```python
1  %matplotlib inline
2  import warnings
3  warnings.filterwarnings("ignore")
4
5  import sqlite3
6  import pandas as pd
7  import numpy as np
8  import nltk
9  import string
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 from sklearn.feature_extraction.text import TfidfTransformer
13 from sklearn.feature_extraction.text import TfidfVectorizer
14
15 from sklearn.feature_extraction.text import CountVectorizer
16 from sklearn.metrics import confusion_matrix
17 from sklearn import metrics
18 from sklearn.metrics import roc_curve, auc
19 from nltk.stem.porter import PorterStemmer
20
21 import re
22 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
23 import string
24 from nltk.corpus import stopwords
25 from nltk.stem import PorterStemmer
26 from nltk.stem.wordnet import WordNetLemmatizer
27
28 from gensim.models import Word2Vec
29 from gensim.models import KeyedVectors
30 import pickle
31
32 from tqdm import tqdm
33 import os
34
35 # from plotly import plotly
36 # import plotly.offline as offline
37 # import plotly.graph_objs as go
38 # offline.init_notebook_mode()
39 from collections import Counter
40 print('done all modules imported ')
```

➪　done all modules imported

```python
1  #getting the file from google drive (test data)
2  import gdown
3
4  url = 'https://drive.google.com/uc?id=1bDLwb_Vq7q2W9S89JB96PgmZG3LsLns9'
5  output = 'train.csv'
```

```
6  # https://drive.google.com/file/d/1bDLwb_Vq7q2W9S89JB96PgmZG3LsLns9/view?usp=sharing
7  gdown.download(url, output, quiet=False)
```

```
Downloading...
From: https://drive.google.com/uc?id=1bDLwb_Vq7q2W9S89JB96PgmZG3LsLns9
To: /content/train.csv
201MB [00:01, 109MB/s]
'train.csv'
```

```
1  #getting the data from google drive (resources data)resources data
2  import gdown
3  url = 'https://drive.google.com/uc?id=14OVXWu_SJU-lJD-jKMOCld14EZ21lYYe'
4  output = 'resources.csv'
5  gdown.download(url, output, quiet=False)
6
```

```
Downloading...
From: https://drive.google.com/uc?id=14OVXWu_SJU-lJD-jKMOCld14EZ21lYYe
To: /content/resources.csv
127MB [00:00, 136MB/s]
'resources.csv'
```

```
1  dft = pd.read_csv('train.csv',nrows=50000)
2  dfr = pd.read_csv('resources.csv')
```

```
1  print("Number of data points in train data", dft.shape)
2  print('^^'*50)
3  print("The attributes of data :", dft.columns.values)
4  print('^^'*50)
5  print(dfr.shape)
6  print(dfr.columns.values)
```

```
Number of data points in train data (50000, 17)
^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
(1541272, 4)
['id' 'description' 'quantity' 'price']
```

```python
1  #sort the datapoints by date and time column
2  # list comprehension python : https://stackoverflow.com/a/2582163/4084039
3  # cols = ['Date' if x=='project_submitted_datetime' else x for x in list(dft.columns)]
4  # #sort dataframe based on time  uisng pandas to_datetime function  : https://stackoverflow.com/a/49702492/4084039
5  # dft['Date'] = pd.to_datetime(dft['project_submitted_datetime'])
6  # dft.drop('project_submitted_datetime', axis=1, inplace=True)# we drop the col
7  # dft.sort_values(by=['Date'], inplace=True)# sort the values y date
8  dft.head(2)
```

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | projec |
|---|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | |

## ▾ 1.1 Text preprocessing

```
1  # merge two column text dataframe:
2  dft["essay"] = dft["project_essay_1"].map(str) +\
3                 dft["project_essay_2"].map(str) + \
4                 dft["project_essay_3"].map(str) + \
5                 dft["project_essay_4"].map(str)
```

```
1  # https://stackoverflow.com/a/47091490/4084039
2  import re
3
4  def decontracted(phrase):
5      # specific
6      phrase = re.sub(r"won't", "will not", phrase)
7      phrase = re.sub(r"can\'t", "can not", phrase)
8
9      # general
10     phrase = re.sub(r"n\'t", " not", phrase)
11     phrase = re.sub(r"\'re", " are", phrase)
12     phrase = re.sub(r"\'s", " is", phrase)
13     phrase = re.sub(r"\'d", " would", phrase)
14     phrase = re.sub(r"\'ll", " will", phrase)
15     phrase = re.sub(r"\'t", " not", phrase)
16     phrase = re.sub(r"\'ve", " have", phrase)
17     phrase = re.sub(r"\'m", " am", phrase)
18     return phrase
```

```
1  # https://gist.github.com/sebleier/554280
2  # we are removing the words from the stop words list so as to get btter prediction : that is , no , not ,etc .
3  stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've",\
4              "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
5              'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their',\
6              'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
7              'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
8              'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
9              'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after',\
10             'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further',\
11             'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more',\
12             'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
13             's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
14             've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn',\
15             "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn',\
16             "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
17             'won', "won't", 'wouldn', "wouldn't"]
```

▼  **Preprocessing of the \*\*project_subject_categories\*\***

```python
 1 categories = list(dft['project_subject_categories'].values)
 2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
 3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
 4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
 5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
 6 cat_list = []
 7 for i in categories:
 8     temp = ""
 9     for j in i.split(','):
10         if 'The' in j.split():
11             j=j.replace('The','')
12         j = j.replace(' ','')
13         temp+=j.strip()+" "
14         temp = temp.replace('&','_')
15     cat_list.append(temp.strip())
16
17 dft['clean_categories'] = cat_list
18 dft.drop(['project_subject_categories'], axis=1, inplace=True)
19
20 from collections import Counter
21 my_counter = Counter()
22 for word in dft['clean_categories'].values:
23     my_counter.update(word.split())
24
25 cat_dict = dict(my_counter)
26 project_subject_categories_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
27 print(project_subject_categories_dict)
```

```
{'Warmth': 643, 'Care_Hunger': 643, 'History_Civics': 2689, 'Music_Arts': 4699, 'AppliedLearning': 5569, 'SpecialNeeds'
```

▼ **preprocessing of the **project_subject_subcategories****

```python
 1 sub_catogories = list(dft['project_subject_subcategories'].values)
 2 # # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
 3 # # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
 4 # # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
 5 # # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
 6
 7 sub_cat_list = []
 8 for i in sub_catogories:
 9     temp = ""
10     for j in i.split(','):
11         if 'The' in j.split():
12             j=j.replace('The','')
13         j = j.replace(' ','')
14         temp +=j.strip()+" "
15         temp = temp.replace('&','_')
```

```
16        sub_cat_list.append(temp.strip())
17
18 dft['clean_subcategories'] = sub_cat_list
19 dft.drop(['project_subject_subcategories'], axis=1, inplace=True)
20
21 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
22 my_counter = Counter()
23 for word in dft['clean_subcategories'].values:
24     my_counter.update(word.split())
25
26 sub_cat_dict = dict(my_counter)
27 project_subject_subcategories_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
28 project_subject_subcategories_dict
```

⊐→

{'AppliedSciences': 4901.

