

credit_card_fraud_detection

May 9, 2025

0.0.1 Problem statement:-

The aim of the project is to predict fraudulent credit card transactions using machine learning models. This is crucial from the bank's as well as customer's perspective. The banks cannot afford to lose their customers' money to fraudsters. Every fraud is a loss to the bank as the bank is responsible for the fraud transactions.

The dataset contains transactions made over a period of two days in September 2013 by European credit cardholders. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions. We need to take care of the data imbalance while building the model and come up with the best model by trying various algorithms.

0.1 Steps:-

The steps are broadly divided into below steps. The sub steps are also listed while we approach each of the steps. 1. Reading, understanding and visualising the data 2. Preparing the data for modelling 3. Building the model 4. Evaluate the model

```
[2]: # This was used while running the model in Google Colab
      # from google.colab import drive
      # drive.mount('/content/drive')
```

```
[1]: # Importing the libraries
      import pandas as pd
      import numpy as np

      import matplotlib.pyplot as plt
      %matplotlib inline
      import seaborn as sns

      import warnings
      warnings.filterwarnings('ignore')
```

```
[2]: pd.set_option('display.max_columns', 500)
```

1 Exploratory data analysis

1.1 Reading and understanding the data

```
[3]: # Reading the dataset
```

```
df = pd.read_csv('creditcard.csv')
df.head()
```

```
[3]:
```

	Time	V1	V2	V3	V4	V5	V6	V7 \
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941

	V8	V9	V10	V11	V12	V13	V14 \
0	0.098698	0.363787	0.090794	-0.551600	-0.617801	-0.991390	-0.311169
1	0.085102	-0.255425	-0.166974	1.612727	1.065235	0.489095	-0.143772
2	0.247676	-1.514654	0.207643	0.624501	0.066084	0.717293	-0.165946
3	0.377436	-1.387024	-0.054952	-0.226487	0.178228	0.507757	-0.287924
4	-0.270533	0.817739	0.753074	-0.822843	0.538196	1.345852	-1.119670

	V15	V16	V17	V18	V19	V20	V21 \
0	1.468177	-0.470401	0.207971	0.025791	0.403993	0.251412	-0.018307
1	0.635558	0.463917	-0.114805	-0.183361	-0.145783	-0.069083	-0.225775
2	2.345865	-2.890083	1.109969	-0.121359	-2.261857	0.524980	0.247998
3	-0.631418	-1.059647	-0.684093	1.965775	-1.232622	-0.208038	-0.108300
4	0.175121	-0.451449	-0.237033	-0.038195	0.803487	0.408542	-0.009431

	V22	V23	V24	V25	V26	V27	V28 \
0	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.021053
1	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.014724
2	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055353	-0.059752
3	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.061458
4	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.215153

	Amount	Class
0	149.62	0
1	2.69	0
2	378.66	0
3	123.50	0
4	69.99	0

```
[4]: df.shape
```

```
[4]: (284807, 31)
```

```
[5]: df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Time        284807 non-null float64
1   V1          284807 non-null float64
2   V2          284807 non-null float64
3   V3          284807 non-null float64
4   V4          284807 non-null float64
5   V5          284807 non-null float64
6   V6          284807 non-null float64
7   V7          284807 non-null float64
8   V8          284807 non-null float64
9   V9          284807 non-null float64
10  V10         284807 non-null float64
11  V11         284807 non-null float64
12  V12         284807 non-null float64
13  V13         284807 non-null float64
14  V14         284807 non-null float64
15  V15         284807 non-null float64
16  V16         284807 non-null float64
17  V17         284807 non-null float64
18  V18         284807 non-null float64
19  V19         284807 non-null float64
20  V20         284807 non-null float64
21  V21         284807 non-null float64
22  V22         284807 non-null float64
23  V23         284807 non-null float64
24  V24         284807 non-null float64
25  V25         284807 non-null float64
26  V26         284807 non-null float64
27  V27         284807 non-null float64
28  V28         284807 non-null float64
29  Amount      284807 non-null float64
30  Class       284807 non-null int64
dtypes: float64(30), int64(1)
memory usage: 67.4 MB

```

```
[6]: df.describe()
```

```

[6]:
      count      Time      V1      V2      V3      V4  \
count  284807.000000  2.848070e+05  2.848070e+05  2.848070e+05  2.848070e+05
mean    94813.859575  1.168375e-15  3.416908e-16 -1.379537e-15  2.074095e-15
std     47488.145955  1.958696e+00  1.651309e+00  1.516255e+00  1.415869e+00
min         0.000000 -5.640751e+01 -7.271573e+01 -4.832559e+01 -5.683171e+00
25%     54201.500000 -9.203734e-01 -5.985499e-01 -8.903648e-01 -8.486401e-01

```

50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01

	V5	V6	V7	V8	V9 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	9.604066e-16	1.487313e-15	-5.556467e-16	1.213481e-16	-2.406331e-15
std	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00	1.098632e+00
min	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01	-1.343407e+01
25%	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01	-6.430976e-01
50%	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02	-5.142873e-02
75%	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01	5.971390e-01
max	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01	1.559499e+01

	V10	V11	V12	V13	V14 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	2.239053e-15	1.673327e-15	-1.247012e-15	8.190001e-16	1.207294e-15
std	1.088850e+00	1.020713e+00	9.992014e-01	9.952742e-01	9.585956e-01
min	-2.458826e+01	-4.797473e+00	-1.868371e+01	-5.791881e+00	-1.921433e+01
25%	-5.354257e-01	-7.624942e-01	-4.055715e-01	-6.485393e-01	-4.255740e-01
50%	-9.291738e-02	-3.275735e-02	1.400326e-01	-1.356806e-02	5.060132e-02
75%	4.539234e-01	7.395934e-01	6.182380e-01	6.625050e-01	4.931498e-01
max	2.374514e+01	1.201891e+01	7.848392e+00	7.126883e+00	1.052677e+01

	V15	V16	V17	V18	V19 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	4.887456e-15	1.437716e-15	-3.772171e-16	9.564149e-16	1.039917e-15
std	9.153160e-01	8.762529e-01	8.493371e-01	8.381762e-01	8.140405e-01
min	-4.498945e+00	-1.412985e+01	-2.516280e+01	-9.498746e+00	-7.213527e+00
25%	-5.828843e-01	-4.680368e-01	-4.837483e-01	-4.988498e-01	-4.562989e-01
50%	4.807155e-02	6.641332e-02	-6.567575e-02	-3.636312e-03	3.734823e-03
75%	6.488208e-01	5.232963e-01	3.996750e-01	5.008067e-01	4.589494e-01
max	8.877742e+00	1.731511e+01	9.253526e+00	5.041069e+00	5.591971e+00

	V20	V21	V22	V23	V24 \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	6.406204e-16	1.654067e-16	-3.568593e-16	2.578648e-16	4.473266e-15
std	7.709250e-01	7.345240e-01	7.257016e-01	6.244603e-01	6.056471e-01
min	-5.449772e+01	-3.483038e+01	-1.093314e+01	-4.480774e+01	-2.836627e+00
25%	-2.117214e-01	-2.283949e-01	-5.423504e-01	-1.618463e-01	-3.545861e-01
50%	-6.248109e-02	-2.945017e-02	6.781943e-03	-1.119293e-02	4.097606e-02
75%	1.330408e-01	1.863772e-01	5.285536e-01	1.476421e-01	4.395266e-01
max	3.942090e+01	2.720284e+01	1.050309e+01	2.252841e+01	4.584549e+00

	V25	V26	V27	V28	Amount \
count	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	284807.000000
mean	5.340915e-16	1.683437e-15	-3.660091e-16	-1.227390e-16	88.349619

std	5.212781e-01	4.822270e-01	4.036325e-01	3.300833e-01	250.120109
min	-1.029540e+01	-2.604551e+00	-2.256568e+01	-1.543008e+01	0.000000
25%	-3.171451e-01	-3.269839e-01	-7.083953e-02	-5.295979e-02	5.600000
50%	1.659350e-02	-5.213911e-02	1.342146e-03	1.124383e-02	22.000000
75%	3.507156e-01	2.409522e-01	9.104512e-02	7.827995e-02	77.165000
max	7.519589e+00	3.517346e+00	3.161220e+01	3.384781e+01	25691.160000

	Class
count	284807.000000
mean	0.001727
std	0.041527
min	0.000000
25%	0.000000
50%	0.000000
75%	0.000000
max	1.000000

1.2 Handling missing values

Handling missing values in columns

```
[7]: # Cheking percent of missing values in columns
df_missing_columns = (round(((df.isnull()).sum()/len(df.index))*100),2).
    to_frame('null').sort_values('null', ascending=False)
df_missing_columns
```

```
[7]:      null
Time    0.0
V1      0.0
V2      0.0
V3      0.0
V4      0.0
V5      0.0
V6      0.0
V7      0.0
V8      0.0
V9      0.0
V10     0.0
V11     0.0
V12     0.0
V13     0.0
V14     0.0
V15     0.0
V16     0.0
V17     0.0
V18     0.0
V19     0.0
V20     0.0
```

```
V21      0.0
V22      0.0
V23      0.0
V24      0.0
V25      0.0
V26      0.0
V27      0.0
V28      0.0
Amount   0.0
Class    0.0
```

We can see that there is no missing values in any of the columns. Hence, there is no problem with null values in the entire dataset.

1.2.1 Checking the distribution of the classes

```
[8]: classes = df['Class'].value_counts()
      classes
```

```
[8]: Class
0     284315
1         492
Name: count, dtype: int64
```

```
[9]: normal_share = round((classes[0]/df['Class'].count()*100),2)
      normal_share
```

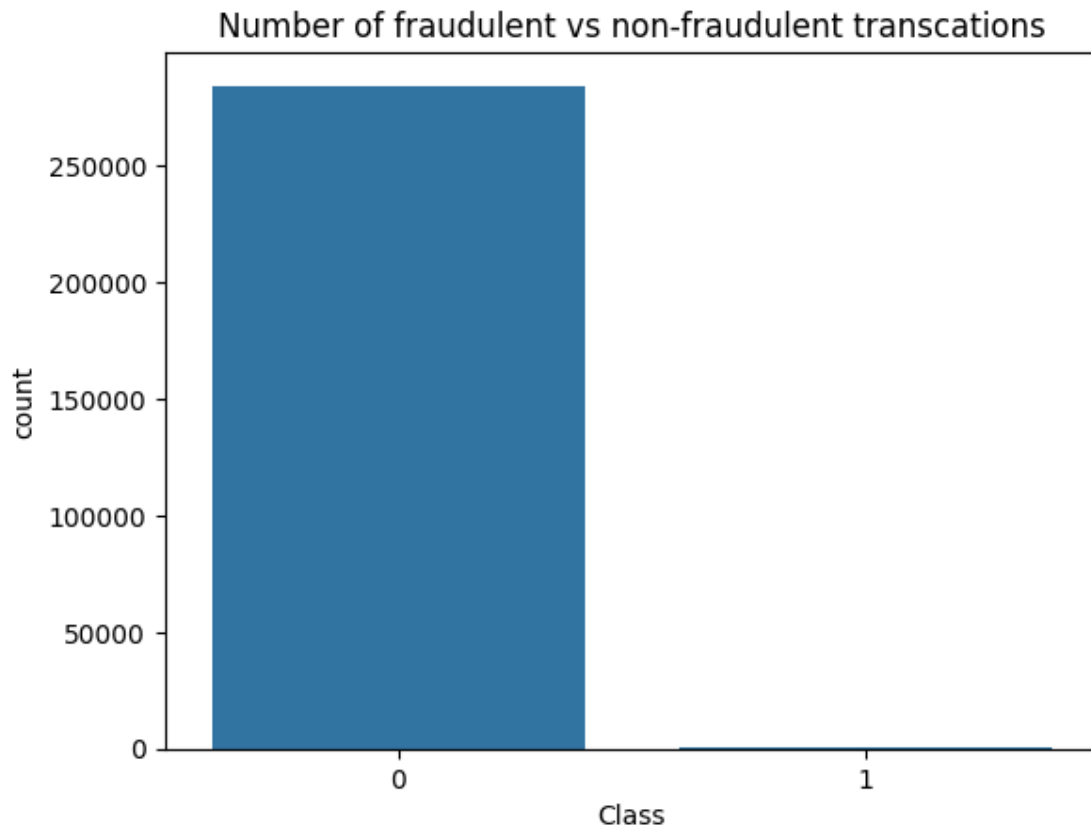
```
[9]: np.float64(99.83)
```

```
[10]: fraud_share = round((classes[1]/df['Class'].count()*100),2)
      fraud_share
```

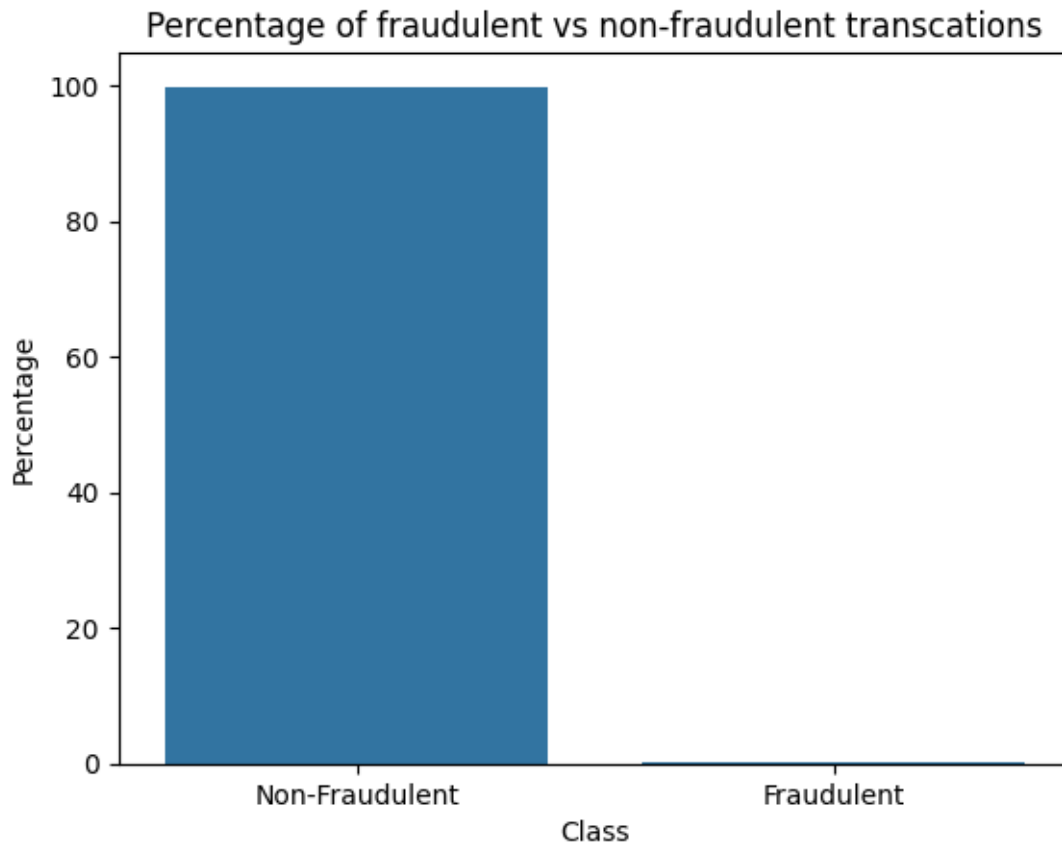
```
[10]: np.float64(0.17)
```

We can see that there is only 0.17% frauds. We will take care of the class imbalance later.

```
[11]: # Bar plot for the number of fraudulent vs non-fraudulent transctions
      sns.countplot(x='Class', data=df)
      plt.title('Number of fraudulent vs non-fraudulent transctions')
      plt.show()
```



```
[12]: # Bar plot for the percentage of fraudulent vs non-fraudulent transctions
fraud_percentage = {'Class':['Non-Fraudulent', 'Fraudulent'], 'Percentage':
    ↳ [normal_share, fraud_share]}
df_fraud_percentage = pd.DataFrame(fraud_percentage)
sns.barplot(x='Class',y='Percentage', data=df_fraud_percentage)
plt.title('Percentage of fraudulent vs non-fraudulent transctions')
plt.show()
```



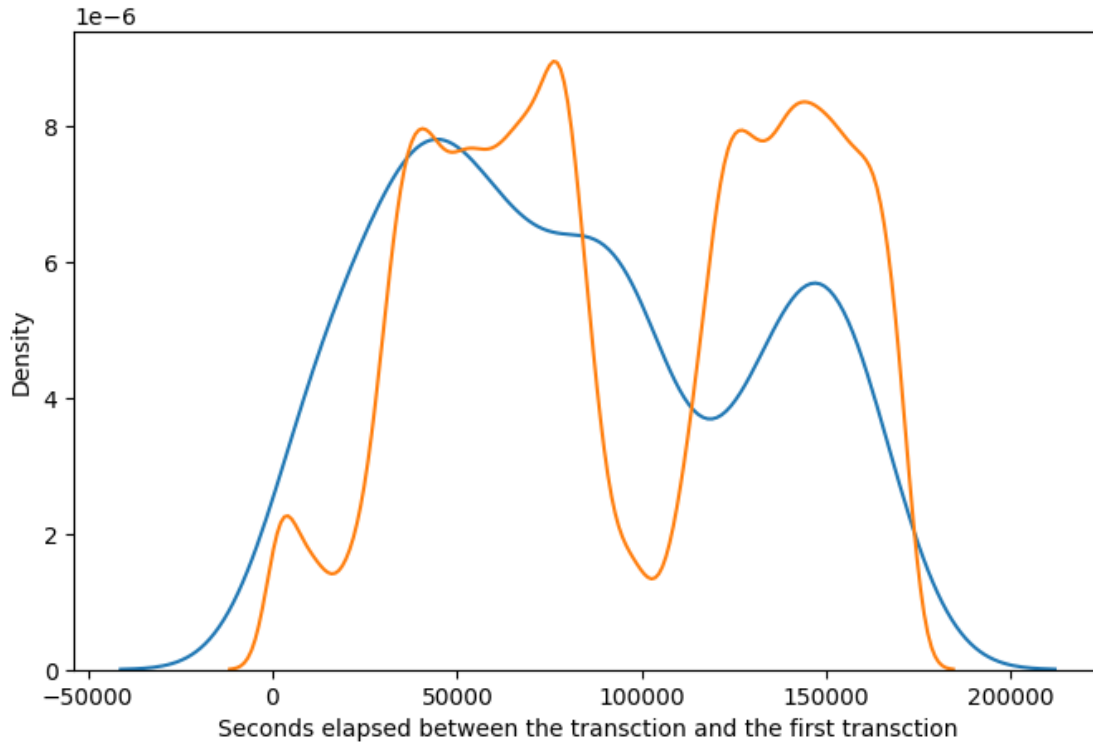
1.3 Outliers treatment

We are not performing any outliers treatment for this particular dataset. Because all the columns are already PCA transformed, which assumed that the outlier values are taken care while transforming the data.

1.3.1 Observe the distribution of classes with time

```
[13]: # Creating fraudulent dataframe
data_fraud = df[df['Class'] == 1]
# Creating non fraudulent dataframe
data_non_fraud = df[df['Class'] == 0]

[14]: # Distribution plot
plt.figure(figsize=(8,5))
ax = sns.distplot(data_fraud['Time'],label='fraudulent',hist=False)
ax = sns.distplot(data_non_fraud['Time'],label='non fraudulent',hist=False)
ax.set(xlabel='Seconds elapsed between the transction and the first transction')
plt.show()
```

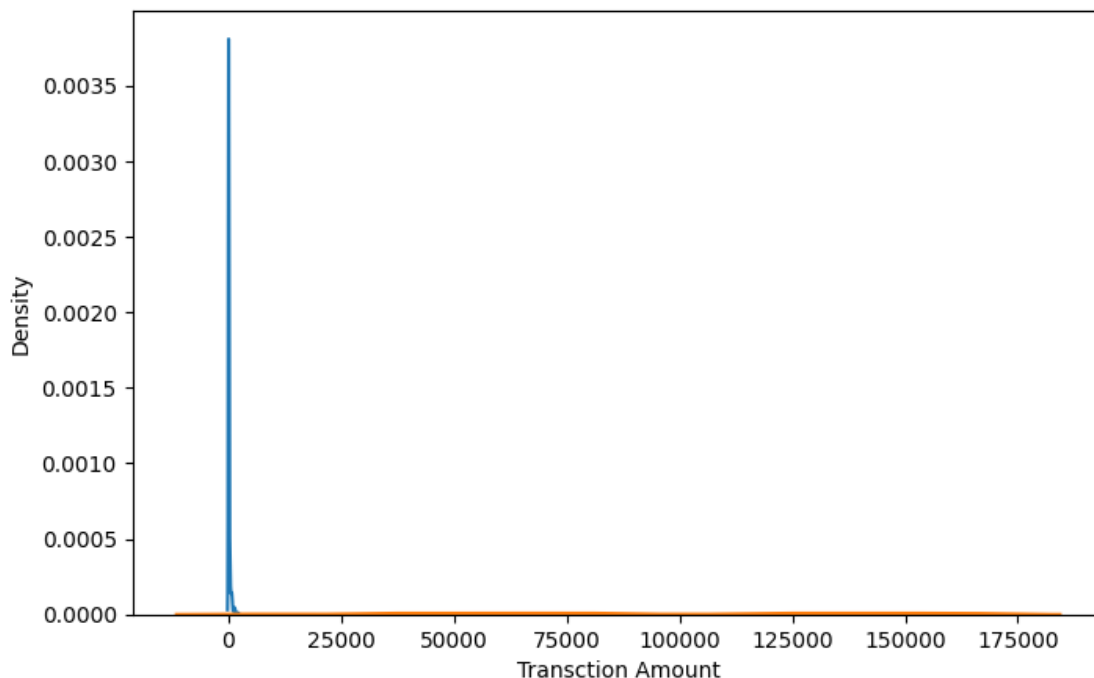



Analysis We do not see any specific pattern for the fraudulent and non-fraudulent transactions with respect to Time. Hence, we can drop the Time column.

```
[15]: # Dropping the Time column
df.drop('Time', axis=1, inplace=True)
```

1.3.2 Observe the distribution of classes with amount

```
[16]: # Distribution plot
plt.figure(figsize=(8,5))
ax = sns.distplot(data_fraud['Amount'],label='fraudulent',hist=False)
ax = sns.distplot(data_non_fraud['Amount'],label='non fraudulent',hist=False)
ax.set(xlabel='Transaction Amount')
plt.show()
```



Analysis We can see that the fraudulent transactions are mostly dense in the lower range of amount, whereas the non-fraudulent transactions are spreaded throughout low to high range of amount.

1.4 Train-Test Split

```
[19]: # Import library
      from sklearn.model_selection import train_test_split
```

```
[20]: # Putting feature variables into X
      X = df.drop(['Class'], axis=1)
```

```
[21]: # Putting target variable to y
      y = df['Class']
```

```
[22]: # Splitting data into train and test set 80:20
      X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8,
      ↪test_size=0.2, random_state=100)
```

1.5 Feature Scaling

We need to scale only the Amount column as all other columns are already scaled by the PCA transformation.

```
[23]: # Standardization method
from sklearn.preprocessing import StandardScaler
```

```
[24]: # Instantiate the Scaler
scaler = StandardScaler()
```

```
[25]: # Fit the data into scaler and transform
X_train['Amount'] = scaler.fit_transform(X_train[['Amount']])
```

```
[26]: X_train.head()
```

```
[26]:
```

	V1	V2	V3	V4	V5	V6	V7	\
201788	2.023734	-0.429219	-0.691061	-0.201461	-0.162486	0.283718	-0.674694	
179369	-0.145286	0.736735	0.543226	0.892662	0.350846	0.089253	0.626708	
73138	-3.015846	-1.920606	1.229574	0.721577	1.089918	-0.195727	-0.462586	
208679	1.851980	-1.007445	-1.499762	-0.220770	-0.568376	-1.232633	0.248573	
206534	2.237844	-0.551513	-1.426515	-0.924369	-0.401734	-1.438232	-0.119942	

	V8	V9	V10	V11	V12	V13	V14	\
201788	0.192230	1.124319	-0.037763	0.308648	0.875063	-0.009562	0.116038	
179369	-0.049137	-0.732566	0.297692	0.519027	0.041275	-0.690783	0.647121	
73138	0.919341	-0.612193	-0.966197	1.106534	1.026421	-0.474229	0.641488	
208679	-0.539483	-0.813368	0.785431	-0.784316	0.673626	1.428269	0.043937	
206534	-0.449263	-0.717258	0.851668	-0.497634	-0.445482	0.324575	0.125543	

	V15	V16	V17	V18	V19	V20	V21	\
201788	0.086537	0.628337	-0.997868	0.482547	0.576077	-0.171390	-0.195207	
179369	0.526333	-1.098558	0.511739	0.243984	3.349611	0.206709	-0.124288	
73138	-0.430684	-0.631257	0.634633	-0.718062	-0.039929	0.842838	0.274911	
208679	-0.309507	-1.805728	-0.012118	0.377096	-0.658353	-0.196551	-0.406722	
206534	0.266588	0.802640	0.225312	-1.865494	0.621879	-0.045417	0.050447	

	V22	V23	V24	V25	V26	V27	V28	\
201788	-0.477813	0.340513	0.059174	-0.431015	-0.297028	-0.000063	-0.046947	
179369	-0.263560	-0.110568	-0.434224	-0.509076	0.719784	-0.006357	0.146053	
73138	-0.319550	0.212891	-0.268792	0.241190	0.318445	-0.100726	-0.365257	
208679	-0.899081	0.137370	0.075894	-0.244027	0.455618	-0.094066	-0.031488	
206534	0.125601	0.215531	-0.080485	-0.063975	-0.307176	-0.042838	-0.063872	

	Amount
201788	-0.345273
179369	-0.206439
73138	0.358043
208679	0.362400
206534	-0.316109

Scaling the test set We don't fit scaler on the test set. We only transform the test set.

```
[27]: # Transform the test set
X_test['Amount'] = scaler.transform(X_test[['Amount']])
X_test.head()
```

```
[27]:
```

	V1	V2	V3	V4	V5	V6	V7	\
49089	1.229452	-0.235478	-0.627166	0.419877	1.797014	4.069574	-0.896223	
154704	2.016893	-0.088751	-2.989257	-0.142575	2.675427	3.332289	-0.652336	
67247	0.535093	-1.469185	0.868279	0.385462	-1.439135	0.368118	-0.499370	
251657	2.128486	-0.117215	-1.513910	0.166456	0.359070	-0.540072	0.116023	
201903	0.558593	1.587908	-2.368767	5.124413	2.171788	-0.500419	1.059829	

	V8	V9	V10	V11	V12	V13	V14	\
49089	1.036103	0.745991	-0.147304	-0.850459	0.397845	-0.259849	-0.277065	
154704	0.752811	1.962566	-1.025024	1.126976	-2.418093	1.250341	-0.056209	
67247	0.303698	1.042073	-0.437209	1.145725	0.907573	-1.095634	-0.055080	
251657	-0.216140	0.680314	0.079977	-1.705327	-0.127579	-0.207945	0.307878	
201903	-0.254233	-1.959060	0.948915	-0.288169	-1.007647	0.470316	-2.771902	

	V15	V16	V17	V18	V19	V20	V21	\
49089	-0.766810	-0.200946	-0.338122	0.006032	0.477431	-0.057922	-0.170060	
154704	-0.736695	0.014783	1.890249	0.333755	-0.450398	-0.147619	-0.184153	
67247	-0.621880	-0.191066	0.311988	-0.478635	0.231159	0.437685	0.028010	
251657	0.213491	0.163032	-0.587029	-0.561292	0.472667	-0.227278	-0.357993	
201903	0.221958	0.354333	2.603189	1.092576	0.668084	0.249457	-0.035049	

	V22	V23	V24	V25	V26	V27	V28	\
49089	-0.288750	-0.130270	1.025935	0.847990	-0.271476	0.060052	0.018104	
154704	-0.089661	0.087188	0.570679	0.101899	0.620842	-0.048958	-0.042831	
67247	-0.384708	-0.128376	0.286638	-0.136700	0.913904	-0.083364	0.052485	
251657	-0.905085	0.223474	-1.075605	-0.188519	0.267672	-0.071733	-0.072238	
201903	0.271455	0.381606	0.332001	-0.334757	0.448890	0.168585	0.004955	

	Amount
49089	-0.340485
154704	-0.320859
67247	0.853442
251657	-0.344410
201903	-0.229480

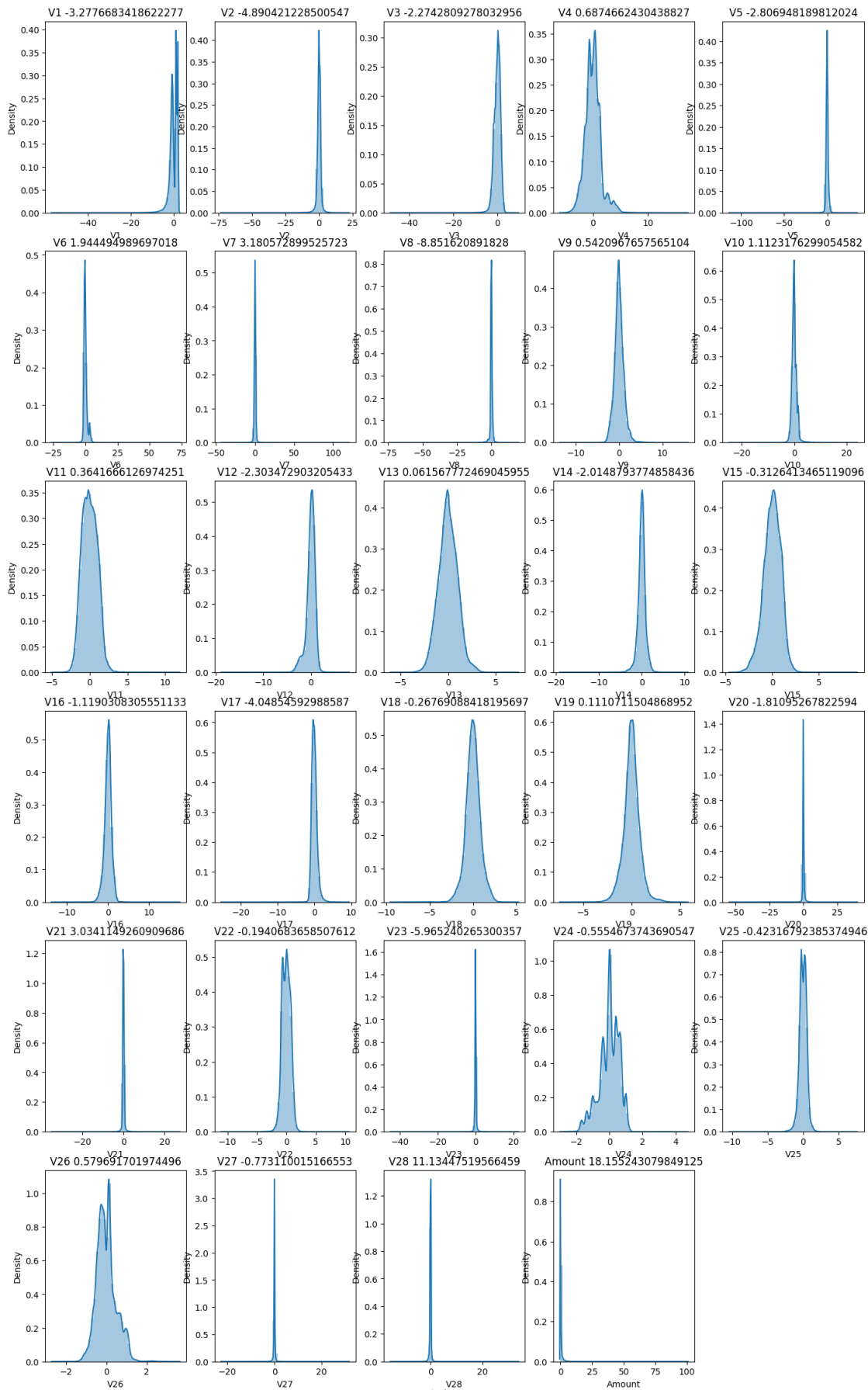
1.6 Checking the Skewness

```
[28]: # Listing the columns
cols = X_train.columns
cols
```

```
[28]: Index(['V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11',
          'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20', 'V21',
```

```
    'V22', 'V23', 'V24', 'V25', 'V26', 'V27', 'V28', 'Amount'],  
    dtype='object')
```

```
[29]: # Plotting the distribution of the variables (skewness) of all the columns  
k=0  
plt.figure(figsize=(17,28))  
for col in cols :  
    k=k+1  
    plt.subplot(6, 5,k)  
    sns.distplot(X_train[col])  
    plt.title(col+' '+str(X_train[col].skew()))
```



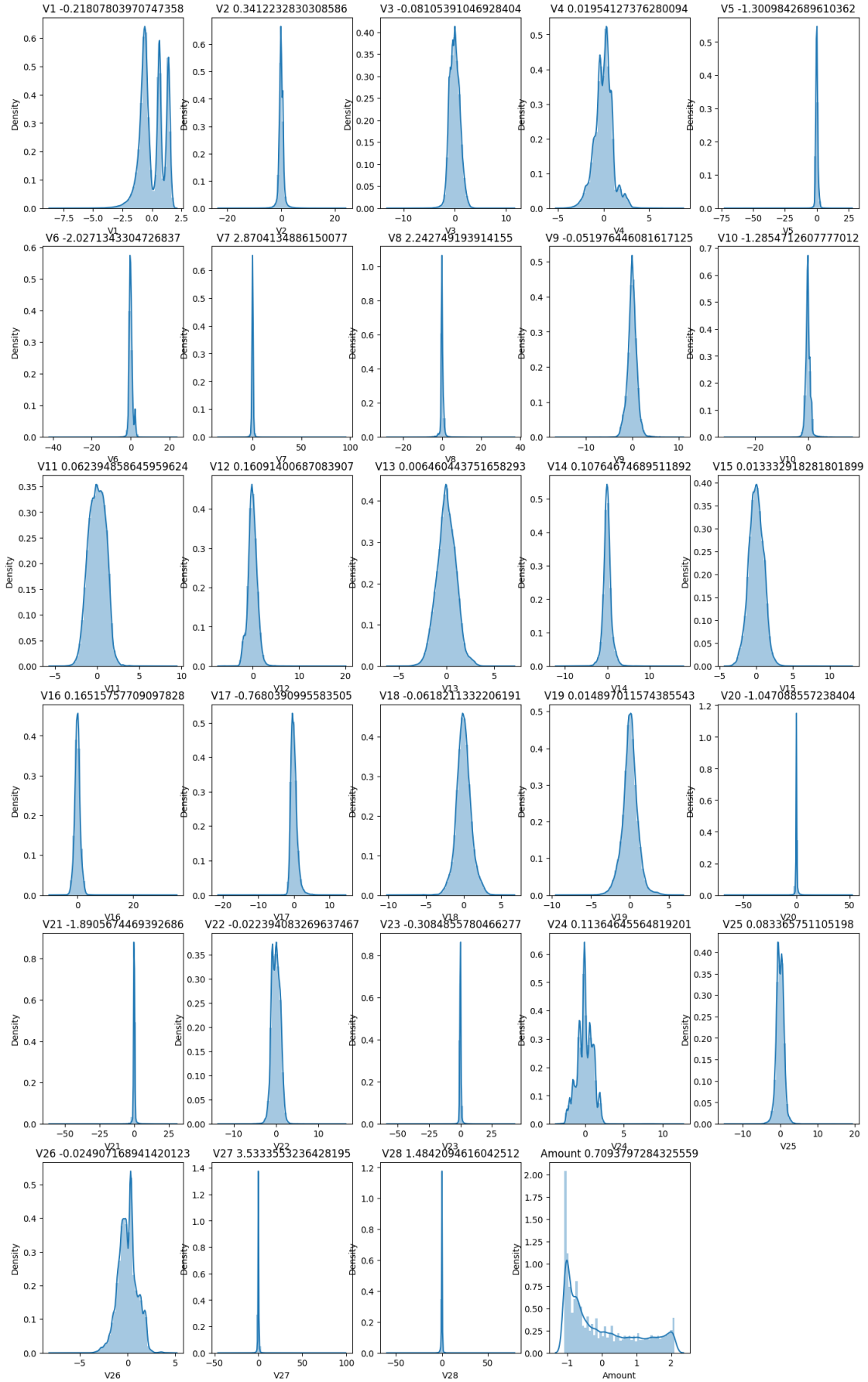
We see that there are many variables, which are heavily skewed. We will mitigate the skewness only for those variables for bringing them into normal distribution.

