

# DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now, a large number of volunteers is needed to manually screen each submission before it's approved to be posted on the DonorsChoose.org website.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are three main problems they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be posted as quickly and as efficiently as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience for teachers
- How to focus volunteer time on the applications that need the most assistance

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school. DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

# About the DonorsChoose Data Set

The `train.csv` data set provided by DonorsChoose contains the following features:

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. <b>Example:</b> 123456789
<code>project_title</code>		Title of the project. <b>Example:</b> Art Will Make You a Better Person
<code>project_grade_category</code>		Grade level of students for which the project is targeted. One of the following enumerated list of categories: <ul style="list-style-type: none"><li>• Grades K-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>
<code>project_subject_categories</code>		One or more (comma-separated) subject categories for the project from the following enumerated list of categories: <ul style="list-style-type: none"><li>• Applied &amp; Design</li><li>• Care &amp; Safety</li><li>• Health &amp; Physical Education</li><li>• History &amp; Social Studies</li><li>• Literacy &amp; Language</li><li>• Math &amp; Science</li><li>• Music &amp; Arts</li><li>• Special Education</li></ul>
<code>project_subject_subcategories</code>		One or more (comma-separated) subject subcategories for the project from the following enumerated list of categories: <ul style="list-style-type: none"><li>• Music &amp; Arts</li><li>• Literacy &amp; Language, Math &amp; Science</li></ul>
<code>school_state</code>		State where school is located ( <a href="https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_abbreviations">Two-letter U.S. postal abbreviations</a> ) ( <a href="https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_abbreviations">https://en.wikipedia.org/wiki/List of U.S. state abbreviations#Postal abbreviations</a> )
<code>project_resource_summary</code>		An explanation of the resources needed for the project. <b>Example:</b> My students need hands on literacy materials to enhance their sensory
<code>project_essay_1</code>		First application essay
<code>project_essay_2</code>		Second application essay
<code>project_essay_3</code>		Third application essay
<code>project_essay_4</code>		Fourth application essay
<code>project_submitted_datetime</code>		Datetime when project application was submitted. <b>Example:</b> 2018-01-12T12:43:21Z
<code>teacher_id</code>		A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4f1

teacher\_prefix

- .....

Number of project applications previously submitted by the sam

## Exe

Additionally, the `resources.csv` data set provides more data about the resources required for each project. Each line in this file represents a resource required by a project:

**Note:** Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- \_\_project\_essay\_1: \_\_ "Introduce us to your classroom"
- \_\_project\_essay\_2: \_\_ "Tell us more about your students"
- \_\_project\_essay\_3: \_\_ "Describe how your students will use the materials you're requesting"
- project\_essay\_3: \_\_ "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- \_\_project\_essay\_1: \_\_ "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- \_\_project\_essay\_2: \_\_ "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with project\_submitted\_datetime of 2016-05-17 and later, the values of project\_essay\_3 and project\_essay\_4 will be NaN.

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
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.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
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 chart_studio import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
```

## 1.1 Reading Data

In [2]:

```
project_data = pd.read_csv('../train_data.csv')
resource_data = pd.read_csv('../resources.csv')
```

In [3]:

```
print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

Number of data points in train data (109248, 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']

In [4]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

Number of data points in train data (1541272, 4)

['id' 'description' 'quantity' 'price']

Out[4]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

## 1.2 Preprocessing Categorical Data

### 1.2.1 preprocessing project\_subject\_categories

In [5]:

```
categories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty) ex: "Math & Science"=> "Math&Science"
            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
    temp = temp.replace('&', '_') # we are replacing the & value into
    cat_list.append(temp.strip())

project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)

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

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

In [6]:

```
sorted_cat_dict.keys()
```

Out[6]:

```
dict_keys(['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language'])
```

## 1.2.2 preprocessing of project\_subject\_subcategories

In [7]:

```
sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " #"
    temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
```

In [8]:

```
sorted_sub_cat_dict.keys()
```

Out[8]:

```
dict_keys(['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy'])
```

### 1.2.3 preprocessing of School State

In [9]:

```
project_data['school_state'].unique()
```

Out[9]:

```
array(['IN', 'FL', 'AZ', 'KY', 'TX', 'CT', 'GA', 'SC', 'NC', 'CA', 'NY',  
      'OK', 'MA', 'NV', 'OH', 'PA', 'AL', 'LA', 'VA', 'AR', 'WA', 'WV',  
      'ID', 'TN', 'MS', 'CO', 'UT', 'IL', 'MI', 'HI', 'IA', 'RI', 'NJ',  
      'MO', 'DE', 'MN', 'ME', 'WY', 'ND', 'OR', 'AK', 'MD', 'WI', 'SD',  
      'NE', 'NM', 'DC', 'KS', 'MT', 'NH', 'VT'], dtype=object)
```

In [10]:

```
project_data['school_state'][project_data['school_state'].isnull()==True]
```

Out[10]:

```
Series([], Name: school_state, dtype: object)
```

In [11]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039  
my_counter = Counter()  
for word in project_data['school_state'].values:  
    my_counter.update(word.split())  
  
school_state_dict = dict(my_counter)  
sorted_school_state_dict = dict(sorted(school_state_dict.items(), key=lambda kv: kv[1]))
```

In [12]:

```
sorted_school_state_dict.keys()
```

Out[12]:

```
dict_keys(['VT', 'WY', 'ND', 'MT', 'RI', 'SD', 'NE', 'DE', 'AK', 'NH', 'W  
V', 'ME', 'HI', 'DC', 'NM', 'KS', 'IA', 'ID', 'AR', 'CO', 'MN', 'OR', 'K  
Y', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'AL', 'WI', 'VA', 'AZ', 'NJ', 'O  
K', 'WA', 'MA', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'N  
C', 'FL', 'NY', 'TX', 'CA'])
```

## 1.2.4 preprocessing of Teacher Prefix

In [13]:

```
project_data.groupby(['teacher_prefix'])['teacher_prefix'].count()
```

Out[13]:

```
teacher_prefix  
Dr.          13  
Mr.         10648  
Mrs.         57269  
Ms.          38955  
Teacher      2360  
Name: teacher_prefix, dtype: int64
```



In [14]:

```
project_data['teacher_prefix'][project_data['teacher_prefix'].isnull()==True]
```

Out[14]:

```
7820      NaN
30368     NaN
57654     NaN
Name: teacher_prefix, dtype: object
```

In [15]:

```
project_data['teacher_prefix'].fillna(project_data['teacher_prefix'].mode()[0],inplace=True)
```

In [16]:

```
project_data['teacher_prefix'][project_data['teacher_prefix'].isnull()==True]
```

Out[16]:

```
Series([], Name: teacher_prefix, dtype: object)
```

In [17]:

```
project_data['teacher_prefix'].unique()
```

Out[17]:

```
array(['Mrs.', 'Mr.', 'Ms.', 'Teacher', 'Dr.'], dtype=object)
```

In [18]:

```
teacher_prefix = list(project_data['teacher_prefix'].values)

teacher_prefix_list = []
for i in teacher_prefix:
    temp = ""
    temp = i.split('.')
    temp = i.replace('.', '')
    teacher_prefix_list.append(temp)

project_data['clean_teacher_prefix'] = teacher_prefix_list
project_data.drop(['teacher_prefix'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_teacher_prefix'].values:
    my_counter.update(word.split())

teacher_prefix_dict = dict(my_counter)
sorted_teacher_prefix_dict = dict(sorted(teacher_prefix_dict.items(), key=lambda kv: kv[1]))
```

In [19]:

```
sorted_teacher_prefix_dict.keys()
```

Out[19]:

```
dict_keys(['Dr', 'Teacher', 'Mr', 'Ms', 'Mrs'])
```

In [20]:

```
project_data.groupby(['clean_teacher_prefix'])['clean_teacher_prefix'].count()
```

Out[20]:

```
clean_teacher_prefix
Dr                13
Mr             10648
Mrs            57272
Ms             38955
Teacher         2360
Name: clean_teacher_prefix, dtype: int64
```

## 1.2.5 preprocessing of Project Grade Category

In [21]:

```
project_data.groupby(['project_grade_category'])['project_grade_category'].count()
```

Out[21]:

```
project_grade_category
Grades 3-5          37137
Grades 6-8          16923
Grades 9-12         10963
Grades PreK-2       44225
Name: project_grade_category, dtype: int64
```

In [22]:

```
project_data['project_grade_category'][project_data['project_grade_category'].isnull()==True]
```

Out[22]:

```
Series([], Name: project_grade_category, dtype: object)
```

In [23]:

```
project_grade_category = list(project_data['project_grade_category'].values)

project_grade_category_list = []
for i in project_grade_category:
    temp = ""
    temp = i.split(' ')
    temp = i.replace('Grades ', '')
    project_grade_category_list.append(temp)

project_data['clean_project_grade_category'] = project_grade_category_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_project_grade_category'].values:
    my_counter.update(word.split())

project_grade_category_dict = dict(my_counter)
sorted_project_grade_category_dict = dict(sorted(project_grade_category_dict.items(), key=lambda kv: kv[1]))
```

In [24]:

```
sorted_project_grade_category_dict.keys()
```

Out[24]:

```
dict_keys(['9-12', '6-8', '3-5', 'PreK-2'])
```

In [25]:

```
project_data.groupby(['clean_project_grade_category'])['clean_project_grade_category'].count()
```

Out[25]:

```
clean_project_grade_category
3-5          37137
6-8          16923
9-12         10963
PreK-2       44225
Name: clean_project_grade_category, dtype: int64
```

In [ ]:

## 1.3 Text preprocessing

In [26]:

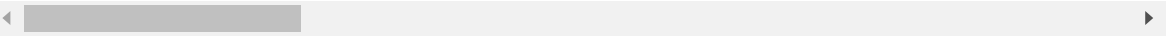
```
# merge two column text dataframe:
project_data["essay"] = project_data["project_essay_1"].map(str) + \
    project_data["project_essay_2"].map(str) + \
    project_data["project_essay_3"].map(str) + \
    project_data["project_essay_4"].map(str)
```

In [27]:

```
project_data.head(2)
```

Out[27]:

	Unnamed: 0	id	teacher_id	school_state	project_submitted_date
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	IN	2016-12-05 1
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	FL	2016-10-25 0



In [28]:

```
#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

In [29]:

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\nWe have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\r\n\r\nThe limits of your language are the limits of your world.\r\n\r\n-Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English along side of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills.\r\n\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\n\r\nnannan

=====

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\n\r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity. My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still. nannan

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed r

aces in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

=====  
My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\n\r\nThey also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

=====  
The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\n\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.nannan

=====

In [30]:

```
# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\'re", " are", phrase)
    phrase = re.sub(r"\'s", " is", phrase)
    phrase = re.sub(r"\'d", " would", phrase)
    phrase = re.sub(r"\'ll", " will", phrase)
    phrase = re.sub(r"\'t", " not", phrase)
    phrase = re.sub(r"\'ve", " have", phrase)
    phrase = re.sub(r"\'m", " am", phrase)
    return phrase
```

In [31]:

```
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\n\r\nThey also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

=====



In [32]:

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. They also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.   
nannan

In [33]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays cognitive delays gross fine motor delays to autism They are eager beavers and always strive to work their hardest working past their limitations The materials we have are the ones I seek out for my students I teach in a Title I school where most of the students receive free or reduced price lunch Despite their disabilities and limitations my students love coming to school and come eager to learn and explore Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting This is how my kids feel all the time They want to be able to move as they learn or so they say Wobble chairs are the answer and I love them because they develop their core which enhances gross motor and in turn fine motor skills They also want to learn through games my kids do not want to sit and do worksheets They want to learn to count by jumping and playing Physical engagement is the key to our success The number toss and color and shape mats can make that happen My students will forget they are doing work and just have the fun a 6 year old deserves   
nannan

In [34]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you'r
e", "you've",\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't
hey', 'them', 'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "th
at'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'ha
d', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as'
, 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through'
, 'during', 'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'ov
er', 'under', 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'an
y', 'both', 'each', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too'
, 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'no
w', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't",
'doesn', "doesn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
tn', "mightn't", 'mustn',\
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'w
asn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [35]:

```
# Combining all the above stundents
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

100%|██████████| 109248/109248 [01:21<00:00, 1341.45it/s]

In [36]:

```
# after preprocessing
preprocessed_essays[20000]
```

Out[36]:

