

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
import sqlite3
import pandas
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import GridSearchCV
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import Normalizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer
import re
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle
from tqdm import tqdm
import os
from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
from sklearn.naive_bayes import MultinomialNB
from prettytable import PrettyTable
```

```
import csv
import math as m
from sklearn.preprocessing import Normalizer
```



```
data = pandas.read_csv('/content/drive/MyDrive/preprocessed_data (1).csv')
```

```
!pip install plotly==3.10.0
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Collecting plotly==3.10.0

Downloading plotly-3.10.0-py2.py3-none-any.whl (41.5 MB)

41.5 MB 58.5 MB/s

Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from plotly==3.10.0) (2.23.0)

Collecting retrying>=1.3.3

Downloading retrying-1.3.3.tar.gz (10 kB)

Requirement already satisfied: decorator>=4.0.6 in /usr/local/lib/python3.7/dist-packages (from plotly==3.10.0) (4.4.2)

Requirement already satisfied: pytz in /usr/local/lib/python3.7/dist-packages (from plotly==3.10.0) (2022.4)

Requirement already satisfied: nbformat>=4.2 in /usr/local/lib/python3.7/dist-packages (from plotly==3.10.0) (5.7.0)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from plotly==3.10.0) (1.15.0)

Requirement already satisfied: traitlets>=5.1 in /usr/local/lib/python3.7/dist-packages (from nbformat>=4.2->plotly==3.10.0) (5

Requirement already satisfied: jsonschema>=2.6 in /usr/local/lib/python3.7/dist-packages (from nbformat>=4.2->plotly==3.10.0) (4.1.0)

Requirement already satisfied: fastjsonschema in /usr/local/lib/python3.7/dist-packages (from nbformat>=4.2->plotly==3.10.0) (2

Requirement already satisfied: importlib-metadata>=3.6 in /usr/local/lib/python3.7/dist-packages (from nbformat>=4.2->plotly==3

```
Requirement already satisfied: setuptools in /usr/local/lib/python2.7/dist-packages (from ipformats==4.2.0) 0.4.1
```

```
y = data['project_is_approved'].values
```

```
data.drop(['project_is_approved'], axis='columns', inplace=True)
```

X = data

```
Requirement already satisfied: numpy<1.19.0, >=1.17.0 in /usr/local/lib/python3.7/dist-packages (from isort==4.3.21)
```

```
data.columns
```

```
Index(['school_state', 'teacher_prefix', 'project_grade_category',
      'teacher_number_of_previously_posted_projects', 'clean_categories',
      'clean_subcategories', 'essay', 'price'],
      dtype='object')
```

Stored in directory: /root/.cache/pip/wheels/f0/8d/8d/f6a2f2f7f00a2552bc7f06d52a1b287d3d18b06a861a656dddf

```
# train test split
```

```
from sklearn.model_selection import train_test_split
```

```
# from sklearn importing train_test_split to split the data
```

```
# test size 0.3 mean 70% splitted as train and 30% as test
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, stratify=y)
```

WARNING: The following packages were previously imported in this runtime:

```
vectorizer_bow = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=50000)
```



```
print(X_train_category_ohe.shape, y_train.shape)
# print(X_cv_category_ohe.shape, y_cv.shape)
print(X_test_category_ohe.shape, y_test.shape)
print("=="*100)
```

After vectorizations

(76473, 9) (76473,)

 $(32775, 9) \quad (32775,)$

=====

```
vectorizer_clean_subcategories = CountVectorizer()  
vectorizer_clean_subcategories.fit(X_train['clean_subcategories'].values) # fitting
```

```
X_train_subcategory_oh = vectorizer_clean_subcategories.transform(X_train['clean_subcategories'].values)
#X_cv_subcategory_oh = vectorizer_sub.transform(X_cv['clean_subcategories'].values) #transform
X_test_subcategory_oh = vectorizer_clean_subcategories.transform(X_test['clean_subcategories'].values)
```

```
print("After vectorizations")
print(X_train_subcategory_ohe.shape, y_train.shape)
# print(X_cv_subcategory_ohe.shape, y_cv.shape)
print(X_test_subcategory_ohe.shape, y_test.shape)
print("==>"*100)
```