1.6.1 Mitigate skewness with PowerTransformer

```
[30]: # Importing PowerTransformer
from sklearn.preprocessing import PowerTransformer
# Instantiate the powertransformer
pt = PowerTransformer(method='yeo-johnson', standardize=True, copy=False)
# Fit and transform the PT on training data
X_train[cols] = pt.fit_transform(X_train)

[32]: # Transform the test set
X_test[cols] = pt.transform(X_test)

[34]: # Plotting the distribution of the variables (skewness) of all the columns
k=0
plt.figure(figsize=(17,28))
for col in cols :
    k=k+1
    plt.subplot(6, 5,k)
    sns.distplot(X_train[col])
    plt.title(col+' '+str(X_train[col].skew()))
```



Now we can see that all the variables are normally distributed after the transformation.

2 Model building on imbalanced data

2.0.1 Metric selection for heavily imbalanced data

As we have seen that the data is heavily imbalanced, where only 0.17% transactions are fraudulent, we should not consider Accuracy as a good measure for evaluating the model. Because in the case of all the datapoints return a particular class(1/0) irrespective of any prediction, still the model will result more than 99% Accuracy.

Hence, we have to measure the ROC-AUC score for fair evaluation of the model. The ROC curve is used to understand the strength of the model by evaluating the performance of the model at all the classification thresholds. The default threshold of 0.5 is not always the ideal threshold to find the best classification label of the test point. Because the ROC curve is measured at all thresholds, the best threshold would be one at which the TPR is high and FPR is low, i.e., misclassifications are low. After determining the optimal threshold, we can calculate the F1 score of the classifier to measure the precision and recall at the selected threshold.

Why SVM was not tried for model building and Random Forest was not tried for few cases? In the dataset we have 284807 datapoints and in the case of Oversampling we would have even more number of datapoints. SVM is not very efficient with large number of datapoints because it takes lot of computational power and resources to make the transformation. When we perform the cross validation with K-Fold for hyperparameter tuning, it takes lot of computational resources and it is very time consuming. Hence, because of the unavailability of the required resources and time SVM was not tried.

For the same reason Random forest was also not tried for model building in few of the hyperparameter tuning for oversampling technique.

Why KNN was not used for model building? KNN is not memory efficient. It becomes very slow as the number of datapoints increases as the model needs to store all the data points. It is computationally heavy because for a single datapoint the algorithm has to calculate the distance of all the datapoints and find the nearest neighbors.

2.0.2 Logistic regression

```
[35]: # Importing scikit logistic regression module
      from sklearn.linear_model import LogisticRegression
```

```
[36]: # Importing metrics
      from sklearn import metrics
      from sklearn.metrics import confusion_matrix
      from sklearn.metrics import f1_score
      from sklearn.metrics import classification_report
```

Tuning hyperparameter C C is the the inverse of regularization strength in Logistic Regression. Higher values of C correspond to less regularization.

```
[37]: # Importing libraries for cross validation
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV

[39]: # Creating KFold object with 5 splits
folds = KFold(n_splits=5, shuffle=True, random_state=4)

# Specify params
params = {"C": [0.01, 0.1, 1, 10, 100, 1000]}

# Specifying score as recall as we are more focused on acheiving the higher
↪sensitivity than the accuracy
model_cv = GridSearchCV(estimator = LogisticRegression(),
                        param_grid = params,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# Fit the model
model_cv.fit(X_train, y_train)
```

Fitting 5 folds for each of 6 candidates, totalling 30 fits

```
[39]: GridSearchCV(cv=KFold(n_splits=5, random_state=4, shuffle=True),
                  estimator=LogisticRegression(),
                  param_grid={'C': [0.01, 0.1, 1, 10, 100, 1000]},
                  return_train_score=True, scoring='roc_auc', verbose=1)
```

```
[41]: # results of grid search CV
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[41]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  param_C  \
0      0.579831    0.085067      0.044459      0.015167      0.01
1      0.671977    0.051240      0.055647      0.007821     0.10
2      0.550463    0.058531      0.031705      0.005627      1.00
3      0.566925    0.093087      0.040286      0.007552     10.00
4      0.517871    0.072010      0.030372      0.010828    100.00
5      0.541186    0.061354      0.028787      0.004566   1000.00

      params  split0_test_score  split1_test_score  split2_test_score  \
0  {'C': 0.01}             0.986595             0.987068             0.969244
1  {'C': 0.1}              0.985593             0.987368             0.966190
```

2	{'C': 1}	0.985601	0.987346	0.960695
3	{'C': 10}	0.985580	0.987338	0.961110
4	{'C': 100}	0.985578	0.987338	0.959647
5	{'C': 1000}	0.985578	0.987338	0.959637

	split3_test_score	split4_test_score	mean_test_score	std_test_score \
0	0.981472	0.993990	0.983674	0.008241
1	0.980005	0.994159	0.982663	0.009395
2	0.979551	0.994229	0.981484	0.011399
3	0.979525	0.991787	0.981068	0.010726
4	0.979517	0.991783	0.980772	0.011272
5	0.979519	0.991782	0.980771	0.011275

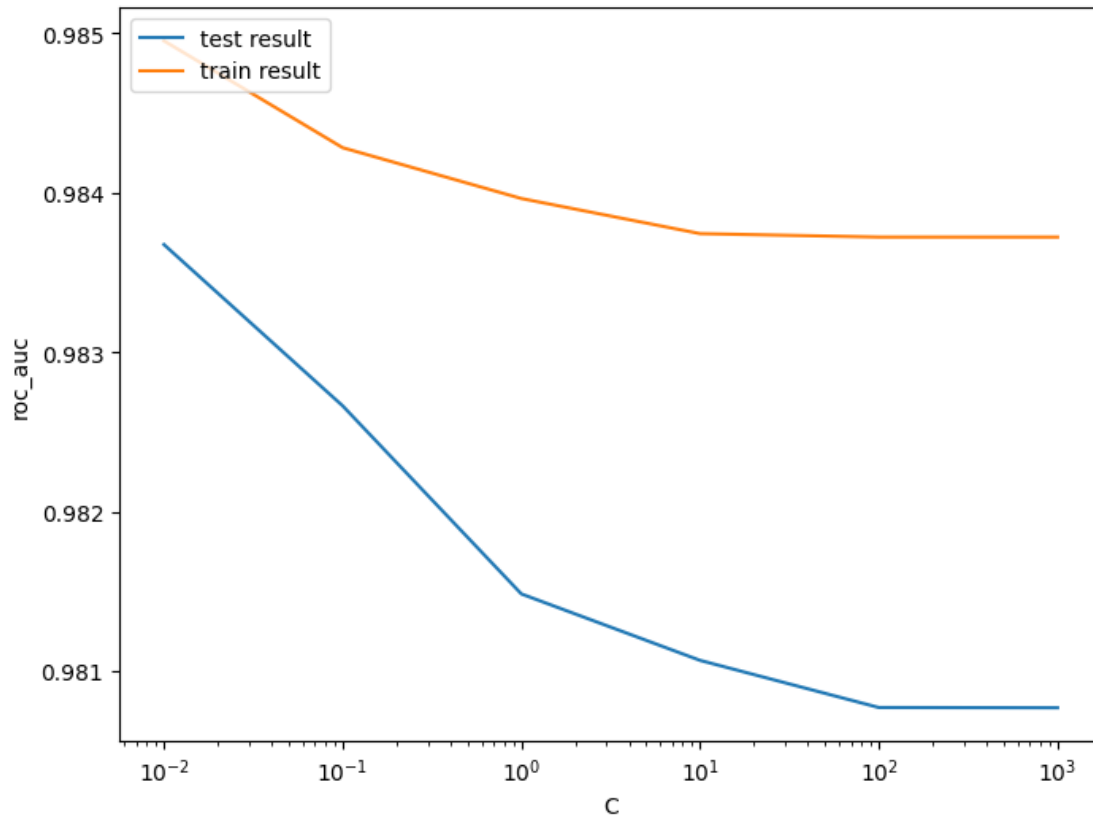
	rank_test_score	split0_train_score	split1_train_score \
0	1	0.983877	0.984106
1	2	0.982962	0.983607
2	3	0.982770	0.983390
3	4	0.982758	0.983365
4	5	0.982757	0.983362
5	6	0.982757	0.983362

	split2_train_score	split3_train_score	split4_train_score \
0	0.988321	0.985739	0.982709
1	0.988169	0.984679	0.981988
2	0.987509	0.984222	0.981921
3	0.987466	0.984354	0.980767
4	0.987354	0.984366	0.980764
5	0.987352	0.984367	0.980764

	mean_train_score	std_train_score
0	0.984950	0.001943
1	0.984281	0.002132
2	0.983962	0.001927
3	0.983742	0.002200
4	0.983721	0.002164
5	0.983721	0.002164

[42]: *# plot of C versus train and validation scores*

```
plt.figure(figsize=(8, 6))
plt.plot(cv_results['param_C'], cv_results['mean_test_score'])
plt.plot(cv_results['param_C'], cv_results['mean_train_score'])
plt.xlabel('C')
plt.ylabel('roc_auc')
plt.legend(['test result', 'train result'], loc='upper left')
plt.xscale('log')
```



```
[43]: # Best score with best C
best_score = model_cv.best_score_
best_C = model_cv.best_params_['C']

print(" The highest test roc_auc is {0} at C = {1}".format(best_score, best_C))
```

The highest test roc_auc is 0.9836736960858568 at C = 0.01

Logistic regression with optimal C

```
[44]: # Instantiate the model with best C
logistic_imb = LogisticRegression(C=0.01)
```

```
[45]: # Fit the model on the train set
logistic_imb_model = logistic_imb.fit(X_train, y_train)
```

Prediction on the train set

```
[46]: # Predictions on the train set
y_train_pred = logistic_imb_model.predict(X_train)
```

```
[47]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train, y_train_pred)
```

```
print(confusion)
```

```
[[227427    22]
 [   138    258]]
```

```
[48]: TP = confusion[1,1] # true positive
      TN = confusion[0,0] # true negatives
      FP = confusion[0,1] # false positives
      FN = confusion[1,0] # false negatives
```

```
[49]: # Accuracy
      print("Accuracy:-",metrics.accuracy_score(y_train, y_train_pred))

      # Sensitivity
      print("Sensitivity:-",TP / float(TP+FN))

      # Specificity
      print("Specificity:-", TN / float(TN+FP))

      # F1 score
      print("F1-Score:-", f1_score(y_train, y_train_pred))
```

Accuracy:- 0.999297768219623

Sensitivity:- 0.6515151515151515

Specificity:- 0.9999032750198946

F1-Score:- 0.7633136094674556

```
[50]: # classification_report
      print(classification_report(y_train, y_train_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	0.92	0.65	0.76	396
accuracy			1.00	227845
macro avg	0.96	0.83	0.88	227845
weighted avg	1.00	1.00	1.00	227845

ROC on the train set

```
[51]: # ROC Curve function

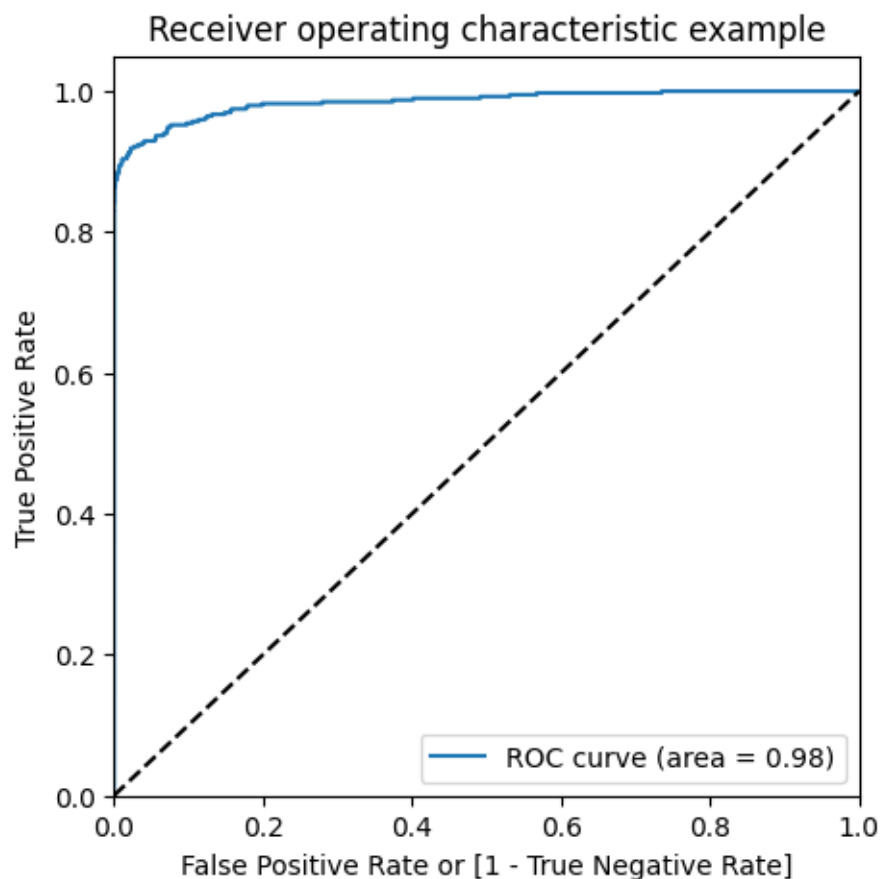
def draw_roc( actual, probs ):
    fpr, tpr, thresholds = metrics.roc_curve( actual, probs,
                                              drop_intermediate = False )
    auc_score = metrics.roc_auc_score( actual, probs )
    plt.figure(figsize=(5, 5))
```

```
plt.plot( fpr, tpr, label='ROC curve (area = %0.2f)' % auc_score )
plt.plot([0, 1], [0, 1], 'k--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate or [1 - True Negative Rate]')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.legend(loc="lower right")
plt.show()

return None
```

```
[52]: # Predicted probability
y_train_pred_proba = logistic_imb_model.predict_proba(X_train)[:,-1]
```

```
[53]: # Plot the ROC curve
draw_roc(y_train, y_train_pred_proba)
```



We achieved very good ROC 0.99 on the train set.

Prediction on the test set

```
[54]: # Prediction on the test set
y_test_pred = logistic_imb_model.predict(X_test)
```

```
[55]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56855   11]
 [   57   39]]
```

```
[56]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[57]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-",TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_test, y_test_pred))
```

```
Accuracy:- 0.9988062216916541
Sensitivity:- 0.40625
Specificity:- 0.9998065627967503
F1-Score:- 0.5342465753424658
```

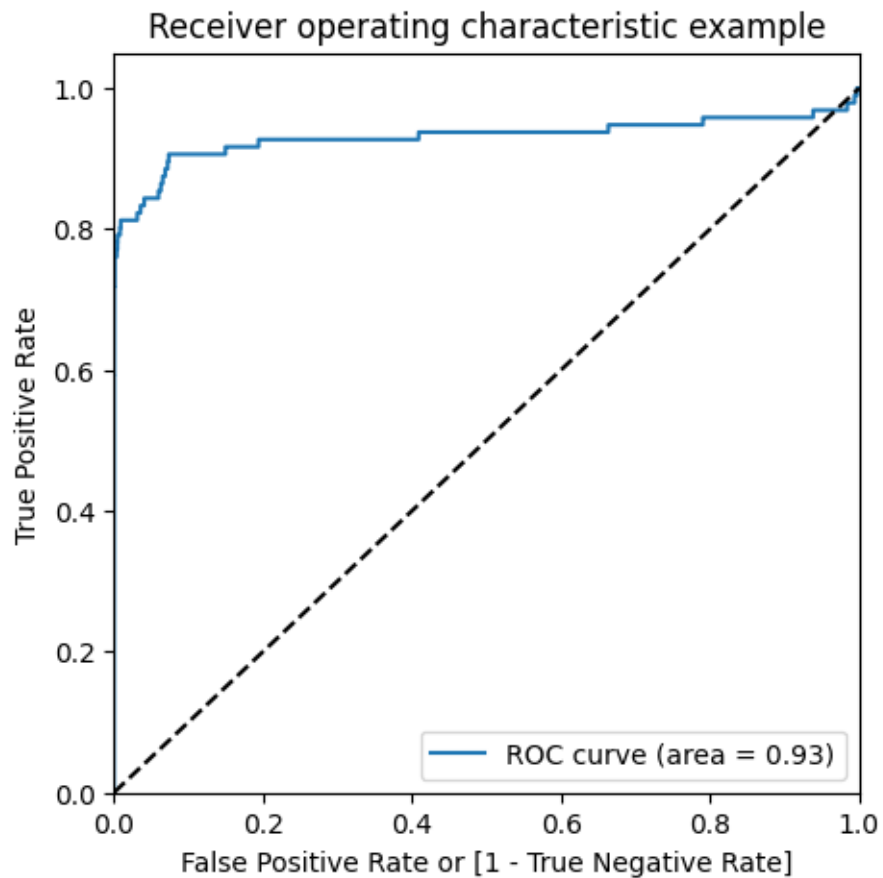
```
[58]: # classification_report
print(classification_report(y_test, y_test_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.78	0.41	0.53	96
accuracy			1.00	56962
macro avg	0.89	0.70	0.77	56962
weighted avg	1.00	1.00	1.00	56962

ROC on the test set

```
[59]: # Predicted probability
y_test_pred_proba = logistic_imb_model.predict_proba(X_test)[: ,1]
```

```
[60]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



We can see that we have very good ROC on the test set 0.97, which is almost close to 1.

Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 0.70
 - Specificity = 0.99
 - F1-Score = 0.76
 - ROC = 0.99
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.77
 - Specificity = 0.99
 - F1-Score = 0.65
 - ROC = 0.97

Overall, the model is performing well in the test set, what it had learnt from the train set.

2.0.3 XGBoost

```
[62]: # Importing XGBoost
      from xgboost import XGBClassifier
```

Tuning the hyperparameters

```
[63]: # hyperparameter tuning with XGBoost

      # creating a KFold object
      folds = 3

      # specify range of hyperparameters
      param_grid = {'learning_rate': [0.2, 0.6],
                    'subsample': [0.3, 0.6, 0.9]}

      # specify model
      xgb_model = XGBClassifier(max_depth=2, n_estimators=200)

      # set up GridSearchCV()
      model_cv = GridSearchCV(estimator = xgb_model,
                              param_grid = param_grid,
                              scoring= 'roc_auc',
                              cv = folds,
                              verbose = 1,
                              return_train_score=True)

      # fit the model
      model_cv.fit(X_train, y_train)
```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

```
[63]: GridSearchCV(cv=3,
                  estimator=XGBClassifier(base_score=None, booster=None,
                                          callbacks=None, colsample_bylevel=None,
                                          colsample_bynode=None,
                                          colsample_bytree=None, device=None,
                                          early_stopping_rounds=None,
                                          enable_categorical=False, eval_metric=None,
                                          feature_types=None, feature_weights=None,
                                          gamma=None, grow_policy=None,
                                          importance_type=None,
                                          interaction_constraints=None,
                                          learning_rate=None, max_bin=None,
                                          max_cat_threshold=None,
                                          max_cat_to_onehot=None,
                                          max_delta_step=None, max_depth=2,
                                          max_leaves=None, min_child_weight=None,
```

```

        missing=nan, monotone_constraints=None,
        multi_strategy=None, n_estimators=200,
        n_jobs=None, num_parallel_tree=None, ...),
    param_grid={'learning_rate': [0.2, 0.6],
                'subsample': [0.3, 0.6, 0.9]},
    return_train_score=True, scoring='roc_auc', verbose=1)

```

```

[64]: # cv results
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results

```

```

[64]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \

0      1.205345    0.379356      0.033322      0.000941
1      0.874265    0.062975      0.032574      0.000085
2      0.775534    0.173189      0.019743      0.002628
3      0.593227    0.014679      0.017233      0.000512
4      0.599941    0.009383      0.017455      0.001003
5      0.628660    0.043908      0.018282      0.000620

   param_learning_rate  param_subsample  \

0                    0.2                0.3
1                    0.2                0.6
2                    0.2                0.9
3                    0.6                0.3
4                    0.6                0.6
5                    0.6                0.9

   params  split0_test_score  \

0 {'learning_rate': 0.2, 'subsample': 0.3}      0.964960
1 {'learning_rate': 0.2, 'subsample': 0.6}      0.954333
2 {'learning_rate': 0.2, 'subsample': 0.9}      0.500000
3 {'learning_rate': 0.6, 'subsample': 0.3}      0.464846
4 {'learning_rate': 0.6, 'subsample': 0.6}      0.500000
5 {'learning_rate': 0.6, 'subsample': 0.9}      0.488406

   split1_test_score  split2_test_score  mean_test_score  std_test_score  \

0      0.966867      0.979262      0.970363      0.006340
1      0.911532      0.933843      0.933236      0.017479
2      0.540591      0.500000      0.513530      0.019135
3      0.500000      0.474165      0.479670      0.014870
4      0.482963      0.500000      0.494321      0.008031
5      0.486019      0.471900      0.482108      0.007284

   rank_test_score  split0_train_score  split1_train_score  \

0                1      0.999691      0.999767
1                2      0.993920      0.968862
2                3      0.500000      0.527574

```

3	6	0.462853	0.500000
4	4	0.500000	0.480886
5	5	0.488261	0.486293

	split2_train_score	mean_train_score	std_train_score
0	0.999219	0.999559	0.000242
1	0.977446	0.980076	0.010397
2	0.500000	0.509191	0.012998
3	0.470746	0.477866	0.015979
4	0.500000	0.493629	0.009010
5	0.468767	0.481107	0.008762

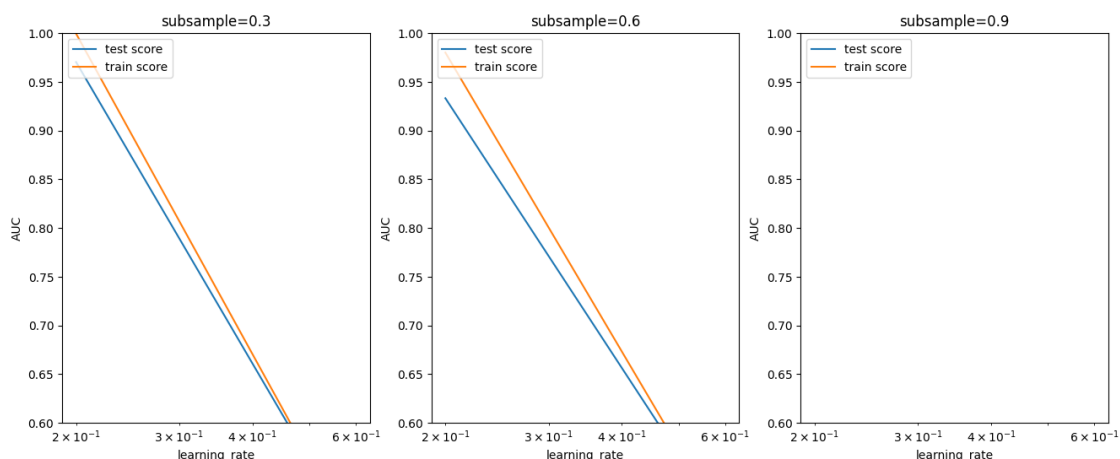
```
[65]: ## plotting
plt.figure(figsize=(16,6))

param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

for n, subsample in enumerate(param_grid['subsample']):

    # subplot 1/n
    plt.subplot(1,len(param_grid['subsample']), n+1)
    df = cv_results[cv_results['param_subsample']==subsample]

    plt.plot(df["param_learning_rate"], df["mean_test_score"])
    plt.plot(df["param_learning_rate"], df["mean_train_score"])
    plt.xlabel('learning_rate')
    plt.ylabel('AUC')
    plt.title("subsample={0}".format(subsample))
    plt.ylim([0.60, 1])
    plt.legend(['test score', 'train score'], loc='upper left')
    plt.xscale('log')
```



Model with optimal hyperparameters We see that the train score almost touches to 1. Among the hyperparameters, we can choose the best parameters as `learning_rate : 0.2` and `subsample: 0.3`

```
[66]: model_cv.best_params_
```

```
[66]: {'learning_rate': 0.2, 'subsample': 0.3}
```

```
[67]: # chosen hyperparameters
# 'objective':'binary:logistic' outputs probability rather than label, which we
      ↪ need for calculating auc
params = {'learning_rate': 0.2,
          'max_depth': 2,
          'n_estimators':200,
          'subsample':0.9,
          'objective':'binary:logistic'}

# fit model on training data
xgb_imb_model = XGBClassifier(params = params)
xgb_imb_model.fit(X_train, y_train)
```

```
[67]: XGBClassifier(base_score=None, booster=None, callbacks=None,
                  colsample_bylevel=None, colsample_bynode=None,
                  colsample_bytree=None, device=None, early_stopping_rounds=None,
                  enable_categorical=False, eval_metric=None, feature_types=None,
                  feature_weights=None, gamma=None, grow_policy=None,
                  importance_type=None, interaction_constraints=None,
                  learning_rate=None, max_bin=None, max_cat_threshold=None,
                  max_cat_to_onehot=None, max_delta_step=None, max_depth=None,
                  max_leaves=None, min_child_weight=None, missing=nan,
                  monotone_constraints=None, multi_strategy=None, n_estimators=None,
                  n_jobs=None, num_parallel_tree=None, ...)
```

Prediction on the train set

```
[68]: # Predictions on the train set
y_train_pred = xgb_imb_model.predict(X_train)
```

```
[69]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train, y_train_pred)
print(confusion)
```

```
[[227449    0]
 [      4   392]]
```

```
[70]: TP = confusion[1,1] # true positive
      TN = confusion[0,0] # true negatives
      FP = confusion[0,1] # false positives
      FN = confusion[1,0] # false negatives
```

```
[71]: # Accuracy
      print("Accuracy:-", metrics.accuracy_score(y_train, y_train_pred))

      # Sensitivity
      print("Sensitivity:-", TP / float(TP+FN))

      # Specificity
      print("Specificity:-", TN / float(TN+FP))

      # F1 score
      print("F1-Score:-", f1_score(y_train, y_train_pred))
```

```
Accuracy:- 0.9999824442054905
Sensitivity:- 0.98989898989899
Specificity:- 1.0
F1-Score:- 0.9949238578680203
```

```
[72]: # classification_report
      print(classification_report(y_train, y_train_pred))
```

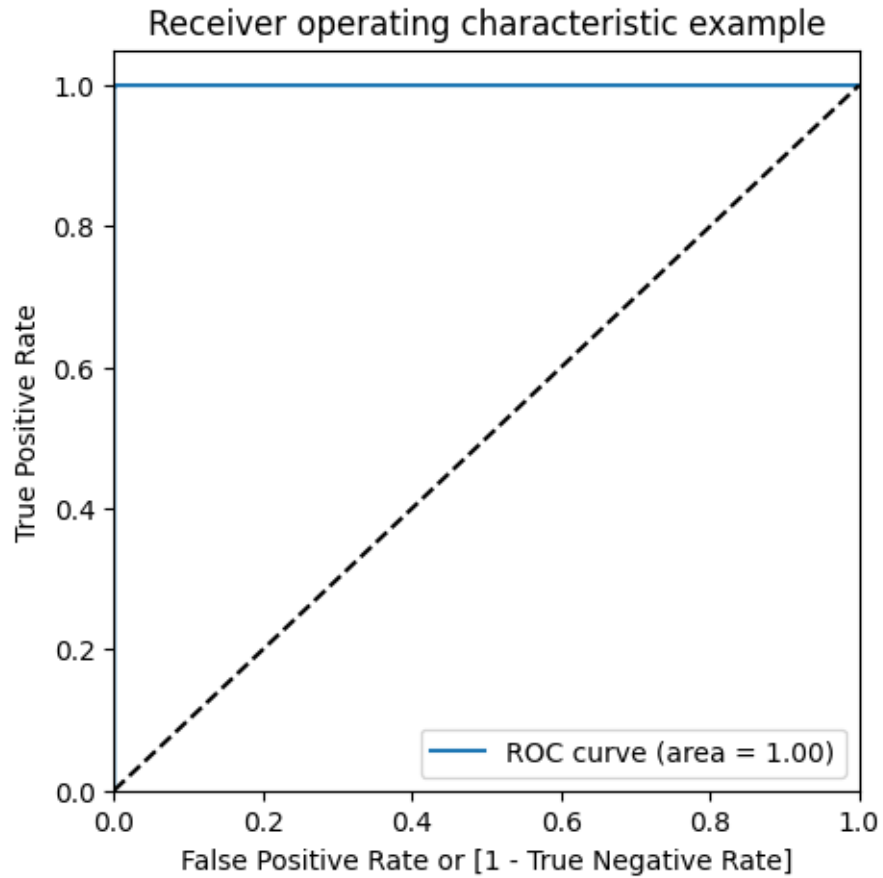
	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	1.00	0.99	0.99	396
accuracy			1.00	227845
macro avg	1.00	0.99	1.00	227845
weighted avg	1.00	1.00	1.00	227845

```
[73]: # Predicted probability
      y_train_pred_proba_imb_xgb = xgb_imb_model.predict_proba(X_train)[: ,1]
```

```
[74]: # roc_auc
      auc = metrics.roc_auc_score(y_train, y_train_pred_proba_imb_xgb)
      auc
```

```
[74]: np.float64(0.999998079267498)
```

```
[75]: # Plot the ROC curve
      draw_roc(y_train, y_train_pred_proba_imb_xgb)
```



Prediction on the test set

```
[76]: # Predictions on the test set
y_test_pred = xgb_imb_model.predict(X_test)
```

```
[77]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56851   15]
 [   28   68]]
```

```
[78]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[79]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))
```

```

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_test, y_test_pred))

```

```

Accuracy:- 0.9992451107756047
Sensitivity:- 0.7083333333333334
Specificity:- 0.9997362219955686
F1-Score:- 0.7597765363128491

```

```

[80]: # classification_report
print(classification_report(y_test, y_test_pred))

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.82	0.71	0.76	96
accuracy			1.00	56962
macro avg	0.91	0.85	0.88	56962
weighted avg	1.00	1.00	1.00	56962

```

[81]: # Predicted probability
y_test_pred_proba = xgb_imb_model.predict_proba(X_test)[: ,1]

```

```

[82]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc

```

```

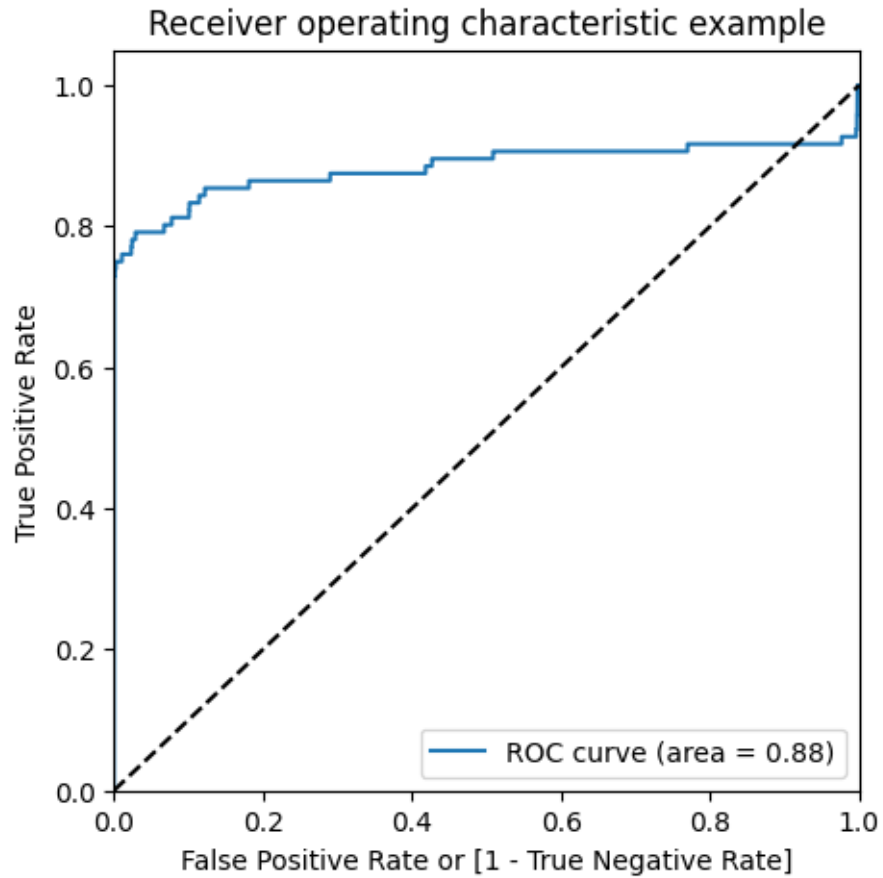
[82]: np.float64(0.882980017350731)

```

```

[83]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)

```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 0.85
 - Specificity = 0.99
 - ROC-AUC = 0.99
 - F1-Score = 0.90
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.75
 - Specificity = 0.99
 - ROC-AUC = 0.98
 - F-Score = 0.79

Overall, the model is performing well in the test set, what it had learnt from the train set.

