

```
In [12]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')

In [12]: df=pd.read_csv(r'C:\Users\sharad\Downloads\l1\classification dataset\loan_data.csv')

In [13]: df.head()

Out[13]:
   credit_policy  purpose  int.rate  installment  log.annual.inc  dt  fico  days.with.cr.line  revol.bal  revol.util  lng.last.6mths  delinq.2yrs  pub.rec  not.fully.paid
0      1  debt_consolidation  0.1380      529.10  11.300407  19.48  737      9639.968333      28864      52.1      0      0      0      0
1      1      credit_card  0.1071      228.22  11.021243  14.29  707      2760.000000      23823      76.7      0      0      0      0
2      1  debt_consolidation  0.1357      366.86  10.373491  11.63  682      4710.000000      3511      25.6      1      0      0      0
3      1  debt_consolidation  0.1008      162.34  11.356407  8.10  712      2699.668333      23867      73.2      1      0      0      0
4      1      credit_card  0.1426      102.92  11.299732  14.97  667      4066.000000      4740      39.5      0      1      0      0

In [14]: df['paid']=df['not.fully.paid']

In [15]: df.drop('not.fully.paid',inplace=True,axis=1)

In [21]: df_cat=df.purpose
from sklearn.preprocessing import LabelEncoder as le

In [22]: le=[]
df_cat=le.fit_transform(df_cat)

In [27]: sns.heatmap(df.isnull())
plt.show()

Out[17]:
<seaborn.axisgrid.FacetGrid at 0xcab8757760>

In [40]: df.paid.value_counts()

Out[40]:
0    8945
1    1033
Name: paid, dtype: int64

In [36]: sns.pairgrid(df hue='paid')

Out[16]:
<seaborn.axisgrid.PairGrid at 0xcab8757760>

In [12]: plt.figure(figsize=(15,12))
sns.heatmap(df.corr(),annot=True)

credit_policy  1  0.018  -0.29  0.005  0.005  -0.001  0.15  0.099  -0.13  0.1  -0.54  -0.976  -0.954  -0.16
purpose        0.018  1  0.12  0.05  0.089  0.053  0.057  0.048  0.042  -0.003  0.038  -0.012  0.0048  0.048
int.rate       -0.29  0.12  1  0.38  0.054  0.22  0.71  0.12  0.033  0.45  0.2  0.16  0.088  0.16
installment    -0.005  0.05  0.28  1  0.45  0.05  0.086  0.18  0.23  0.081  -0.01  -0.0044  -0.033  0.05
log.annual.inc -0.005  0.009  0.056  0.45  1  -0.044  0.13  0.34  0.17  0.095  0.029  0.019  0.037  -0.033
dt             -0.001  -0.033  0.22  0.05  0.054  1  0.24  0.06  0.19  0.34  0.029  -0.022  0.0002  0.017
fico           0.15  0.097  -0.71  0.088  0.11  -0.14  1  0.36  0.026  0.54  0.19  0.12  -0.15  -0.15
days.with.cr.line -0.099  0.048  -0.12  0.18  0.34  0.06  0.26  1  0.73  0.024  0.542  0.081  0.072  -0.029
revol.bal      0.19  0.042  0.001  0.23  0.37  0.19  0.038  0.27  1  0.2  0.022  -0.033  0.035  0.054
revol.util     -0.1  -0.053  0.46  0.081  0.055  0.38  -0.54  0.024  0.2  1  -0.014  -0.043  0.067  0.062
lng.last.6mths -0.54  0.036  0.2  -0.05  0.029  0.029  -0.19  -0.042  0.022  -0.014  1  0.021  0.075  0.15
delinq.2yrs    -0.076  -0.032  0.16  -0.064  0.029  -0.022  0.22  0.081  0.033  -0.043  0.021  1  0.0002  0.0089
pub.rec        -0.054  0.048  0.088  -0.033  0.017  0.0002  -0.15  0.072  -0.031  0.067  0.073  0.0002  1  0.045
paid           0.16  0.048  0.16  0.05  -0.033  0.037  -0.15  -0.029  0.054  0.082  0.15  0.0089  0.048  1

In [45]: from scipy.stats import skew

In [51]: for col in df:
print(col, ' :-- ',end='')
print(skew(df[col]))
sns.distplot(df[col])
plt.show()

credit_policy  :-- -1.5393803114408613

purpose        :-- 0.9018541173261189

int.rate       :-- 0.164339416289317964

installment    :-- 0.9123795474881175

log.annual.inc :-- 0.928653616694390843

dt             :-- 0.02393727346996256

fico           :-- 0.711859325420144

days.with.cr.line :-- 1.1555672187692227

revol.bal      :-- 11.15531949753894

revol.util     :-- 0.95997684792027563

lng.last.6mths :-- 3.583589522653665

delinq.2yrs    :-- 6.869843937622224

pub.rec        :-- 5.125631575999643

paid           :-- 1.8543012464275515

In [78]: df.purpose=df_cat
df=df.iloc[:,1:]
y=df['paid']

In [79]: from sklearn.preprocessing import StandardScaler as sc

In [80]: s=sc()
X=sc.fit(X)
X=sc.transform(X)

In [105]: from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report

In [156]: models=[(Logistic Regression, LogisticRegression()),
                  ('Decision Tree', DecisionTreeClassifier()),
                  ('Support Vector', SVC()),
                  ('Knearest Neighbors', KNeighborsClassifier()),
                  ('Random Forest', RandomForestClassifier()),
                  ('AdaBoostClassifier', AdaBoostClassifier()),
                  ('Gradient Boosting', GradientBoostingClassifier()),
                  ('Xtreme Gradient Boosting', XGBClassifier())
                ]

In [157]: xtrain,xtest,ytrain,ytest=train_test_split(x,y,random_state=0,test_size=0.30)