**preprocessing of the \*\*project_grade_category\*\***

civics_government . 586,

```
1  Grade= list(dft['project_grade_category'].values)
2  # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
4  # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5  # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6  # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7
8  grade_cat_list = []
9  for i in Grade:
10     # consider we have text like this:
11     for j in i.split(' '): #     # split by spae
12         j=j.replace('Grades','')# clean grades from the row
13     grade_cat_list.append(j.strip())
14
15
16
17 dft['clean_grade'] = grade_cat_list
18 dft.drop(['project_grade_category'], axis=1, inplace=True)
19
20 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
21 my_counter = Counter()
22 for word in dft['clean_grade'].values:
23     my_counter.update(word.split())
24
25 project_grade_category_dict= dict(my_counter)
26 sorted_project_grade_category_dict = dict(sorted(project_grade_category_dict.items(), key=lambda kv: kv[1]))
27 sorted_project_grade_category_dict
```

> {'3-5': 16968, '6-8': 7750, '9-12': 4966, 'PreK-2': 20316}

# Preparing data matrix for the models for Model

▶ **Spliting the Data for test , train and cv using the SKlearn Module train_test_split**

> ↳ *3 cells hidden*

▶ **Pre-processing the Text Features ,**

here we are independtenly doing the preprocessing as i observed that if we pre-processes and then divide the data , it is casuing in lower AUC score

⌄ 6 cells hidden

## ▾ 2.1 vectorizing the Categorical Data

The vectorizing Processes used here is One-hot Encoding

```python
from collections import Counter
my_counter = Counter()
for word in X_train['clean_categories'].values:
    my_counter.update(word.split())

cat_dict = dict(my_counter)
project_subject_categories_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
print(project_subject_categories_dict)

#  vectorizing Processes of clean_categories
from sklearn.feature_extraction.text import CountVectorizer
vectorizer1 = CountVectorizer(vocabulary=list(project_subject_categories_dict.keys()), lowercase=False, binary=True)
vectorizer1.fit(X_train['clean_categories'].values)

# The fit is only to be happen on the train data orelse it may casue data leakage
# we use the fitted CountVectorizer on X_train data for converting the text to vector form (i.e using One- hot encoding )
X_train_cat = vectorizer1.transform(X_train['clean_categories'].values)
X_cv_cat = vectorizer1.transform(X_cv['clean_categories'].values)
X_test_cat = vectorizer1.transform(X_test['clean_categories'].values)

print(vectorizer1.get_feature_names())
```

```
{'Warmth': 289, 'Care_Hunger': 289, 'History_Civics': 1246, 'Music_Arts': 2103, 'AppliedLearning': 2485, 'SpecialNeeds'
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Scie
```

```python
print("After vectorizations the shape of the data changes : it would be a spare matix form")
print(X_train_cat.shape, y_train.shape)
print("="*100)
print(X_cv_cat.shape, y_cv.shape)
print("="*100)
print(X_test_cat.shape, y_test.shape)
```

```
After vectorizations the shape of the data changes : it would be a spare matix form
(22445, 9) (22445,)
====================================================================================================
(11055, 9) (11055,)
====================================================================================================
(16500, 9) (16500,)
```

```python
1  # The fit is only to be happen on the train data orelse it may casue data leakage
2  # we use the fitted CountVectorizer on X_train data for converting the text to vector form (i.e using One- hot encoding )
3  from sklearn.feature_extraction.text import CountVectorizer
4  vectorizer2 = CountVectorizer(vocabulary=list(project_subject_subcategories_dict.keys()), lowercase=False, binary=True)
5  vectorizer2.fit(X_train['clean_subcategories'].values)
6
7
8  X_train_subcat = vectorizer2.transform(X_train['clean_subcategories'].values)
9  X_cv_subcat = vectorizer2.transform(X_cv['clean_subcategories'].values)
10 X_test_subcat = vectorizer2.transform(X_test['clean_subcategories'].values)
11
12 print(vectorizer2.get_feature_names())
```

⊢→  ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'Fo

```python
1  print("After vectorizations the shape of the data changes : it would be a spare matix form")
2  print(X_train_subcat.shape, y_train.shape)
3  print("="*100)
4  print(X_cv_subcat.shape, y_cv.shape)
5  print("="*100)
6  print(X_test_subcat.shape, y_test.shape)
7
```

⊢→  After vectorizations the shape of the data changes : it would be a spare matix form
```
(22445, 30) (22445,)
====================================================================================================
(11055, 30) (11055,)
====================================================================================================
(16500, 30) (16500,)
```

```python
1  #school_state data column convert categorical to vectors
2  # 1st the vocabulary of the words are taken into account then it is used to do one - hot - encodng
3  # then sort the dictonary and apply the one - hot - encoding
4  from collections import Counter
```

```
 5 my_counter = Counter()
 6 for word in dft['school_state'].values:
 7     my_counter.update(word.split())
 8
 9 school_state_dict = dict(my_counter)
10 sorted_school_state_dict = dict(sorted(school_state_dict.items(), key=lambda kv: kv[1]))
11
12 # here the data is converted into ector form
13 from sklearn.feature_extraction.text import CountVectorizer
14 vectorizer3 = CountVectorizer(vocabulary=list(sorted_school_state_dict.keys()), lowercase=False, binary=True)
15 vectorizer3.fit(dft['school_state'].values)
16
17
18 X_train_school_state = vectorizer3.transform(X_train['school_state'].values)
19 X_cv_school_state = vectorizer3.transform(X_cv['school_state'].values)
20 X_test_school_state = vectorizer3.transform(X_test['school_state'].values)
21
22 print(vectorizer3.get_feature_names())
```

⇨  ['VT', 'WY', 'ND', 'MT', 'RI', 'NH', 'SD', 'NE', 'AK', 'DE', 'WV', 'ME', 'NM', 'HI', 'DC', 'KS', 'ID', 'IA', 'AR', 'CO'

```
1 print("After vectorizations the shape of the data changes : it would be a spare matix form")
2 print(X_train_school_state .shape, y_train.shape)
3 print(X_cv_school_state .shape, y_cv.shape)
4 print(X_test_school_state .shape, y_test.shape)
5
```

⇨  After vectorizations the shape of the data changes : it would be a spare matix form
   (22445, 51) (22445,)
   (11055, 51) (11055,)
   (16500, 51) (16500,)

```
 1 # here i am converting the clean_grade into vectors
 2 #Fillna is used to fill all he spaces that are present in the data coloumn
 3 #https://stackoverflow.com/questions/42224700/attributeerror-float-object-has-no-attribute-split
 4 dft['clean_grade']=dft['clean_grade'].fillna("")# fill the nulll values with space
 5
 6 ## here i am using the data from the unsplitted data as there as small number of values and may not present in the X_train data
 7 #  so i am fitting it withe original data then converting into vectors
 8 from sklearn.feature_extraction.text import CountVectorizer
 9 vectorizer4 = CountVectorizer(vocabulary=list(sorted_project_grade_category_dict.keys()), lowercase=False, binary=True)
10 vectorizer4.fit(dft['clean_grade'].values)
11
12 # firstly convert fit the train data into the vectoriaer then it learn hte vocablery
```

```
13
14  # we use the fitted CountVectorizer to convert the text to vector
15  X_train_project_grade_category = vectorizer4.transform(X_train['clean_grade'].values)
16  X_cv_project_grade_category = vectorizer4.transform(X_cv['clean_grade'].values)
17  X_test_project_grade_category = vectorizer4.transform(X_test['clean_grade'].values)
18
19  print(vectorizer4.get_feature_names())
```

⤷  ['9-12', '6-8', '3-5', 'PreK-2']

```
1  print("After vectorizations the shape of the data changes : it would be a spare matix form")
2  print(X_train_project_grade_category  .shape, y_train.shape)
3  print(X_cv_project_grade_category  .shape, y_cv.shape)
4  print(X_test_project_grade_category  .shape, y_test.shape)
5
```

⤷  After vectorizations the shape of the data changes : it would be a spare matix form
    (22445, 4) (22445,)
    (11055, 4) (11055,)
    (16500, 4) (16500,)

```
1  ##https://stackoverflow.com/questions/42224700/attributeerror-float-object-has-no-attribute-split
2  dft['teacher_prefix']=dft['teacher_prefix'].fillna(" ")# filll the null values with space
3  my_counter = Counter()
4  for word in dft['teacher_prefix'].values:
5      my_counter.update(word.split())
6
7  # dict sort by value python: https://stackoverflow.com/a/613218/4084039
8  teacher_cat_dict = dict(my_counter)
9  sorted_teacher_prefix_dict = dict(sorted(teacher_cat_dict.items(), key=lambda kv: kv[1]))
10 ## here i am using the data from the unsplitted data as there as small number of values and may not present in the X_train data
11 #  so i am fitting it withe original data then converting into vectors
12 # after counting the frequency of the words we then transform into vectors
13 from sklearn.feature_extraction.text import CountVectorizer
14 vectorizer5 = CountVectorizer(vocabulary=list(sorted_teacher_prefix_dict.keys()), lowercase=False, binary=True)
15 vectorizer5.fit(dft['teacher_prefix'].values.astype('U'))
16
17 #  the fitted CountVectorizer to convert the text to vector
18 X_train_teacher_prefix = vectorizer5.transform(X_train['teacher_prefix'].values.astype('U'))
19 X_cv_teacher_prefix= vectorizer5.transform(X_cv['teacher_prefix'].values.astype('U'))
20 X_test_teacher_prefix = vectorizer5.transform(X_test['teacher_prefix'].values.astype('U'))
21
22 print(vectorizer5.get_feature_names())
23
24
25 # when i executeed this error comes
26 #np.nan is an invalid document, expected byte or unicode string.
```

```
27  # then iconvert to unicode just write .astype('U') after the .values in fit and transform
28  #https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-document
```