2.0.4 Decision Tree

```
[84]: # Importing decision tree classifier
from sklearn.tree import DecisionTreeClassifier
```

```
[85]: # Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
}

# Instantiate the grid search model
dtree = DecisionTreeClassifier()

grid_search = GridSearchCV(estimator = dtree,
                           param_grid = param_grid,
                           scoring= 'roc_auc',
                           cv = 3,
                           verbose = 1)

# Fit the grid search to the data
grid_search.fit(X_train,y_train)
```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```
[85]: GridSearchCV(cv=3, estimator=DecisionTreeClassifier(),
                  param_grid={'max_depth': range(5, 15, 5),
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50)},
                  scoring='roc_auc', verbose=1)
```

```
[89]: # cv results
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results
```

```
[89]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      3.998808    0.055323      0.016658      0.000187
1      3.858425    0.005329      0.016564      0.000062
2      3.834432    0.020078      0.016550      0.000254
3      3.842935    0.054652      0.016533      0.000231
4      7.532441    0.031690      0.017664      0.000240
5      7.424438    0.195222      0.187010      0.246578
6      7.369755    0.066111      0.014435      0.002657
7      7.366655    0.030611      0.017967      0.000107

   param_max_depth  param_min_samples_leaf  param_min_samples_split  \
```

0	5	50	50
1	5	50	100
2	5	100	50
3	5	100	100
4	10	50	50
5	10	50	100
6	10	100	50
7	10	100	100

	params	split0_test_score \
0	{'max_depth': 5, 'min_samples_leaf': 50, 'min...	0.933337
1	{'max_depth': 5, 'min_samples_leaf': 50, 'min...	0.933337
2	{'max_depth': 5, 'min_samples_leaf': 100, 'min...	0.933279
3	{'max_depth': 5, 'min_samples_leaf': 100, 'min...	0.933282
4	{'max_depth': 10, 'min_samples_leaf': 50, 'min...	0.917537
5	{'max_depth': 10, 'min_samples_leaf': 50, 'min...	0.917523
6	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.933448
7	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.933446

	split1_test_score	split2_test_score	mean_test_score	std_test_score \
0	0.933183	0.923866	0.930129	0.004429
1	0.933183	0.923871	0.930130	0.004427
2	0.936720	0.944598	0.938199	0.004738
3	0.936720	0.944598	0.938200	0.004737
4	0.916479	0.938009	0.924008	0.009909
5	0.916442	0.930472	0.921479	0.006374
6	0.919684	0.944093	0.932408	0.009992
7	0.919673	0.921774	0.924965	0.006058

	rank_test_score
0	5
1	4
2	2
3	1
4	7
5	8
6	3
7	6

```
[90]: # Printing the optimal sensitivity score and hyperparameters
print("Best roc_auc:-", grid_search.best_score_)
print(grid_search.best_estimator_)
```

Best roc_auc:- 0.9382001202914115

DecisionTreeClassifier(max_depth=5, min_samples_leaf=100, min_samples_split=100)

```
[91]: # Model with optimal hyperparameters
dt_imb_model = DecisionTreeClassifier(criterion = "gini",
                                     random_state = 100,
                                     max_depth=5,
                                     min_samples_leaf=100,
                                     min_samples_split=100)

dt_imb_model.fit(X_train, y_train)
```

```
[91]: DecisionTreeClassifier(max_depth=5, min_samples_leaf=100, min_samples_split=100,
                             random_state=100)
```

Prediction on the train set

```
[92]: # Predictions on the train set
y_train_pred = dt_imb_model.predict(X_train)
```

```
[93]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train, y_train)
print(confusion)
```

```
[[227449    0]
 [      0   396]]
```

```
[94]: TP = confusion[1,1] # true positive
      TN = confusion[0,0] # true negatives
      FP = confusion[0,1] # false positives
      FN = confusion[1,0] # false negatives
```

```
[95]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train, y_train_pred))
```

```
Accuracy:- 0.9991704887094297
Sensitivity:- 1.0
Specificity:- 1.0
F1-Score:- 0.749003984063745
```

```
[96]: # classification_report
print(classification_report(y_train, y_train_pred))
```

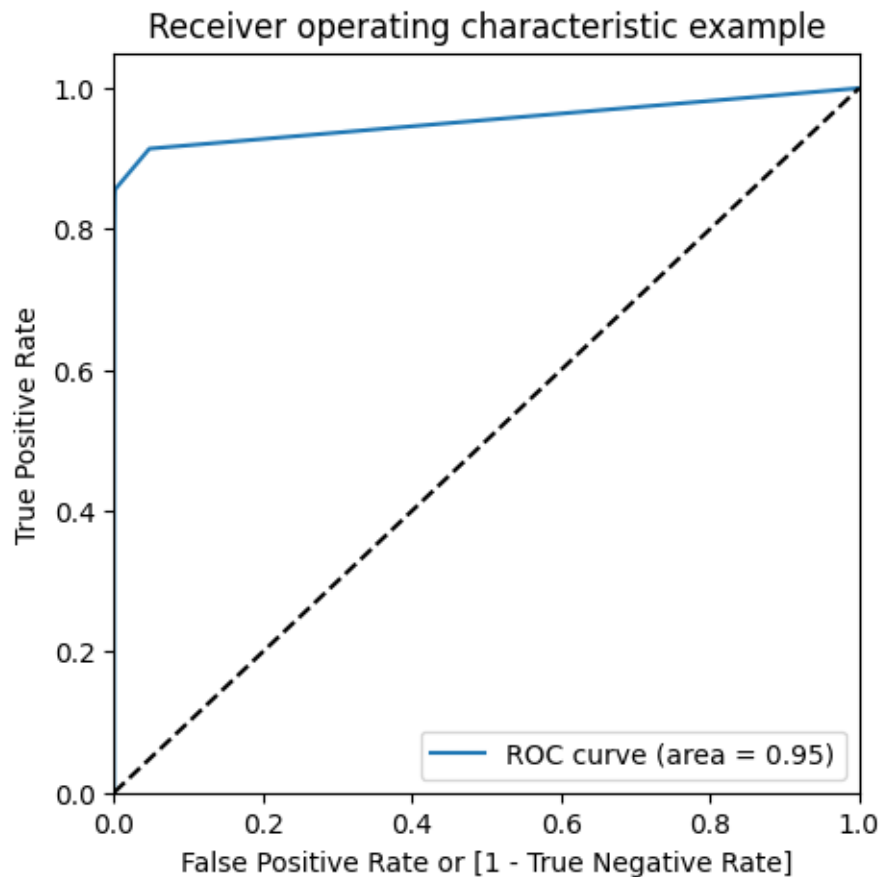
	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	0.79	0.71	0.75	396
accuracy			1.00	227845
macro avg	0.89	0.86	0.87	227845
weighted avg	1.00	1.00	1.00	227845

```
[97]: # Predicted probability
y_train_pred_proba = dt_imb_model.predict_proba(X_train)[:,-1]
```

```
[98]: # roc_auc
auc = metrics.roc_auc_score(y_train, y_train_pred_proba)
auc
```

```
[98]: np.float64(0.9534547393930157)
```

```
[99]: # Plot the ROC curve
draw_roc(y_train, y_train_pred_proba)
```



Prediction on the test set

```
[100]: # Predictions on the test set
y_test_pred = dt_imb_model.predict(X_test)
```

```
[101]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56834   32]
 [   38   58]]
```

```
[102]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[103]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train, y_train_pred))
```

```
Accuracy:- 0.9987711105649381
Sensitivity:- 0.6041666666666666
Specificity:- 0.9994372735905462
F1-Score:- 0.749003984063745
```

```
[104]: # classification_report
print(classification_report(y_test, y_test_pred))
```

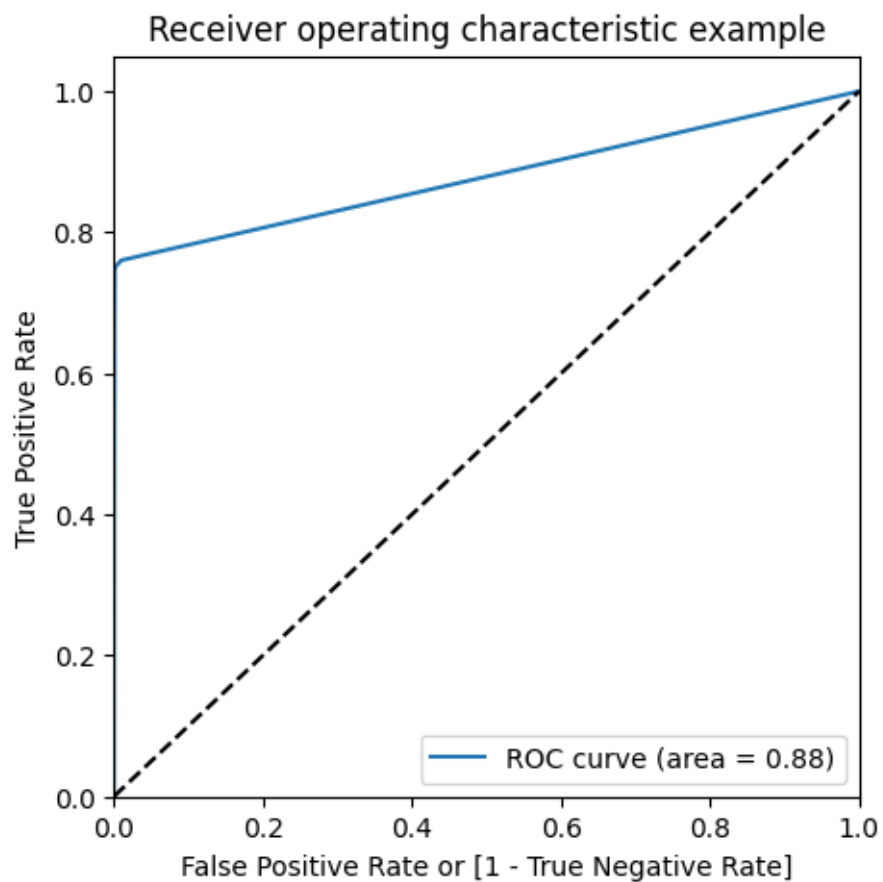
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.64	0.60	0.62	96
accuracy			1.00	56962
macro avg	0.82	0.80	0.81	56962
weighted avg	1.00	1.00	1.00	56962

```
[105]: # Predicted probability
y_test_pred_proba = dt_imb_model.predict_proba(X_test)[:,-1]
```

```
[106]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

```
[106]: np.float64(0.8787224754979542)
```

```
[107]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 1.0
 - Specificity = 1.0
 - F1-Score = 0.75
 - ROC-AUC = 0.95

- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.58
 - Specificity = 0.99
 - F-1 Score = 0.75
 - ROC-AUC = 0.92

2.0.5 Random forest

```
[108]: # Importing random forest classifier
from sklearn.ensemble import RandomForestClassifier
```

```
[ ]: param_grid = {
    'max_depth': range(5,10,5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
    'n_estimators': [100,200,300],
    'max_features': [10, 20]
}
# Create a based model
rf = RandomForestClassifier()
# Instantiate the grid search model
grid_search = GridSearchCV(estimator = rf,
                           param_grid = param_grid,
                           cv = 2,
                           n_jobs = -1,
                           verbose = 1,
                           return_train_score=True)

# Fit the model
grid_search.fit(X_train, y_train)
```

Fitting 2 folds for each of 24 candidates, totalling 48 fits

```
[ ]: # printing the optimal accuracy score and hyperparameters
print('We can get accuracy of',grid_search.best_score_, 'using',grid_search.
      ↪best_params_)
```

```
[95]: # model with the best hyperparameters

rfc_imb_model = RandomForestClassifier(bootstrap=True,
                                       max_depth=5,
                                       min_samples_leaf=50,
                                       min_samples_split=50,
                                       max_features=10,
                                       n_estimators=100)
```

```
[96]: # Fit the model
rfc_imb_model.fit(X_train, y_train)
```

```
[96]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                             criterion='gini', max_depth=5, max_features=10,
                             max_leaf_nodes=None, max_samples=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=50, min_samples_split=50,
                             min_weight_fraction_leaf=0.0, n_estimators=100,
                             n_jobs=None, oob_score=False, random_state=None,
                             verbose=0, warm_start=False)
```

Prediction on the train set

```
[97]: # Predictions on the train set
y_train_pred = rfc_imb_model.predict(X_train)
```

```
[98]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train, y_train_pred)
print(confusion)
```

```
[[227449    0]
 [      0   396]]
```

```
[99]: TP = confusion[1,1] # true positive
      TN = confusion[0,0] # true negatives
      FP = confusion[0,1] # false positives
      FN = confusion[1,0] # false negatives
```

```
[100]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train, y_train_pred))
```

```
Accuracy:- 0.9993460466545239
Sensitivity:- 1.0
Specificity:- 1.0
F1-Score:- 0.7983761840324763
```

```
[101]: # classification_report
print(classification_report(y_train, y_train_pred))
```

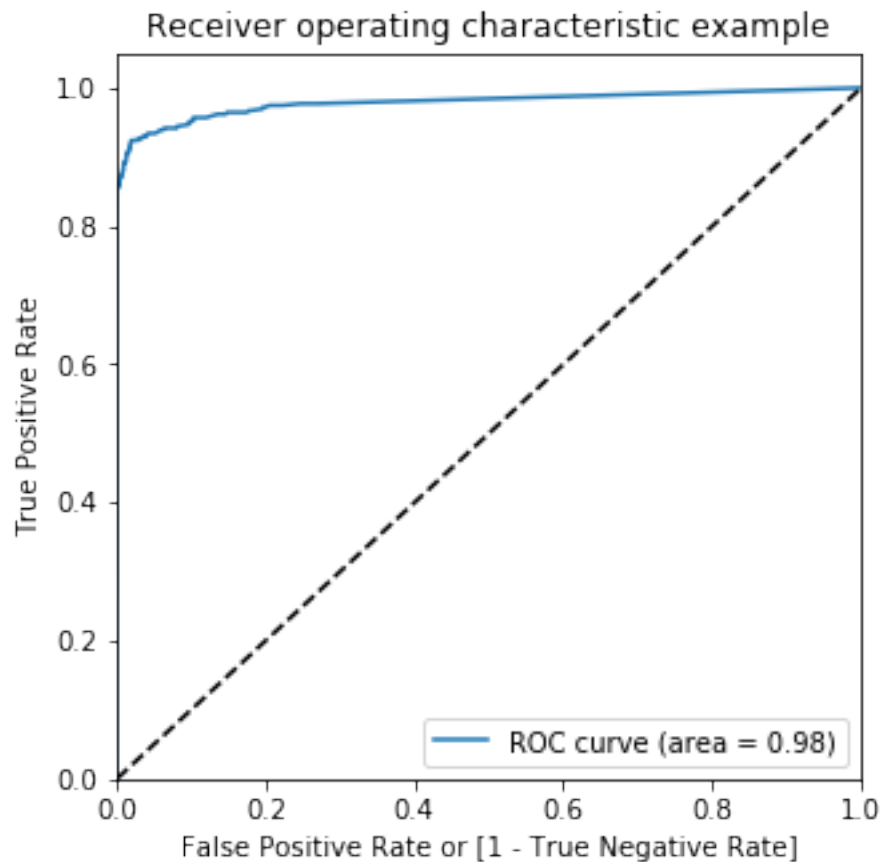

	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	0.86	0.74	0.80	396
accuracy			1.00	227845
macro avg	0.93	0.87	0.90	227845
weighted avg	1.00	1.00	1.00	227845

```
[102]: # Predicted probability
y_train_pred_proba = rfc_imb_model.predict_proba(X_train)[: ,1]
```

```
[103]: # roc_auc
auc = metrics.roc_auc_score(y_train, y_train_pred_proba)
auc
```

```
[103]: 0.9791822295960585
```

```
[104]: # Plot the ROC curve
draw_roc(y_train, y_train_pred_proba)
```



Prediction on the test set

```
[105]: # Predictions on the test set
y_test_pred = rfc_imb_model.predict(X_test)
```

```
[106]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56841    25]
 [   36    60]]
```

```
[107]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[108]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-",TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train, y_train_pred))
```

```
Accuracy:- 0.9989291106351603
Sensitivity:- 0.625
Specificity:- 0.9995603699926142
F1-Score:- 0.7983761840324763
```

```
[109]: # classification_report
print(classification_report(y_test, y_test_pred))
```

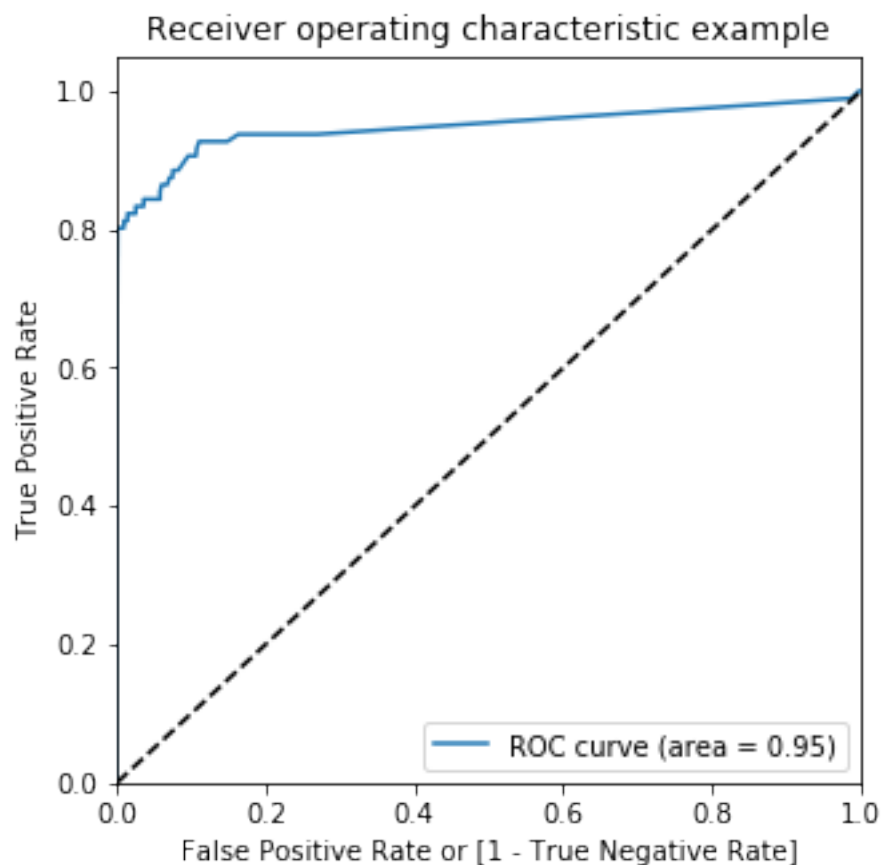
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.71	0.62	0.66	96
accuracy			1.00	56962
macro avg	0.85	0.81	0.83	56962
weighted avg	1.00	1.00	1.00	56962

```
[110]: # Predicted probability
y_test_pred_proba = rfc_imb_model.predict_proba(X_test)[: ,1]
```

```
[111]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

```
[111]: 0.9474696179029063
```

```
[112]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 1.0
 - Specificity = 1.0
 - F1-Score = 0.80
 - ROC-AUC = 0.98

- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.62
 - Specificity = 0.99
 - F-1 Score = 0.75
 - ROC-AUC = 0.96

2.0.6 Choosing best model on the imbalanced data

We can see that among all the models we tried (Logistic, XGBoost, Decision Tree, and Random Forest), almost all of them have performed well. More specifically Logistic regression and XGBoost performed best in terms of ROC-AUC score.

But as we have to choose one of them, we can go for the best as **XGBoost**, which gives us ROC score of 1.0 on the train data and 0.98 on the test data.

Keep in mind that XGBoost requires more resource utilization than Logistic model. Hence building XGBoost model is more costlier than the Logistic model. But XGBoost having ROC score 0.98, which is 0.01 more than the Logistic model. The 0.01 increase of score may convert into huge amount of saving for the bank.

Print the important features of the best model to understand the dataset

- This will not give much explanation on the already transformed dataset
- But it will help us in understanding if the dataset is not PCA transformed

```
[57]: # Features of XGBoost model

var_imp = []
for i in xgb_imb_model.feature_importances_:
    var_imp.append(i)
print('Top var =', var_imp.index(np.sort(xgb_imb_model.
    ↪feature_importances_) [-1]) + 1)
print('2nd Top var =', var_imp.index(np.sort(xgb_imb_model.
    ↪feature_importances_) [-2]) + 1)
print('3rd Top var =', var_imp.index(np.sort(xgb_imb_model.
    ↪feature_importances_) [-3]) + 1)
# Variable on Index-16 and Index-13 seems to be the top 2 variables
top_var_index = var_imp.index(np.sort(xgb_imb_model.feature_importances_) [-1])
second_top_var_index = var_imp.index(np.sort(xgb_imb_model.
    ↪feature_importances_) [-2])

X_train_1 = X_train.to_numpy()[np.where(y_train==1.0)]
X_train_0 = X_train.to_numpy()[np.where(y_train==0.0)]

np.random.shuffle(X_train_0)

import matplotlib.pyplot as plt
%matplotlib inline
```

```
plt.rcParams['figure.figsize'] = [20, 20]

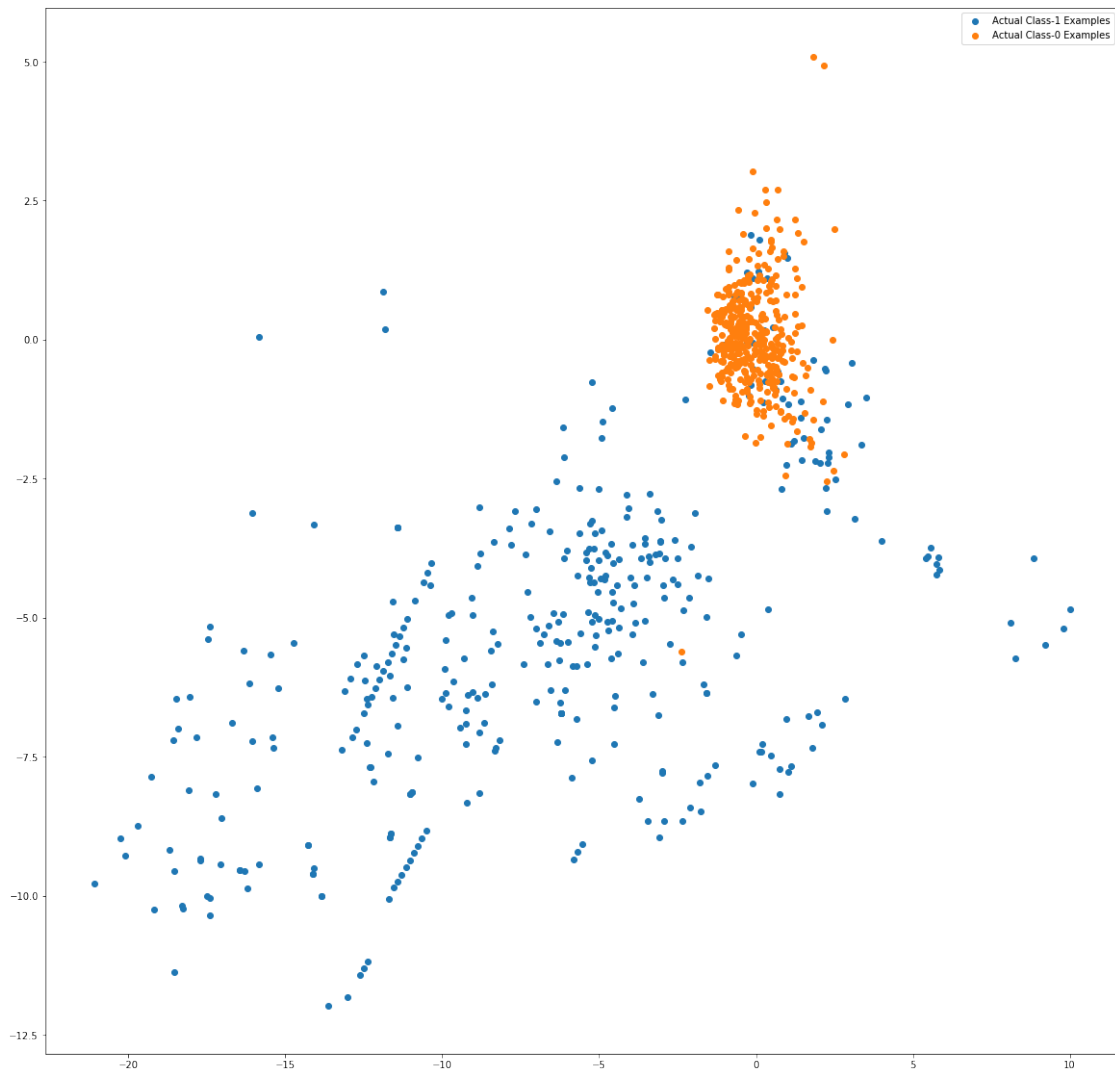
plt.scatter(X_train_1[:, top_var_index], X_train_1[:, second_top_var_index],
            ↪label='Actual Class-1 Examples')
plt.scatter(X_train_0[:X_train_1.shape[0], top_var_index], X_train_0[:X_train_1.
            ↪shape[0], second_top_var_index],
            label='Actual Class-0 Examples')
plt.legend()
```

Top var = 17

2nd Top var = 14

3rd Top var = 10

[57]: <matplotlib.legend.Legend at 0x11887c88>



Print the FPR,TPR & select the best threshold from the roc curve for the best model

```
[66]: print('Train auc =', metrics.roc_auc_score(y_train, y_train_pred_proba_imb_xgb))
      fpr, tpr, thresholds = metrics.roc_curve(y_train, y_train_pred_proba_imb_xgb)
      threshold = thresholds[np.argmax(tpr-fpr)]
      print("Threshold=",threshold)
```

Train auc = 1.0

Threshold= 0.8474788

We can see that the threshold is 0.85, for which the TPR is the highest and FPR is the lowest and we got the best ROC score.

3 Handling data imbalance

As we see that the data is heavily imbalanced, We will try several approaches for handling data imbalance.

- Undersampling :- Here for balancing the class distribution, the non-fraudulent transactions count will be reduced to 396 (similar count of fraudulent transactions)
- Oversampling :- Here we will make the same count of non-fraudulent transactions as fraudulent transactions.
- SMOTE :- Synthetic minority oversampling technique. It is another oversampling technique, which uses nearest neighbor algorithm to create synthetic data.
- Adasyn:- This is similar to SMOTE with minor changes that the new synthetic data is generated on the region of low density of imbalanced data points.

3.1 Undersampling

```
[116]: # Importing undersampler library
      from imblearn.under_sampling import RandomUnderSampler
      from collections import Counter
```

```
[117]: # instantiating the random undersampler
      rus = RandomUnderSampler()
      # resampling X, y
      X_train_rus, y_train_rus = rus.fit_resample(X_train, y_train)
```

```
[118]: # Before sampling class distribution
      print('Before sampling class distribution:-',Counter(y_train))
      # new class distribution
      print('New class distribution:-',Counter(y_train_rus))
```

Before sampling class distribution:- Counter({0: 227449, 1: 396})

New class distribution:- Counter({0: 396, 1: 396})

3.2 Model building on balanced data with Undersampling

3.2.1 Logistic Regression

```
[50]: # Creating KFold object with 5 splits
folds = KFold(n_splits=5, shuffle=True, random_state=4)

# Specify params
params = {"C": [0.01, 0.1, 1, 10, 100, 1000]}

# Specifying score as roc_auc
model_cv = GridSearchCV(estimator = LogisticRegression(),
                        param_grid = params,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# Fit the model
model_cv.fit(X_train_rus, y_train_rus)
```

Fitting 5 folds for each of 6 candidates, totalling 30 fits

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 30 out of 30 | elapsed: 0.7s finished
```

```
[50]: GridSearchCV(cv=KFold(n_splits=5, random_state=4, shuffle=True),
                  error_score=nan,
                  estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
                                                fit_intercept=True,
                                                intercept_scaling=1, l1_ratio=None,
                                                max_iter=100, multi_class='auto',
                                                n_jobs=None, penalty='l2',
                                                random_state=None, solver='lbfgs',
                                                tol=0.0001, verbose=0,
                                                warm_start=False),
                  iid='deprecated', n_jobs=None,
                  param_grid={'C': [0.01, 0.1, 1, 10, 100, 1000]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)
```

```
[51]: # results of grid search CV
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[51]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  param_C \
0      0.016201    0.009065      0.0040      1.095540e-03    0.01
1      0.016801    0.002136      0.0042      7.483665e-04    0.1
2      0.026201    0.004118      0.0040      1.095453e-03    1
```

3	0.020201	0.002786	0.0030	9.536743e-08	10
4	0.020801	0.002561	0.0030	6.324097e-04	100
5	0.021601	0.001497	0.0026	4.898624e-04	1000

	params	split0_test_score	split1_test_score	split2_test_score	\
0	{'C': 0.01}	0.983943	0.995410	0.972276	
1	{'C': 0.1}	0.981240	0.995568	0.976122	
2	{'C': 1}	0.981081	0.994302	0.978365	
3	{'C': 10}	0.975199	0.994777	0.978846	
4	{'C': 100}	0.972496	0.994619	0.978846	
5	{'C': 1000}	0.972178	0.994619	0.978686	

	split3_test_score	split4_test_score	mean_test_score	std_test_score	\
0	0.976110	0.987913	0.983130	0.008264	
1	0.974186	0.989202	0.983264	0.008052	
2	0.971621	0.989363	0.982947	0.008036	
3	0.966330	0.990169	0.981064	0.010270	
4	0.965368	0.990169	0.980300	0.010848	
5	0.965368	0.990169	0.980204	0.010899	

	rank_test_score	split0_train_score	split1_train_score	\
0	2	0.988177	0.985624	
1	1	0.990903	0.987521	
2	3	0.991911	0.988759	
3	4	0.992371	0.988999	
4	5	0.992431	0.989049	
5	6	0.992401	0.989059	

	split2_train_score	split3_train_score	split4_train_score	\
0	0.989641	0.988655	0.987128	
1	0.992268	0.991223	0.989408	
2	0.993094	0.992646	0.990384	
3	0.993392	0.992904	0.990553	
4	0.993412	0.993113	0.990563	
5	0.993502	0.993113	0.990533	

