

Ganeshan M

AM.

EN.U4CSE19320

Phase 2

Problem Definition

Banks play a crucial role in market economies. They decide who can get finance and on what terms and can make or break investment decisions. For markets and society to function, individuals and companies need access to credit.

Credit scoring algorithms, which make a guess at the probability of default, are the method banks use to determine whether or not a loan should be granted. This competition requires participants to improve on the state of the art in credit scoring, by predicting the probability that somebody will experience financial distress in the next two years.

The goal is to build a model that borrowers can use to help make the best financial decisions.

Improve on the state of the art in credit scoring by predicting the probability that somebody will experience financial distress in the next two years.

[GitHub \(https://github.com/mganeshan29/Credit-Prediction-model\)](https://github.com/mganeshan29/Credit-Prediction-model)

Dataset

Give Me Some Credit

In this Model, credit scoring data sets from Kaggle competition called 'Give me some credit' is used to build classifiers

Training dataset [cs-training.csv file \(https://www.kaggle.com/brycecf/give-me-some-credit-dataset?select=cs-training.csv\)](https://www.kaggle.com/brycecf/give-me-some-credit-dataset?select=cs-training.csv) which will be used for model training and test data [cs-test.csv file \(https://www.kaggle.com/brycecf/give-me-some-credit-dataset?select=cs-test.csv\)](https://www.kaggle.com/brycecf/give-me-some-credit-dataset?select=cs-test.csv).

Prepare Data

In this part, the data and variables is done. Target ratio, variables values distributions are also be investigated. Null value analysis is done and null values in variables, filled with statistical approach to prevent them to influence modelling in a bad way.

After the describing inputs, missing values are also investigated, and for 2 inputs ('MonthlyIncome' , 'NumberOfDependents') some missing values are observed. Handling with missing values is very important to create accurate models.

Data Dictionary :

| Variable Name | Description | Type |
|--------------------------------------|--|------------|
| SeriousDlqin2yrs | Person experienced 90 days past due delinquency or worse | Y/N |
| RevolvingUtilizationOfUnsecuredLines | Total balance on credit cards and personal lines of credit except real estate and no installment debt like car loans divided by the sum of credit limits | percentage |
| age | Age of borrower in years | integer |
| NumberOfTime30-59DaysPastDueNotWorse | Number of times borrower has been 30-59 days past due but no worse in the last 2 years. | integer |
| DebtRatio | Monthly debt payments, alimony, living costs divided by monthly gross income | percentage |
| MonthlyIncome | Monthly income | real |
| NumberOfOpenCreditLinesAndLoans | Number of Open loans (installment like car loan or mortgage) and Lines of credit (e.g. credit cards) | integer |
| NumberOfTimes90DaysLate | Number of times borrower has been 90 days or more past due. | integer |
| NumberRealEstateLoansOrLines | Number of mortgage and real estate loans including home equity lines of credit | integer |
| NumberOfTime60-89DaysPastDueNotWorse | Number of times borrower has been 60-89 days past due but no worse in the last 2 years. | integer |
| NumberOfDependents | Number of dependents in family excluding themselves (spouse, children etc.) | integer |

Python packages

NumPy

NumPy is a well known general-purpose array-processing package. An extensive collection of high complexity mathematical functions make NumPy powerful to process large multi-dimensional arrays and matrices. NumPy is very useful for handling linear algebra, Fourier transforms, and random numbers.

Scikit-learn

The Python library, Scikit-Learn, is built on top of the matplotlib, NumPy, and SciPy libraries. This Python ML library has several tools for data analysis and data mining tasks.

Pandas

Pandas is a popular Python library for data analysis. It is not directly related to Machine Learning. As we know that the dataset must be prepared before training. In this case, Pandas comes handy as it was developed specifically for data extraction and preparation.

Matplotlib

Matplotlib is a very popular Python library for data visualization. Like Pandas, it is not directly related to Machine Learning. It particularly comes in handy when a programmer wants to visualize the patterns in the data. It is a 2D plotting library used for creating 2D graphs and plots.

seaborn

Seaborn is a library for making statistical graphs in Python. It is built on top of matplotlib and also integrated with pandas data structures.

Scikit-plot

Scikit-plot is the result of an unartistic data scientist's dreadful realization that visualization is one of the most crucial components in the data science process, not just a mere afterthought.

Importing Libraries

```
In [17]: import itertools
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

import scikitplot as skplt
from sklearn import metrics
import plotly.graph_objs as go
from plotly.offline import iplot
from sklearn.metrics import f1_score
from sklearn.metrics import precision_score
from sklearn.metrics import confusion_matrix

from sklearn import model_selection
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeRegressor
from sklearn.datasets import make_classification
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, auc, recall_score, roc_auc_score

from sklearn.metrics import recall_score
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
```

Data analysis

```
In [18]: train_data = pd.read_csv("cs-training.csv")
test_data = pd.read_csv("cs-test.csv")
```

Review data

```
In [19]: print(train_data.describe)
print(train_data.head)
```

```

2          3          0          0.658180
3          4          0          0.233810
4          5          0          0.907239
...      ...      ...      ...
149995    149996          0          0.040674
149996    149997          0          0.299745
149997    149998          0          0.246044
149998    149999          0          0.000000
149999    150000          0          0.850283

```

```

      age  NumberOfTime30-59DaysPastDueNotWorse  DebtRatio  MonthlyIncome  \
0      45                                   2      0.802982         9120.0
1      40                                   0      0.121876         2600.0
2      38                                   1      0.085113         3042.0
3      30                                   0      0.036050         3300.0
4      49                                   1      0.024926        63588.0
...      ...      ...      ...      ...      ...
149995    74                                   0      0.225131         2100.0
149996    44                                   0      0.716562         5584.0
149997    58                                   0      0.387000         3000.0

```

```
In [20]: train_data.shape
```

```
Out[20]: (150000, 12)
```

```
In [21]: y = train_data.iloc[:,0].values #Taking first col (credit worthiness value)
X = train_data.iloc[:, 1:11].values #Taking the rest of the cols
```

Print columns headers of the dataset

```
In [22]: train_data.columns
```

```
Out[22]: Index(['Unnamed: 0', 'SeriousDlqin2yrs',  
              'RevolvingUtilizationOfUnsecuredLines', 'age',  
              'NumberOfTime30-59DaysPastDueNotWorse', 'DebtRatio', 'MonthlyIncome',  
              'NumberOfOpenCreditLinesAndLoans', 'NumberOfTimes90DaysLate',  
              'NumberRealEstateLoansOrLines', 'NumberOfTime60-89DaysPastDueNotWorse',  
              'NumberOfDependents'],  
             dtype='object')
```