```
'my kindergarten students varied disabilities ranging speech language dela
ys cognitive delays gross fine motor delays autism they eager beavers alwa
ys strive work hardest working past limitations the materials ones i seek
students i teach title i school students receive free reduced price lunch
despite disabilities limitations students love coming school come eager le
arn explore have ever felt like ants pants needed groove move meeting this
kids feel time the want able move learn say wobble chairs answer i love de
velop core enhances gross motor turn fine motor skills they also want lear
n games kids not want sit worksheets they want learn count jumping playing
physical engagement key success the number toss color shape mats make happ
en my students forget work fun 6 year old deserves nannan'
```

In [37]:

```
project_data['preprocessed_essays'] = preprocessed_essays
project_data.drop(['essay'], axis=1, inplace=True)
```

## 1.4 Preprocessing of `project\_title`

In [38]:

```
# similarly you can preprocess the titles also
```

In [39]:

```
project_data['project_title'][2000:2010]
```

Out[39]:

```
2000          Steady Stools for Active Learning
2001          Classroom Supplies
2002  Kindergarten Students Deserve Quality Books a...
2003          Listen to Understand!
2004          iPads to iGnite Learning
2005          Tablets For Learning
2006          Go P.E.!
2007          Making Learning Fun!
2008  Empowerment Through Silk Screen Designed Tee S...
2009          Let's Play Together!
Name: project_title, dtype: object
```

In [40]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_titles.append(sent.lower().strip())
```

100%|██████████| 109248/109248 [00:03<00:00, 28057.16it/s]

In [41]:

```
preprocessed_titles[2000:2010]
```

Out[41]:

```
['steady stools active learning',
 'classroom supplies',
 'kindergarten students deserve quality books vibrant rug',
 'listen understand',
 'ipads ignite learning',
 'tablets for learning',
 'go p e',
 'making learning fun',
 'empowerment through silk screen designed tee shirts',
 'let play together']
```

In [42]:

```
project_data['preprocessed_titles'] = preprocessed_titles
project_data.drop(['project_title'], axis=1, inplace=True)
```

## 1.5 Preparing data for models

In [43]:

```
project_data.columns
```

Out[43]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
      'project_essay_3', 'project_essay_4', 'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_teacher_prefix',
      'clean_project_grade_category', 'preprocessed_essays',
      'preprocessed_titles'],
      dtype='object')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
- project\_title : text data
- text : text data
- project\_resource\_summary: text data (optional)
- quantity : numerical
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

## 1.5.1 Vectorizing Categorical data

- <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/> (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>)

```
# we use count vectorizer to convert the values into one from sklearn.feature_extraction.text import
CountVectorizer vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
binary=True) categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names()) print("Shape of matrix after one hot encoding ",categories_one_hot.shape)#
we use count vectorizer to convert the values into one vectorizer =
CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names()) print("Shape of matrix after one hot encoding
",sub_categories_one_hot.shape)# you can do the similar thing with state, teacher_prefix and
project_grade_category also
```

## 1.5.2 Vectorizing Text data

### 1.5.2.1 Bag of words

```
# We are considering only the words which appeared in at least 10 documents(rows or projects). vectorizer =
CountVectorizer(min_df=10) text_bow = vectorizer.fit_transform(preprocessed_essays) print("Shape of matrix
after one hot encoding ",text_bow.shape)# you can vectorize the title also # before you vectorize the title make
sure you preprocess it
```

### 1.5.2.2 TFIDF vectorizer

```
from sklearn.feature_extraction.text import TfidfVectorizer vectorizer = TfidfVectorizer(min_df=10) text_tfidf =
vectorizer.fit_transform(preprocessed_essays) print("Shape of matrix after one hot encoding ",text_tfidf.shape)
```

### 1.5.2.3 Using Pretrained Models: Avg W2V

```
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print("Loading Glove Model")
    f = open(gloveFile, 'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print("Done.", len(model), " words loaded!")
    return model

model = loadGloveModel('glove.42B.300d.txt')
# ===== Output: Loading Glove Model
1917495it [06:32, 4879.69it/s] Done. 1917495 words loaded!
# ===== words = []
for i in preprocod_texts:
    words.extend(i.split(' '))
for i in preprocod_titles:
    words.extend(i.split(' '))
print("all the words in the corpus", len(words))
words = set(words)
print("the unique words in the corpus", len(words))
inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our corpus", len(inter_words),
      "(", np.round(len(inter_words)/len(words)*100, 3), "%)")
words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file with
open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
glove_words = set(model.keys())
# average Word2Vec #
compute average word2vec for each review.
avg_w2v_vectors = []
# the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays):
    # for each review/sentence vector = np.zeros(300)
    # as word vectors are of zero length cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split():
        # for each word in a review/sentence if word in glove_words:
        vector += model[word]
    cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors.append(vector)
print(len(avg_w2v_vectors))
print(len(avg_w2v_vectors[0]))
```

### 1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(preprocessed_essays)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
# average Word2Vec #
compute average word2vec for each review.
tfidf_w2v_vectors = []
# the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays):
    # for each review/sentence vector = np.zeros(300)
    # as word vectors are of zero length tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split():
        # for each word in a review/sentence if (word in glove_words) and (word in tfidf_words):
        vec = model[word]
        # getting the vector for each word # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
        tf_idf = dictionary[word] * (sentence.count(word)/len(sentence.split()))
        # getting the tfidf value for each word vector += (vec * tf_idf)
        # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)
print(len(tfidf_w2v_vectors))
print(len(tfidf_w2v_vectors[0]))
# Similarly you can vectorize for title also
```

### 1.5.3 Vectorizing Numerical features

```
price_data = resource_data.groupby('id').agg({'price': 'sum', 'quantity': 'sum'}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler
# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error # ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)
price_scalar = StandardScaler()
price_scalar.fit(project_data['price'].values.reshape(-1, 1))
# finding the mean and standard deviation of this data
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
# Now standardize
```

the data with above mean and variance. price\_standardized =  
price\_scalar.transform(project\_data['price'].values.reshape(-1, 1))price\_standardized

### 1.5.4 Merging all the above features

- we need to merge all the numerical vectors i.e categorical, text, numerical vectors

```
print(categories_one_hot.shape) print(sub_categories_one_hot.shape) print(text_bow.shape)
print(price_standardized.shape)# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack # with the same hstack function we are concatenating a sparse matrix and a
dense matrix :) X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

In [ ]:

## 1.6 Merging Numerical data in Resources to project\_data

In [44]:

```
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

In [ ]:

In [ ]:

### Computing Sentiment Scores

In [45]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# import nltk
# nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest
students with the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multi
ple intelligences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety
of different backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school
is a caring community of successful \
learners which can be seen through collaborative student project based learning in and
out of the classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities
to practice a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspec
t of the kindergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love
to role play in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with
real food i will take their idea \
and create common core cooking lessons where we learn important math and writing concep
ts while cooking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that wen
t into making the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this p
roject would expand our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make hom
emade applesauce make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create o
ur own cookbooks to be printed and \
shared with families students will gain math and literature skills as well as a life lo
ng enjoyment for healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

for k in ss:
    print('{0}: {1}, '.format(k, ss[k]), end='')

# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
```

neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,

In [ ]:

In [ ]:



In [ ]:

## Assignment 5: Logistic Regression




1. **[Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets**

- **Set 1:** categorical, numerical features + project\_title(BOW) + preprocessed\_eassay (`BOW with bi-grams` with `min\_df=10` and `max\_features=5000`)
- **Set 2:** categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (`TFIDF with bi-grams` with `min\_df=10` and `max\_features=5000`)
- **Set 3:** categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_eassay (TFIDF W2V)

2. **Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)**

- Find the best hyper parameter which will give the maximum [AUC](https://www.applidaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. **Representation of results**

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.  

- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.  

- Along with plotting ROC curve, you need to print the [confusion matrix](https://www.applidaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](https://seaborn.pydata.org/generated/seaborn.heatmap.html).  
  
(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)  
(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)  
(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

4. **[Task-2] Apply Logistic Regression on the below feature set **Set 5** by finding the best hyper parameter as suggested in step 2 and step 3.**

5. **Consider these set of features **Set 5** :**

- **school\_state** : categorical data
- **clean\_categories** : categorical data
- **clean\_subcategories** : categorical data
- **project\_grade\_category** :categorical data
- **teacher\_prefix** : categorical data
- **quantity** : numerical data
- **teacher\_number\_of\_previously\_posted\_projects** : numerical data
- **price** : numerical data
- **sentiment score's of each of the essay** : numerical data
- **number of words in the title** : numerical data
- **number of words in the combine essays** : numerical data

And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3

(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>).

6. **Conclusion** (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this [prettytable](#) library.

(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>) [link](#)

(<http://zetcode.com/python/prettytable/>)



### Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf). (<https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf>)

## 2. Logistic Regression

### 2.1 Splitting data into Train and cross validation(or test): Stratified Sampling

In [46]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [47]:

```
project_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 109248 entries, 0 to 109247
```

```
Data columns (total 20 columns):
```

#	Column	Non-Null Count	Dtype
0	Unnamed: 0	109248 non-null	int64
1	id	109248 non-null	object
2	teacher_id	109248 non-null	object
3	school_state	109248 non-null	object
4	project_submitted_datetime	109248 non-null	object
5	project_essay_1	109248 non-null	object
6	project_essay_2	109248 non-null	object
7	project_essay_3	3758 non-null	object
8	project_essay_4	3758 non-null	object
9	project_resource_summary	109248 non-null	object
10	teacher_number_of_previously_posted_projects	109248 non-null	int64
11	project_is_approved	109248 non-null	int64
12	clean_categories	109248 non-null	object
13	clean_subcategories	109248 non-null	object
14	clean_teacher_prefix	109248 non-null	object
15	clean_project_grade_category	109248 non-null	object
16	preprocessed_essays	109248 non-null	object
17	preprocessed_titles	109248 non-null	object
18	price	109248 non-null	float6

```
4  
19 quantity 109248 non-null int64
```

```
dtypes: float64(1), int64(4), object(15)
```

```
memory usage: 17.5+ MB
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
- project\_title : text data
- Essay : text data
- quantity : numerical
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

In [48]:

```
data1 = project_data.drop(['Unnamed: 0', 'id', 'project_submitted_datetime', 'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4', 'project_resource_summary', 'teacher_id'], axis = 1)
```

In [49]:

```
data1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 109248 entries, 0 to 109247
Data columns (total 11 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   school_state                          109248 non-null  object
1   teacher_number_of_previously_posted_projects  109248 non-null  int64
2   project_is_approved                  109248 non-null  int64
3   clean_categories                      109248 non-null  object
4   clean_subcategories                  109248 non-null  object
5   clean_teacher_prefix                 109248 non-null  object
6   clean_project_grade_category         109248 non-null  object
7   preprocessed_essays                  109248 non-null  object
8   preprocessed_titles                  109248 non-null  object
9   price                                109248 non-null  float64
4
10  quantity                             109248 non-null  int64
dtypes: float64(1), int64(3), object(7)
memory usage: 10.0+ MB
```

In [50]:

```
data1 = data1[:50000]
```

In [51]:

```
y = data1['project_is_approved'].values
X = data1.drop(['project_is_approved'], axis=1)
X.head(1)
```

Out[51]:

	school_state	teacher_number_of_previously_posted_projects	clean_categories	clean_subc
0	IN	0	Literacy_Language	ES

In [52]:

```
# train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)
```

In [ ]:

## 2.2 Make Data Model Ready: encoding numerical, categorical features

In [53]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

### 2.2.1 Numerical features

1. teacher\_number\_of\_previously\_posted\_projects
2. price
3. quantity

#### 2.2.1.1 Teacher number of previously posted projects

In [54]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))

X_train_TPPP_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
X_cv_TPPP_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
X_test_TPPP_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))

print("After vectorizations")
print(X_train_TPPP_norm.shape, y_train.shape)
print(X_cv_TPPP_norm.shape, y_cv.shape)
print(X_test_TPPP_norm.shape, y_test.shape)
print("=="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====  
=====

In [55]:

```
print("Transpose of teacher number of previously posted projects")

X_train_TPPP_norm = X_train_TPPP_norm.transpose()
X_cv_TPPP_norm = X_cv_TPPP_norm.transpose()
X_test_TPPP_norm = X_test_TPPP_norm.transpose()

print("After transpose")
print(X_train_TPPP_norm.shape, y_train.shape)
print(X_cv_TPPP_norm.shape, y_cv.shape)
print(X_test_TPPP_norm.shape, y_test.shape)
print("=="*100)
```

Transpose of teacher number of previously posted projects

After transpose

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====  
=====

### 2.2.1.2 price

In [56]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(1,-1))

X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1))
X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====

In [57]:

```
print("Transpose of price")

X_train_price_norm = X_train_price_norm.transpose()
X_cv_price_norm = X_cv_price_norm.transpose()
X_test_price_norm = X_test_price_norm.transpose()

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

Transpose of price

After vectorizations

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====

### 2.2.1.3 quantity



In [58]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['quantity'].values.reshape(1,-1))

X_train_quantity_norm = normalizer.transform(X_train['quantity'].values.reshape(1,-1))
X_cv_quantity_norm = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
X_test_quantity_norm = normalizer.transform(X_test['quantity'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====

In [59]:

```
print("Transpose of Quantity")

X_train_quantity_norm = X_train_quantity_norm.transpose()
X_cv_quantity_norm = X_cv_quantity_norm.transpose()
X_test_quantity_norm = X_test_quantity_norm.transpose()

print("After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("="*100)
```

Transpose of Quantity

After vectorizations

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====

In [ ]:

## 2.2.2 Categorical Data

## Categorical Features for vectorization

1. Clean Categories
2. Clean Sub Categories
3. School State
4. Teacher Prefix
5. Project grade category

### 2.2.2.1 Clean Categories

In [60]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
                             binary=True)
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_CC_ohe = vectorizer.transform(X_train['clean_categories'].values)
X_cv_CC_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_CC_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_CC_ohe.shape, y_train.shape)
print(X_cv_CC_ohe.shape, y_cv.shape)
print(X_test_CC_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(22445, 9) (22445,)
(11055, 9) (11055,)
(16500, 9) (16500,)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
=====
=====
```

### 2.2.2.2 Clean Sub Categories

In [61]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_CSC_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
X_cv_CSC_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_CSC_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_CSC_ohe.shape, y_train.shape)
print(X_cv_CSC_ohe.shape, y_cv.shape)
print(X_test_CSC_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(22445, 30) (22445,)
(11055, 30) (11055,)
(16500, 30) (16500,)
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'Nutrition Education', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
=====
=====
```

### 2.2.2.3 School State

In [62]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_school_state_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_ohe = vectorizer.transform(X_train['school_state'].values)
X_cv_state_ohe = vectorizer.transform(X_cv['school_state'].values)
X_test_state_ohe = vectorizer.transform(X_test['school_state'].values)

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(22445, 51) (22445,)
(11055, 51) (11055,)
(16500, 51) (16500,)
['VT', 'WY', 'ND', 'MT', 'RI', 'SD', 'NE', 'DE', 'AK', 'NH', 'WV', 'ME',
 'HI', 'DC', 'NM', 'KS', 'IA', 'ID', 'AR', 'CO', 'MN', 'OR', 'KY', 'MS', 'N
V', 'MD', 'CT', 'TN', 'UT', 'AL', 'WI', 'VA', 'AZ', 'NJ', 'OK', 'WA', 'M
A', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'N
Y', 'TX', 'CA']
=====
=====
```

#### 2.2.2.4 Teacher prefix

In [63]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_teacher_prefix_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_teacher_prefix'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer.transform(X_train['clean_teacher_prefix'].values)
X_cv_teacher_ohe = vectorizer.transform(X_cv['clean_teacher_prefix'].values)
X_test_teacher_ohe = vectorizer.transform(X_test['clean_teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(22445, 5) (22445,)
(11055, 5) (11055,)
(16500, 5) (16500,)
['Dr', 'Teacher', 'Mr', 'Ms', 'Mrs']
=====
=====
```

### 2.2.2.5 Project Grade category

In [64]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_project_grade_category_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(X_train['clean_project_grade_category'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer.transform(X_train['clean_project_grade_category'].values)
X_cv_grade_ohe = vectorizer.transform(X_cv['clean_project_grade_category'].values)
X_test_grade_ohe = vectorizer.transform(X_test['clean_project_grade_category'].values)

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(22445, 4) (22445,)
(11055, 4) (11055,)
(16500, 4) (16500,)
['9-12', '6-8', '3-5', 'PreK-2']
=====
=====
```

In [65]:

```
data1['clean_project_grade_category'].unique()
```

Out[65]:

```
array(['PreK-2', '6-8', '3-5', '9-12'], dtype=object)
```

In [ ]:

## 2.3 Make Data Model Ready: encoding eassay, and project\_title

In [66]:

```
# please write all the code with proper documentation, and proper titles for each subsection  
# go through documentations and blogs before you start coding  
# first figure out what to do, and then think about how to do.  
# reading and understanding error messages will be very much helpfull in debugging your code  
# make sure you featurize train and test data separatly  
  
# when you plot any graph make sure you use  
    # a. Title, that describes your plot, this will be very helpful to the reader  
    # b. Legends if needed  
    # c. X-axis label  
    # d. Y-axis label
```

## **Ecoding Essay and Project title**

- 2.3.1 BOW
- 2.3.2 TFIDF
- 2.3.3 AVG W2V
- 2.3.4 TFIDF W2V

## **2.3.1 BOW Essays and Title**

### **2.3.1.1 BOW Essay**

In [67]:

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train
data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_bow = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("="*100)
```

```
(22445, 10) (22445,)
(11055, 10) (11055,)
(16500, 10) (16500,)
=====
=====
After vectorizations
(22445, 5000) (22445,)
(11055, 5000) (11055,)
(16500, 5000) (16500,)
=====
=====
```

### 2.3.1.2 BOW Title

In [68]:

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['preprocessed_titles'].values) # fit has to happen only on train
data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['preprocessed_titles'].values)
X_cv_title_bow = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_title_bow = vectorizer.transform(X_test['preprocessed_titles'].values)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
print("="*100)
```

```
(22445, 10) (22445,)
(11055, 10) (11055,)
(16500, 10) (16500,)
=====
=====
After vectorizations
(22445, 638) (22445,)
(11055, 638) (11055,)
(16500, 638) (16500,)
=====
=====
```

## 2.3.2 TF IDF Essay and Title

### 2.3.2.1 TF IDF Essay



In [69]:

```
from sklearn.feature_extraction.text import TfidfVectorizer

print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = TfidfVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
(22445, 10) (22445,)
(11055, 10) (11055,)
(16500, 10) (16500,)
=====
=====
After vectorizations
(22445, 5000) (22445,)
(11055, 5000) (11055,)
(16500, 5000) (16500,)
=====
=====
```

### 2.3.2.2 TF IDF Title

In [70]:

```
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = TfidfVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['preprocessed_titles'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['preprocessed_titles'].values)
X_cv_title_tfidf = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_title_tfidf = vectorizer.transform(X_test['preprocessed_titles'].values)

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
print("="*100)
```

```
(22445, 10) (22445,)
(11055, 10) (11055,)
(16500, 10) (16500,)
=====
=====
After vectorizations
(22445, 638) (22445,)
(11055, 638) (11055,)
(16500, 638) (16500,)
=====
=====
```

## 2.3.3 AVG W2V Essay and Title

### 2.3.3.1 AVG W2V Essay

In [71]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('../glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

In [72]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train.append(vector)

print(len(avg_w2v_vectors_train))
print(len(avg_w2v_vectors_train[0]))
print(avg_w2v_vectors_train[0])
```

100%|██████████| 22445/22445 [00:09<00:00, 2485.41it/s]

22445

300

[ -2.63049958e-03 1.01326269e-01 2.14073322e-02 -9.33893986e-02  
-5.45645636e-02 4.43942304e-02 -2.97693575e+00 -2.82121399e-03  
5.02780587e-02 -3.46241902e-02 1.15298055e-01 2.85400596e-02  
7.60204951e-02 -9.64316292e-02 2.87635594e-03 -7.81713159e-02  
4.84611175e-02 -4.28344874e-02 5.73115021e-02 -4.28735874e-02  
7.51226888e-02 4.47968671e-02 -7.57592440e-02 1.07475944e-03  
-3.44299245e-02 -9.48578336e-02 6.95911958e-02 -1.11368661e-01  
-8.70413070e-02 -6.27257992e-02 -1.93357450e-01 -1.32888049e-01  
5.02082909e-02 1.18718632e-01 -8.88458322e-02 5.97352154e-03  
3.87395301e-02 -2.98265385e-03 -7.92925895e-02 3.50873706e-03  
-1.79557546e-02 5.33112622e-02 -1.32022130e-01 -1.61365500e-01  
5.72652056e-02 -9.97882972e-02 6.26813497e-02 -6.84446937e-02  
6.64249385e-02 -4.16751976e-02 -3.85423917e-02 -9.57764315e-02  
-3.04711078e-02 1.24294119e-02 -3.18809413e-02 -6.96509587e-02  
1.89437336e-01 -1.58679441e-02 4.29350571e-02 6.32968121e-02  
1.84236965e-02 -3.07158071e-02 1.20443094e-01 -6.93545678e-02  
2.22726183e-02 8.39173119e-02 1.09189648e-01 -2.80926552e-02  
1.06219076e-01 -1.54424871e-01 -1.02668163e-01 -1.84745552e-02  
5.14139537e-02 -4.83479042e-02 -4.71362755e-02 -2.40004372e-01  
3.58018608e-02 -2.35372252e-02 2.76235329e-02 -7.94012657e-03  
1.22990964e-01 -4.78112559e-01 -5.15079168e-02 -6.99692091e-02  
-9.28900203e-02 2.47376141e-02 4.61389476e-02 -8.35468601e-02  
7.32190490e-02 6.75255987e-02 -1.08839940e-02 -3.84342573e-02  
5.31420000e-02 -9.45432168e-04 -5.44425105e-02 -7.30479070e-02  
-2.17616559e+00 -1.84292386e-02 9.21621245e-02 9.09319470e-02  
-1.13920134e-01 -7.82581145e-02 -2.84582867e-02 -6.04816224e-02  
9.18341462e-02 6.68893063e-02 9.39850133e-02 -2.61757520e-01  
4.23157176e-02 7.55481906e-02 -1.01372092e-01 6.87611189e-05  
3.65148706e-02 1.82554635e-01 -1.41275264e-02 1.34914116e-01  
-8.57449103e-02 3.31565895e-02 2.57250037e-02 -2.60603154e-02  
6.26168273e-02 5.53392014e-02 -1.98074938e-02 -4.41010825e-02  
8.55120225e-02 5.22427782e-02 4.02925927e-02 4.14764741e-02  
4.69772783e-03 1.31570976e-01 -7.19044945e-02 2.68019280e-02  
8.29310266e-02 -6.21744643e-02 8.12752944e-02 -1.85843515e-01  
1.66466963e-01 -5.77297820e-02 5.94340497e-02 3.36576082e-01  
6.58336965e-02 5.44033406e-02 9.81260091e-02 -7.81550601e-02  
-6.36602657e-02 -1.31027225e-01 -1.53789580e-03 -2.03464643e-02  
1.96997541e-01 -9.19918881e-04 1.38182820e-02 3.50967979e-02  
2.31736685e-02 -1.95459057e-02 -1.22509785e-02 1.53267507e-01  
5.83722290e-02 7.96085175e-02 -9.64473014e-02 6.71618245e-02  
4.65918217e-02 1.04959768e-02 2.45279643e-02 -6.42463399e-02  
-2.92751818e-03 1.20166814e-01 -9.44673724e-02 -5.43407112e-02  
1.03853131e-01 -6.29492175e-02 -1.68300406e-02 5.69743839e-02  
-6.13026876e-02 -1.35731406e-01 -1.24452251e-02 -4.04451678e-02  
-5.43396503e-02 1.44625664e-02 -2.09576669e-01 -4.36322869e-02  
3.93598888e-03 2.68220153e-01 -6.05159734e-02 9.66184685e-03  
-2.30277126e-02 -5.15654343e-02 -6.79129021e-02 -4.90027552e-02  
6.92583895e-02 1.22248101e-01 8.25541881e-02 -6.86395874e-02  
-1.00052119e-01 -3.70363517e-03 6.15272308e-03 -2.14001888e-03  
-3.63745720e-02 -1.39407699e-02 6.10691069e-02 2.15775697e-02  
7.09274098e-02 9.66571329e-03 -5.75933336e-02 1.08145894e-01  
-7.58177621e-02 3.04165000e-02 3.72449077e-02 -9.77580517e-02  
2.12848118e-01 4.80300315e-02 2.84020729e-02 -4.10073755e-02  
-4.07838587e-03 -1.48846647e-01 -5.15218741e-02 -7.85891399e-02  
-4.62468189e-02 -6.32938552e-02 1.01920534e-01 -3.57182154e-02  
-8.36898629e-02 -1.87375051e-01 -1.09338667e-01 -8.91805315e-02  
-1.86425181e+00 -3.60894287e-02 -2.85615020e-02 -4.05551524e-02  
1.88361049e-03 -1.94255806e-01 -1.14249224e-01 -9.42727287e-02  
-4.06751329e-03 -5.97658042e-02 -4.63775322e-02 -1.47818077e-02

```

9.91436965e-02 -2.33500119e-02 -8.16509021e-03 1.32395959e-01
-5.83603357e-03 8.19939112e-02 -1.73516856e-01 1.59771244e-01
-5.50540783e-02 -1.02703245e-02 -1.21353357e-01 3.31691420e-03
-1.82844990e-03 -2.05141373e-02 2.57810832e-03 1.85551733e-01
-2.28334322e-02 6.79984196e-03 1.67078119e-01 2.65139522e-02
1.04275126e-01 -3.83207063e-02 2.55472134e-01 -4.74595322e-02
1.95538483e-02 3.25160986e-02 -5.55684741e-02 1.42274029e-02
-1.30398126e-02 -5.00695916e-02 -5.83655391e-02 -2.22274399e-02
1.17459176e-01 7.40340594e-02 -1.29300783e-01 2.93711266e-02
-1.90066265e-01 5.39769483e-02 -1.22539214e-01 7.00169105e-02
1.19876752e-01 2.31548322e-02 -3.93150587e-02 2.19802646e-02
8.05595664e-02 -8.21964965e-03 3.36027636e-02 1.39029828e-01
2.18894266e-03 1.52799821e-01 1.21472588e-02 3.02761427e-02
-6.73986874e-02 -3.70049014e-02 3.63809091e-02 -1.92621770e-02
2.41643007e-02 -6.47919007e-02 9.65560566e-02 -3.91337093e-02
-2.41569650e-02 1.86488885e-01 9.19699130e-02 5.08127490e-02]