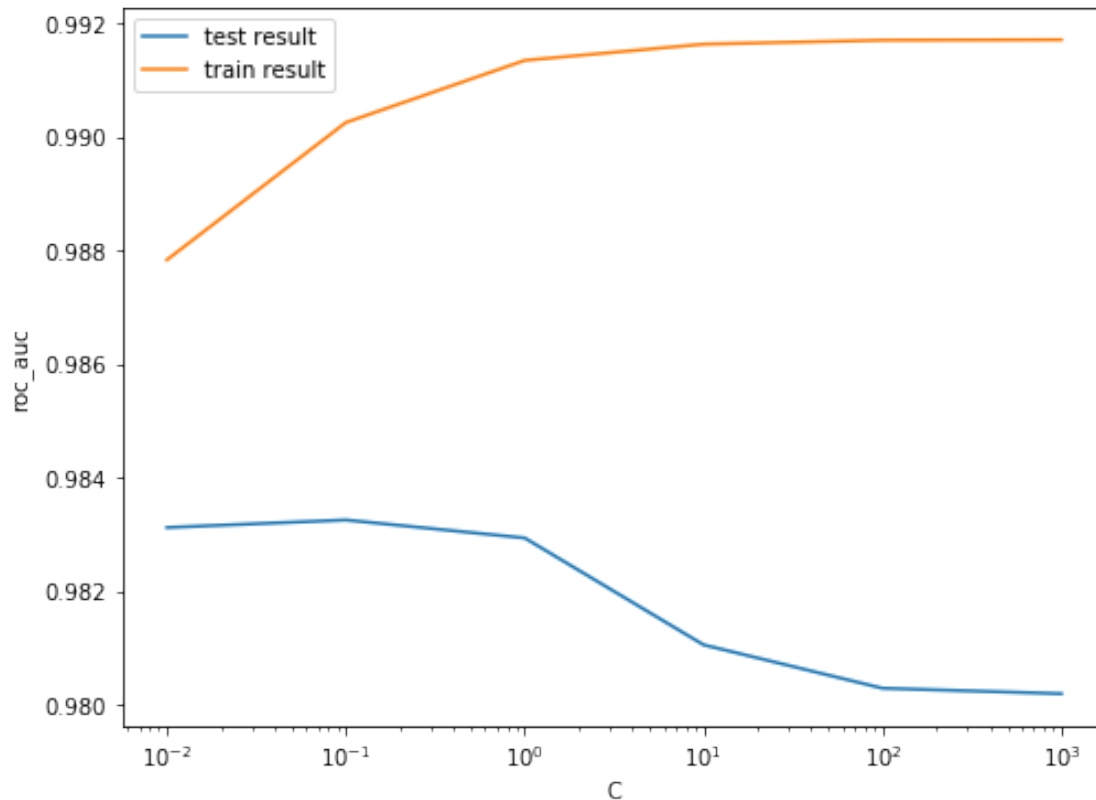
	mean_train_score	std_train_score
0	0.987845	0.001374
1	0.990264	0.001649
2	0.991359	0.001593
3	0.991644	0.001635
4	0.991714	0.001660
5	0.991721	0.001678

[52]: *# plot of C versus train and validation scores*

```
plt.figure(figsize=(8, 6))
```



```
plt.plot(cv_results['param_C'], cv_results['mean_test_score'])
plt.plot(cv_results['param_C'], cv_results['mean_train_score'])
plt.xlabel('C')
plt.ylabel('roc_auc')
plt.legend(['test result', 'train result'], loc='upper left')
plt.xscale('log')
```



```
[53]: # Best score with best C
best_score = model_cv.best_score_
best_C = model_cv.best_params_['C']

print(" The highest test roc_auc is {0} at C = {1}".format(best_score, best_C))
```

The highest test roc_auc is 0.9832637280039689 at C = 0.1

Logistic regression with optimal C

```
[119]: # Instantiate the model with best C
logistic_bal_rus = LogisticRegression(C=0.1)
```

```
[120]: # Fit the model on the train set
logistic_bal_rus_model = logistic_bal_rus.fit(X_train_rus, y_train_rus)
```

Prediction on the train set

```
[121]: # Predictions on the train set
y_train_pred = logistic_bal_rus_model.predict(X_train_rus)
```

```
[122]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_rus, y_train_pred)
print(confusion)
```

```
[[391   5]
 [ 32 364]]
```

```
[123]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[124]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_train_rus, y_train_pred))

# Sensitivity
print("Sensitivity:-",TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train_rus, y_train_pred))
```

```
Accuracy:- 0.9532828282828283
Sensitivity:- 0.9191919191919192
Specificity:- 0.9873737373737373
F1-Score:- 0.9516339869281046
```

```
[125]: # classification_report
print(classification_report(y_train_rus, y_train_pred))
```

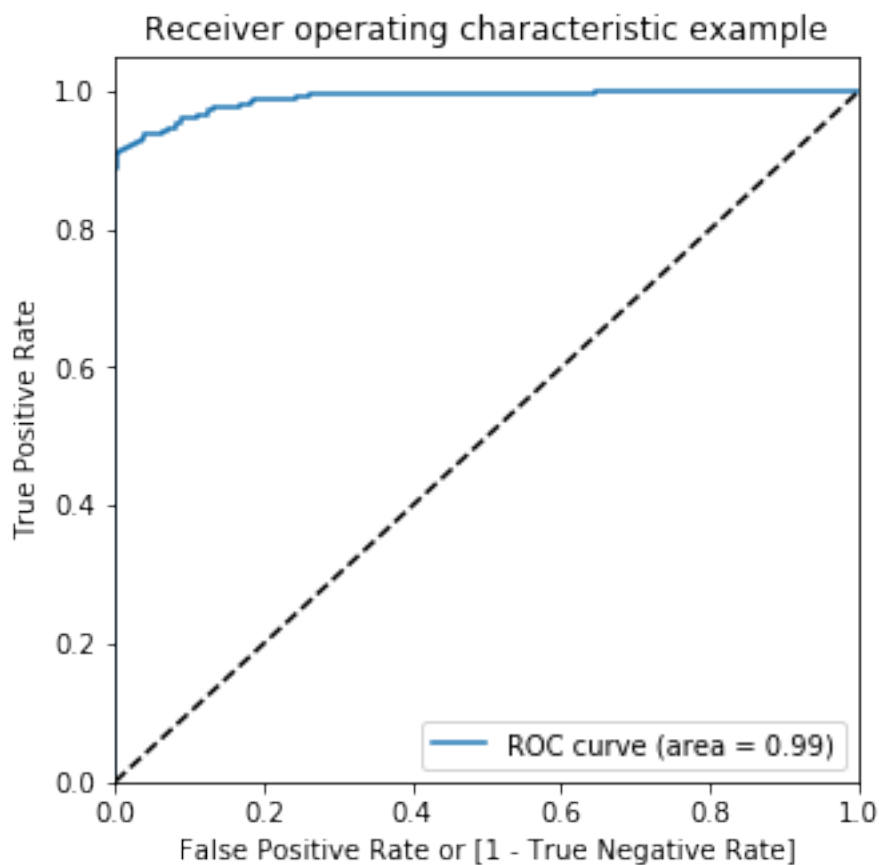
	precision	recall	f1-score	support
0	0.92	0.99	0.95	396
1	0.99	0.92	0.95	396
accuracy			0.95	792
macro avg	0.96	0.95	0.95	792
weighted avg	0.96	0.95	0.95	792

```
[126]: # Predicted probability
y_train_pred_proba = logistic_bal_rus_model.predict_proba(X_train_rus)[: ,1]
```

```
[127]: # roc_auc
auc = metrics.roc_auc_score(y_train_rus, y_train_pred_proba)
auc
```

```
[127]: 0.9892230384654627
```

```
[128]: # Plot the ROC curve
draw_roc(y_train_rus, y_train_pred_proba)
```



Prediction on the test set

```
[129]: # Prediction on the test set
y_test_pred = logistic_bal_rus_model.predict(X_test)
```

```
[130]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55658 1208]
 [   13   83]]
```

```
[131]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[132]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 0.9785646571398476
Sensitivity:- 0.8645833333333334
Specificity:- 0.978757078043119
```

```
[133]: # classification_report
print(classification_report(y_test, y_test_pred))
```

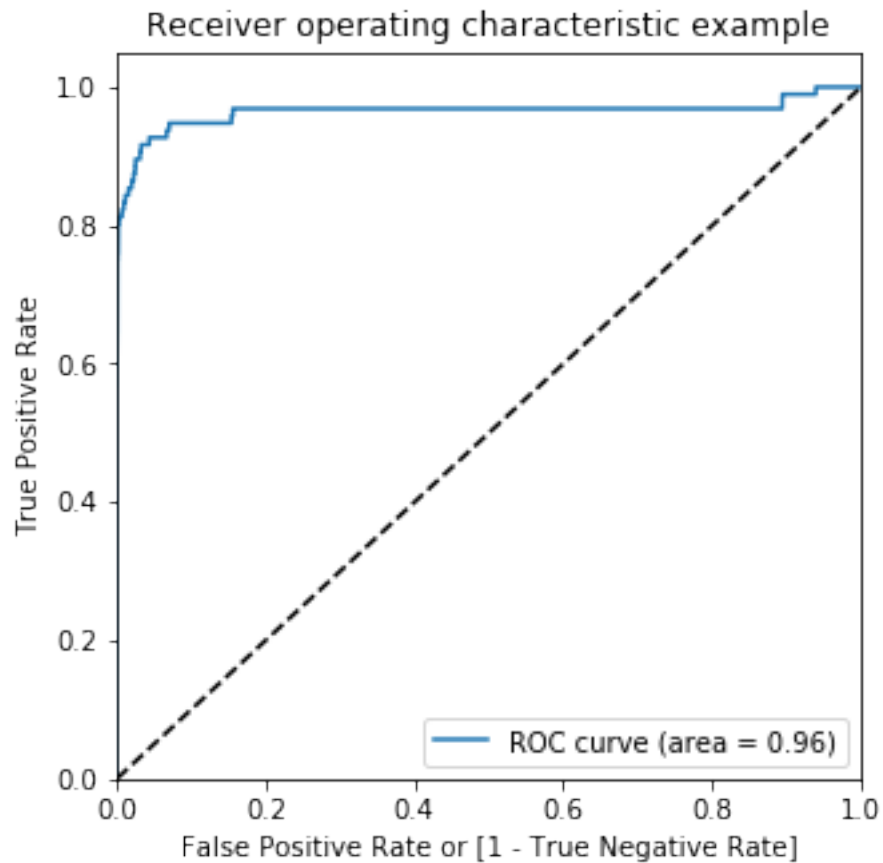
	precision	recall	f1-score	support
0	1.00	0.98	0.99	56866
1	0.06	0.86	0.12	96
accuracy			0.98	56962
macro avg	0.53	0.92	0.55	56962
weighted avg	1.00	0.98	0.99	56962

```
[134]: # Predicted probability
y_test_pred_proba = logistic_bal_rus_model.predict_proba(X_test)[: ,1]
```

```
[135]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

```
[135]: 0.9639748854031114
```

```
[136]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.95
 - Sensitivity = 0.92
 - Specificity = 0.98
 - ROC = 0.99
- Test set
 - Accuracy = 0.97
 - Sensitivity = 0.86
 - Specificity = 0.97
 - ROC = 0.96

3.2.2 XGBoost

```
[73]: # hyperparameter tuning with XGBoost
```

```
# creating a KFold object  
folds = 3
```

```

# specify range of hyperparameters
param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

# specify model
xgb_model = XGBClassifier(max_depth=2, n_estimators=200)

# set up GridSearchCV()
model_cv = GridSearchCV(estimator = xgb_model,
                        param_grid = param_grid,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# fit the model
model_cv.fit(X_train_rus, y_train_rus)

```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 3.9s finished

```

```

[73]: GridSearchCV(cv=3, error_score=nan,
                  estimator=XGBClassifier(base_score=None, booster=None,
                                          colsample_bylevel=None,
                                          colsample_bynode=None,
                                          colsample_bytree=None, gamma=None,
                                          gpu_id=None, importance_type='gain',
                                          interaction_constraints=None,
                                          learning_rate=None, max_delta_step=None,
                                          max_depth=2, min_child_weight=None,
                                          missing=nan, monotone_constraints=None,
                                          n_estimators=200,
                                          objective='binary:logistic',
                                          random_state=None, reg_alpha=None,
                                          reg_lambda=None, scale_pos_weight=None,
                                          subsample=None, tree_method=None,
                                          validate_parameters=False,
                                          verbosity=None),
                  iid='deprecated', n_jobs=None,
                  param_grid={'learning_rate': [0.2, 0.6],
                              'subsample': [0.3, 0.6, 0.9]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)

```

```
[74]: # cv results
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[74]:
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	\
0	0.210345	0.069442	0.016668	0.014384	
1	0.168343	0.003300	0.006334	0.000471	
2	0.247348	0.025370	0.006334	0.000471	
3	0.247347	0.116300	0.011667	0.005437	
4	0.188344	0.026248	0.007001	0.000817	
5	0.171343	0.023115	0.006334	0.000471	

	param_learning_rate	param_subsample	\
0	0.2	0.3	
1	0.2	0.6	
2	0.2	0.9	
3	0.6	0.3	
4	0.6	0.6	
5	0.6	0.9	

	params	split0_test_score	\
0	{'learning_rate': 0.2, 'subsample': 0.3}	0.967172	
1	{'learning_rate': 0.2, 'subsample': 0.6}	0.969295	
2	{'learning_rate': 0.2, 'subsample': 0.9}	0.969238	
3	{'learning_rate': 0.6, 'subsample': 0.3}	0.967172	
4	{'learning_rate': 0.6, 'subsample': 0.6}	0.964073	
5	{'learning_rate': 0.6, 'subsample': 0.9}	0.970500	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.973714	0.982381	0.974422	0.006229	
1	0.974518	0.981749	0.975187	0.005106	
2	0.974690	0.981061	0.974996	0.004831	
3	0.969754	0.978478	0.971801	0.004837	
4	0.976297	0.976010	0.972127	0.005696	
5	0.968951	0.976928	0.972127	0.003454	

	rank_test_score	split0_train_score	split1_train_score	\
0	3	0.999986	1.0	
1	1	1.000000	1.0	
2	2	1.000000	1.0	
3	6	1.000000	1.0	
4	4	1.000000	1.0	
5	4	1.000000	1.0	

	split2_train_score	mean_train_score	std_train_score
0	0.999971	0.999986	0.000012
1	1.000000	1.000000	0.000000

2	1.000000	1.000000	0.000000
3	1.000000	1.000000	0.000000
4	1.000000	1.000000	0.000000
5	1.000000	1.000000	0.000000

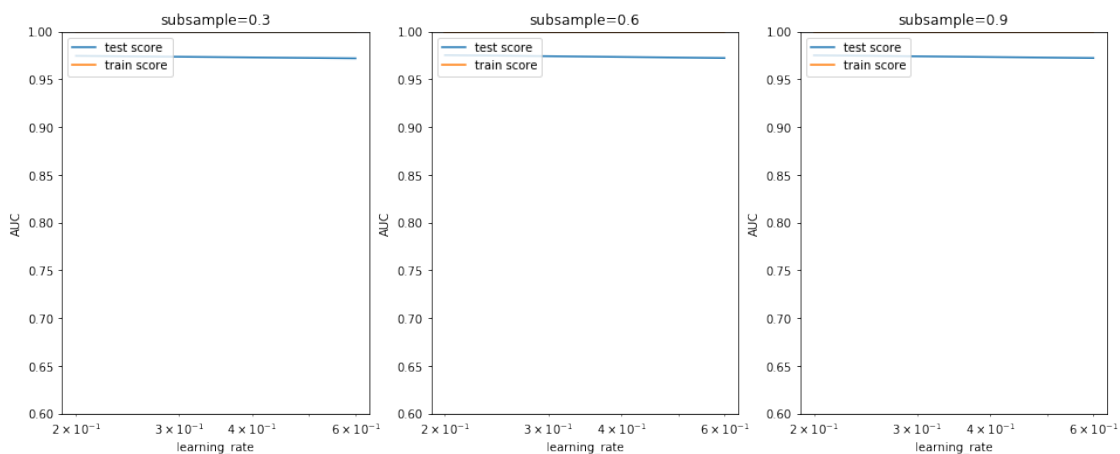
```
[75]: ## plotting
plt.figure(figsize=(16,6))

param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

for n, subsample in enumerate(param_grid['subsample']):

    # subplot 1/n
    plt.subplot(1,len(param_grid['subsample']), n+1)
    df = cv_results[cv_results['param_subsample']==subsample]

    plt.plot(df["param_learning_rate"], df["mean_test_score"])
    plt.plot(df["param_learning_rate"], df["mean_train_score"])
    plt.xlabel('learning_rate')
    plt.ylabel('AUC')
    plt.title("subsample={0}".format(subsample))
    plt.ylim([0.60, 1])
    plt.legend(['test score', 'train score'], loc='upper left')
    plt.xscale('log')
```



Model with optimal hyperparameters We see that the train score almost touches to 1. Among the hyperparameters, we can choose the best parameters as learning_rate : 0.2 and subsample: 0.3


```
[76]: model_cv.best_params_
```

```
[76]: {'learning_rate': 0.2, 'subsample': 0.6}
```

```
[137]: # chosen hyperparameters
# 'objective': 'binary:logistic' outputs probability rather than label, which we
# need for calculating auc
params = {'learning_rate': 0.2,
          'max_depth': 2,
          'n_estimators': 200,
          'subsample': 0.6,
          'objective': 'binary:logistic'}

# fit model on training data
xgb_bal_rus_model = XGBClassifier(params = params)
xgb_bal_rus_model.fit(X_train_rus, y_train_rus)
```

```
[137]: XGBClassifier(base_score=0.5, booster=None, colsample_bylevel=1,
                  colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                  importance_type='gain', interaction_constraints=None,
                  learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                  min_child_weight=1, missing=nan, monotone_constraints=None,
                  n_estimators=100, n_jobs=0, num_parallel_tree=1,
                  objective='binary:logistic',
                  params={'learning_rate': 0.2, 'max_depth': 2, 'n_estimators': 200,
                          'objective': 'binary:logistic', 'subsample': 0.6},
                  random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
                  subsample=1, tree_method=None, validate_parameters=False,
                  verbosity=None)
```

Prediction on the train set

```
[138]: # Predictions on the train set
y_train_pred = xgb_bal_rus_model.predict(X_train_rus)
```

```
[139]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_rus, y_train_pred)
print(confusion)
```

```
[[396   0]
 [   0 396]]
```

```
[140]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[141]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_rus, y_train_pred))
```

```
# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 1.0
Sensitivity:- 1.0
Specificity:- 1.0

```
[142]: # classification_report
print(classification_report(y_train_rus, y_train_pred))
```

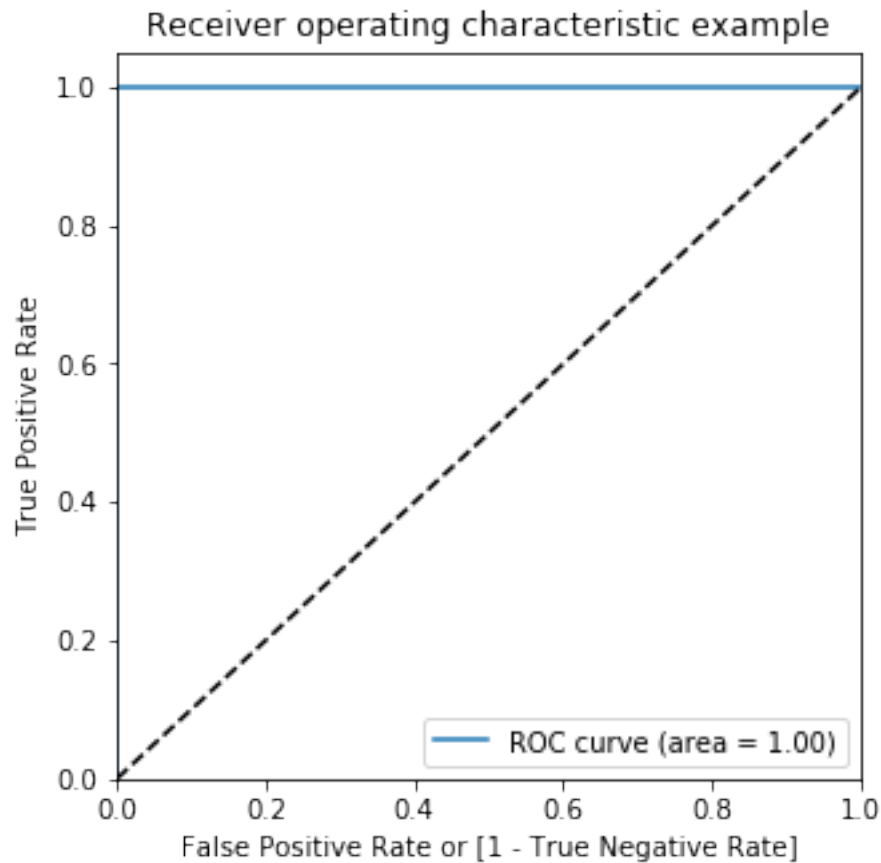
	precision	recall	f1-score	support
0	1.00	1.00	1.00	396
1	1.00	1.00	1.00	396
accuracy			1.00	792
macro avg	1.00	1.00	1.00	792
weighted avg	1.00	1.00	1.00	792

```
[143]: # Predicted probability
y_train_pred_proba = xgb_bal_rus_model.predict_proba(X_train_rus)[:,-1]
```

```
[144]: # roc_auc
auc = metrics.roc_auc_score(y_train_rus, y_train_pred_proba)
auc
```

[144]: 1.0

```
[146]: # Plot the ROC curve
draw_roc(y_train_rus, y_train_pred_proba)
```



Prediction on the test set

```
[147]: # Predictions on the test set
y_test_pred = xgb_bal_rus_model.predict(X_test)
```

```
[148]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[54810  2056]
 [    11    85]]
```

```
[149]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[150]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9637126505389558
Sensitivity:- 0.8854166666666666
Specificity:- 0.9638448281925931

```
[151]: # classification_report
print(classification_report(y_test, y_test_pred))
```

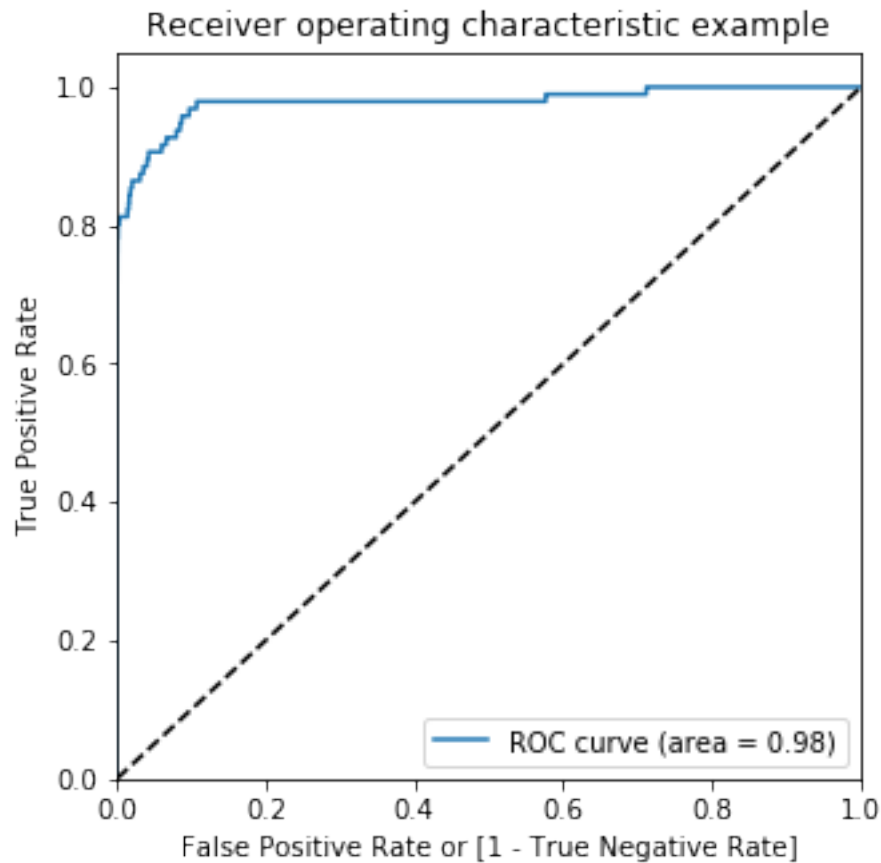
	precision	recall	f1-score	support
0	1.00	0.96	0.98	56866
1	0.04	0.89	0.08	96
accuracy			0.96	56962
macro avg	0.52	0.92	0.53	56962
weighted avg	1.00	0.96	0.98	56962

```
[152]: # Predicted probability
y_test_pred_proba = xgb_bal_rus_model.predict_proba(X_test)[: ,1]
```

```
[153]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[153]: 0.9777381439114174

```
[154]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 1.0
 - Sensitivity = 1.0
 - Specificity = 1.0
 - ROC-AUC = 1.0
- Test set
 - Accuracy = 0.96
 - Sensitivity = 0.92
 - Specificity = 0.96
 - ROC-AUC = 0.98

3.2.3 Decision Tree

```
[105]: # Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
```

```

}

# Instantiate the grid search model
dtree = DecisionTreeClassifier()

grid_search = GridSearchCV(estimator = dtree,
                           param_grid = param_grid,
                           scoring= 'roc_auc',
                           cv = 3,
                           verbose = 1)

# Fit the grid search to the data
grid_search.fit(X_train_rus,y_train_rus)

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 24 out of 24 | elapsed: 0.2s finished

```

```

[105]: GridSearchCV(cv=3, error_score=nan,
                  estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features=None,
                                                    max_leaf_nodes=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    presort='deprecated',
                                                    random_state=None,
                                                    splitter='best'),
                  iid='deprecated', n_jobs=None,
                  param_grid={'max_depth': range(5, 15, 5),
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50)},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                  scoring='roc_auc', verbose=1)

```

```

[106]: # cv results
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results

```

```

[106]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      0.012001  1.632972e-03      0.004000      8.166321e-04
1      0.009334  4.714266e-04      0.003000      1.123916e-07
2      0.007334  4.712580e-04      0.004000      8.165347e-04

```

3	0.007000	1.123916e-07	0.003334	4.714827e-04
4	0.008667	4.715951e-04	0.003333	4.714266e-04
5	0.013334	4.989110e-03	0.004000	8.165347e-04
6	0.007334	1.247235e-03	0.004000	8.165347e-04
7	0.007334	4.714827e-04	0.004000	1.414392e-03

	param_max_depth	param_min_samples_leaf	param_min_samples_split	\
0	5	50	50	
1	5	50	100	
2	5	100	50	
3	5	100	100	
4	10	50	50	
5	10	50	100	
6	10	100	50	
7	10	100	100	

	params	split0_test_score	\
0	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.954345	
1	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.951217	
2	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.948806	
3	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.947773	
4	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.954459	
5	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.954345	
6	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.947773	
7	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.947544	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.963671	0.968607	0.962207	0.005914	
1	0.961920	0.968033	0.960390	0.006950	
2	0.935950	0.960772	0.948510	0.010136	
3	0.935950	0.960772	0.948165	0.010137	
4	0.964532	0.966454	0.961815	0.005260	
5	0.959510	0.966311	0.960055	0.004900	
6	0.935950	0.960772	0.948165	0.010137	
7	0.935950	0.959338	0.947611	0.009548	

	rank_test_score
0	1
1	3
2	5
3	6
4	2
5	4
6	6
7	8

```
[107]: # Printing the optimal sensitivity score and hyperparameters
print("Best roc_auc:-", grid_search.best_score_)
print(grid_search.best_estimator_)
```

```
Best roc_auc:- 0.9622073002754821
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                        max_depth=5, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=50, min_samples_split=50,
                        min_weight_fraction_leaf=0.0, presort='deprecated',
                        random_state=None, splitter='best')
```

```
[155]: # Model with optimal hyperparameters
dt_bal_rus_model = DecisionTreeClassifier(criterion = "gini",
                                          random_state = 100,
                                          max_depth=5,
                                          min_samples_leaf=50,
                                          min_samples_split=50)

dt_bal_rus_model.fit(X_train_rus, y_train_rus)
```

```
[155]: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                              max_depth=5, max_features=None, max_leaf_nodes=None,
                              min_impurity_decrease=0.0, min_impurity_split=None,
                              min_samples_leaf=50, min_samples_split=50,
                              min_weight_fraction_leaf=0.0, presort='deprecated',
                              random_state=100, splitter='best')
```

Prediction on the train set

```
[156]: # Predictions on the train set
y_train_pred = dt_bal_rus_model.predict(X_train_rus)
```

```
[157]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_rus, y_train_pred)
print(confusion)
```

```
[[391   5]
 [ 53 343]]
```

```
[158]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[159]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_rus, y_train_pred))

# Sensitivity
```



```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9267676767676768
Sensitivity:- 0.8661616161616161
Specificity:- 0.9873737373737373

```
[160]: # classification_report
print(classification_report(y_train_rus, y_train_pred))
```

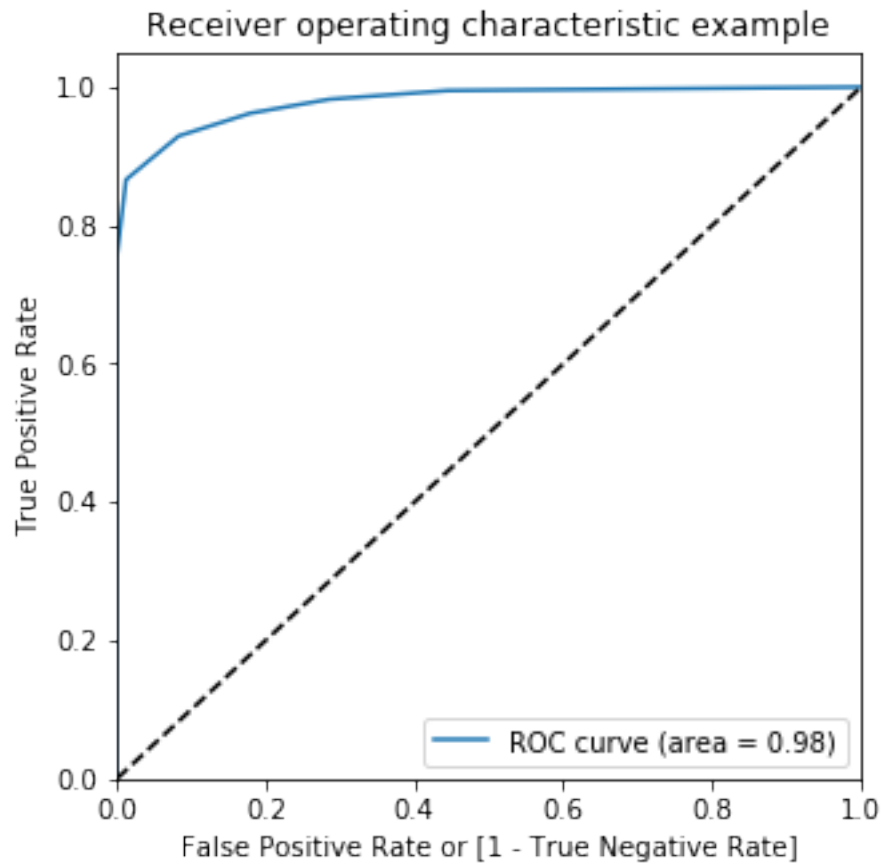
	precision	recall	f1-score	support
0	0.88	0.99	0.93	396
1	0.99	0.87	0.92	396
accuracy			0.93	792
macro avg	0.93	0.93	0.93	792
weighted avg	0.93	0.93	0.93	792

```
[161]: # Predicted probability
y_train_pred_proba = dt_bal_rus_model.predict_proba(X_train_rus)[: ,1]
```

```
[162]: # roc_auc
auc = metrics.roc_auc_score(y_train_rus, y_train_pred_proba)
auc
```

[162]: 0.9789944903581267

```
[163]: # Plot the ROC curve
draw_roc(y_train_rus, y_train_pred_proba)
```



Prediction on the test set

```
[164]: # Predictions on the test set
y_test_pred = dt_bal_rus_model.predict(X_test)
```

```
[165]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55851  1015]
 [   19    77]]
```

```
[166]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[167]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9818475474877989
Sensitivity:- 0.8020833333333334
Specificity:- 0.9821510217001371

```
[168]: # classification_report
print(classification_report(y_test, y_test_pred))
```

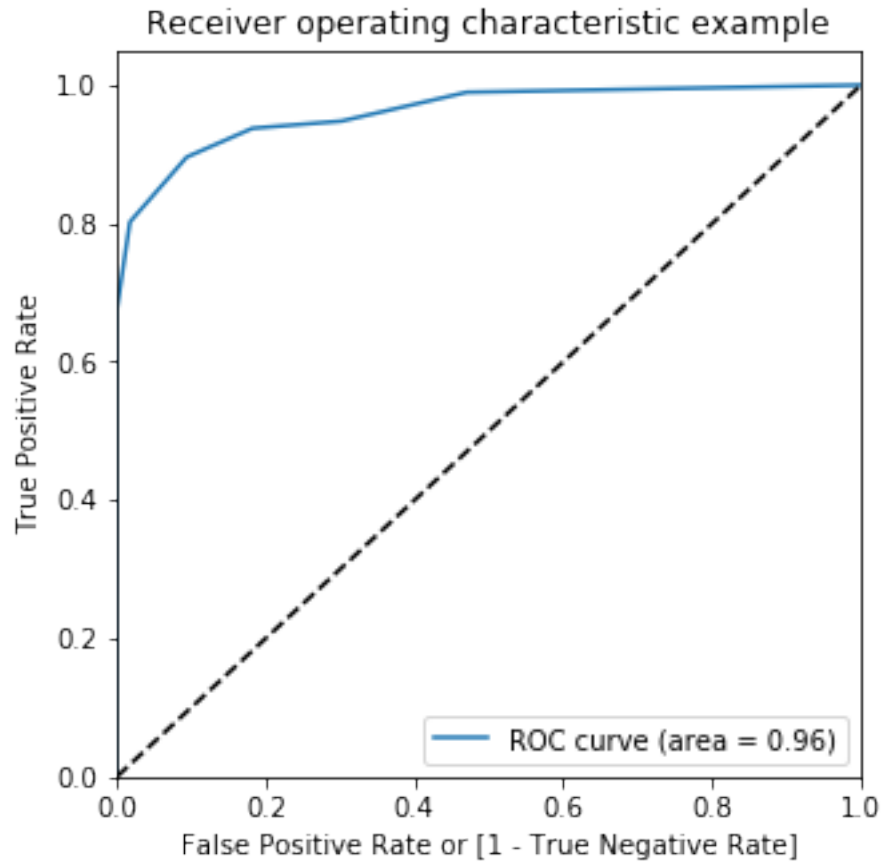
	precision	recall	f1-score	support
0	1.00	0.98	0.99	56866
1	0.07	0.80	0.13	96
accuracy			0.98	56962
macro avg	0.54	0.89	0.56	56962
weighted avg	1.00	0.98	0.99	56962

```
[169]: # Predicted probability
y_test_pred_proba = dt_bal_rus_model.predict_proba(X_test)[: ,1]
```

```
[170]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[170]: 0.9613739243719154

```
[171]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.93
 - Sensitivity = 0.88
 - Specificity = 0.97
 - ROC-AUC = 0.98
- Test set
 - Accuracy = 0.96
 - Sensitivity = 0.85
 - Specificity = 0.96
 - ROC-AUC = 0.96

3.2.4 Random forest

```
[123]: param_grid = {
    'max_depth': range(5,10,5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
    'n_estimators': [100,200,300],
```

```

        'max_features': [10, 20]
    }
    # Create a based model
    rf = RandomForestClassifier()
    # Instantiate the grid search model
    grid_search = GridSearchCV(estimator = rf,
                               param_grid = param_grid,
                               scoring= 'roc_auc',
                               cv = 2,
                               n_jobs = -1,
                               verbose = 1,
                               return_train_score=True)

    # Fit the model
    grid_search.fit(X_train_rus, y_train_rus)

```

Fitting 2 folds for each of 24 candidates, totalling 48 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.