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
```

```
1  print("After vectorizations the shape of the data changes : it would be a spare matix form")
2  print(X_train_teacher_prefix.shape, y_train.shape)
3  print(X_cv_teacher_prefix.shape, y_cv.shape)
4  print(X_test_teacher_prefix.shape, y_test.shape)
5  print("="*100)
```

```
After vectorizations the shape of the data changes : it would be a spare matix form
(22445, 5) (22445,)
(11055, 5) (11055,)
(16500, 5) (16500,)
====================================================================================================
```

## 3.1 BOW Vectorization of the data : Processed essay and Project titles

```
1  # i am assigning the values of diffrent preprocessing to avoid confusion
2  # Changing the termonilogy
3  X_train_essay=preprocessed_essays_train
4  X_cv_essay=preprocessed_essays_cv
5  X_test_essay=preprocessed_essays_test
```

```
1   vectorizer6 = CountVectorizer(min_df=10,max_features=5000,ngram_range=(1, 2))
2   # here i am using the DF to be 10 and max features to be 5000 and uni and bi grams , refrence given below
3   # https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html
4   vectorizer6.fit(X_train_essay)
5
6   # we use the fitted CountVectorizer to convert the text to vector
7   X_train_bow = vectorizer6.transform(X_train_essay)
8   X_cv_bow = vectorizer6.transform(X_cv_essay)
9   X_test_bow = vectorizer6.transform(X_test_essay)
10
11  print("After vectorizations the shape of the data changes : it would be a spare matix form")
12  print(X_train_bow.shape, y_train.shape)
13  print(X_cv_bow.shape, y_cv.shape)
14  print(X_test_bow.shape, y_test.shape)
15  print("so the dimension of all are the same and there will be no data leakage")
16  # print(vectorizer6.get_feature_names())
```

```
After vectorizations the shape of the data changes : it would be a spare matix form
(22445, 5000) (22445,)
(11055, 5000) (11055,)
(16500, 5000) (16500,)
so the dimension of all are the same and there will be no data leakage
```

```
1 # i am assigning the values of diffrent preprocessing to avoid confusion
2 # Changing the termonilogy
3 X_train_title=preprocessed_titles_train
4 X_cv_title=preprocessed_titles_cv
5 X_test_title=preprocessed_titles_test
```

```
1 # Project title BOW
2 vectorizer7 = CountVectorizer(min_df=10,max_features=5000,ngram_range=(1, 2))
3 vectorizer7.fit(X_train_title)
4 # here i am using the DF to be 10 and max features to be 5000 and uni and bi grams , refrence given below
5 # https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html
6 X_train_bow_title = vectorizer7.transform(X_train_title)
7 X_cv_bow_title= vectorizer7.transform(X_cv_title)
8 X_test_bow_title = vectorizer7.transform(X_test_title)
9
10 print("After vectorizations the shape of the data changes : it would be a spare matix form")
11 print(X_train_bow_title.shape, y_train.shape)
12 print(X_cv_bow_title.shape, y_cv.shape)
13 print(X_test_bow_title.shape, y_test.shape)
14 print("so the dimension of all are the same and there will be no data leakage")
15
```

```
After vectorizations the shape of the data changes : it would be a spare matix form
(22445, 1619) (22445,)
(11055, 1619) (11055,)
(16500, 1619) (16500,)
so the dimension of all are the same and there will be no data leakage
```

## ▾ 3.2 TFIDF Vectorization of the data : Processed essay and Project titles

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 # We are considering only the words which appeared in at least 10 documents and max of 5000 features .
3 vectorizer8 = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(1, 2))
4 vectorizer8.fit(X_train_title)# that is learned from trainned  data
```

```
 5
 6
 7  # we use the fitted CountVectorizer to convert the text to vector
 8  X_train_tf_title = vectorizer8.transform(X_train_title)
 9  X_cv_tf_title= vectorizer8.transform(X_cv_title)
10  X_test_tf_title = vectorizer8.transform(X_test_title)
11
12  print("After vectorizations the shape of the data changes : it would be a spare matix form")
13  print(X_train_tf_title.shape, y_train.shape)
14  print(X_cv_tf_title.shape, y_cv.shape)
15  print(X_test_tf_title.shape, y_test.shape)
16  print("so the dimension of all are the same and there will be no data leakage")
```

�localhost�localhost    After vectorizations the shape of the data changes : it would be a spare matix form
         (22445, 1619) (22445,)
         (11055, 1619) (11055,)
         (16500, 1619) (16500,)
         so the dimension of all are the same and there will be no data leakage

```
 1  #for essay
 2  from sklearn.feature_extraction.text import TfidfVectorizer
 3  # We are considering only the words which appeared in at least 10 documents(rows or projects).
 4  vectorizer9 = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(1, 2))# its a countvectors used for convert text to vectors
 5  vectorizer9.fit(X_train_essay)# that is learned from trainned  data
 6
 7
 8
 9  # we use the fitted CountVectorizer to convert the text to vector
10  X_train_tf_essay = vectorizer9.transform(X_train_essay)
11  X_cv_tf_essay= vectorizer9.transform(X_cv_essay)
12  X_test_tf_essay = vectorizer9.transform(X_test_essay)
13
14
15
16  print("After vectorizations")
17  print(X_train_tf_essay.shape, y_train.shape)
18  print(X_cv_tf_essay.shape, y_cv.shape)
19  print(X_test_tf_essay.shape, y_test.shape)
20  print("="*100)
21  # so the dimension of alll are the same by using first fit and then transform
22
```

�localhost�localhost

## ▼ 3.3 AVG W2V Vectorization of the data : Processed essay and Project titles

```
1 import gdown
2
3 url = 'https://drive.google.com/uc?id=1MqUasf7jYoPbG35MJ28VQcOjjNp-ZDDp'
4 output = 'glove_vectors'
5 gdown.download(url, output, quiet=False)
```

⬡→  Downloading...
    From: https://drive.google.com/uc?id=1MqUasf7jYoPbG35MJ28VQcOjjNp-ZDDp
    To: /content/glove_vectors
    128MB [00:02, 63.3MB/s]
    'glove_vectors'

```
1 import pickle
2 with open('glove_vectors', 'rb') as f:
3     model = pickle.load(f)
4     glove_words =  set(model.keys())
```

```
 1 def AVG_W2V(values):
 2 # AVG W2V Vectorization.
 3
 4   train_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
 5   for sentence in tqdm(values): # for each review/sentence
 6     vector = np.zeros(300) # as word vectors are of zero length    # we are taking the 300 dimensions  very large
 7     cnt_words =0; # num of words with a valid vector in the sentence/review
 8     for word in sentence.split(): # for each word in a review/sentence
 9         if word in glove_words:
10             vector += model[word]
11             cnt_words += 1
12     if cnt_words != 0:
13         vector /= cnt_words
14     train_avg_w2v_vectors.append(vector)
15
16   print(len(train_avg_w2v_vectors))
17   print(len(train_avg_w2v_vectors[0]))
18   return train_avg_w2v_vectors
19
```

```
1 train_avg_w2v_vectors=AVG_W2V(preprocessed_essays_train)
2 test_avg_w2v_vectors=AVG_W2V(preprocessed_essays_test)
3 cv_avg_w2v_vectors=AVG_W2V(preprocessed_essays_cv)
```

```
100%|████████| 22445/22445 [00:05<00:00, 3863.97it/s]
  2%||          | 363/16500 [00:00<00:04, 3627.11it/s]22445
300
100%|████████| 16500/16500 [00:04<00:00, 3752.85it/s]
  3%|█         | 382/11055 [00:00<00:02, 3810.13it/s]16500
300
100%|████████| 11055/11055 [00:02<00:00, 3764.67it/s]11055
300
```

```
1  #  AVG W2V for preprocessed_titles
2  train_avg_w2v_vectors_title=AVG_W2V(preprocessed_titles_train)
3  test_avg_w2v_vectors_title=AVG_W2V(preprocessed_titles_test)
4  cv_avg_w2v_vectors_title=AVG_W2V(preprocessed_titles_cv)
5
```

```
100%|████████| 22445/22445 [00:00<00:00, 64309.34it/s]
 41%|███       | 6771/16500 [00:00<00:00, 67703.77it/s]22445
300
100%|████████| 16500/16500 [00:00<00:00, 66732.51it/s]
100%|████████| 11055/11055 [00:00<00:00, 65830.20it/s]16500
300
11055
300
```