```

In [73]:

```

avg_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv.append(vector)

```

100%|██████████| 11055/11055 [00:04<00:00, 2585.56it/s]

In [74]:

```

avg_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_test.append(vector)

```

100%|██████████| 16500/16500 [00:06<00:00, 2416.21it/s]

### 2.3.3.2 AVG W2V Title

In [75]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_train_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train_title.append(vector)

print(len(avg_w2v_vectors_train_title))
print(len(avg_w2v_vectors_train_title[0]))
print(avg_w2v_vectors_train_title[0])
```

100%|██████████| 22445/22445 [00:00<00:00, 47357.54it/s]



22445

300

```
[ 1.17780600e-01  7.15982500e-02 -2.67117500e-01 -1.93448750e-01
  2.86788000e-01  1.48055000e-01 -4.11172500e+00  4.23460000e-02
  5.68140000e-02 -4.92536250e-01  2.11305750e-01 -1.12011750e-01
  8.86917500e-02  1.01014250e-01 -7.60000000e-02  1.52705000e-01
 -2.46551500e-01 -2.53445000e-02  2.81269000e-01  6.83382500e-02
  3.94912500e-02  7.51310000e-02  1.86274000e-02  2.63997500e-02
 -4.98172500e-02 -1.51413750e-01 -1.90702000e-01 -2.74921750e-01
 -1.45461500e-02 -4.45915000e-02 -2.77381750e-01  6.23530000e-02
  2.97665000e-01  6.91200000e-02 -4.79922500e-02  5.93025000e-02
 -1.06390000e-01 -2.30610000e-01  1.67173750e-01 -3.26642500e-02
 -9.45325000e-02  1.57470250e-01 -1.92692250e-01 -1.33901750e-01
  2.44445000e-02  1.05870750e-01 -3.08360000e-02  7.11475000e-02
 -2.06392500e-01 -2.42745000e-01 -1.00461750e-01 -5.12222500e-02
 -1.15115500e-01 -5.33822500e-02 -1.11373500e-01 -3.88745000e-01
  4.37018000e-01 -6.12987500e-02 -2.75342500e-02  1.60496500e-01
  2.46222500e-01  1.85017500e-01  2.33500000e-02  1.05022500e-03
 -4.92655000e-02  7.93835000e-02  1.96346500e-01  2.34230000e-01
  8.51237500e-02 -1.15237500e-02 -2.19271325e-01  8.18665000e-02
  1.35597500e-01  8.20340000e-02 -2.81492500e-01 -1.80840000e-01
  2.39938500e-01  2.08652250e-01  1.17468500e-01 -2.27622500e-02
 -1.01198750e-01 -7.12545000e-01 -2.06820000e-01 -1.05817850e-01
  2.16367000e-01 -7.40875000e-02  2.48649750e-01 -2.00483000e-01
 -6.59552500e-02  1.54570000e-02 -1.75320000e-01 -6.55700000e-03
 -8.87790000e-02  1.97055600e-01  3.16150000e-02 -2.42640000e-01
 -2.38502500e+00 -7.06150000e-03  1.25750000e-01  2.73495000e-02
 -3.46662500e-01  2.17787750e-01  1.38550500e-01  1.04405750e-01
 -2.49292500e-02 -5.29277500e-02  2.07531500e-01 -2.09304500e-01
 -5.81787500e-02 -1.59612750e-01 -2.72469250e-01 -1.89234000e-02
 -8.40590925e-02  8.62552500e-02 -3.43032500e-01  2.64025000e-01
  2.83316250e-01 -4.75867500e-02 -1.07607750e-01 -4.24662500e-02
  6.80600000e-02  4.73297500e-02  7.52400000e-02 -8.20575000e-02
  9.56375000e-03 -1.99825000e-02  9.79906850e-02 -7.13017500e-02
 -7.77857500e-02  2.79362250e-01 -2.47410000e-01 -9.36655000e-02
  8.23455000e-02  3.44450000e-03  1.47108750e-01 -3.80358250e-01
  1.72950000e-01  2.62912500e-01 -1.36275000e-02  7.37382500e-01
  1.66877325e-01  6.87625000e-02  3.36120000e-01 -7.47575000e-02
  1.33837575e-01 -4.28388250e-01  1.25790350e-01 -4.13965000e-01
  2.42095000e-01  1.63443750e-01  8.62320000e-02  1.02269250e-01
 -1.12372750e-01 -1.25827500e-01  4.26607500e-02 -2.48510000e-01
 -7.29107500e-02  1.53974000e-01 -3.41075000e-01 -7.66792500e-02
  2.35622500e-01  9.70135000e-02 -3.34010000e-01 -1.32977750e-01
  3.17492500e-01  3.24510750e-01  3.52577500e-02  2.20924250e-01
  3.15532500e-01 -3.63862500e-01  5.87250000e-04 -1.67925000e-01
 -1.13573320e-01 -4.48000000e-02 -2.49375000e-01  1.67225000e-01
 -7.36630000e-02  6.03365000e-02 -1.36600000e-02 -3.54292500e-01
  3.76989250e-02  1.78502500e-01 -2.25468750e-01  6.29670000e-02
 -8.41142500e-02 -3.40182500e-02 -1.59786000e-01 -2.28472500e-01
  1.58404750e-01  1.34965500e-01  8.06125000e-02 -8.46107000e-02
  1.84056250e-01 -3.51887500e-02  2.40275000e-03 -1.66521250e-01
  9.38625000e-03 -2.91807500e-01  1.21067500e-01 -3.14824250e-01
  3.45400000e-02  2.86350000e-03 -2.05664250e-01  2.50352500e-01
  3.47165000e-02  2.29065000e-02 -8.68782500e-02 -2.13417725e-01
  3.95921750e-01 -5.43172500e-02  4.55385000e-02  1.10769250e-01
 -1.52537250e-01 -2.69349250e-01  6.84335000e-02 -3.02767500e-02
  7.93077500e-02  1.38307500e-02  1.76405000e-02  1.18704250e-01
 -4.80046000e-02 -3.18817500e-01  1.00539150e-01 -2.37682750e-01
 -2.96170000e+00  3.06062500e-01  9.11795000e-02  1.45079250e-01
 -1.48237500e-01 -2.39912250e-01 -1.01294000e-01  2.15333500e-01
  2.56685450e-01  3.88532500e-02 -1.27020250e-01  1.02536155e-01
```

```

7.68747500e-02  1.09518500e-01 -1.74922500e-01  7.35931500e-02
1.54175000e-02  1.09654500e-01 -3.99497500e-01  8.98642500e-02
1.58950000e-02  -2.19881500e-01 -1.96467500e-02  1.26125000e-02
-1.84446250e-02  1.06317250e-01  5.49365000e-02  2.59225000e-02
6.59708250e-02  -8.18650000e-02  9.20415000e-02  9.92500000e-03
1.46294250e-01  -9.50930000e-02 -1.62382500e-01  2.38265000e-01
-2.94755000e-01  4.18597500e-02 -2.21197500e-01 -2.03560000e-01
-9.88950000e-04  -6.63242500e-02 -1.19730475e-01  2.30615000e-01
8.21365000e-02  1.07446400e-01  1.34780000e-01  1.06406125e-01
-4.01485000e-01  -6.95065000e-02 -2.40685600e-01 -2.27020000e-02
2.12290000e-01  1.52085000e-01  2.02850000e-02 -2.42348750e-01
4.63105000e-01  1.05055250e-01  2.98950000e-03  2.52230000e-02
-2.66121250e-02  4.17792500e-02 -3.01722500e-01 -1.54871750e-01
-1.65280000e-02  -3.47767500e-02  1.60712500e-02 -2.66458000e-02
5.93525000e-02  -1.01528000e-01  1.35918850e-01  1.04576500e-01
-2.57370000e-02  1.97851250e-01 -4.17980750e-02  5.56925000e-02]

```

In [76]:

```

avg_w2v_vectors_cv_title = []; # the avg-w2v for each sentence/review is stored in this
list
for sentence in tqdm(X_cv['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv_title.append(vector)

```

100%|██████████| 11055/11055 [00:00<00:00, 40909.26it/s]

In [77]:

```

avg_w2v_vectors_test_title = []; # the avg-w2v for each sentence/review is stored in th
is list
for sentence in tqdm(X_test['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_test_title.append(vector)

```

100%|██████████| 16500/16500 [00:00<00:00, 42115.28it/s]

## 2.3.4 TF IDF W2V Essay and Title

### 2.3.4.1 TF IDF W2V Essay

In [78]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_essays'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
```

In [79]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this
list
for sentence in tqdm(X_train['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_train.append(vector)

print(len(tfidf_w2v_vectors_train))
print(len(tfidf_w2v_vectors_train[0]))
```

100%|██████████| 22445/22445 [01:03<00:00, 354.00it/s]

22445

300

In [80]:

```
tfidf_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_cv.append(vector)
```

100%|██████████| 11055/11055 [00:32<00:00, 343.49it/s]

In [81]:

```
tfidf_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_test.append(vector)
```

100%|██████████| 16500/16500 [00:47<00:00, 347.07it/s]

### 2.3.4.2 TF IDF W2V Title

In [82]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_titles'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
```

In [83]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_train_title = []; # the avg-w2v for each sentence/review is stored in
this list
for sentence in tqdm(X_train['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_train_title.append(vector)

print(len(tfidf_w2v_vectors_train_title))
print(len(tfidf_w2v_vectors_train_title[0]))
```

100%|██████████| 22445/22445 [00:00<00:00, 23157.58it/s]

22445

300

In [84]:

```
tfidf_w2v_vectors_cv_title = []; # the avg-w2v for each sentence/review is stored in th
is list
for sentence in tqdm(X_cv['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_cv_title.append(vector)
```

100%|██████████| 11055/11055 [00:00<00:00, 21299.91it/s]

In [85]:

```
tfidf_w2v_vectors_test_title = []; # the avg-w2v for each sentence/review is stored in
this list
for sentence in tqdm(X_test['preprocessed_titles'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sen
            tence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # ge
            tting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_test_title.append(vector)
```

100%|██████████| 16500/16500 [00:00<00:00, 23749.52it/s]

In [ ]:

## Concatinating all the features

### 1. SET 1 BOW

In [86]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_BOW = hstack((X_train_essay_bow, X_train_title_bow, X_train_state_ohe, X_train_tea
cher_ohe, X_train_grade_ohe, X_train_CSC_ohe, X_train_CC_ohe, X_train_price_norm, X_tra
in_quantity_norm, X_train_TPPP_norm)).tocsr()
X_cr_BOW = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_state_ohe, X_cv_teacher_ohe, X_
cv_grade_ohe, X_cv_CSC_ohe, X_cv_CC_ohe, X_cv_price_norm, X_cv_quantity_norm, X_cv_TPPP
_norm)).tocsr()
X_te_BOW = hstack((X_test_essay_bow, X_test_title_bow, X_test_state_ohe, X_test_teacher
_ohe, X_test_grade_ohe, X_test_CSC_ohe, X_test_CC_ohe, X_test_price_norm, X_test quanti
ty_norm, X_test_TPPP_norm)).tocsr()

print("Final Data matrix")
print(X_tr_BOW.shape, y_train.shape)
print(X_cr_BOW.shape, y_cv.shape)
print(X_te_BOW.shape, y_test.shape)
print("=*100)
```

Final Data matrix

```
(22445, 5740) (22445,)
(11055, 5740) (11055,)
(16500, 5740) (16500,)
```

```
=====
=====
```

## 2. SET 2 TF IDF

In [87]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_TFIDF = hstack((X_train_essay_tfidf, X_train_title_tfidf, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_CSC_ohe, X_train_CC_ohe, X_train_price_norm, X_train_quantity_norm, X_train_TPPP_norm)).tocsr()
X_cr_TFIDF = hstack((X_cv_essay_tfidf, X_cv_title_tfidf, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_CSC_ohe, X_cv_CC_ohe, X_cv_price_norm, X_cv_quantity_norm, X_cv_TPPP_norm)).tocsr()
X_te_TFIDF = hstack((X_test_essay_tfidf, X_test_title_tfidf, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_CSC_ohe, X_test_CC_ohe, X_test_price_norm, X_test_quantity_norm, X_test_TPPP_norm)).tocsr()

print("Final Data matrix")
print(X_tr_TFIDF.shape, y_train.shape)
print(X_cr_TFIDF.shape, y_cv.shape)
print(X_te_TFIDF.shape, y_test.shape)
print("=="*100)
```

```
Final Data matrix
(22445, 5740) (22445,)
(11055, 5740) (11055,)
(16500, 5740) (16500,)
=====
=====
```

## 3. SET 3 AVG W2V

In [88]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_AVG_W2V = hstack((avg_w2v_vectors_train, avg_w2v_vectors_train_title, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_CSC_ohe, X_train_CC_ohe, X_train_price_norm, X_train_quantity_norm, X_train_TPPP_norm)).tocsr()
X_cr_AVG_W2V = hstack((avg_w2v_vectors_cv, avg_w2v_vectors_cv_title, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_CSC_ohe, X_cv_CC_ohe, X_cv_price_norm, X_cv_quantity_norm, X_cv_TPPP_norm)).tocsr()
X_te_AVG_W2V = hstack((avg_w2v_vectors_test, avg_w2v_vectors_test_title, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_CSC_ohe, X_test_CC_ohe, X_test_price_norm, X_test_quantity_norm, X_test_TPPP_norm)).tocsr()

print("Final Data matrix")
print(X_tr_AVG_W2V.shape, y_train.shape)
print(X_cr_AVG_W2V.shape, y_cv.shape)
print(X_te_AVG_W2V.shape, y_test.shape)
print("=="*100)
```

```
Final Data matrix
(22445, 702) (22445,)
(11055, 702) (11055,)
(16500, 702) (16500,)
=====
=====
```

#### 4. SET 4 TF IDF W2V

In [89]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_TFIDF_W2V = hstack((tfidf_w2v_vectors_train, tfidf_w2v_vectors_train_title, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_CSC_ohe, X_train_CC_ohe, X_train_price_norm, X_train_quantity_norm, X_train_TPPP_norm)).tocsr()
X_cr_TFIDF_W2V = hstack((tfidf_w2v_vectors_cv, tfidf_w2v_vectors_cv_title, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_CSC_ohe, X_cv_CC_ohe, X_cv_price_norm, X_cv_quantity_norm, X_cv_TPPP_norm)).tocsr()
X_te_TFIDF_W2V = hstack((tfidf_w2v_vectors_test, tfidf_w2v_vectors_test_title, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_CSC_ohe, X_test_CC_ohe, X_test_price_norm, X_test_quantity_norm, X_test_TPPP_norm)).tocsr()

print("Final Data matrix")
print(X_tr_TFIDF_W2V.shape, y_train.shape)
print(X_cr_TFIDF_W2V.shape, y_cv.shape)
print(X_te_TFIDF_W2V.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(22445, 702) (22445,)
(11055, 702) (11055,)
(16500, 702) (16500,)
=====
=====
```

In [ ]:

## 2.4 Applying Logistic Regression on different kind of featurization as mentioned in the instructions

Apply Logistic Regression on different kind of featurization as mentioned in the instructions  
For Every model that you work on make sure you do the step 2 and step 3 of instructions

In [90]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

### 2.4.1 Applying logistic regression on BOW, SET 1



In [91]:

```
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
#from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence va
lues, or non-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

clf = LogisticRegression(class_weight='balanced');
parameters={'C':[10**-4, 10**-3,10**-2,1,10,100,1000,500,1000,10000]}
sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc',return_train_score=True)
sd.fit(X_tr_BOW, y_train);

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

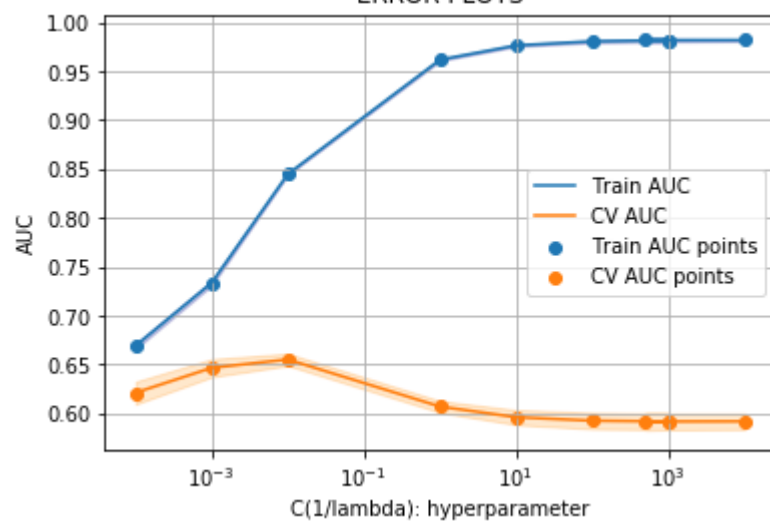
plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc + train_auc_
std,alpha=0.2,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.
2,color='darkorange')

plt.scatter(parameters['C'], train_auc, label='Train AUC points')
plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
plt.xscale('log')

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

# ERROR PLOTS



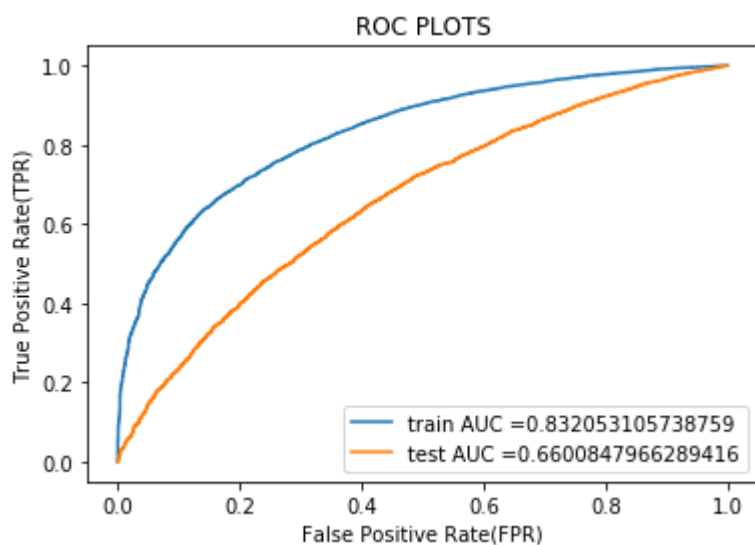
In [92]:

```
##Fitting Model to Hyper-Parameter Curve
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

neigh = LogisticRegression(C=10**-2,class_weight='balanced');
neigh.fit(X_tr_BOW ,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_tr_BOW)[:,-1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_te_BOW)[:,-1])

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()
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")
plt.title("ROC PLOTS")
plt.show()
```



In [93]:

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

```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_tr_BOW )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
```

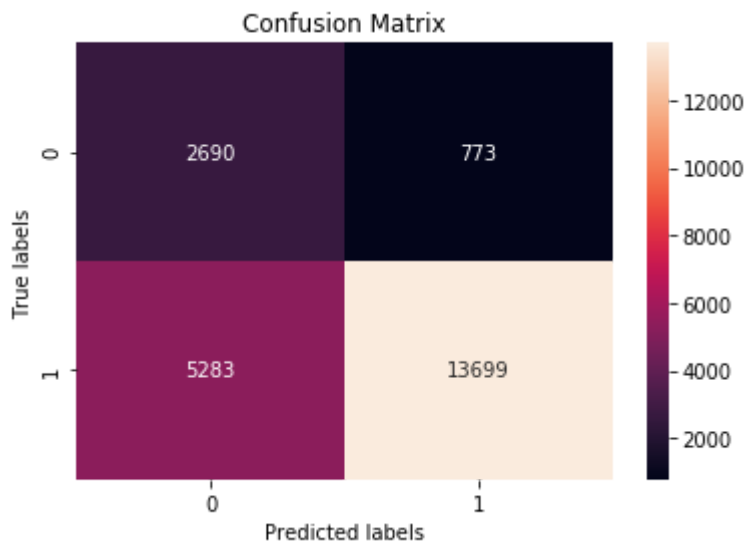
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```

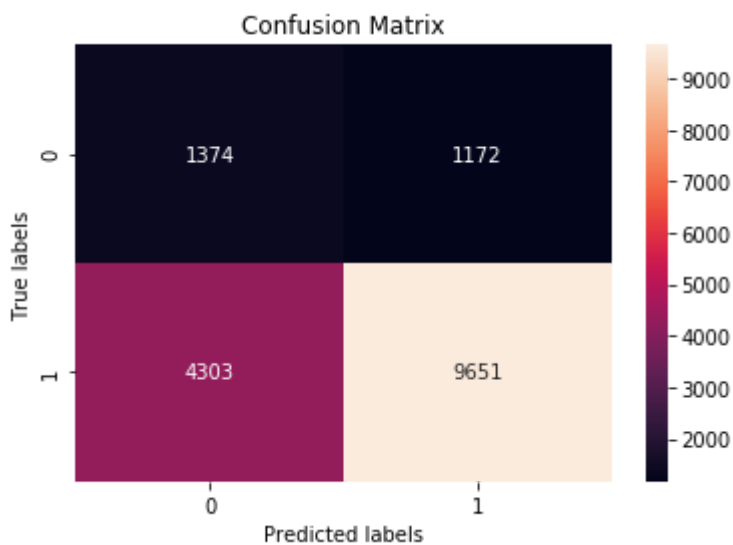


In [94]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_te_BOW )), annot=True, ax = ax,fmt
='g'); #annot=True to annotate cells

# Labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b
usiness']);
```



## 2.4.2 Applying Logistic regression on TFIDF, SET 2

In [95]:

```
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
#from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence va
lues, or non-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

clf = LogisticRegression(class_weight='balanced');
parameters = {'C':[10**-4, 10**-3, 10**-2, 1, 10, 100, 1000, 500, 1000, 10000]}
sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc', return_train_score=True)
sd.fit(X_tr_TFIDF, y_train);

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

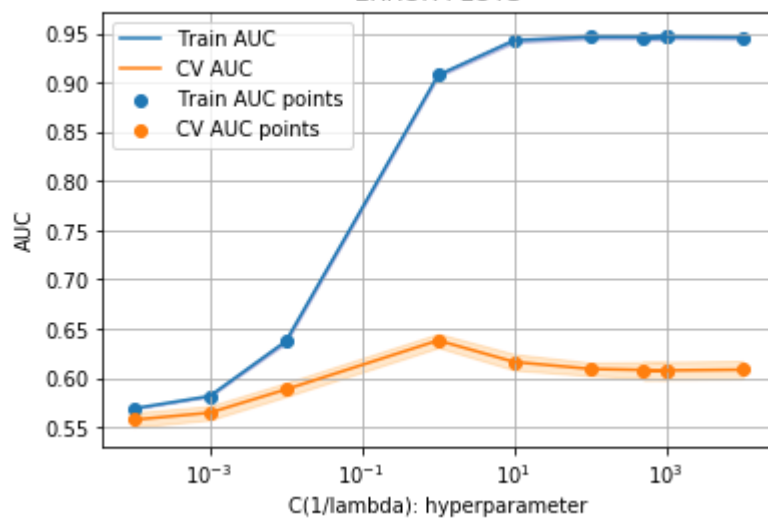
plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], train_auc - train_auc_std, train_auc + train_auc_std, alpha=0.2, color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.2, color='darkorange')

plt.scatter(parameters['C'], train_auc, label='Train AUC points')
plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
plt.xscale('log')