[Parallel(n_jobs=-1)]: Done 48 out of 48 | elapsed: 11.4s finished

```

[123]: GridSearchCV(cv=2, error_score=nan,
                  estimator=RandomForestClassifier(bootstrap=True, ccp_alpha=0.0,
                                                    class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features='auto',
                                                    max_leaf_nodes=None,
                                                    max_samples=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    n_estimators=100, n_jobs=None,
                                                    oob_score=False,
                                                    random_state=None, verbose=0,
                                                    warm_start=False),
                  iid='deprecated', n_jobs=-1,
                  param_grid={'max_depth': range(5, 10, 5), 'max_features': [10, 20],
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50),
                              'n_estimators': [100, 200, 300]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)

```

```

[124]: # printing the optimal accuracy score and hyperparameters
print('We can get roc-auc of',grid_search.best_score_, 'using',grid_search.
      ↪best_params_)

```

We can get roc-auc of 0.976788082848689 using {'max_depth': 5, 'max_features': 10, 'min_samples_leaf': 50, 'min_samples_split': 50, 'n_estimators': 200}

```
[172]: # model with the best hyperparameters

rfc_bal_rus_model = RandomForestClassifier(bootstrap=True,
                                          max_depth=5,
                                          min_samples_leaf=50,
                                          min_samples_split=50,
                                          max_features=10,
                                          n_estimators=200)
```

```
[173]: # Fit the model
rfc_bal_rus_model.fit(X_train_rus, y_train_rus)
```

```
[173]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                             criterion='gini', max_depth=5, max_features=10,
                             max_leaf_nodes=None, max_samples=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=50, min_samples_split=50,
                             min_weight_fraction_leaf=0.0, n_estimators=200,
                             n_jobs=None, oob_score=False, random_state=None,
                             verbose=0, warm_start=False)
```

Prediction on the train set

```
[174]: # Predictions on the train set
y_train_pred = rfc_bal_rus_model.predict(X_train_rus)
```

```
[175]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_rus, y_train_pred)
print(confusion)
```

```
[[391   5]
 [ 44 352]]
```

```
[176]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[177]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_rus, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

```
# F1 score
print("F1-Score:-", f1_score(y_train_rus, y_train_pred))
```

Accuracy:- 0.9381313131313131
 Sensitivity:- 0.8888888888888888
 Specificity:- 0.9873737373737373
 F1-Score:- 0.9349269588313412

```
[178]: # classification_report
print(classification_report(y_train_rus, y_train_pred))
```

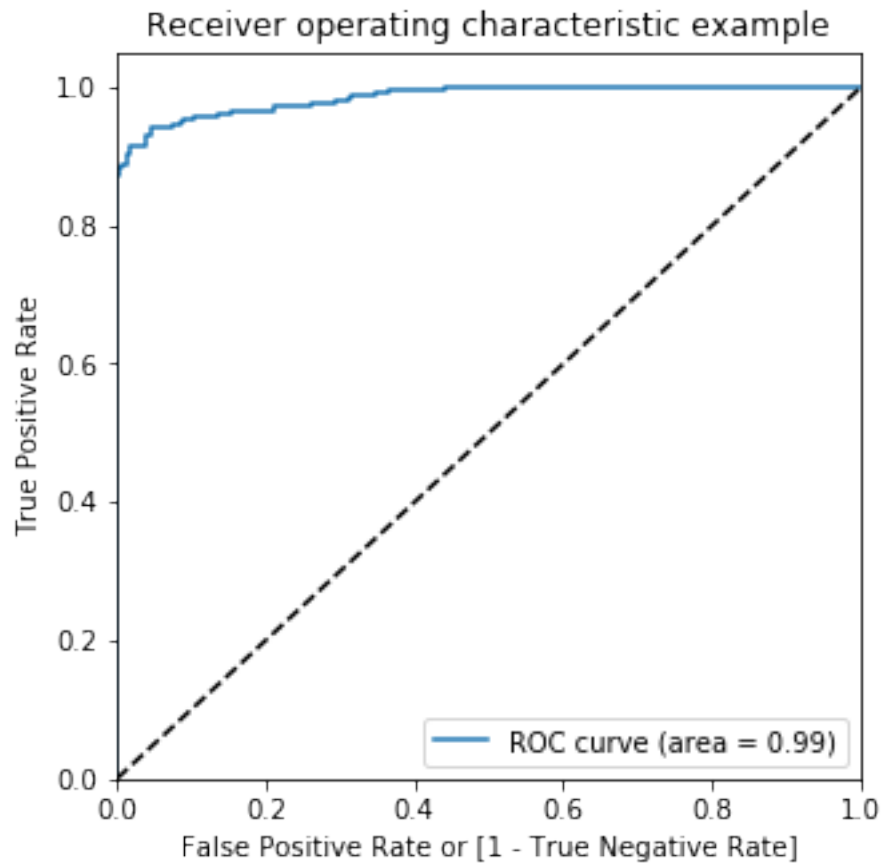
	precision	recall	f1-score	support
0	0.90	0.99	0.94	396
1	0.99	0.89	0.93	396
accuracy			0.94	792
macro avg	0.94	0.94	0.94	792
weighted avg	0.94	0.94	0.94	792

```
[179]: # Predicted probability
y_train_pred_proba = rfc_bal_rus_model.predict_proba(X_train_rus)[: ,1]
```

```
[180]: # roc_auc
auc = metrics.roc_auc_score(y_train_rus, y_train_pred_proba)
auc
```

```
[180]: 0.9851099377614528
```

```
[181]: # Plot the ROC curve
draw_roc(y_train_rus, y_train_pred_proba)
```



Prediction on the test set

```
[182]: # Predictions on the test set
y_test_pred = rfc_bal_rus_model.predict(X_test)
```

```
[183]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55832  1034]
 [   18    78]]
```

```
[184]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[185]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```



```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9815315473473544
Sensitivity:- 0.8125
Specificity:- 0.981816902894524

```
[186]: # classification_report
print(classification_report(y_test, y_test_pred))
```

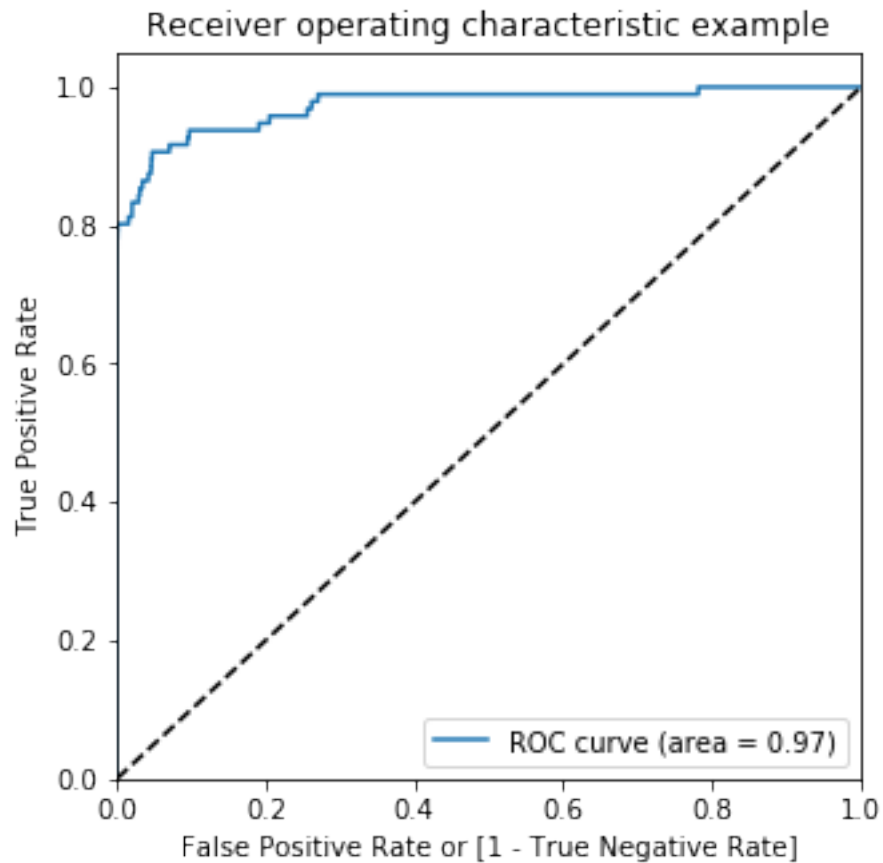
	precision	recall	f1-score	support
0	1.00	0.98	0.99	56866
1	0.07	0.81	0.13	96
accuracy			0.98	56962
macro avg	0.53	0.90	0.56	56962
weighted avg	1.00	0.98	0.99	56962

```
[187]: # Predicted probability
y_test_pred_proba = rfc_bal_rus_model.predict_proba(X_test)[: ,1]
```

```
[188]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[188]: 0.9730361178032567

```
[189]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.94
 - Sensitivity = 0.89
 - Specificity = 0.98
 - ROC-AUC = 0.98
- Test set
 - Accuracy = 0.98
 - Sensitivity = 0.83
 - Specificity = 0.98
 - ROC-AUC = 0.97

4 Oversampling

```
[190]: # Importing oversampler library
from imblearn.over_sampling import RandomOverSampler
```

```
[191]: # instantiating the random oversampler
ros = RandomOverSampler()
# resampling X, y
X_train_ros, y_train_ros = ros.fit_resample(X_train, y_train)
```

```
[192]: # Befor sampling class distribution
print('Before sampling class distribution:-',Counter(y_train))
# new class distribution
print('New class distribution:-',Counter(y_train_ros))
```

Before sampling class distribution:- Counter({0: 227449, 1: 396})
New class distribution:- Counter({0: 227449, 1: 227449})

4.0.1 Logistic Regression

```
[145]: # Creating KFold object with 5 splits
folds = KFold(n_splits=5, shuffle=True, random_state=4)

# Specify params
params = {"C": [0.01, 0.1, 1, 10, 100, 1000]}

# Specifying score as roc-auc
model_cv = GridSearchCV(estimator = LogisticRegression(),
                        param_grid = params,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# Fit the model
model_cv.fit(X_train_ros, y_train_ros)
```

Fitting 5 folds for each of 6 candidates, totalling 30 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 30 out of 30 | elapsed: 1.4min finished

```
[145]: GridSearchCV(cv=KFold(n_splits=5, random_state=4, shuffle=True),
                  error_score=nan,
                  estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
                                                fit_intercept=True,
                                                intercept_scaling=1, l1_ratio=None,
                                                max_iter=100, multi_class='auto',
                                                n_jobs=None, penalty='l2',
                                                random_state=None, solver='lbfgs',
                                                tol=0.0001, verbose=0,
                                                warm_start=False),
                  iid='deprecated', n_jobs=None,
                  param_grid={'C': [0.01, 0.1, 1, 10, 100, 1000]}),
```

```
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)
```

```
[146]: # results of grid search CV
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[146]:
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C \
0	2.392937	0.133817	0.052003	0.003847	0.01
1	2.366276	0.096595	0.048522	0.003303	0.1
2	2.725587	0.393503	0.056963	0.008990	1
3	2.949569	0.306817	0.061003	0.008391	10
4	2.584676	0.096526	0.056722	0.007632	100
5	2.384325	0.060643	0.050203	0.003371	1000

	params	split0_test_score	split1_test_score	split2_test_score \
0	{'C': 0.01}	0.988802	0.988039	0.988728
1	{'C': 0.1}	0.988821	0.988048	0.988751
2	{'C': 1}	0.988819	0.988049	0.988751
3	{'C': 10}	0.988820	0.988049	0.988751
4	{'C': 100}	0.988820	0.988050	0.988751
5	{'C': 1000}	0.988820	0.988050	0.988751

	split3_test_score	split4_test_score	mean_test_score	std_test_score \
0	0.988207	0.988824	0.988520	0.000330
1	0.988206	0.988834	0.988532	0.000336
2	0.988202	0.988837	0.988532	0.000336
3	0.988202	0.988837	0.988532	0.000337
4	0.988201	0.988837	0.988532	0.000336
5	0.988201	0.988837	0.988532	0.000336

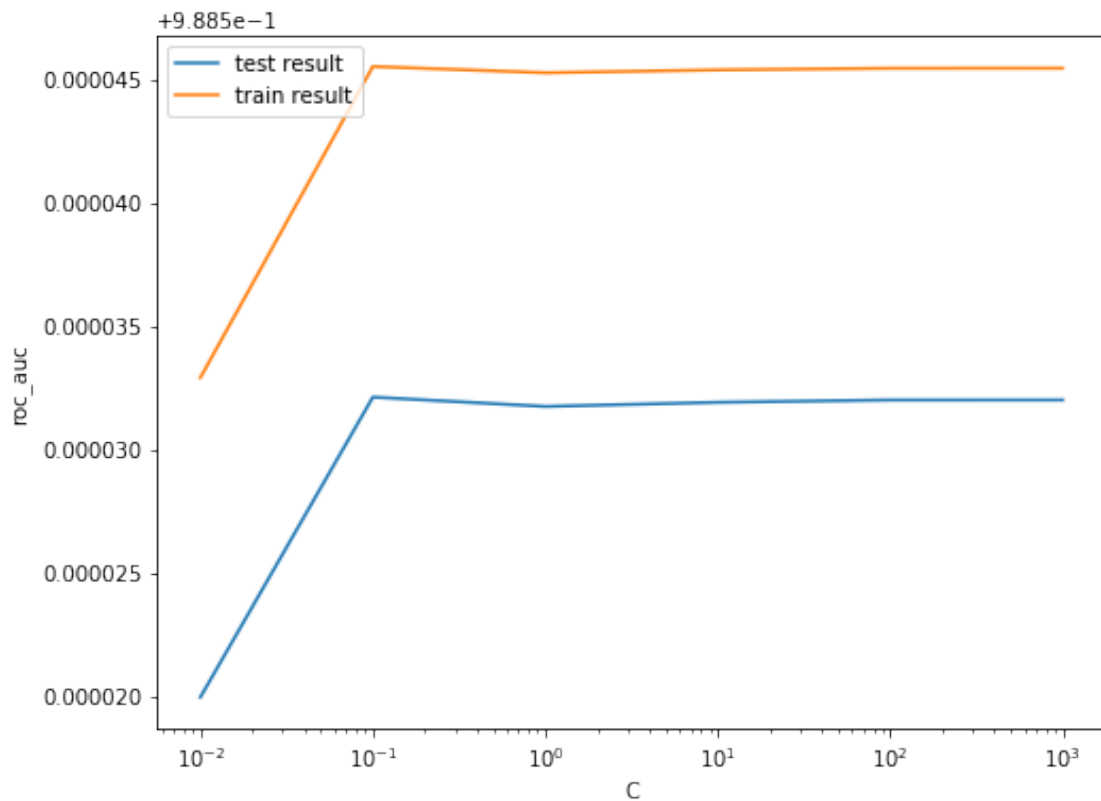
	rank_test_score	split0_train_score	split1_train_score \
0	6	0.988501	0.988558
1	1	0.988516	0.988572
2	5	0.988517	0.988571
3	4	0.988517	0.988571
4	3	0.988517	0.988571
5	2	0.988517	0.988571

	split2_train_score	split3_train_score	split4_train_score \
0	0.988503	0.988633	0.988469
1	0.988514	0.988643	0.988483
2	0.988514	0.988641	0.988484
3	0.988514	0.988641	0.988484
4	0.988514	0.988641	0.988484
5	0.988514	0.988641	0.988484

	mean_train_score	std_train_score
0	0.988533	0.000058
1	0.988546	0.000056
2	0.988545	0.000056
3	0.988545	0.000056
4	0.988545	0.000056
5	0.988545	0.000056

```
[147]: # plot of C versus train and validation scores
```

```
plt.figure(figsize=(8, 6))
plt.plot(cv_results['param_C'], cv_results['mean_test_score'])
plt.plot(cv_results['param_C'], cv_results['mean_train_score'])
plt.xlabel('C')
plt.ylabel('roc_auc')
plt.legend(['test result', 'train result'], loc='upper left')
plt.xscale('log')
```



```
[148]: # Best score with best C
```

```
best_score = model_cv.best_score_
best_C = model_cv.best_params_['C']
```

```
print(" The highest test roc_auc is {0} at C = {1}".format(best_score, best_C))
```

The highest test roc_auc is 0.9885321193436821 at C = 0.1

Logistic regression with optimal C

```
[193]: # Instantiate the model with best C
logistic_bal_ros = LogisticRegression(C=0.1)
```

```
[194]: # Fit the model on the train set
logistic_bal_ros_model = logistic_bal_ros.fit(X_train_ros, y_train_ros)
```

Prediction on the train set

```
[195]: # Predictions on the train set
y_train_pred = logistic_bal_ros_model.predict(X_train_ros)
```

```
[196]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_ros, y_train_pred)
print(confusion)
```

```
[[222261  5188]
 [ 17649 209800]]
```

```
[197]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[198]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_ros, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train_ros, y_train_pred))
```

```
Accuracy:- 0.9497975370302794
Sensitivity:- 0.9224045830054209
Specificity:- 0.9771904910551377
F1-Score:- 0.9483836116780467
```

```
[199]: # classification_report
print(classification_report(y_train_ros, y_train_pred))
```

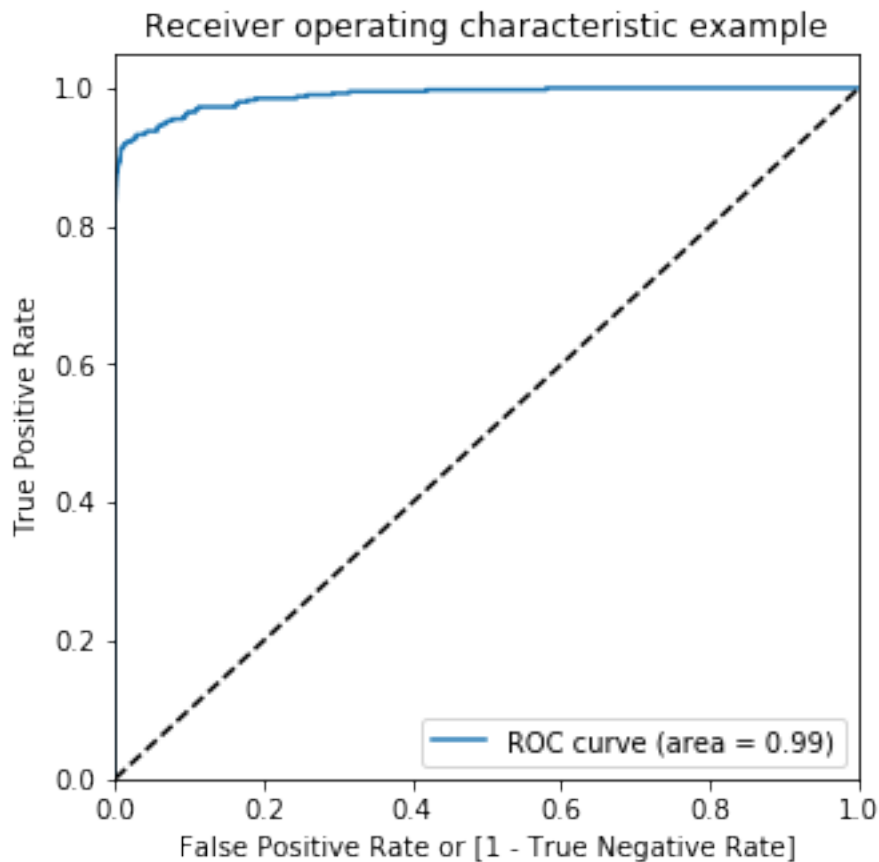
	precision	recall	f1-score	support
0	0.93	0.98	0.95	227449
1	0.98	0.92	0.95	227449
accuracy			0.95	454898
macro avg	0.95	0.95	0.95	454898
weighted avg	0.95	0.95	0.95	454898

```
[200]: # Predicted probability
y_train_pred_proba = logistic_bal_ros_model.predict_proba(X_train_ros)[: ,1]
```

```
[201]: # roc_auc
auc = metrics.roc_auc_score(y_train_ros, y_train_pred_proba)
auc
```

```
[201]: 0.9886578544816166
```

```
[202]: # Plot the ROC curve
draw_roc(y_train_ros, y_train_pred_proba)
```



Prediction on the test set

```
[203]: # Prediction on the test set
y_test_pred = logistic_bal_ros_model.predict(X_test)
```

```
[204]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55540  1326]
 [   11    85]]
```

```
[205]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[206]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-",TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 0.9765282117903163
Sensitivity:- 0.8854166666666666
Specificity:- 0.976682024408258
```

```
[207]: # classification_report
print(classification_report(y_test, y_test_pred))
```

	precision	recall	f1-score	support
0	1.00	0.98	0.99	56866
1	0.06	0.89	0.11	96
accuracy			0.98	56962
macro avg	0.53	0.93	0.55	56962
weighted avg	1.00	0.98	0.99	56962

```
[208]: # Predicted probability
y_test_pred_proba = logistic_bal_ros_model.predict_proba(X_test)[:,-1]
```

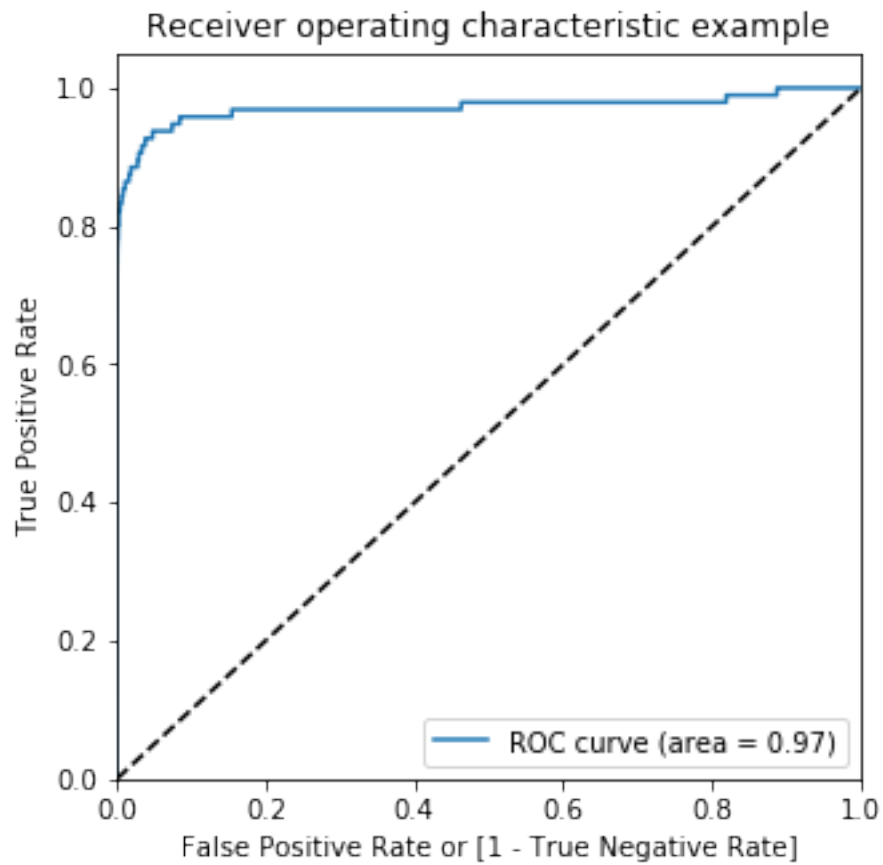
```
[209]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
```



```
auc
```

```
[209]: 0.9712808034091842
```

```
[210]: # Plot the ROC curve  
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.95
 - Sensitivity = 0.92
 - Specificity = 0.97
 - ROC = 0.98
- Test set
 - Accuracy = 0.97
 - Sensitivity = 0.89
 - Specificity = 0.97
 - ROC = 0.97

4.0.2 XGBoost

```
[222]: # hyperparameter tuning with XGBoost

# creating a KFold object
folds = 3

# specify range of hyperparameters
param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

# specify model
xgb_model = XGBClassifier(max_depth=2, n_estimators=200)

# set up GridSearchCV()
model_cv = GridSearchCV(estimator = xgb_model,
                        param_grid = param_grid,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# fit the model
model_cv.fit(X_train_ros, y_train_ros)
```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 33.2min finished

```
[222]: GridSearchCV(cv=3, error_score=nan,
                  estimator=XGBClassifier(base_score=None, booster=None,
                                          colsample_bylevel=None,
                                          colsample_bynode=None,
                                          colsample_bytree=None, gamma=None,
                                          gpu_id=None, importance_type='gain',
                                          interaction_constraints=None,
                                          learning_rate=None, max_delta_step=None,
                                          max_depth=2, min_child_weight=None,
                                          missing=nan, monotone_constraints=None,
                                          n_estimato...
                                          objective='binary:logistic',
                                          random_state=None, reg_alpha=None,
                                          reg_lambda=None, scale_pos_weight=None,
                                          subsample=None, tree_method=None,
                                          validate_parameters=False,
                                          verbosity=None),
                  iid='deprecated', n_jobs=None,
```

```

param_grid={'learning_rate': [0.2, 0.6],
            'subsample': [0.3, 0.6, 0.9]},
pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
scoring='roc_auc', verbose=1)

```

```

[164]: # cv results
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results

```

```

[164]:
   mean_fit_time  std_fit_time  mean_score_time  std_score_time \
0      89.934024    4.977601      0.749709      0.031054
1     117.291133    1.948062      0.781700      0.011541
2     133.174869    3.055986      0.774044      0.020544
3     108.884205    2.397979      0.861049      0.051926
4     126.067211    3.452522      0.857716      0.020887
5     134.505360    0.828143      0.842048      0.029339

   param_learning_rate  param_subsample \
0                   0.2                0.3
1                   0.2                0.6
2                   0.2                0.9
3                   0.6                0.3
4                   0.6                0.6
5                   0.6                0.9

           params  split0_test_score \
0  {'learning_rate': 0.2, 'subsample': 0.3}      0.999911
1  {'learning_rate': 0.2, 'subsample': 0.6}      0.999919
2  {'learning_rate': 0.2, 'subsample': 0.9}      0.999915
3  {'learning_rate': 0.6, 'subsample': 0.3}      0.999985
4  {'learning_rate': 0.6, 'subsample': 0.6}      0.999993
5  {'learning_rate': 0.6, 'subsample': 0.9}      0.999998

   split1_test_score  split2_test_score  mean_test_score  std_test_score \
0      0.999917      0.999916      0.999915      0.000003
1      0.999927      0.999913      0.999919      0.000006
2      0.999931      0.999909      0.999919      0.000009
3      0.999989      0.999976      0.999983      0.000006
4      0.999986      0.999979      0.999986      0.000006
5      0.999990      0.999973      0.999987      0.000010

   rank_test_score  split0_train_score  split1_train_score \
0                6      0.999927      0.999927
1                4      0.999924      0.999935
2                5      0.999921      0.999932
3                3      0.999998      0.999994
4                2      0.999999      0.999998

```

5	1	0.999999	0.999999
---	---	----------	----------

	split2_train_score	mean_train_score	std_train_score
0	0.999925	0.999926	1.298571e-06
1	0.999932	0.999930	4.822947e-06
2	0.999924	0.999925	4.543692e-06
3	0.999999	0.999997	2.188322e-06
4	1.000000	0.999999	7.647958e-07
5	1.000000	1.000000	3.284538e-07

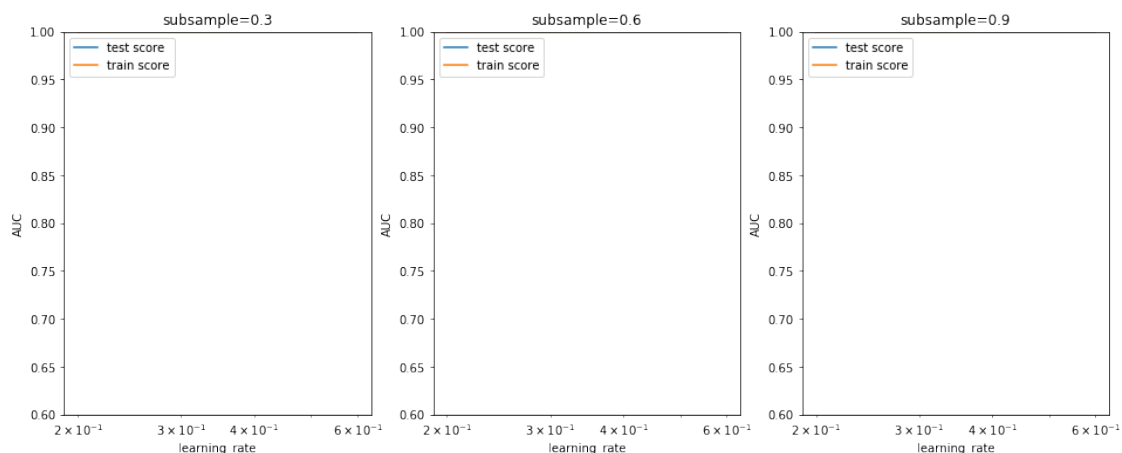
```
[165]: # # plotting
plt.figure(figsize=(16,6))

param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

for n, subsample in enumerate(param_grid['subsample']):

    # subplot 1/n
    plt.subplot(1,len(param_grid['subsample']), n+1)
    df = cv_results[cv_results['param_subsample']==subsample]

    plt.plot(df["param_learning_rate"], df["mean_test_score"])
    plt.plot(df["param_learning_rate"], df["mean_train_score"])
    plt.xlabel('learning_rate')
    plt.ylabel('AUC')
    plt.title("subsample={0}".format(subsample))
    plt.ylim([0.60, 1])
    plt.legend(['test score', 'train score'], loc='upper left')
    plt.xscale('log')
```



Model with optimal hyperparameters We see that the train score almost touches to 1. Among the hyperparameters, we can choose the best parameters as `learning_rate : 0.2` and `subsample: 0.3`

```
[166]: model_cv.best_params_
```

```
[166]: {'learning_rate': 0.6, 'subsample': 0.9}
```

```
[211]: # chosen hyperparameters
params = {'learning_rate': 0.6,
          'max_depth': 2,
          'n_estimators': 200,
          'subsample': 0.9,
          'objective': 'binary:logistic'}

# fit model on training data
xgb_bal_ros_model = XGBClassifier(params = params)
xgb_bal_ros_model.fit(X_train_ros, y_train_ros)
```

```
[211]: XGBClassifier(base_score=0.5, booster=None, colsample_bylevel=1,
                    colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                    importance_type='gain', interaction_constraints=None,
                    learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                    min_child_weight=1, missing=nan, monotone_constraints=None,
                    n_estimators=100, n_jobs=0, num_parallel_tree=1,
                    objective='binary:logistic',
                    params={'learning_rate': 0.6, 'max_depth': 2, 'n_estimators': 200,
                           'objective': 'binary:logistic', 'subsample': 0.9},
                    random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
                    subsample=1, tree_method=None, validate_parameters=False,
                    verbosity=None)
```

Prediction on the train set

```
[212]: # Predictions on the train set
y_train_pred = xgb_bal_ros_model.predict(X_train_ros)
```