## 3.4 TFIDF weighted W2V vectorization using the Pretrained Models

```
1  #here i am converting the dictionary with word as a key, and the idf as a value for the vectorization processes
2  tfidf_model = TfidfVectorizer()
3  tfidf_model.fit(preprocessed_essays_train)
4  dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
5  tfidf_words = set(tfidf_model.get_feature_names())
```

```
1  # average Word2Vec
2  # compute average word2vec for each review.
3  def tf_idf_done(word_list):
4
5    train_title_tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
```

```
 6    for sentence in tqdm(word_list): # for each review/sentence
 7      vector = np.zeros(300) # as word vectors are of zero length
 8      tf_idf_weight =0; # num of words with a valid vector in the sentence/review
 9      for word in sentence.split():#.split(): # for each word in a review/sentence
10          if (word in glove_words) and (word in tfidf_words):
11              vec = model[word] # getting the vector for each word
12              # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
13              tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
14              vector += (vec * tf_idf) # calculating tfidf weighted w2v
15              tf_idf_weight += tf_idf
16      if tf_idf_weight != 0:
17          vector /= tf_idf_weight
18      train_title_tfidf_w2v_vectors.append(vector)
19
20    print(len(train_title_tfidf_w2v_vectors))
21    print(len(train_title_tfidf_w2v_vectors[0]))
22    return train_title_tfidf_w2v_vectors
```

```
1 # Vectorization the Processed titles
2 train_title_tfidf_w2v_vectors=tf_idf_done(preprocessed_titles_train)
3 test_title_tfidf_w2v_vectors=tf_idf_done(preprocessed_titles_test)
4 cv_title_tfidf_w2v_vectors=tf_idf_done(preprocessed_titles_cv)
```

```
100%|██████████| 22445/22445 [00:00<00:00, 33230.60it/s]
 23%|██        |        | 3765/16500 [00:00<00:00, 37645.91it/s]22445
300
100%|██████████| 16500/16500 [00:00<00:00, 35520.68it/s]
 33%|███       |        | 3682/11055 [00:00<00:00, 36810.82it/s]16500
300
100%|██████████| 11055/11055 [00:00<00:00, 35964.09it/s]11055
300
```

```
1 # Vectrozation of the preprocessed_essays
2 train_tfidf_w2v_vectors=tf_idf_done(preprocessed_essays_train)
3 test_tfidf_w2v_vectors=tf_idf_done(preprocessed_essays_test)
4 cv_tfidf_w2v_vectors=tf_idf_done(preprocessed_essays_cv)
```

```
100%|████████| 22445/22445 [00:37<00:00, 602.21it/s]
  0%|          |   59/16500 [00:00<00:28, 579.82it/s]22445
300
100%|████████| 16500/16500 [00:27<00:00, 605.52it/s]
  1%|          |   59/11055 [00:00<00:18, 583.36it/s]16500
300
100%|████████| 11055/11055 [00:18<00:00, 597.67it/s]11055
300
```

## 4. 1 Vectorization of the Numerical features

```
1  #  there are 2 numerical features in the Resources data.csv so have to merge th dataframes
2  price_data = dfr.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
3  dft = pd.merge(dft, price_data, on='id', how='left')
4  print(price_data.head(7))
5
6  # same has to be done foe the train and test data
7  X_train = pd.merge(X_train, price_data, on = "id", how = "left")
8  X_test = pd.merge(X_test, price_data, on = "id", how = "left")
9  X_cv = pd.merge(X_cv, price_data, on = "id", how = "left")
```

```
          id     price   quantity
0  p000001    459.56          7
1  p000002    515.89         21
2  p000003    298.97          4
3  p000004   1113.69         98
4  p000005    485.99          8
5  p000006    130.62          5
6  p000007    157.98          6
```

```
1  from sklearn.preprocessing import StandardScaler
2  # https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
3
4  scalar =  StandardScaler()
5  scalar.fit(X_train['price'].values.reshape(-1,1))
6
7  train_price_standar = scalar.transform(X_train['price'].values.reshape(-1, 1))
8  test_price_standar = scalar.transform(X_test['price'].values.reshape(-1, 1))
9  cv_price_standar = scalar.transform(X_cv['price'].values.reshape(-1, 1))
```

```
1  scalar.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
2  train_prev_proj_standar = scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
3  test_prev_proj_standar = scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
4  cv_prev_proj_standar = scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
5
6
```

```
1  scalar.fit(X_train['quantity'].values.reshape(-1,1))
2  train_qnty_standar = scalar.transform(X_train['quantity'].values.reshape(-1, 1))
3  cv_qnty_standar = scalar.transform(X_cv['quantity'].values.reshape(-1, 1))
4  test_qnty_standar = scalar.transform(X_test['quantity'].values.reshape(-1, 1))
5
```

## ▼ *Mergaing of all the Data matrix for diffrent sets of opeations *

1. Set 1: categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW with bi-grams with min_df=10 and max_features=5000`)
2. Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF with bi-grams with min_df=10 and max_features=50(
3. Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
4. Set 4: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

Set 1: categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW with bi-grams with min_df=10 and max_features=5000`)

```
1  from scipy.sparse import hstack
2  # with the same hstack function we are concatinating a sparse matrix and a dense matirx :)
3  X_set1_train = hstack((X_train_bow_title,X_train_bow,X_train_teacher_prefix,X_train_cat,X_train_subcat ,X_train_project_grade_cate
4                     X_train_school_state,train_qnty_standar,train_price_standar,train_prev_proj_standar))
5  # printing the shape of X_set1_train data metrix
6  print(" printing the shape of X_set1_train data metrix",X_set1_train.shape, y_train.shape)
7
8  X_set1_cv = hstack((X_cv_bow_title,X_cv_bow,X_cv_teacher_prefix,X_cv_cat,X_cv_subcat,X_cv_project_grade_category,X_cv_school_state
9                   cv_qnty_standar,cv_price_standar,cv_prev_proj_standar))
10 # printing the shape of X_set1_cv data metrix
11 print("printing the shape of X_set1_cv data metrix",X_set1_cv.shape, y_cv.shape)
12
13 X_set1_test = hstack((X_test_bow_title,X_test_bow,X_test_teacher_prefix,X_test_cat,X_test_subcat,X_test_project_grade_category,X_t
14                    test_qnty_standar,test_price_standar,test_prev_proj_standar))
15
16 # printing the shape of X_set1_test data metrix
17 print("printing the shape of X_set1_test data metrix ",X_set1_test.shape, y_test.shape)
18
```

```
     printing the shape of X_set1_train data metrix (22445, 6721) (22445,)
     printing the shape of X_set1_cv data metrix (11055, 6721) (11055,)
     printing the shape of X_set1_test data metrix  (16500, 6721) (16500,)
```

Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF with bi-grams with min_df=10 and max_features=5000)

```python
1  X_set2_train = hstack((X_train_tf_essay,X_train_tf_title,X_train_teacher_prefix,X_train_cat,X_train_subcat,X_train_project_grade_
2                     train_qnty_standar,train_price_standar,train_prev_proj_standar))
3
4
5  print("printing the shape of X_set2_train data metrix",X_set2_train.shape, y_train.shape)
6
7  print("*"*50)
8  X_set2_cv = hstack((X_cv_tf_essay,X_cv_tf_title,X_cv_teacher_prefix,X_cv_cat,X_cv_subcat,X_cv_project_grade_category,X_cv_school_
9                 cv_qnty_standar,cv_price_standar,cv_prev_proj_standar))
10
11 print("printing the shape of X_set2_cv data metrix",X_set2_cv.shape, y_cv.shape)
12
13 print("*"*50)
14 X_set2_test = hstack((X_test_tf_essay,X_test_tf_title,X_test_teacher_prefix,X_test_cat,X_test_subcat, X_test_project_grade_categor
15                   test_qnty_standar,test_price_standar,test_prev_proj_standar))
16
17 print("printing the shape of X_set2_test data metrix",X_set2_test.shape, y_test.shape)
18
19
```

```
     printing the shape of X_set2_train data metrix (22445, 6721) (22445,)
     **************************************************
     printing the shape of X_set2_cv data metrix (11055, 6721) (11055,)
     **************************************************
     printing the shape of X_set2_test data metrix (16500, 6721) (16500,)
```

Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)

```python
1  X_set3_train = hstack((train_avg_w2v_vectors,train_avg_w2v_vectors_title,train_prev_proj_standar,train_price_standar,train_qnty_st
2                     X_train_teacher_prefix,X_train_cat,X_train_subcat,
3                     X_train_project_grade_category,X_train_school_state))
4
5
6  print("printing the shape of X_set3_train data metrix",X_set3_train.shape, y_train.shape)
7  print("*"*50)
8
```

```
 9 X_set3_cv = hstack((cv_avg_w2v_vectors,cv_avg_w2v_vectors_title,cv_prev_proj_standar,cv_price_standar,cv_qnty_standar,
10                     X_cv_teacher_prefix,X_cv_cat,X_cv_subcat,
11                     X_cv_project_grade_category,X_cv_school_state))
12
13
14 print("printing the shape of X_set3_cv data metrix",X_set3_cv.shape, y_cv.shape)
15 print("*"*50)
16
17 X_set3_test = hstack((test_avg_w2v_vectors,test_avg_w2v_vectors_title,test_prev_proj_standar,test_price_standar,test_qnty_standar
18                       X_test_teacher_prefix,X_test_cat,X_test_subcat,
19                       X_test_project_grade_category,X_test_school_state))
20
21
22 print("printing the shape of X_set3_test data metrix",X_set3_test.shape, y_test.shape)
23
24
```

```
printing the shape of X_set3_train data metrix (22445, 702) (22445,)
****************************************************
printing the shape of X_set3_cv data metrix (11055, 702) (11055,)
****************************************************
printing the shape of X_set3_test data metrix (16500, 702) (16500,)
```

Set 4: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

```
 1 X_set4_train = hstack((train_tfidf_w2v_vectors,train_title_tfidf_w2v_vectors,train_prev_proj_standar,train_price_standar,train_qnt
 2                        X_train_teacher_prefix,X_train_cat,X_train_subcat,
 3                        X_train_project_grade_category,X_train_school_state))
 4
 5
 6 print("printing the shape of X_set4_train data metrix",X_set4_train.shape, y_train.shape)
 7 print("*"*50)
 8
 9 X_set4_cv = hstack((cv_tfidf_w2v_vectors,cv_title_tfidf_w2v_vectors,cv_prev_proj_standar,cv_price_standar,cv_qnty_standar,
10                     X_cv_teacher_prefix,X_cv_cat,X_cv_subcat,
11                     X_cv_project_grade_category,X_cv_school_state))
12
13
14 print("printing the shape of X_set4_CV data metrix",X_set4_cv.shape, y_cv.shape)
15
16 print("*"*50)
17 X_set4_test = hstack((test_title_tfidf_w2v_vectors,test_tfidf_w2v_vectors,test_prev_proj_standar,test_price_standar,test_qnty_star
18                       X_test_teacher_prefix,X_test_cat,X_test_subcat,
19                       X_test_project_grade_category,X_test_school_state))
20
21
22 print("printing the shape of X_set4_test data metrix",X_set4_test.shape, y_test.shape)
```

```
printing the shape of X_set4_train data metrix (22445, 702) (22445,)
****************************************************
printing the shape of X_set4_CV data metrix (11055, 702) (11055,)
****************************************************
printing the shape of X_set4_test data metrix (16500, 702) (16500,)
```
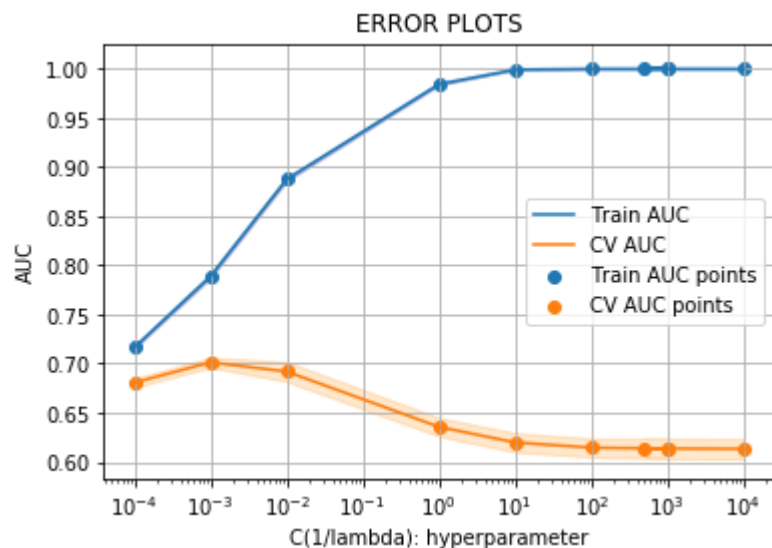
## ▾ Applying the Logistic Regression model on SET:1

```python
1  import warnings
2  warnings.filterwarnings('ignore')
3  from sklearn.metrics import roc_auc_score
4  import matplotlib.pyplot as plt
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import learning_curve, GridSearchCV
7
8
9  clf = LogisticRegression(class_weight='balanced');
10 parameters ={'C':[10**-4, 10**-3,10**-2,1,10,100,1000,500,1000,10000]}
11 sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc',return_train_score=True)
12 sd.fit(X_set1_train, y_train);
13
14
15 train_auc= sd.cv_results_['mean_train_score']
16 train_auc_std= sd.cv_results_['std_train_score']
17 cv_auc = sd.cv_results_['mean_test_score']
18 cv_auc_std= sd.cv_results_['std_test_score']
19
20 plt.plot(parameters['C'], train_auc, label='Train AUC')
21 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
22 plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
23
24 plt.plot(parameters['C'], cv_auc, label='CV AUC')
25 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
26 plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
27
28 plt.scatter(parameters['C'], train_auc, label='Train AUC points')
29 plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
30 plt.xscale('log')
31
32 plt.legend()
33 plt.xlabel("C(1/lambda): hyperparameter")
34 plt.ylabel("AUC")
35 plt.title("ERROR PLOTS")
36 plt.grid()
37 plt.show()
```
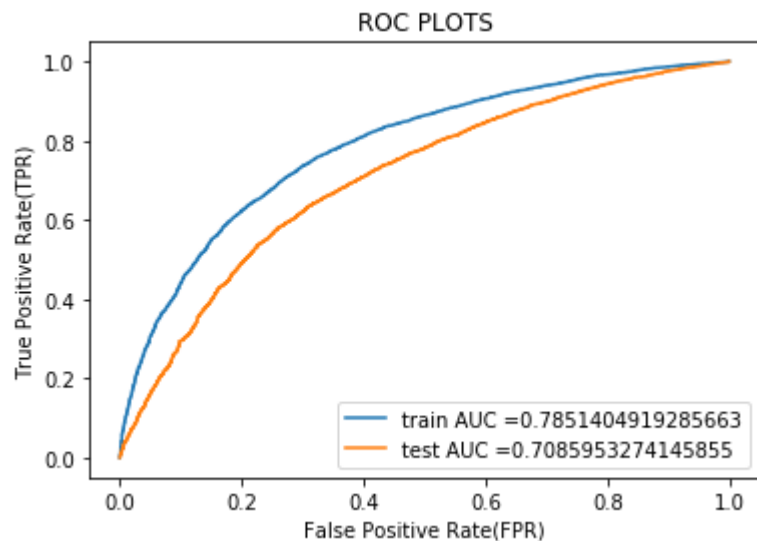
ERROR PLOTS

```
1  ##Fitting Model to Hyper-Parameter Curve
2  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
3  from sklearn.metrics import roc_curve, auc
4
5
6  neigh = LogisticRegression(C=10**-3,class_weight='balanced');
7  neigh.fit(X_set1_train ,y_train)
8  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
9  # not the predicted outputs
10
11 train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set1_train)[:,1])
12 test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set1_test)[:,1])
13
14 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
15 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
16 plt.legend()
17 plt.ylabel("True Positive Rate(TPR)")
18 plt.xlabel("False Positive Rate(FPR)")
19 plt.title("ROC PLOTS")
20 plt.show()
```