After vectorizations

 $(76473, 30) \quad (76473,)$ $(32775, 30) \quad (32775,)$

=====

```
vectorizer_project_grade_category = CountVectorizer()  
vectorizer_project_grade_category.fit(X_train['project_grade_category'].values) # fitting
```

```
X_train_grade_ohe = vectorizer_project_grade_category.transform(X_train['project_grade_category'].values)
#X_cv_grade_ohe = vectorizer_grade.transform(X_cv['project_grade_category'].values) #transform
X_test_grade_ohe = vectorizer_project_grade_category.transform(X_test['project_grade_category'].values)
```


SET-1

Set 1: categorical, numerical features + preprocessed_essay (BOW)

```
from scipy.sparse import hstack
X_tr_set_one = hstack((X_train_essay_bow, X_train_state_ohe,
                        X_train_teacher_ohe, X_train_grade_ohe,
                        X_train_price,
                        X_train_category_ohe,
                        X_train_subcategory_ohe,
                        X_train_teacher_number_of_previously_posted_projects)).tocsr()

X_te_set_one = hstack((X_test_essay_bow,
                        X_test_state_ohe,
                        X_test_teacher_ohe,
                        X_test_grade_ohe,
                        X_test_price,
                        X_test_category_ohe,
                        X_test_subcategory_ohe,
                        X_test_teacher_number_of_previously_posted_projects)).tocsr()

print("SHAPE OF TRAIN AND TEST AFTER STACKING")
print(X_tr_set_one.shape)
print(X_te_set_one.shape)
```

```
SHAPE OF TRAIN AND TEST AFTER STACKING
(76473, 50101)
(32775, 50101)
```

Double-click (or enter) to edit

```
from scipy.sparse import hstack
```

```
X_tr_set_two = hstack((X_train_essay_bow,  
                        X_train_state_ohe,  
                        X_train_teacher_ohe,  
                        X_train_grade_ohe,  
                        X_train_price,  
                        X_train_category_ohe,  
                        X_train_subcategory_ohe,  
                        X_train_teacher_number_of_previously_posted_projects )).tocsr()
```

```
X_te_set_two = hstack((X_test_essay_bow,  
                        X_test_state_ohe,  
                        X_test_teacher_ohe,  
                        X_test_grade_ohe,  
                        X_test_price,  
                        X_test_category_ohe,  
                        X_test_subcategory_ohe,  
                        X_test_teacher_number_of_previously_posted_projects)).tocsr()
```

```
print("SHAPE OF TRAIN AND TEST AFTER STACKING")  
print(X_tr_set_two.shape)  
print(X_te_set_two.shape)
```

```
SHAPE OF TRAIN AND TEST AFTER STACKING  
(76473, 50101)  
(32775, 50101)
```

➤ MULTINOMIAL NAIVE BAYES USING GRID SEARCH CROSS VALIDATION (SET - 1)

```
# REFER : http://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes.MultinomialNB.html
```

```
from sklearn.naive_bayes import MultinomialNB  
from sklearn.model_selection import GridSearchCV
```


