

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	Title of the project. Examples: <ul style="list-style-type: none">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: <ul style="list-style-type: none">Grades PreK-2Grades 3-5Grades 6-8Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none">Applied LearningCare & HungerHealth & SportsHistory & CivicsLiteracy & LanguageMath & ScienceMusic & The ArtsSpecial NeedsWarmth Examples: <ul style="list-style-type: none">Music & The ArtsLiteracy & Language, Math & Science

school_state	State where school is located (Two-letter U.S. postal code). Example: WY
project_subject_subcategories	One or more (comma-separated) subject subcategories for the project. Examples: <ul style="list-style-type: none"> Literacy Literature & Writing, Social Sciences
project_resource_summary	An explanation of the resources needed for the project. Example: <ul style="list-style-type: none"> My students need hands on literacy materials to manage sensory needs!
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	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: <code>p036502</code>
<code>description</code>	Description of the resource. Example: <code>Tenor Saxophone Reeds, Box of 25</code>
<code>quantity</code>	Quantity of the resource required. Example: <code>3</code>
<code>price</code>	Price of the resource required. Example: <code>9.95</code>

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 <code>0</code> indicates the project was not approved, and a value of <code>1</code> 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

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

1.1 Reading Data

In [2]:

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

In [3]:

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

```
Number of data points in train data (109248, 17)
```

```
-----
----
```

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

In [4]:

```
print("Number of data points in train data", resource_data.shape)
```

```
print(resource_data.columns.values)
resource_data.head(2)
```

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

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

Out[4]:

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

1.2 preprocessing of project_subject_categories

In [5]:

```
categories = list(project_data['project_subject_categories'].
values)
# remove special characters from list of strings python: http://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 placing 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 [6]:

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

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

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth , Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"

        j=j.replace('The', '') # if we have the words "The " we are going to replace it with ''(i.e removing 'The')
        j = j.replace(' ', '') # we are placing 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 preprocessing of school_state

In [7]:

```
school_states = list(project_data['school_state'].values)
# remove special characters from list of strings python: http://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

school_states_list = []
for i in school_states:
    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('&', '_')
    school_states_list.append(temp.strip())
```

```
project_data.drop(['school_state'], axis=1, inplace=True)
project_data['school_state'] = school_states_list
```

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

```
my_counter = Counter()
```

```
for word in project_data['school_state'].values:
    my_counter.update(word.split())
```

```
school_states_dict = dict(my_counter)
```

```
sorted_school_states_dict = dict(sorted(school_states_dict.items(), key=lambda kv: kv[1]))
```

1.3 preprocessing of teacher_prefix

In [8]:

```
#NaN values in teacher prefix will create a problem while encoding, so we replace NaN values with the mode of that particular column  
#removing dot(.) since it is a special character  
mode_of_teacher_prefix = project_data['teacher_prefix'].value_counts().index[0]  
  
project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna(mode_of_teacher_prefix)
```

In [9]:

```
prefixes = []  
  
for i in range(len(project_data)):  
    a = project_data["teacher_prefix"][i].replace(".", "")  
    prefixes.append(a)
```

In [10]:

```
project_data.drop(['teacher_prefix'], axis = 1, inplace = True)  
project_data["teacher_prefix"] = prefixes  
print("After removing the special characters ,Column values:  
")  
np.unique(project_data["teacher_prefix"].values)
```

After removing the special characters ,Column values:

Out[10]:

```
array(['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher'], dt  
type=object)
```


1.3 preprocessing of project_grade_category

In [11]:

```
# We need to get rid of The spaces between the text and the h  
yphens because they're special characters.  
#Rmoving multiple characters from a string in Python  
#https://stackoverflow.com/questions/3411771/multiple-charact  
er-replace-with-python  
  
project_grade_category = []  
  
for i in range(len(project_data)):  
    a = project_data["project_grade_category"][i].replace(" "  
    , "_").replace("-", "_")  
    project_grade_category.append(a)
```

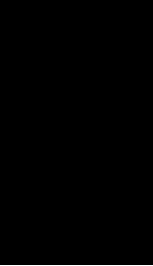
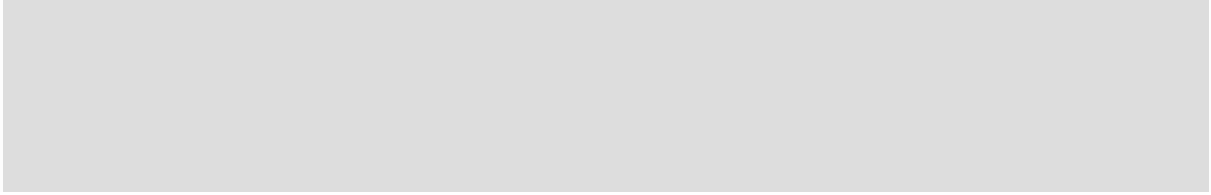
In [12]:

```
project_data.drop(['project_grade_category'], axis = 1, inplace = True)  
project_data["project_grade_category"] = project_grade_category  
print("After removing the special characters ,Column values:  
")  
np.unique(project_data["project_grade_category"].values)
```

After removing the special characters ,Column
values:

Out[12]:

```
array(['Grades_3_5', 'Grades_6_8', 'Grades_9_1  
2', 'Grades_PreK_2'], dtype=object)
```



1.3 Text preprocessing

In [13]:

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

In [14]:

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

Out[14]:

Unnamed: 0	id	teacher_id
0		
	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc
1		
	140945 p258326	897464ce9ddc600bced1151f324dd63a



In [15]:

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

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school.

Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English along side of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All fami

lies with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch.

The videos are to help the child develop early reading skills.

Parents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.

=====

====

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students.

The school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity.

My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will ut

ilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.

Whenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them.

We ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.

=====
=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.

My class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas. They attend a Title I school.

hool, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.

Your generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.

It costs a lot of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!

nannan

=====
=====

My kindergarten students have varied disabilities ranging from speech and language delays, c

ognitive delays, gross/fine motor delays, to a
utism. They are eager beavers and always striv
e to work their hardest working past their lim
itations. \r\n\r\nThe materials we have are th
e ones I seek out for my students. I teach in
a Title I school where most of the students re
ceive free or reduced price lunch. Despite th
eir disabilities and limitations, my students
love coming to school and come eager to learn
and explore. Have you ever felt like you had an
ts in your pants and you needed to groove and
move as you were in a meeting? This is how my
kids feel all the time. They want to be able to
move as they learn or so they say. Wobble chai
rs are the answer and I love them because they
develop their core, which enhances gross moto
r and in turn fine motor skills. \r\nThey also
want to learn through games, my kids don't wa
nt to sit and do worksheets. They want to lear
n to count by jumping and playing. Physical en
gagement is the key to our success. The number
toss and color and shape mats can make that h
appen. My students will forget they are doing
work and just have the fun a 6 year old deserv
es. nannan

=====
=====

The mediocre teacher tells. The good teacher e
xplains. The superior teacher demonstrates. Th
e great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup
is 97.6% African-American, making up the larg
est segment of the student body. A typical sch
ool in Dallas is made up of 23.2% African-Amer
ican students. Most of the students are on fre
e or reduced lunch. We aren't receiving doctor
s, lawyers, or engineers children from rich ba

ckgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.

The cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.

nannan
=====

In [16]:

```
# 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 [17]:

```

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

```

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

ves.nannan

=====
=====

In [18]:

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

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

and just have the fun a 6 year old deserves.nannan

In [19]:

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

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

In [20]:

```
# 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', 'b
etween', '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', "co
uldn'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 [21]:

```
# Combining all the above stundents
```

```

from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\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_essays.append(sent.lower().strip())

```

```

100%|██████████| 109248/109248 [00:45<00:00, 2
414.31it/s]

```

In [22]:

```

# after preprocesing
preprocessed_essays[20000]

```

Out[22]:

```

'my kindergarten students varied disabilities
ranging speech language delays cognitive delay
s gross fine motor delays autism they eager be
avers always strive work hardest working past
limitations the materials ones i seek students
i teach title i school students receive free
reduced price lunch despite disabilities limit
ations students love coming school come eager
learn explore have ever felt like ants pants n
eeded groove move meeting this kids feel time
the want able move learn say wobble chairs ans
wer i love develop core enhances gross motor t
urn fine motor skills they also want learn gam
es kids not want sit worksheets they want lear
n count jumping playing physical engagement ke

```

y success the number toss color shape mats mak
e happen my students forget work fun 6 year ol
d deserves nannan'

In [23]:

```
#creating a new column with the preprocessed essays and repla  
cing it with the original columns  
project_data['preprocessed_essays'] = preprocessed_essays  
project_data.drop(['project_essay_1'], axis=1, inplace=True)  
project_data.drop(['project_essay_2'], axis=1, inplace=True)  
project_data.drop(['project_essay_3'], axis=1, inplace=True)  
project_data.drop(['project_essay_4'], axis=1, inplace=True)
```

In [24]:

```
essay_word_count=[]  
for i in range(len(project_data['preprocessed_essays'])):  
    essay_word_count.append(len(project_data['preprocessed_es  
says'][i].split()))
```

In [25]:

```
project_data['preprocessed_essays'][1].split()
```

Out[25]:

```
['our',  
 'students',  
 'arrive',  
 'school',  
 'eager',  
 'learn',  
 'they',  
 'polite',  
 'generous',  
 'strive',  
 'best',  
 'they',
```

'know',
'education',
'succeed',
'life',
'help',
'improve',
'lives',
'our',
'school',
'focuses',
'families',
'low',
'incomes',
'tries',
'give',
'student',
'education',
'deserve',
'while',
'not',
'much',
'students',
'use',
'materials',
'given',
'best',
'the',
'projector',
'need',
'school',
'crucial',
'academic',
'improvement',
'students',
'as',
'technology',
'continues',

'grow',
'many',
'resources',
'internet',
'teachers',
'use',
'growth',
'students',
'however',
'school',
'limited',
'resources',
'particularly',
'technology',
'without',
'disadvantage',
'one',
'things',
'could',
'really',
'help',
'classrooms',
'projector',
'with',
'projector',
'not',
'crucial',
'instruction',
'also',
'growth',
'students',
'with',
'projector',
'show',
'presentations',
'documentaries',
'photos',

```
'historical',  
'land',  
'sites',  
'math',  
'problems',  
'much',  
'with',  
'projector',  
'make',  
'teaching',  
'learning',  
'easier',  
'also',  
'targeting',  
'different',  
'types',  
'learners',  
'classrooms',  
'auditory',  
'visual',  
'kinesthetic',  
'etc',  
'nannan']
```

In [26]:

```
len(project_data['preprocessed_essays'][1].split())
```

Out[26]:

109

In [27]:

```
essay_word_count[1]
```

Out[27]:

109

In [28]:

```
project_data['essay_word_count'] = essay_word_count
```

In [29]:

```
#simple example on how to use nltk's sentiment intensity analyzer
```

```
import nltk
```

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

```
#if it shows an error, then do the following:
```

```
# import nltk
```

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

```
sid = SentimentIntensityAnalyzer()
```

```
for_sentiment = 'a person is a person no matter how small dr  
seuss i teach the smallest students with the biggest enthusia  
sm \
```

```
for learning my students learn in many different ways using a  
ll of our senses and multiple intelligences i use a wide rang  
e\
```

```
of techniques to help all my students succeed students in my  
class come from a variety of different backgrounds which make  
s\
```

```
for wonderful sharing of experiences and cultures including n  
ative americans our school is a caring community of successfu  
l \
```

```
learners which can be seen through collaborative student proj  
ect based learning in and out of the classroom kindergartener  
s \
```

```
in my class love to work with hands on materials and have man  
y different opportunities to practice a skill before it is\  
mastered having the social skills to work cooperatively with  
friends is a crucial aspect of the kindergarten curriculum\  
montana is the perfect place to learn about agriculture and n  
utrition my students love to role play in our pretend kitchen  
\
```

in the early childhood classroom i have had several kids ask me can we try cooking with real food i will take their idea \ and create common core cooking lessons where we learn important math and writing concepts while cooking delicious healthy \ food for snack time my students will have a grounded appreciation for the work that went into making the food and knowledge \ of where the ingredients came from as well as how it is healthy for their bodies this project would expand our learning of \ nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce make our own bread \ and mix up healthy plants from our classroom garden in the spring we will also create our own cookbooks to be printed and \ shared with families students will gain math and literature skills as well as a life long enjoyment for healthy cooking \
nannan'

```
ss = sid.polarity_scores(for_sentiment)
```

```
ss
```

Out[29]:

```
{'neg': 0.01, 'neu': 0.745, 'pos': 0.245, 'compound': 0.9975}
```

In [30]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer
analyzer = SentimentIntensityAnalyzer()
neg=[];pos=[];neu=[]; compound = []

for i in tqdm(range(len(project_data['preprocessed_essays']))):
```

```

sentiment_scores = analyzer.polarity_scores(project_data[
'preprocessed_essays'][i])
neg.append(sentiment_scores['neg'])
pos.append(sentiment_scores['pos'])
neu.append(sentiment_scores['neu'])
compound.append(sentiment_scores['compound'])

```

```

100%|██████████| 109248/109248 [02:20<00:00, 7
76.19it/s]

```

In [31]:

```

#new columns indicating the sentiment score of each project e
ssay
project_data['neg'] = neg
project_data['neu'] = neu
project_data['pos'] = pos
project_data['compound'] = compound

```

In [32]:

```
project_data.head(5)
```

Out[32]:

Unnamed: 0	id	teacher_id	project_submitted_date
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	2016-12-05 13:43:57
1	140945 p258326	897464ce9ddc600bced1151f324dd63a	2016-10-25 09:2

2

21895 p182444 3465aaf82da834c0582ebd0ef8040ca0

2016-08-31 12:03:56

3

45 p246581 f3cb9bffbba169bef1a77b243e620b60

2016-10-06 21:1

4

172407 p104768 be1f7507a41f8479dc06f047086a39ec

2016-07-11 01:10:09



1.4 Preprocessing of $project_{it} \leq$

In [33]:

```
# similarly you can preprocess the titles also
preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\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_titles.append(sent.lower().strip())
```

```
100%|██████████| 109248/109248 [00:01<00:00, 5
6598.82it/s]
```

In [34]:

```
#creating a new column with the preprocessed titles, useful for analysis
project_data['preprocessed_titles'] = preprocessed_titles
```

In [35]:

```
title_word_count=[]
for i in range(len(project_data['preprocessed_titles'])):
    title_word_count.append(len(project_data['preprocessed_titles'][i].split()))
```

In [36]:

```
project_data['title_word_count'] = title_word_count
```


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

In [37]:

```
from sklearn.model_selection import train_test_split
#How to split whole dataset into Train,CV and test
#https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html#sklearn.model_selection.train_test_split
project_data_train, project_data_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'], test_size=0.33, stratify = project_data['project_is_approved'])
project_data_train, project_data_cv, y_train, y_cv = train_test_split(project_data_train, y_train, test_size=0.33, stratify=y_train)
```