Drop lines with Missing data

```
In [23]: train_data.dropna(axis=1)  
cols_with_missing = [col for col in train_data.columns  
                     if train_data[col].isnull().any()]  
reduced_train_data = train_data.drop(cols_with_missing, axis=1)  
reduced_test_data = test_data.drop(cols_with_missing, axis=1)
```

```
In [24]: reduced_train_data.columns
```

```
Out[24]: Index(['Unnamed: 0', 'SeriousDlqin2yrs',  
              'RevolvingUtilizationOfUnsecuredLines', 'age',  
              'NumberOfTime30-59DaysPastDueNotWorse', 'DebtRatio',  
              'NumberOfOpenCreditLinesAndLoans', 'NumberOfTimes90DaysLate',  
              'NumberRealEstateLoansOrLines', 'NumberOfTime60-89DaysPastDueNotWorse'],  
             dtype='object')
```

Applying Machine Learning Algorithms for Classification Problem

logistic-regression

Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables.

specify the target variable

```
In [25]: train_y = reduced_train_data.SeriousDlqin2yrs  
test_y = reduced_test_data.SeriousDlqin2yrs
```

create list of features

```
In [26]: feature_names = ['RevolvingUtilizationOfUnsecuredLines', 'age',  
                          'NumberOfTime30-59DaysPastDueNotWorse', 'DebtRatio',  
                          'NumberOfOpenCreditLinesAndLoans', 'NumberOfTimes90DaysLate',  
                          'NumberRealEstateLoansOrLines', 'NumberOfTime60-89DaysPastDueNotWorse']  
#create data corresponding to the features  
train_X = reduced_train_data[feature_names]  
test_X = reduced_test_data[feature_names]
```

Spilting Dataset into Train and Test set.

```
In [27]: x_train, x_test, y_train, y_test = train_test_split(train_X, train_y, test_size = 0.3, random_state
```

Logistic regression User Defined

In [28]:

```

class logistic_regression:
    def __init__(self,x,y):
        self.intercept = np.ones((x.shape[0], 1))
        self.x = np.concatenate((self.intercept, x), axis=1)
        self.weight = np.zeros(self.x.shape[1])
        self.y = y

    def sigmoid(self, x, weight):
        z = np.dot(x, weight)
        return 1 / (1 + np.exp(-z))

    def loss(self, h, y):
        return (-y * np.log(h) - (1 - y) * np.log(1 - h)).mean()

    def gradient_descent(self, X, h, y):
        return np.dot(X.T, (h - y)) / y.shape[0]

    def fit(self, lr , iterations):
        for i in range(iterations):
            sigma = self.sigmoid(self.x, self.weight)
            loss = self.loss(sigma,self.y)
            dW = self.gradient_descent(self.x , sigma, self.y)
            #Updating the weights
            self.weight -= lr * dW
        return print('Working successfully')

    def predict(self, x_new , treshhold):
        x_new = np.concatenate((self.intercept, x_new), axis=1)
        result = self.sigmoid(x_new, self.weight)
        result = result >= treshhold
        y_pred = np.zeros(result.shape[0])
        for i in range(len(y_pred)):
            if result[i].any() == True:
                y_pred[i] = 1
            else:
                continue

        return y_pred

```



```
In [29]: regressor = logistic_regression(X,y)
regressor.fit(0.1 , 5000)
pred_lr = regressor.predict(X,0.5)
```

Working successfully

```
In [30]: logmodel = LogisticRegression(random_state=42)
logmodel.fit(x_train,y_train)
y_pred_l = logmodel.predict(x_test)
pred_lr = y_pred_l.astype(int)
```

```
In [31]: roc=roc_auc_score(y_test, pred_lr)
acc = accuracy_score(y_test, pred_lr)
prec = precision_score(y_test, pred_lr)
results = pd.DataFrame(['Logistic Regression', acc,prec,roc]),
columns = ['Model', 'Accuracy', 'Precision', 'ROC_AUC'])
results
```

Out[31]:

| | Model | Accuracy | Precision | ROC_AUC |
|---|---------------------|----------|-----------|----------|
| 0 | Logistic Regression | 0.932778 | 0.532468 | 0.512651 |

Built in Logistic Regression

```
In [32]: logreg = LogisticRegression(random_state=1)

#fit the model
logreg.fit(x_train,y_train)
```

Out[32]: LogisticRegression(random_state=1)

Make predictions

```
In [33]: predictions_train = logreg.predict(x_train)
y_pred = logreg.predict(x_test)
t_pred = logreg.predict(x_test)
```

```
In [34]: # I had to do this because confusion matrix was throwing errors
print(predictions_train)
print(y_pred)
print(y_pred.shape)
print(y_pred.dtype)
print(y_test.shape)
print(y_test.dtype)

[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
(45000,)
int64
(45000,)
int64
```

```
In [35]: y_predi = y_pred.astype(int) #converting float to int
```

create confuson matrices in text view

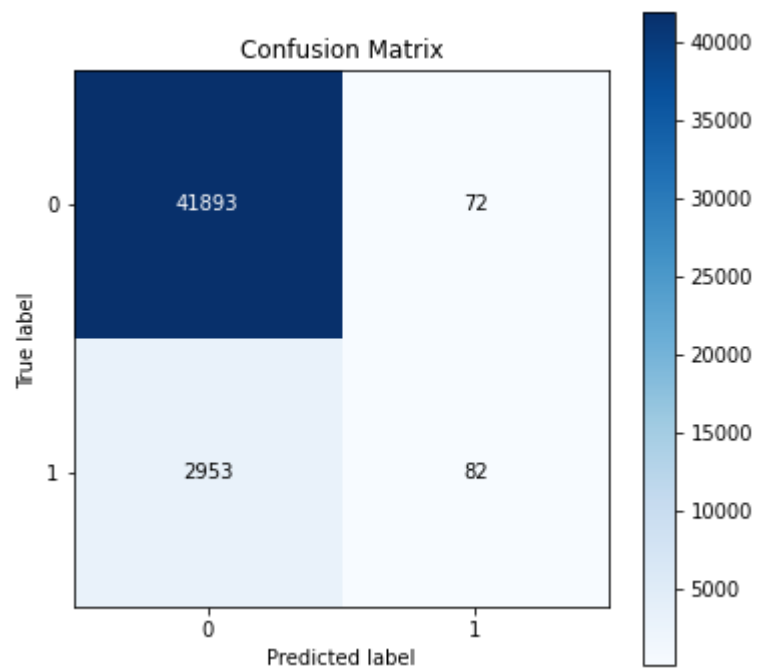
```
In [36]: tn, fp, fn, tp = confusion_matrix( y_test,y_predi).ravel()
(tn, fp, fn, tp)
```

```
Out[36]: (41893, 72, 2953, 82)
```

consusion matrix

```
In [37]: skplt.metrics.plot_confusion_matrix(y_test,y_predi,figsize=(6,6))
```

```
Out[37]: <AxesSubplot:title={'center':'Confusion Matrix'}, xlabel='Predicted label', ylabel='True label'>
```



calculating ROC

```
In [38]: fpr, tpr, thresholds = metrics.roc_curve(y_test, y_predi)
print(fpr)
print(fpr.shape)
print(tpr)
print(tpr.shape)
print(thresholds)
```

```
[0.          0.00171572 1.          ]
(3,)
[0.          0.02701812 1.          ]
(3,)
[2 1 0]
```