▸ plotting hyperparameter v/s auc

```
log_param=[]
for i in param['alpha']: # converting alpha into log- alpha
    log_param.append(m.log(i))

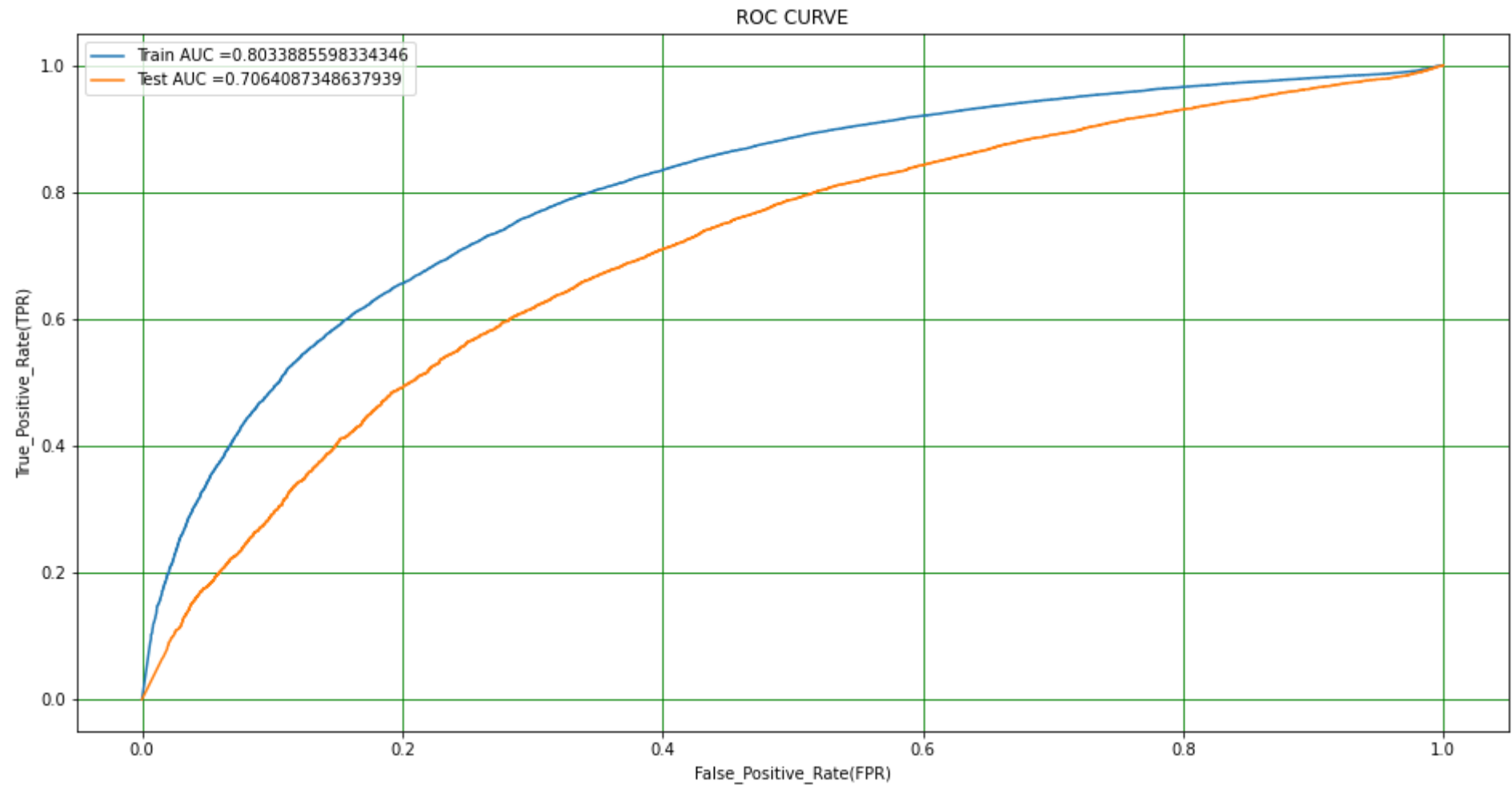
plt.figure(figsize=(12,8))
plt.grid()
plt.plot(log_param, train_auc, label='Train AUC')
plt.plot(log_param, cv_auc, label='CV AUC')
plt.scatter(log_param,train_auc)
plt.scatter(log_param,cv_auc)
plt.legend()
plt.xlabel("log alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC")
plt.grid(color='red', linestyle='-', linewidth=0.5)
```



```
plt.title("ROC CURVE")
```

```
# DEFINING THE GRID PARAMETERS
```

```
plt.grid(color='green',lw=0.8)
```



▼ confusion matrix

```
# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

def find_best_threshold(threshold, fpr, tpr):
    Thres = threshold[np.argmax(tpr*(1-fpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(Thres,3))
    return Thres

# DEFINING THE THRESHOLD VALUE IF VALUE GREATER THAN THRESHOLD THEN 1 AND IF ITS LESS THAN THEN ITS 0
def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

▼ Train data

```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
CONFUSION_MATRIX=metrics.confusion_matrix(y_train,predict_with_best_t(y_train_probs, best_t))

# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

print("CONFUSION MATRIX OF TRAIN DATA")
print("\n")
print('==>'*50)
print(CONFUSION_MATRIX)
sns.heatmap(CONFUSION_MATRIX, annot=True, fmt='d',cmap='RdGy',annot_kws = {"size":16})
```

the maximum value of tpr*(1-fpr) 0.5371868823700227 for threshold 0.577

CONFUSION MATRIX OF TRAIN DATA

Heatmap of Train Data Confusion Matrix. The color scale ranges from 5000 (dark red) to 45000 (black).

	Predicted Label 0	Predicted Label 1
True Label 0	8510	3069
True Label 1	17462	47432

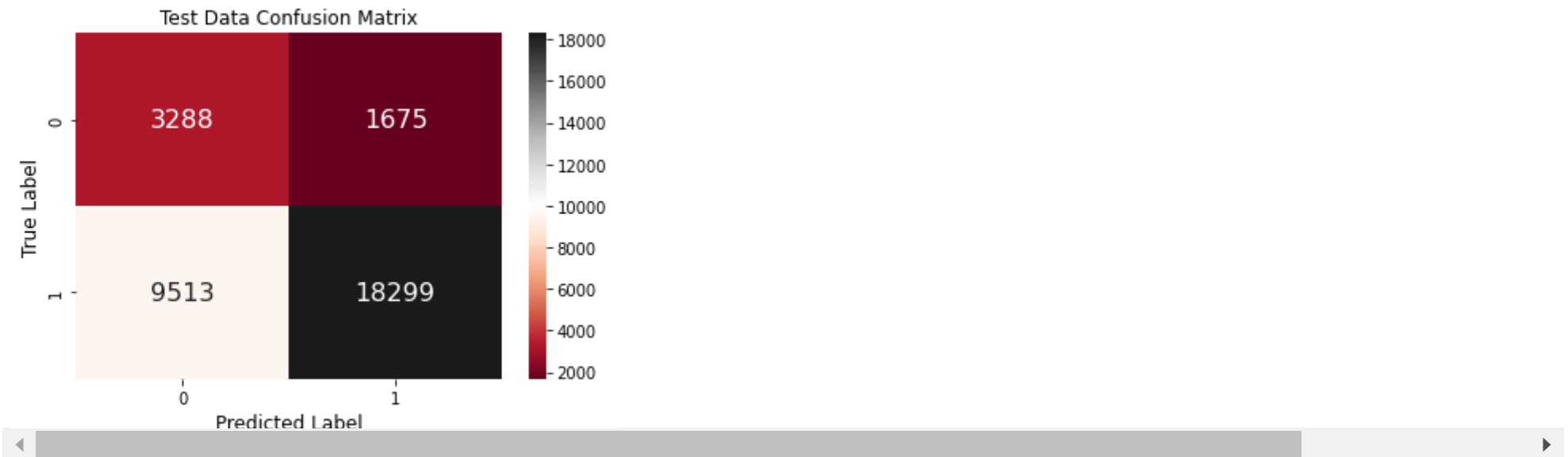
```
# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
best_t = find_best_threshold(te_thresholds, test_fpr, test_tpr)
confusion_matrix = metrics.confusion_matrix(y_test, predict_with_best_t(y_test_probs, best_t))
```



```
print("CONFUSION MATRIX OF TEST DATA")
print('\n')
print('==>'*50)
print(confusion_matrix)
sns.heatmap(confusion_matrix, annot=True, fmt='d', cmap='RdGy', annot_kws = {"size":16})
plt.ylabel('True Label',size=12)
plt.xlabel('Predicted Label',size=12)
plt.title('Test Data Confusion Matrix',size=12)
```

```
the maximum value of tpr*(1-fpr) 0.43589578558015635 for threshold 0.88
CONFUSION MATRIX OF TEST DATA
```

```
[ [ 3288  1675]
  [ 9513 18299]]
Text(0.5, 1.0, 'Test Data Confusion Matrix')
```



▼ MULTINOMIAL NAIVE BAYES USING GRID SEARCH CROSS VALIDATION (SET - 2)

```

from sklearn.naive_bayes import MultinomialNB
from sklearn.model_selection import GridSearchCV

model=MultinomialNB(class_prior=[0.5,0.5])
param={'alpha': [0.00001, 0.0001, 0.001, 0.1, 1, 10, 100,1000]}
clf=GridSearchCV(model,param,scoring='roc_auc',cv=10,return_train_score=True) # running 10 fold cross validation grid search
clf.fit(X_tr_set_two,y_train) #fitting

GridSearchCV(cv=10, estimator=MultinomialNB(class_prior=[0.5, 0.5]),
             param_grid={'alpha': [1e-05, 0.0001, 0.001, 0.1, 1, 10, 100,
                                   1000]}),
             return_train_score=True, scoring='roc_auc')

train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

print('Best score: ',clf.best_score_)
print('alpha value with best score: ',clf.best_params_)
print('='*40)

Best score:  0.7006215706152047
alpha value with best score:  {'alpha': 0.1}
=====

```

▼ plotting hyperparameter v/s auc

```

log_param=[]
for i in param['alpha']: # converting alpha into log- alpha
    log_param.append(m.log(i))

```

```
plt.figure(figsize=(12,8))
plt.grid()
plt.plot(log_param, train_auc, label='Train AUC')
plt.plot(log_param, cv_auc, label='CV AUC')
plt.scatter(log_param,train_auc)
plt.scatter(log_param,cv_auc)
plt.legend()
plt.xlabel("log alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC")
plt.grid(color='red', linestyle='-', linewidth=0.5)
```

alpha: hyperparameter v/s AUC



▼ roc plot of train and test data

```

model_set2=MultinomialNB(alpha=1e-05,class_prior=[0.5,0.5])
model_set2.fit(X_tr_set_two,y_train)

# converting train and test output into probability
y_train_probs = clf.predict_proba(X_tr_set_two)[:,-1] # converting train and test output into probability
y_test_probs= clf.predict_proba(X_te_set_two )[:,-1]

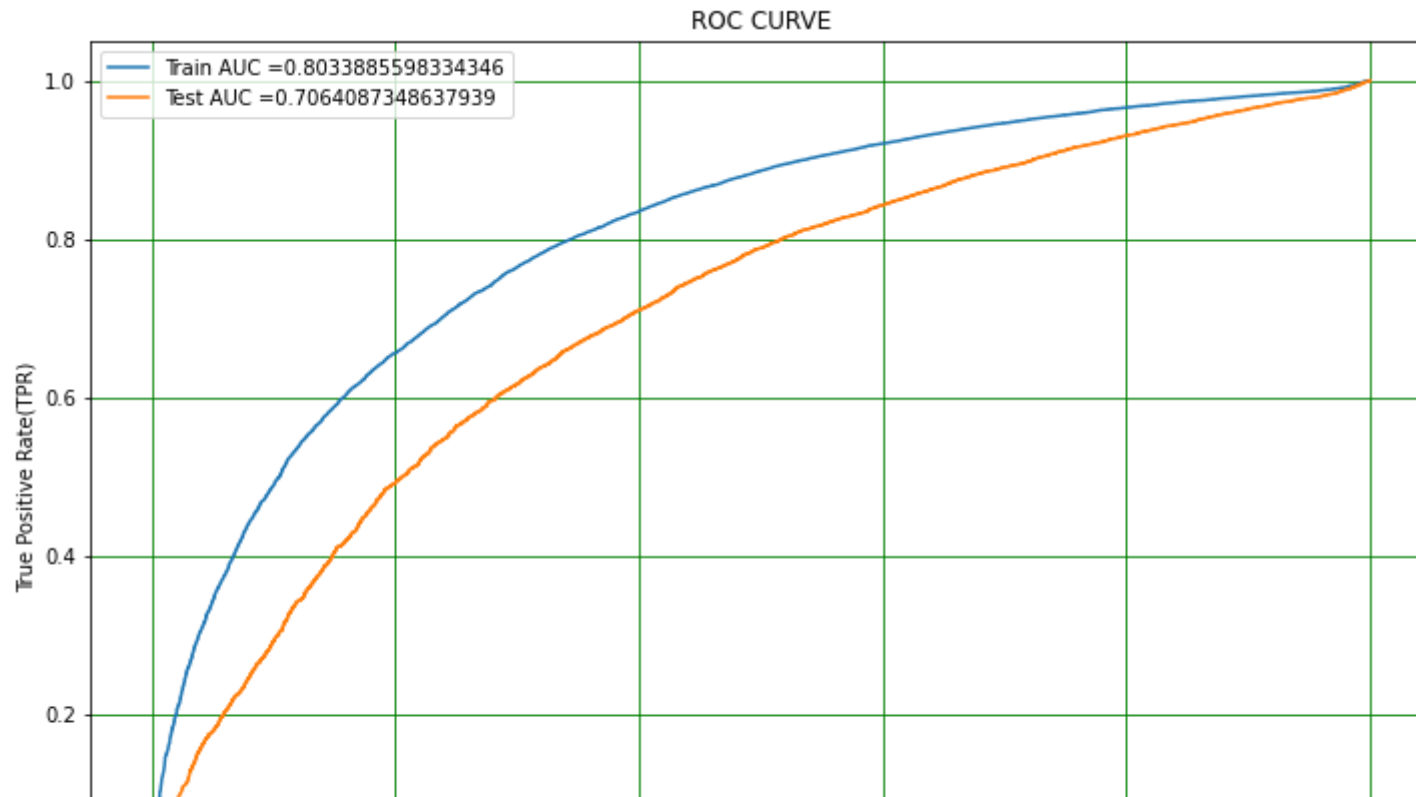
# storing values of fpr and tpr
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_probs) # storing values of fpr and tpr
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_probs)

