

# Credit card default prediction

## Table of Contents

- 1.[Introduction](#)
- 2.[Dataset Overview](#)
  - 2.1[Dataset Information](#)
  - 2.2[Variable Info](#)
- 3.[Import libraries](#)
- 4.[Import Dataset](#)
- 5.[Visualize the data](#)
  - 5.1[finding outlier](#)
  - 5.2[univariate analysis](#)
  - 5.3[univariate with numerical data](#)
- 6.[Feature vector and Target variable](#)
  - 6.2[train and test the model](#)
- 7.[Feature scaling](#)
- 8.[Model Development](#)
- 9.[Dealing with imbalanced data](#)
- 10.[Model development with imbalance data](#)
  - 10.1[XGBOOST](#)
  - 10.2[RandomForest](#)
  - 10.3[LGBM](#)
- 11.[checking overfitting and underfitting](#)
- 12.[K-Fold cross validation](#)
- 13.[ROC-AUC](#)
- 14.[Conclusion](#)

## 1.Introduction

The primary objective of this project is to leverage machine learning to predict whether a credit card user is likely to default on their payments. By evaluating past financial behaviors and patterns, we aim to provide credit decisions that are both responsible and sustainable

## 2.Dataset Overview

### 2.1 Dataset Information

This dataset contains information on default payments, demographic factors, credit data, history of payment, and bill statements of credit card clients in Taiwan from April 2005 to September 2005.

### 2.2 Variable Info

1	There are 25 features:
2	
3	ID: ID of each client
4	LIMIT_BAL: Amount of given credit in NT dollars (includes individual and family/supplementary credit
5	SEX: Gender (1=male, 2=female)
6	EDUCATION: (1=graduate school, 2=university, 3=high school, 4=others, 5=unknown, 6=unknown)
7	MARRIAGE: Marital status (1=married, 2=single, 3=others)
8	AGE: Age in years
9	PAY_0: Repayment status in September, 2005 (-1=pay duly, 1=payment delay for one month, 2=payment delay for two months, ... 8=payment delay for eight months, 9=payment delay for nine months and above)
10	PAY_2: Repayment status in August, 2005 (scale same as above)
11	PAY_3: Repayment status in July, 2005 (scale same as above)
12	PAY_4: Repayment status in June, 2005 (scale same as above)
13	PAY_5: Repayment status in May, 2005 (scale same as above)
14	PAY_6: Repayment status in April, 2005 (scale same as above)
15	BILL_AMT1: Amount of bill statement in September, 2005 (NT dollar)
16	BILL_AMT2: Amount of bill statement in August, 2005 (NT dollar)
17	BILL_AMT3: Amount of bill statement in July, 2005 (NT dollar)
18	BILL_AMT4: Amount of bill statement in June, 2005 (NT dollar)
19	BILL_AMT5: Amount of bill statement in May, 2005 (NT dollar)
20	BILL_AMT6: Amount of bill statement in April, 2005 (NT dollar)
21	PAY_AMT1: Amount of previous payment in September, 2005 (NT dollar)
22	PAY_AMT2: Amount of previous payment in August, 2005 (NT dollar)
23	PAY_AMT3: Amount of previous payment in July, 2005 (NT dollar)
24	PAY_AMT4: Amount of previous payment in June, 2005 (NT dollar)
25	PAY_AMT5: Amount of previous payment in May, 2005 (NT dollar)
26	PAY_AMT6: Amount of previous payment in April, 2005 (NT dollar)
27	default.payment.next.month: Default payment (1=yes, 0=no)
28	

## 3.Import libraries

```
In [1]: 1 # This Python 3 environment comes with many helpful analytics libraries installed
2 # It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
3 # For example, here's several helpful packages to load in
4
5 import numpy as np # linear algebra
6 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
7 import matplotlib.pyplot as plt # for data visualization purposes
8 import seaborn as sns # for data visualization
9 %matplotlib inline
```

```
In [2]: 1 import warnings
2
3 warnings.filterwarnings('ignore')
```

```
In [3]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
In [4]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.svm import SVC
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.ensemble import RandomForestClassifier
6 from xgboost import XGBClassifier
7 from lightgbm import LGBMClassifier
```

```
In [5]: 1 # 4. import Dataset
```

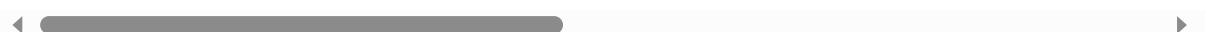
```
In [6]: 1 df = pd.read_csv(r"D:\Capstone real time project\Internship\default-of-credit-card-payment.csv")
```

```
In [7]: 1 df.head()
```

Out[7]:

	ID	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	...	BILL_AMT
0	1	20000	2	2	1	24	2	2	-1	-1	...	1
1	2	120000	2	2	2	26	-1	2	0	0	...	1
2	3	90000	2	2	2	34	0	0	0	0	...	1
3	4	50000	2	2	1	37	0	0	0	0	...	1
4	5	50000	1	2	1	57	-1	0	-1	0	...	1

5 rows × 25 columns



```
In [8]: 1 df.drop(['ID'],axis=1,inplace=True) # removing unnecessary attribute
2
```

```
In [9]: 1 # Rename the column name
2 df = df.rename(columns={'default payment next month':'default'})
```

```
In [10]: 1 df.tail()
```

```
Out[10]:
```

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	..
29995	220000	1	3	1	39	0	0	0	0	0	..
29996	150000	1	3	2	43	-1	-1	-1	-1	0	..
29997	30000	1	2	2	37	4	3	2	-1	0	..
29998	80000	1	3	1	41	1	-1	0	0	0	..
29999	50000	1	2	1	46	0	0	0	0	0	..

5 rows × 24 columns



```
In [11]: 1 df.shape
```

```
Out[11]: (30000, 24)
```

```
In [12]: 1 df.isnull().sum() # there is no missing values
```

```
Out[12]: LIMIT_BAL      0
SEX                  0
EDUCATION           0
MARRIAGE            0
AGE                 0
PAY_0               0
PAY_2               0
PAY_3               0
PAY_4               0
PAY_5               0
PAY_6               0
BILL_AMT1           0
BILL_AMT2           0
BILL_AMT3           0
BILL_AMT4           0
BILL_AMT5           0
BILL_AMT6           0
PAY_AMT1            0
PAY_AMT2            0
PAY_AMT3            0
PAY_AMT4            0
PAY_AMT5            0
PAY_AMT6            0
default             0
dtype: int64
```

```
In [13]: 1 df.info() # summary of dataset
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 24 columns):
#   Column      Non-Null Count  Dtype
---  -
0   LIMIT_BAL    30000 non-null  int64
1   SEX          30000 non-null  int64
2   EDUCATION    30000 non-null  int64
3   MARRIAGE     30000 non-null  int64
4   AGE          30000 non-null  int64
5   PAY_0        30000 non-null  int64
6   PAY_2        30000 non-null  int64
7   PAY_3        30000 non-null  int64
8   PAY_4        30000 non-null  int64
9   PAY_5        30000 non-null  int64
10  PAY_6        30000 non-null  int64
11  BILL_AMT1    30000 non-null  int64
12  BILL_AMT2    30000 non-null  int64
13  BILL_AMT3    30000 non-null  int64
14  BILL_AMT4    30000 non-null  int64
15  BILL_AMT5    30000 non-null  int64
16  BILL_AMT6    30000 non-null  int64
17  PAY_AMT1     30000 non-null  int64
18  PAY_AMT2     30000 non-null  int64
19  PAY_AMT3     30000 non-null  int64
20  PAY_AMT4     30000 non-null  int64
21  PAY_AMT5     30000 non-null  int64
22  PAY_AMT6     30000 non-null  int64
23  default      30000 non-null  int64
dtypes: int64(24)
memory usage: 5.5 MB
```

```
In [14]: 1 # Above shows all variables are numerical but some variable categorical we
```

```
In [15]: 1 df['SEX'].value_counts()
```

```
Out[15]: 2 18112
1 11888
Name: SEX, dtype: int64
```

```
In [16]: 1 df['EDUCATION'].value_counts()
```

```
Out[16]: 2 14030
1 10585
3 4917
5 280
4 123
6 51
0 14
Name: EDUCATION, dtype: int64
```

```
In [17]: 1 # Put 0 , 5, 6 into category 4(others)
2 df['EDUCATION'].replace({0:4,5:4,6:4}, inplace=True)
3 # 1=graduation school ,2= university,3=high school,4=others
```

```
In [18]: 1 df['EDUCATION'].value_counts()
```

```
Out[18]: 2    14030
          1    10585
          3     4917
          4     468
          Name: EDUCATION, dtype: int64
```

```
In [19]: 1 df['MARRIAGE'].value_counts()
```

```
Out[19]: 2    15964
          1    13659
          3     323
          0      54
          Name: MARRIAGE, dtype: int64
```

```
In [20]: 1 df['MARRIAGE'].replace({0 : 3},inplace = True)
          2 df['MARRIAGE'].value_counts()
          3
```

```
Out[20]: 2    15964
          1    13659
          3     377
          Name: MARRIAGE, dtype: int64
```

```
In [21]: 1 df_columns = df[['PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6']]
          2 df_columns.nunique()
          3
```

```
Out[21]: PAY_0    11
          PAY_2    11
          PAY_3    11
          PAY_4    11
          PAY_5    10
          PAY_6    10
          dtype: int64
```

```
In [22]: 1 #Out of 23 attribute we have 9 categorical value
```

```
In [23]: 1 df1 = df.copy()
```