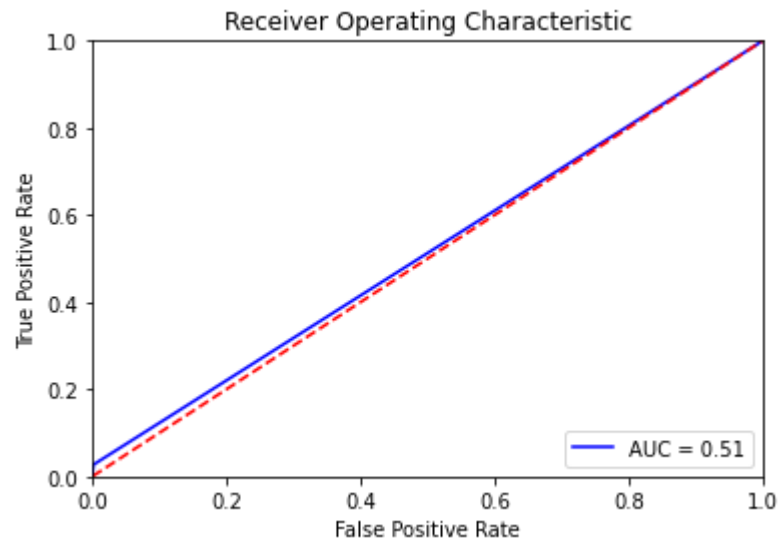
calculating AUC

```
In [39]: from sklearn import metrics
# AUC
roc_auc = metrics.auc(fpr, tpr)
print(roc_auc)
```

```
0.5126512032169273
```

Plotting ROC

```
In [40]: plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



Calculate F1 score

```
In [41]: #calculate F1 score
from sklearn.metrics import f1_score
f1_score(y_test, y_predi)
```

Out[41]: 0.051426779554719346

Accuracy score

```
In [42]: accuracy_score(y_test, y_predi)
```

Out[42]: 0.9327777777777778

```
In [43]: form = []
form.append(['Logistic Regression', accuracy_score(y_test, y_predi)])
```

Precision

```
In [44]: precision_score(y_test, y_predi)
```

Out[44]: 0.5324675324675324

Recall

```
In [45]: recall_score(y_test, y_predi)
```

Out[45]: 0.02701812191103789

Cost-sensitive accuracy

```
In [46]: fp_cost = 1
fn_cost = 0
cost_sensitive_accuracy = (tp + tn) / (tp + tn + fp*fp_cost + fn*fn_cost)
print(cost_sensitive_accuracy)
```

0.9982876305087164

```
In [47]: results = pd.DataFrame(['Logistic Regression', accuracy_score(y_test, y_predi) ,precision_score(y_t
columns = ['Model', 'Accuracy', 'Precision', 'Recall', 'F1 Score', 'ROC_AUC'])
results
```

Out[47]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------------|----------|-----------|----------|----------|----------|
| 0 | Logistic Regression | 0.932778 | 0.532468 | 0.027018 | 0.051427 | 0.512651 |

Decision Tree

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems.

```
In [48]: from sklearn.tree import DecisionTreeClassifier

#specify the model, set any numeric valye as parameter to ensure reproducibility
dTree = DecisionTreeClassifier(random_state=1)

#fit the model
dTree.fit(x_train,y_train)
```

Out[48]: DecisionTreeClassifier(random_state=1)

Make predictions

```
In [49]: predictions_train = dTree.predict(x_train)
y_pred = dTree.predict(x_test)
```

```
In [50]: #this section investigates resulting data  
print(predictions_train)  
print(y_pred)  
print(y_pred.shape)  
print(y_pred.dtype)  
print(y_test.shape)  
print(y_test.dtype)
```

```
[0 0 0 ... 1 0 1]  
[0 0 0 ... 0 0 0]  
(45000,)  
int64  
(45000,)  
int64
```

```
In [51]: #conver float to int  
y_predi = y_pred.astype(int)
```

Creating confusion matrices in text view

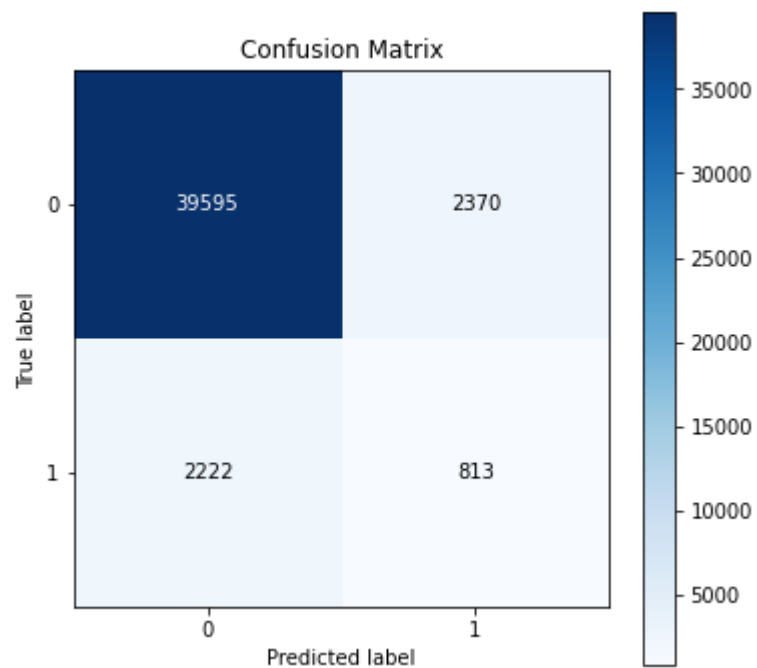
```
In [52]: tn, fp, fn, tp = confusion_matrix( y_test,y_predi).ravel()  
(tn, fp, fn, tp)
```

```
Out[52]: (39595, 2370, 2222, 813)
```

Plotting Confusion Matrix


```
In [53]: skplt.metrics.plot_confusion_matrix(y_test,y_predi,figsize=(6,6))
```

```
Out[53]: <AxesSubplot:title={'center':'Confusion Matrix'}, xlabel='Predicted label', ylabel='True label'>
```



ROC

```
In [54]: fpr, tpr, thresholds = metrics.roc_curve(y_test, y_predi)
print(fpr)
print(fpr.shape)
print(tpr)
print(tpr.shape)
print(thresholds)
```