# PLOTTING THE ROC CURVE
plt.figure(figsize=(12,8))
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()

# NOW WE LABEL THE PLOT X AND Y
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")

# TITLE OF THE CURVE IS ROC CURVE
plt.title("ROC CURVE")
plt.grid(color='green',lw=0.8)

```



▼ Confusion matrix

Train data

```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
CONFUSION_MATRIX=metrics.confusion_matrix(y_train,predict_with_best_t(y_train_probs, best_t))
```

```
# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
# REFER : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
```

```
print("CONFUSION MATRIX OF TRAIN DATA")
print("\n")
```

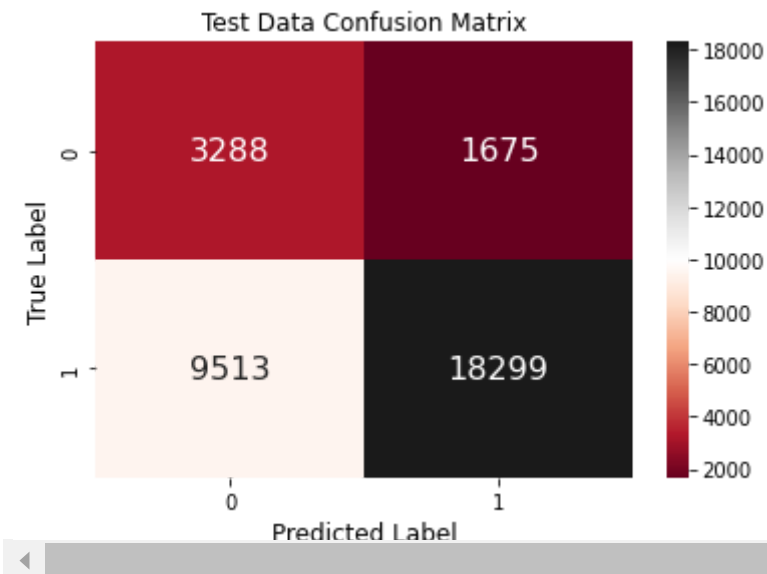


```
print("CONFUSION MATRIX OF TEST DATA")
print('\n')
print('==>'*50)
print(confusion_matrix)
sns.heatmap(confusion_matrix, annot=True, fmt='d', cmap='RdGy', annot_kws = {"size":16})
plt.ylabel('True Label',size=12)
plt.xlabel('Predicted Label',size=12)
plt.title('Test Data Confusion Matrix',size=12)
```

```

=====
[[ 3288  1675]
 [ 9513 18299]]
Text(0.5, 1.0, 'Test Data Confusion Matrix')

```



▼ Top 20 features from set-1

```
from scipy.sparse import hstack
```

```
X_tr_set_one = hstack((X_train_essay_bow,  
                        X_train_state_ohe,  
                        X_train_teacher_ohe,  
                        X_train_grade_ohe,  
                        X_train_price,  
                        X_train_category_ohe,  
                        X_train_subcategory_ohe,  
                        X_train_teacher_number_of_previously_posted_projects )).tocsr()
```