In [24]: 1 df.describe()

Out[24]:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	
count	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000
mean	167484.322667	1.603733	1.842267	1.557267	35.485500	-0.016700	
std	129747.661567	0.489129	0.744494	0.521405	9.217904	1.123802	
min	10000.000000	1.000000	1.000000	1.000000	21.000000	-2.000000	
25%	50000.000000	1.000000	1.000000	1.000000	28.000000	-1.000000	
50%	140000.000000	2.000000	2.000000	2.000000	34.000000	0.000000	
75%	240000.000000	2.000000	2.000000	2.000000	41.000000	0.000000	
max	1000000.000000	2.000000	4.000000	3.000000	79.000000	8.000000	

8 rows × 24 columns



```
In [25]: 1 categorical = []
2 numerical = []
3 for col in df.columns:
4     if df[col].nunique() > 12:
5         numerical.append(col)
6
7     else:
8         categorical.append(col)
9 print('categorical_columns:', categorical)
10 print('numerical_columns:', numerical)
```

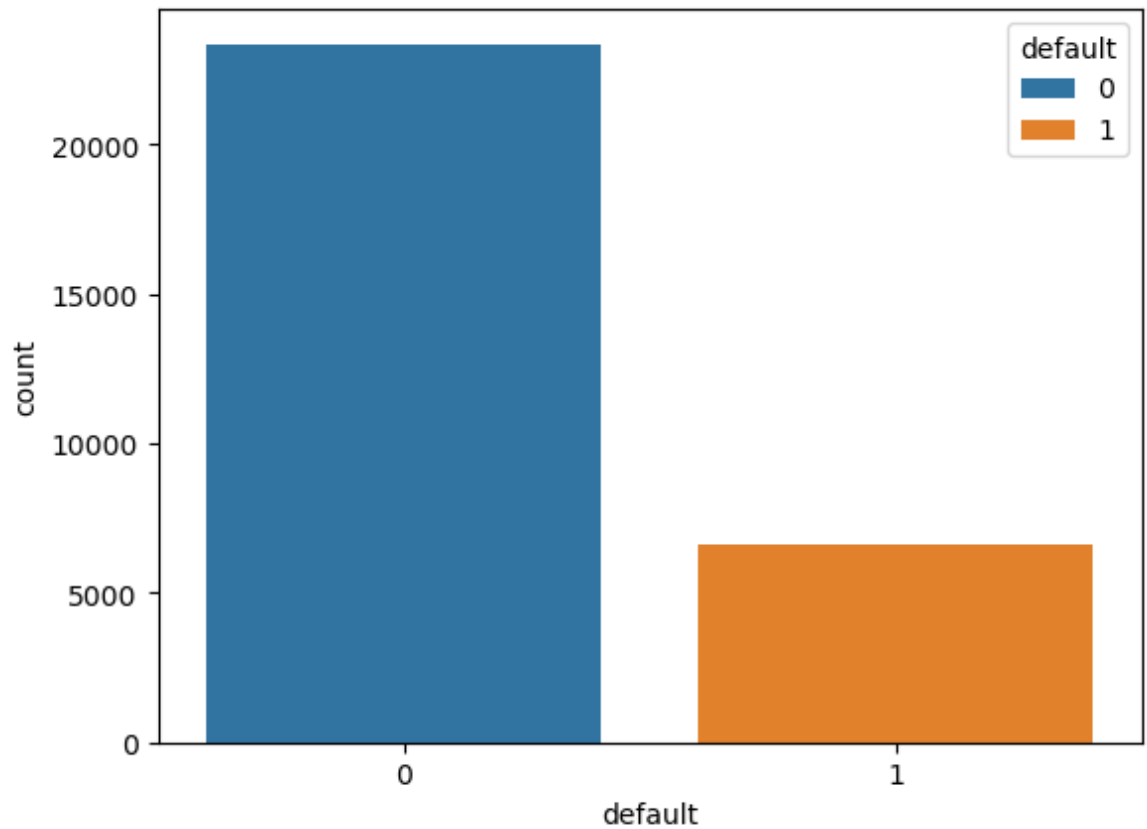
```
categorical_columns: ['SEX', 'EDUCATION', 'MARRIAGE', 'PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6', 'default']
numerical_columns: ['LIMIT_BAL', 'AGE', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5', 'BILL_AMT6', 'PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6']
```

In [26]: 1 categorical.pop()

Out[26]: 'default'

## 5. Visualize the data

```
In [27]: 1 vis = sns.countplot(x='default',data= df,hue='default')
```



```
In [28]: 1 # its seems like imbalance data
```

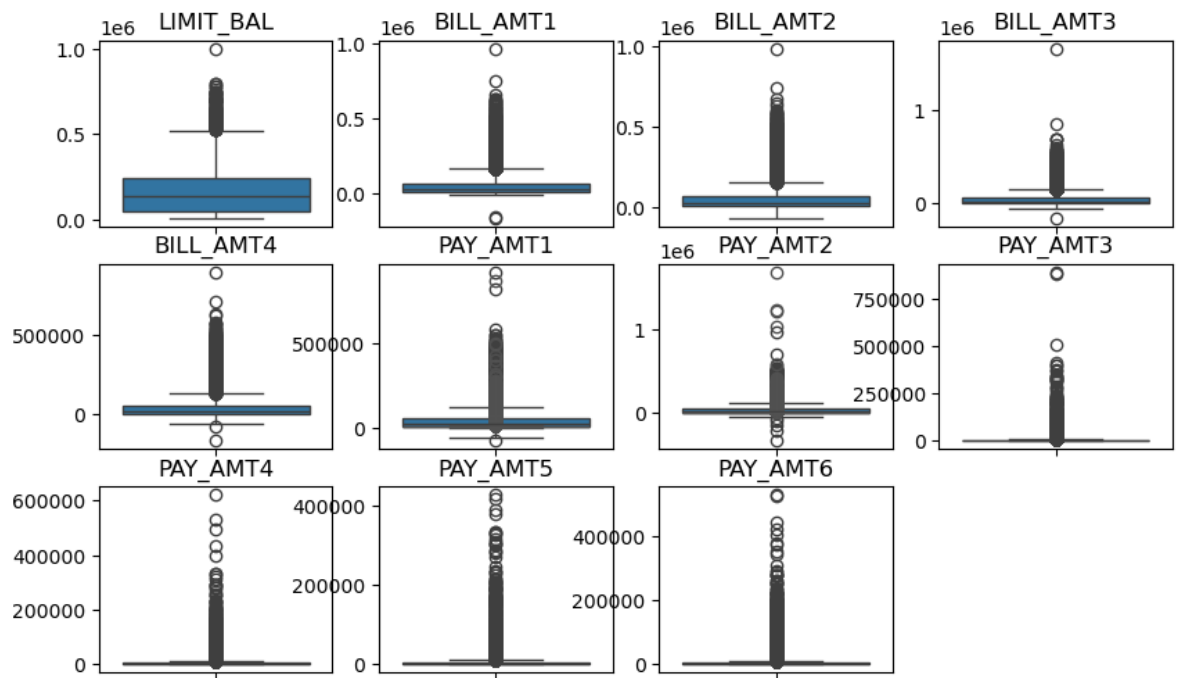
## 5.1 finding outliers



```

In [29]: 1 plt.figure(figsize=(10,8))
2 outlier_col_1 = ['LIMIT_BAL', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_',
3               'BILL_AMT5', 'BILL_AMT6']
4 outlier_col_2 = ['PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5']
5 for i in enumerate(outlier_col_1):
6     plt.subplot(4,4,i[0]+1)
7     sns.boxplot(y=df[i[1]])
8     plt.title(i[1])
9     plt.ylabel("")
10
11 for i in enumerate(outlier_col_2):
12     plt.subplot(4,4,i[0]+6)
13     sns.boxplot(y=df[i[1]])
14     plt.title(i[1])
15     plt.ylabel("")

