Negative', 'Predicted Positive'],
35         index=['Actual Negative',
36             'Actual Positive'])
37     print(confusion_matrix)
38     print("\n")
39     print("BEST ROC AUC SCORE is ", best_roc_score, "at the
threshold", best_threshold)
40     return best_roc_score, best_threshold, best_recall_score

```

Draw ROC

```
In [ ]: 1 def draw_roc( actual, probs ):
2         fpr, tpr, thresholds = metrics.roc_curve( actual, probs,
3                                                     drop_intermediate =
4             False )
5         auc_score = metrics.roc_auc_score( actual, probs )
6         plt.figure(figsize=(5, 5))
7         plt.plot( fpr, tpr, label='ROC curve (area = %0.2f)' % auc_score
8             )
9         plt.plot([0, 1], [0, 1], 'k--')
10        plt.xlim([0.0, 1.0])
11        plt.ylim([0.0, 1.0])
12        plt.xlabel('False Positive Rate or [1 - True Negative Rate]')
13        plt.ylabel('True Positive Rate')
14        plt.title('Receiver operating characteristic')
15        plt.legend(loc="lower right")
16        plt.show()
```

XG BOOST

In []:

```
1  # Import necessary libraries
2  import xgboost as xgb
3  from sklearn.model_selection import GridSearchCV, KFold
4  from sklearn.pipeline import make_pipeline
5  import pandas as pd
6
7  # Initialize the XGBoost model
8  def initialize_xgboost():
9      return xgb.XGBClassifier(eval_metric='logloss', n_estimators=100,
10                             max_depth=3)
11
12 # Train the model using GridSearchCV with pipeline
13 def train_xgboost(X_train, y_train, cv):
14     xgb_model = initialize_xgboost()
15     pipeline = make_pipeline(xgb_model)
16     param_grid = {
17         'xgbclassifier__n_estimators': [50, 100],
18         'xgbclassifier__max_depth': [3, 5]
19     }
20     grid = GridSearchCV(pipeline, param_grid=param_grid,
21                        scoring='roc_auc', cv=cv, n_jobs=-1, verbose=1000)
22     grid.fit(X_train, y_train)
23     print("BEST GRID SCORE", grid.best_score_)
24     print("BEST GRID PARAMS")
25     print(grid.best_params_)
26     return grid.best_estimator_.named_steps['xgbclassifier']
27
28 # Evaluate the model
29 def evaluate_xgboost(model, X_test, y_test):
30     # Get the predicted probabilities for both classes
31     y_pred_proba_test = model.predict_proba(X_test)
32     # Pass the full probability array to model_evaluation
33     best_roc_score, best_threshold, best_recall_score =
34     model_evaluation(y_pred_proba_test)
35     return best_roc_score, best_threshold, best_recall_score
36
37 # Example usage
38 cv = KFold(3)
39 best_xgb_model = train_xgboost(X_train, y_train, cv)
40 best_roc_score, best_threshold, best_recall_score =
41     evaluate_xgboost(best_xgb_model, X_test, y_test)
42
43 # Update results dataframe
44 data = pd.DataFrame(['XGBOOST', best_threshold, best_recall_score,
45                     best_roc_score], columns=results.columns)
46 results = pd.concat([results, data], ignore_index=False)
47
48 # Optional: Plot ROC curve
49 draw_roc(y_test, best_xgb_model.predict_proba(X_test)[: , 1])
```

Fitting 3 folds for each of 4 candidates, totalling 12 fits
BEST GRID SCORE nan
BEST GRID PARAMS
{'xgbclassifier__max_depth': 3, 'xgbclassifier__n_estimators': 50}

-----Results-----

-----for Test with threshold 0.0 -----

accuracy	precision	recall	roc_auc
0.002	0.002	1.0	0.5

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	0	71079
Actual Positive	0	123

-----for Test with threshold 0.1 -----

accuracy	precision	recall	roc_auc
0.999	0.75	0.829	0.914

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71045	34
Actual Positive	21	102

-----for Test with threshold 0.2 -----

accuracy	precision	recall	roc_auc
0.999	0.852	0.797	0.898

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71062	17
Actual Positive	25	98

-----for Test with threshold 0.3 -----

accuracy	precision	recall	roc_auc
0.999	0.882	0.789	0.894

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71066	13
Actual Positive	26	97

-----for Test with threshold 0.4 -----

accuracy	precision	recall	roc_auc
1.0	0.923	0.78	0.89

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71071	8
Actual Positive	27	96

-----for Test with threshold 0.5 -----

accuracy	precision	recall	roc_auc
1.0	0.931	0.772	0.886

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71072	7
Actual Positive	28	95

-----for Test with threshold 0.6 -----

accuracy	precision	recall	roc_auc
1.0	0.949	0.756	0.878

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71074	5
Actual Positive	30	93

-----for Test with threshold 0.7 -----

accuracy	precision	recall	roc_auc
0.999	0.968	0.732	0.866

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71076	3
Actual Positive	33	90

-----for Test with threshold 0.8 -----

accuracy	precision	recall	roc_auc
0.999	0.966	0.691	0.846

