

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. Example: p036502
<code>project_title</code>	Title of the project. Examples: • Art Will Make You Happy! • First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: • Grades PreK-2 • Grades 3-5 • Grades 6-8 • Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth Examples: • Music & The Arts • Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: • Literacy • Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: • My students need hands on literacy materials to manage sensory needs!
<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. Example: 2016-04-28 12:43:56.245

Feature	Description
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher.
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. Example: 2

* See the section **Notes on the Essay Data** for more details about these features.

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:

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	Price of the resource required. Example: 9.95

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):

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project was not approved, and a value of 1 indicates the project was approved.

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 plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter

from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_validate

```

```

C:\Users\myuri\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

```

1.1 Reading Data

In [2]:

```
project_data = pd.read_csv('train_data.csv',nrows=50000)
```

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 (50000, 17)
```

```

-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']

```

In [4]:

```
resource_data = pd.read_csv('resources.csv')
```

In [5]:

```

print("Number of data points in train data", resource_data.shape)
print('-'*50)
print("The attributes of data :", resource_data.columns.values)
resource_data.head(2)

```

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

The attributes of data : ['id' 'description' 'quantity' 'price']

Out[5]:

	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

In [6]:

```
# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]

#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
project_data.drop('project_submitted_datetime', axis=1, inplace=True)
project_data.sort_values(by=['Date'], inplace=True)

# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
project_data = project_data[cols]
project_data.head(2)
```

Out[6]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_s
473	100660 p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2	
41558	33679 p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5	L

1.2 preprocessing of project_subject_categories

In [7]:

```
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]))

```

1.3 preprocessing of project_subject_subcategories

In [8]:

```

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 & H
unger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Scienc
e"=> "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 placeing 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('&','_')
            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]))

```

1.3 Text preprocessing

In [9]:

```

# 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 [10]:

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

Out[10]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
473	100660 p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2	Flexi Seating Flexi Learn

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	Project Description
41558	33679	p137682	06f0e62e17de34fc01020c77549e1d5	Mrs.	WA	Grades 3-5	The Art Thinki

2016-04-27 01:05:25

In [11]:

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

I recently read an article about giving students a choice about how they learn. We already set goals; why not let them choose where to sit, and give them options of what to sit on? I teach at a low-income (Title 1) school. Every year, I have a class with a range of abilities, yet they are all the same age. They learn differently, and they have different interests. Some have ADHD, and some are fast learners. Yet they are eager and active learners that want and need to be able to move around the room, yet have a place that they can be comfortable to complete their work. We need a classroom rug that we can use as a class for reading time, and students can use during other learning times. I have also requested four Kore Kids wobble chairs and four Back Jack padded portable chairs so that students can still move during whole group lessons without disrupting the class. Having these areas will provide these little ones with a way to wiggle while working. Benjamin Franklin once said, "Tell me and I forget, teach me and I may remember, involve me and I learn." I want these children to be involved in their learning by having a choice on where to sit and how to learn, all by giving them options for comfortable flexible seating.

In [12]:

```
# 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 [13]:

```
sent = decontracted(project_data['essay'].values[0])
print(sent)
print("="*50)
```

I recently read an article about giving students a choice about how they learn. We already set goals; why not let them choose where to sit, and give them options of what to sit on? I teach at a low-income (Title 1) school. Every year, I have a class with a range of abilities, yet they are all the same age. They learn differently, and they have different interests. Some have ADHD, and some are fast learners. Yet they are eager and active learners that want and need to be able to move around the room, yet have a place that they can be comfortable to complete their work. We need a classroom rug that we can use as a class for reading time, and students can use during other learning times. I have also requested four Kore Kids wobble chairs and four Back Jack padded portable chairs so that students can still move during whole group lessons without disrupting the class. Having these areas will provide these little ones with a way to wiggle while working. Benjamin Franklin once said, "Tell me and I forget, teach me and I may remember, involve me and I learn." I want these children to be involved in their learning by having a choice on where to sit and how to learn, all by giving them options for comfortable flexible seating.

In [14]:

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
```

```
# If you want to remove from string python: http://texendrieter.com/info/remove-line-breaks-python/
sent = sent.replace('\\r', ' ')
sent = sent.replace('\\n', ' ')
sent = sent.replace('\\n', ' ')
print(sent)
```

I recently read an article about giving students a choice about how they learn. We already set goals; why not let them choose where to sit, and give them options of what to sit on? I teach at a low-income (Title 1) school. Every year, I have a class with a range of abilities, yet they are all the same age. They learn differently, and they have different interests. Some have ADHD, and some are fast learners. Yet they are eager and active learners that want and need to be able to move around the room, yet have a place that they can be comfortable to complete their work. We need a classroom rug that we can use as a class for reading time, and students can use during other learning times. I have also requested four Kore Kids wobble chairs and four Back Jack padded portable chairs so that students can still move during whole group lessons without disrupting the class. Having these areas will provide these little ones with a way to wiggle while working. Benjamin Franklin once said, Tell me and I forget, teach me and I may remember, involve me and I learn. I want these children to be involved in their learning by having a choice on where to sit and how to learn, all by giving them options for comfortable flexible seating.

In [15]:

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

I recently read an article about giving students a choice about how they learn. We already set goals; why not let them choose where to sit and give them options of what to sit on? I teach at a low income Title 1 school. Every year I have a class with a range of abilities yet they are all the same age. They learn differently and they have different interests. Some have ADHD and some are fast learners. Yet they are eager and active learners that want and need to be able to move around the room yet have a place that they can be comfortable to complete their work. We need a classroom rug that we can use as a class for reading time and students can use during other learning times. I have also requested four Kore Kids wobble chairs and four Back Jack padded portable chairs so that students can still move during whole group lessons without disrupting the class. Having these areas will provide these little ones with a way to wiggle while working. Benjamin Franklin once said Tell me and I forget teach me and I may remember involve me and I learn. I want these children to be involved in their learning by having a choice on where to sit and how to learn all by giving them options for comfortable flexible seating.

In [16]:

```
# 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're", "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', 'they', 'them',
'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', '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', 'over', 'under',
, 'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', '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", 'now', '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', 'mightn',
"mightn't", 'mustn', \
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
"wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [17]:

```
# Combining all the above students
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('\\n', ' ')
    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.lower() not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

```
100%|███████████████████████████████████████████████████████████████| 50000/50000 [01:  
33<00:00, 535.24it/s]
```

In [18]:

```
# after preprocessing
preprocessed_essays[0]
```

Out [18] :

'recently read article giving students choice learn already set goals not let choose sit give options sit teach low income title 1 school every year class range abilities yet age learn differently different interests adhd fast learners yet eager active learners want need able move around room yet place comfortable complete work need classroom rug use class reading time students use learning times also requested four kore kids wobble chairs four back jack padded portable chairs students still move whole group lessons without disrupting class areas provide little ones way wiggle working benjamin franklin said tell forget teach may remember involve learn want children involved learning choice sit learn giving options comfortable flexible seating'

1.4 Preprocessing of `project_title`

In [19]:

```
# preprocessing of project title
```

In [20]:

```
sent = decontracted(project_data['project_title'].values[0])
# \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', ' ')
#remove special character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

Flexible Seating for Flexible Learning

In [21]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_project_title = []
# 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('\n', ' ')
    sent = sent.replace('\t', ' ')
    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_project_title.append(sent.lower().strip())
```


100% | 50000/50000
[00:03<00:00, 13663.10it/s]

1.5 Preparing data for models

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 (optional)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

In [22]:

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

Out[22]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2 Flexi Seating Flexi Learn
41558	33679	p137682	06f6e62e17de34fc81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5 Going De The Art In Thinki

In [23]:

```
#number of words in project titlefor set 5 new feature
```

In [24]:

```
new_title = []  
for i in tqdm(project_data['project_title']):  
    j = decontracted(i)  
    new_title.append(j)
```

100% | 50000/50000
[00:01<00:00, 34411.00it/s]

In [25]:

```
#Introducing New Features  
title_word_count = []  
#for i in project_data['project_title']:  
for i in tqdm(new_title):  
    j = len(i.split())  
    title_word_count.append(j)  
    #print(j)  
project_data['title word count'] = title_word_count
```

```
100%|██████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:00<00:00, 355579.07it/s]
```

In [26]:

```
#number of words in project title for set 5 new feature
```

In [27]:

```
new_essay = []
for i in tqdm(project_data['essay']):
    j = decontracted(i)
    new_essay.append(j)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:02<00:00, 19633.26it/s]
```

In [28]:

```
#Introducing New Features
essay_word_count = []
#for i in project_data['project_title']:
for i in tqdm(new_essay ):
    j = len(i.split())
    essay_word_count.append(j)
    #print(j)
project_data['essay word count'] = essay_word_count
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:02<00:00, 22441.33it/s]
```

In [29]:

```
#split the data into train ,test and cross validation
```

In [30]:

```
y = project_data['project_is_approved'].values
project_data.drop(['project_is_approved'], axis=1, inplace=True)
print(project_data.shape)
```

(50000, 19)

In [31]:

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

Out [31]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2 Flexi Seating Flexi Learn
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5 Going De The Ar In Thinki



Computing Sentiment Scores

In [32]:

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

# 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 w
ith the biggest enthusiasm'
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
```

```
[nltk_data] Downloading package vader_lexicon to
[nltk_data] C:\Users\myuri\AppData\Roaming\nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
```

neg: 0.109, neu: 0.693, pos: 0.198, compound: 0.2023,

In [33]:

```
SID = SentimentIntensityAnalyzer()
#There is NEGITIVE and POSITIVE and NEUTRAL and COMPUND SCORES
#http://www.nltk.org/howto/sentiment.html

negative = []
positive = []
neutral = []
compound = []

for i in tqdm(project_data['essay']):
    j = SID.polarity_scores(i)['neg']
    k = SID.polarity_scores(i)['neu']
    l = SID.polarity_scores(i)['pos']
    m = SID.polarity_scores(i)['compound']
    negative.append(j)
    positive.append(k)
    neutral.append(l)
    compound.append(m)
```

```
100%|███████████████████████████████████████████████████████████| 50000/50000 [28  
:09<00:00, 29.59it/s]
```