In [158]: accuracy=[]

for name,model in models:
print(name,' :-- ')
print()
model.fit(xtrain,ytrain)
ypred=model.predict(xtest)
print(confusion_matrix(ytest,ypred))
print()
print(classification_report(ytest,ypred))
accuracy.append(accuracy_score(ytest,ypred))

Logistic Regression :--

[[2282  34]
 [ 465  3]]

precision recall f1-score support
0      0.84      0.99      0.91      2486
1      0.39      0.03      0.05      468
accuracy      0.83      2874
macro avg      0.61      0.59      0.46      2874
weighted avg      0.74      0.83      0.77      2874

Decision Tree :--

[[2807 289]
 [384 136]]

precision recall f1-score support
0      0.85      0.83      0.84      2486
1      0.22      0.24      0.23      468
accuracy      0.83      2874
macro avg      0.54      0.54      0.54      2874
weighted avg      0.75      0.74      0.77      2874

Support vector :--

[[2484  2]
 [ 466  2]]

precision recall f1-score support
0      0.84      0.99      0.91      2486
1      0.59      0.09      0.01      468
accuracy      0.87      0.59      0.46      2874
macro avg      0.78      0.84      0.76      2874

Nearest Neighbors :--

[[2324  82]
 [ 441  27]]

precision recall f1-score support
0      0.84      0.97      0.90      2486
1      0.25      0.06      0.09      468
accuracy      0.82      2874
macro avg      0.54      0.51      0.59      2874
weighted avg      0.74      0.82      0.77      2874

Random Forest :--

[[2386  20]
 [ 458  12]]

precision recall f1-score support
0      0.84      0.99      0.91      2486
1      0.36      0.03      0.05      468
accuracy      0.61      0.51      0.48      2874
weighted avg      0.76      0.83      0.77      2874

AdaBoostClassifier :--

[[2388  18]
 [ 458  12]]

precision recall f1-score support
0      0.84      0.99      0.91      2486
1      0.42      0.03      0.05      468
accuracy      0.83      0.51      0.48      2874
weighted avg      0.77      0.84      0.77      2874

Gradient Boosting :--

[[2386  20]
 [ 455  13]]

precision recall f1-score support
0      0.84      0.99      0.91      2486
1      0.39      0.03      0.05      468
accuracy      0.82      0.51      0.48      2874
weighted avg      0.77      0.83      0.77      2874

Xtreme Gradient Boosting :--

[18:02:10] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:10] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:10] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:10] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

precision recall f1-score support
0      0.84      0.97      0.90      2486
1      0.34      0.09      0.14      468
accuracy      0.82      2874
macro avg      0.59      0.53      0.52      2874
weighted avg      0.76      0.82      0.78      2874

In [168]: print('the mean accuracy of the model :-- ',np.array(accuracy).mean()*100)

the mean accuracy of the model :-- 81.94580422407795

In [143]: xgb=XGBClassifier(learning_rate=0.16, gamma=20)
xgb.fit(xtrain,ytrain)
ypred=xgb.predict(xtest)
print('training score :-- ',xgb.score(xtrain,ytrain))*100
print('testing score :-- ',xgb.score(xtest,ytest))*100

[18:58:06] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:58:07] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:58:07] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:58:07] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

training score :-- 84.1206181384248
testing score :-- 83.7360751605762

In [147]: skf=StratifiedKFold(n_splits=5)
skf.get_n_splits(X,y)

Out[147]: 5

In [161]: accuracy=[]

for train,test in skf.split(x,y):
xtrain,xtest=train,y[train],y[test]

xgb=XGBClassifier(learning_rate=0.16, gamma=20)
xgb.fit(xtrain,ytrain)
ypred=xgb.predict(xtest)
accuracy.append(accuracy_score(ytest,ypred))

[18:02:02] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:03] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:03] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[18:02:03] WARNING: C:\Users\Administrator\workspace\xgboost-win64-release-1.5.1\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'er' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

precision recall f1-score support
0      0.84      0.97      0.90      2486
1      0.34      0.09      0.14      468
accuracy      0.82      2874
macro avg      0.59      0.53      0.52      2874
weighted avg      0.76      0.82      0.78      2874

Because the dataset is imbalance we are getting less accuracy

In [173]: print('this is the minimum accuracy we can get from xtreme boosting classifier :-- ',np.array(accuracy).mean()*100)

this is the minimum accuracy we can get from xtreme boosting classifier :-- 73.10162599410215

In [ ]:
```