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71076	3
Actual Positive	38	85

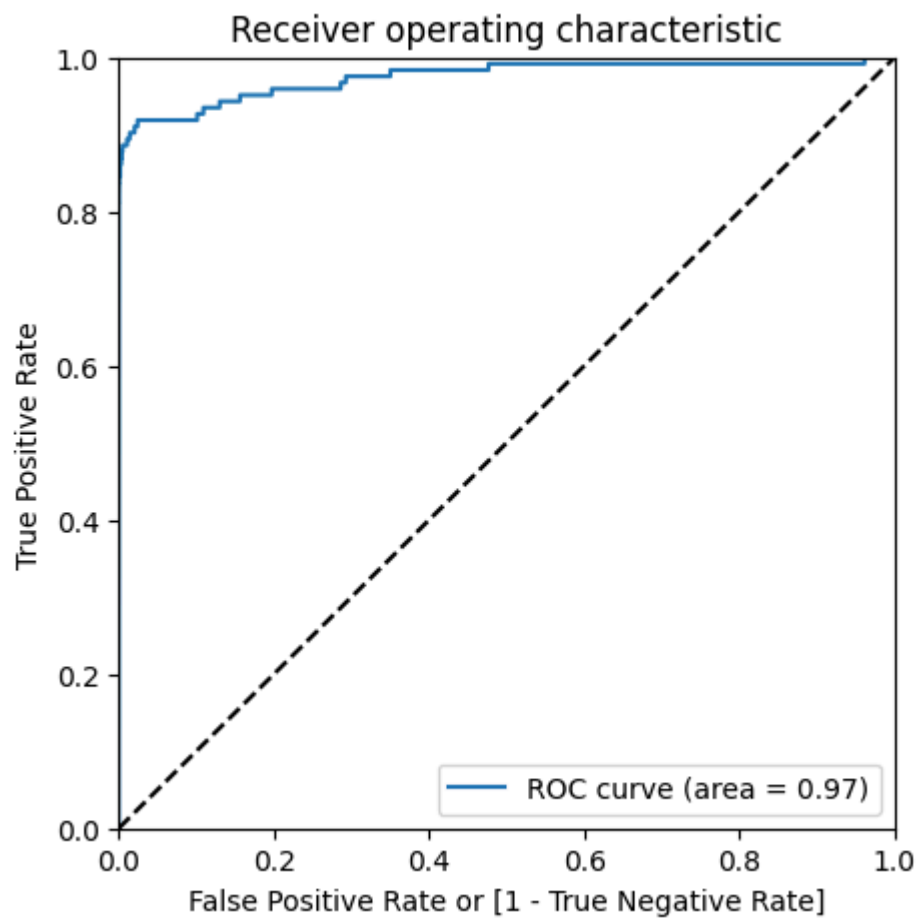
-----for Test with threshold 0.9 -----

accuracy	precision	recall	roc_auc
0.999	0.986	0.593	0.797

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71078	1
Actual Positive	50	73

BEST ROC AUC SCORE is 0.914 at the threshold 0.1



Random Forest

In []:

```
1  # Import necessary libraries
2  from sklearn.ensemble import RandomForestClassifier
3  from sklearn.model_selection import train_test_split
4  import pandas as pd
5
6  # Train a simple Random Forest model
7  def train_random_forest(X_train, y_train):
8      rf_model = RandomForestClassifier(n_estimators=50, max_depth=10,
9      random_state=42)
10     rf_model.fit(X_train, y_train)
11     return rf_model
12
13 # Example usage
14 X_train, X_test, y_train, y_test = train_test_split(X, y,
15     test_size=0.25, stratify=y, random_state=42)
16 best_rf_model = train_random_forest(X_train, y_train)
17
18 # Evaluate the model using the existing function
19 y_pred_proba_test = best_rf_model.predict_proba(X_test)
20 best_roc_score, best_threshold, best_recall_score =
21     model_evaluation(y_pred_proba_test)
22
23 # Update results dataframe
24 data = pd.DataFrame(['RANDOM FOREST', best_threshold,
25     best_recall_score, best_roc_score], columns=results.columns)
26 results = pd.concat([results, data], ignore_index=False)
```

-----Results-----

-----for Test with threshold 0.0 -----

accuracy	precision	recall	roc_auc
0.002	0.002	1.0	0.5

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	1	71078
Actual Positive	0	123

-----for Test with threshold 0.1 -----

accuracy	precision	recall	roc_auc
0.999	0.707	0.846	0.922

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71036	43
Actual Positive	19	104

-----for Test with threshold 0.2 -----

accuracy	precision	recall	roc_auc
0.999	0.811	0.837	0.919

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71055	24
Actual Positive	20	103

-----for Test with threshold 0.3 -----

accuracy	precision	recall	roc_auc
0.999	0.877	0.813	0.906

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71065	14
Actual Positive	23	100

-----for Test with threshold 0.4 -----

accuracy	precision	recall	roc_auc
1.0	0.933	0.789	0.894

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	71072	7
Actual Positive	26	97

-----for Test with threshold 0.5 -----

accuracy	precision	recall	roc_auc
1.0	0.94	0.764	0.882