In [34]:

```
project_data['negative'] = negative
```

In [35]:

```
project_data['positive'] = positive
project_data['neutral'] = neutral
project_data['compound'] = compound
project_data.head(2)
```

Out[35]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title	
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27	Grades PreK-2	Flexi Seating Flexi

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5

Going De
The Ar
In
Thinki

2 rows × 23 columns

In [36]:

```
#https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
#Splitting data into Train and cross validation
# split the data set into train and test
X_train, X_test, y_train, y_test = train_test_split(project_data, y, test_size=0.33, stratify=y)
# split the train data set into cross validation train and cross validation test
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)
```

In [37]:

```
print(X_train.columns)
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_grade_category', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'clean_categories',
      'clean_subcategories', 'essay', 'title_word_count', 'essay_word_count',
      'negative', 'positive', 'neutral', 'compound'],
      dtype='object')
```

In [38]:

```
print(X_train.shape)
print(X_test.shape)
print(X_cv.shape)
```

```
(22445, 23)
(16500, 23)
(11055, 23)
```

In [39]:

```
print(y_train.shape)
print(y_test.shape)
print(y_cv.shape)
```

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

In [40]:

```
print(X_train['essay'].values[0])
```

If you were to walk in my classroom, you would see some students working at their seats, some would be standing at a table, some would be laying on the floor, while others would be working with me or working collaboratively with other students.\r\n\r\nThe students at my school and in my classroom are amazing! They come from diverse backgrounds. My students are very talented and capable. However, some lack experiences and opportunity. Technology keeps students in engaged in learning! I would like the Apple TV to keep them motivated, engaged and up to date! This donation would ensure that they receive quality learning through technology!\r\n\r\nPlease consider helping our classroom!!My classroom is has several pieces of technology but I am always looking for more innovative and creative ways to engage my students and offer them as many learning platforms as possible! Fortunately, I have had many technology donations over the past few years and love teaching and sharing this technology with my students. The only thing we are missing is an Apple TV. Technology engages children while learning. If we had an Apple TV, students would be able to project their work from a computer or device onto our class Smartboard. We will be able to watch t

Project shall work from a computer or device once our class commences. We will be able to watch the news and current live events. We will also be able to access many apps offered by Apple TV. Students will also have access to many more math programs and games.

In [41]:

```
#preprocessing of train ,cross validation and test essay data
```

In [42]:

```
#preprocess the X_train essay
```

In [43]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_essay_train_data = []
# tqdm is for printing the status bar
for sentence in tqdm(X_train['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\n', ' ')
    sent = sent.replace('\\t', ' ')
    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_essay_train_data.append(sent.lower().strip())
```

```
100%|███████████████████████████████████████████████████████| 22445/22445 [00:  
40<00:00, 557.34it/s]
```

In [44]:

```
#preprocess the X cv essay
```

In [45]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_essay_cv_data = []
# tqdm is for printing the status bar
for sentence in tqdm(X_cv['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\n', ' ')
    sent = sent.replace('\\t', ' ')
    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_essay_cv_data.append(sent.lower().strip())
```

```
100%|███████████████████████████████████████████████████████████| 11055/11055 [00:  
20<00:00, 543.11it/s]
```

In [46]:

```
#preprocess the X test essay
```

In [47]:

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_essay_test_data = []
# tqdm is for printing the status bar
for sentence in tqdm(X_test['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\n', ' ')
    sent = sent.replace('\n', ' ')

```

[illegible]

```
#preprocessing of x_train,x_cv and x_test of project title
```

```
100% |██████████████████████████████████████████████████████████████████████████| 22445/22445  
[00:01<00:00, 13059.81it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:00<00:00, 11055.76it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:01<00:00, 11234.84it/s]
```

1.5.1 Vectorizing Categorical data

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

In [52]:

```
#vectorisation of clean categories
```

In [53]:

```
# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vect_categories= CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vect_categories.fit(project_data['clean_categories'].values)

train_categories_one_hot=vect_categories.transform(X_train['clean_categories'].values)
cv_categories_one_hot=vect_categories.transform(X_cv['clean_categories'].values)
test_categories_one_hot=vect_categories.transform(X_test['clean_categories'].values)

print(vect_categories.get_feature_names())
print("Shape of train matrix after one hot encodig ",train_categories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",cv_categories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",test_categories_one_hot.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of train matrix after one hot encodig (22445, 9)
Shape of train matrix after one hot encodig (11055, 9)
Shape of train matrix after one hot encodig (16500, 9)
```

In [54]:

```
#vectorisation of clean subcategories
```

In [55]:

```
# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_subcategories= CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_subcategories.fit(project_data['clean_subcategories'].values)

train_subcategories_one_hot=vectorizer_subcategories.transform(X_train['clean_subcategories'].values)
cv_subcategories_one_hot=vectorizer_subcategories.transform(X_cv['clean_subcategories'].values)
test_subcategories_one_hot=vectorizer_subcategories.transform(X_test['clean_subcategories'].values)

print(vectorizer_subcategories.get_feature_names())
print("Shape of train matrix after one hot encodig ",train_subcategories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",cv_subcategories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",test_subcategories_one_hot.shape)
```

```
['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']
Shape of train matrix after one hot encodig (22445, 30)
Shape of train matrix after one hot encodig (11055, 30)
Shape of train matrix after one hot encodig (16500, 30)
```

In [56]:

```
# Build the data matrix using these features-- school_state : categorical data (one hot encoding)
##Encoding for school state
```

In [57]:

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict_state = dict(my_counter)
sorted_cat_dict_state = dict(sorted(cat_dict_state.items(), key=lambda kv: kv[1]))

from sklearn.feature_extraction.text import CountVectorizer
vectorizer_state = CountVectorizer(vocabulary=list(sorted_cat_dict_state.keys()), lowercase=False,
binary=True)
vectorizer_state.fit(project_data['school_state'].values)

train_state_one_hot=vectorizer_state.transform(X_train['school_state'].values)
cv_state_one_hot=vectorizer_state.transform(X_cv['school_state'].values)
test_state_one_hot=vectorizer_state.transform(X_test['school_state'].values)

print(vectorizer_state.get_feature_names())
print("Shape of train matrix after one hot encodig ",train_state_one_hot.shape)
print("Shape of train matrix after one hot encodig ",cv_state_one_hot.shape)
print("Shape of train matrix after one hot encodig ",test_state_one_hot.shape)

['VT', 'WY', 'ND', 'MT', 'RI', 'NH', 'SD', 'NE', 'AK', 'DE', 'WV', 'ME', 'NM', 'HI', 'DC', 'KS', 'I
D', 'IA', 'AR', 'CO', 'MN', 'OR', 'MS', 'KY', 'NV', 'MD', 'TN', 'CT', 'AL', 'UT', 'WI', 'VA', 'AZ',
'NJ', 'OK', 'MA', 'LA', 'WA', 'MO', 'IN', 'OH', 'PA', 'MI', 'GA', 'SC', 'IL', 'NC', 'FL', 'TX', 'NY
', 'CA']
Shape of train matrix after one hot encodig (22445, 51)
Shape of train matrix after one hot encodig (11055, 51)
Shape of train matrix after one hot encodig (16500, 51)
```

In [58]:

```
#Encoding for project_grade_category
```

In [59]:

```
project_data.project_grade_category = project_data.project_grade_category.str.replace('\s+', '_')
project_data.project_grade_category = project_data.project_grade_category.str.replace('-', '_')
project_data['project_grade_category'].value_counts()
```

Out[59]:

```
Grades_PreK_2      20316
Grades_3_5         16968
Grades_6_8         7750
Grades_9_12        4966
Name: project_grade_category, dtype: int64
```

In [60]:

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict_grade = dict(my_counter)
sorted_cat_dict_grade = dict(sorted(cat_dict_grade.items(), key=lambda kv: kv[1]))
print(sorted_cat_dict_grade)

# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_grade_cat = CountVectorizer(vocabulary=list(sorted_cat_dict_grade.keys()), lowercase=False,
binary=True)
vectorizer_grade_cat.fit(project_data['project_grade_category'].values)
```



```

train_grade_one_hot=vectorizer_grade_cat.transform(X_train['project_grade_category'].values)
cv_grade_one_hot=vectorizer_grade_cat.transform(X_cv['project_grade_category'].values)
test_grade_one_hot=vectorizer_grade_cat.transform(X_test['project_grade_category'].values)

print(vectorizer_grade_cat.get_feature_names())
print("Shape of train matrix after one hot encoding ",train_grade_one_hot.shape)
print("Shape of train matrix after one hot encoding ",cv_grade_one_hot.shape)
print("Shape of train matrix after one hot encoding ",test_grade_one_hot.shape)

```

```

{'Grades_9_12': 4966, 'Grades_6_8': 7750, 'Grades_3_5': 16968, 'Grades_PreK_2': 20316}
['Grades_9_12', 'Grades_6_8', 'Grades_3_5', 'Grades_PreK_2']
Shape of train matrix after one hot encoding (22445, 4)
Shape of train matrix after one hot encoding (11055, 4)
Shape of train matrix after one hot encoding (16500, 4)

```

In [61]:

```
#Encoding for teacher_prefix
```

In [62]:

```

project_data.teacher_prefix = project_data.teacher_prefix.str.replace('\s+', '_')
project_data.teacher_prefix = project_data.teacher_prefix.str.replace('-', '_')
project_data['teacher_prefix'].value_counts()

```

Out[62]:

```

Mrs.      26140
Ms.       17936
Mr.       4859
Teacher   1061
Dr.        2
Name: teacher_prefix, dtype: int64

```

In [63]:

```

#https://stackoverflow.com/questions/42224700/attributeerror-float-object-has-no-attribute-split
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna("")

```

In [64]:

```

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict_prefix = dict(my_counter)
sorted_cat_dict_prefix = dict(sorted(cat_dict_prefix.items(), key=lambda kv: kv[1]))

# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_teacher_prefix = CountVectorizer(vocabulary=list(sorted_cat_dict_prefix.keys()), lowerc
ase=False, binary=True)
vectorizer_teacher_prefix.fit(project_data['teacher_prefix'].values.astype('U'))

train_teacher_prefix_one_hot=vectorizer_teacher_prefix.transform(X_train['teacher_prefix'].values.
astype('U'))
cv_teacher_prefix_one_hot=vectorizer_teacher_prefix.transform(X_cv['teacher_prefix'].values.as
type('U'))
test_teacher_prefix_one_hot=vectorizer_teacher_prefix.transform(X_test['teacher_prefix'].values.as
type('U'))

print(vectorizer_teacher_prefix.get_feature_names())
print("Shape of train matrix after one hot encoding ",train_teacher_prefix_one_hot.shape)
print("Shape of train matrix after one hot encoding ",cv_teacher_prefix_one_hot.shape)
print("Shape of train matrix after one hot encoding ",test_teacher_prefix_one_hot.shape)

```

```
['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']  
Shape of train matrix after one hot encodig (22445, 5)  
Shape of train matrix after one hot encodig (11055, 5)  
Shape of train matrix after one hot encodig (16500, 5)
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [65]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).  
vectorizer_essay_bow = CountVectorizer(min_df=10,max_features=5000,ngram_range = (2,2))  
vectorizer_essay_bow.fit(preprocessed_essay_train_data)
```

```
text_bow_essays_train = vectorizer_essay_bow.transform(preprocessed_essay_train_data)  
print("Shape of matrix after one hot encodig ",text_bow_essays_train.shape)
```

Shape of matrix after one hot encodig (22445, 5000)

In [66]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).  
text_bow_essays_cv = vectorizer_essay_bow.transform(preprocessed_essay_cv_data)  
print("Shape of matrix after one hot encodig ",text_bow_essays_cv.shape)
```

Shape of matrix after one hot encodig (11055, 5000)

In [67]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).  
text_bow_essays_test= vectorizer_essay_bow.transform(preprocessed_essay_test_data)  
print("Shape of matrix after one hot encodig ",text_bow_essays_test.shape)
```

Shape of matrix after one hot encodig (16500, 5000)

In [68]:

```
# you can vectorize the title also  
# before you vectorize the title make sure you preprocess it
```

In [69]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).  
vectorizer_bow_title = CountVectorizer(ngram_range = (2,2),min_df=10,max_features=5000)  
vectorizer_bow_title.fit(train_preprocessed_project_title)
```

```
text_bow_title_train= vectorizer_bow_title.transform(train_preprocessed_project_title)  
print("Shape of matrix after one hot encodig ",text_bow_title_train.shape)
```

Shape of matrix after one hot encodig (22445, 636)

In [70]:

```
text_bow_title_cv=vectorizer_bow_title.transform(cv_preprocessed_project_title)  
print("Shape of matrix after one hot encodig ",text_bow_title_cv.shape)
```

Shape of matrix after one hot encodig (11055, 636)

In [71]:

```
text_bow_title_test= vectorizer_bow_title.transform(test_preprocessed_project_title)  
print("Shape of matrix after one hot encodig ",text_bow_title_test.shape)
```

Shape of matrix after one hot encodig (16500, 636)

TFIDF Vectorizer on project_title

In [72]:

```
# Similarly you can vectorize for title also
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(2,2))
vectorizer.fit(train_preprocessed_project_title)

text_tfidf_title_train = vectorizer.transform(train_preprocessed_project_title)
print("Shape of matrix after one hot encodig ",text_tfidf_title_train.shape)
```

Shape of matrix after one hot encodig (22445, 636)

In [73]:

```
# Similarly you can vectorize for title
text_tfidf_title_cv = vectorizer.transform(cv_preprocessed_project_title)
print("Shape of matrix after one hot encodig ",text_tfidf_title_cv.shape)
```

Shape of matrix after one hot encodig (11055, 636)

In [74]:

```
# Similarly you can vectorize for title also
text_tfidf_title_test = vectorizer.transform(test_preprocessed_project_title)
print("Shape of matrix after one hot encodig ",text_tfidf_title_test.shape)
```

Shape of matrix after one hot encodig (16500, 636)

TFIDF Vectorizer on preprocessed essay

In [75]:

```
# Similarly you can vectorize for title also
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(2,2))
vectorizer.fit(preprocessed_essay_train_data)

tfidf_essay_train = vectorizer.transform(preprocessed_essay_train_data)
print("Shape of matrix after one hot encodig ",tfidf_essay_train.shape)
```

Shape of matrix after one hot encodig (22445, 5000)

In [76]:

```
# Similarly you can vectorize for title also
tfidf_essay_cv = vectorizer.transform(preprocessed_essay_cv_data)
print("Shape of matrix after one hot encodig ",tfidf_essay_cv.shape)
```

Shape of matrix after one hot encodig (11055, 5000)

In [77]:

```
tfidf_essay_test = vectorizer.transform(preprocessed_essay_test_data)
print("Shape of matrix after one hot encodig ",tfidf_essay_test.shape)
```

Shape of matrix after one hot encodig (16500, 5000)

1.5.2.3 Using Pretrained Models: Avg W2V

In [78]:

```
'''
# 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 coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      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-sa
ve-and-load-variables-in-python/

import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)

'''
```

Out[78]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\nencoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\nword = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n    model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =\nloadGloveModel('glove.42B.300d.txt')\n\n# =====\n\nOutput:\n\nLoading G\nlove Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n# =====\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split('\n\n'))\nfor i in preprocod_titles:\n    words.extend(i.split('\n\n'))\nprint("all the words in the\ncoupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus",\nlen(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha\nt are present in both glove vectors and our coupus",\n      len(inter_words),\n      "(" ,np.round(len(inter_words)/len(words)*100,3), "%) ") \n\nwords_courpus = {}\nwords_glove =\nset(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\nprint("word 2 vec length". len(words_courpus))\n\n\n# stronging variables into pickle files python
```

◀ ▶

```
# 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_essay_train_data = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essay_train_data): # 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_essay_train_data.append(vector)

print(len(avg_w2v_essay_train_data))
print(len(avg_w2v_essay_train_data[0]))
```

100% |██| 22445/22445
[00:20<00:00, 1109.29it/s]

```
In [81]:  
  
# average Word2Vec  
# compute average word2vec for each review.  
avg_w2v_essay_cv_data = []; # the avg-w2v for each sentence/review is stored in this list  
for sentence in tqdm(preprocessed_essay_cv_data): # 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_essay_cv_data.append(vector)  
  
print(len(avg_w2v_essay_cv_data))  
print(len(avg_w2v_essay_cv_data[0]))
```

100% |██| 11055/11055
[00:09<00:00, 1120.09it/s]

11055
300

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay_test_data = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essay_test_data): # 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
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:15<00:00, 1085.28it/s]
```

In [83]:

In [84]:

```
100%|██████████████████████████████████████████████████████████████████████████| 22445/22445  
[00:01<00:00, 20233.37it/s]
```

In [85]:

```
100%|██████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:00<00:00, 19654.74it/s]
```

11055

In [89]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essay_cv_data = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essay_cv_data): # 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 tf
            idf 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_essay_cv_data.append(vector)

print(len(tfidf_w2v_essay_cv_data))
print(len(tfidf_w2v_essay_cv_data[0]))
```

```
100%|███████████████████████████████████████████| 11055/11055 [01:  
11<00:00, 154.06it/s]
```

11055
300

In [90]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essay_test_data = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essay_test_data): # 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 tf
            idf 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_essay_test_data.append(vector)

print(len(tfidf_w2v_essay_test_data))
print(len(tfidf_w2v_essay_test_data[0]))
```

```
100%|██████████████████████████████████████████████████████████████| 16500/16500 [01:  
46<00:00, 154.91it/s]
```

16500
300

Using Pretrained Models: TFIDF weighted W2V on project_title

In [91]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(train_preprocessed_project_title)
# 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 [92]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_train_project_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(train_preprocessed_project_title): # 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 tf
            idf 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_train_project_title.append(vector)

print(len(tfidf_w2v_train_project_title))
print(len(tfidf_w2v_train_project_title[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 22445/22445  
[00:02<00:00, 9772.65it/s]
```

22445
300

In [93]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_cv_project_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(cv_preprocessed_project_title): # 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 tf
            idf 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_cv_project_title.append(vector)

print(len(tfidf_w2v_cv_project_title))
print(len(tfidf_w2v_cv_project_title[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:01<00:00, 8227.55it/s]
```

11055
300

In [94]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_test_project_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(test_preprocessed_project_title): # 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
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:01<00:00, 8585.96it/s]
```

In [95]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_grade_category', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'clean_categories',
      'clean_subcategories', 'essay', 'title_word_count', 'essay_word_count',
      'negative', 'positive', 'neutral', 'compound'],
      dtype='object')
```

```
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')
```

```
x_train = pd.merge(X_train, price_data, on = "id", how = "left")
x_test = pd.merge(X_test, price_data, on = "id", how = "left")
x_cv = pd.merge(X_cv, price_data, on = "id", how = "left")
```

 $(22445, 25)$

In [99]:

```
# check this one: https://www.youtube.com/watch?v=0HQoCln3Z4&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(x_train['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.
train_price_standar = price_scalar.transform(x_train['price'].values.reshape(-1, 1))
train_price_standar

```

Mean : 301.22037202049455, Standard deviation : 385.68361456591435

Out[99]:

```

array([[ -0.39210992],
       [  0.10897437],
       [ -0.13539692],
       ...,
       [ -0.58662169],
       [ -0.47728336],
       [  0.49641629]])

```

In [100]:

```

price_scalar.fit(x_test['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.
test_price_standar = price_scalar.transform(x_test['price'].values.reshape(-1, 1))
test_price_standar

```

Mean : 297.4745654545455, Standard deviation : 370.1066540554988

Out[100]:

```

array([[ -0.35572061],
       [ -0.3524783 ],
       [ -0.02716667],
       ...,
       [ -0.41208274],
       [ -0.27949934],
       [ -0.47741526]])