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

ERROR PLOTS



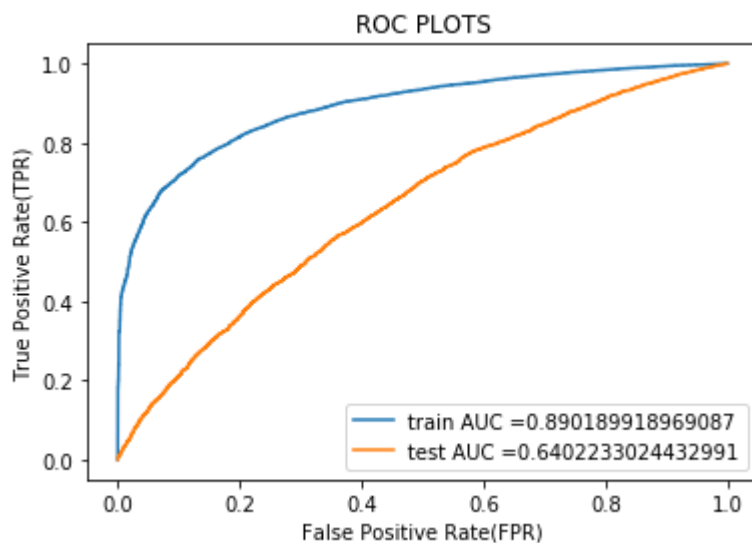
In [96]:

```
##Fitting Model to Hyper-Parameter Curve
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

neigh = LogisticRegression(C=10**0, class_weight='balanced');
neigh.fit(X_tr_TFIDF , y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_tr_TFIDF)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_te_TFIDF)[:,1])

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()
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")
plt.title("ROC PLOTS")
plt.show()
```





In [97]:

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

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
```

```
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_tr_TFIDF )), annot=True, ax = ax,  
fmt='g'); #annot=True to annotate cells
```

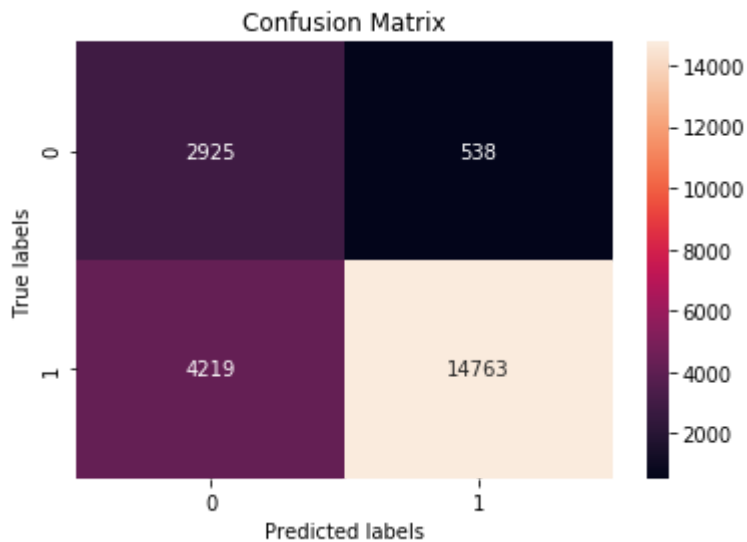
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b  
usiness']);
```

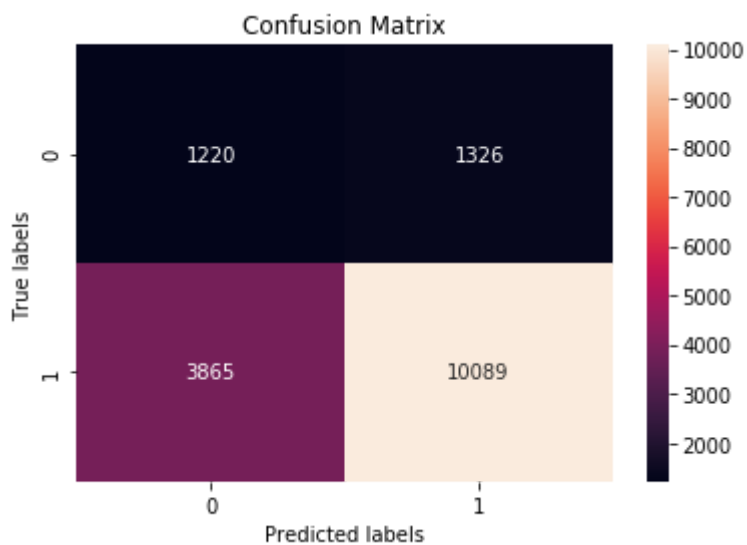


In [98]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_te_TFIDF )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells

# Labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```



In [ ]:

### 2.4.3 Applying Logistic regression on AVG W2V, SET 3

In [99]:

```
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
#from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence va
lues, or non-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

clf = LogisticRegression(class_weight='balanced');
parameters = {'C':[10**-4, 10**-3, 10**-2, 1, 10, 100, 1000, 500, 1000, 10000]}
sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc', return_train_score=True)
sd.fit(X_tr_AVG_W2V, y_train);

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

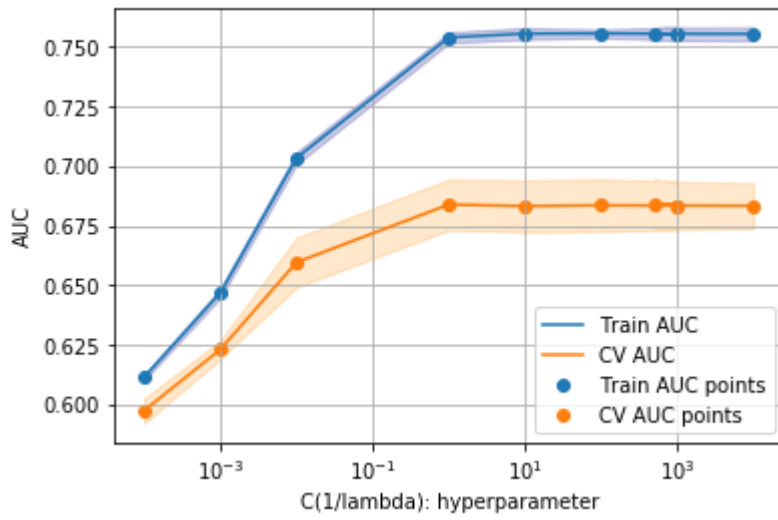
plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], train_auc - train_auc_std, train_auc + train_auc_std, alpha=0.2, color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.2, color='darkorange')

plt.scatter(parameters['C'], train_auc, label='Train AUC points')
plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
plt.xscale('log')

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

ERROR PLOTS



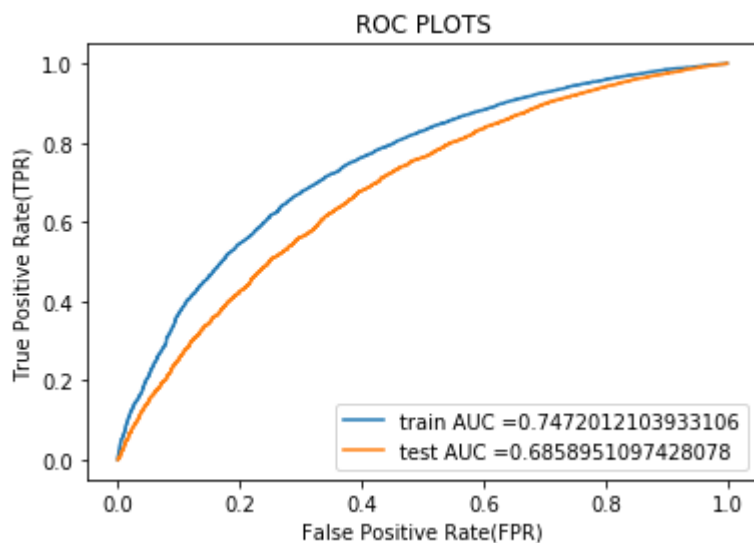
In [100]:

```
##Fitting Model to Hyper-Parameter Curve
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

neigh = LogisticRegression(C=10**0, class_weight='balanced');
neigh.fit(X_tr_AVG_W2V , y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_tr_AVG_W2V)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_te_AVG_W2V)[:,1])

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()
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")
plt.title("ROC PLOTS")
plt.show()
```



In [101]:

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

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
```

```
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_tr_AVG_W2V )), annot=True, ax = ax,fmt='g'); #annot=True to annotate cells
```

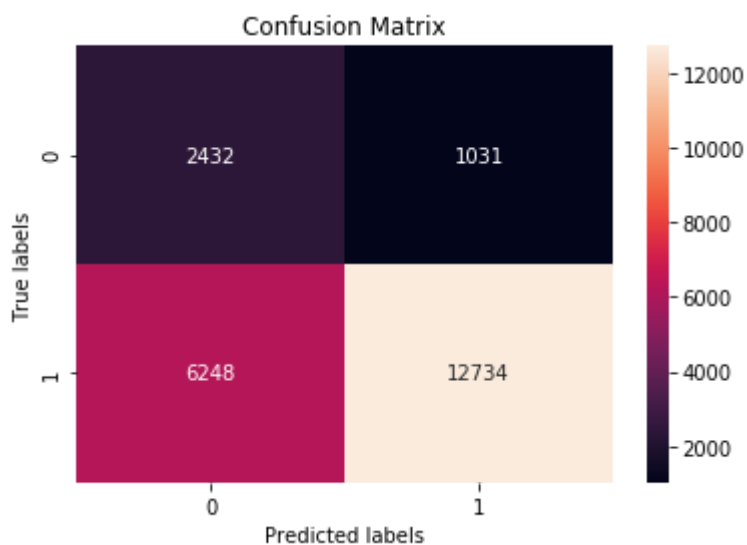
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```



In [102]:

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

```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_te_AVG_W2V )), annot=True, ax = ax
,fmt='g'); #annot=True to annotate cells
```

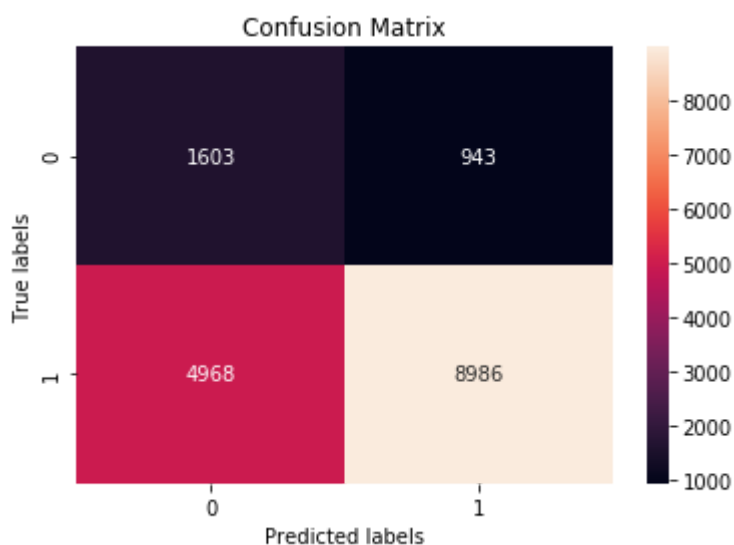
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b
usiness']);
```



In [ ]:

In [ ]:

## 2.4.4 Applying Logistic regression on TFIDF W2v, SET 4

In [103]:

```
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
#from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence va
lues, or non-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.
"""

clf = LogisticRegression(class_weight='balanced');
parameters = {'C':[10**-4, 10**-3, 10**-2, 1, 10, 100, 1000, 500, 1000, 10000]}
sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc', return_train_score=True)
sd.fit(X_tr_TFIDF_W2V, y_train);