```
[0.          0.05647563 1.          ]
(3,)
[0.          0.26787479 1.          ]
(3,)
[2 1 0]
```

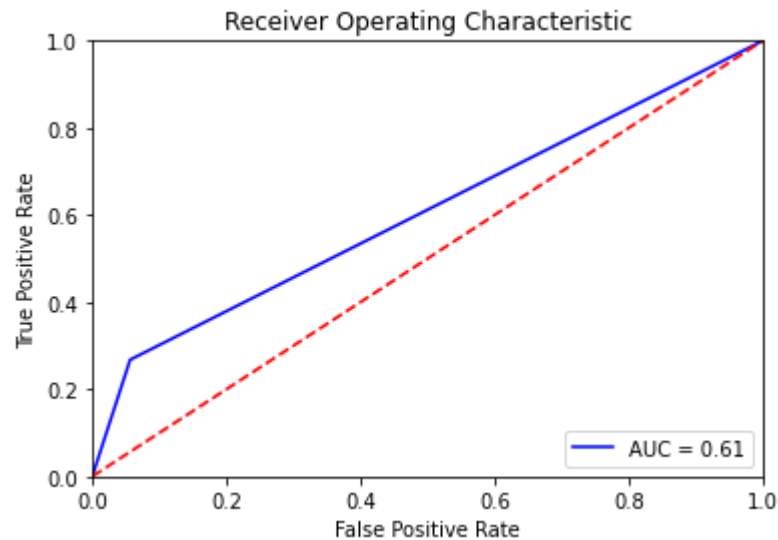
AUC

```
In [55]: roc_auc = metrics.auc(fpr, tpr)
roc_auc
```

```
Out[55]: 0.6056995798059535
```

Plotting ROC

```
In [56]: plt.title('Receiver Operating Characteristic')
plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



F1 score

```
In [57]: f1_score=f1_score(y_test, y_predi)
         f1_score
```

```
Out[57]: 0.26149887423608875
```

Accuracy score

```
In [58]: a_score=accuracy_score(y_test, y_predi)
         a_score
```

```
Out[58]: 0.8979555555555555
```

Precision

```
In [59]: p_score=precision_score(y_test, y_predi)
         p_score
```

```
Out[59]: 0.2554194156456173
```

Recall

```
In [60]: recall_score=recall_score(y_test, y_predi)
         recall_score
```

```
Out[60]: 0.26787479406919273
```

cost-sensitive accuracy

```
In [61]: fp_cost = 1
fn_cost = 0
cost_sensitive_accuracy = (tp + tn) / (tp + tn + fp*fp_cost + fn*fn_cost)
print(cost_sensitive_accuracy)
```

0.9445976904016083

```
In [62]: dt = pd.DataFrame([['Decision Tree', a_score, p_score, recall_score, f1_score, roc_auc]],
columns = ['Model', 'Accuracy', 'Precision', 'Recall', 'F1 Score', 'ROC_AUC'])
dt
```

Out[62]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------|----------|-----------|----------|----------|---------|
| 0 | Decision Tree | 0.897956 | 0.255419 | 0.267875 | 0.261499 | 0.6057 |

```
In [63]: results = results.append(dt, ignore_index = True)
```

```
In [64]: results
```

Out[64]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------------|----------|-----------|----------|----------|----------|
| 0 | Logistic Regression | 0.932778 | 0.532468 | 0.027018 | 0.051427 | 0.512651 |
| 1 | Decision Tree | 0.897956 | 0.255419 | 0.267875 | 0.261499 | 0.605700 |

K-Nearest Neighbor(KNN) Algorithm

K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.

K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

```
In [65]: model_KNN = KNeighborsClassifier()

neighbors = [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20]
model_KNN.fit(x_train, y_train)
knn = KNeighborsClassifier(n_neighbors = 20)
knn.fit(x_train, y_train)
knn_scores_proba = knn.predict_proba(x_train)
knn_scores = knn_scores_proba[:,1]
print("AUC Score :", roc_auc_score(y_train,knn_scores))
```

AUC Score : 0.8886195756715168

```
In [66]: model_KNN.fit(x_train,y_train)
y_pred_knn = logmodel.predict(x_test)
pred_knn = y_pred_knn.astype(int)
```

```
In [67]: scorelist = []
n_neighbors, maxauc = -1, 0
for k in range(100, 1000+1, 100):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(x_train, y_train)
    y_pred = knn.predict_proba(x_test)[:,-1]
    score = roc_auc_score(y_test, y_pred)
    print(k, score)
    if score > maxauc:
        n_neighbors, maxauc = k, score
    scorelist.append(score)
print()
print(n_neighbors, maxauc)
```

```
100 0.7738365716625469
200 0.7751776319444048
300 0.7679608507207013
400 0.7659084657313275
500 0.760544731027327
600 0.7546051182920732
700 0.750963969150569
800 0.7484402963087424
900 0.7465734585834943
1000 0.7445061871006885

200 0.7751776319444048
```

```
In [68]: scorelist = []
n_neighbors, maxauc = -1, 0
for k in range(320, 400+1, 10):
    knn = KNeighborsClassifier(n_neighbors=k)
    knn.fit(x_train, y_train)
    y_pred = knn.predict_proba(x_test)[: ,1]
    score = roc_auc_score(y_test, y_pred)
    print(k, score)
    if score > maxauc:
        n_neighbors, maxauc = k, score
    scorelist.append(score)
print()
print(n_neighbors, maxauc)
```

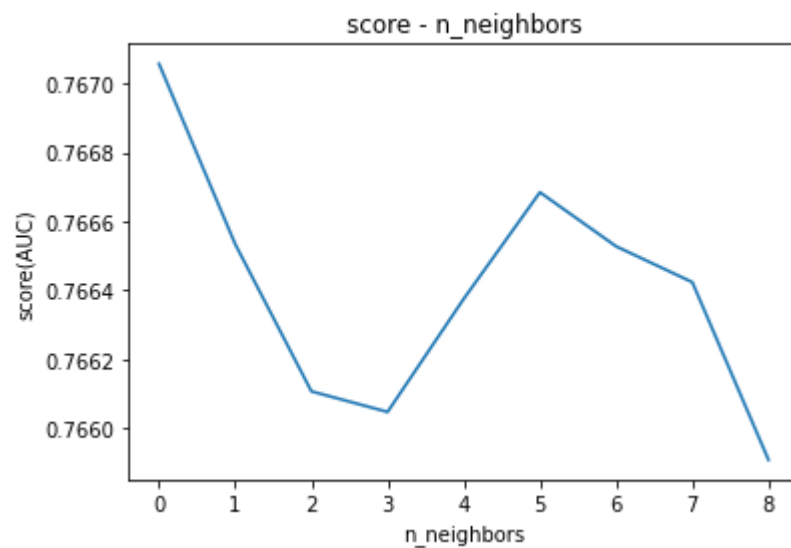
```
320 0.7670568377860976
330 0.7665339104466713
340 0.7661070700833106
350 0.766047414188218
360 0.7663753449518907
370 0.7666842632451809
380 0.7665274996756338
390 0.7664235965053643
400 0.7659084657313275
```

```
320 0.7670568377860976
```