```

In [101]:

```

price_scalar.fit(x_cv['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.
cv_price_standar = price_scalar.transform(x_cv['price'].values.reshape(-1, 1))
test_price_standar

```

Mean : 298.27790411578474, Standard deviation : 374.8149794131345

Out[101]:

```

array([[ -0.35572061],
       [ -0.3524783 ],
       [ -0.02716667],
       ...,
       [ -0.41208274],
       [ -0.27949934],
       [ -0.47741526]])

```

In [102]:

```

print(train_price_standar.shape, y_train.shape)
print(test_price_standar.shape, y_test.shape)
print(cv_price_standar.shape, y_cv.shape)

```

```

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

```

```
(11055, 1) (11055,)
```

Vectorizing teacher_number_of_previously_posted_projects

In [103]:

```
price_scalar.fit(x_train['teacher_number_of_previously_posted_projects'].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.
train_prev_proj_standar =
price_scalar.transform(x_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
train_prev_proj_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 11.328714635776342, Standard deviation : 27.86197142523261

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

Out[103]:

```
array([[ -0.3348189 ],
       [ -0.40660133],
       [ -0.40660133],
       ...,
       [ -0.29892769],
       [ -0.37071012],
       [ -0.40660133]])
```

In [104]:

```
price_scalar.fit(x_test['teacher_number_of_previously_posted_projects'].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.
test_prev_proj_standar =
price_scalar.transform(x_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
test_prev_proj_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 11.024909090909091, Standard deviation : 27.689953065133775

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

Out[104]:

```
array([[ -0.21758466],
       [ -0.28981303],
       [ -0.1453563 ],
       ...,
       [ -0.39815557],
```

```
[-0.39815557],  
[-0.32592721]])
```

In [105]:

```
price_scalar.fit(x_cv['teacher_number_of_previously_posted_projects'].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.  
cv_prev_proj_standar = price_scalar.transform(x_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))  
cv_prev_proj_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 11.423156942559928, Standard deviation : 29.413326546749165

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Out[105]:

```
array([[ -0.28637213],  
       [ -0.35436852],  
       [ -0.38836671],  
       ...,  
       [  0.29159718],  
       [ -0.18437755],  
       [ -0.38836671]])
```

Standardize Quantity

In [106]:

```
price_scalar.fit(x_train['quantity'].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.  
train_quantity_standar = price_scalar.transform(x_train['quantity'].values.reshape(-1, 1))  
train_quantity_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 16.932056137224325, Standard deviation : 26.27109332797795

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Out[106]:

```
array([[ -0.60644816],  
       [ -0.4922542 ],  
       [ -0.56838351],  
       ...,  
       [  0.26903882],  
       [ -0.60644816],  
       [  0.23097417]])
```

In [107]:

```
price_scalar.fit(x_test['quantity'].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.
test_quantity_standar = price_scalar.transform(x_test['quantity'].values.reshape(-1, 1))
test_quantity_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 17.18072727272727, Standard deviation : 26.832227424535287

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Out[107]:

```
array([[ 0.25414486],
       [-0.45395886],
       [-0.19307854],
       ...,
       [-0.52849609],
       [ 1.52127783],
       [ 0.47775656]])
```

In [108]:

```
price_scalar.fit(x_cv['quantity'].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.
cv_quantity_standar = price_scalar.transform(x_cv['quantity'].values.reshape(-1, 1))
cv_quantity_standar
```

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Mean : 17.123383084577114, Standard deviation : 27.82871903968372

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

Out[108]:

```
array([[ -0.14817006],
       [ 0.46270965],
       [-0.43564287],
       ...,
       [ 0.13930274],
       [-0.22003827],
       [-0.57937928]])
```

In [109]:

```
print(train_quantity_standar.shape, y_train.shape)
print(test_quantity_standar.shape, y_test.shape)
print(cv_quantity_standar.shape, y_cv.shape)
```

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

```
new_title = []
for i in tqdm(project_data['project_title']):
    j = decontracted(i)
    new_title.append(j)
```

In [111]:

```
#Introducing New Features
title_word_count = []
#for i in project_data['project_title']:
for i in tqdm(new_title):
    j = len(i.split())
    title_word_count.append(j)
    #print(j)
project_data['title word count'] = title_word_count
```

Standardize Title word count

```

title_scalar = StandardScaler()
title_scalar.fit(X_train['title_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {title_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
train_title_word_count_standar = title_scalar.transform(X_train['title_word_count'].values.reshape
(-1, 1))

title_scalar.fit(X_test['title_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {title_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
test_title_word_count_standar = title_scalar.transform(X_test['title_word_count'].values.reshape(-1
, 1))

title_scalar.fit(X_cv['title_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {title_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
cv_title_word_count_standar = title_scalar.transform(X_cv['title_word_count'].values.reshape(-1, 1)
)

print(train_title_word_count_standar.shape, y_train.shape)
print(test_title_word_count_standar.shape, y_test.shape)
print(cv_title_word_count_standar.shape, y_cv.shape)

```

Mean : 5.208687903764758, Standard deviation : 27.82871903968372

C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
Mean : 5.22630303030303, Standard deviation : 27.82871903968372
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
Mean : 5.215196743554952, Standard deviation : 27.82871903968372
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

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

Standardize Title_word_count

```
In [113]:
```

```
essay_scalar = StandardScaler()

essay_scalar.fit(X_train['essay_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
train_essay_word_count_standar = essay_scalar.transform(X_train['essay_word_count'].values.reshape
(-1, 1))

essay_scalar.fit(X_train['essay_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
test_essay_word_count_standar = essay_scalar.transform(X_test['essay_word_count'].values.reshape(-1
, 1))

essay_scalar.fit(X_cv['essay_word_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
cv_essay_word_count_standar = essay_scalar.transform(X_cv['essay_word_count'].values.reshape(-1, 1)
)

print(train_essay_word_count_standar.shape, y_train.shape)
print(test_essay_word_count_standar.shape, y_test.shape)
print(cv_essay_word_count_standar.shape, y_cv.shape)
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```



```
C:\Users\myuri\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

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

Standardize POSITIVE

In [114]:

```
essay_scalar.fit(X_train['positive'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
train_positive_standar = essay_scalar.transform(X_train['positive'].values.reshape(-1, 1))

essay_scalar.fit(X_train['positive'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
test_positive_standar = essay_scalar.transform(X_test['positive'].values.reshape(-1, 1))

essay_scalar.fit(X_cv['positive'].values.reshape(-1,1)) # finding the mean and standard deviation
of this data
cv_positive_standar= essay_scalar.transform(X_cv['positive'].values.reshape(-1, 1))

print(train_positive_standar.shape, y_train.shape)
print(test_positive_standar.shape, y_test.shape)
print(cv_positive_standar.shape, y_cv.shape)
```

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

Standardize NEGATIVE

In [115]:

```
essay_scalar.fit(X_train['negative'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
train_negative_standar = essay_scalar.transform(X_train['negative'].values.reshape(-1, 1))

essay_scalar.fit(X_train['negative'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
test_negative_standar = essay_scalar.transform(X_test['negative'].values.reshape(-1, 1))

essay_scalar.fit(X_cv['negative'].values.reshape(-1,1)) # finding the mean and standard deviation
of this data
cv_negative_standar = essay_scalar.transform(X_cv['negative'].values.reshape(-1, 1))

print(train_negative_standar.shape, y_train.shape)
print(test_negative_standar.shape, y_test.shape)
print(cv_negative_standar.shape, y_cv.shape)
```

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

Standardize neutral

In [116]:

```
essay_scalar.fit(X_train['neutral'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
train_neutral_standar = essay_scalar.transform(X_train['neutral'].values.reshape(-1, 1))

essay_scalar.fit(X_train['neutral'].values.reshape(-1,1)) # finding the mean and standard
```

```

essay_scalar.fit(X_train['neutral'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
test_neutral_standar = essay_scalar.transform(X_test['neutral'].values.reshape(-1, 1))

essay_scalar.fit(X_cv['neutral'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
cv_neutral_standar = essay_scalar.transform(X_cv['neutral'].values.reshape(-1, 1))

print(train_neutral_standar.shape, y_train.shape)
print(test_neutral_standar.shape, y_test.shape)
print(cv_neutral_standar.shape, y_cv.shape)

```

```

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

```

Standardize compound

In [117]:

```

essay_scalar.fit(X_train['compound'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
train_compound_standar = essay_scalar.transform(X_train['compound'].values.reshape(-1, 1))

essay_scalar.fit(X_train['compound'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
test_compound_standar = essay_scalar.transform(X_test['compound'].values.reshape(-1, 1))

essay_scalar.fit(X_cv['compound'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
cv_compound_standar = essay_scalar.transform(X_cv['compound'].values.reshape(-1, 1))

print(train_compound_standar.shape, y_train.shape)
print(test_compound_standar.shape, y_test.shape)
print(cv_compound_standar.shape, y_cv.shape)

```

```

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

```

1.5.4 Merging all the above features

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

In [118]:

```
# combine all the numerical data together
```

In [119]:

```

#categrical data --category
print("Shape of train matrix after one hot encodig ",train_categories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",cv_categories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",test_categories_one_hot.shape)

```

```

Shape of train matrix after one hot encodig (22445, 9)
Shape of train matrix after one hot encodig (11055, 9)
Shape of train matrix after one hot encodig (16500, 9)

```

In [120]:

```

#categrical data --subcategory
print("Shape of train matrix after one hot encodig ",train_subcategories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",cv_subcategories_one_hot.shape)
print("Shape of train matrix after one hot encodig ",test_subcategories_one_hot.shape)

```

```
Shape of train matrix after one hot encodig (22445, 20)
```

```
Shape of train matrix after one hot encoding (22445, 30)
Shape of train matrix after one hot encoding (11055, 30)
Shape of train matrix after one hot encoding (16500, 30)
```

In [121]:

```
#category --state
print("Shape of train matrix after one hot encoding ",train_state_one_hot.shape)
print("Shape of train matrix after one hot encoding ",cv_state_one_hot.shape)
print("Shape of train matrix after one hot encoding ",test_state_one_hot.shape)
```