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

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], train_auc - train_auc_std, train_auc + train_auc_
std, alpha=0.2, color='darkblue')

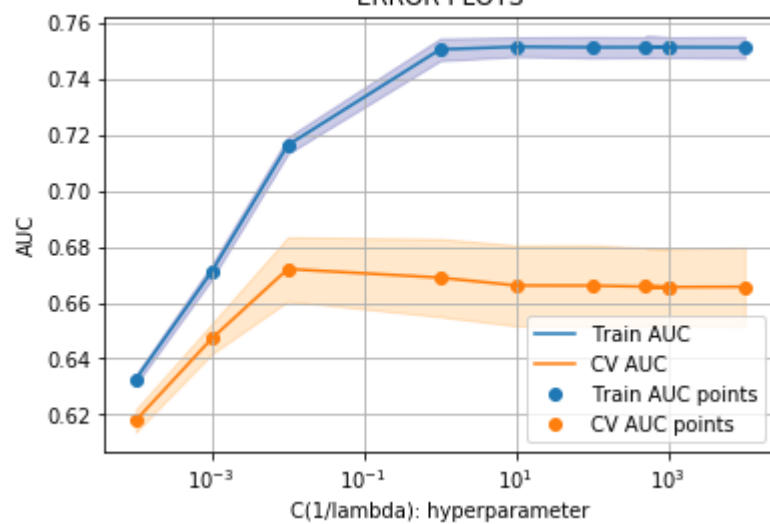
plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.
2, color='darkorange')

plt.scatter(parameters['C'], train_auc, label='Train AUC points')
plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
plt.xscale('log')

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



# ERROR PLOTS



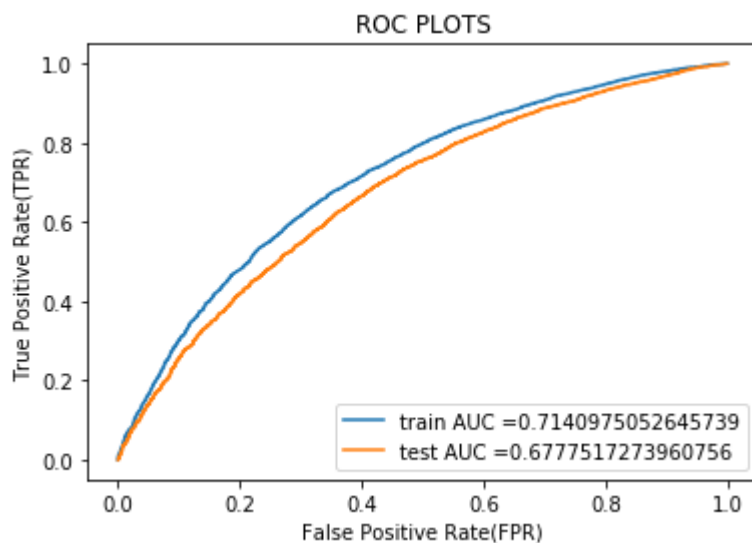
In [104]:

```
##Fitting Model to Hyper-Parameter Curve
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

neigh = LogisticRegression(C=10**-2,class_weight='balanced');
neigh.fit(X_tr_TFIDF_W2V ,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_tr_TFIDF_W2V)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_te_TFIDF_W2V)[:,1])

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()
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")
plt.title("ROC PLOTS")
plt.show()
```



In [105]:

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

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
```

```
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_tr_TFIDF_W2V )), annot=True, ax =  
ax,fmt='g'); #annot=True to annotate cells
```

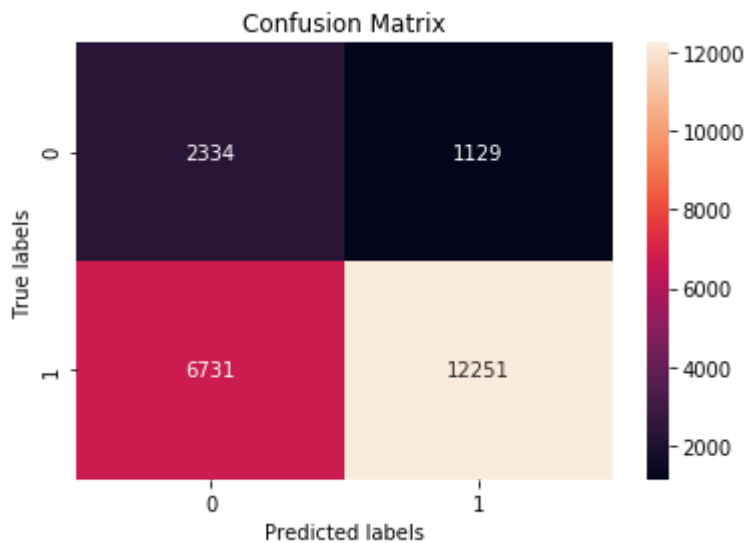
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b  
usiness']);
```

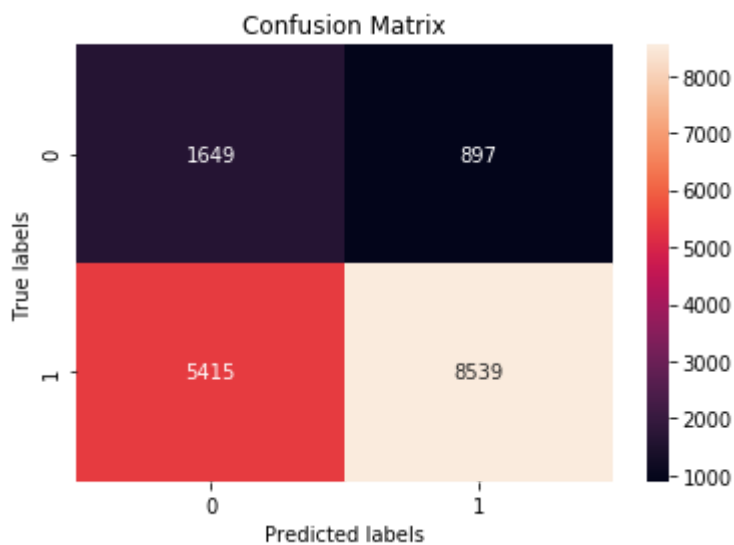


In [106]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_te_TFIDF_W2V )), annot=True, ax =
ax,fmt='g'); #annot=True to annotate cells

# Labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b
usiness']);
```



In [ ]:

## 2.5 Logistic Regression with added Features `Set 5`

In [107]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

In [108]:

```
data1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 50000 entries, 0 to 49999
```

```
Data columns (total 11 columns):
```

#	Column	Non-Null Count	Dtype
0	school_state	50000 non-null	object
1	teacher_number_of_previously_posted_projects	50000 non-null	int64
2	project_is_approved	50000 non-null	int64
3	clean_categories	50000 non-null	object
4	clean_subcategories	50000 non-null	object
5	clean_teacher_prefix	50000 non-null	object
6	clean_project_grade_category	50000 non-null	object
7	preprocessed_essays	50000 non-null	object
8	preprocessed_titles	50000 non-null	object
9	price	50000 non-null	float64
10	quantity	50000 non-null	int64

```
dtypes: float64(1), int64(3), object(7)
```

```
memory usage: 4.6+ MB
```

In [109]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# import nltk
nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

neg = []
pos = []
neu = []
compound = []

for a in tqdm(data1["preprocessed_essays"]) :
    b = sid.polarity_scores(a)['neg']
    c = sid.polarity_scores(a)['pos']
    d = sid.polarity_scores(a)['neu']
    e = sid.polarity_scores(a)['compound']
    neg.append(b)
    pos.append(c)
    neu.append(d)
    compound.append(e)
```

```
[nltk_data] Downloading package vader_lexicon to C:\Users\KALYAN
[nltk_data]       SRINIVAS\AppData\Roaming\nltk_data...
[nltk_data]   Package vader_lexicon is already up-to-date!
100%|██████████| 50000/50000 [08:20<00:00, 99.84it/s]
```

In [ ]:

In [110]:

```
data1["pos"] = pos
data1["neg"] = neg
data1["neu"] = neu
data1["compound"] = compound
```

## Essays and title word count

In [111]:

```
essay_word_count = []
for ess in data1["preprocessed_essays"] :
    c = len(ess.split())
    essay_word_count.append(c)
```

In [112]:

```
data1['essay_word_count'] = essay_word_count
```

In [113]:

```
title_word_count = []
for ess in data1["preprocessed_titles"] :
    c = len(ess.split())
    title_word_count.append(c)
```

In [114]:

```
data1['title_word_count'] = title_word_count
```

In [115]:

```
y = data1['project_is_approved'].values
X = data1.drop(['project_is_approved'], axis=1)
X.head(1)
```

Out[115]:

	school_state	teacher_number_of_previously_posted_projects	clean_categories	clean_subc
0	IN	0	Literacy_Language	ES



In [116]:

```
# train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)
```

## pos vectorization

In [117]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['pos'].values.reshape(1,-1))

X_train_pos_norm = normalizer.transform(X_train['pos'].values.reshape(1,-1))
X_cv_pos_norm = normalizer.transform(X_cv['pos'].values.reshape(1,-1))
X_test_pos_norm = normalizer.transform(X_test['pos'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_pos_norm.shape, y_train.shape)
print(X_cv_pos_norm.shape, y_cv.shape)
print(X_test_pos_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

```
=====
=====
```

In [118]:

```
print("Transpose of pos")

X_train_pos_norm = X_train_pos_norm.transpose()
X_cv_pos_norm = X_cv_pos_norm.transpose()
X_test_pos_norm = X_test_pos_norm.transpose()

print("After vectorizations")
print(X_train_pos_norm.shape, y_train.shape)
print(X_cv_pos_norm.shape, y_cv.shape)
print(X_test_pos_norm.shape, y_test.shape)
print("="*100)
```

Transpose of pos

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

```
=====
=====
```

## neg vectorization



In [119]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['neg'].values.reshape(1,-1))

X_train_neg_norm = normalizer.transform(X_train['neg'].values.reshape(1,-1))
X_cv_neg_norm = normalizer.transform(X_cv['neg'].values.reshape(1,-1))
X_test_neg_norm = normalizer.transform(X_test['neg'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_neg_norm.shape, y_train.shape)
print(X_cv_neg_norm.shape, y_cv.shape)
print(X_test_neg_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====

In [120]:

```
print("Transpose of neg")

X_train_neg_norm = X_train_neg_norm.transpose()
X_cv_neg_norm = X_cv_neg_norm.transpose()
X_test_neg_norm = X_test_neg_norm.transpose()

print("After vectorizations")
print(X_train_neg_norm.shape, y_train.shape)
print(X_cv_neg_norm.shape, y_cv.shape)
print(X_test_neg_norm.shape, y_test.shape)
print("="*100)
```

Transpose of neg

After vectorizations

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====

**neutral vectorization**

In [121]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['neu'].values.reshape(1,-1))

X_train_neu_norm = normalizer.transform(X_train['neu'].values.reshape(1,-1))
X_cv_neu_norm = normalizer.transform(X_cv['neu'].values.reshape(1,-1))
X_test_neu_norm = normalizer.transform(X_test['neu'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_neu_norm.shape, y_train.shape)
print(X_cv_neu_norm.shape, y_cv.shape)
print(X_test_neu_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

```
=====
=====
```

In [122]:

```
print("Transpose of neutral")

X_train_neu_norm = X_train_neu_norm.transpose()
X_cv_neu_norm = X_cv_neu_norm.transpose()
X_test_neu_norm = X_test_neu_norm.transpose()

print("After vectorizations")
print(X_train_neu_norm.shape, y_train.shape)
print(X_cv_neu_norm.shape, y_cv.shape)
print(X_test_neu_norm.shape, y_test.shape)
print("="*100)
```

Transpose of neutral

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

```
=====
=====
```

## compound vectorization

In [123]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['compound'].values.reshape(1,-1))

X_train_compound_norm = normalizer.transform(X_train['compound'].values.reshape(1,-1))
X_cv_compound_norm = normalizer.transform(X_cv['compound'].values.reshape(1,-1))
X_test_compound_norm = normalizer.transform(X_test['compound'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_compound_norm.shape, y_train.shape)
print(X_cv_compound_norm.shape, y_cv.shape)
print(X_test_compound_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====