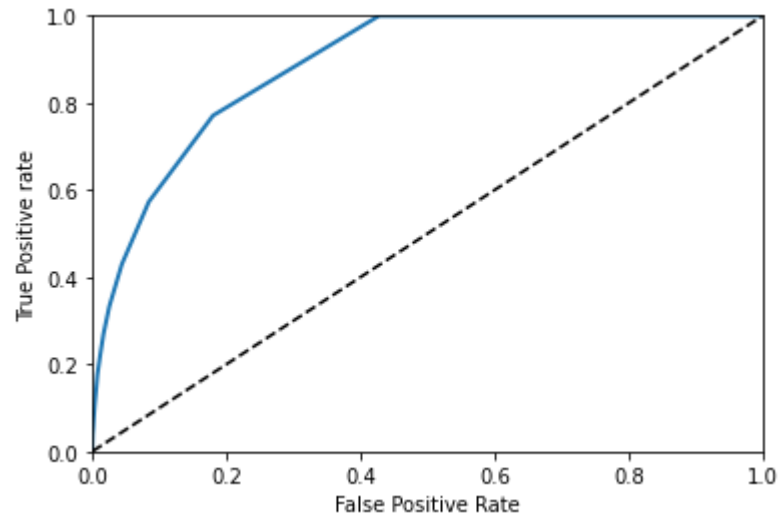
```
In [69]: plt.plot(scorelist)
plt.title('score - n_neighbors')
plt.xlabel('n_neighbors')
plt.ylabel('score(AUC)')
```

```
Out[69]: Text(0, 0.5, 'score(AUC)')
```



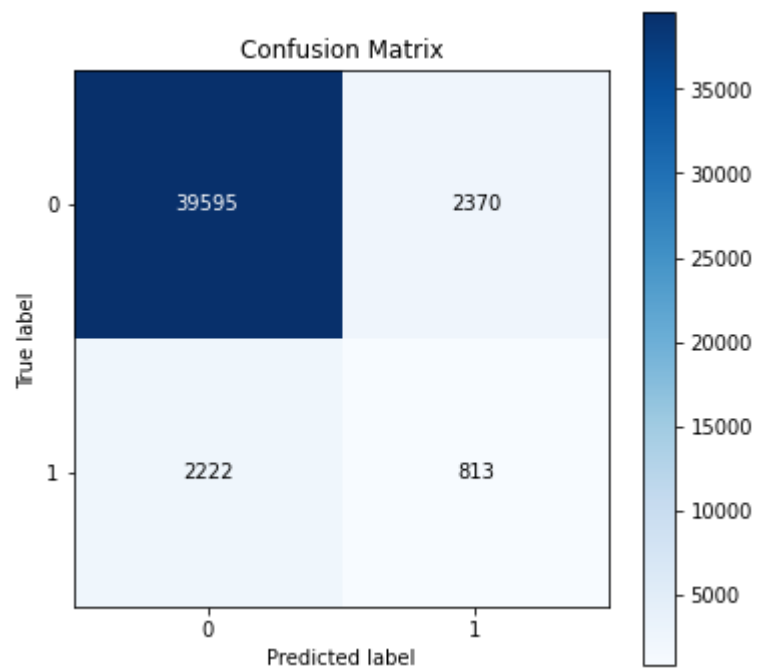
Plotting ROC

```
In [70]: def plot_roc_curve(fpr, tpr, label=None):  
    plt.plot(fpr, tpr, linewidth=2, label=label)  
    plt.plot([0,1],[0,1], "k--")  
    plt.axis([0,1,0,1])  
    plt.xlabel("False Positive Rate")  
    plt.ylabel("True Positive rate")  
fpr_gbc, tpr_gbc, thresh_gbc = roc_curve(y_train, knn_scores)  
plot_roc_curve(fpr_gbc, tpr_gbc)
```



```
In [71]: skplt.metrics.plot_confusion_matrix(y_test,y_predi,figsize=(6,6))
```

```
Out[71]: <AxesSubplot:title={'center':'Confusion Matrix'}, xlabel='Predicted label', ylabel='True label'>
```



ROC

```
In [72]: roc=roc_auc_score(y_test, y_pred)
roc
```

```
Out[72]: 0.7659084657313275
```

AUC

```
In [73]: auc = metrics.auc(fpr, tpr)
auc
```

```
Out[73]: 0.6056995798059535
```

Accuracy

```
In [74]: acc=accuracy_score(y_test, y_pred.round(), normalize=True)
acc
```

```
Out[74]: 0.9331333333333334
```

Precision

```
In [75]: prec = precision_score(y_test, y_pred.round())
prec
```

```
Out[75]: 0.6444444444444445
```

```
In [76]: # f1_scorea=f1_score(pred_knn, pred_knn)
f1_score
```

```
Out[76]: 0.26149887423608875
```

```
In [77]: recall_score
```

```
Out[77]: 0.26787479406919273
```

```
In [78]: k_results = pd.DataFrame([[ 'K-Nearest Neighbour ', acc ,prec,recall_score,f1_score ,roc]],
columns = [ 'Model', 'Accuracy', 'Precision', 'Recall', 'F1 Score','ROC_AUC'])
k_results
```

Out[78]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------------|----------|-----------|----------|----------|----------|
| 0 | K-Nearest Neighbour | 0.933133 | 0.644444 | 0.267875 | 0.261499 | 0.765908 |

```
In [79]: results = results.append(k_results, ignore_index = True)
```

```
In [80]: results
```

Out[80]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------------|----------|-----------|----------|----------|----------|
| 0 | Logistic Regression | 0.932778 | 0.532468 | 0.027018 | 0.051427 | 0.512651 |
| 1 | Decision Tree | 0.897956 | 0.255419 | 0.267875 | 0.261499 | 0.605700 |
| 2 | K-Nearest Neighbour | 0.933133 | 0.644444 | 0.267875 | 0.261499 | 0.765908 |

Gaussian Naive Bayes

Gaussian Naive Bayes is a variant of Naive Bayes that follows Gaussian normal distribution and supports continuous data. We have explored the idea behind Gaussian Naive Bayes along with an example.

Before going into it, we shall go through a brief overview of Naive Bayes.