In [38]:

```
print("Split ratio")
print('- '*50)
print('Train dataset:',len(project_data_train)/len(project_data)*100, '%\n', 'size:', len(project_data_train))
print('Cross validation dataset:', len(project_data_cv)/len(project_data)*100, '%\n', 'size:', len(project_data_cv))
print('Test dataset:', len(project_data_test)/len(project_data)*100, '%\n', 'size:', len(project_data_test))
```

Split ratio

Train dataset: 44.889608963093146 %
size: 49041

Cross validation dataset: 22.110244581136495 %
size: 24155
Test dataset: 33.000146455770356 %
size: 36052

In [39]:

```
#Features  
project_data_train.drop(['project_is_approved'], axis=1, inplace=True)  
project_data_cv.drop(['project_is_approved'], axis=1, inplace=True)  
project_data_test.drop(['project_is_approved'], axis=1, inplace=True)
```

1.5 Preparing data for models

In [40]:

```
project_data.columns
```

Out[40]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'project_submitted_datetime',  
      'project_title', 'project_resource_summary',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'clean_categories', 'clean_subcategories', 'school_state',  
      'teacher_prefix', 'project_grade_category', 'essay',  
      'preprocessed_essays', 'essay_word_count', 'neg', 'neu', 'pos',  
      'compound', 'preprocessed_titles', 'title_word_count'],  
      dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optional)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

Make Data Model Ready: encoding numerical, categorical features

Vectorizing Categorical data

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

In [41]:

```
# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_cat.fit(project_data_train['clean_categories'].values) #fitting has to be on Train data

train_categories_one_hot = vectorizer_cat.transform(project_data_train['clean_categories'].values)
cv_categories_one_hot = vectorizer_cat.transform(project_data_cv['clean_categories'].values)
test_categories_one_hot = vectorizer_cat.transform(project_data_test['clean_categories'].values)

print(vectorizer_cat.get_feature_names())
print("Shape of training data matrix after one hot encoding ", train_categories_one_hot.shape)
print("Shape of cross validation data matrix after one hot encoding ", cv_categories_one_hot.shape)
print("Shape of test data matrix after one hot encoding ", test_categories_one_hot.shape)
```

```
t_categories_one_hot.shape)
```

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

Shape of training data matrix after one hot encoding (49041, 9)

Shape of cross validation data matrix after one hot encoding (24155, 9)

Shape of test data matrix after one hot encoding (36052, 9)

In [42]:

```
# we use count vectorizer to convert the values into one
vectorizer_subcat = CountVectorizer(vocabulary=list(sorted_subcat_dict.keys()), lowercase=False, binary=True)
vectorizer_subcat.fit(project_data_train['clean_subcategories'].values)
```

```
train_subcategories_one_hot = vectorizer_subcat.transform(project_data_train['clean_subcategories'].values)
cv_subcategories_one_hot = vectorizer_subcat.transform(project_data_cv['clean_subcategories'].values)
test_subcategories_one_hot = vectorizer_subcat.transform(project_data_test['clean_subcategories'].values)
```

```
print(vectorizer_subcat.get_feature_names())
```

```
print("Shape of train data matrix after one hot encoding ", train_subcategories_one_hot.shape)
```

```
print("Shape of cross validation data matrix after one hot encoding ", cv_subcategories_one_hot.shape)
```

```
print("Shape of test data matrix after one hot encoding ", test_subcategories_one_hot.shape)
```

```
t_subcategories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLi  
teracy', 'ParentInvolvement', 'Extracurricular  
, 'Civics_Government', 'ForeignLanguages', 'N  
utritionEducation', 'Warmth', 'Care_Hunger', '  
SocialSciences', 'PerformingArts', 'CharacterE  
ducation', 'TeamSports', 'Other', 'College_Car  
eerPrep', 'Music', 'History_Geography', 'Healt  
h_LifeScience', 'EarlyDevelopment', 'ESL', 'Gy  
m_Fitness', 'EnvironmentalScience', 'VisualArt  
s', 'Health_Wellness', 'AppliedSciences', 'Spe  
cialNeeds', 'Literature_Writing', 'Mathematics  
, 'Literacy']
```

```
Shape of train data matrix after one hot encod  
ing (49041, 30)
```

```
Shape of cross validation data matrix after on  
e hot encoding (24155, 30)
```

```
Shape of test data matrix after one hot encodi  
ng (36052, 30)
```

In [43]:

```
## we use count vectorizer to convert the values into one hot  
encoded features
```

```
vectorizer_school_state = CountVectorizer()  
vectorizer_school_state.fit(project_data_train['school_state']  
.values)
```

```
print(vectorizer_school_state.get_feature_names())
```

```
train_school_state_category_one_hot = vectorizer_school_state  
.transform(project_data_train['school_state'].values)
```

```
cv_school_state_category_one_hot = vectorizer_school_state.transform(project_data_cv['school_state'].values)
test_school_state_category_one_hot = vectorizer_school_state.transform(project_data_test['school_state'].values)
```

```
print("Shape of train data matrix after one hot encoding ",train_school_state_category_one_hot.shape)
print("Shape of cross validation data matrix after one hot encoding ",cv_school_state_category_one_hot.shape)
print("Shape of test data matrix after one hot encoding ",test_school_state_category_one_hot.shape)
```

```
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
```

```
Shape of train data matrix after one hot encoding (49041, 51)
```

```
Shape of cross validation data matrix after one hot encoding (24155, 51)
```

```
Shape of test data matrix after one hot encoding (36052, 51)
```

In [44]:

```
#This step is to initialize a vectorizer with vocab from train data
my_counter = Counter()
for project_grade in project_data_train['project_grade_category'].values:
    my_counter.update(project_grade.split())
```

In [45]:


```
project_grade_cat_dict = dict(my_counter)
sorted_project_grade_cat_dict = dict(sorted(project_grade_cat
_dict.items(), key=lambda kv: kv[1]))
```

In [46]:

```
## we use count vectorizer to convert the values into one hot
encoded features
```

```
vectorizer_grade = CountVectorizer(vocabulary=list(sorted_pro
ject_grade_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_grade.fit(project_data_train['project_grade_catego
ry'].values)
```

```
print(vectorizer_grade.get_feature_names())
```

```
train_project_grade_category_one_hot = vectorizer_grade.trans
form(project_data_train['project_grade_category'].values)
cv_project_grade_category_one_hot = vectorizer_grade.transfor
m(project_data_cv['project_grade_category'].values)
test_project_grade_category_one_hot = vectorizer_grade.transf
orm(project_data_test['project_grade_category'].values)
```

```
print("Shape of train data matrix after one hot encoding ",tr
ain_project_grade_category_one_hot.shape)
print("Shape of cross validation data matrix after one hot en
coding ",cv_project_grade_category_one_hot.shape)
print("Shape of test data matrix after one hot encoding ",tes
t_project_grade_category_one_hot.shape)
```

```
['Grades_9_12', 'Grades_6_8', 'Grades_3_5', 'G
rades_PreK_2']
```

```
Shape of train data matrix after one hot encod
ing (49041, 4)
```

```
Shape of cross validation data matrix after on
```

e hot encoding (24155, 4)
Shape of test data matrix after one hot encoding (36052, 4)