ROC PLOTS

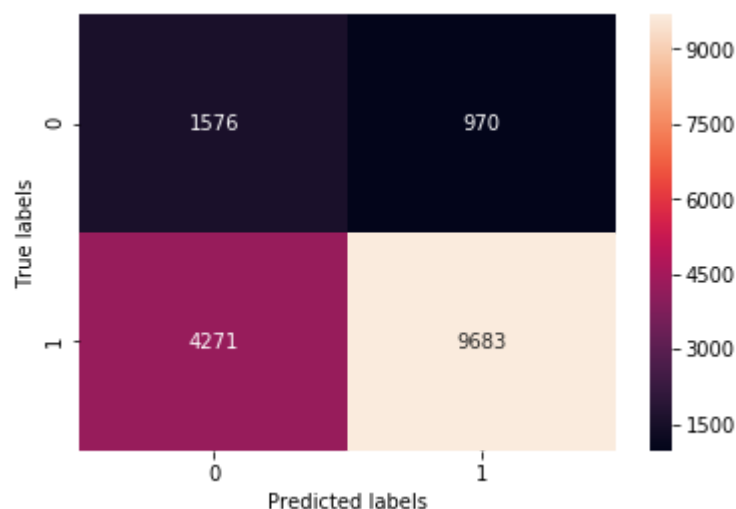

```
1  # Heat map representaion of the Confusion metrix
2  # https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3  import seaborn as sns
4  import matplotlib.pyplot as plt
5
6  ax= plt.subplot()
7  sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set1_train )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
8
9  # labels, title and ticks
10 ax.set_xlabel('Predicted labels')
11 ax.set_ylabel('True labels');
12 print("The Confusion metrix for Train data ")
```

The Confusion metrix for Train data



```
1  #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4
5  ax= plt.subplot()
6  sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set1_test )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
7
8  # labels, title and ticks
9  ax.set_xlabel('Predicted labels')
10 ax.set_ylabel('True labels');
11 print("The Confusion metrix for Test data ")
```

⊏→

The Confusion metrix for Test data



## Applying the Logistic Regression model on SET:2

```
1  from sklearn.metrics import roc_auc_score
2  import matplotlib.pyplot as plt
3  #from sklearn.grid_search import GridSearchCV
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.model_selection import learning_curve, GridSearchCV
6
7  clf = LogisticRegression(class_weight='balanced');
8  parameters ={'C':[10**-4, 10**-3,10**-2,1,10,100,1000,500,1000,10000]}
9  sd = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc',return_train_score=True)
10 sd.fit(X_set2_train, y_train);
11
12 train_auc= sd.cv_results_['mean_train_score']
13 train_auc_std= sd.cv_results_['std_train_score']
14 cv_auc =sd.cv_results_['mean_test_score']
15 cv_auc_std=sd.cv_results_['std_test_score']
16
17 plt.plot(parameters['C'], train_auc, label='Train AUC')
18 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
19 plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
20
21 plt.plot(parameters['C'], cv_auc, label='CV AUC')
22 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
23 plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
24
```

```
25 plt.scatter(parameters['C'], train_auc, label='Train AUC points')
26 plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
27 plt.xscale('log')
28
29 plt.legend()
30 plt.xlabel("C : hyperparameter")
31 plt.ylabel("AUC")
32 plt.title("ERROR PLOTS")
33 plt.grid()
34 plt.show()
```
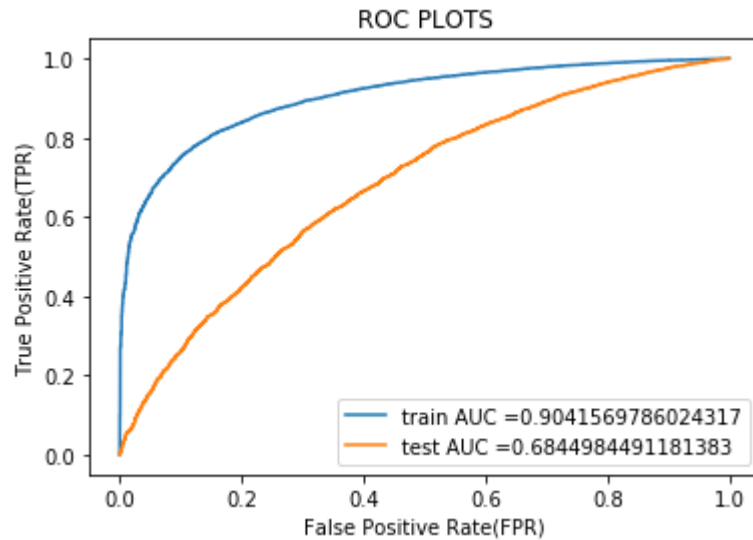


```
1  #Fitting Model to Hyper-Parameter Curve
2  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
3  from sklearn.metrics import roc_curve, auc
4
5
6  neigh = LogisticRegression(C=1,class_weight='balanced');
7  neigh.fit(X_set2_train ,y_train)
8  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
9  # not the predicted outputs
10
11 train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set2_train)[:,1])
12 test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set2_test)[:,1])
13
14 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
15 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
16 plt.legend()
17 plt.ylabel("True Positive Rate(TPR)")
18 plt.xlabel("False Positive Rate(FPR)")
```

```
19  plt.title("ROC PLOTS")
20  plt.show()
21
```
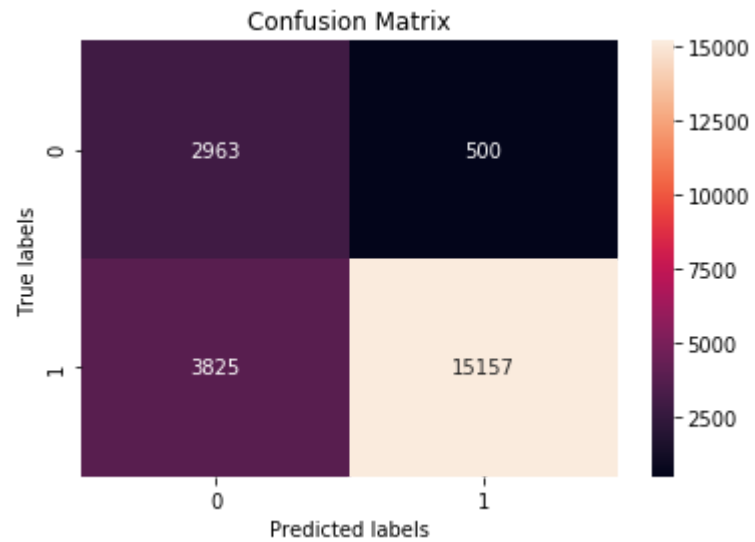


```
1  #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
2  ax= plt.subplot()
3  sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set2_train )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
4  ax.set_xlabel('Predicted labels')
5  ax.set_ylabel('True labels');
6  ax.set_title('Confusion Matrix')
7  print("The Confusion metrix for Train data ")
```
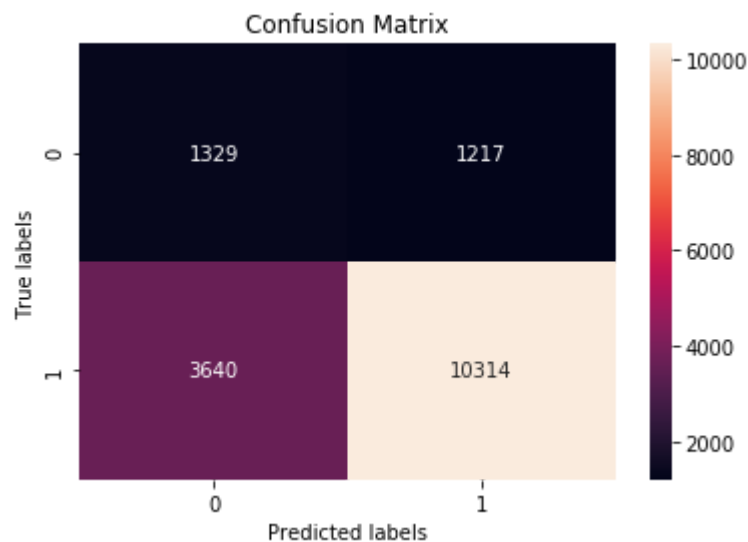
The Confusion metrix for Train data



```
1  ax= plt.subplot()
2  sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set2_test )), annot=True, ax = ax,fmt='g')
3  ax.set_xlabel('Predicted labels')
4  ax.set_ylabel('True labels')
5  ax.set_title('Confusion Matrix')
6  print("The Confusion metrix for Test data ")
```