Naive Bayes are a group of supervised machine learning classification algorithms based on the Bayes theorem. It is a simple classification technique, but has high functionality. They find use when the dimensionality of the inputs is high.

```
In [81]: from sklearn.naive_bayes import GaussianNB

gaussian = GaussianNB()
gaussian.fit(x_train, y_train)
gaussian_scores_proba = gaussian.predict_proba(x_train)
gaussian_scores = gaussian_scores_proba[:,1]
y_pred = gaussian.predict_proba(x_test)[:,1]
y_predi=gaussian.predict(x_test)

gaussian_scores = roc_auc_score(y_test, y_pred)
print(gaussian_scores)
```

0.7097319979719509

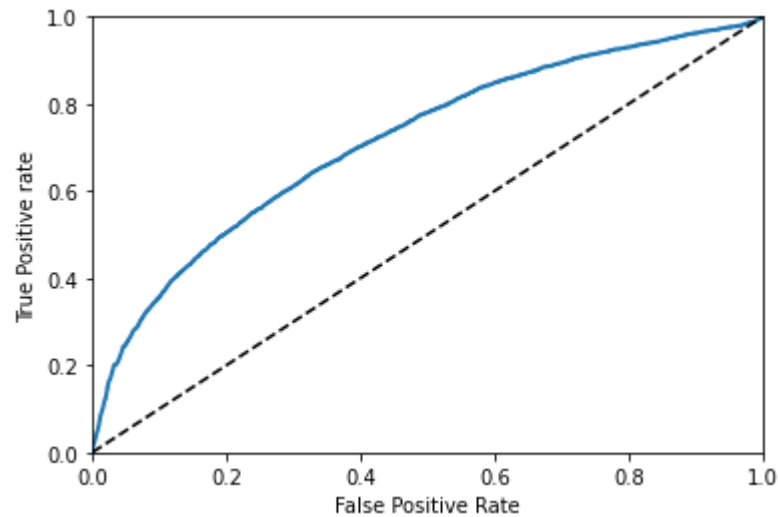
```
In [82]: print(y_test)
```

```
16269      0
140471     0
78683      0
2605       0
81156      0
..
148024     0
59238      0
111773     0
107702     0
89084      0
Name: SeriousDlqin2yrs, Length: 45000, dtype: int64
```

```
In [83]: gaussian = GaussianNB()
gaussian.fit(x_train, y_train)
gaussian_scores_proba = gaussian.predict_proba(x_train)
gaussian_scores = gaussian_scores_proba[:,1]
roc_auc=roc_auc_score(y_train,gaussian_scores)
roc_auc
```

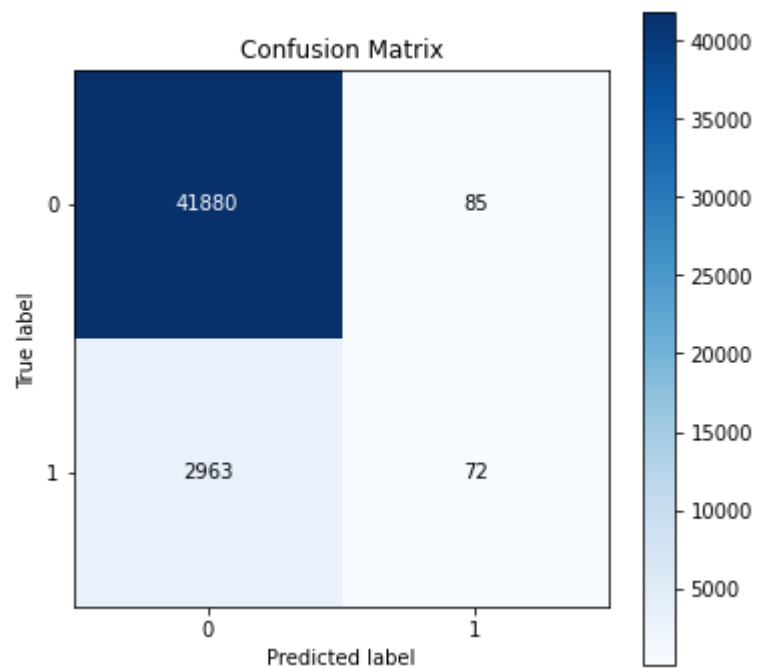
Out[83]: 0.714745677265423

```
In [84]: fpr_gbc, tpr_gbc, thresh_gbc = roc_curve(y_train, gaussian_scores)
plot_roc_curve(fpr_gbc, tpr_gbc)
```



```
In [85]: skplt.metrics.plot_confusion_matrix(y_test,y_predi,figsize=(6,6))
```

```
Out[85]: <AxesSubplot:title={'center':'Confusion Matrix'}, xlabel='Predicted label', ylabel='True label'>
```




```
In [86]: prec = precision_score(y_test, y_predi)
prec
```

```
Out[86]: 0.4585987261146497
```

```
In [87]: acc=accuracy_score(y_test, y_predi.round(), normalize=True)
acc
```

```
Out[87]: 0.9322666666666667
```

```
In [88]: roc=roc_auc_score(y_test, y_predi)
roc
```

```
Out[88]: 0.5108488657783581
```

```
In [89]: f1=f1_score.round(4)
f1
```

```
Out[89]: 0.2615
```

```
In [90]:
rs=recall_score.round(4)
rs
roc_auc.round(3)
```

```
Out[90]: 0.715
```

```
In [91]: gaussian = pd.DataFrame(['Gaussian Naive Bayes', acc ,prec,rs, f1 ,roc_auc]],
columns = ['Model', 'Accuracy', 'Precision', 'Recall', 'F1 Score','ROC_AUC'])
dt
```

```
Out[91]:
```

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|---------------|----------|-----------|----------|----------|---------|
| 0 | Decision Tree | 0.897956 | 0.255419 | 0.267875 | 0.261499 | 0.6057 |

```
In [92]: results = results.append(gaussian, ignore_index = True)
```

In [93]: results

Out[93]:

| | Model | Accuracy | Precision | Recall | F1 Score | ROC_AUC |
|---|----------------------|----------|-----------|----------|----------|----------|
| 0 | Logistic Regression | 0.932778 | 0.532468 | 0.027018 | 0.051427 | 0.512651 |
| 1 | Decision Tree | 0.897956 | 0.255419 | 0.267875 | 0.261499 | 0.605700 |
| 2 | K-Nearest Neighbour | 0.933133 | 0.644444 | 0.267875 | 0.261499 | 0.765908 |
| 3 | Gaussian Naive Bayes | 0.932267 | 0.458599 | 0.267900 | 0.261500 | 0.714746 |

Model Tuning using KFold Validation

Cross-validation is a statistical method used to estimate the skill of machine learning models.