In [47]:

```
#https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-document
#ValueError: np.nan is an invalid document, expected byte or unicode string.
vectorizer_prefix = CountVectorizer()
vectorizer_prefix.fit(project_data_train['teacher_prefix'].values.astype("U"))

print(vectorizer_prefix.get_feature_names())

train_teacher_prefix_categories_one_hot = vectorizer_prefix.transform(project_data_train['teacher_prefix'].values.astype("U"))
cv_teacher_prefix_categories_one_hot = vectorizer_prefix.transform(project_data_cv['teacher_prefix'].values.astype("U"))
test_teacher_prefix_categories_one_hot = vectorizer_prefix.transform(project_data_test['teacher_prefix'].values.astype("U"))

print("Shape of train data matrix after one hot encoding ", train_teacher_prefix_categories_one_hot.shape)
print("Shape of cross validation data matrix after one hot encoding ", cv_teacher_prefix_categories_one_hot.shape)
print("Shape of test data matrix after one hot encoding ", test_teacher_prefix_categories_one_hot.shape)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
Shape of train data matrix after one hot encoding (49041, 5)
Shape of cross validation data matrix after on
```

e hot encoding (24155, 5)

Shape of test data matrix after one hot encoding (36052, 5)

Make Data Model Ready: encoding essay, and project_title

Vectorizing Text data

1.5.2.1 Bag of words

In [48]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer_bow_essay = CountVectorizer(ngram_range=(2,2), min_df=10, max_features = 5000)
vectorizer_bow_essay.fit(project_data_train['preprocessed_essays'].values) #Fitting has to be on Train data

train_essay_bow = vectorizer_bow_essay.transform(project_data_train['essay'].values)
cv_essay_bow = vectorizer_bow_essay.transform(project_data_cv['essay'].values)
test_essay_bow = vectorizer_bow_essay.transform(project_data_test['essay'].values)

print("Shape of train data matrix after one hot encoding ", train_essay_bow.shape)
print("Shape of cross validation data matrix after one hot encoding ", cv_essay_bow.shape)
print("Shape of test data matrix after one hot encoding ", test_essay_bow.shape)
```

```
t_essay_bow.shape)
```

Shape of train data matrix after one hot encoding (49041, 5000)

Shape of cross validation data matrix after one hot encoding (24155, 5000)

Shape of test data matrix after one hot encoding (36052, 5000)

In [49]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer_bow_title = CountVectorizer(ngram_range=(2,2), min_
_df=10, max_features = 5000)
vectorizer_bow_title.fit_transform(project_data_train['prepro
cessed_titles'].values)    #Fitting has to be on Train data

train_title_bow = vectorizer_bow_title.transform(project_data
_train['preprocessed_titles'].values)
cv_title_bow = vectorizer_bow_title.transform(project_data_cv
['preprocessed_titles'].values)
test_title_bow = vectorizer_bow_title.transform(project_data_
test['preprocessed_titles'].values)

print("Shape of train data matrix after one hot encoding ",tr
ain_title_bow.shape)
print("Shape of cross validation data matrix after one hot en
coding ",cv_title_bow.shape)
print("Shape of test data matrix after one hot encoding ",tes
t_title_bow.shape)
```

Shape of train data matrix after one hot encoding (49041, 1663)

Shape of cross validation data matrix after on

e hot encoding (24155, 1663)
Shape of test data matrix after one hot encoding (36052, 1663)

1.5.2.2 TFIDF vectorizer

In [50]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essay = TfidfVectorizer(ngram_range=(2,2), min_df=10, max_features = 5000)
vectorizer_tfidf_essay.fit(project_data_train['preprocessed_essays']) #Fitting has to be on Train data

train_essay_tfidf = vectorizer_tfidf_essay.transform(project_data_train['preprocessed_essays'].values)
cv_essay_tfidf = vectorizer_tfidf_essay.transform(project_data_cv['preprocessed_essays'].values)
test_essay_tfidf = vectorizer_tfidf_essay.transform(project_data_test['preprocessed_essays'].values)

print("Shape of train data matrix after one hot encoding ", train_essay_tfidf.shape)
print("Shape of cross validation data matrix after one hot encoding ", cv_essay_tfidf.shape)
print("Shape of test data matrix after one hot encoding ", test_essay_tfidf.shape)
```

Shape of train data matrix after one hot encoding (49041, 5000)
Shape of cross validation data matrix after one hot encoding (24155, 5000)
Shape of test data matrix after one hot encoding (36052, 5000)

In [51]:

```

vectorizer_tfidf_title = TfidfVectorizer(ngram_range=(2,2), m
in_df=10, max_features = 5000)
vectorizer_tfidf_title.fit(project_data_train['preprocessed_t
itles'])      #Fitting has to be on Train data

train_title_tfidf = vectorizer_tfidf_title.transform(project_
data_train['preprocessed_titles'].values)
cv_title_tfidf = vectorizer_tfidf_title.transform(project_dat
a_cv['preprocessed_titles'].values)
test_title_tfidf = vectorizer_tfidf_title.transform(project_d
ata_test['preprocessed_titles'].values)

print("Shape of train data matrix after one hot encoding ",tr
ain_title_tfidf.shape)
print("Shape of cross validation data matrix after one hot en
coding ",cv_title_tfidf.shape)
print("Shape of test data matrix after one hot encoding ",tes
t_title_tfidf.shape)

```

Shape of train data matrix after one hot encoding (49041, 1663)

Shape of cross validation data matrix after one hot encoding (24155, 1663)

Shape of test data matrix after one hot encoding (36052, 1663)

1.5.2.3 Using Pretrained Models: Avg W2V

In [52]:

```

'''
# 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 preproced_texts:
    words.extend(i.split(' '))

for i in preproced_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 vec
tors and our coupus", \
    len(inter_words), "(" ,np.round(len(inter_words)/len(word
s)*100,3), "%)")

words_courpus = {}

```



```

words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/

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

'''

```

Out[52]:

```

'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef load
GloveModel(gloveFile):\n    print ("Loading Gl
ove Model")\n    f = open(gloveFile,\r', enc
oding="utf8")\n    model = {}\n    for line in
tqdm(f):\n        splitLine = line.split()\n
        word = splitLine[0]\n        embedding
= np.array([float(val) for val in splitLine[1:
]])\n        model[word] = embedding\n    prin
t ("Done.",len(model)," words loaded!")\n    r
eturn model\nmodel = loadGloveModel(\'glove.42
B.300d.txt\')\n\n# =====
\n\nOutput:\n    \nLoading Glove Model\n1917495
it [06:32, 4879.69it/s]\nDone. 1917495 words
loaded!\n\n# =====\n\nw
ords = []\nfor i in preproced_texts:\n    word
s.extend(i.split(\' \'))\n\nfor i in preproced
_titles:\n    words.extend(i.split(\' \'))\npr

```

```

int("all the words in the coupus", len(words))
\nwords = set(words)\nprint("the unique words
in the coupus", len(words))\n\ninter_words = s
et(model.keys()).intersection(words)\nprint("T
he number of words that are present in both gl
ove vectors and our coupus",      len(inter_w
ords), "(" ,np.round(len(inter_words)/len(words)
*100,3), "%)")\n\nwords_courpus = {}\nwords_glo
ve = set(model.keys())\nfor i in words:\n    i
f i in words_glove:\n        words_courpus[i]
= model[i]\nprint("word 2 vec length", len(wor
ds_courpus))\n\n\n# stronging variables into p
ickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
\n\nimport pickle\nwith open(\'glove
_vectors\', \'wb\') as f:\n    pickle.dump(wor
ds_courpus, f)\n\n\n'

```

In [53]:

```

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

```

In [54]:

```

# average Word2Vec
# compute average word2vec for each review.
train_avg_w2v_essays = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(project_data_train['preprocessed_essays']
): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length

```

```

    cnt_words =0; # num of words with a valid vector in the s
entence/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
    train_avg_w2v_essays.append(vector)

print(len(train_avg_w2v_essays))
print(len(train_avg_w2v_essays[0]))

```

```

100%|██████████| 49041/49041 [00:08<00:00, 576
2.28it/s]

```

```
49041
```

```
300
```

In [55]:

```

# average Word2Vec
# compute average word2vec for each review.
cv_avg_w2v_essays = []; # the avg-w2v for each sentence/revie
w is stored in this list
for sentence in tqdm(project_data_cv['preprocessed_essays']):
    # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero leng
th
    cnt_words =0; # num of words with a valid vector in the s
entence/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
    cv_avg_w2v_essays.append(vector)

print(len(cv_avg_w2v_essays))
print(len(cv_avg_w2v_essays[0]))

```

```

100%|██████████| 24155/24155 [00:04<00:00, 572
3.17it/s]

```

```

24155
300

```

In [56]:

```

# average Word2Vec
# compute average word2vec for each review.
test_avg_w2v_essays = []; # the avg-w2v for each sentence/review
                           # is stored in this list
for sentence in tqdm(project_data_test['preprocessed_essays']):
    # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    test_avg_w2v_essays.append(vector)

print(len(test_avg_w2v_essays))
print(len(test_avg_w2v_essays[0]))

```

100%|██████████| 36052/36052 [00:06<00:00, 576
0.36it/s]

36052

300

In [57]:

```
# average Word2Vec
# compute average word2vec for each review.
train_avg_w2v_titles = []; # the avg-w2v for each sentence/re
view is stored in this list
for sentence in tqdm(project_data_train['preprocessed_titles'
]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero leng
th
    cnt_words = 0; # num of words with a valid vector in the s
entence/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
    train_avg_w2v_titles.append(vector)

print(len(train_avg_w2v_titles))
print(len(train_avg_w2v_titles[0]))
```

100%|██████████| 49041/49041 [00:00<00:00, 100
132.30it/s]

49041

300

In [58]:

```
# average Word2Vec
# compute average word2vec for each review.
cv_avg_w2v_titles = []; # the avg-w2v for each sentence/review
                        # is stored in this list
for sentence in tqdm(project_data_cv['preprocessed_titles']):
    # 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
    cv_avg_w2v_titles.append(vector)

print(len(cv_avg_w2v_titles))
print(len(cv_avg_w2v_titles[0]))
```

```
100%|██████████| 24155/24155 [00:00<00:00, 971
72.04it/s]
```

24155

300

In [59]:

```
# average Word2Vec
# compute average word2vec for each review.
test_avg_w2v_titles = []; # the avg-w2v for each sentence/review
                           # is stored in this list
for sentence in tqdm(project_data_test['preprocessed_titles']
```

```

): # 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
    test_avg_w2v_titles.append(vector)

print(len(test_avg_w2v_titles))
print(len(test_avg_w2v_titles[0]))

```

```

100%|██████████| 36052/36052 [00:00<00:00, 100
078.59it/s]

```

```

36052
300

```

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [60]:

```

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

```

```
tfidf_words = set(tfidf_model.get_feature_names())
```

In [61]:

```
# average Word2Vec
# compute average word2vec for each review.
train_tfidf_w2v_essays = []; # the avg-w2v for each sentence/
review is stored in this list
for sentence in tqdm(project_data_train['preprocessed_essays'
]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero leng
th
    tf_idf_weight = 0; # num of words with a valid vector in t
he 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 w
ord
            # here we are multiplying idf value(dictionary[wo
rd]) and the tf value((sentence.count(word)/len(sentence.spli
t()))))
            tf_idf = dictionary[word]*(sentence.count(word)/l
en(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weig
hted w2v
            tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
    train_tfidf_w2v_essays.append(vector)

print(len(train_tfidf_w2v_essays))
print(len(train_tfidf_w2v_essays[0]))
```

```
100%|██████████| 49041/49041 [01:03<00:00, 770
.37it/s]
```

49041

In [62]:

```

# average Word2Vec
# compute average word2vec for each review.
cv_tfidf_w2v_essays = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(project_data_cv['preprocessed_essays']):
    # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    cv_tfidf_w2v_essays.append(vector)

print(len(cv_tfidf_w2v_essays))
print(len(cv_tfidf_w2v_essays[0]))

```

```

100%|██████████| 24155/24155 [00:31<00:00, 775.02it/s]

```

24155

300

In [63]:

```
# average Word2Vec
# compute average word2vec for each review.
test_tfidf_w2v_essays = []; # the avg-w2v for each sentence/
review is stored in this list
for sentence in tqdm(project_data_test['preprocessed_essays']
): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero leng
th
    tf_idf_weight = 0; # num of words with a valid vector in t
he 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 w
ord
            # here we are multiplying idf value(dictionary[wo
rd]) and the tf value((sentence.count(word)/len(sentence.spli
t()))))
            tf_idf = dictionary[word]*(sentence.count(word)/l
en(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weig
hted w2v
            tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
    test_tfidf_w2v_essays.append(vector)

print(len(test_tfidf_w2v_essays))
print(len(test_tfidf_w2v_essays[0]))
```

100%|██████████| 36052/36052 [00:46<00:00, 773
.66it/s]

36052

300

In [64]:

```
# Similarly you can vectorize for title also
tfidf_model = TfidfVectorizer()
tfidf_model.fit(project_data_train['preprocessed_titles'])
# 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 [65]:

```
# average Word2Vec
# compute average word2vec for each review.
train_tfidf_w2v_titles = []; # the avg-w2v for each sentence/
review is stored in this list
for sentence in tqdm(project_data_train['preprocessed_titles'
]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero leng
th
    tf_idf_weight = 0; # num of words with a valid vector in t
he 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 w
ord
            # here we are multiplying idf value(dictionary[word]
rd]) and the tf value((sentence.count(word)/len(sentence.spli
t())))
            tf_idf = dictionary[word]*(sentence.count(word)/l
en(sentence.split())) # getting the tfidf value for each word
```

```

        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    train_tfidf_w2v_titles.append(vector)

print(len(train_tfidf_w2v_titles))
print(len(train_tfidf_w2v_titles[0]))

```

```

100%|██████████| 49041/49041 [00:01<00:00, 43757.86it/s]

```

```
49041
```

```
300
```

In [66]:

```

# average Word2Vec
# compute average word2vec for each review.
cv_tfidf_w2v_titles = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(project_data_cv['preprocessed_titles']):
    # 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()))

```

```

en(sentence.split())) # getting the tfidf value for each word
    vector += (vec * tf_idf) # calculating tfidf weighted w2v
    tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    cv_tfidf_w2v_titles.append(vector)

print(len(cv_tfidf_w2v_titles))
print(len(cv_tfidf_w2v_titles[0]))

```

```

100%|██████████| 24155/24155 [00:00<00:00, 436
24.30it/s]

```

```

24155
300

```

In [67]:

```

# average Word2Vec
# compute average word2vec for each review.
test_tfidf_w2v_titles = []; # the avg-w2v for each sentence/
review is stored in this list
for sentence in tqdm(project_data_test['preprocessed_titles']
): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the
sentence/review
    for word in sentence.split(): # for each word in a review
/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word])
and the tf value((sentence.count(word)/len(sentence.split())))

```

```

        tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
        vector += (vec * tf_idf) # calculating tfidf weighted w2v

        tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
        test_tfidf_w2v_titles.append(vector)

print(len(test_tfidf_w2v_titles))
print(len(test_tfidf_w2v_titles[0]))

```

```

100%|██████████| 36052/36052 [00:00<00:00, 45580.12it/s]

```

```

36052
300

```

1.5.3 Vectorizing Numerical features

In [68]:

```

price_data = resource_data.groupby('id').agg({'price': 'sum',
'quantity': 'sum'}).reset_index()

```

In [69]:

```

project_data_train = pd.merge(project_data_train, price_data,
on='id', how='left')
project_data_cv = pd.merge(project_data_cv, price_data, on='id', how='left')
project_data_test = pd.merge(project_data_test, price_data, on='id', how='left')

```

In [70]:

```

from sklearn.preprocessing import Normalizer
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array ins
tead:
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.

normalizer = Normalizer()
normalizer.fit(project_data_train['price'].values.reshape(-1,1
))

price_normalized_train = normalizer.transform(project_data_tr
ain['price'].values.reshape(-1, 1))
price_normalized_cv = normalizer.transform(project_data_cv['p
rice'].values.reshape(-1, 1))
price_normalized_test = normalizer.transform(project_data_tes
t['price'].values.reshape(-1, 1))

print('After normalization')
print(price_normalized_train.shape)
print(price_normalized_cv.shape)
print(price_normalized_test.shape)

```

After normalization

```

(49041, 1)
(24155, 1)
(36052, 1)

```

In [71]:

```

normalizer = Normalizer()
normalizer.fit(project_data_train['quantity'].values.reshape(-
1,1))

```

```
quantity_normalized_train = normalizer.transform(project_data_train['quantity'].values.reshape(-1, 1))
quantity_normalized_cv = normalizer.transform(project_data_cv['quantity'].values.reshape(-1, 1))
quantity_normalized_test = normalizer.transform(project_data_test['quantity'].values.reshape(-1, 1))

print('After normalization')
print(quantity_normalized_train.shape)
print(quantity_normalized_cv.shape)
print(quantity_normalized_test.shape)
```

After normalization

```
(49041, 1)
(24155, 1)
(36052, 1)
```

In [72]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

previously_posted_projects_normalized_train = normalizer.transform(project_data_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
previously_posted_projects_normalized_cv = normalizer.transform(project_data_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
previously_posted_projects_normalized_test = normalizer.transform(project_data_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

print('After normalization')
```



```
print(previously_posted_projects_normalized_train.shape)
print(previously_posted_projects_normalized_cv.shape)
print(previously_posted_projects_normalized_test.shape)
```

After normalization

```
(49041, 1)
(24155, 1)
(36052, 1)
```

In [73]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['essay_word_count'].values.
               reshape(-1,1))
```

```
essay_word_count_normalized_train = normalizer.transform(proj
ect_data_train['essay_word_count'].values.reshape(-1, 1))
essay_word_count_normalized_cv = normalizer.transform(project
_data_cv['essay_word_count'].values.reshape(-1, 1))
essay_word_count_normalized_test = normalizer.transform(proje
ct_data_test['essay_word_count'].values.reshape(-1, 1))
```

```
print('After normalization')
print(essay_word_count_normalized_train.shape)
print(essay_word_count_normalized_cv.shape)
print(essay_word_count_normalized_test.shape)
```

After normalization

```
(49041, 1)
(24155, 1)
(36052, 1)
```

In [74]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['title_word_count'].values.
```

```
reshape(-1,1))
```

```
title_word_count_normalized_train = normalizer.transform(project_data_train['title_word_count'].values.reshape(-1, 1))
title_word_count_normalized_cv = normalizer.transform(project_data_cv['title_word_count'].values.reshape(-1, 1))
title_word_count_normalized_test = normalizer.transform(project_data_test['title_word_count'].values.reshape(-1, 1))
```

```
print('After normalization')
print(title_word_count_normalized_train.shape)
print(title_word_count_normalized_cv.shape)
print(title_word_count_normalized_test.shape)
```

After normalization

(49041, 1)

(24155, 1)

(36052, 1)

In [75]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['neg'].values.reshape(-1,1))
)
```

```
sent_neg_train = normalizer.transform(project_data_train['neg'].values.reshape(-1, 1))
sent_neg_cv = normalizer.transform(project_data_cv['neg'].values.reshape(-1, 1))
sent_neg_test = normalizer.transform(project_data_test['neg'].values.reshape(-1, 1))
```

```
print('After normalization')
print(sent_neg_train.shape)
```

```
print(sent_neg_cv.shape)
print(sent_neg_test.shape)
```

After normalization

```
(49041, 1)
(24155, 1)
(36052, 1)
```

In [76]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['pos'].values.reshape(-1,1)
)

sent_pos_train = normalizer.transform(project_data_train['pos
'].values.reshape(-1, 1))
sent_pos_cv = normalizer.transform(project_data_cv['pos'].val
ues.reshape(-1, 1))
sent_pos_test = normalizer.transform(project_data_test['pos']
.values.reshape(-1, 1))

print('After normalization')
print(sent_pos_train.shape)
print(sent_pos_cv.shape)
print(sent_pos_test.shape)
```

After normalization

```
(49041, 1)
(24155, 1)
(36052, 1)
```

In [77]:

```
normalizer = Normalizer()
normalizer.fit(project_data_train['neu'].values.reshape(-1,1)
)
```

```

sent_neu_train = normalizer.transform(project_data_train['neu
'].values.reshape(-1, 1))
sent_neu_cv = normalizer.transform(project_data_cv['neu'].val
ues.reshape(-1, 1))
sent_neu_test = normalizer.transform(project_data_test['neu']
.values.reshape(-1, 1))

print('After normalization')
print(sent_neu_train.shape)
print(sent_neu_cv.shape)
print(sent_neu_test.shape)

```

After normalization

```

(49041, 1)
(24155, 1)
(36052, 1)

```

In [78]:

```

normalizer = Normalizer()
normalizer.fit(project_data_train['compound'].values.reshape(-
1,1))

sent_compound_train = normalizer.transform(project_data_train
['compound'].values.reshape(-1, 1))
sent_compound_cv = normalizer.transform(project_data_cv['comp
ound'].values.reshape(-1, 1))
sent_compound_test = normalizer.transform(project_data_test['
compound'].values.reshape(-1, 1))

print('After normalization')
print(sent_compound_train.shape)
print(sent_compound_cv.shape)

```

```
print(sent_compound_test.shape)
```

After normalization

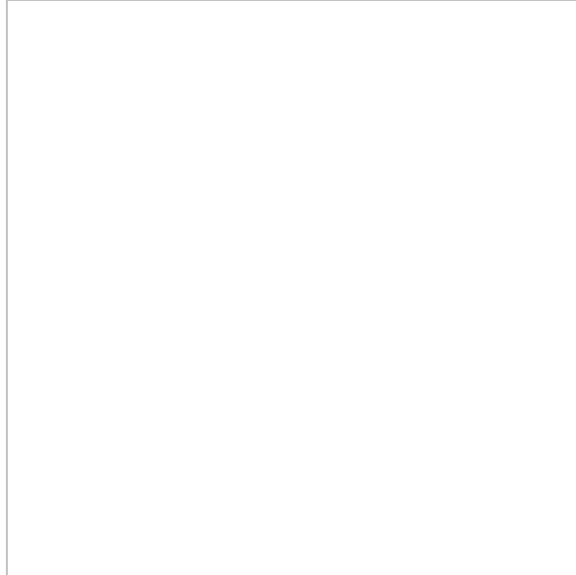
(49041, 1)

(24155, 1)

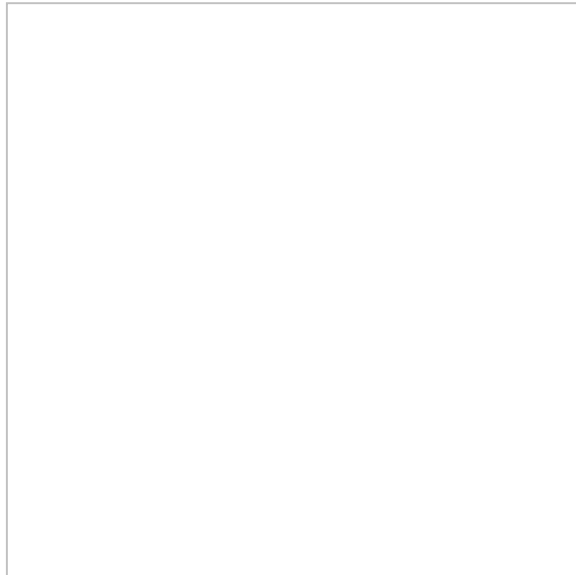
(36052, 1)

Assignment 5: Logistic Regression

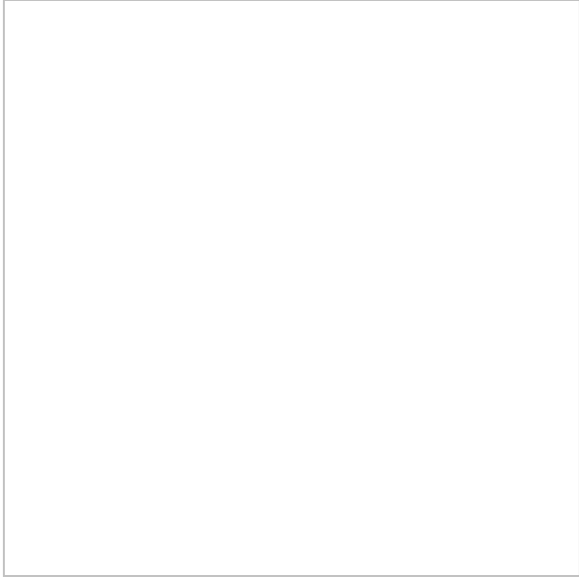
1. **[Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets**
 - **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay (`BOW with bi-grams` with `min_df=10` and `max_features=5000`)
 - **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (`TFIDF with bi-grams` with `min_df=10` and `max_features=5000`)
 - **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
 - **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)
2. **Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)**
 - Find the best hyper parameter which will give the maximum [AUC](#) value
 - Find the best hyper paramter using k-fold cross validation or simple cross validation data
 - Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning
3. **Representation of results**
 - You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.



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



- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

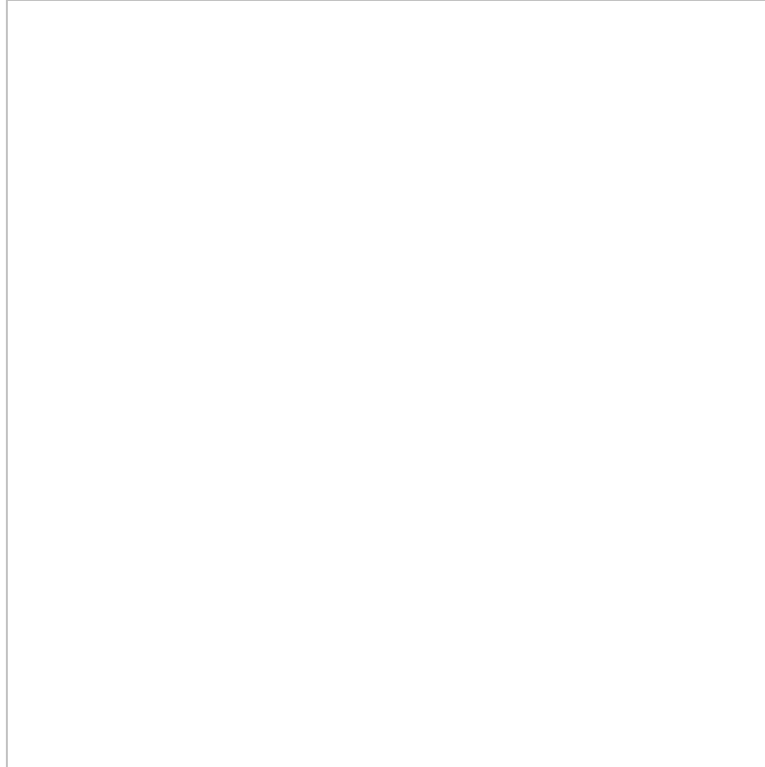
- 
4. [Task-2] Apply Logistic Regression on the below feature set **Set 5** by finding the best hyper parameter as suggested in step 2 and step 3.
 5. Consider these set of features **Set 5** :

- **school state** : categorical data
- **clean categories** : categorical data
- **clean subcategories** : categorical data
- **project grade category** :categorical data
- **teacher prefix** : categorical data
- **quantity** : numerical data
- **teacher number of previously posted projects** : numerical data
- **price** : numerical data
- **sentiment score's of each of the essay** : numerical data
- **number of words in the title** : numerical data
- **number of words in the combine essays** : numerical data

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

6. 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](#).

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

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

Logistic Regression

Set 1: Categorical, Numerical features + Project_title(BOW) + Preprocessed_essay (BOW with bi-grams with min_df=10 and max_features=5000)

In [79]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
```

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

```
X_train = hstack((train_categories_one_hot, train_subcategories_one_hot, train_school_state_category_one_hot , train_project_grade_category_one_hot, train_teacher_prefix_categories_one_hot, price_normalized_train, quantity_normalized_train, previously_posted_projects_normalized_train, title_word_count_normalized_train, essay_word_count_normalized_train, train_title_bow, train_essay_bow)).tocsr()
```

```
X_test = hstack((test_categories_one_hot, test_subcategories_one_hot, test_school_state_category_one_hot , test_project_grade_category_one_hot, test_teacher_prefix_categories_one_hot, price_normalized_test, quantity_normalized_test, previously_posted_projects_normalized_test, title_word_count_normalized_test, essay_word_count_normalized_test, test_title_bow, test_essay_bow)).tocsr()
```

```
X_cv = hstack((cv_categories_one_hot, cv_subcategories_one_hot, cv_school_state_category_one_hot , cv_project_grade_category_one_hot, cv_teacher_prefix_categories_one_hot, price_normalized_cv, quantity_normalized_cv, previously_posted_projects_normalized_cv, title_word_count_normalized_cv, essay_word_count_normalized_cv, cv_title_bow, cv_essay_bow)).tocsr()
```

In [80]:

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

```
(49041, 6767)
```

```
(24155, 6767)
```

```
(36052, 6767)
```

Using GridSearchCV (K fold Cross Validation) to determine the best hyperparameter

In [81]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

parameters = {'C':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10*
*1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(lr, parameters, cv= 10, scoring='roc_auc')

clf.fit(X_train, y_train)