```
[213]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_ros, y_train_ros)
print(confusion)
```

```
[[227449    0]
 [    0 227449]]
```

```
[214]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
```

```
FN = confusion[1,0] # false negatives
```

```
[215]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_ros, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 1.0
Sensitivity:- 1.0
Specificity:- 1.0
```

```
[216]: # classification_report
print(classification_report(y_train_ros, y_train_pred))
```

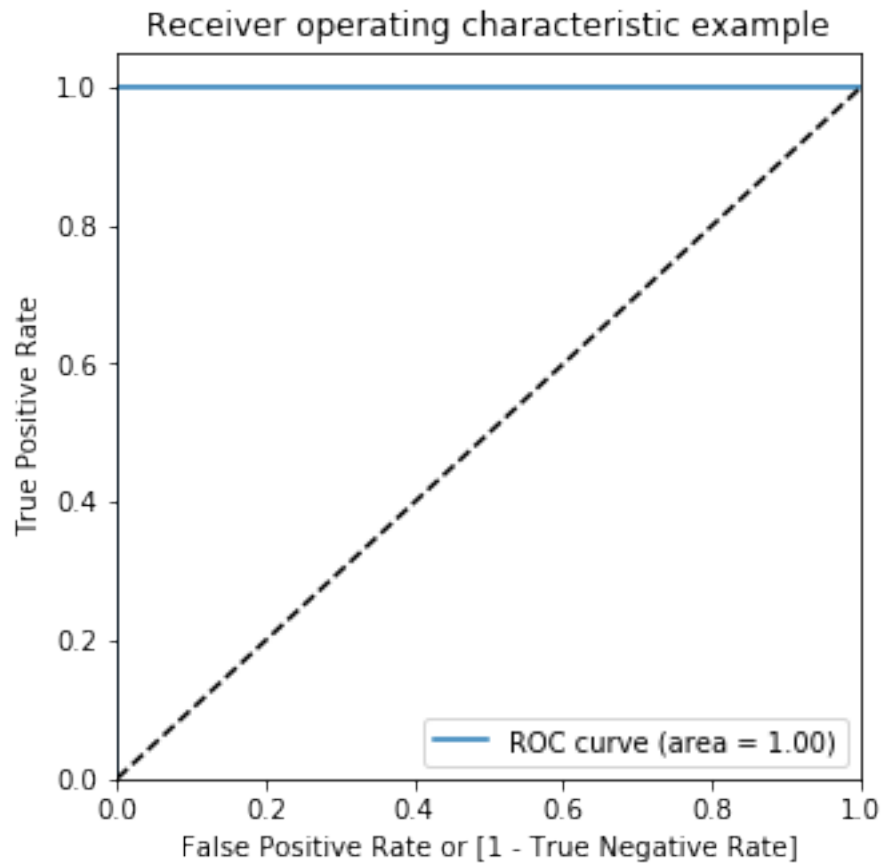
	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	1.00	1.00	1.00	227449
accuracy			1.00	454898
macro avg	1.00	1.00	1.00	454898
weighted avg	1.00	1.00	1.00	454898

```
[217]: # Predicted probability
y_train_pred_proba = xgb_bal_ros_model.predict_proba(X_train_ros)[: ,1]
```

```
[218]: # roc_auc
auc = metrics.roc_auc_score(y_train_ros, y_train_pred_proba)
auc
```

```
[218]: 1.0
```

```
[219]: # Plot the ROC curve
draw_roc(y_train_ros, y_train_pred_proba)
```



Prediction on the test set

```
[223]: # Predictions on the test set
y_test_pred = xgb_bal_ros_model.predict(X_test)
```

```
[224]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56857    9]
 [   19   77]]
```

```
[225]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[226]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9995084442259752
Sensitivity:- 0.8020833333333334
Specificity:- 0.9998417331973412

```
[227]: # classification_report
print(classification_report(y_test, y_test_pred))
```

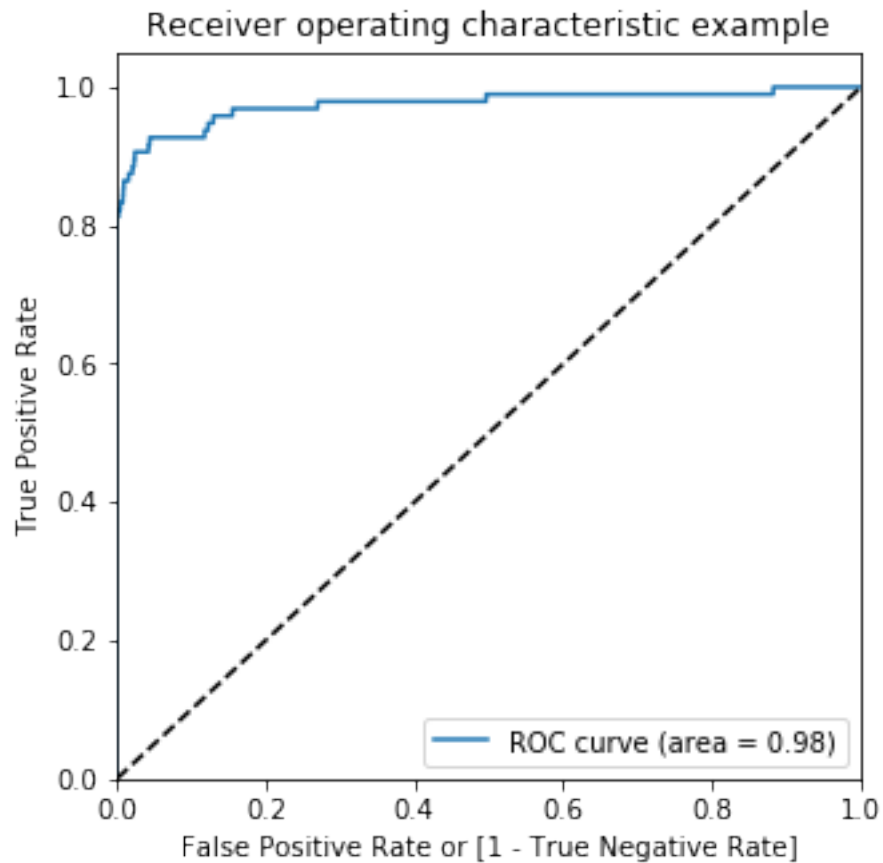
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.90	0.80	0.85	96
accuracy			1.00	56962
macro avg	0.95	0.90	0.92	56962
weighted avg	1.00	1.00	1.00	56962

```
[228]: # Predicted probability
y_test_pred_proba = xgb_bal_ros_model.predict_proba(X_test)[: ,1]
```

```
[229]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[229]: 0.9751521119825555

```
[230]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```

Model summary

- Train set
 - Accuracy = 1.0
 - Sensitivity = 1.0
 - Specificity = 1.0
 - ROC-AUC = 1.0
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.80
 - Specificity = 0.99
 - ROC-AUC = 0.97

4.0.3 Decision Tree

```
[180]: # Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
```

```

}

# Instantiate the grid search model
dtree = DecisionTreeClassifier()

grid_search = GridSearchCV(estimator = dtree,
                           param_grid = param_grid,
                           scoring= 'roc_auc',
                           cv = 3,
                           verbose = 1)

# Fit the grid search to the data
grid_search.fit(X_train_ros,y_train_ros)

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 24 out of 24 | elapsed: 3.2min finished

```

```

[180]: GridSearchCV(cv=3, error_score=nan,
                  estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features=None,
                                                    max_leaf_nodes=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    presort='deprecated',
                                                    random_state=None,
                                                    splitter='best'),
                  iid='deprecated', n_jobs=None,
                  param_grid={'max_depth': range(5, 15, 5),
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50)},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                  scoring='roc_auc', verbose=1)

```

```

[181]: # cv results
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results

```

```

[181]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      6.194021    0.117651      0.089005      1.123916e-07
1      6.149352    0.019800      0.095672      9.428756e-03
2      6.112016    0.013696      0.089672      4.715390e-04

```

3	6.115016	0.025955	0.089672	4.715951e-04
4	9.501210	0.124013	0.093672	1.247235e-03
5	9.553962	0.181357	0.093402	2.803609e-04
6	9.538348	0.164482	0.096734	4.431263e-03
7	9.481282	0.091798	0.092400	1.697113e-03

	param_max_depth	param_min_samples_leaf	param_min_samples_split	\
0	5	50	50	
1	5	50	100	
2	5	100	50	
3	5	100	100	
4	10	50	50	
5	10	50	100	
6	10	100	50	
7	10	100	100	

	params	split0_test_score	\
0	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.990413	
1	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.990413	
2	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.990354	
3	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.990361	
4	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.999553	
5	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.999505	
6	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.999641	
7	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.999567	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.983910	0.991274	0.988532	0.003287	
1	0.983913	0.991299	0.988542	0.003293	
2	0.983844	0.991209	0.988469	0.003289	
3	0.983844	0.991199	0.988468	0.003288	
4	0.999622	0.999577	0.999584	0.000029	
5	0.999611	0.999628	0.999582	0.000054	
6	0.999581	0.999588	0.999603	0.000027	
7	0.999591	0.999586	0.999581	0.000010	

	rank_test_score
0	6
1	5
2	7
3	8
4	2
5	3
6	1
7	4

```
[182]: # Printing the optimal sensitivity score and hyperparameters
print("Best roc_auc:-", grid_search.best_score_)
print(grid_search.best_estimator_)
```

```
Best roc_auc:- 0.9996033327168891
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                        max_depth=10, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=100, min_samples_split=50,
                        min_weight_fraction_leaf=0.0, presort='deprecated',
                        random_state=None, splitter='best')
```

```
[231]: # Model with optimal hyperparameters
dt_bal_ros_model = DecisionTreeClassifier(criterion = "gini",
                                          random_state = 100,
                                          max_depth=10,
                                          min_samples_leaf=100,
                                          min_samples_split=50)

dt_bal_ros_model.fit(X_train_ros, y_train_ros)
```

```
[231]: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                              max_depth=10, max_features=None, max_leaf_nodes=None,
                              min_impurity_decrease=0.0, min_impurity_split=None,
                              min_samples_leaf=100, min_samples_split=50,
                              min_weight_fraction_leaf=0.0, presort='deprecated',
                              random_state=100, splitter='best')
```

Prediction on the train set

```
[232]: # Predictions on the train set
y_train_pred = dt_bal_ros_model.predict(X_train_ros)
```

```
[233]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_ros, y_train_pred)
print(confusion)
```

```
[[225914  1535]
 [      0 227449]]
```

```
[234]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[235]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_ros, y_train_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9966256171713219
Sensitivity:- 1.0
Specificity:- 0.9932512343426438

```
[236]: # classification_report
print(classification_report(y_train_ros, y_train_pred))
```

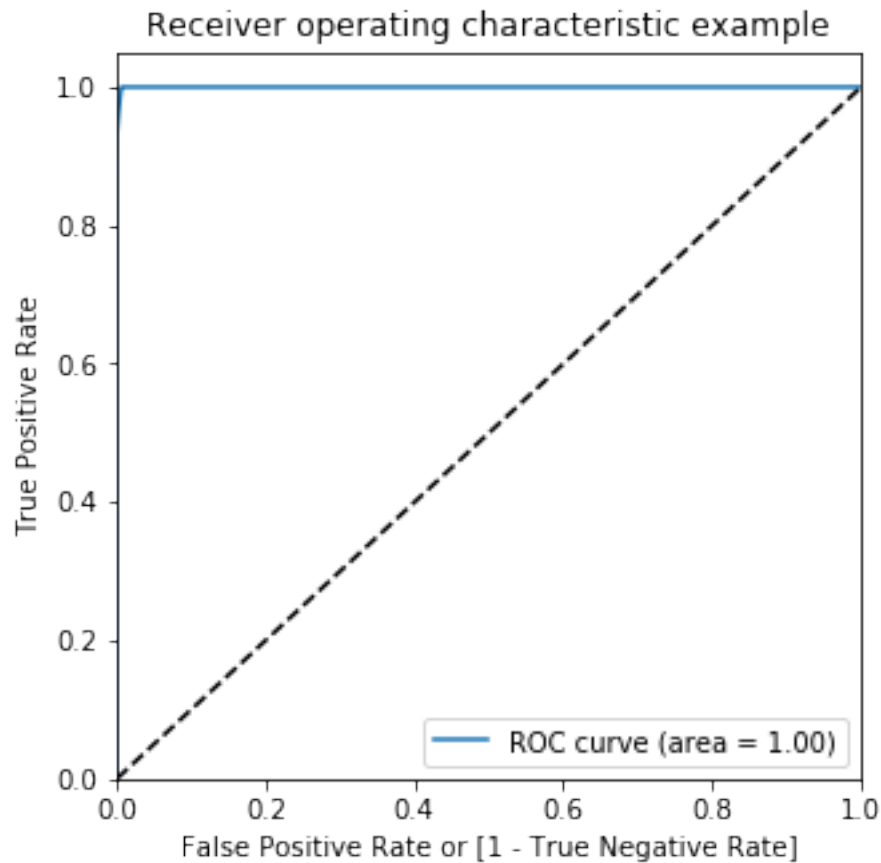
	precision	recall	f1-score	support
0	1.00	0.99	1.00	227449
1	0.99	1.00	1.00	227449
accuracy			1.00	454898
macro avg	1.00	1.00	1.00	454898
weighted avg	1.00	1.00	1.00	454898

```
[237]: # Predicted probability
y_train_pred_proba = dt_bal_ros_model.predict_proba(X_train_ros)[: ,1]
```

```
[238]: # roc_auc
auc = metrics.roc_auc_score(y_train_ros, y_train_pred_proba)
auc
```

[238]: 0.9997642505020377

```
[239]: # Plot the ROC curve
draw_roc(y_train_ros, y_train_pred_proba)
```



Prediction on the test set

```
[240]: # Predictions on the test set
y_test_pred = dt_bal_ros_model.predict(X_test)
```

```
[241]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56431  435]
 [   20   76]]
```

```
[242]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[243]: # Accuracy
print("Accuracy:-",metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9920122186720972
 Sensitivity:- 0.7916666666666666
 Specificity:- 0.9923504378714874

```
[244]: # classification_report
print(classification_report(y_test, y_test_pred))
```

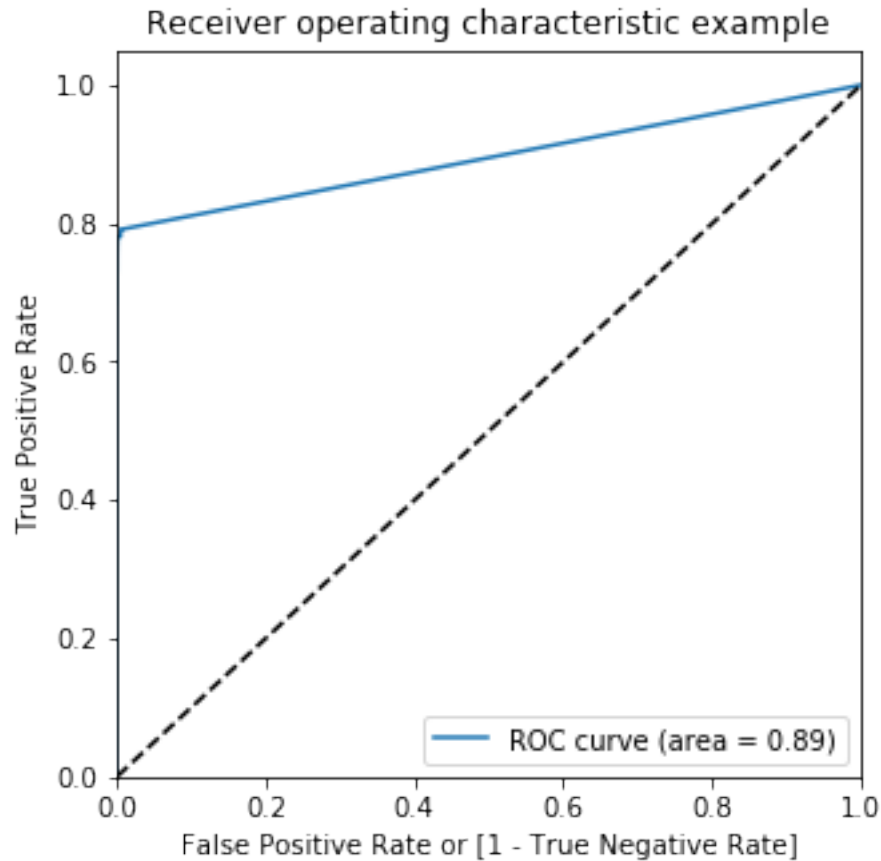
	precision	recall	f1-score	support
0	1.00	0.99	1.00	56866
1	0.15	0.79	0.25	96
accuracy			0.99	56962
macro avg	0.57	0.89	0.62	56962
weighted avg	1.00	0.99	0.99	56962

```
[245]: # Predicted probability
y_test_pred_proba = dt_bal_ros_model.predict_proba(X_test)[: ,1]
```

```
[246]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[246]: 0.8948251151830618

```
[247]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 1.0
 - Specificity = 0.99
 - ROC-AUC = 0.99
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.79
 - Specificity = 0.99
 - ROC-AUC = 0.90

4.1 SMOTE (Synthetic Minority Oversampling Technique)

We are creating synthetic samples by doing upsampling using SMOTE(Synthetic Minority Oversampling Technique).

```
[68]: # Importing SMOTE
      from imblearn.over_sampling import SMOTE
```



```
[69]: # Instantiate SMOTE
sm = SMOTE(random_state=27)
# Fitting SMOTE to the train set
X_train_smote, y_train_smote = sm.fit_sample(X_train, y_train)
```

```
[70]: print('Before SMOTE oversampling X_train shape=',X_train.shape)
print('After SMOTE oversampling X_train shape=',X_train_smote.shape)
```

Before SMOTE oversampling X_train shape= (227845, 29)

After SMOTE oversampling X_train shape= (454898, 29)

4.1.1 Logistic Regression

```
[ ]: # Creating KFold object with 5 splits
folds = KFold(n_splits=5, shuffle=True, random_state=4)

# Specify params
params = {"C": [0.01, 0.1, 1, 10, 100, 1000]}

# Specifying score as roc-auc
model_cv = GridSearchCV(estimator = LogisticRegression(),
                        param_grid = params,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# Fit the model
model_cv.fit(X_train_smote, y_train_smote)
```

Fitting 5 folds for each of 6 candidates, totalling 30 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 30 out of 30 | elapsed: 1.5min finished

```
[ ]: GridSearchCV(cv=KFold(n_splits=5, random_state=4, shuffle=True),
                error_score=nan,
                estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
                                              fit_intercept=True,
                                              intercept_scaling=1, l1_ratio=None,
                                              max_iter=100, multi_class='auto',
                                              n_jobs=None, penalty='l2',
                                              random_state=None, solver='lbfgs',
                                              tol=0.0001, verbose=0,
                                              warm_start=False),
                iid='deprecated', n_jobs=None,
                param_grid={'C': [0.01, 0.1, 1, 10, 100, 1000]},
                pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                scoring='roc_auc', verbose=1)
```

```
[ ]: # results of grid search CV
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[ ]:  mean_fit_time  std_fit_time  mean_score_time  std_score_time  param_C  \
0      2.517373    0.063613    0.068640    0.012480    0.01
1      2.580687    0.139180    0.065641    0.006184    0.1
2      2.673410    0.107579    0.065920    0.006089     1
3      2.650617    0.100909    0.065520    0.006240    10
4      2.693168    0.148317    0.065520    0.006240   100
5      2.745892    0.157325    0.056160    0.007642  1000

      params  split0_test_score  split1_test_score  split2_test_score  \
0  {'C': 0.01}          0.989805          0.989796          0.989484
1  {'C': 0.1}          0.989834          0.989807          0.989488
2   {'C': 1}          0.989836          0.989807          0.989486
3   {'C': 10}         0.989836          0.989807          0.989486
4   {'C': 100}        0.989836          0.989807          0.989486
5   {'C': 1000}       0.989836          0.989807          0.989486

      split3_test_score  split4_test_score  mean_test_score  std_test_score  \
0          0.989631          0.989910          0.989725          0.000150
1          0.989632          0.989942          0.989741          0.000161
2          0.989630          0.989944          0.989741          0.000162
3          0.989630          0.989945          0.989741          0.000163
4          0.989630          0.989945          0.989741          0.000163
5          0.989630          0.989945          0.989741          0.000163

      rank_test_score  split0_train_score  split1_train_score  \
0              6          0.989758          0.989666
1              1          0.989780          0.989686
2              2          0.989781          0.989687
3              5          0.989781          0.989687
4              3          0.989781          0.989687
5              4          0.989781          0.989687

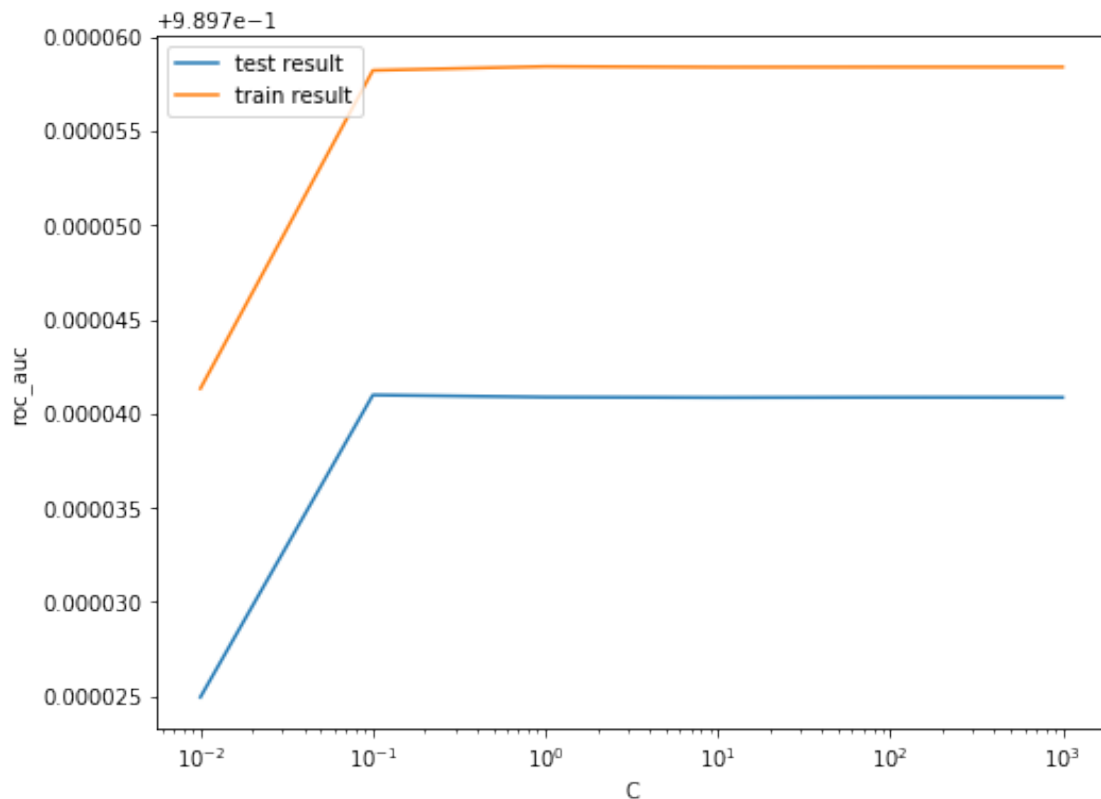
      split2_train_score  split3_train_score  split4_train_score  \
0          0.989760          0.989841          0.989682
1          0.989772          0.989853          0.989700
2          0.989772          0.989852          0.989701
3          0.989772          0.989852          0.989701
4          0.989772          0.989852          0.989701
5          0.989772          0.989852          0.989701

      mean_train_score  std_train_score
0          0.989741          0.000063
1          0.989758          0.000060
```

2	0.989758	0.000060
3	0.989758	0.000060
4	0.989758	0.000060
5	0.989758	0.000060

```
[ ]: # plot of C versus train and validation scores
```

```
plt.figure(figsize=(8, 6))
plt.plot(cv_results['param_C'], cv_results['mean_test_score'])
plt.plot(cv_results['param_C'], cv_results['mean_train_score'])
plt.xlabel('C')
plt.ylabel('roc_auc')
plt.legend(['test result', 'train result'], loc='upper left')
plt.xscale('log')
```



```
[ ]: # Best score with best C
```

```
best_score = model_cv.best_score_
best_C = model_cv.best_params_['C']
```

```
print(" The highest test roc_auc is {0} at C = {1}".format(best_score, best_C))
```

The highest test roc_auc is 0.9897409900830768 at C = 0.1

Logistic regression with optimal C

```
[71]: # Instantiate the model with best C
logistic_bal_smote = LogisticRegression(C=0.1)
```

```
[72]: # Fit the model on the train set
logistic_bal_smote_model = logistic_bal_smote.fit(X_train_smote, y_train_smote)
```

Prediction on the train set

```
[73]: # Predictions on the train set
y_train_pred = logistic_bal_smote_model.predict(X_train_smote)
```

```
[74]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_smote, y_train_pred)
print(confusion)
```

```
[[221911  5538]
 [ 17693 209756]]
```

```
[75]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[76]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_smote, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

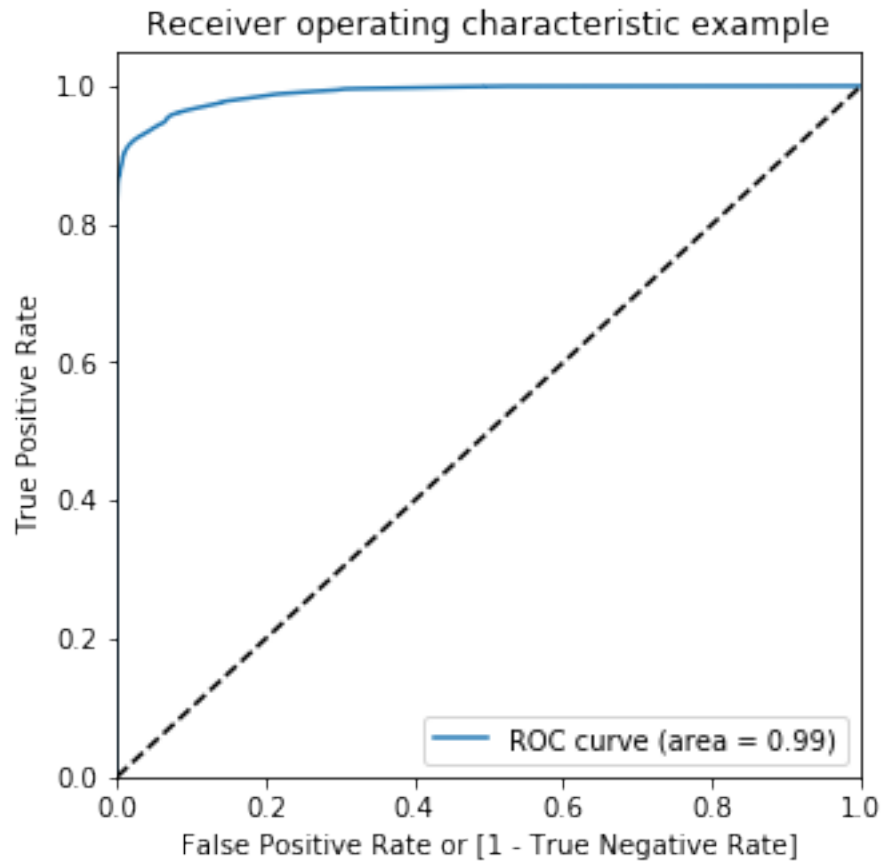
```
Accuracy:- 0.9489314087993352
Sensitivity:- 0.9222111330452102
Specificity:- 0.9756516845534603
```

```
[77]: # classification_report
print(classification_report(y_train_smote, y_train_pred))
```

	precision	recall	f1-score	support
0	0.93	0.98	0.95	227449
1	0.97	0.92	0.95	227449
accuracy			0.95	454898
macro avg	0.95	0.95	0.95	454898
weighted avg	0.95	0.95	0.95	454898

```
[89]: # Predicted probability
y_train_pred_proba_log_bal_smote = logistic_bal_smote_model.
      ↪predict_proba(X_train_smote)[: ,1]
```

```
[90]: # Plot the ROC curve
draw_roc(y_train_smote, y_train_pred_proba_log_bal_smote)
```



Prediction on the test set

```
[80]: # Prediction on the test set
y_test_pred = logistic_bal_smote_model.predict(X_test)
```

```
[81]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55416 1450]
 [  10   86]]
```

```
[82]: TP = confusion[1,1] # true positive
      TN = confusion[0,0] # true negatives
      FP = confusion[0,1] # false positives
      FN = confusion[1,0] # false negatives
```

```
[83]: # Accuracy
      print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

      # Sensitivity
      print("Sensitivity:-", TP / float(TP+FN))

      # Specificity
      print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9743688774972789
 Sensitivity:- 0.8958333333333334
 Specificity:- 0.9745014595716245

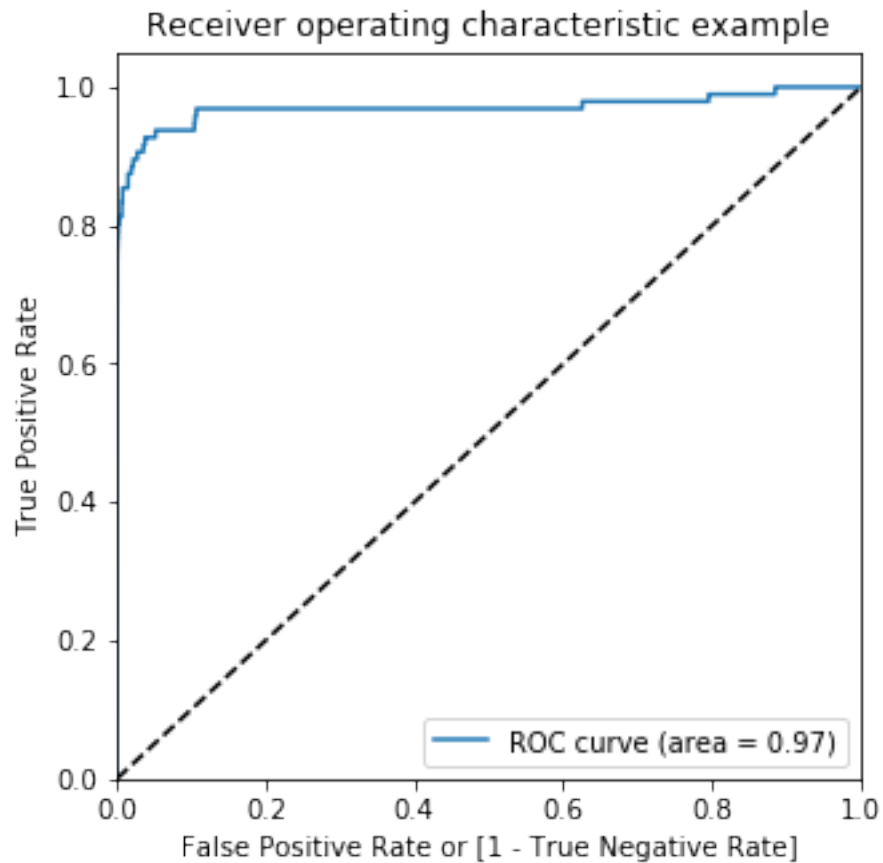
```
[84]: # classification_report
      print(classification_report(y_test, y_test_pred))
```

	precision	recall	f1-score	support
0	1.00	0.97	0.99	56866
1	0.06	0.90	0.11	96
accuracy			0.97	56962
macro avg	0.53	0.94	0.55	56962
weighted avg	1.00	0.97	0.99	56962

ROC on the test set

```
[85]: # Predicted probability
      y_test_pred_proba = logistic_bal_smote_model.predict_proba(X_test)[: ,1]
```

```
[86]: # Plot the ROC curve
      draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.95
 - Sensitivity = 0.92
 - Specificity = 0.98
 - ROC = 0.99
- Test set
 - Accuracy = 0.97
 - Sensitivity = 0.90
 - Specificity = 0.99
 - ROC = 0.97

4.1.2 XGBoost

```
[ ]: # hyperparameter tuning with XGBoost  
  
# creating a KFold object  
folds = 3
```

```

# specify range of hyperparameters
param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

# specify model
xgb_model = XGBClassifier(max_depth=2, n_estimators=200)

# set up GridSearchCV()
model_cv = GridSearchCV(estimator = xgb_model,
                        param_grid = param_grid,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# fit the model
model_cv.fit(X_train_smote, y_train_smote)

```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 45.8min finished

```

```

[ ]: GridSearchCV(cv=3, error_score=nan,
                  estimator=XGBClassifier(base_score=None, booster=None,
                                          colsample_bylevel=None,
                                          colsample_bynode=None,
                                          colsample_bytree=None, gamma=None,
                                          gpu_id=None, importance_type='gain',
                                          interaction_constraints=None,
                                          learning_rate=None, max_delta_step=None,
                                          max_depth=2, min_child_weight=None,
                                          missing=nan, monotone_constraints=None,
                                          n_estimators=200,
                                          objective='binary:logistic',
                                          random_state=None, reg_alpha=None,
                                          reg_lambda=None, scale_pos_weight=None,
                                          subsample=None, tree_method=None,
                                          validate_parameters=False,
                                          verbosity=None),
                  iid='deprecated', n_jobs=None,
                  param_grid={'learning_rate': [0.2, 0.6],
                              'subsample': [0.3, 0.6, 0.9]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)