It is commonly used in applied machine learning to compare and select a model for a given predictive modeling problem because it is easy to understand, easy to implement, and results in skill estimates that generally have a lower bias than other methods.

```
In [94]: from sklearn.model_selection import KFold
X, y = make_classification(n_samples=100, n_features=20, n_informative=15, n_redundant=5, random_state=42)
# prepare the cross-validation procedure
cv = KFold(n_splits=10, random_state=1, shuffle=True)
kfold = KFold(n_splits=10, random_state=None)
```

K-fold Cross-Validation on Logistic Regression Model

In [95]: form = []

```

In [96]: # evaluate a logistic regression model using k-fold cross-validation
from numpy import mean
from numpy import std
from sklearn.datasets import make_classification
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
# create dataset
X, y = make_classification(n_samples=100, n_features=20, n_informative=15, n_redundant=5, random_state=1)
# prepare the cross-validation procedure
cv = KFold(n_splits=10, random_state=1, shuffle=True)
# create model
lr_model = LogisticRegression()
# evaluate model
lr_scores = cross_val_score(logreg, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report performance
print('Accuracy: %.3f (%.3f)' % (mean(lr_scores), std(lr_scores)))
LRscore=mean(lr_scores)

```

Accuracy: 0.850 (0.128)

```

In [97]: k_results = pd.DataFrame(['logistic regression Tuned', LRscore ,std(lr_scores)]),
columns = ['Model', 'Accuracy', 'std'])

```

```

In [98]: form.append(['logistic regression Tuned'.format(lr_scores),mean(lr_scores)])
# results1 = results1.append(k_results, ignore_index = True)
k_results

```

Out[98]:

| | Model | Accuracy | std |
|---|---------------------------|----------|----------|
| 0 | logistic regression Tuned | 0.85 | 0.128452 |

K-fold Cross-Validation on Decision Tree

In [99]:

```

a, b = make_classification(n_samples=100, n_features=20, n_informative=15, n_redundant=5, random_state=42)
# prepare the cross-validation procedure
cv = KFold(n_splits=10, random_state=1, shuffle=True)
# create model
Dt_model = DecisionTreeClassifier()
# evaluate model
scores = cross_val_score(Dt_model, a, b, scoring='accuracy', cv=cv, n_jobs=-1)
# report performance
print('Accuracy: %.3f (%.3f)' % (mean(scores), std(scores)))
DTCscore=mean(scores)

```

Accuracy: 0.650 (0.143)

In [100]:

```

form.append(['Decision Tree Tuned'].format(scores),mean(scores))
# results

```

In [101]:

```

k_model = pd.DataFrame(['Decision Tree Tuned', DTCscore ,std(scores)]),
columns = ['Model', 'Accuracy', 'std'])
k_results = k_results.append(k_model, ignore_index = True)

```

In [102]:

k_results

Out[102]:

| | Model | Accuracy | std |
|---|---------------------------|----------|----------|
| 0 | logistic regression Tuned | 0.85 | 0.128452 |
| 1 | Decision Tree Tuned | 0.65 | 0.143178 |

K-fold Cross-Validation on KNN

```
In [103]: knn_model = KNeighborsClassifier()
# evaluate model
knn_scores = cross_val_score(model_KNN, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report performance
print('Accuracy: %.3f (%.3f)' % (mean(knn_scores), std(knn_scores)))
KNNscore=mean(knn_scores)

Accuracy: 0.770 (0.110)
```

```
In [104]: k_model = pd.DataFrame(['KNN Tuned', KNNscore, std(knn_scores)],
columns = ['Model', 'Accuracy', 'std'])
k_results = k_results.append(k_model, ignore_index = True)
```

```
In [105]: form.append(['KNN Tuned'.format(knn_scores), mean(knn_scores)])
```

K-fold Cross-Validation Gaussian Naive Bayes

```
In [106]: gnb_model = GaussianNB()
# evaluate model
gnb_scores = cross_val_score(gnb_model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
# report performance
print('Accuracy: %.3f (%.3f)' % (mean(gnb_scores), std(gnb_scores)))
form.append(['Gaussian Naive Bayes Tuned'.format(gnb_scores), mean(gnb_scores)])
NBscore=mean(gnb_scores)

Accuracy: 0.720 (0.117)
```

```
In [107]: k_model = pd.DataFrame(['Gaussian Naive Bayes', NBscore, mean(gnb_scores)],
columns = ['Model', 'Accuracy', 'std'])
k_results = k_results.append(k_model, ignore_index = True)
```

In [108]: k_results

Out[108]:

| | Model | Accuracy | std |
|---|---------------------------|----------|----------|
| 0 | logistic regression Tuned | 0.85 | 0.128452 |
| 1 | Decision Tree Tuned | 0.65 | 0.143178 |
| 2 | KNN Tuned | 0.77 | 0.110000 |
| 3 | Gaussian Naive Bayes | 0.72 | 0.720000 |

ROC Curve Comparison of Classifiers

```
In [117]: logreg.fit(x_train, y_train)
dTree.fit(x_train, y_train)
knn_model.fit(x_train, y_train)
# gaussian.fit(x_train, y_train)
log_pred = logreg.predict(x_test)
dt_pred = dTree.predict(x_test)
knn_pred = knn_model.predict(x_test)
```

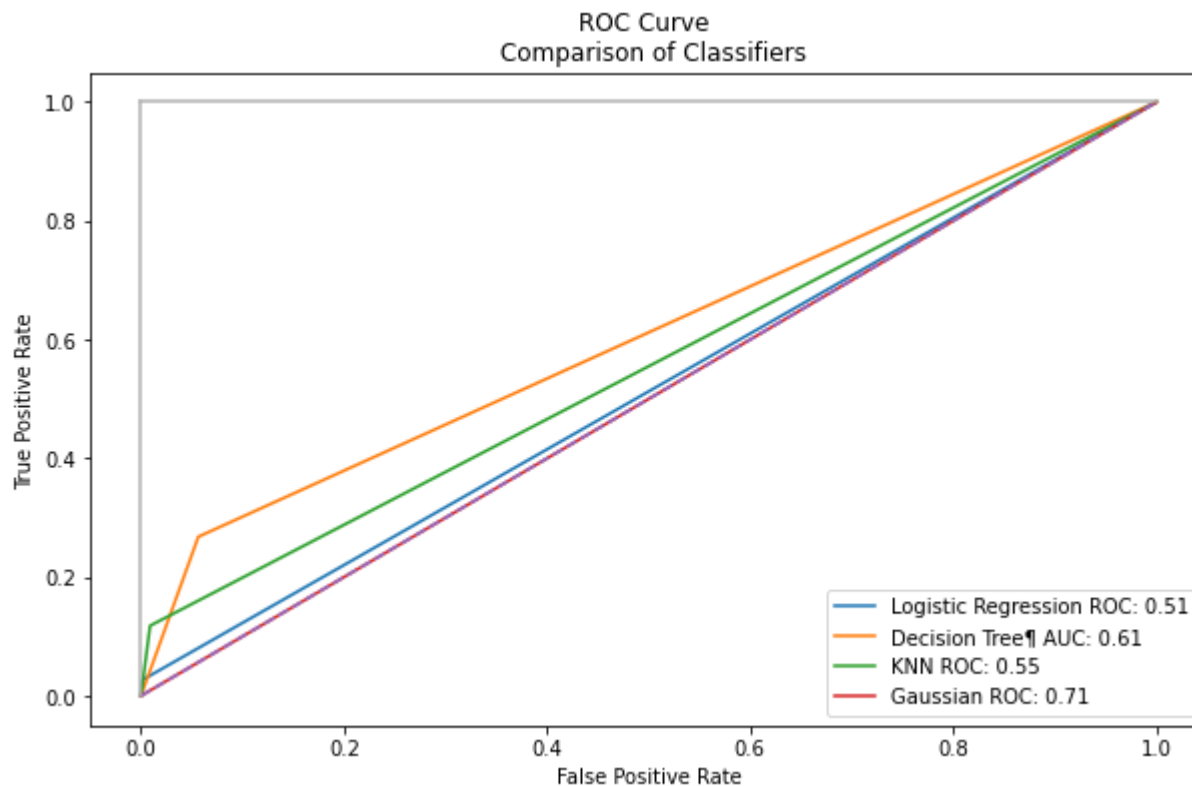
```
In [118]: log_fpr, log_tpr, log_threshold = roc_curve(y_test, log_pred)
dt_fpr, dt_tpr, dt_threshold = roc_curve(y_test, dt_pred)
knn_fpr, knn_tpr, knn_threshold = roc_curve(y_test, knn_pred)
gn_fpr, gn_tpr, gn_threshold = roc_curve(y_test, gn_pred)
```

```

In [119]: fig = plt.figure(figsize=(10,6))
plt.title('ROC Curve \n Comparison of Classifiers')
plt.plot(log_fpr, log_tpr, label='Logistic Regression ROC: {:.2f}'.format(roc_auc_score(y_test, log
plt.plot(dt_fpr, rfc_tpr, label='Decision Tree AUC: {:.2f}'.format(roc_auc_score(y_test, dt_pred))
plt.plot(knn_fpr, knn_tpr, label='KNN ROC: {:.2f}'.format(roc_auc_score(y_test, knn_pred)))
plt.plot(gn_fpr, gn_tpr, label='Gaussian ROC: {:.2f}'.format(roc_auc_score(y_test, y_pred)))

plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0], c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()

```



Scatter Plot For Algorithm Comparison

```
In [120]: scores=[KNNscore,NBscore,LRscore,DTCscore,]
AlgorithmsName=["K-NN","Naive Bayes","Logistic Regression","Decision Tree"]

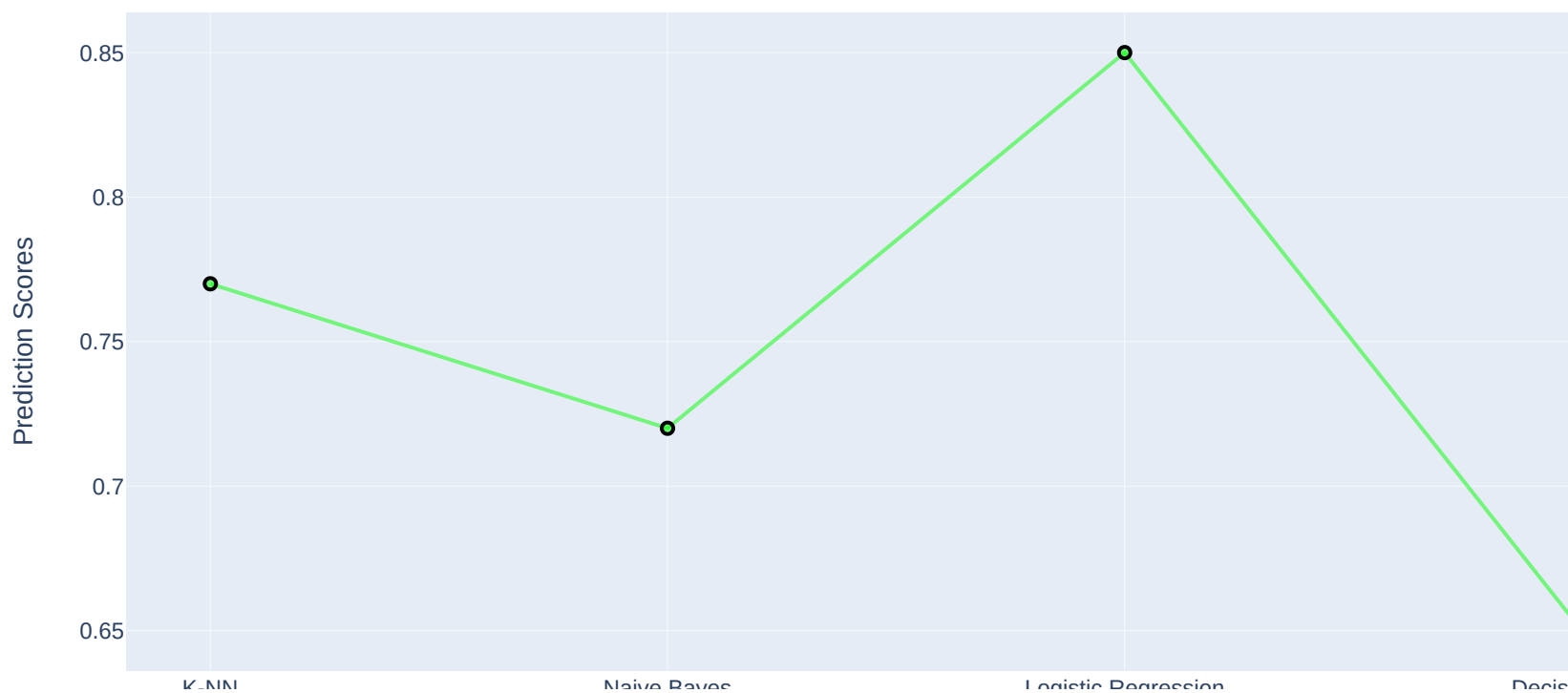
#create traces

tracel = go.Scatter(
    x = AlgorithmsName,
    y= scores,
    name='Algorithms Name',
    marker =dict(color='rgba(0,255,0,0.5)',
                line =dict(color='rgb(0,0,0)',width=2)),
    text=AlgorithmsName
)
data = [tracel]

layout = go.Layout(barmode = "group",
                    xaxis= dict(title= 'ML Algorithms',ticklen= 5,zeroline= False),
                    yaxis= dict(title= 'Prediction Scores',ticklen= 5,zeroline= False))
fig = go.Figure(data = data, layout = layout)
```



```
In [122]: ipplot(fig)
```



Conclusion

We found best result with **Logistic Regression**.

Logistic Regression has the best performance with **0.85 Accuracy** compared to other three classifiers KNN, Gaussian Naive Bayes, and decision tree.