```
Shape of train matrix after one hot encoding (22445, 51)
Shape of train matrix after one hot encoding (11055, 51)
Shape of train matrix after one hot encoding (16500, 51)
```

In [122]:

```
#category ----grade
print("Shape of train matrix after one hot encoding ",train_grade_one_hot.shape)
print("Shape of train matrix after one hot encoding ",cv_grade_one_hot.shape)
print("Shape of train matrix after one hot encoding ",test_grade_one_hot.shape)
```

```
Shape of train matrix after one hot encoding (22445, 4)
Shape of train matrix after one hot encoding (11055, 4)
Shape of train matrix after one hot encoding (16500, 4)
```

In [123]:

```
#category ----teacher
print("Shape of train matrix after one hot encoding ",train_teacher_prefix_one_hot.shape)
print("Shape of train matrix after one hot encoding ",cv_teacher_prefix_one_hot.shape)
print("Shape of train matrix after one hot encoding ",test_teacher_prefix_one_hot.shape)
```

```
Shape of train matrix after one hot encoding (22445, 5)
Shape of train matrix after one hot encoding (11055, 5)
Shape of train matrix after one hot encoding (16500, 5)
```

In [124]:

```
#bow essay
print("Shape of matrix after one hot encoding ",text_bow_essays_train.shape)
print("Shape of matrix after one hot encoding ",text_bow_essays_cv.shape)
print("Shape of matrix after one hot encoding ",text_bow_essays_test.shape)

#bow project title
print("Shape of matrix after one hot encoding ",text_bow_title_train.shape)
print("Shape of matrix after one hot encoding ",text_bow_title_cv.shape)
print("Shape of matrix after one hot encoding ",text_bow_title_test.shape)
```

```
Shape of matrix after one hot encoding (22445, 5000)
Shape of matrix after one hot encoding (11055, 5000)
Shape of matrix after one hot encoding (16500, 5000)
Shape of matrix after one hot encoding (22445, 636)
Shape of matrix after one hot encoding (11055, 636)
Shape of matrix after one hot encoding (16500, 636)
```

In [125]:

```
#bow essay tfidf
print("Shape of matrix after one hot encoding ",tfidf_essay_train.shape)
print("Shape of matrix after one hot encoding ",tfidf_essay_cv.shape)
print("Shape of matrix after one hot encoding ",tfidf_essay_test.shape)

#bow project title
print("Shape of matrix after one hot encoding ",text_tfidf_title_train.shape)
print("Shape of matrix after one hot encoding ",text_tfidf_title_cv.shape)
print("Shape of matrix after one hot encoding ",text_tfidf_title_test.shape)
```

```
Shape of matrix after one hot encoding (22445, 5000)
```

```
Shape of matrix after one hot encoding (11055, 5000)
Shape of matrix after one hot encoding (16500, 5000)
Shape of matrix after one hot encoding (22445, 636)
Shape of matrix after one hot encoding (11055, 636)
Shape of matrix after one hot encoding (16500, 636)
```

Assignment 5: Logistic Regression

[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) Hyper parameter tuning (find best hyper parameters corresponding the algorithm that you choose) Find the best hyper parameter which will give the maximum AUC value Find the best hyper parameter 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 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 with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps. [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. 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 parameter as suggested in step 2 and step 3 Conclusion 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 link](#)

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](#).

2. Logistic Regression

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

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

In [126]:

```
# 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 [127]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train1=hstack(((train_categories_one_hot,train_subcategories_one_hot,train_state_one_hot,train_gr
```

```
ade_one_hot,
        train_teacher_prefix_one_hot,text_bow_essays_train,text_bow_title_train,train_pri
ce_standar,
        train_quantity_standar,train_prev_proj_standar))).tocsr()
X_train1.shape
```

Out[127]:

(22445, 5738)

In [128]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_cv1=hstack(((cv_categories_one_hot,cv_subcategories_one_hot,cv_state_one_hot,cv_grade_one_hot,
               cv_teacher_prefix_one_hot,text_bow_essays_cv,text_bow_title_cv,cv_price_standar,
               cv_quantity_standar,cv_prev_proj_standar))).tocsr()
X_cv1.shape
```

Out[128]:

(11055, 5738)

In [129]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_test1=hstack(((test_categories_one_hot,test_subcategories_one_hot,test_state_one_hot,test_grade_o
ne_hot,
               test_teacher_prefix_one_hot,text_bow_essays_test,text_bow_title_test,
               test_price_standar,test_quantity_standar,test_prev_proj_standar
               ))).tocsr()
X_test1.shape
```

Out[129]:

(16500, 5738)

In [130]:

```
print(X_train1.shape,y_train.shape)
print(X_cv1.shape,y_cv.shape)
print(X_test1.shape,y_test.shape)
```

```
(22445, 5738) (22445,)
(11055, 5738) (11055,)
(16500, 5738) (16500,)
```

In [131]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
tive class
    # not the predicted outputs

    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
    # we will be predicting for the last data points
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

In [132]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
```

```

train_auc = []
cv_auc = []
C= [10**-4, 10**-3, 10**-2, 10**-1, 1, 10**1, 10**2, 10**3, 10**4]
for i in C:
    model = LogisticRegression(C=i,class_weight='balanced')
    model.fit(X_train1, y_train)

    y_train_pred = batch_predict(model,X_train1)
    y_cv_pred = batch_predict(model,X_cv1)

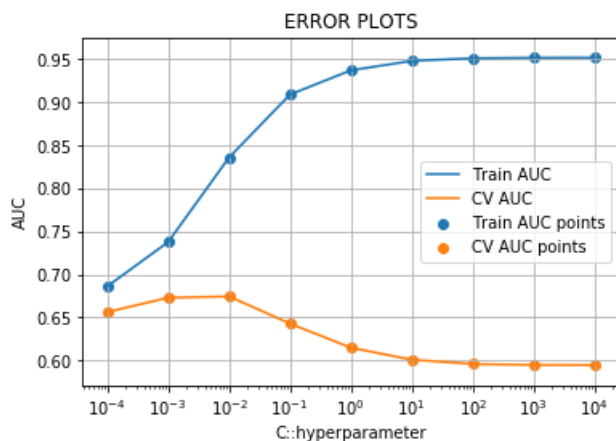
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    tive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(C, train_auc, label='Train AUC')
plt.plot(C, cv_auc, label='CV AUC')

plt.scatter(C, train_auc, label='Train AUC points')
plt.scatter(C, cv_auc, label='CV AUC points')

plt.legend()
plt.xscale('log')
plt.xlabel("C:hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [133]:

```
best_C = 0.01
```

In [134]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression

model=LogisticRegression(C=best_C,class_weight='balanced')
model.fit(X_train1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

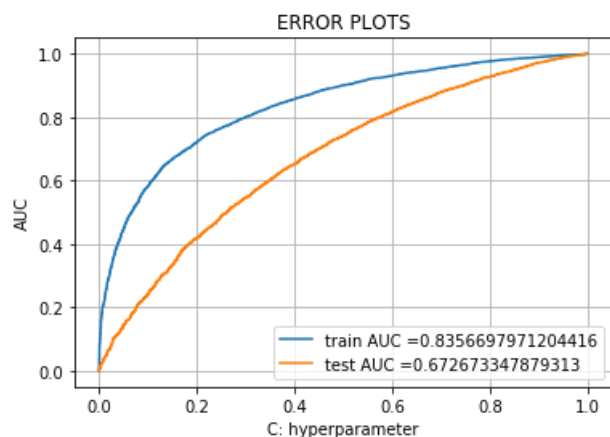
y_train_pred = batch_predict(model,X_train1)
y_test_pred = batch_predict(model, X_test1)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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



In [135]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(fpr*(1-tpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [136]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399
[[ 1732  1731]
 [ 1888 17094]]
```

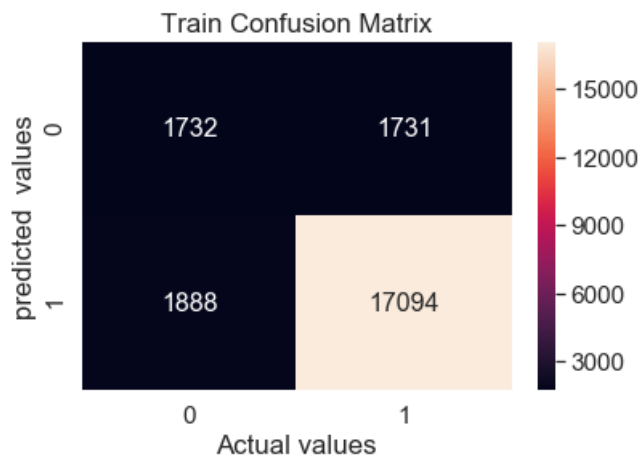
In [137]:

```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred,
tr_thresholds, train_fpr, train_fpr)),
                                     range(2),range(2))
sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix , annot = True, annot_kws={"size":16}, fmt = 'd')# font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Train Confusion Matrix')
```

```
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399
```

Out[137]:

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



```
In [138]:
```

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.443
[[ 1035  1511]
 [ 2637 11317]]
```

```
In [139]:
```

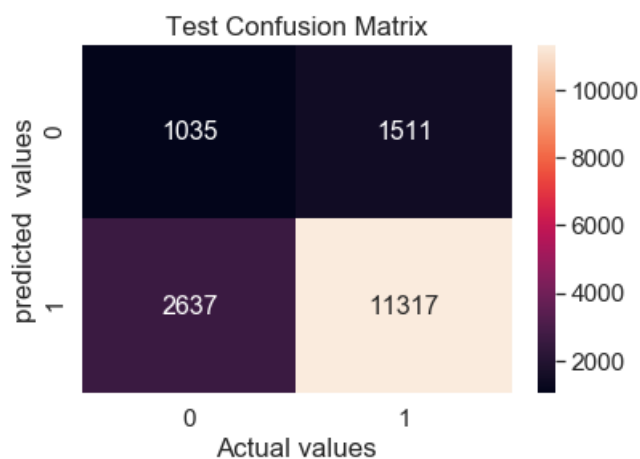
```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
                                                                    range(2), range(2)),
                                                                    test_fpr, test_fpr)),
                                     range(2), range(2))

sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Test Confusion Matrix')
```

```
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.443
```

```
Out[139]:
```

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



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


```
and max_features=5000,
```

In [140]:

```
# 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_train2=hstack(((train_categories_one_hot,train_subcategories_one_hot,train_state_one_hot,train_grade_one_hot,
                    train_teacher_prefix_one_hot,train_price_standard,train_quantity_standard,
                    train_prev_proj_standard
                    ,avg_w2v_essay_train_data,avg_w2v_project_title_train_data))).tocsr()
X_train2.shape
```

Out[140]:

```
(22445, 702)
```

In [141]:

```
# 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_cv2=hstack(((cv_categories_one_hot,cv_subcategories_one_hot,cv_state_one_hot,cv_grade_one_hot,
cv_teacher_prefix_one_hot
                    ,cv_price_standard,cv_quantity_standard,cv_prev_proj_standard,
                    avg_w2v_essay_cv_data,avg_w2v_project_title_cv_data))).tocsr()
X_cv2.shape
```

Out[141]:

```
(11055, 702)
```

In [142]:

```
# 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_test2=hstack(((test_categories_one_hot,test_subcategories_one_hot,test_state_one_hot,test_grade_one_hot,
                    test_teacher_prefix_one_hot,
test_price_standard,test_quantity_standard,test_prev_proj_standard,
                    avg_w2v_essay_test_data,avg_w2v_project_title_test_data))).tocsr()
X_test2.shape
```

Out[142]:

```
(16500, 702)
```

In [143]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs

    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
    # we will be predicting for the last data points
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

In [144]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
```

```

train_auc = []
cv_auc = []
C= [10**-4, 10**-3, 10**-2, 10**-1, 1, 10**1, 10**2, 10**3, 10**4]
for i in C:
    model = LogisticRegression(C=i,class_weight='balanced')
    model.fit(X_train2, y_train)

    y_train_pred = batch_predict(model,X_train2)
    y_cv_pred = batch_predict(model,X_cv2)

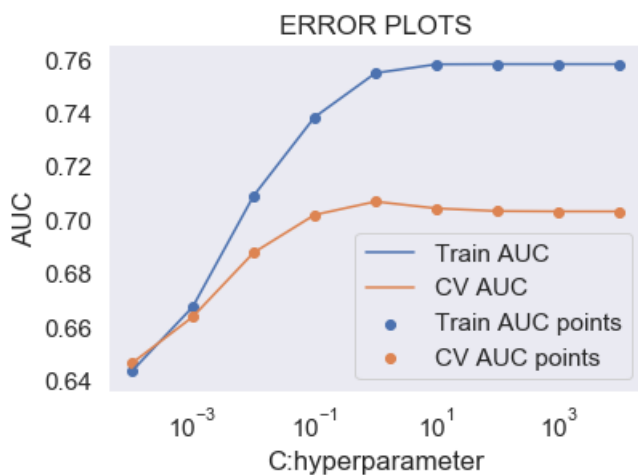
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(C, train_auc, label='Train AUC')
plt.plot(C, cv_auc, label='CV AUC')

plt.scatter(C, train_auc, label='Train AUC points')
plt.scatter(C, cv_auc, label='CV AUC points')

plt.legend()
plt.xscale('log')
plt.xlabel("C:hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [145]:

```
best_C = 0.01
```

In [146]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression

model=LogisticRegression(C=best_C,class_weight='balanced')
model.fit(X_train1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

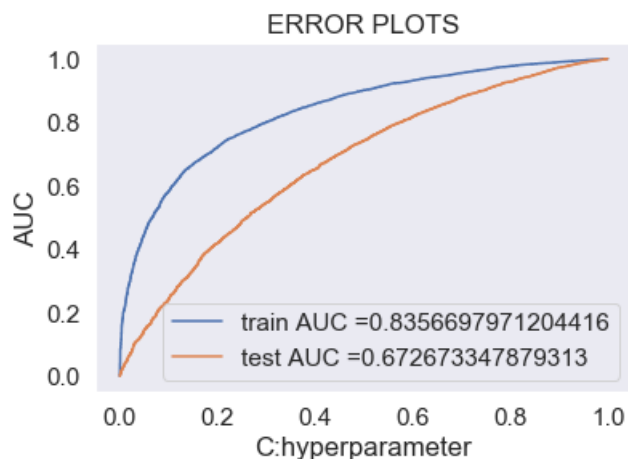
y_train_pred = batch_predict(model,X_train1)
y_test_pred = batch_predict(model, X_test1)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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



In [147]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(fpr*(1-tpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [148]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399
[[ 1732  1731]
 [ 1888 17094]]
```

In [149]:

```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred,
                                                                    tr_thresholds, train_fpr, train_fpr)),
                                      range(2), range(2))

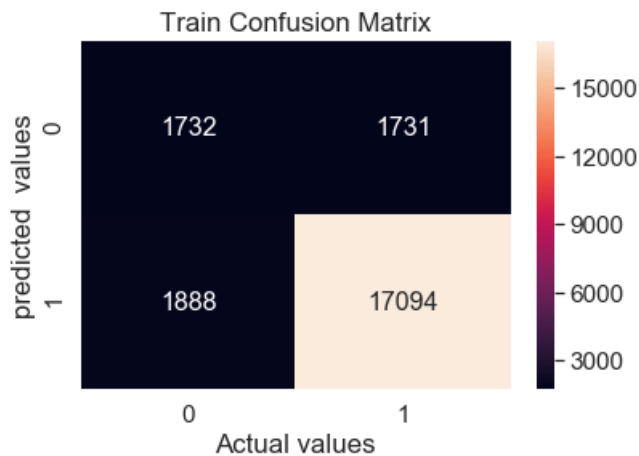
sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')# font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Train Confusion Matrix')
```

```
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399
```

Out [149]:

Out[149]:

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



In [150]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 0.443
[[1035 1511]
 [2637 11317]]

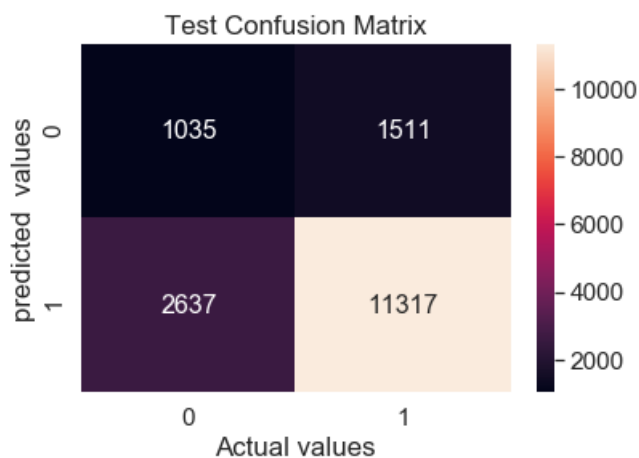
In [151]:

```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
                                                                    test_fpr, test_fpr)),
                                      range(2), range(2))
sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Test Confusion Matrix')
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 0.443

Out[151]:

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



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

In [152]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train3=hstack(((train_categories_one_hot,train_subcategories_one_hot,train_state_one_hot,train_grade_one_hot,
                    train_teacher_prefix_one_hot,train_price_standard,train_quantity_standard,
                    train_prev_proj_standard,
                    tfidf_w2v_essay_train_data,tfidf_w2v_train_project_title))).tocsr()
X_train3.shape
```

Out[152]:

(22445, 702)

In [153]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_cv3=hstack(((cv_categories_one_hot,cv_subcategories_one_hot,cv_state_one_hot,cv_grade_one_hot,cv_teacher_prefix_one_hot,
                cv_price_standard,cv_quantity_standard,cv_prev_proj_standard,
                tfidf_w2v_essay_cv_data,tfidf_w2v_cv_project_title))).tocsr()
X_cv3.shape
```

Out[153]:

(11055, 702)

In [154]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_test3=hstack(((test_categories_one_hot,test_subcategories_one_hot,test_state_one_hot,test_grade_one_hot,
                  test_teacher_prefix_one_hot,
                  test_price_standard,test_quantity_standard,test_prev_proj_standard,
                  tfidf_w2v_essay_test_data,tfidf_w2v_test_project_title))).tocsr()
X_test3.shape
```

Out[154]:

(16500, 702)

In [155]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression

train_auc = []
cv_auc = []
C = [10**-4, 10**-3, 10**-2, 10**-1, 1, 10**1, 10**2, 10**3, 10**4]
for i in C:
    model = LogisticRegression(C=i,class_weight='balanced')
    model.fit(X_train3, y_train)

    y_train_pred =batch_predict(model,X_train3)
    y_cv_pred =batch_predict(model,X_cv3)

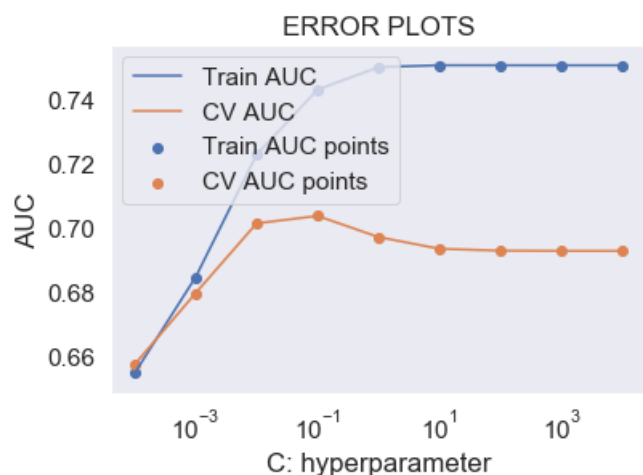
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(C, train_auc, label='Train AUC')
plt.plot(C, cv_auc, label='CV AUC')

plt.scatter(C, train_auc, label='Train AUC points')
plt.scatter(C, cv_auc, label='CV AUC points')

plt.legend()
plt.xscale('log')
```

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



In [156]:

```
best_C = 0.01
```

In [157]:

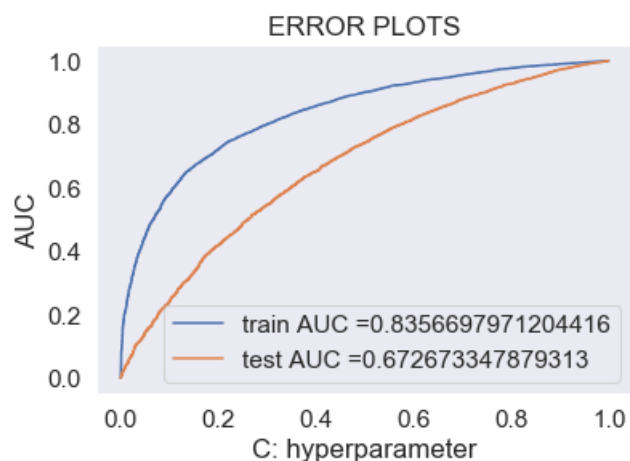
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression

model=LogisticRegression(C=best_C,class_weight='balanced')
model.fit(X_train1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(model,X_train1)
y_test_pred = batch_predict(model, X_test1)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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



In [158]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(fpr*(1-tpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [159]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399
[[ 1732  1731]
 [ 1888 17094]]
```

In [160]:

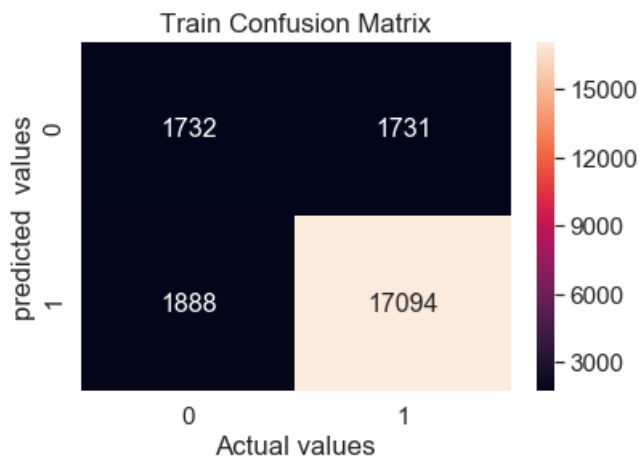
```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred,
                                                                    range(2), range(2))),
                                      range(2), range(2))

sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Train Confusion Matrix')
```

the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.399

Out[160]:

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



In [161]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.443
[[1035 1511]
[2637 11317]]

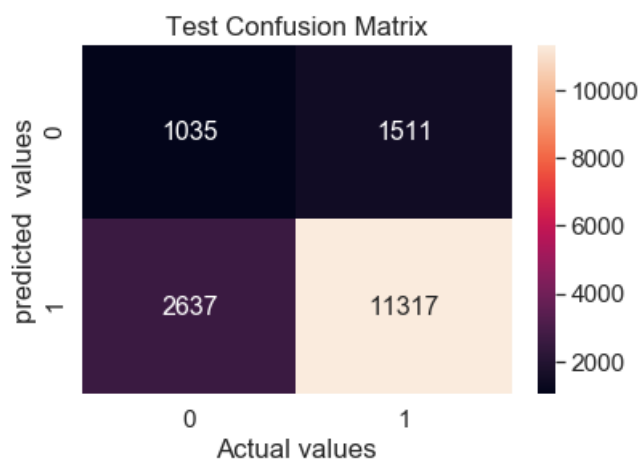
In [162]:

```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
test_fpr, test_fpr)),
                                     range(2), range(2))
sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Test Confusion Matrix')
```

the maximum value of tpr*(1-fpr) 0.25 for threshold 0.443

Out[162]:

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



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

In [163]:

```
# 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_train4=hstack(((train_categories_one_hot,train_subcategories_one_hot,train_state_one_hot,train_grade_one_hot,
                    train_teacher_prefix_one_hot,train_price_standar,train_quantity_standar,
                    train_prev_proj_standar,
                    tfidf_w2v_essay_train_data,tfidf_w2v_train_project_title))).tocsr()
X_train4.shape
```

Out[163]:

(22445, 702)

In [164]:

```
# 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_cv4=hstack(((cv_categories_one_hot,cv_subcategories_one_hot,cv_state_one_hot,cv_grade_one_hot,
                cv_teacher_prefix_one_hot,cv_price_standar,cv_quantity_standar,
```



```

cv_teacher_prefix_one_hot,cv_price_standar,cv_quantity_standar,
cv_prev_proj_standar,
tfidf_w2v_essay_cv_data,tfidf_w2v_cv_project_title))).tocsr()
X_cv4.shape

```

Out[164]:

```

(11055, 702)

```

In [165]:

```

# 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_test4=hstack(((test_categories_one_hot,test_subcategories_one_hot,test_state_one_hot,test_grade_c
ne_hot,
test_teacher_prefix_one_hot,test_price_standar,test_quantity_standar,
test_prev_proj_standar,
tfidf_w2v_essay_test_data,tfidf_w2v_test_project_title))).tocsr()
X_test4.shape

```

Out[165]:

```

(16500, 702)

```

In [166]:

```

import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
C = [10**-4, 10**-3, 10**-2, 10**-1, 1, 10**1, 10**2, 10**3, 10**4]
for i in C:
    model = LogisticRegression(C=i,class_weight='balanced')
    model.fit(X_train4, y_train)

    y_train_pred = batch_predict(model,X_train4)
    y_cv_pred = batch_predict(model,X_cv4 )

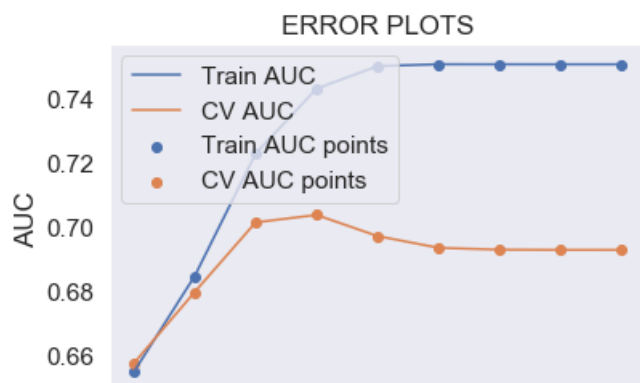
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
tive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(C, train_auc, label='Train AUC')
plt.plot(C, cv_auc, label='CV AUC')

plt.scatter(C, train_auc, label='Train AUC points')
plt.scatter(C, cv_auc, label='CV AUC points')

plt.legend()
plt.xscale('log')
plt.xlabel("C ::hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



10^{-3} 10^{-1} 10^1 10^3
 C::hyperparameter

In [167]:

```
best_C = 0.1
```

In [168]:

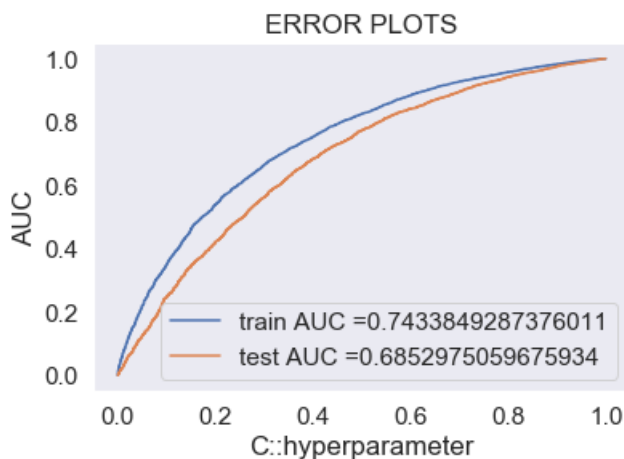
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression

model=LogisticRegression(C=best_C,class_weight='balanced')
model.fit(X_train4, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(model,X_train4)
y_test_pred = batch_predict(model, X_test4)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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



In [169]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(fpr*(1-tpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [170]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.406
[[ 1732  1731]
 [ 3329 15653]]
```

In [171]:

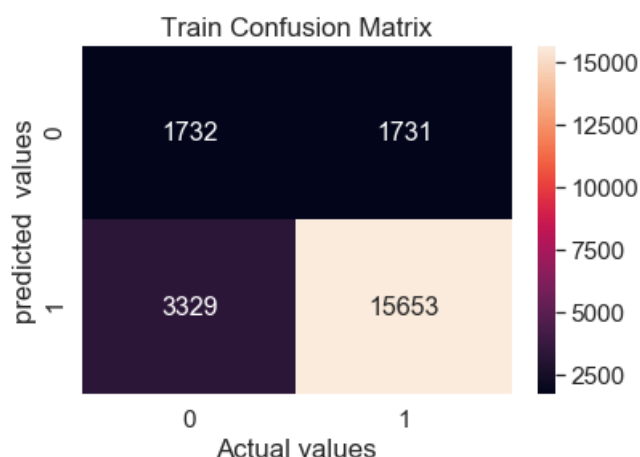
```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred,
                                                                    tr_thresholds, train_fpr, train_fpr)),
                                      range(2), range(2))

sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Train Confusion Matrix')
```

the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.406

Out[171]:

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



In [172]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.466
[[1460 1086]
 [4042 9912]]
```

In [173]:

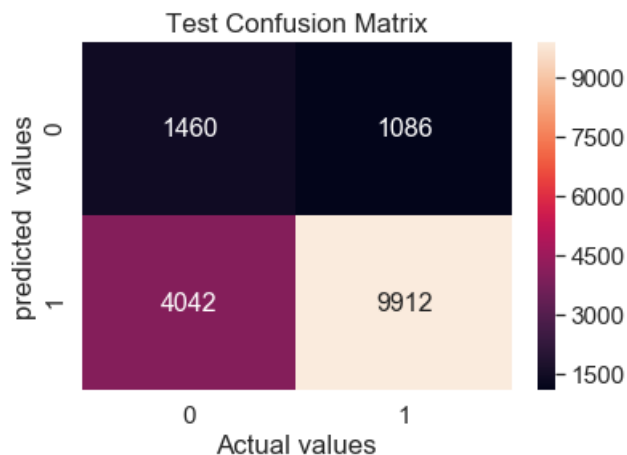
```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
                                                                    test_fpr, test_fpr)),
                                      range(2), range(2))

sns.set(font_scale=1.4) #for label size
sns.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd') # font size
plt.xlabel('Actual values')
plt.ylabel('predicted values')
plt.title('Test Confusion Matrix')
```

the maximum value of $\text{tpr} \cdot (1 - \text{fpr})$ 0.25 for threshold 0.466

Out[173]:

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



Apply Logistic Regression on the below feature set Set 5 by finding the best hyper parameter

In [174]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train5=
hstack((train_categories_one_hot,train_subcategories_one_hot,train_state_one_hot,train_grade_one_hot,
train_teacher_prefix_one_hot, text_tfidf_title_train, tfidf_essay_train,
train_quantity_standar, train_prev_proj_standar, train_price_standar,train_positive_standar,
train_negative_standar, train_neutral_standar,train_compound_standar,
train_title_word_count_standar,
train_essay_word_count_standar)).tocsr()
print(X_train5.shape, y_train.shape)
print(type(X_train5))
```

```
(22445, 5744) (22445,)
<class 'scipy.sparse.csr.csr_matrix'>
```

In [175]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_cv5=hstack((cv_categories_one_hot,cv_subcategories_one_hot,cv_state_one_hot,cv_grade_one_hot,
cv_teacher_prefix_one_hot, text_tfidf_title_cv, tfidf_essay_cv,
cv_quantity_standar, cv_prev_proj_standar, cv_price_standar,cv_positive_standar,
cv_negative_standar, cv_neutral_standar,cv_compound_standar, cv_title_word_count_standar,
cv_essay_word_count_standar)).tocsr()
print(X_cv5.shape,y_cv.shape)
print(type(X_cv5))
```

```
(11055, 5744) (11055,)
<class 'scipy.sparse.csr.csr_matrix'>
```

In [176]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_test5=hstack((test_categories_one_hot,test_subcategories_one_hot,test_state_one_hot,test_grade_one_hot,
test_teacher_prefix_one_hot, text_tfidf_title_test, tfidf_essay_test,
test_quantity_standar, test_prev_proj_standar, test_price_standar,test_positive_standar,
test negative standar, test neutral standar,test compound standar, test title word count standar,
```

```
test_essay_word_count_standar)).tocsr()
print(X_test5.shape,y_test.shape)
print(type(X_test5))
```

```
(16500, 5744) (16500,)
<class 'scipy.sparse.csr.csr_matrix'>
```

Hyperparameter Tunning

In [177]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
C = [10**-4, 10**-3, 10**-2, 10**-1, 1, 10**1, 10**2, 10**3, 10**4]
for i in C:
    model = LogisticRegression(C=i,class_weight='balanced')
    model.fit(X_train5, y_train)

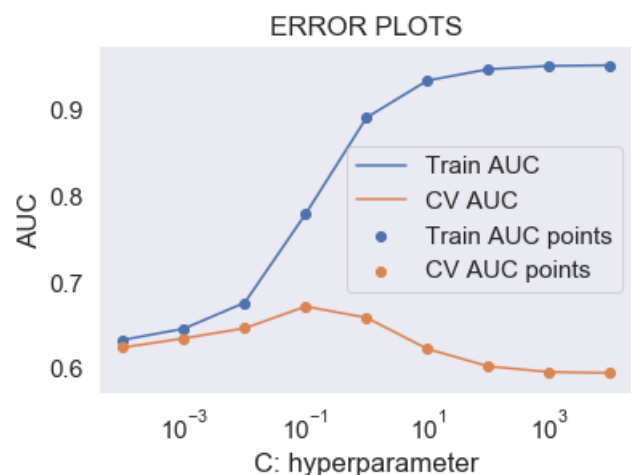
    y_train_pred = batch_predict(model,X_train5)
    y_cv_pred = batch_predict(model,X_cv5)

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(C, train_auc, label='Train AUC')
plt.plot(C, cv_auc, label='CV AUC')

plt.scatter(C, train_auc, label='Train AUC points')
plt.scatter(C, cv_auc, label='CV AUC points')

plt.legend()
plt.xscale('log')
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [178]:

```
best_c5 = 0.01
```

In [179]:

```
# https://scikit-
```

```

learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression

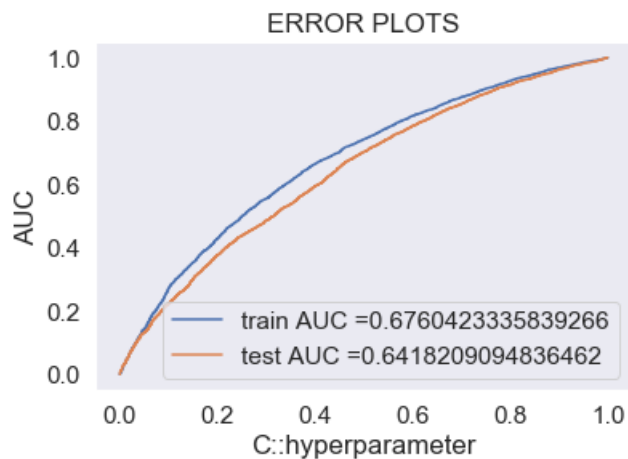
model=LogisticRegression(C=best_c5,class_weight='balanced')
model.fit(X_train5, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

y_train_pred = batch_predict(model,X_train5)
y_test_pred = batch_predict(model, X_test5)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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.xlabel("C::hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



Train Confusion Matrix

In [180]:

```

from sklearn.metrics import confusion_matrix
import seaborn as sea
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))

```

```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.466
[[ 1732  1731]
 [ 4907 14075]]

```

In [181]:

```

train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test,predict(y_test_pred, tr_thresholds,
                                                                    test_fpr,test_fpr)),
range(2),range(2))
sea.set(font_scale=1.4)
sea.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")

```

```

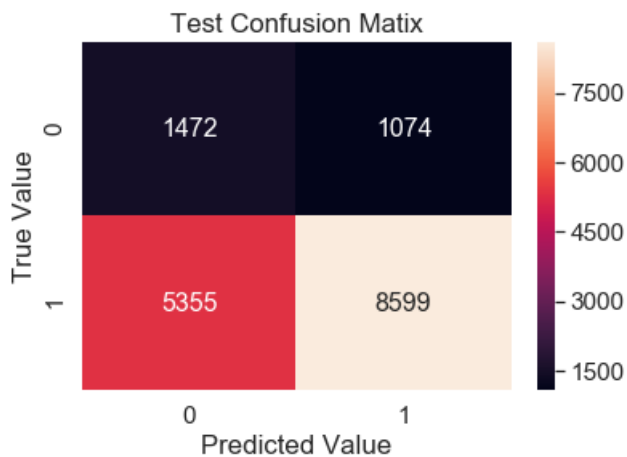
the maximum value of tpr*(1-fpr) 0.24999984572938835 for threshold 0.5

```

Out[181]:

Out[181]:

Text(0.5, 1.0, 'Test Confusion Matix')



Test Confusion Matrix

In [182]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.24999984572938835 for threshold 0.5
[[1472 1074]
 [5355 8599]]

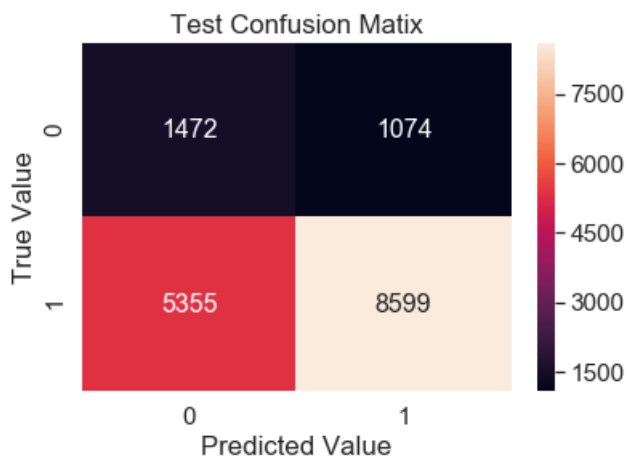
In [183]:

```
train_confusion_matrix = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds,
                                                                    test_fpr, test_fpr)),
                                     range(2), range(2))
sea.set(font_scale=1.4)
sea.heatmap(train_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999984572938835 for threshold 0.5

Out[183]:

Text(0.5, 1.0, 'Test Confusion Matix')



2.2 Make Data Model Ready: encoding numerical, categorical features

In [184]:

```
# 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
```

2.3 Make Data Model Ready: encoding eassay, and project_title

In [185]:

```
# 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
```

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

In [186]:

```
# 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
```

2.5 Logistic Regression with added Features `Set 5

In [187]:

```
# 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
```

Summary

Summary

In [190]:

```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model", " hyperParameter", "AUC"]
x.add_row(["BOW", "Auto", "0.01", "0.83"])
x.add_row(["TFIDF", "Auto", "0.01", "0.83"])
x.add_row(["AVGW2V", "Auto", "0.01", "0.83"])
x.add_row(["TFIF-2V", "Auto", "0.5", "0.74"])
x.add_row(["SET5", "Auto", "0.01", "0.67"])
print(x)
```

Vectorizer	Model	hyperParameter	AUC
BOW	Auto	0.01	0.83
TFIDF	Auto	0.01	0.83
AVGW2V	Auto	0.01	0.83
TFIF-2V	Auto	0.5	0.74
SET5	Auto	0.01	0.67

Observations

In [189]:

From above we can the BOW **and** TFIDF encoding contain the Project_Essays **and** Project_titles **in** those models.
So, that we can say that Text Data also plays major role **in** predicting the output.
The Set 5 which we built model on numerical features only performs badly compared to **all** the Models which having text data.

File "<ipython-input-189-a5fb3ba3fa79>", line 1

From above we can the BOW and TFIDF encoding contain the Project_Essays and Project_titles in those models.

SyntaxError: invalid syntax