```



```
[ ]: # cv results
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[ ]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      120.679338      2.195478      0.847048      0.012833
1      154.374830      1.994512      0.847715      0.007040
2      173.991618      0.672700      0.816047      0.020834
3      122.247942      0.657327      0.870037      0.022553
4      153.590355      1.449181      0.857216      0.032023
5      175.839443      2.079717      0.829144      0.007553

   param_learning_rate  param_subsample  \
0                0.2                0.3
1                0.2                0.6
2                0.2                0.9
3                0.6                0.3
4                0.6                0.6
5                0.6                0.9

                                params  split0_test_score  \
0  {'learning_rate': 0.2, 'subsample': 0.3}      0.999645
1  {'learning_rate': 0.2, 'subsample': 0.6}      0.999671
2  {'learning_rate': 0.2, 'subsample': 0.9}      0.999665
3  {'learning_rate': 0.6, 'subsample': 0.3}      0.999956
4  {'learning_rate': 0.6, 'subsample': 0.6}      0.999953
5  {'learning_rate': 0.6, 'subsample': 0.9}      0.999970

   split1_test_score  split2_test_score  mean_test_score  std_test_score  \
0      0.999753      0.999685      0.999694      0.000045
1      0.999738      0.999652      0.999687      0.000037
2      0.999735      0.999648      0.999683      0.000038
3      0.999950      0.999953      0.999953      0.000002
4      0.999962      0.999959      0.999958      0.000004
5      0.999958      0.999951      0.999960      0.000008

   rank_test_score  split0_train_score  split1_train_score  \
0                4      0.999718      0.999736
1                5      0.999733      0.999731
2                6      0.999720      0.999723
3                3      0.999979      0.999972
4                2      0.999980      0.999981
5                1      0.999985      0.999981

   split2_train_score  mean_train_score  std_train_score
0      0.999720      0.999725      0.000008
1      0.999697      0.999721      0.000017
```

2	0.999720	0.999721	0.000001
3	0.999977	0.999976	0.000003
4	0.999984	0.999982	0.000002
5	0.999977	0.999981	0.000003

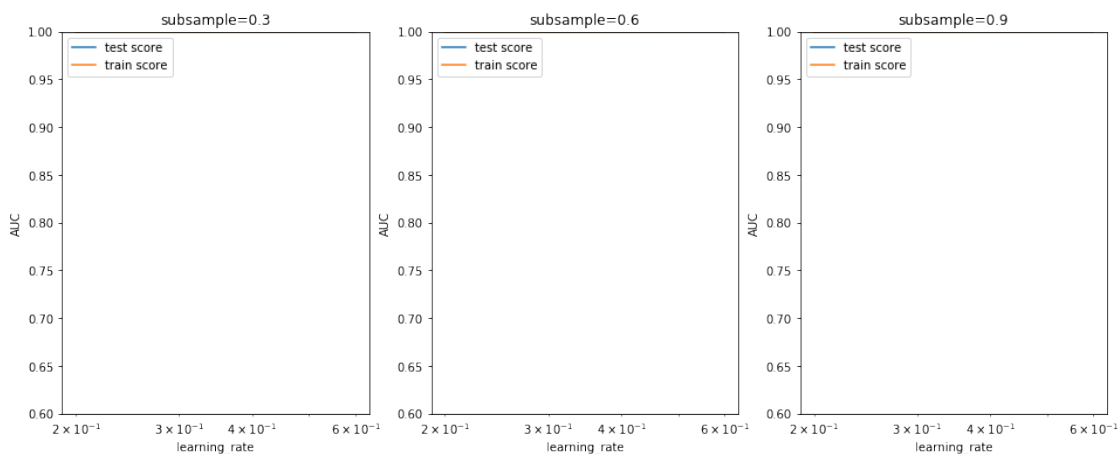
```
[ ]: ## plotting
plt.figure(figsize=(16,6))

param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

for n, subsample in enumerate(param_grid['subsample']):

    # subplot 1/n
    plt.subplot(1,len(param_grid['subsample']), n+1)
    df = cv_results[cv_results['param_subsample']==subsample]

    plt.plot(df["param_learning_rate"], df["mean_test_score"])
    plt.plot(df["param_learning_rate"], df["mean_train_score"])
    plt.xlabel('learning_rate')
    plt.ylabel('AUC')
    plt.title("subsample={0}".format(subsample))
    plt.ylim([0.60, 1])
    plt.legend(['test score', 'train score'], loc='upper left')
    plt.xscale('log')
```



Model with optimal hyperparameters We see that the train score almost touches to 1. Among the hyperparameters, we can choose the best parameters as learning_rate : 0.2 and subsample: 0.3

```
[ ]: model_cv.best_params_
```

```
[ ]: {'learning_rate': 0.6, 'subsample': 0.9}
```

```
[267]: # chosen hyperparameters
# 'objective': 'binary:logistic' outputs probability rather than label, which we
# need for calculating auc
params = {'learning_rate': 0.6,
          'max_depth': 2,
          'n_estimators': 200,
          'subsample': 0.9,
          'objective': 'binary:logistic'}

# fit model on training data
xgb_bal_smote_model = XGBClassifier(params = params)
xgb_bal_smote_model.fit(X_train_smote, y_train_smote)
```

```
[267]: XGBClassifier(base_score=0.5, booster=None, colsample_bylevel=1,
                  colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                  importance_type='gain', interaction_constraints=None,
                  learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                  min_child_weight=1, missing=nan, monotone_constraints=None,
                  n_estimators=100, n_jobs=0, num_parallel_tree=1,
                  objective='binary:logistic',
                  params={'learning_rate': 0.6, 'max_depth': 2, 'n_estimators': 200,
                          'objective': 'binary:logistic', 'subsample': 0.9},
                  random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
                  subsample=1, tree_method=None, validate_parameters=False,
                  verbosity=None)
```

Prediction on the train set

```
[268]: # Predictions on the train set
y_train_pred = xgb_bal_smote_model.predict(X_train_smote)
```

```
[269]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_smote, y_train_pred)
print(confusion)
```

```
[[227447      2]
 [      0 227449]]
```

```
[270]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[271]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_smote, y_train_pred))
```

```
# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9999956034099952
Sensitivity:- 1.0
Specificity:- 0.9999912068199904

```
[272]: # classification_report
print(classification_report(y_train_smote, y_train_pred))
```

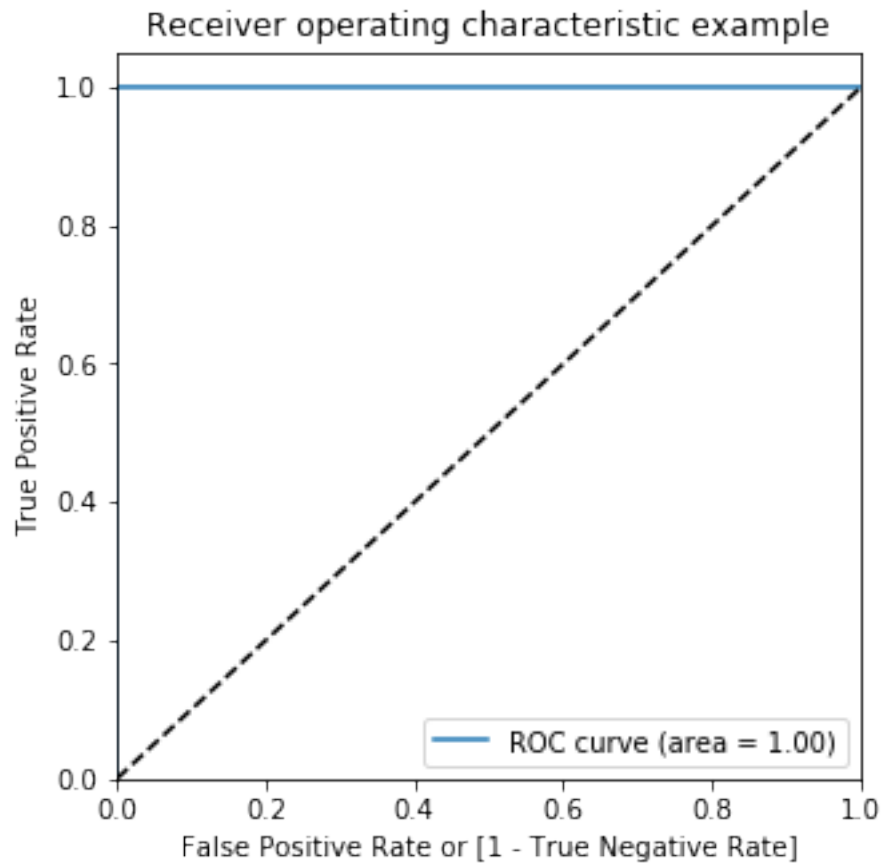
	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	1.00	1.00	1.00	227449
accuracy			1.00	454898
macro avg	1.00	1.00	1.00	454898
weighted avg	1.00	1.00	1.00	454898

```
[273]: # Predicted probability
y_train_pred_proba = xgb_bal_smote_model.predict_proba(X_train_smote)[: ,1]
```

```
[274]: # roc_auc
auc = metrics.roc_auc_score(y_train_smote, y_train_pred_proba)
auc
```

[274]: 1.0

```
[275]: # Plot the ROC curve
draw_roc(y_train_smote, y_train_pred_proba)
```



Prediction on the test set

```
[276]: # Predictions on the test set
y_test_pred = xgb_bal_smote_model.predict(X_test)
```

```
[277]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56839   27]
 [   20   76]]
```

```
[278]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[279]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9991748885221726
Sensitivity:- 0.7916666666666666
Specificity:- 0.9995251995920234

```
[280]: # classification_report
print(classification_report(y_test, y_test_pred))
```

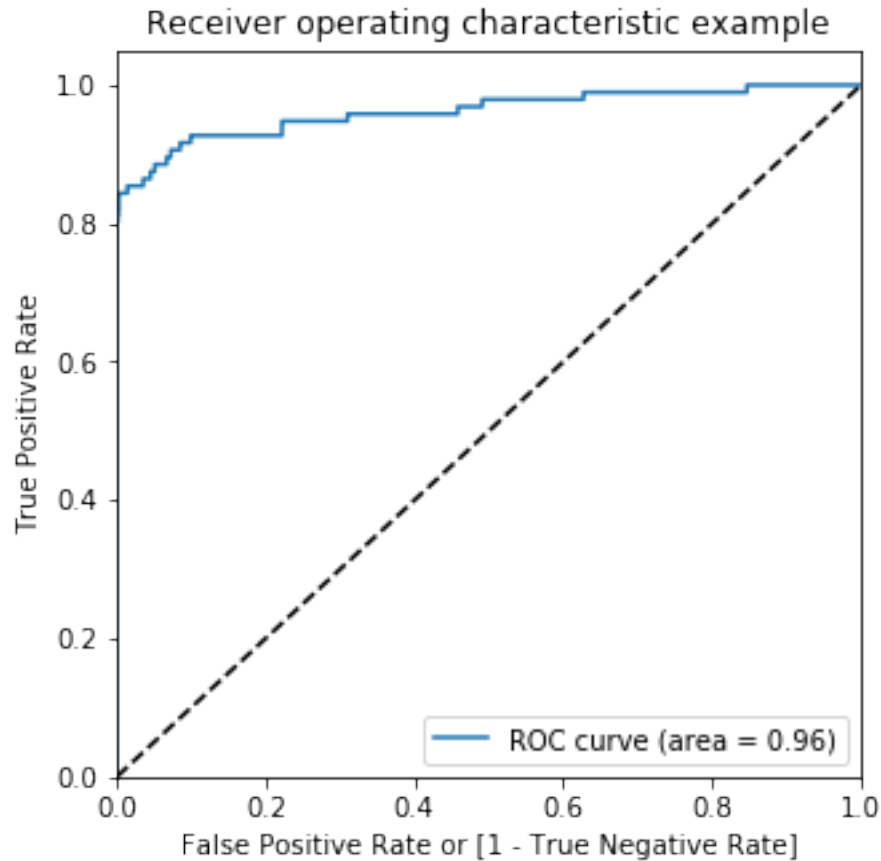
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.74	0.79	0.76	96
accuracy			1.00	56962
macro avg	0.87	0.90	0.88	56962
weighted avg	1.00	1.00	1.00	56962

```
[281]: # Predicted probability
y_test_pred_proba = xgb_bal_smote_model.predict_proba(X_test)[: ,1]
```

```
[282]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[282]: 0.9618437789423088

```
[283]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 1.0
 - Specificity = 0.99
 - ROC-AUC = 1.0
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.79
 - Specificity = 0.99
 - ROC-AUC = 0.96

Overall, the model is performing well in the test set, what it had learnt from the train set.

4.1.3 Decision Tree

```
[ ]: # Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
```

```
'min_samples_split': range(50, 150, 50),
}
```

```
# Instantiate the grid search model
```

```
dtree = DecisionTreeClassifier()
```

```
grid_search = GridSearchCV(estimator = dtree,
                           param_grid = param_grid,
                           scoring= 'roc_auc',
                           cv = 3,
                           verbose = 1)
```

```
# Fit the grid search to the data
```

```
grid_search.fit(X_train_smote,y_train_smote)
```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 24 out of 24 | elapsed: 5.9min finished

```
[ ]: GridSearchCV(cv=3, error_score=nan,
                  estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features=None,
                                                    max_leaf_nodes=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    presort='deprecated',
                                                    random_state=None,
                                                    splitter='best'),
                  iid='deprecated', n_jobs=None,
                  param_grid={'max_depth': range(5, 15, 5),
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50)},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                  scoring='roc_auc', verbose=1)
```

```
[ ]: # cv results
```

```
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results
```

```
[ ]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      9.971941      0.245916      0.091339      0.001700
1      9.798227      0.043148      0.090672      0.002495
```


2	9.804227	0.061308	0.090672	0.001247
3	10.006914	0.237793	0.093204	0.000280
4	22.208819	2.856627	0.096871	0.002617
5	19.618377	0.749607	0.102538	0.008521
6	18.075125	0.167592	0.096537	0.002231
7	18.079367	0.254541	0.103004	0.002159

	param_max_depth	param_min_samples_leaf	param_min_samples_split	\
0	5	50	50	
1	5	50	100	
2	5	100	50	
3	5	100	100	
4	10	50	50	
5	10	50	100	
6	10	100	50	
7	10	100	100	

	params	split0_test_score	\
0	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.986116	
1	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.986116	
2	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.986081	
3	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.986069	
4	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.998130	
5	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.998157	
6	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.998118	
7	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.998088	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.985690	0.984837	0.985548	0.000532	
1	0.985690	0.984839	0.985548	0.000531	
2	0.985637	0.984771	0.985496	0.000544	
3	0.985639	0.984771	0.985493	0.000540	
4	0.998119	0.997982	0.998077	0.000067	
5	0.998135	0.997940	0.998077	0.000097	
6	0.998048	0.997934	0.998033	0.000076	
7	0.998075	0.997923	0.998029	0.000075	

	rank_test_score
0	6
1	5
2	7
3	8
4	2
5	1
6	3
7	4

```
[ ]: # Printing the optimal sensitivity score and hyperparameters
print("Best roc_auc:-", grid_search.best_score_)
print(grid_search.best_estimator_)
```

```
Best roc_auc:- 0.9980773622123168
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                      max_depth=10, max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=50, min_samples_split=100,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')
```

```
[284]: # Model with optimal hyperparameters
dt_bal_smote_model = DecisionTreeClassifier(criterion = "gini",
                                           random_state = 100,
                                           max_depth=10,
                                           min_samples_leaf=50,
                                           min_samples_split=100)

dt_bal_smote_model.fit(X_train_smote, y_train_smote)
```

```
[284]: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                             max_depth=10, max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=50, min_samples_split=100,
                             min_weight_fraction_leaf=0.0, presort='deprecated',
                             random_state=100, splitter='best')
```

Prediction on the train set

```
[285]: # Predictions on the train set
y_train_pred = dt_bal_smote_model.predict(X_train_smote)
```

```
[286]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_smote, y_train_pred)
print(confusion)
```

```
[[223809   3640]
 [  2374 225075]]
```

```
[287]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[288]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_smote, y_train_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9867794538555896
Sensitivity:- 0.9895624953286232
Specificity:- 0.9839964123825561

```
[290]: # classification_report
print(classification_report(y_train_smote, y_train_pred))
```

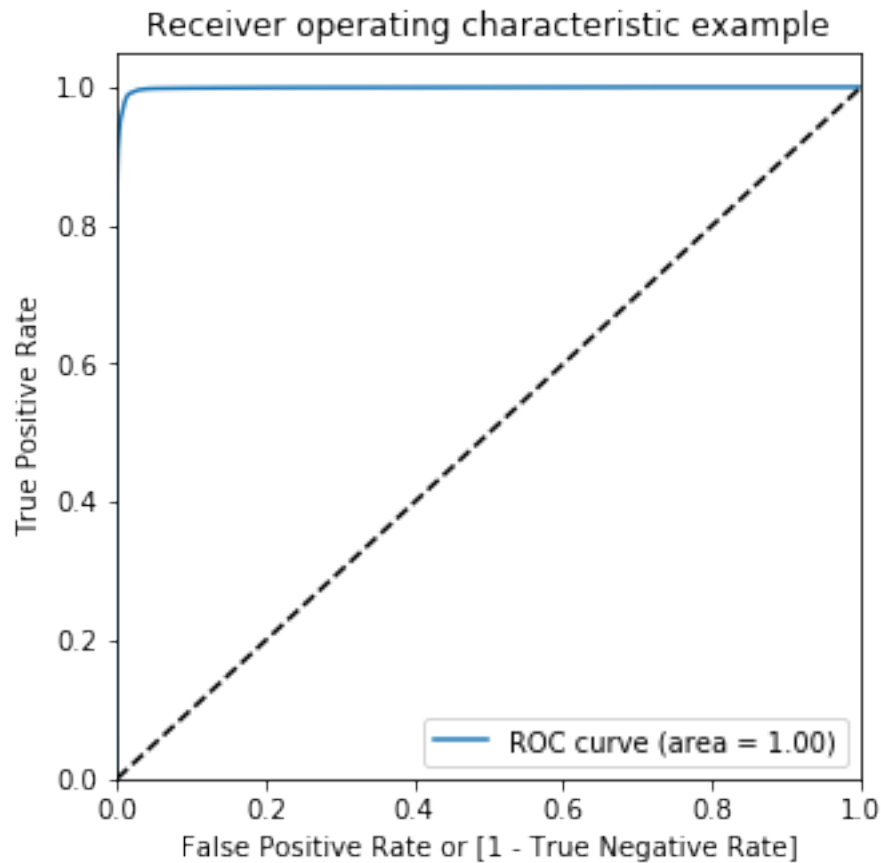
	precision	recall	f1-score	support
0	0.99	0.98	0.99	227449
1	0.98	0.99	0.99	227449
accuracy			0.99	454898
macro avg	0.99	0.99	0.99	454898
weighted avg	0.99	0.99	0.99	454898

```
[291]: # Predicted probability
y_train_pred_proba = dt_bal_smote_model.predict_proba(X_train_smote)[: ,1]
```

```
[292]: # roc_auc
auc = metrics.roc_auc_score(y_train_smote, y_train_pred_proba)
auc
```

[292]: 0.9986355757920081

```
[294]: # Plot the ROC curve
draw_roc(y_train_smote, y_train_pred_proba)
```



Prediction on the test set

```
[295]: # Predictions on the test set
y_test_pred = dt_bal_smote_model.predict(X_test)
```

```
[296]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[55852  1014]
 [    19    77]]
```

```
[297]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[298]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9818651030511569
Sensitivity:- 0.8020833333333334
Specificity:- 0.9821686069004326

```
[299]: # classification_report
print(classification_report(y_test, y_test_pred))
```

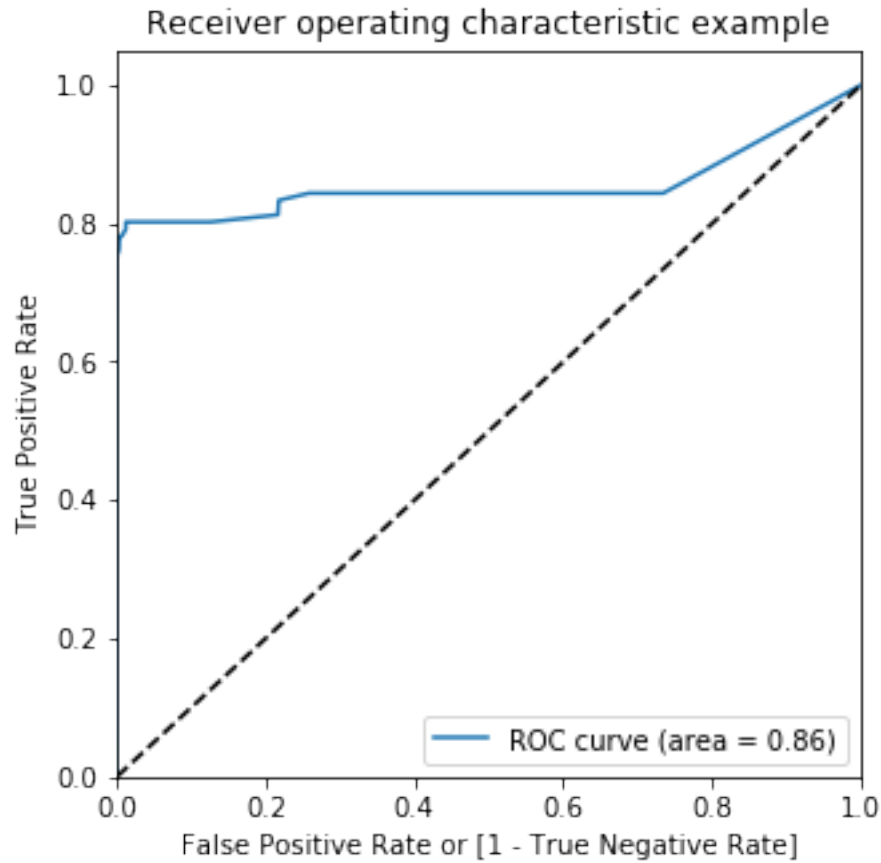
	precision	recall	f1-score	support
0	1.00	0.98	0.99	56866
1	0.07	0.80	0.13	96
accuracy			0.98	56962
macro avg	0.54	0.89	0.56	56962
weighted avg	1.00	0.98	0.99	56962

```
[300]: # Predicted probability
y_test_pred_proba = dt_bal_smote_model.predict_proba(X_test)[: ,1]
```

```
[301]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[301]: 0.8551876157692353

```
[302]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 0.99
 - Specificity = 0.98
 - ROC-AUC = 0.99
- Test set
 - Accuracy = 0.98
 - Sensitivity = 0.80
 - Specificity = 0.98
 - ROC-AUC = 0.86

4.2 AdaSyn (Adaptive Synthetic Sampling)

```
[303]: # Importing adasyn
from imblearn.over_sampling import ADASYN
```

```
[304]: # Instantiate adasyn
ada = ADASYN(random_state=0)
```

```
X_train_adasyn, y_train_adasyn = ada.fit_resample(X_train, y_train)
```

```
[305]: # Before sampling class distribution
print('Before sampling class distribution:-', Counter(y_train))
# new class distribution
print('New class distribution:-', Counter(y_train_adasyn))
```

Before sampling class distribution:- Counter({0: 227449, 1: 396})

New class distribution:- Counter({0: 227449, 1: 227448})

4.2.1 Logistic Regression

```
[238]: # Creating KFold object with 3 splits
folds = KFold(n_splits=3, shuffle=True, random_state=4)

# Specify params
params = {"C": [0.01, 0.1, 1, 10, 100, 1000]}

# Specifying score as roc-auc
model_cv = GridSearchCV(estimator = LogisticRegression(),
                        param_grid = params,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# Fit the model
model_cv.fit(X_train_adasyn, y_train_adasyn)
```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 42.9s finished

```
[238]: GridSearchCV(cv=KFold(n_splits=3, random_state=4, shuffle=True),
                  error_score=nan,
                  estimator=LogisticRegression(C=1.0, class_weight=None, dual=False,
                                                fit_intercept=True,
                                                intercept_scaling=1, l1_ratio=None,
                                                max_iter=100, multi_class='auto',
                                                n_jobs=None, penalty='l2',
                                                random_state=None, solver='lbfgs',
                                                tol=0.0001, verbose=0,
                                                warm_start=False),
                  iid='deprecated', n_jobs=None,
                  param_grid={'C': [0.01, 0.1, 1, 10, 100, 1000]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)
```

```
[239]: # results of grid search CV
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[239]:
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C \
0	1.919777	0.172030	0.106673	0.021640	0.01
1	1.930110	0.079116	0.090672	0.002495	0.1
2	2.188786	0.195766	0.093672	0.004497	1
3	2.356865	0.238141	0.103206	0.009386	10
4	2.035450	0.107419	0.091672	0.002357	100
5	2.046450	0.091060	0.094005	0.007119	1000

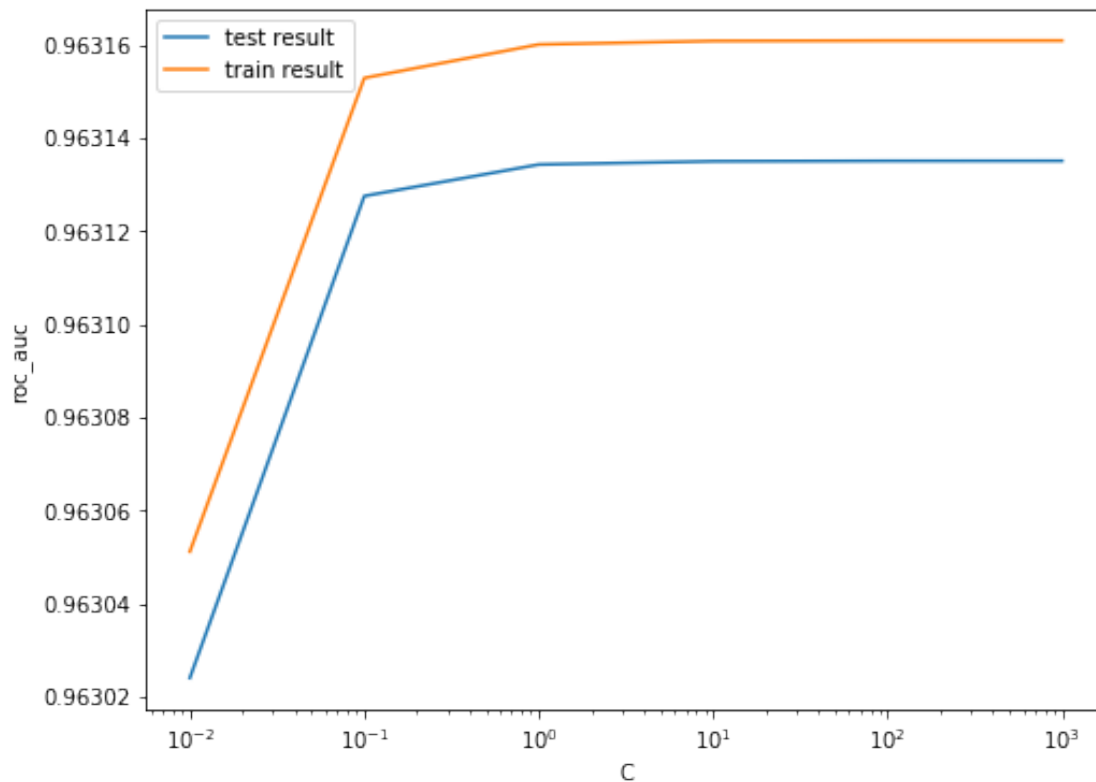
	params	split0_test_score	split1_test_score	split2_test_score \
0	{'C': 0.01}	0.963472	0.962327	0.963273
1	{'C': 0.1}	0.963578	0.962435	0.963370
2	{'C': 1}	0.963585	0.962442	0.963376
3	{'C': 10}	0.963585	0.962443	0.963377
4	{'C': 100}	0.963585	0.962443	0.963377
5	{'C': 1000}	0.963585	0.962443	0.963377

	mean_test_score	std_test_score	rank_test_score	split0_train_score \
0	0.963024	0.000499	6	0.962770
1	0.963128	0.000497	5	0.962881
2	0.963134	0.000497	4	0.962890
3	0.963135	0.000496	3	0.962891
4	0.963135	0.000496	2	0.962891
5	0.963135	0.000496	1	0.962891

	split1_train_score	split2_train_score	mean_train_score	std_train_score
0	0.963211	0.963172	0.963051	0.000199
1	0.963305	0.963272	0.963153	0.000192
2	0.963312	0.963278	0.963160	0.000191
3	0.963312	0.963279	0.963161	0.000191
4	0.963312	0.963279	0.963161	0.000191
5	0.963312	0.963279	0.963161	0.000191

```
[240]: # plot of C versus train and validation scores

plt.figure(figsize=(8, 6))
plt.plot(cv_results['param_C'], cv_results['mean_test_score'])
plt.plot(cv_results['param_C'], cv_results['mean_train_score'])
plt.xlabel('C')
plt.ylabel('roc_auc')
plt.legend(['test result', 'train result'], loc='upper left')
plt.xscale('log')
```

```
[241]: # Best score with best C
best_score = model_cv.best_score_
best_C = model_cv.best_params_['C']

print(" The highest test roc_auc is {0} at C = {1}".format(best_score, best_C))
```

The highest test roc_auc is 0.9631351482818916 at C = 1000

Logistic regression with optimal C

```
[306]: # Instantiate the model with best C
logistic_bal_adasyn = LogisticRegression(C=1000)
```

```
[307]: # Fit the model on the train set
logistic_bal_adasyn_model = logistic_bal_adasyn.fit(X_train_adasyn,
↳ y_train_adasyn)
```

Prediction on the train set

```
[308]: # Predictions on the train set
y_train_pred = logistic_bal_adasyn_model.predict(X_train_adasyn)
```

```
[309]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_adasyn, y_train_pred)
print(confusion)
```

```
[[207019  20430]
 [ 31286 196162]]
```

```
[310]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[311]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_adasyn, y_train_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))

# F1 score
print("F1-Score:-", f1_score(y_train_adasyn, y_train_pred))
```

```
Accuracy:- 0.8863127257379143
Sensitivity:- 0.862447680348915
Specificity:- 0.9101776662020936
F1-Score:- 0.8835330150436899
```

```
[312]: # classification_report
print(classification_report(y_train_adasyn, y_train_pred))
```

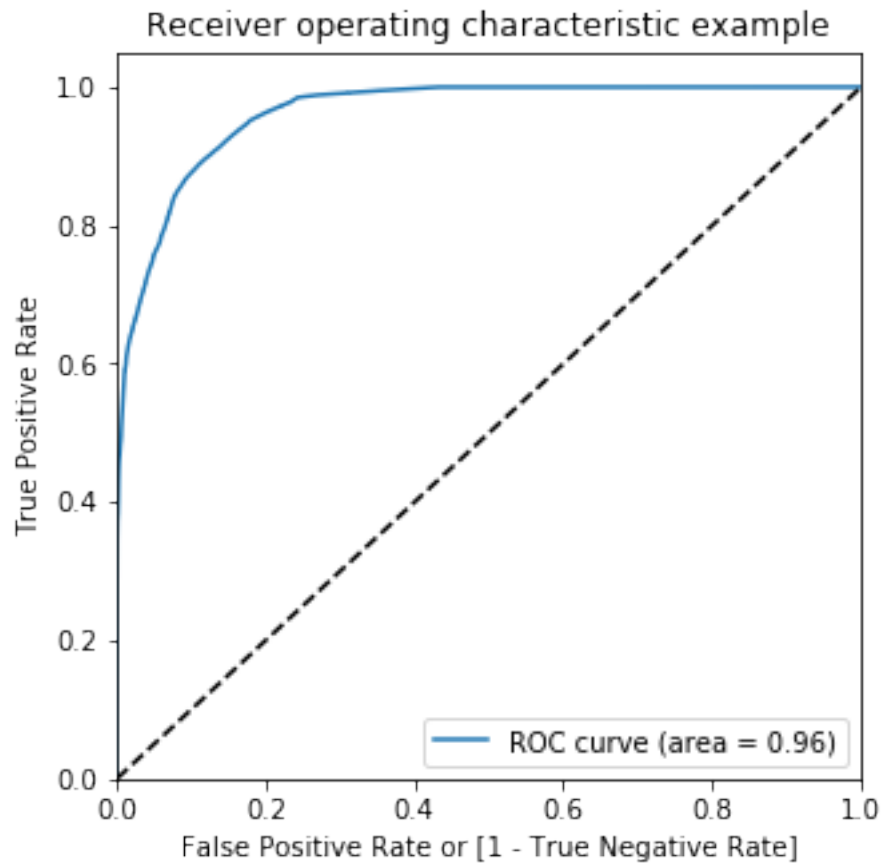
	precision	recall	f1-score	support
0	0.87	0.91	0.89	227449
1	0.91	0.86	0.88	227448
accuracy			0.89	454897
macro avg	0.89	0.89	0.89	454897
weighted avg	0.89	0.89	0.89	454897

```
[313]: # Predicted probability
y_train_pred_proba = logistic_bal_adasyn_model.predict_proba(X_train_adasyn)[:
↪,1]
```

```
[314]: # roc_auc
auc = metrics.roc_auc_score(y_train_adasyn, y_train_pred_proba)
auc
```

```
[314]: 0.9631610161614914
```

```
[315]: # Plot the ROC curve
draw_roc(y_train_adasyn, y_train_pred_proba)
```



Prediction on the test set

```
[316]: # Prediction on the test set
y_test_pred = logistic_bal_adasyn_model.predict(X_test)
```

```
[317]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[51642  5224]
 [    4    92]]
```

```
[318]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[319]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

```
Accuracy:- 0.9082195147642288
Sensitivity:- 0.9583333333333334
Specificity:- 0.9081349136566665
```

```
[320]: # classification_report
print(classification_report(y_test, y_test_pred))
```

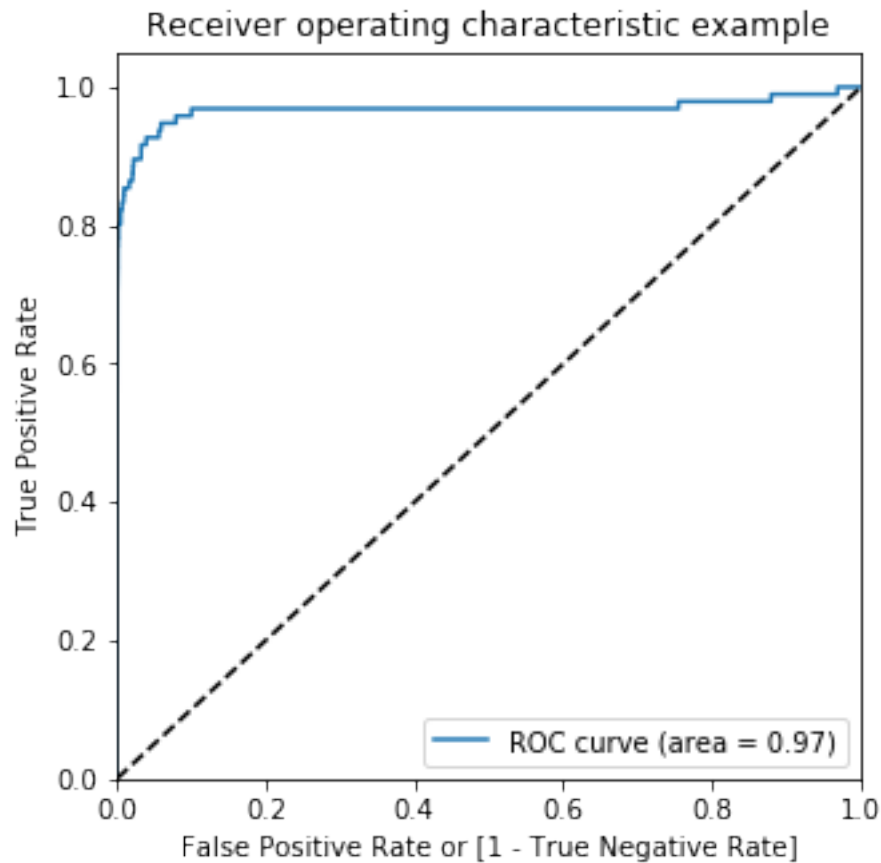
	precision	recall	f1-score	support
0	1.00	0.91	0.95	56866
1	0.02	0.96	0.03	96
accuracy			0.91	56962
macro avg	0.51	0.93	0.49	56962
weighted avg	1.00	0.91	0.95	56962

```
[321]: # Predicted probability
y_test_pred_proba = logistic_bal_adasyn_model.predict_proba(X_test)[: ,1]
```

```
[322]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

```
[322]: 0.9671573487086602
```

```
[323]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.88
 - Sensitivity = 0.86
 - Specificity = 0.91
 - ROC = 0.96
- Test set
 - Accuracy = 0.90
 - Sensitivity = 0.95
 - Specificity = 0.90
 - ROC = 0.97

4.2.2 Decision Tree

```
[205]: # Create the parameter grid
param_grid = {
    'max_depth': range(5, 15, 5),
    'min_samples_leaf': range(50, 150, 50),
    'min_samples_split': range(50, 150, 50),
```

```

}

# Instantiate the grid search model
dtree = DecisionTreeClassifier()

grid_search = GridSearchCV(estimator = dtree,
                           param_grid = param_grid,
                           scoring= 'roc_auc',
                           cv = 3,
                           verbose = 1)

# Fit the grid search to the data
grid_search.fit(X_train_adasyn,y_train_adasyn)

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
 [Parallel(n_jobs=1)]: Done 24 out of 24 | elapsed: 5.4min finished

```

[205]: GridSearchCV(cv=3, error_score=nan,
                  estimator=DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None,
                                                    criterion='gini', max_depth=None,
                                                    max_features=None,
                                                    max_leaf_nodes=None,
                                                    min_impurity_decrease=0.0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=0.0,
                                                    presort='deprecated',
                                                    random_state=None,
                                                    splitter='best'),
                  iid='deprecated', n_jobs=None,
                  param_grid={'max_depth': range(5, 15, 5),
                              'min_samples_leaf': range(50, 150, 50),
                              'min_samples_split': range(50, 150, 50)},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                  scoring='roc_auc', verbose=1)

```

```

[206]: # cv results
cv_results = pd.DataFrame(grid_search.cv_results_)
cv_results

```

```

[206]:   mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
0      9.992159    0.352046    0.098602    0.007498
1      9.723238    0.012229    0.092537    0.001110
2      9.719180    0.064421    0.087535    0.006825

```

3	9.705453	0.014290	0.092735	0.001223
4	17.367458	0.262487	0.104000	0.014708
5	17.314552	0.315418	0.094069	0.000663
6	17.106967	0.206823	0.104667	0.014260
7	17.102270	0.148033	0.099734	0.006711

	param_max_depth	param_min_samples_leaf	param_min_samples_split	\
0	5	50	50	
1	5	50	100	
2	5	100	50	
3	5	100	100	
4	10	50	50	
5	10	50	100	
6	10	100	50	
7	10	100	100	

	params	split0_test_score	\
0	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.902958	
1	{'max_depth': 5, 'min_samples_leaf': 50, 'min_...	0.902958	
2	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.902958	
3	{'max_depth': 5, 'min_samples_leaf': 100, 'min_...	0.902958	
4	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.935282	
5	{'max_depth': 10, 'min_samples_leaf': 50, 'min_...	0.935250	
6	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.937615	
7	{'max_depth': 10, 'min_samples_leaf': 100, 'mi...	0.937226	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.920356	0.909339	0.910884	0.007186	
1	0.920356	0.909330	0.910882	0.007187	
2	0.920310	0.911441	0.911570	0.007085	
3	0.920317	0.911441	0.911572	0.007087	
4	0.946369	0.934662	0.938771	0.005378	
5	0.945371	0.937635	0.939419	0.004320	
6	0.945862	0.940961	0.941479	0.003387	
7	0.945535	0.940369	0.941043	0.003425	

	rank_test_score
0	7
1	8
2	6
3	5
4	4
5	3
6	1
7	2

```
[207]: # Printing the optimal sensitivity score and hyperparameters
print("Best roc_auc:-", grid_search.best_score_)
print(grid_search.best_estimator_)
```

```
Best roc_auc:- 0.9414793563319087
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                      max_depth=10, max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=100, min_samples_split=50,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')
```

```
[324]: # Model with optimal hyperparameters
dt_bal_adasyn_model = DecisionTreeClassifier(criterion = "gini",
                                           random_state = 100,
                                           max_depth=10,
                                           min_samples_leaf=100,
                                           min_samples_split=50)

dt_bal_adasyn_model.fit(X_train_adasyn, y_train_adasyn)
```

```
[324]: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                             max_depth=10, max_features=None, max_leaf_nodes=None,
                             min_impurity_decrease=0.0, min_impurity_split=None,
                             min_samples_leaf=100, min_samples_split=50,
                             min_weight_fraction_leaf=0.0, presort='deprecated',
                             random_state=100, splitter='best')
```

Prediction on the train set

```
[325]: # Predictions on the train set
y_train_pred = dt_bal_adasyn_model.predict(X_train_adasyn)
```

```
[326]: # Confusion matrix
confusion = metrics.confusion_matrix(y_train_adasyn, y_train_pred)
print(confusion)
```

```
[[215929  11520]
 [   1118 226330]]
```

```
[327]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[328]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_train_adasyn, y_train_pred))

# Sensitivity
```



```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9722178866864367
Sensitivity:- 0.9950845907636031
Specificity:- 0.9493512831447929

```
[329]: # classification_report
print(classification_report(y_train_adasyn, y_train_pred))
```

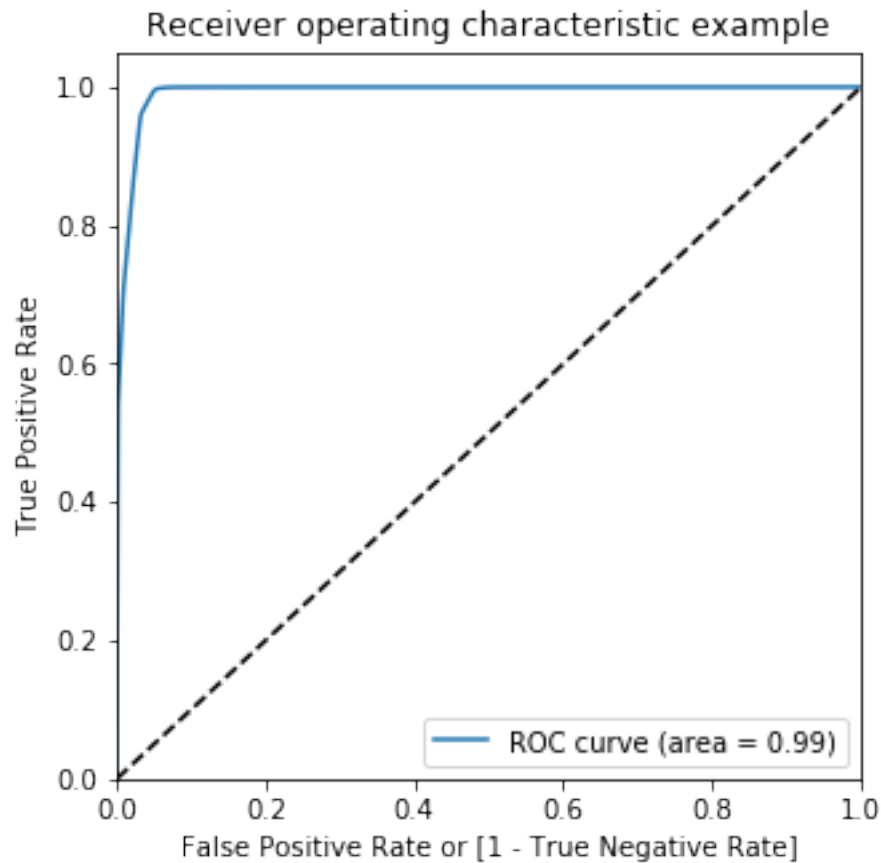
	precision	recall	f1-score	support
0	0.99	0.95	0.97	227449
1	0.95	1.00	0.97	227448
accuracy			0.97	454897
macro avg	0.97	0.97	0.97	454897
weighted avg	0.97	0.97	0.97	454897

```
[330]: # Predicted probability
y_train_pred_proba = dt_bal_adasyn_model.predict_proba(X_train_adasyn)[: ,1]
```

```
[331]: # roc_auc
auc = metrics.roc_auc_score(y_train_adasyn, y_train_pred_proba)
auc
```

[331]: 0.9917591040224101

```
[332]: # Plot the ROC curve
draw_roc(y_train_adasyn, y_train_pred_proba)
```



Prediction on the test set

```
[333]: # Predictions on the test set
y_test_pred = dt_bal_adasyn_model.predict(X_test)
```

```
[334]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[53880 2986]
 [ 15    81]]
```

```
[335]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[336]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9473157543625575
Sensitivity:- 0.84375
Specificity:- 0.9474905919178419

```
[337]: # classification_report
print(classification_report(y_test, y_test_pred))
```

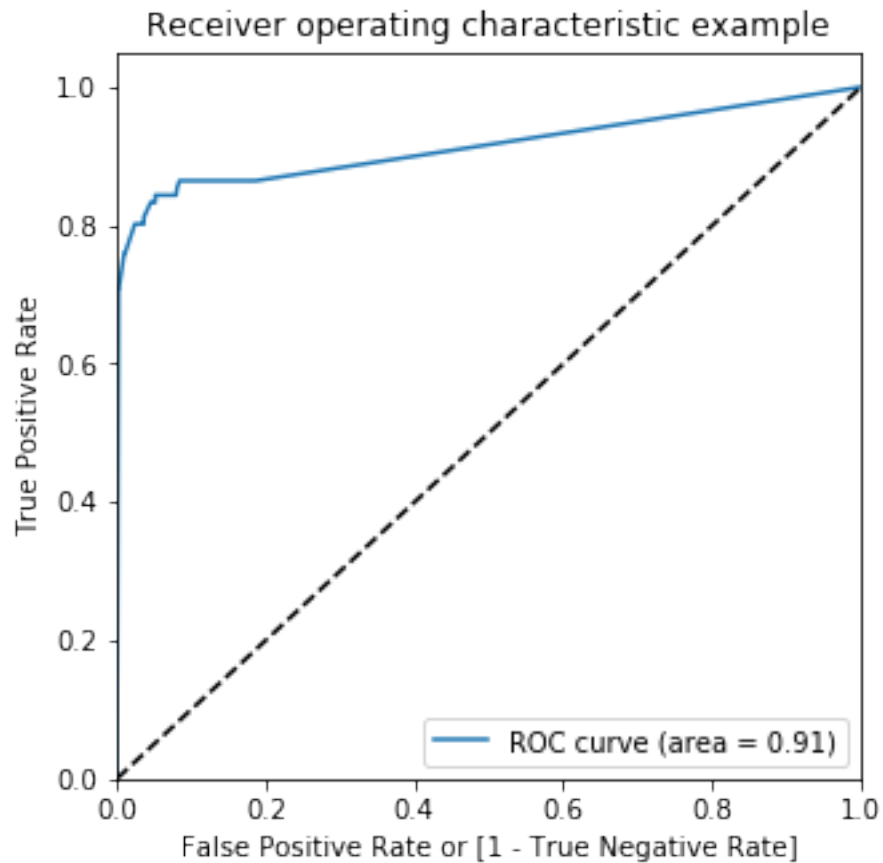
	precision	recall	f1-score	support
0	1.00	0.95	0.97	56866
1	0.03	0.84	0.05	96
accuracy			0.95	56962
macro avg	0.51	0.90	0.51	56962
weighted avg	1.00	0.95	0.97	56962

```
[338]: # Predicted probability
y_test_pred_proba = dt_bal_adasyn_model.predict_proba(X_test)[: ,1]
```

```
[339]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[339]: 0.9141440147305362

```
[340]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.97
 - Sensitivity = 0.99
 - Specificity = 0.95
 - ROC-AUC = 0.99
- Test set
 - Accuracy = 0.95
 - Sensitivity = 0.84
 - Specificity = 0.95
 - ROC-AUC = 0.91

4.2.3 XGBoost

```
[221]: # hyperparameter tuning with XGBoost
```

```
# creating a KFold object  
folds = 3
```

```

# specify range of hyperparameters
param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

# specify model
xgb_model = XGBClassifier(max_depth=2, n_estimators=200)

# set up GridSearchCV()
model_cv = GridSearchCV(estimator = xgb_model,
                        param_grid = param_grid,
                        scoring= 'roc_auc',
                        cv = folds,
                        verbose = 1,
                        return_train_score=True)

# fit the model
model_cv.fit(X_train_adasyn, y_train_adasyn)

```

Fitting 3 folds for each of 6 candidates, totalling 18 fits

```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 18 out of 18 | elapsed: 42.5min finished

```

```

[221]: GridSearchCV(cv=3, error_score=nan,
                  estimator=XGBClassifier(base_score=None, booster=None,
                                          colsample_bylevel=None,
                                          colsample_bynode=None,
                                          colsample_bytree=None, gamma=None,
                                          gpu_id=None, importance_type='gain',
                                          interaction_constraints=None,
                                          learning_rate=None, max_delta_step=None,
                                          max_depth=2, min_child_weight=None,
                                          missing=nan, monotone_constraints=None,
                                          n_estimators=200,
                                          objective='binary:logistic',
                                          random_state=None, reg_alpha=None,
                                          reg_lambda=None, scale_pos_weight=None,
                                          subsample=None, tree_method=None,
                                          validate_parameters=False,
                                          verbosity=None),
                  iid='deprecated', n_jobs=None,
                  param_grid={'learning_rate': [0.2, 0.6],
                              'subsample': [0.3, 0.6, 0.9]},
                  pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                  scoring='roc_auc', verbose=1)

```

```
[222]: # cv results
cv_results = pd.DataFrame(model_cv.cv_results_)
cv_results
```

```
[222]:
```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	\
0	107.725133	10.671068	0.794354	0.057250	
1	138.776001	0.322162	0.785001	0.008165	
2	157.356024	1.177755	0.809046	0.050527	
3	108.153853	0.945556	0.795379	0.047079	
4	139.687656	0.522447	0.793045	0.024834	
5	184.802151	45.551483	0.767030	0.030258	

	param_learning_rate	param_subsample	\
0	0.2	0.3	
1	0.2	0.6	
2	0.2	0.9	
3	0.6	0.3	
4	0.6	0.6	
5	0.6	0.9	

	params	split0_test_score	\
0	{'learning_rate': 0.2, 'subsample': 0.3}	0.975756	
1	{'learning_rate': 0.2, 'subsample': 0.6}	0.978500	
2	{'learning_rate': 0.2, 'subsample': 0.9}	0.977110	
3	{'learning_rate': 0.6, 'subsample': 0.3}	0.979173	
4	{'learning_rate': 0.6, 'subsample': 0.6}	0.971621	
5	{'learning_rate': 0.6, 'subsample': 0.9}	0.977355	

	split1_test_score	split2_test_score	mean_test_score	std_test_score	\
0	0.996202	0.994729	0.988896	0.009310	
1	0.996075	0.993204	0.989260	0.007698	
2	0.996104	0.993729	0.988981	0.008450	
3	0.998146	0.998145	0.991822	0.008944	
4	0.996825	0.997548	0.988664	0.012055	
5	0.998183	0.995571	0.990370	0.009265	

	rank_test_score	split0_train_score	split1_train_score	\
0	5	0.999304	0.999014	
1	3	0.999295	0.999072	
2	4	0.999300	0.999069	
3	1	0.999937	0.999934	
4	6	0.999950	0.999942	
5	2	0.999953	0.999935	

	split2_train_score	mean_train_score	std_train_score
0	0.999315	0.999211	0.000139
1	0.999224	0.999197	0.000093

2	0.999194	0.999188	0.000095
3	0.999942	0.999938	0.000004
4	0.999947	0.999946	0.000003
5	0.999955	0.999948	0.000009

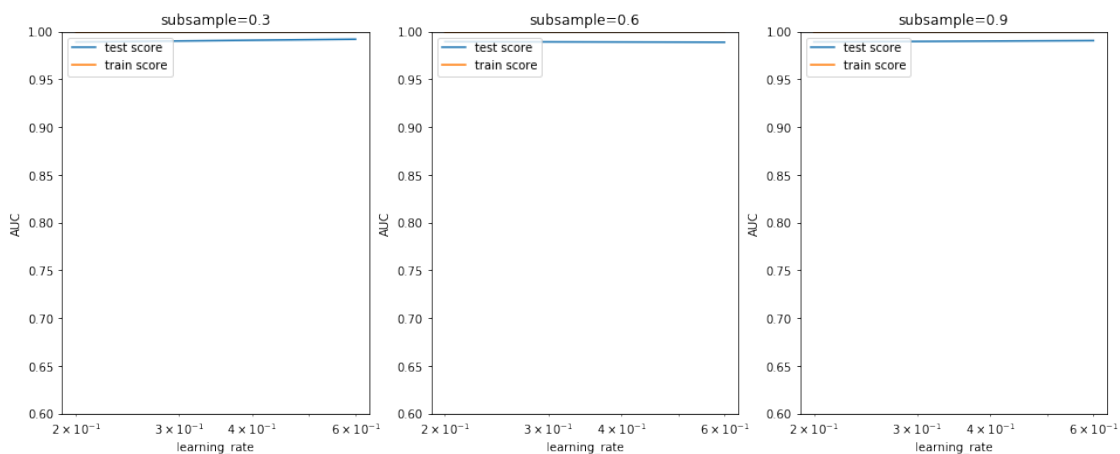
```
[223]: ## plotting
plt.figure(figsize=(16,6))

param_grid = {'learning_rate': [0.2, 0.6],
              'subsample': [0.3, 0.6, 0.9]}

for n, subsample in enumerate(param_grid['subsample']):

    # subplot 1/n
    plt.subplot(1,len(param_grid['subsample']), n+1)
    df = cv_results[cv_results['param_subsample']==subsample]

    plt.plot(df["param_learning_rate"], df["mean_test_score"])
    plt.plot(df["param_learning_rate"], df["mean_train_score"])
    plt.xlabel('learning_rate')
    plt.ylabel('AUC')
    plt.title("subsample={0}".format(subsample))
    plt.ylim([0.60, 1])
    plt.legend(['test score', 'train score'], loc='upper left')
    plt.xscale('log')
```



```
[224]: model_cv.best_params_
```

```
[224]: {'learning_rate': 0.6, 'subsample': 0.3}
```

```
[341]: # chosen hyperparameters
```

```
params = {'learning_rate': 0.6,  
          'max_depth': 2,  
          'n_estimators': 200,  
          'subsample': 0.3,  
          'objective': 'binary:logistic'}  
  
# fit model on training data  
xgb_bal_adasyn_model = XGBClassifier(params = params)  
xgb_bal_adasyn_model.fit(X_train_adasyn, y_train_adasyn)
```

```
[341]: XGBClassifier(base_score=0.5, booster=None, colsample_bylevel=1,  
                    colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,  
                    importance_type='gain', interaction_constraints=None,  
                    learning_rate=0.300000012, max_delta_step=0, max_depth=6,  
                    min_child_weight=1, missing=nan, monotone_constraints=None,  
                    n_estimators=100, n_jobs=0, num_parallel_tree=1,  
                    objective='binary:logistic',  
                    params={'learning_rate': 0.6, 'max_depth': 2, 'n_estimators': 200,  
                             'objective': 'binary:logistic', 'subsample': 0.3},  
                    random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,  
                    subsample=1, tree_method=None, validate_parameters=False,  
                    verbosity=None)
```

Prediction on the train set

```
[342]: # Predictions on the train set
```

```
y_train_pred = xgb_bal_adasyn_model.predict(X_train_adasyn)
```

```
[343]: # Confusion matrix
```

```
confusion = metrics.confusion_matrix(y_train_adasyn, y_train_pred)  
print(confusion)
```

```
[[227449    0]  
 [    0 227448]]
```

```
[344]: TP = confusion[1,1] # true positive  
TN = confusion[0,0] # true negatives  
FP = confusion[0,1] # false positives  
FN = confusion[1,0] # false negatives
```

```
[345]: # Accuracy
```

```
print("Accuracy:-", metrics.accuracy_score(y_train_adasyn, y_train_pred))
```

```
# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))
```

```
# Specificity
```



```
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9999956034003302

Sensitivity:- 1.0

Specificity:- 1.0

```
[346]: # classification_report
print(classification_report(y_train_adasyn, y_train_pred))
```

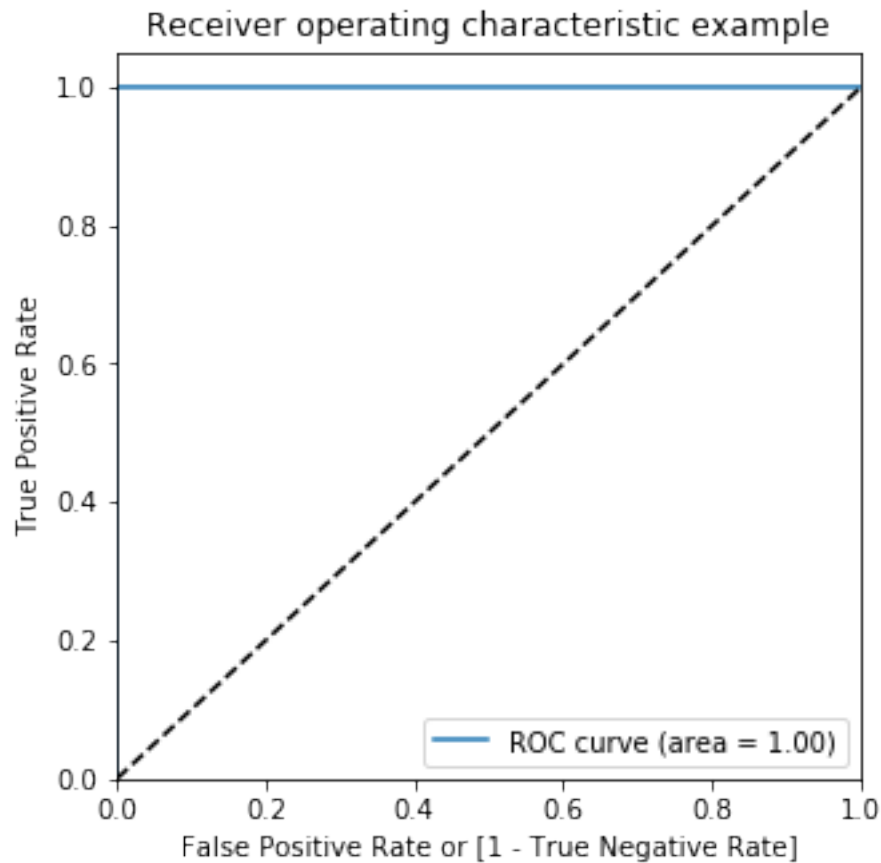
	precision	recall	f1-score	support
0	1.00	1.00	1.00	227449
1	1.00	1.00	1.00	227448
accuracy			1.00	454897
macro avg	1.00	1.00	1.00	454897
weighted avg	1.00	1.00	1.00	454897

```
[347]: # Predicted probability
y_train_pred_proba = xgb_bal_adasyn_model.predict_proba(X_train_adasyn)[: ,1]
```

```
[348]: # roc_auc
auc = metrics.roc_auc_score(y_train_adasyn, y_train_pred_proba)
auc
```

```
[348]: 1.0
```

```
[349]: # Plot the ROC curve
draw_roc(y_train_adasyn, y_train_pred_proba)
```



Prediction on the test set

```
[350]: # Predictions on the test set
y_test_pred = xgb_bal_adasyn_model.predict(X_test)
```

```
[351]: # Confusion matrix
confusion = metrics.confusion_matrix(y_test, y_test_pred)
print(confusion)
```

```
[[56825   41]
 [   21   75]]
```

```
[352]: TP = confusion[1,1] # true positive
TN = confusion[0,0] # true negatives
FP = confusion[0,1] # false positives
FN = confusion[1,0] # false negatives
```

```
[353]: # Accuracy
print("Accuracy:-", metrics.accuracy_score(y_test, y_test_pred))

# Sensitivity
```

```
print("Sensitivity:-", TP / float(TP+FN))

# Specificity
print("Specificity:-", TN / float(TN+FP))
```

Accuracy:- 0.9989115550718023
Sensitivity:- 0.78125
Specificity:- 0.9992790067878873

```
[354]: # classification_report
print(classification_report(y_test, y_test_pred))
```

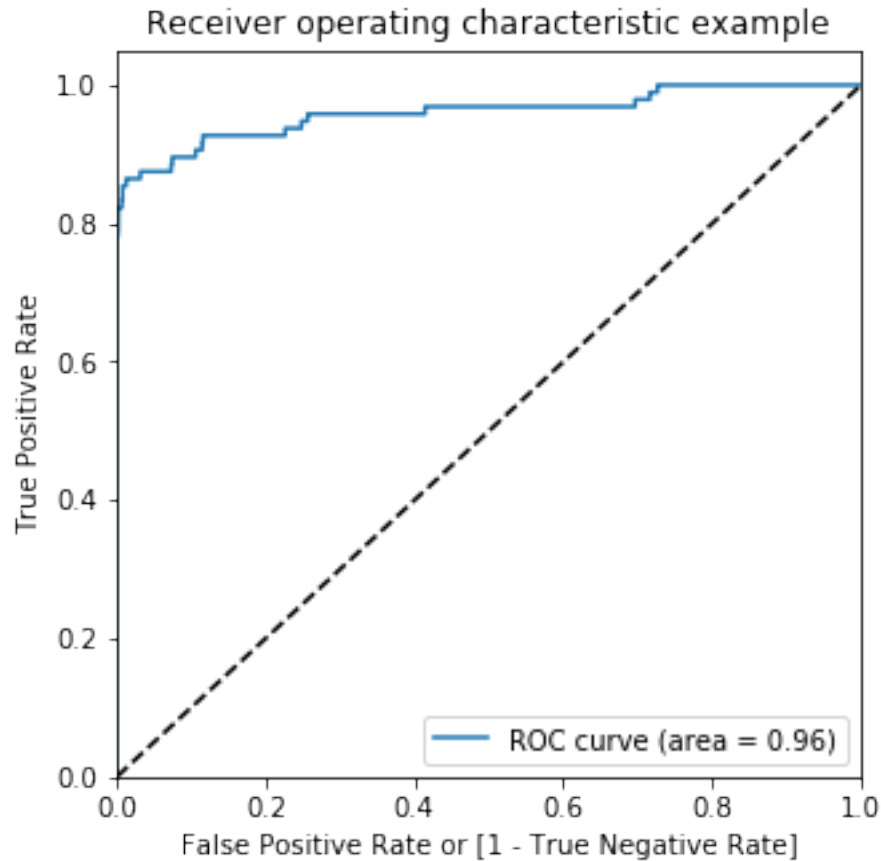
	precision	recall	f1-score	support
0	1.00	1.00	1.00	56866
1	0.65	0.78	0.71	96
accuracy			1.00	56962
macro avg	0.82	0.89	0.85	56962
weighted avg	1.00	1.00	1.00	56962

```
[355]: # Predicted probability
y_test_pred_proba = xgb_bal_adasyn_model.predict_proba(X_test)[: ,1]
```

```
[356]: # roc_auc
auc = metrics.roc_auc_score(y_test, y_test_pred_proba)
auc
```

[356]: 0.9599176499724499

```
[357]: # Plot the ROC curve
draw_roc(y_test, y_test_pred_proba)
```



Model summary

- Train set
 - Accuracy = 0.99
 - Sensitivity = 1.0
 - Specificity = 1.0
 - ROC-AUC = 1.0
- Test set
 - Accuracy = 0.99
 - Sensitivity = 0.78
 - Specificity = 0.99
 - ROC-AUC = 0.96

4.2.4 Choosing best model on the balanced data

He we balanced the data with various approach such as Undersampling, Oversampling, SMOTE and Adasy. With every data balancing thechnique we built several models such as Logistic, XGBoost, Decision Tree, and Random Forest.

We can see that almost all the models performed more or less good. But we should be interested in the best model.

Though the Undersampling technique models performed well, we should keep mind that by doing the undersampling some information were lost. Hence, it is better not to consider the undersampling models.

Whereas the SMOTE and Adasyn models performed well. Among those models the simplest model Logistic regression has ROC score 0.99 in the train set and 0.97 on the test set. We can consider the Logistic model as the best model to choose because of the easy interpretation of the models and also the resource requirements to build the model is lesser than the other heavy models such as Random forest or XGBoost.

Hence, we can conclude that the Logistic regression model with SMOTE is the best model for its simplicity and less resource requirement.

Print the FPR,TPR & select the best threshold from the roc curve for the best model

```
[92]: print('Train auc =', metrics.roc_auc_score(y_train_smote,
        ↪ y_train_pred_proba_log_bal_smote))
fpr, tpr, thresholds = metrics.roc_curve(y_train_smote,
        ↪ y_train_pred_proba_log_bal_smote)
threshold = thresholds[np.argmax(tpr-fpr)]
print("Threshold=", threshold)
```

```
Train auc = 0.9897539730968845
```

```
Threshold= 0.5311563613510013
```

We can see that the threshold is 0.53, for which the TPR is the highest and FPR is the lowest and we got the best ROC score.

4.3 Cost benefit analysis

We have tried several models till now with both balanced and imbalanced data. We have noticed most of the models have performed more or less well in terms of ROC score, Precision and Recall.

But while picking the best model we should consider few things such as whether we have required infrastructure, resources or computational power to run the model or not. For the models such as Random forest, SVM, XGBoost we require heavy computational resources and eventually to build that infrastructure the cost of deploying the model increases. On the other hand the simpler model such as Logistic regression requires less computational resources, so the cost of building the model is less.

We also have to consider that for little change of the ROC score how much monetary loss of gain the bank incur. If the amount is huge then we have to consider building the complex model even though the cost of building the model is high.

4.4 Summary to the business

For banks with smaller average transaction value, we would want high precision because we only want to label relevant transactions as fraudulent. For every transaction that is flagged as fraudulent, we can add the human element to verify whether the transaction was done by calling the customer. However, when precision is low, such tasks are a burden because the human element has to be increased.

For banks having a larger transaction value, if the recall is low, i.e., it is unable to detect transactions that are labelled as non-fraudulent. So we have to consider the losses if the missed transaction was a high-value fraudulent one.

So here, to save the banks from high-value fraudulent transactions, we have to focus on a high recall in order to detect actual fraudulent transactions.

After performing several models, we have seen that in the balanced dataset with SMOTE technique the simplest Logistic regression model has good ROC score and also high Recall. Hence, we can go with the logistic model here. It is also easier to interpret and explain to the business.