⟶

The Confusion metrix for Test data



## Applying the Logistic Regression model on SET:3

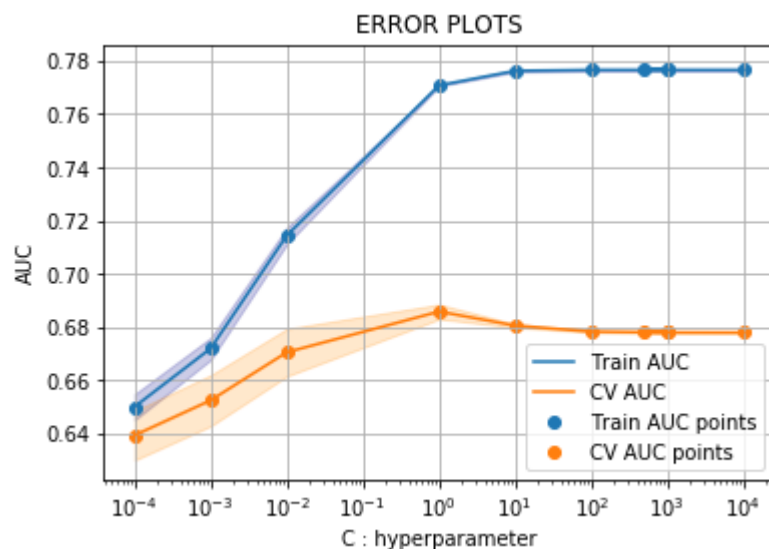Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)

```
 1  clf = LogisticRegression(class_weight='balanced');
 2  parameters ={'C':[10**-4, 10**-3,10**-2,1,10,100,1000,500,1000,10000]}
 3  cl = GridSearchCV(clf , parameters, cv=3, scoring='roc_auc',return_train_score=True)
 4  cl.fit(X_set3_train, y_train);
 5
 6  train_auc= cl.cv_results_['mean_train_score']
 7  train_auc_std= cl.cv_results_['std_train_score']
 8  cv_auc = cl.cv_results_['mean_test_score']
 9  cv_auc_std= cl.cv_results_['std_test_score']
10
11  plt.plot(parameters['C'], train_auc, label='Train AUC')
12  # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
13  plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
14
15  plt.plot(parameters['C'], cv_auc, label='CV AUC')
16  # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
17  plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
18
19  plt.scatter(parameters['C'], train_auc, label='Train AUC points')
20  plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
21  plt.xscale('log')
```

```
22
23  plt.legend()
24  plt.xlabel("C : hyperparameter")
25  plt.ylabel("AUC")
26  plt.title("ERROR PLOTS")
27  plt.grid()
28  plt.show()
```
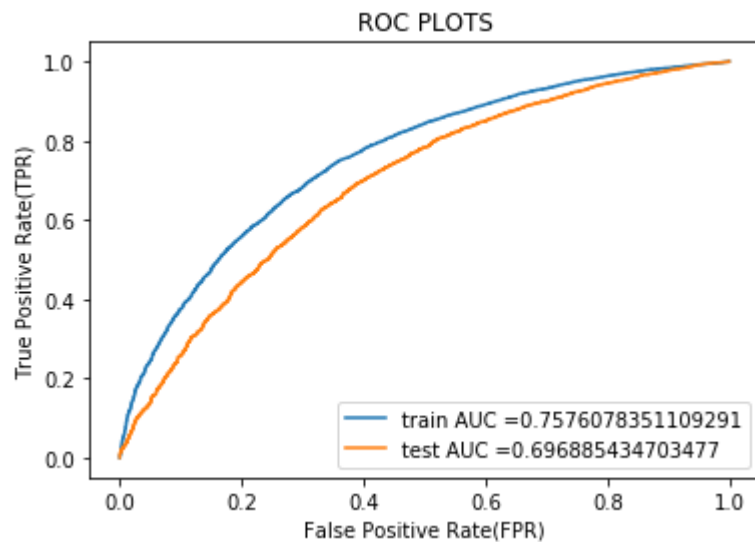


ERROR PLOTS

```
1   #Fitting Model to Hyper-Parameter Curve:
2   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
3   from sklearn.metrics import roc_curve, auc
4
5
6   neigh = LogisticRegression(C=1,class_weight='balanced');
7   neigh.fit(X_set3_train ,y_train)
8   # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
9   # not the predicted outputs
10
11  train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set3_train)[:,1])
12  test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set3_test)[:,1])
13
14  plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
15  plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
16  plt.legend()
17  plt.ylabel("True Positive Rate(TPR)")
18  plt.xlabel("False Positive Rate(FPR)")
19  plt.title("ROC PLOTS")
20  plt.show()
21  print("="*100)
```

===================================================================================================

```
1  ax= plt.subplot()
2  sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set3_train )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
3  ax.set_xlabel('Predicted labels');
4  ax.set_ylabel('True labels');
5  print("The Confusion metrix for Train data ")
```
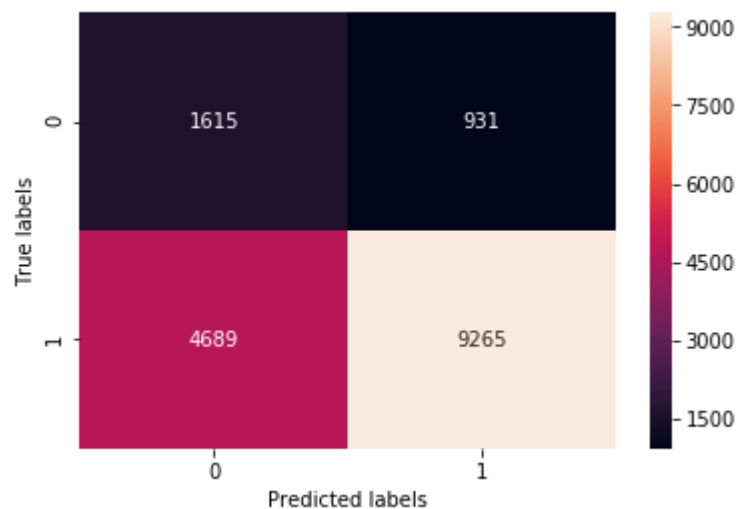
```
1 ax= plt.subplot()
2 sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set3_test )), annot=True, ax = ax,fmt='g')
3 ax.set_xlabel('Predicted labels')
4 ax.set_ylabel('True labels')
5 print("The Confusion metrix for Test data ")
```

⤷   The Confusion metrix for Test data



▶ **Applying the Logistic Regression model on SET:4**

categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

⤓ *4 cells hidden*

▶ **Logistic Regression on SET:5**

1. school_state : categorical data
2. clean_categories : categorical data
3. clean_subcategories : categorical data
4. project_grade_category :categorical data
5. teacher_prefix : categorical data
6. quantity : numerical data