```
X_te_set_one = hstack((X_test_essay_bow,  
                        X_test_state_ohe,  
                        X_test_teacher_ohe,  
                        X_test_grade_ohe,  
                        X_test_price,  
                        X_test_category_ohe,  
                        X_test_subcategory_ohe,  
                        X_test_teacher_number_of_previously_posted_projects)).tocsr()
```

```
print("SHAPE OF TRAIN AND TEST AFTER STACKING")  
print(X_tr_set_one.shape)  
print(X_te_set_one.shape)
```


SHAPE OF TRAIN AND TEST AFTER STACKING

(76473, 50101)

(32775, 50101)

```
features=[]
for fe in vectorizer_bow.get_feature_names() :
    features.append(fe)

for fe in vectorizer_school_state.get_feature_names() :
    features.append(fe)

for fe in vectorizer_teacher_prefix.get_feature_names() : # adding all features into list as the order of data frame
    features.append(fe)

for fe in vectorizer_project_grade_category.get_feature_names() :
    features.append(fe)

features.append("price")

for fe in vectorizer_clean_categories.get_feature_names() :
    features.append(fe)

for fe in vectorizer_clean_subcategories.get_feature_names() :
    features.append(fe)

features.append("teacher_number_of_previously_posted_projects")

class_0=model_set1.feature_log_prob_[0, :].argsort() # finding probability and making argsort for each class
class_1=model_set1.feature_log_prob_[1, :].argsort()
#https://stackoverflow.com/questions/50526898/how-to-get-feature-importance-in-naive-bayes

print("top 20 features of class_0")
print("==>"*100)
print(np.take(features, class_0[-20:])) # since argsort is ascending order
```


+	-----+	-----+	-----+	-----+
	BOW	Naive Bayes	1e-05	0.70648
	TFIDF	Naive Bayes	1e-05	0.70648
+	-----+	-----+	-----+	-----+

▼ optional work

Top 20 features from set-2

```

from scipy.sparse import hstack

X_tr_set_two = hstack((X_train_essay_bow,
                        X_train_state_ohe,
                        X_train_teacher_ohe,
                        X_train_grade_ohe,
                        X_train_price,
                        X_train_category_ohe,
                        X_train_subcategory_ohe,
                        X_train_teacher_number_of_previously_posted_projects)).tocsr()

X_te_set_two = hstack((X_test_essay_bow,
                        X_test_state_ohe,
                        X_test_teacher_ohe,
                        X_test_grade_ohe,
                        X_test_price,
                        X_test_category_ohe,
                        X_test_subcategory_ohe,
                        X_test_teacher_number_of_previously_posted_projects)).tocsr()

print("SHAPE OF TRAIN AND TEST AFTER STACKING")
print(X_tr_set_two.shape)
print(X_te_set_two.shape)

```

```
SHAPE OF TRAIN AND TEST AFTER STACKING
(76473, 50101)
(32775, 50101)
```

```
features=[]
for fe in vectorizer_bow.get_feature_names() :
    features.append(fe)

for fe in vectorizer_school_state.get_feature_names() :
    features.append(fe)

for fe in vectorizer_teacher_prefix.get_feature_names() : # adding all features into list as the order of data frame
    features.append(fe)

for fe in vectorizer_project_grade_category.get_feature_names() :
    features.append(fe)

features.append("price")

for fe in vectorizer_clean_categories.get_feature_names() :
    features.append(fe)

for fe in vectorizer_clean_subcategories.get_feature_names() :
    features.append(fe)

features.append("teacher_number_of_previously_posted_projects")


class_0=model_set2.feature_log_prob_[0, :].argsort() # finding probability and making argsort for each class
class_1=model_set2.feature_log_prob_[1, :].argsort()
#https://stackoverflow.com/questions/50526898/how-to-get-feature-importance-in-naive-bayes


print("top 20 features of class_0")
print("="*80)
print(np.take(features, class_0[-20:])) # since argsort is ascending order
```


◀ [REDACTED] ▶

```
print("top 100 features of class_1")
print("==>"*100)
print(np.take(features, class_1[-100:]))
```

[illegible]

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 9:40 AM