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	71073	6
Actual Positive	29	94

-----for Test with threshold 0.6 -----

	accuracy	precision	recall	roc_auc
	0.999	0.966	0.699	0.85

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	71076	3
Actual Positive	37	86

-----for Test with threshold 0.7 -----

	accuracy	precision	recall	roc_auc
	0.999	0.976	0.65	0.825

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	71077	2
Actual Positive	43	80

-----for Test with threshold 0.8 -----

	accuracy	precision	recall	roc_auc
	0.999	0.971	0.537	0.768

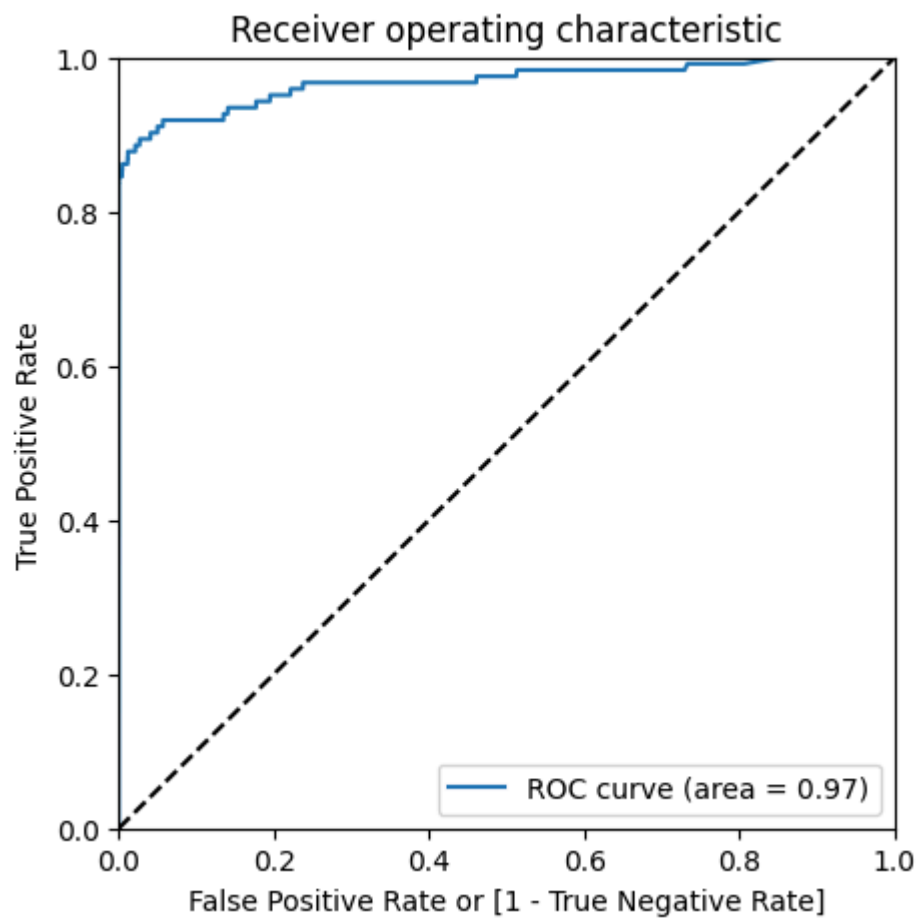
	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	71077	2
Actual Positive	57	66

-----for Test with threshold 0.9 -----

	accuracy	precision	recall	roc_auc
	0.999	0.98	0.407	0.703

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	71078	1
Actual Positive	73	50

BEST ROC AUC SCORE is 0.922 at the threshold 0.1



Using SMOTE analysis

In []:

```
1  # Import necessary libraries
2  from sklearn.linear_model import LogisticRegression
3  from sklearn.model_selection import train_test_split
4  from imblearn.over_sampling import SMOTE
5  import pandas as pd
6
7  # Apply SMOTE to balance the dataset
8  def apply_smote(X_train, y_train):
9      smote = SMOTE(random_state=42)
10     X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
11     return X_resampled, y_resampled
12
13  # Train a simple Logistic Regression model
14  def train_logistic_regression(X_train, y_train):
15      lr_model = LogisticRegression(max_iter=1000, random_state=42)
16      lr_model.fit(X_train, y_train)
17      return lr_model
18
19  # Example usage
20  X_train, X_test, y_train, y_test = train_test_split(X, y,
21      test_size=0.25, stratify=y, random_state=42)
22  X_resampled, y_resampled = apply_smote(X_train, y_train)
23
24  best_lr_model = train_logistic_regression(X_resampled, y_resampled)
25
26  # Evaluate the model using the existing function
27  y_pred_proba_test = best_lr_model.predict_proba(X_test)
28  best_roc_score, best_threshold, best_recall_score =
29  model_evaluation(y_pred_proba_test)
30
31  # Update results dataframe
32  data = pd.DataFrame(['LOGISTIC REGRESSION (SMOTE)', best_threshold,
33      best_recall_score, best_roc_score], columns=results.columns)
34  results = pd.concat([results, data], ignore_index=False)
35
36  # Optional: Plot ROC curve
37  draw_roc(y_test, y_pred_proba_test[:, 1])
```

-----Results-----

-----for Test with threshold 0.0 -----

accuracy	precision	recall	roc_auc
0.002	0.002	1.0	0.5

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	0	71079
Actual Positive	0	123

-----for Test with threshold 0.1 -----

accuracy	precision	recall	roc_auc
0.879	0.013	0.943	0.911

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	62471	8608
Actual Positive	7	116

-----for Test with threshold 0.2 -----

accuracy	precision	recall	roc_auc
0.936	0.024	0.911	0.923

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	66546	4533
Actual Positive	11	112

-----for Test with threshold 0.3 -----

accuracy	precision	recall	roc_auc
0.961	0.038	0.894	0.928

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	68305	2774
Actual Positive	13	110

-----for Test with threshold 0.4 -----

accuracy	precision	recall	roc_auc
0.974	0.056	0.886	0.93

CONFUSION MATRIX

	Predicted Negative	Predicted Positive
Actual Negative	69258	1821
Actual Positive	14	109

-----for Test with threshold 0.5 -----

accuracy	precision	recall	roc_auc
0.982	0.078	0.886	0.934

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	69796	1283
Actual Positive	14	109

-----for Test with threshold 0.6 -----

	precision	recall	roc_auc
accuracy	0.109	0.878	0.933

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	70194	885
Actual Positive	15	108

-----for Test with threshold 0.7 -----

	precision	recall	roc_auc
accuracy	0.156	0.862	0.927

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	70506	573
Actual Positive	17	106

-----for Test with threshold 0.8 -----

	precision	recall	roc_auc
accuracy	0.207	0.862	0.928

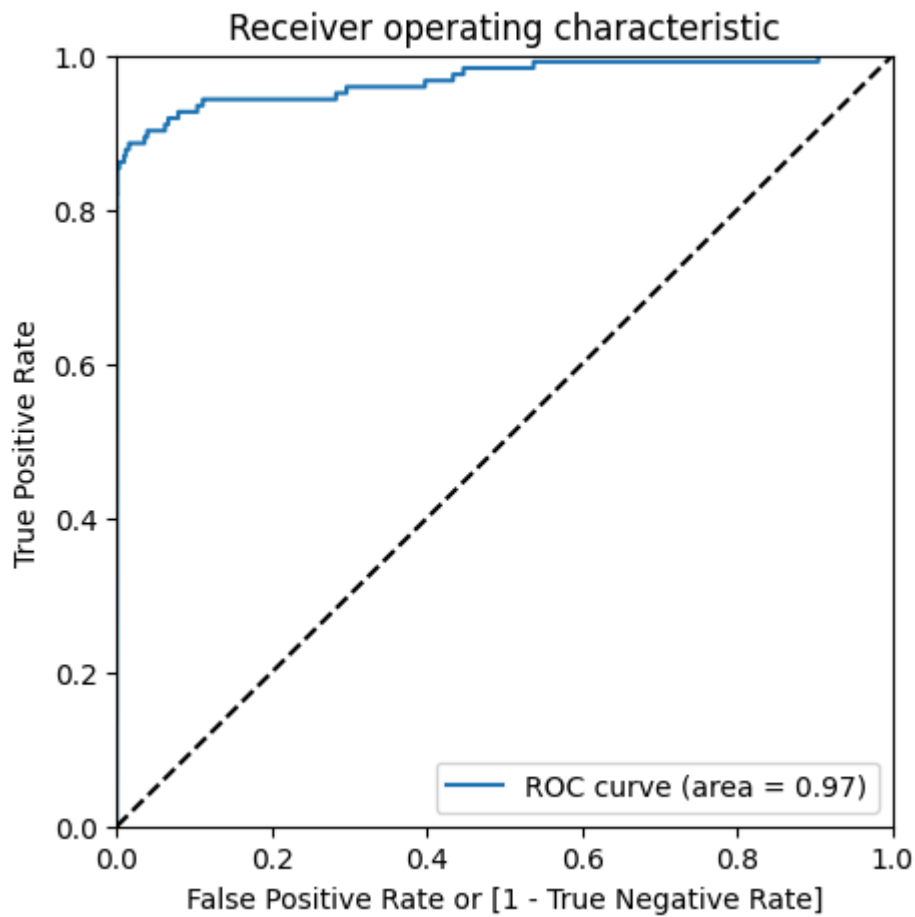
	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	70672	407
Actual Positive	17	106

-----for Test with threshold 0.9 -----

	precision	recall	roc_auc
accuracy	0.268	0.862	0.929

	CONFUSION MATRIX	
	Predicted Negative	Predicted Positive
Actual Negative	70789	290
Actual Positive	17	106

BEST ROC AUC SCORE is 0.934 at the threshold 0.5



```
In [ ]: 1 results.sort_values("recall", ascending=False, inplace=True)
        2 # results.to_csv("results.csv", index = False) #for saving the data
          3 results
```

Out[30]:

	model_name	threshold	recall	roc_auc_score
0	LOGISTIC REGRESSION (SMOTE)	0.5	0.886	0.934000
0	RANDOM FOREST	0.1	0.846	0.922000
0	XGBOOST	0.1	0.829	0.914000
0	RANDOM FOREST	NaN	None	0.968519

Type *Markdown* and LaTeX: α^2