7. teacher_number_of_previously_posted_projects : numerical data

8. price : numerical data

9. sentiment score's of each of the essay : numerical data

10. number of words in the title : numerical data

11. number of words in the combine essays : numerical data

↳ 11 cells hidden

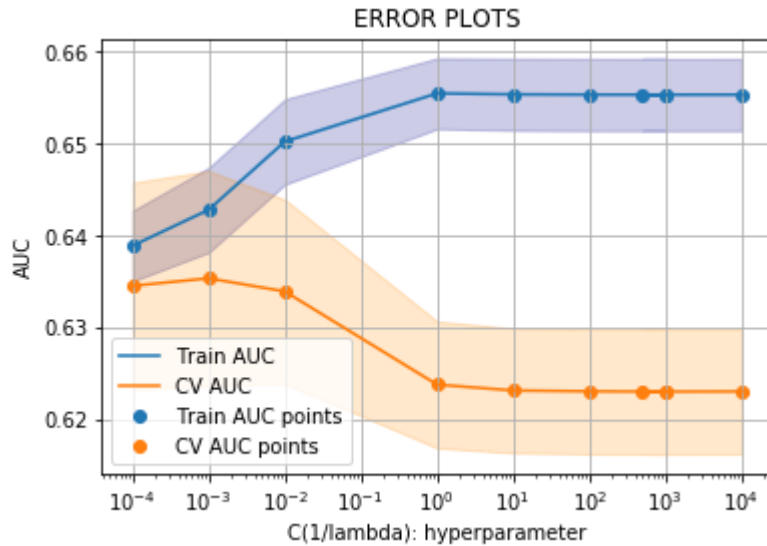## ▾ Applying the Logistic Regression model on SET:5

```
1
2  from sklearn.metrics import roc_auc_score
3  import matplotlib.pyplot as plt
4  #from sklearn.grid_search import GridSearchCV
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import learning_curve, GridSearchCV
7
8  """
9  y_true : array, shape = [n_samples] or [n_samples, n_classes]
10 True binary labels or binary label indicators.
11
12 y_score : array, shape = [n_samples] or [n_samples, n_classes]
13 Target scores, can either be probability estimates of the positive class, confidence values, or non-thresholded measure of
14 decisions (as returned by "decision_function" on some classifiers).
15 For binary y_true, y_score is supposed to be the score of the class with greater label.
16
17 """
18
19 clf = LogisticRegression(class_weight='balanced');
20 parameters ={'C':[10**-4, 10**-3,10**-2,1,10,100,1000,500,1000,10000]}
21 cl = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc',return_train_score=True)
22 cl.fit(X_set5_train, y_train);
23
24 train_auc= cl.cv_results_['mean_train_score']
25 train_auc_std= cl.cv_results_['std_train_score']
26 cv_auc = cl.cv_results_['mean_test_score']
27 cv_auc_std= cl.cv_results_['std_test_score']
28
29 plt.plot(parameters['C'], train_auc, label='Train AUC')
30 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
31 plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
32
33 plt.plot(parameters['C'], cv_auc, label='CV AUC')
34 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
35 plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
36
```

```
37 plt.scatter(parameters['C'], train_auc, label='Train AUC points')
38 plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
39 plt.xscale('log')
40
41 plt.legend()
42 plt.xlabel("C(1/lambda): hyperparameter")
43 plt.ylabel("AUC")
44 plt.title("ERROR PLOTS")
45 plt.grid()
46 plt.show()
```
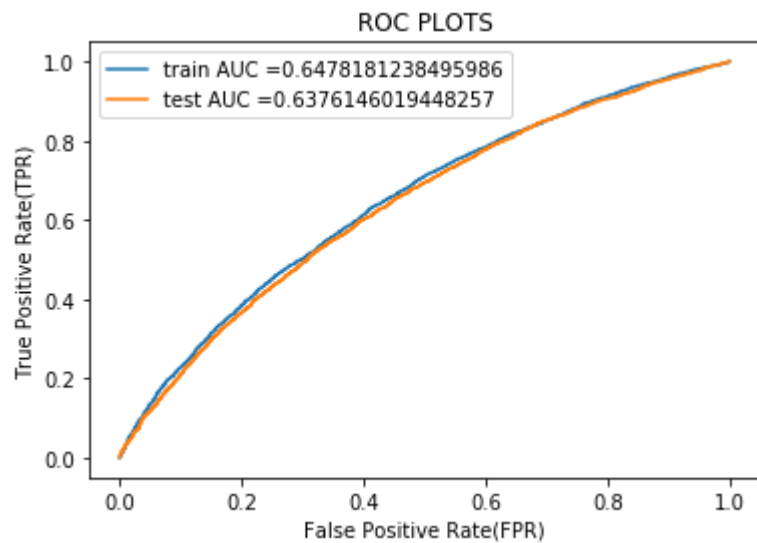


```
1  # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2  from sklearn.metrics import roc_curve, auc
3  neigh = LogisticRegression(C=10**-2,class_weight='balanced');
4  neigh.fit(X_set5_train ,y_train)
5  # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
6  # not the predicted outputs
7
8  train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set5_train)[:,1])
9  test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set5_test)[:,1])
10
11 plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
12 plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
13 plt.legend()
14 plt.ylabel("True Positive Rate(TPR)")
15 plt.xlabel("False Positive Rate(FPR)")
16 plt.title("ROC PLOTS")
17 plt.show()
```

⤷

**ROC PLOTS**



```
1  #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
2  ax= plt.subplot()
3  sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set5_train)), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
4  ax.set_xlabel('Predicted labels')
5  ax.set_ylabel('True labels')
6  ax.set_title('Confusion Matrix')
```

⤷

```
    Text(0.5, 1.0, 'Confusion Matrix')
```

```
1  ax= plt.subplot()
2  sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set5_test )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
3
4  ax.set_xlabel('Predicted labels')
5  ax.set_ylabel('True labels')
6  ax.set_title('Confusion Matrix')
```
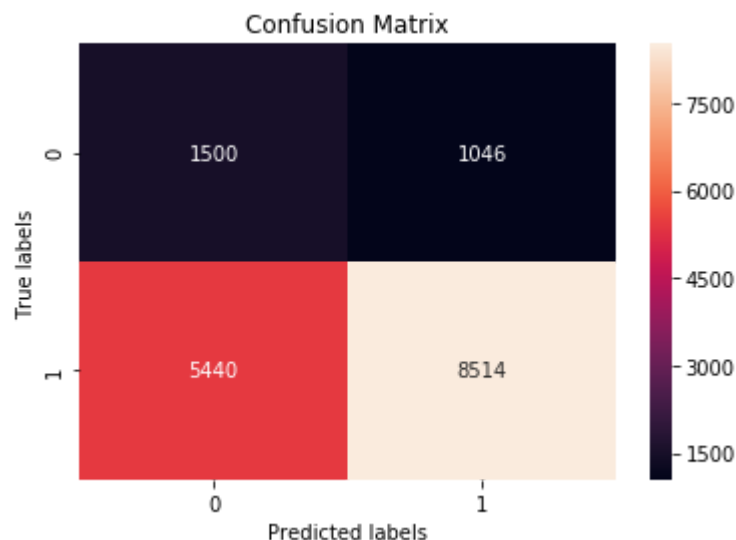
⤷    Text(0.5, 1.0, 'Confusion Matrix')

**Confusion Matrix**

|                | 0    | 1    |
|----------------|------|------|
| **0** (True)   | 1500 | 1046 |
| **1** (True)   | 5440 | 8514 |

Predicted labels

Conclusion :

1. If we compare the roc curves between model with Text features and model without the text features , the model with the text features included is bet
auc and the Confusion metrix gives better results .

```
 1  #how to use pretty table http://zetcode.com/python/prettytable/
 2  from prettytable import PrettyTable
 3
 4  tb = PrettyTable()
 5  tb.field_names= ("Vectorizer", "HyperParameter", "AUC")
 6  tb.add_row(["BOW",10**-3, 70])
 7  tb.add_row(["Tf-Idf",1, 66])
 8  tb.add_row(["AVGW2V",1, 69])
 9  tb.add_row(["Tf-Idf w2v", 10**-2, 60])
10  tb.add_row(["Set 5",10**-3, 64])
11  print(tb.get_string(titles = "Logistic Reg> - Observations"))
```

```
+-----------+---------------+-----+
| Vectorizer | HyperParameter | AUC |
+-----------+---------------+-----+
|    BOW     |     0.001     | 70  |
|   Tf-Idf   |       1       | 66  |
|   AVGW2V   |       1       | 69  |
| Tf-Idf w2v |     0.01      | 60  |
|   Set 5    |     0.001     | 64  |
+-----------+---------------+-----+
```