```

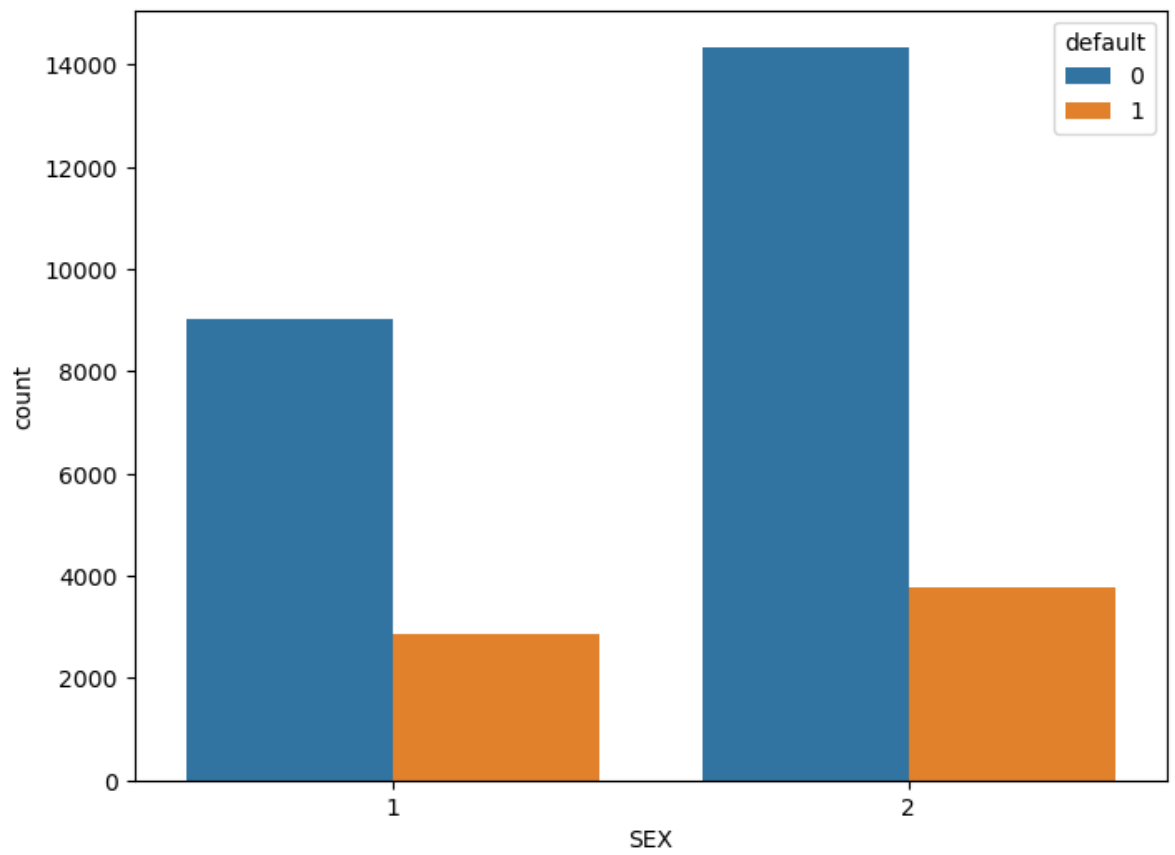


its seems dataset have Many outliers are present in the dataset, and these outliers may hold valuable information for our model.

## 5.2 univariate analysis

In [30]:

```
1 f, ax = plt.subplots(figsize=(8, 6))
2 ax = sns.countplot(x="SEX", data=df, hue='default')
3 plt.show()
```

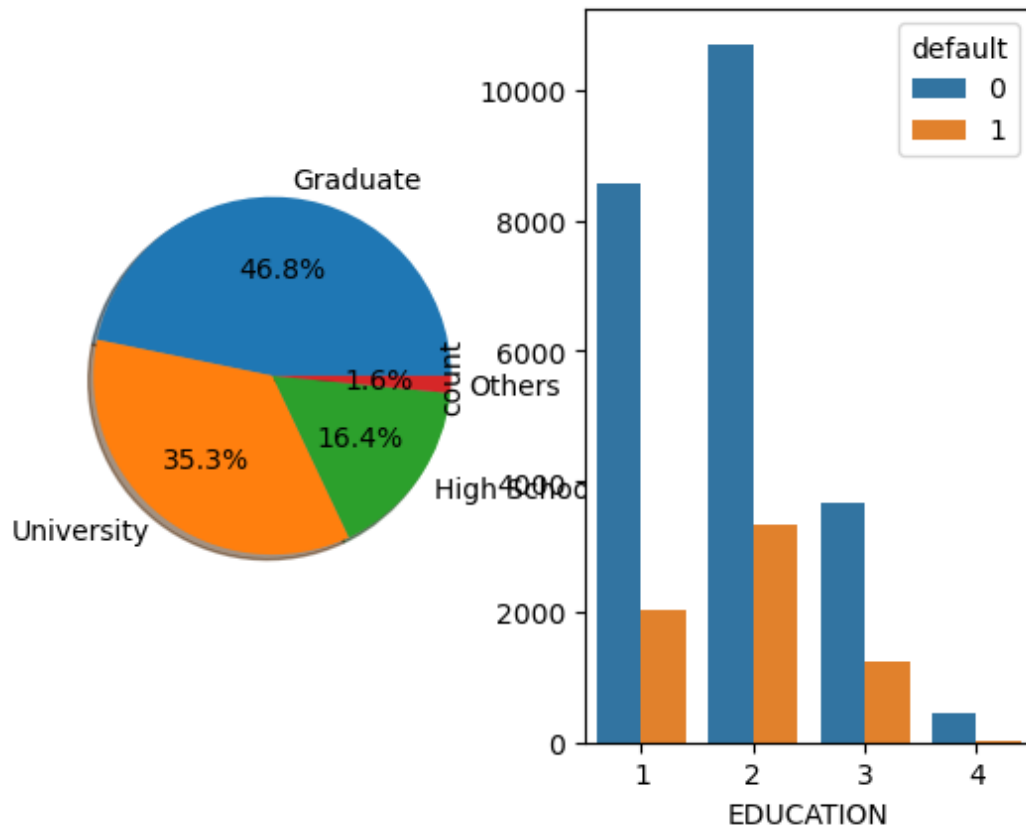


Compare to male female default is high

```

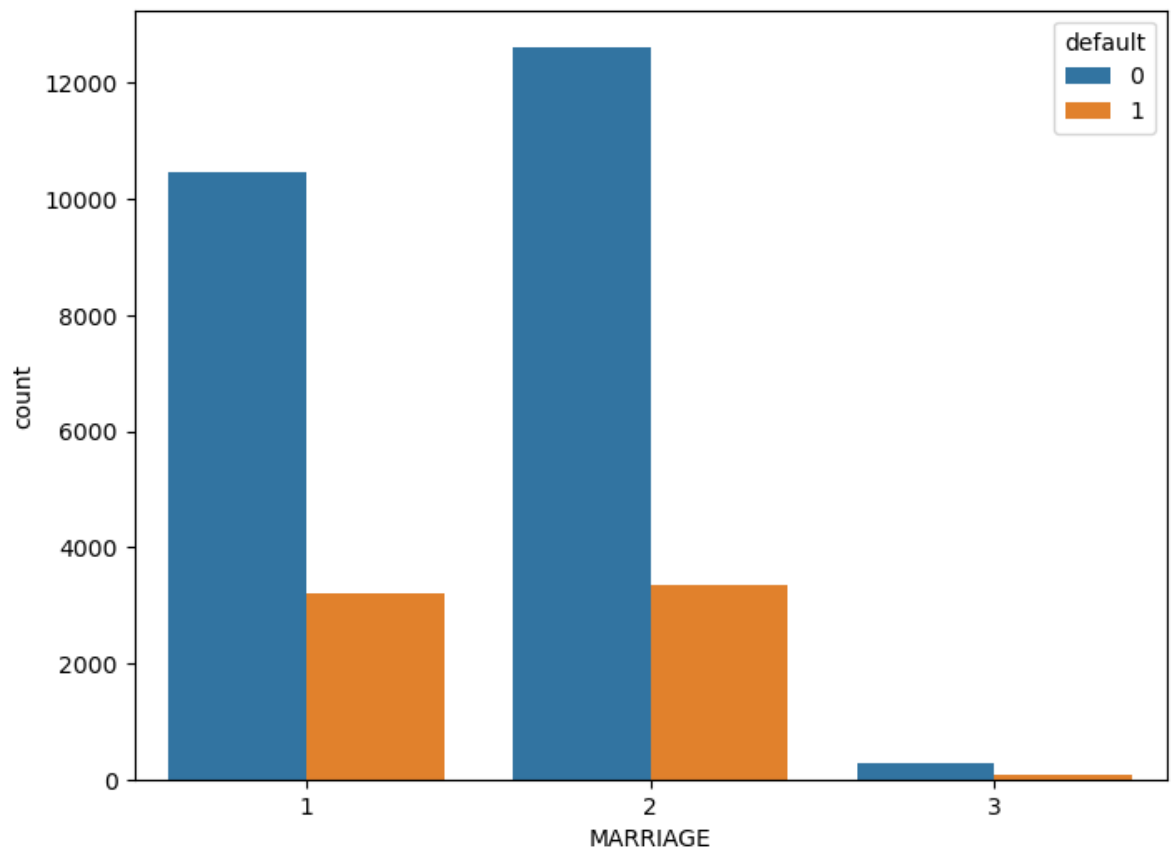
In [31]: 1 df['EDUCATION'].value_counts()
2 labels = ['Graduate', 'University', 'High School', 'Others']
3 values = df['EDUCATION'].value_counts().values
4
5 f, ax = plt.subplots(1,2)
6
7 sns.countplot(x="EDUCATION", data=df, hue='default', ax=ax[1])
8 ax[0].pie(values, labels = labels, autopct='%1.1f%%', shadow = True)
9 plt.show()
10

```



default of credit card holder university degree holder is high compare to others. Graduate users categorized as 'Other' demonstrate a substantial 50% chance of defaulting on their credit card payments.