In [124]:

```
print("Transpose of compound")

X_train_compound_norm = X_train_compound_norm.transpose()
X_cv_compound_norm = X_cv_compound_norm.transpose()
X_test_compound_norm = X_test_compound_norm.transpose()

print("After vectorizations")
print(X_train_compound_norm.shape, y_train.shape)
print(X_cv_compound_norm.shape, y_cv.shape)
print(X_test_compound_norm.shape, y_test.shape)
print("="*100)
```

Transpose of compound

After vectorizations

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====

**essay word count vectorization**

In [125]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['essay_word_count'].values.reshape(1,-1))

X_train_essay_word_count_norm = normalizer.transform(X_train['essay_word_count'].values
.reshape(1,-1))
X_cv_essay_word_count_norm = normalizer.transform(X_cv['essay_word_count'].values.reshape(1,-1))
X_test_essay_word_count_norm = normalizer.transform(X_test['essay_word_count'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_essay_word_count_norm.shape, y_train.shape)
print(X_cv_essay_word_count_norm.shape, y_cv.shape)
print(X_test_essay_word_count_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

```
=====
=====
```

In [126]:

```
print("Transpose of essay word count norm")

X_train_essay_word_count_norm = X_train_essay_word_count_norm.transpose()
X_cv_essay_word_count_norm = X_cv_essay_word_count_norm.transpose()
X_test_essay_word_count_norm = X_test_essay_word_count_norm.transpose()

print("After vectorizations")
print(X_train_essay_word_count_norm.shape, y_train.shape)
print(X_cv_essay_word_count_norm.shape, y_cv.shape)
print(X_test_essay_word_count_norm.shape, y_test.shape)
print("="*100)
```

Transpose of essay word count norm

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

```
=====
=====
```

**title word count vectorization**

In [127]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['title_word_count'].values.reshape(1,-1))

X_train_title_word_count_norm = normalizer.transform(X_train['title_word_count'].values
.reshape(1,-1))
X_cv_title_word_count_norm = normalizer.transform(X_cv['title_word_count'].values.reshape(1,-1))
X_test_title_word_count_norm = normalizer.transform(X_test['title_word_count'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_title_word_count_norm.shape, y_train.shape)
print(X_cv_title_word_count_norm.shape, y_cv.shape)
print(X_test_title_word_count_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

(1, 22445) (22445,)

(1, 11055) (11055,)

(1, 16500) (16500,)

=====

In [128]:

```
print("Transpose of essay word count norm")

X_train_title_word_count_norm = X_train_title_word_count_norm.transpose()
X_cv_title_word_count_norm = X_cv_title_word_count_norm.transpose()
X_test_title_word_count_norm = X_test_title_word_count_norm.transpose()

print("After vectorizations")
print(X_train_title_word_count_norm.shape, y_train.shape)
print(X_cv_title_word_count_norm.shape, y_cv.shape)
print(X_test_title_word_count_norm.shape, y_test.shape)
print("="*100)
```

Transpose of essay word count norm

After vectorizations

(22445, 1) (22445,)

(11055, 1) (11055,)

(16500, 1) (16500,)

=====

## Set 5

In [129]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_NO_TEXT = hstack((X_train_pos_norm, X_train_neg_norm, X_train_neu_norm, X_train_compound_norm, X_train_essay_word_count_norm, X_train_title_word_count_norm, X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_CSC_ohe, X_train_CC_ohe, X_train_price_norm, X_train_quantity_norm, X_train_TPPP_norm)).tocsr()
X_cr_NO_TEXT = hstack((X_cv_pos_norm, X_cv_neg_norm, X_cv_neu_norm, X_cv_compound_norm, X_cv_essay_word_count_norm, X_cv_title_word_count_norm, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_CSC_ohe, X_cv_CC_ohe, X_cv_price_norm, X_cv_quantity_norm, X_cv_TPPP_norm)).tocsr()
X_te_NO_TEXT = hstack((X_test_pos_norm, X_test_neg_norm, X_test_neu_norm, X_test_compound_norm, X_test_essay_word_count_norm, X_test_title_word_count_norm, X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_CSC_ohe, X_test_CC_ohe, X_test_price_norm, X_test_quantity_norm, X_test_TPPP_norm)).tocsr()

print("Final Data matrix")
print(X_tr_NO_TEXT.shape, y_train.shape)
print(X_cr_NO_TEXT.shape, y_cv.shape)
print(X_te_NO_TEXT.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(22445, 108) (22445,)
(11055, 108) (11055,)
(16500, 108) (16500,)
```

```
=====
=====
```

## Logistic regression with no text data

In [130]:

```
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
#from sklearn.grid_search import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence va
lues, or non-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.

"""

clf = LogisticRegression(class_weight='balanced');
parameters = {'C':[10**-4, 10**-3, 10**-2, 1, 10, 100, 1000, 500, 1000, 10000]}
sd=GridSearchCV(clf, parameters, cv=5, scoring='roc_auc', return_train_score=True)
sd.fit(X_tr_NO_TEXT, y_train);

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

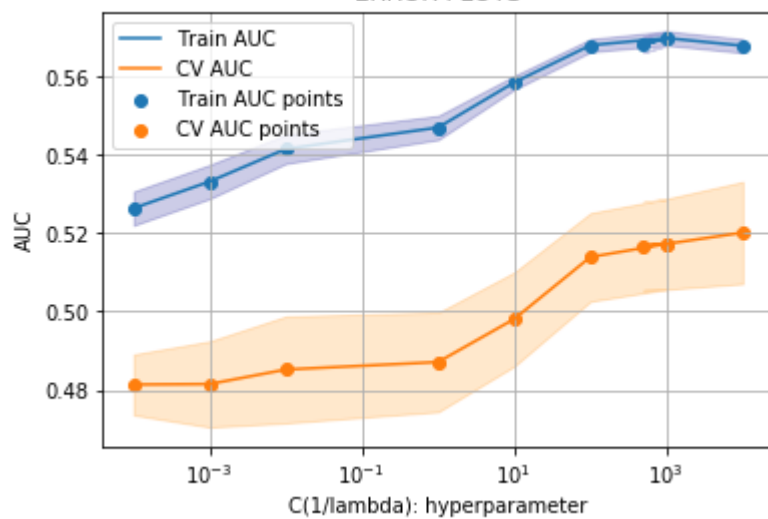
plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], train_auc - train_auc_std, train_auc + train_auc_std, alpha=0.2, color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.2, color='darkorange')

plt.scatter(parameters['C'], train_auc, label='Train AUC points')
plt.scatter(parameters['C'], cv_auc, label='CV AUC points')
plt.xscale('log')

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

# ERROR PLOTS





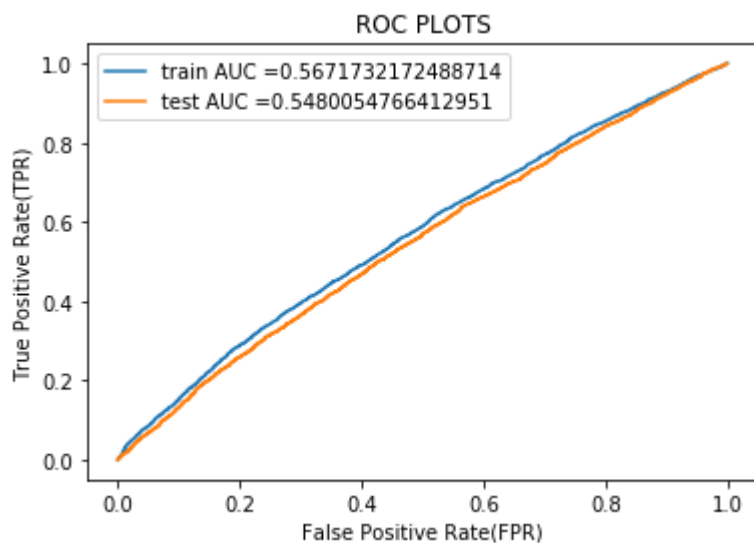
In [131]:

```
##Fitting Model to Hyper-Parameter Curve
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

neigh = LogisticRegression(C=10**3, class_weight='balanced');
neigh.fit(X_tr_NO_TEXT , y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_tr_NO_TEXT)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_te_NO_TEXT)[:,1])

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()
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")
plt.title("ROC PLOTS")
plt.show()
```



In [132]:

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

```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_tr_NO_TEXT )), annot=True, ax = ax,
            fmt='g'); #annot=True to annotate cells
```

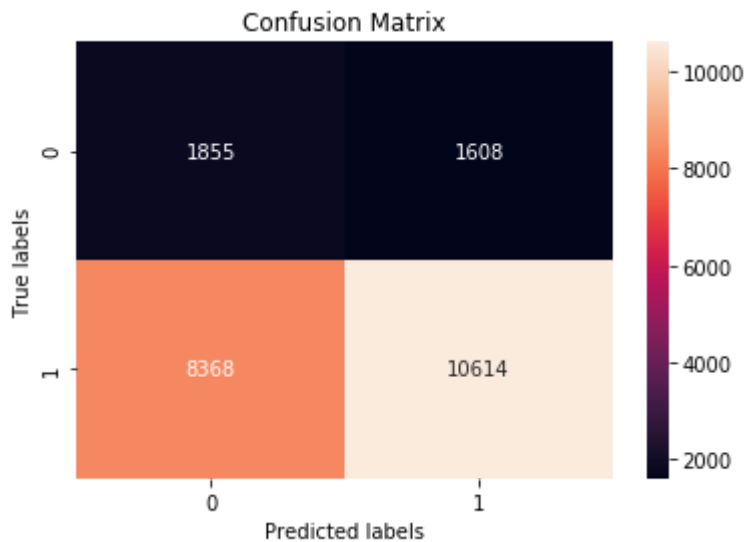
```
# Labels, title and ticks
```

```
ax.set_xlabel('Predicted labels');
```

```
ax.set_ylabel('True labels');
```

```
ax.set_title('Confusion Matrix');
```

```
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```

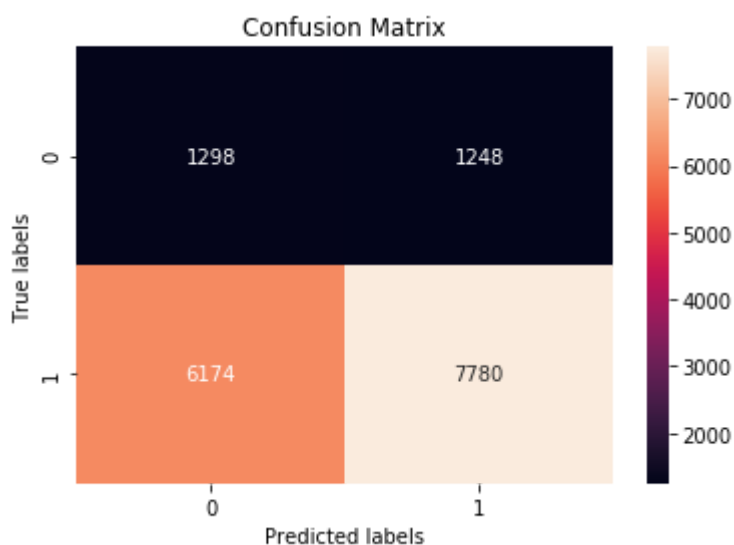


In [133]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_te_NO_TEXT )), annot=True, ax = ax
,fmt='g'); #annot=True to annotate cells

# Labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'b
usiness']);
```



In [ ]:

In [ ]:

### 3. Conclusion

In [134]:

```
# Please compare all your models using Prettytable Library
```

In [135]:

```
# Please compare all your models using Prettytable library
# http://zetcode.com/python/prettytable/

from prettytable import PrettyTable

#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Alpha:Hyper Parameter", "AUC"]

x.add_row(["BOW", "Logistic Regression", 0.01, 0.65])
x.add_row(["TFIDF", "Logistic Regression", 1, 0.64])
x.add_row(["AVG W2V", "Logistic Regression", 1, 0.68])
x.add_row(["TFIDF W2V", "Logistic Regression", 0.01, 0.68])
x.add_row(["NO TEXT", "Logistic Regression", 1000, 0.53])

print(x)
```

Vectorizer	Model	Alpha:Hyper Parameter	AUC
BOW	Logistic Regression	0.01	0.65
TFIDF	Logistic Regression	1	0.64
AVG W2V	Logistic Regression	1	0.68
TFIDF W2V	Logistic Regression	0.01	0.68
NO TEXT	Logistic Regression	1000	0.53

## Observation

1. Without text data there is a noticable difference observed in AUC score.
2. TFIDF W2V and AVG W2V has given same AUC score with different hyper parameters.

In [ ]: