```python
1  import os
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from xgboost import XGBClassifier
6  from sklearn.model_selection import GridSearchCV
7  from sklearn.model_selection import train_test_split
8  from sklearn.preprocessing import StandardScaler
9  from sklearn.preprocessing import PowerTransformer
10 from imblearn.over_sampling import SMOTE
11 from sklearn.metrics import accuracy_score,
   confusion_matrix,classification_report,roc_auc_score,
   roc_curve
12
13 df = pd.read_csv('../creditcard data/creditcard.csv')
14 classes = df['Class'].value_counts()
15 print(classes)
16 normal_share = round((classes[0]/df['Class'].count()*
   100),2)
17 print("Normal share =", normal_share)
18 fraud_share = round((classes[1]/df['Class'].count()*
   100),2)
19 print("Fraud share=", fraud_share)
20 sns.countplot(x='Class', data=df)
21 plt.title('Number of fraudulent vs non-fraudulent
   transactions')
22 #plt.show()
23
24 # Creating fraudulent dataframe
25 data_fraud = df[df['Class'] == 1]
26 # Creating non fraudulent dataframe
27 data_non_fraud = df[df['Class'] == 0]
28 # Distribution plot
29 plt.figure(figsize=(8,5))
30 ax = sns.distplot(data_fraud['Amount'],label='
   fraudulent',hist=False)
31 ax = sns.distplot(data_non_fraud['Amount'],label='non
    fraudulent',hist=False)
32 ax.set(xlabel='Transaction Amount')
33 #plt.show()
34
```

```python
35 # Putting feature variables into X
36 X = df.drop(['Class'], axis=1)
37 # Putting target variable to y
38 y = df['Class']
39 # Splitting data into train and test set 80:20
40 X_train, X_test, y_train, y_test = train_test_split(X
   , y, train_size=0.8, test_size=0.2, random_state=100)
41
42 # Standardization method
43 # Instantiate the Scaler
44 scaler = StandardScaler()
45 # Fit the data into scaler and transform
46 X_train['Amount'] = scaler.fit_transform(X_train[['
   Amount']])
47 print(X_train.head())
48
49 # Transform the test set
50 X_test['Amount'] = scaler.transform(X_test[['Amount'
   ]])
51 print(X_test.head())
52
53 # Plotting the distribution of the variables (
   skewness) of all the columns
54 cols = X_train.columns
55 k = 0
56 plt.figure(figsize=(17, 28))
57 for col in cols:
58  k = k + 1
59  plt.subplot(6, 5, k)
60  sns.distplot(X_train[col])
61  plt.title(col + ' Skewness: ' + str(X_train[col].
   skew()))
62 plt.tight_layout()
63 #plt.show()
64
65 # Instantiate the powertransformer
66 pt = PowerTransformer(method='yeo-johnson',
   standardize=True,
67 copy=False)
68 # Fit and transform the PT on training data
69 X_train[cols] = pt.fit_transform(X_train)
```

```python
70 # Transform the test set
71 X_test[cols] = pt.transform(X_test)
72 # Plotting the distribution of the variables (
   skewness) of all the columns
73 k=0
74 plt.figure(figsize=(17,28))
75 for col in cols :
76  k=k+1
77  plt.subplot(6, 5,k)
78  sns.distplot(X_train[col])
79  plt.title(col+' '+str(X_train[col].skew()))
80 #plt.show()
81
82 # Instantiate SMOTE
83 sm = SMOTE(random_state=27)
84 # Fitting SMOTE to the train set
85 X_train_smote, y_train_smote = sm.fit_resample(
   X_train, y_train)
86 print('Before SMOTE oversampling X_train shape=',
   X_train.shape)
87 print('After SMOTE oversampling X_train shape=',
   X_train_smote.shape)
88
89 # hyperparameter tuning with XGBoost
90 # creating a KFold object
91 folds = 3
92 # specify range of hyperparameters
93 param_grid = {'learning_rate': [0.2, 0.6], '
   subsample': [0.3, 0.6, 0.9]}
94 # specify model
95 xgb_model = XGBClassifier(max_depth=2, n_estimators=
   200)
96 # set up GridSearchCV()
97 model_cv = GridSearchCV(estimator = xgb_model,
   param_grid = param_grid, scoring= 'roc_auc', cv =
   folds, verbose = 1, return_train_score=True)
98 # fit the model
99 model_cv.fit(X_train_smote, y_train_smote)
100
101 # cv results
102 cv_results = pd.DataFrame(model_cv.cv_results_)
```

```python
103 print(cv_results)
104
105 # # plotting
106 plt.figure(figsize=(16,6))
107 param_grid = {'learning_rate': [0.2, 0.6], '
    subsample': [0.3, 0.6, 0.9]}
108 for n, subsample in enumerate(param_grid['subsample'
    ]):
109
110  # subplot 1/n
111  plt.subplot(1,len(param_grid['subsample']), n+1)
112  df = cv_results[cv_results['param_subsample']==
    subsample]
113  plt.plot(df["param_learning_rate"], df["
    mean_test_score"])
114  plt.plot(df["param_learning_rate"], df["
    mean_train_score"])
115  plt.xlabel('learning_rate')
116  plt.ylabel('AUC')
117  plt.title("subsample={0}".format(subsample))
118  plt.ylim([0.60, 1])
119  plt.legend(['test score', 'train score'], loc='
    upper left')
120  plt.xscale('log')
121
122 var = model_cv.best_params_
123 print(var)
124
125 # chosen hyperparameters
126 # 'objective':'binary:logistic' outputs probability
    rather than label, which we need for calculating auc
127 params = {'learning_rate': 0.6, 'max_depth': 2, '
    n_estimators': 200, 'subsample': 0.9, 'objective': '
    binary:logistic'}
128 # fit model on training data
129 xgb_bal_smote_model = XGBClassifier(params=params)
130 xgb_bal_smote_model.fit(X_train_smote, y_train_smote
    )
131 path = '../models'
132 model_filename = 'xgb_bal_smote_model.bin'
133 model_filepath = os.path.join(path, model_filename)
```

```python
134 xgb_bal_smote_model.save_model(model_filepath)
135
136 # Predictions on the train set
137 y_train_pred = xgb_bal_smote_model.predict(
    X_train_smote)
138 # Confusion matrix
139 confusion = confusion_matrix(y_train_smote,
    y_train_pred)
140 print(confusion)
141
142 TP = confusion[1,1] # true positive
143 TN = confusion[0,0] # true negatives
144 FP = confusion[0,1] # false positives
145 FN = confusion[1,0] # false negatives
146 # Accuracy
147 print("Accuracy:-", accuracy_score(y_train_smote,
    y_train_pred))
148 # Sensitivity
149 print("Sensitivity:-",TP / float(TP+FN))
150 # Specificity
151 print("Specificity:-", TN / float(TN+FP))
152
153 print(classification_report(y_train_smote,
    y_train_pred))
154
155 # Predicted probability
156 y_train_pred_proba = xgb_bal_smote_model.
    predict_proba(X_train_smote)[:,1]
157 # roc_auc
158 auc = roc_auc_score(y_train_smote,
    y_train_pred_proba)
159 print(auc)
160
161 # ROC Curve function
162 def draw_roc( actual, probs ):
163  fpr, tpr, thresholds = roc_curve( actual, probs,
    drop_intermediate = False )
164  auc_score = roc_auc_score( actual, probs )
165  plt.figure(figsize=(5, 5))
166  plt.plot( fpr, tpr, label='ROC curve (area = %0.2f
    )' % auc_score )
```

```python
167  plt.plot([0, 1], [0, 1], 'k--')
168  plt.xlim([0.0, 1.0])
169  plt.ylim([0.0, 1.05])
170  plt.xlabel('False Positive Rate or [1 - True
     Negative Rate]')
171  plt.ylabel('True Positive Rate')
172  plt.title('Receiver operating characteristic
     example')
173  plt.legend(loc="lower right")
174  plt.show()
175  return None
176
177  # Plot the ROC curve
178  draw_roc(y_train_smote, y_train_pred_proba)
179
180  # Predictions on the test set
181  y_test_pred = xgb_bal_smote_model.predict(X_test)
182  # Confusion matrix
183  confusion = confusion_matrix(y_test, y_test_pred)
184  print(confusion)
185
186  TP = confusion[1,1] # true positive
187  TN = confusion[0,0] # true negatives
188  FP = confusion[0,1] # false positives
189  FN = confusion[1,0] # false negatives
190  # Accuracy
191  print("Accuracy:-",accuracy_score(y_test,
     y_test_pred))
192  # Sensitivity
193  print("Sensitivity:-",TP / float(TP+FN))
194  # Specificity
195  print("Specificity:-", TN / float(TN+FP))
196
197  # classification_report
198  print(classification_report(y_test, y_test_pred))
199
200  # Predicted probability
201  y_test_pred_proba = xgb_bal_smote_model.
     predict_proba(X_test)[:,1]
202  # roc_auc
203  auc = roc_auc_score(y_test, y_test_pred_proba)
```

```
204 print(auc)
205
206 # Plot the ROC curve
207 draw_roc(y_test, y_test_pred_proba)
```