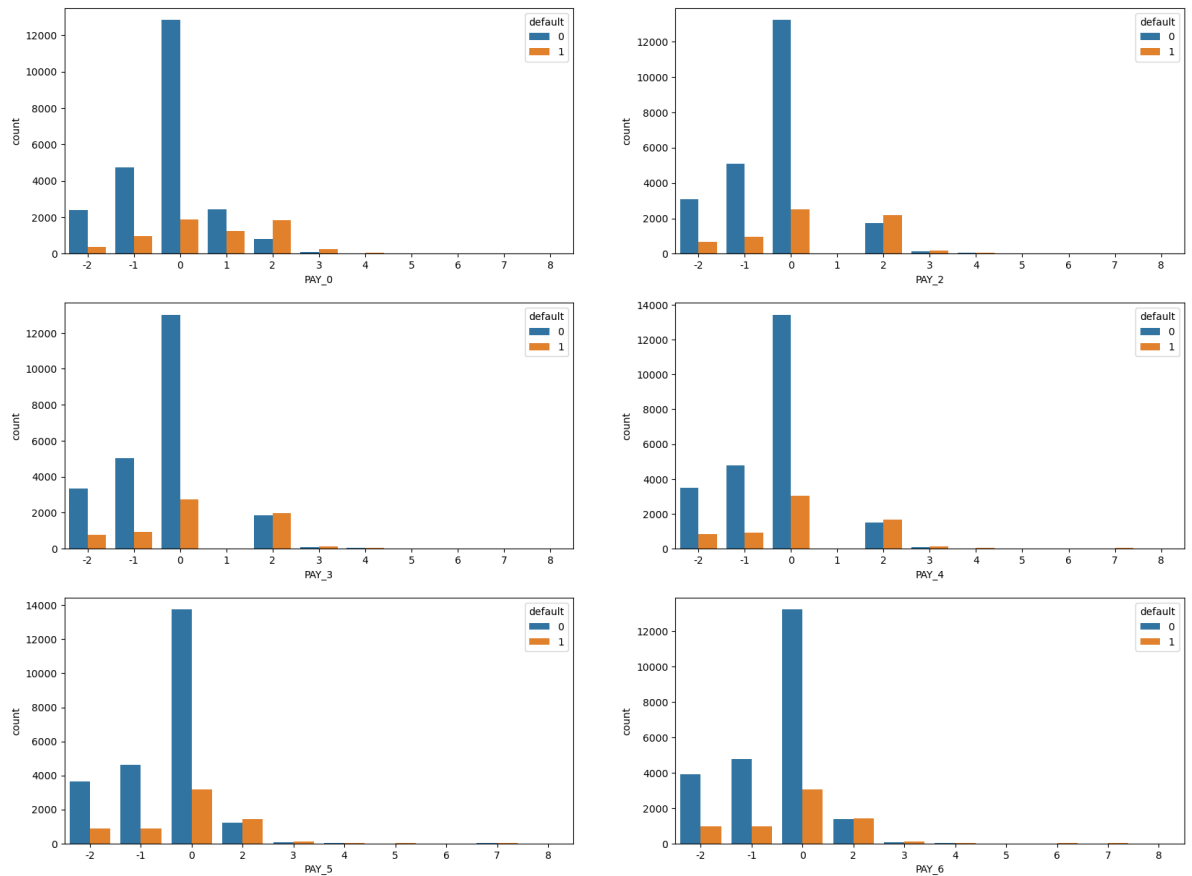
```
In [32]: 1 f, ax = plt.subplots(figsize=(8, 6))
          2 ax = sns.countplot(x="MARRIAGE", data=df, hue='default')
          3 plt.show()
```



```
In [33]: 1 # almost same both married and unmarried as a defaulters
```

In [34]:

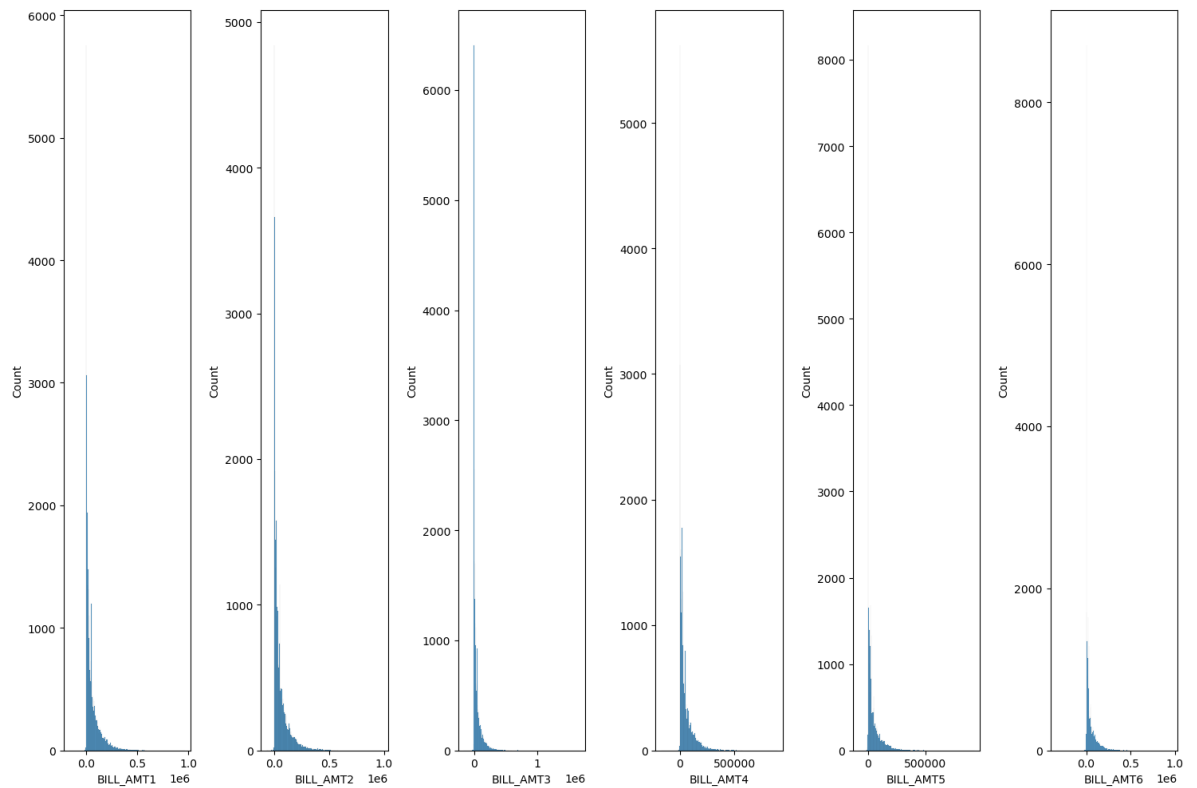
```
1 figure, axes = plt.subplots(3,2,figsize=(20,15))
2 for i,column in enumerate(df_columns):
3     row , col = divmod(i,2)
4     ax=axes[row,col]
5     sns.countplot(data =df, x=column,hue='default',ax=ax)
```



0= revolving credit ,defaulter or non default nut most of users engage with revolving credit

## 5.3 univariate with numerical data

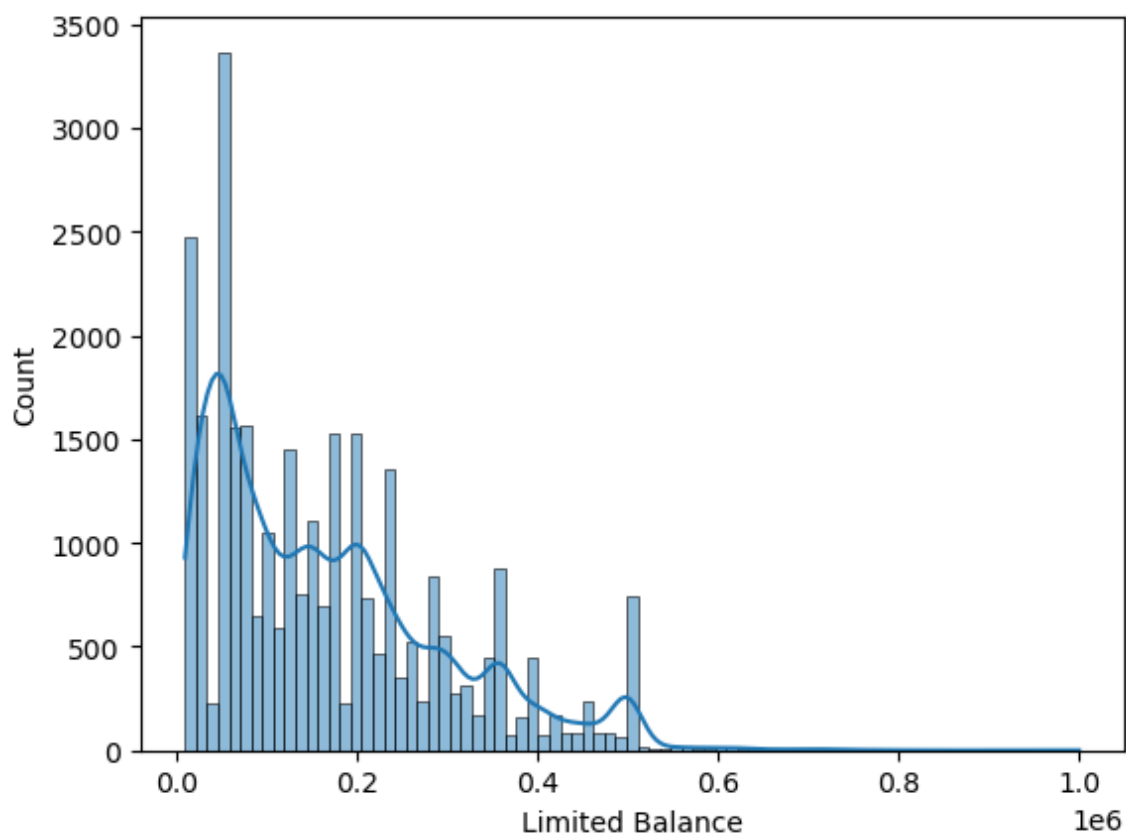
```
In [36]: 1 Previous_bills = ['BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BIL
2 fig, ax = plt.subplots(1,6, figsize=(15,10))
3 for i,bill_columns in enumerate(Previous_bills):
4     sns.histplot(
5         data = df, x = bill_columns,ax=ax[i]
6     )
7     plt.tight_layout()
```



Across all months postively skewness

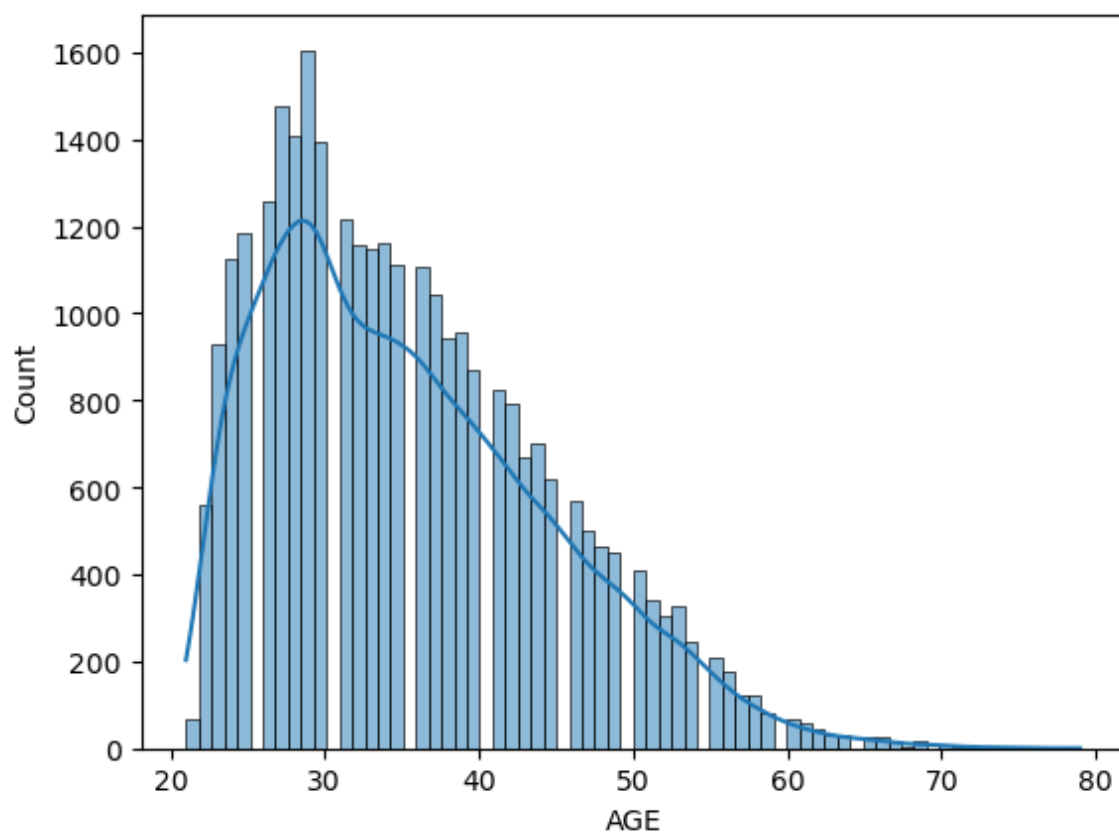
```
In [37]: 1 sns.histplot(df['LIMIT_BAL'], kde=True)
        2 plt.xlabel('Limited Balance')
```

```
Out[37]: Text(0.5, 0, 'Limited Balance')
```



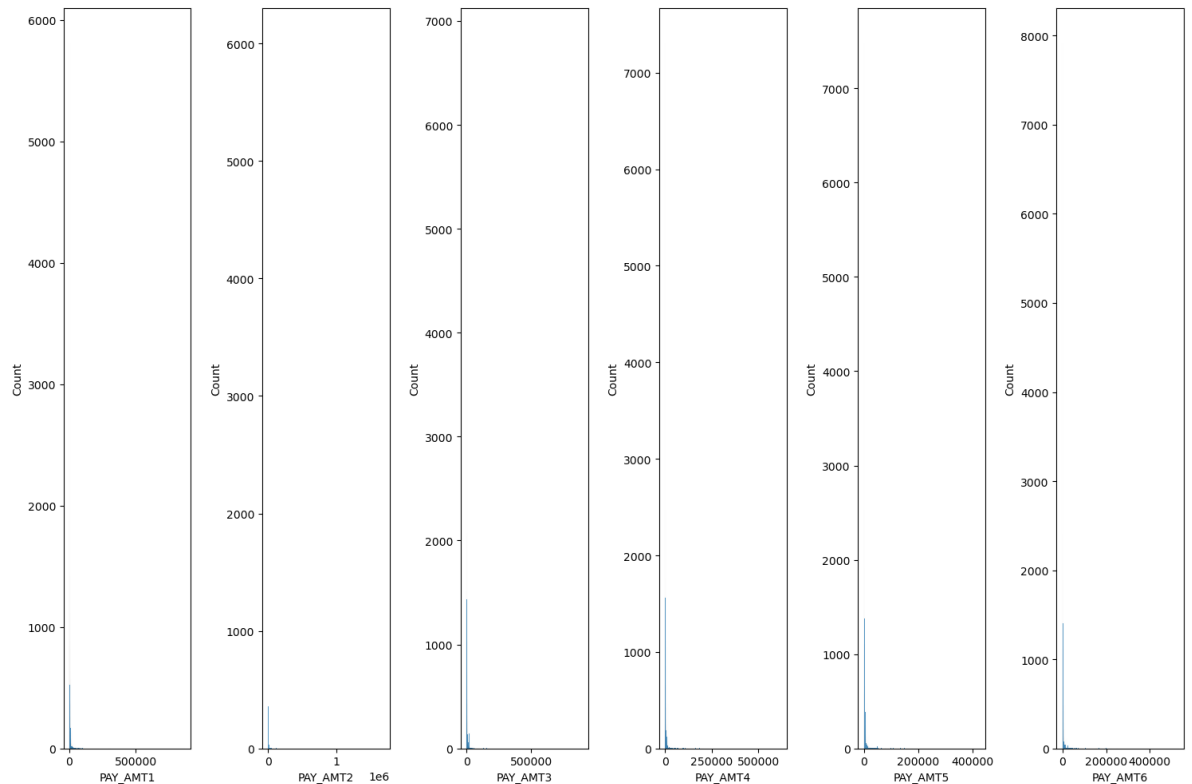
```
In [38]: 1 sns.histplot(df['AGE'], kde=True)
        2 plt.xlabel('AGE')
```

```
Out[38]: Text(0.5, 0, 'AGE')
```





```
In [39]: 1 Previous_payment = ['PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6']
2 fig, ax = plt.subplots(1,6, figsize=(15,10))
3 for i,payment in enumerate(Previous_payment):
4     sns.histplot(
5         data = df, x = payment,ax=ax[i]
6     )
7     plt.tight_layout()
```



## 6.Feature vector and Target variable

```
In [40]: 1
2 X = df.drop(['default'],axis=1)
3 y=df['default']
```

### 6.1 train and test the model

```
In [42]: 1 Xtrain,Xtest,ytrain,ytest = train_test_split(X,y,test_size=0.2,random_state=42)
```

## 7.Feature Scaling

```
In [44]: 1 scaler = StandardScaler()
2 Xtrain = scaler.fit_transform(Xtrain)
3 Xtest= scaler.transform(Xtest)
```

## 8.Model development

```
In [46]: 1 classifiers = [
2         ("Decision Tree", DecisionTreeClassifier()),
3         ("Random Forest", RandomForestClassifier()),
4         ("SVM", SVC()),
5         ("KNN", KNeighborsClassifier()),
6         ("Logistic Regression", LogisticRegression()),
7         ('XGBoost' , XGBClassifier()),
8         ('LGBM' , LGBMClassifier())
9     ]
10
11 # Create an empty DataFrame to store results
12 results_df = pd.DataFrame(columns=["Classifier", "Accuracy"])
13 for method, classifier in classifiers:
14     model = classifier.fit(Xtrain, ytrain)
15     ypred = model.predict(Xtest)
16     accuracy = accuracy_score(ytest, ypred)
17     # Append results to DataFrame
18     results_df = results_df.append({"Classifier": method, "Accuracy": accuracy})
19 print(results_df)
```

```
[LightGBM] [Info] Number of positive: 5339, number of negative: 18661
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.011733 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 3266
[LightGBM] [Info] Number of data points in the train set: 24000, number of used features: 23
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.222458 -> initscore=-1.251397
[LightGBM] [Info] Start training from score -1.251397
```

	Classifier	Accuracy
0	Decision Tree	0.731333
1	Random Forest	0.822333
2	SVM	0.825667
3	KNN	0.789667
4	Logistic Regression	0.819167
5	XGBoost	0.818333
6	LGBM	0.828000

#LGBM algorithm model is good accuracy 82.8% with imbalanced data #svm algorithm model is 82.5% #randomforest model 82.2%

## 9. Dealing with imbalanced data

```
In [47]: 1 df_smote = df.copy()
```

```
In [48]: 1 from imblearn.over_sampling import SMOTE
2 x = df_smote.drop(['default'], axis=1)
3 y = df_smote['default']
```

```
In [49]: 1 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.2, random_state=42)
```

```
In [50]: 1 scaler = StandardScaler()
2 xtrain = scaler.fit_transform(xtrain)
3 xtest= scaler.transform(xtest)
```

```
In [57]: 1 smote = SMOTE(random_state=0)
2 x_sample,y_sample = smote.fit_resample(xtrain,ytrain)
3
```

```
In [58]: 1 ytest.shape,xtest.shape
```

```
Out[58]: ((6000,), (6000, 23))
```

```
In [59]: 1 x_sample.shape,y_sample.shape # after sampling the data shape changed
```

```
Out[59]: ((37322, 23), (37322,))
```

```
In [60]: 1 y_sample.value_counts() # convert imbalance data into balanced data using S
```

```
Out[60]: 0    18661
1    18661
Name: default, dtype: int64
```

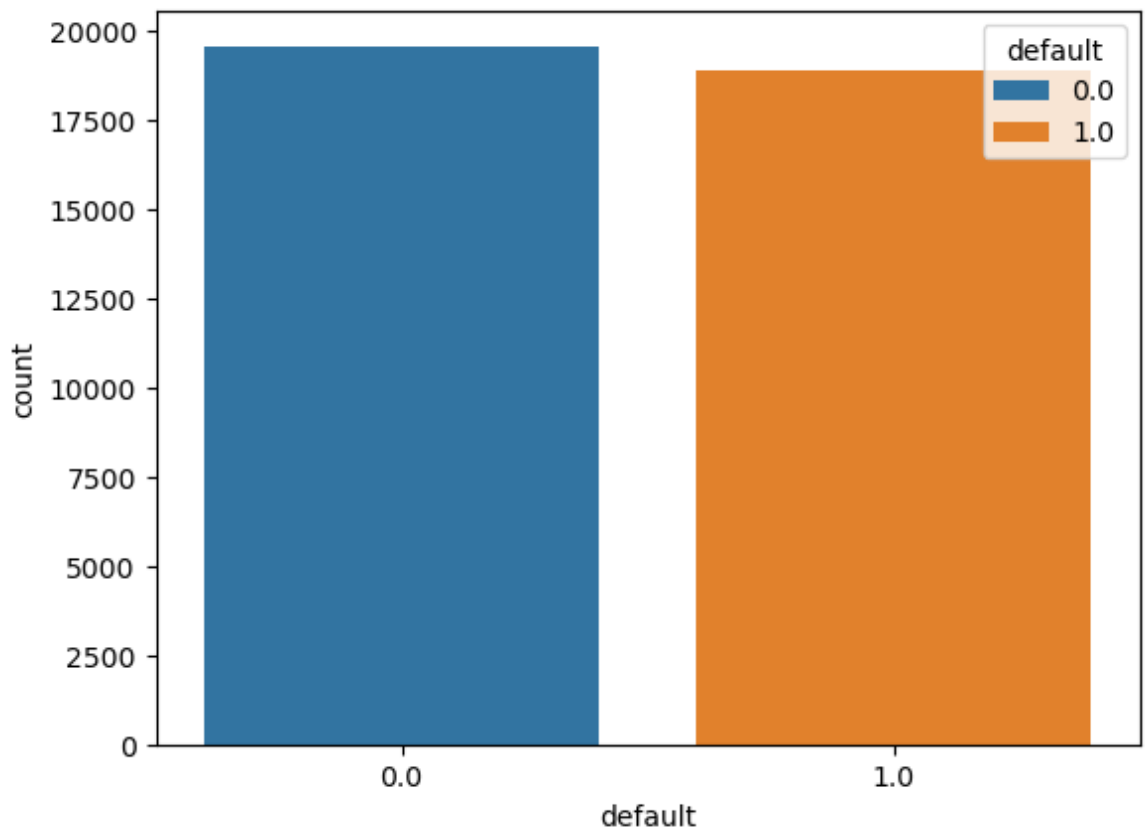
```
In [70]: 1 # Create new DataFrames with the resampled training set and the original te
2 train_resampled_df = pd.DataFrame(x_sample, columns=x.columns)
3 train_resampled_df['default'] = y_sample
4 test_df = pd.DataFrame(xtest, columns=x.columns)
5 test_df['default'] = ytest
```

```
In [71]: 1
2 # Combine the resampled training DataFrame with the original test DataFrame
3 combined_df = pd.concat([train_resampled_df, test_df], ignore_index=True)
4
5 # Save the combined DataFrame as a CSV file
6 combined_df.to_csv('combined_data.csv', index=False)
7
```

```
In [72]: 1 df_smote1 = combined_df.copy()
```

```
In [73]: 1 sns.countplot(x='default',data=df_smote1,hue='default') #see data balance
```

```
Out[73]: <Axes: xlabel='default', ylabel='count'>
```



## 10. Model development with balance data

```
In [74]: 1 classifiers = [
2         ("Decision Tree", DecisionTreeClassifier()),
3         ("Random Forest", RandomForestClassifier()),
4         ("SVM", SVC()),
5         ("KNN", KNeighborsClassifier()),
6         ("Logistic Regression", LogisticRegression()),
7         ('XGBoost' , XGBClassifier()),
8         ('LGBM' , LGBMClassifier())
9     ]
10
11 # Create an empty DataFrame to store results
12 results_df = pd.DataFrame(columns=["Classifier", "Accuracy"])
13 for method, classifier in classifiers:
14     model = classifier.fit(x_sample, y_sample)
15     ypred_smote = model.predict(xtest)
16     accuracy = accuracy_score(ytest, ypred_smote)
17     # Append results to DataFrame
18     results_df = results_df.append({"Classifier": method, "Accuracy": accuracy})
19 print(results_df)
```

```
[LightGBM] [Info] Number of positive: 18661, number of negative: 18661
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.017742 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 5735
[LightGBM] [Info] Number of data points in the train set: 37322, number of used features: 23
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.500000 -> initscore=0.000000
0
```

	Classifier	Accuracy
0	Decision Tree	0.703167
1	Random Forest	0.802000
2	SVM	0.780167
3	KNN	0.652000
4	Logistic Regression	0.674667
5	XGBoost	0.811667
6	LGBM	0.824333

```
In [66]: 1 # here LGBM is Good accuracy 82.4%,XGBoost is 81.1% and RandomForest 80.2%
```

## 10.1 XGBOOST

```
In [76]: 1 # import packages for hyperparameters tuning
2 from hyperopt import STATUS_OK, Trials, fmin, hp, tpe
3 space={ 'max_depth': hp.quniform("max_depth", 3, 18, 1),
4         'gamma': hp.uniform('gamma', 1, 9),
5         'reg_alpha' : hp.quniform('reg_alpha', 40, 180, 1),
6         'reg_lambda' : hp.uniform('reg_lambda', 0, 1),
7         'colsample_bytree' : hp.uniform('colsample_bytree', 0.5, 1),
8         'min_child_weight' : hp.quniform('min_child_weight', 0, 10, 1),
9         'n_estimators': 180,
10         'seed': 0
11     }
```

```
In [77]: 1 import xgboost as xgb
2
3
4 def objective(space):
5     clf=xgb.XGBClassifier(
6         n_estimators =space['n_estimators'], max_depth = int(s
7         reg_alpha = int(space['reg_alpha']),min_child_weight=i
8         colsample_bytree=int(space['colsample_bytree']))
9
10    evaluation = [( x_sample, y_sample), ( xtest, ytest)]
11
12    clf.fit(x_sample, y_sample,
13            eval_set=evaluation, eval_metric="auc",
14            early_stopping_rounds=10,verbose=False)
15
16
17    pred = clf.predict(xtest)
18    accuracy = accuracy_score(ytest, pred)
19    print ("SCORE:", accuracy)
20    return {'loss': -accuracy, 'status': STATUS_OK }
```

```
In [80]: 1 trials = Trials()
          2
          3 best_hyperparams = fmin(fn = objective,
          4                           space = space,
          5                           algo = tpe.suggest,
          6                           max_evals = 100,
          7                           trials = trials)
          8
```

0.0010555555555555

0.8

0.7976

0.8016

0.7941

0.791

0.7891

0.796

0.8

0.7975

```
In [81]: 1 print("The best hyperparameters are : ", "\n")
2 print(best_hyperparams)
3 best_hyperparams['max_depth'] = int(best_hyperparams['max_depth'])
4
```

The best hyperparameters are :

```
{'colsample_bytree': 0.7125085163110019, 'gamma': 2.916107216242168, 'max_depth': 13.0, 'min_child_weight': 4.0, 'reg_alpha': 51.0, 'reg_lambda': 0.2384520699705056}
```

```
In [82]: 1 xgb = XGBClassifier(**best_hyperparams)
2 xgb.fit(x_sample, y_sample)
3 ypredtion=xgb.predict(xtest)
4 print('Final accuracy of the model is {}'.format(accuracy_score(ytest, ypredtion)))
5 print('Classification Report \n{}'.format(classification_report(ytest, ypredtion)))
6 print('Confusion Matrix \n{}'.format(confusion_matrix(ytest, ypredtion)))
7
```

Final accuracy of the model is 79.4%

Classification Report

	precision	recall	f1-score	support
0	0.87	0.87	0.87	4703
1	0.52	0.53	0.53	1297
accuracy			0.79	6000
macro avg	0.70	0.70	0.70	6000
weighted avg	0.80	0.79	0.79	6000

Confusion Matrix

```
[[4076  627]
 [ 609  688]]
```

## 10.2 RandomForest

```
In [85]: 1 from sklearn.model_selection import RandomizedSearchCV
2 rf_classifier = RandomForestClassifier()
3
4 # Define the hyperparameter search space
5 param_dist = {
6     'n_estimators': [int(x) for x in range(10, 200)],
7     'max_features': ['auto', 'sqrt', 'log2', None],
8     'max_depth': [int(x) for x in range(1, 20)],
9     'min_samples_split': [2, 5, 10],
10    'min_samples_leaf': [1, 2, 4],
11
12 }
13
14 # Perform RandomizedSearchCV
15 random_search = RandomizedSearchCV(
16     rf_classifier,
17     param_distributions=param_dist,
18     n_iter=20, # Number of parameter settings that are sampled
19     cv=5,      # Number of cross-validation folds
20     n_jobs=-1, # Use all available CPU cores
21     verbose=1, # Print progress
22     random_state=0
23 )
24
25 # Fit the model to the training data
26 random_search.fit(x_sample, y_sample)
27 best_parameters = random_search.best_params_
28
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

```
In [86]: 1 rf = RandomForestClassifier(**best_parameters)
2 rf.fit(x_sample,y_sample)
3 y_predrf = rf.predict(xtest)
4 print('Final accuracy of the model is {}'.format(accuracy_score(ytest,y_predrf)))
5 print('Classification Report \n{}'.format(classification_report(ytest,y_predrf)))
6 print('Confusion Matrix \n{}'.format(confusion_matrix(ytest,y_predrf)))
```

Final accuracy of the model is 79.7%

Classification Report

	precision	recall	f1-score	support
0	0.86	0.88	0.87	4703
1	0.53	0.48	0.51	1297
accuracy			0.80	6000
macro avg	0.70	0.68	0.69	6000
weighted avg	0.79	0.80	0.79	6000

Confusion Matrix

```
[[4159  544]
 [ 674  623]]
```

## 10.3 LGBM



```

In [ ]: 1 # import packages for hyperparameters tuning
2 from hyperopt import STATUS_OK, Trials, fmin, hp, tpe
3 space={ 'max_depth': hp.uniform("max_depth", 4, 18),
4         'learning_rate': hp.loguniform('learning_rate', np.log(0.01), np.
5         'num_leaves': hp.uniform('num_leaves', 20, 150),
6         'min_child_samples': hp.uniform('min_child_samples', 5, 100),
7         'colsample_bytree' : hp.uniform('colsample_bytree', 0.5,0.7),
8         'subsample':hp.uniform('subsample', 0.5,0.9),
9         'subsample_freq': hp.uniform('subsample_freq', 1,2),
10        'reg_alpha' : hp.uniform('reg_alpha', 40,180),
11        'reg_lambda' : hp.uniform('reg_lambda', 0,1),
12        'max_bin': hp.uniform('max_bin', 75,100),
13        'feature_fraction': hp.uniform('feature_fraction', 0.5, 1.0),
14        'bagging_fraction': hp.uniform('bagging_fraction', 0.5, 1.0),
15        'bagging_freq': hp.uniform('bagging_freq', 1, 10) ,
16        'nthread': 8,
17        'verbose': 0
18
19    }

```

```

In [94]: 1 import lightgbm as lgb
2
3
4 def objective(space):
5
6     clf=lgb.LGBMClassifier(
7         num_leaves =int(space['num_leaves']), max_depth = int(
8         objective = 'binary',min_child_samples=int(space['min_ch
9         colsample_bytree=space['colsample_bytree'],subsample =
10        subsample_freq = space['subsample_freq'],bagging_freq=
11        n_estimators=150)
12    evaluation = [( x_sample, y_sample), ( xtest, ytest)]
13
14    clf.fit(x_sample, y_sample,
15            eval_set=evaluation, eval_metric="auc")
16    pred = clf.predict(xtest)
17    accuracy = accuracy_score(ytest, pred)
18    print ("SCORE:", accuracy)
19    return { 'loss': -accuracy, 'status': STATUS_OK }

```

In [95]:

```
1 trials = Trials()
2
3 best_hyperparams = fmin(fn = objective,
4                          space = space,
5                          algo = tpe.suggest,
6                          max_evals = 100,
7                          trials = trials)

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] bagging_freq is set=2, subsample_freq=1.52384056321214
57 will be ignored. Current value: bagging_freq=2
SCORE:
0.814
[LightGBM] [Warning] bagging_freq is set=1, subsample_freq=1.17992371880536
97 will be ignored. Current value: bagging_freq=1
[LightGBM] [Warning] bagging_freq is set=1, subsample_freq=1.17992371880536
97 will be ignored. Current value: bagging_freq=1
[LightGBM] [Info] Number of positive: 18661, number of negative: 18661
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of t
esting was 0.013212 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 5735
```

In [103]:

```
1 print("The best hyperparameters are : ", "\n")
2 print(best_hyperparams)
3 best_hyperparams['num_leaves'] = int(best_hyperparams['num_leaves'])
4 best_hyperparams['max_depth'] = int(best_hyperparams['max_depth'])
5 best_hyperparams['bagging_freq'] = int(best_hyperparams['bagging_freq'])
6 best_hyperparams['max_bin'] = int(best_hyperparams['max_bin'])
```

The best hyperparameters are :

```
{'bagging_fraction': 0.6828580064307468, 'bagging_freq': 2, 'colsample_bytre
e': 0.5130684539621747, 'feature_fraction': 0.9684276671090114, 'learning_rat
e': 0.015223939585351605, 'max_bin': 78.63726078924715, 'max_depth': 14, 'min
_child_samples': 24.75726287535893, 'num_leaves': 24, 'reg_alpha': 129.651697
35764212, 'reg_lambda': 0.5362968135570422, 'subsample': 0.8134663323186887,
'subsample_freq': 1.3343624216548728}
```

In [104]:

```
1 lgb = LGBMClassifier(**best_hyperparams,min_data_in_leaf=int(24.7))
2 lgb.fit(x_sample,y_sample)
3 ypredtion=lgb.predict(xtest)
4 print('Final accuracy of the model is {}'.format(accuracy_score(ytest,ypredtion)))
5 print('Classification Report \n{}'.format(classification_report(ytest,ypredtion)))
6 print('Confusion Matrix \n{}'.format(confusion_matrix(ytest,ypredtion)))
```

```

[LightGBM] [Warning] min_data_in_leaf is set=24, min_child_samples=24.7572628
7535893 will be ignored. Current value: min_data_in_leaf=24
[LightGBM] [Warning] feature_fraction is set=0.9684276671090114, colsample_by
tree=0.5130684539621747 will be ignored. Current value: feature_fraction=0.96
84276671090114
[LightGBM] [Warning] bagging_fraction is set=0.6828580064307468, subsample=0.
8134663323186887 will be ignored. Current value: bagging_fraction=0.682858006
4307468
[LightGBM] [Warning] bagging_freq is set=2, subsample_freq=1.3343624216548728
will be ignored. Current value: bagging_freq=2
[LightGBM] [Warning] min_data_in_leaf is set=24, min_child_samples=24.7572628
7535893 will be ignored. Current value: min_data_in_leaf=24
[LightGBM] [Warning] feature_fraction is set=0.9684276671090114, colsample_by
tree=0.5130684539621747 will be ignored. Current value: feature_fraction=0.96
84276671090114
[LightGBM] [Warning] bagging_fraction is set=0.6828580064307468, subsample=0.
8134663323186887 will be ignored. Current value: bagging_fraction=0.682858006
4307468
[LightGBM] [Warning] bagging_freq is set=2, subsample_freq=1.3343624216548728
will be ignored. Current value: bagging_freq=2
[LightGBM] [Info] Number of positive: 18661, number of negative: 18661
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of tes
ting was 0.008318 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 1730
[LightGBM] [Info] Number of data points in the train set: 37322, number of us
ed features: 23
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.500000 -> initscore=0.00000
0
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] min_data_in_leaf is set=24, min_child_samples=24.7572628
7535893 will be ignored. Current value: min_data_in_leaf=24
[LightGBM] [Warning] feature_fraction is set=0.9684276671090114, colsample_by
tree=0.5130684539621747 will be ignored. Current value: feature_fraction=0.96
84276671090114
[LightGBM] [Warning] bagging_fraction is set=0.6828580064307468, subsample=0.
8134663323186887 will be ignored. Current value: bagging_fraction=0.682858006
4307468
[LightGBM] [Warning] bagging_freq is set=2, subsample_freq=1.3343624216548728
will be ignored. Current value: bagging_freq=2
Final accuracy of the model is 77.21666666666667%
Classification Report

```

	precision	recall	f1-score	support
0	0.88	0.83	0.85	4703
1	0.48	0.58	0.52	1297
accuracy			0.77	6000
macro avg	0.68	0.70	0.69	6000
weighted avg	0.79	0.77	0.78	6000

Confusion Matrix

```
[[3880  823]
 [ 544  753]]
```

## 11.checking overfitting and underfitting

```
In [91]: 1 # print the scores on training and test set
          2
          3 print('Training set score: {:.4f}'.format(rf.score(x_sample, y_sample)))
          4
          5 print('Test set score: {:.4f}'.format(rf.score(xtest, ytest)))
```

Training set score: 0.9473

Test set score: 0.7970

#The training set score is significantly higher than the test set score. This suggests that the model is likely overfitting.

## 12.K-Fold cross validation

```
In [93]: 1 from sklearn.model_selection import cross_val_score
          2 accuracies = cross_val_score(estimator = rf, X = x_sample,y = y_sample, cv
          3 print("Accuracy: {:.2f} %".format(accuracies.mean()*100))
          4 print("Standard Deviation: {:.2f} %".format(accuracies.std()*100))
```

Accuracy: 83.68 %

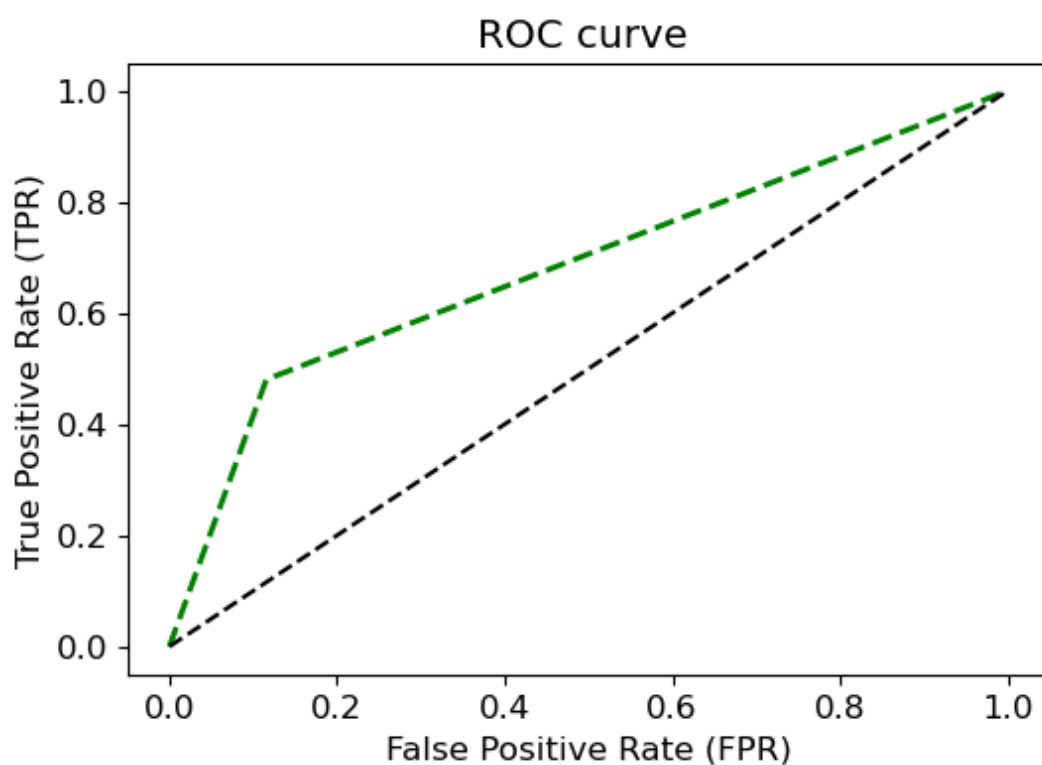
Standard Deviation: 6.81 %

## 13.ROC-AUC

```

In [101]: 1 from sklearn.metrics import roc_curve
2
3 fpr, tpr, thresholds = roc_curve(ytest, y_predrf)
4
5 plt.figure(figsize=(6,4))
6
7 plt.plot(fpr, tpr, linestyle='--', color = 'green',linewidth=2)
8
9 plt.plot([0,1], [0,1], 'k--' )
10
11 plt.rcParams['font.size'] = 12
12
13 plt.title('ROC curve ')
14
15 plt.xlabel('False Positive Rate (FPR)')
16
17 plt.ylabel('True Positive Rate (TPR)')
18
19 plt.show()
20

```



```

In [98]: 1 from sklearn.metrics import roc_auc_score
2
3 ROC_AUC = roc_auc_score(ytest, y_predrf)
4
5 print('ROC AUC : {:.4f}'.format(ROC_AUC))

```

ROC AUC : 0.6823

```
In [99]: 1
2 from sklearn.model_selection import cross_val_score
3
4 Cross_validated_ROC_AUC = cross_val_score(rf, x_sample, y_sample, cv=5, sc
5
6 print('Cross validated ROC AUC : {:.4f}'.format(Cross_validated_ROC_AUC))
```

Cross validated ROC AUC : 0.9097

## 14.Conclusion

```
1 In conclusion, In this project, I build a all classification models to
2 classify thecredit card defaulter.
3 we get 3 better models RandomForest,XGBoost,LGBM .
4 In XGBoost hyperparameter tuning using HyperOPT ,the model accuracy is
5 79.4%
6 In RandomForest the model accuracy is 79.7% using the model tuning
7 technique RandomizedSearchCV
8 In LGBM hyperparameter tuning using HyperOPT, the model acuracy is 77.2%
9 Finally go with RandomForest Classifier
10 The training-set accuracy score is 0.9473 while the test-set accuracy to
11 be 0.797. So, the model is overfitting.
12 If we look at all the 10 scores produced by the 10-fold cross-validation,
13 we can also conclude that there is a relatively low variance in the
14 accuracy between folds, ranging from 100% accuracy to 83.68% accuracy.
15 ROC AUC score is 68.23%
16 Cross validated ROC AUC score is 90.97%
```

[GotoTOP](#)

In [ ]:

```
1
```