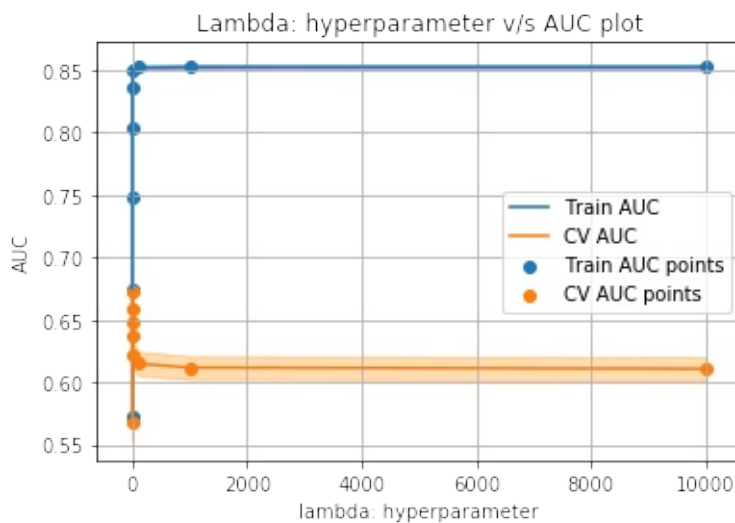
train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/
48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_
std,train_auc + train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/
48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv
_auc + cv_auc_std,alpha=0.3,color='darkorange')
```

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

plt.legend()
plt.xlabel("lambda: hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC plot")
plt.grid()
plt.show()
```



In [82]:

```
#https://datascience.stackexchange.com/questions/21877/how-to-use-the-output-of-gridsearch
#choosing the best hyperparameter
clf.best_params_
```

Out[82]:

```
{'C': 0.01}
```

Train the model using the best hyperparameter value

In [84]:

```
best_c1 = clf.best_params_['C']
```

In [83]:

```
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 mul  
    tiplier  
    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 [86]:

```
#Train the model using the best hyperparameter found from GridSearch/RandomSearch/SimpleCV
```



```
from sklearn.metrics import roc_curve, auc

model = LogisticRegression(C = best_c1)

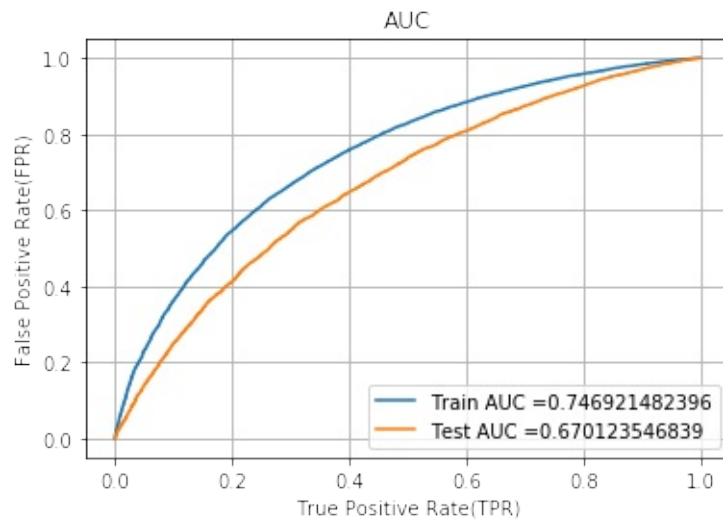
model.fit(X_train, y_train)

# 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_train)
y_test_pred = batch_predict(model, X_test)

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("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Confusion Matrix of Train and Test Data

In [87]:

```
# we are writing our own function for predict, with defined t  
hreshold  
# we will pick a threshold that will give the least fpr  
def find_best_threshold(threshold, fpr, tpr):  
    t = threshold[np.argmax(tpr*(1-fpr))]  
    # (tpr*(1-fpr)) will be maximum if your fpr is very low a  
nd tpr is very high  
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)  
) , "for threshold", np.round(t,3))  
    return t  
  
def predict_with_best_t(proba, threshold):  
    predictions = []  
    for i in proba:  
        if i>=threshold:  
            predictions.append(1)  
        else:  
            predictions.append(0)  
    return predictions
```

In [88]:

```
#our objective here is to make auc the maximum  
#so we find the best threshold that will give the least fpr  
best_t = find_best_threshold(tr_thresholds, train_fpr, train_  
tpr)  
print("Train confusion matrix")  
print(confusion_matrix(y_train, predict_with_best_t(y_train_p  
red, best_t)))
```

the maximum value of tpr*(1-fpr) 0.46781475337

```
for threshold 0.832
Train confusion matrix
[[ 4966  2460]
 [12503 29112]]
```

In [89]:

```
#plotting confusion matrix using seaborn's heatmap
# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

print("Train data confusion matrix")

confusion_matrix_df_train = pd.DataFrame(confusion_matrix(y_train,
predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_train, annot=True, annot_kws={
"size": 16}, fmt='g')
```

Train data confusion matrix

Out[89]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05e1c54a58>
```



In [90]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

Test confusion matrix

```
[[ 3076  2383]
 [ 9807 20786]]
```

In [91]:

```
print("Test data confusion matrix")

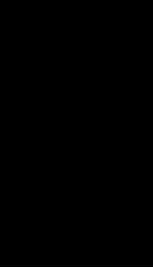
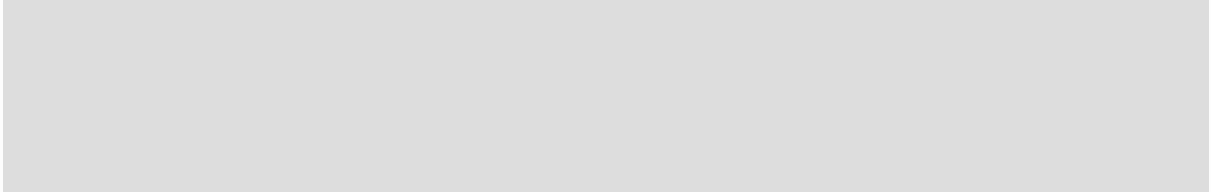
confusion_matrix_df_test = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_test, annot=True, annot_kws={"size": 16}, fmt='g')
```

Test data confusion matrix

Out[91]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f05e0a0e198>





Set 2 : Categorical, Numerical features + Project_title(TFIDF) + Preprocessed_essay (TFIDF with bi-grams with min_df=10 and max_features=5000)

In [92]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
```

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

```
X_train = hstack((train_categories_one_hot, train_subcategories_one_hot, train_school_state_category_one_hot , train_project_grade_category_one_hot, train_teacher_prefix_categories_one_hot, price_normalized_train, quantity_normalized_train, previously_posted_projects_normalized_train, title_word_count_normalized_train, essay_word_count_normalized_train, train_title_tfidf, train_essay_tfidf)).tocsr()
```

```
X_test = hstack((test_categories_one_hot, test_subcategories_one_hot, test_school_state_category_one_hot , test_project_grade_category_one_hot, test_teacher_prefix_categories_one_hot, price_normalized_test, quantity_normalized_test, previously_posted_projects_normalized_test, title_word_count_normalized_test, essay_word_count_normalized_test, test_title_tfidf, test_essay_tfidf)).tocsr()
```

```
X_cv = hstack((cv_categories_one_hot, cv_subcategories_one_hot, cv_school_state_category_one_hot , cv_project_grade_category_one_hot, cv_teacher_prefix_categories_one_hot, price_normalized_cv, quantity_normalized_cv, previously_posted_projects_normalized_cv, title_word_count_normalized_cv, essay_word_count_normalized_cv, cv_title_tfidf, cv_essay_tfidf)).tocsr()
```

In [93]:

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

```
(49041, 6767)
```

```
(24155, 6767)
```

```
(36052, 6767)
```

Using GridSearchCV (K fold Cross Validation) to determine the best hyperparameter

In [94]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

parameters = {'C':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10*
*1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(lr, parameters, cv= 10, scoring='roc_auc')

clf.fit(X_train, y_train)

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

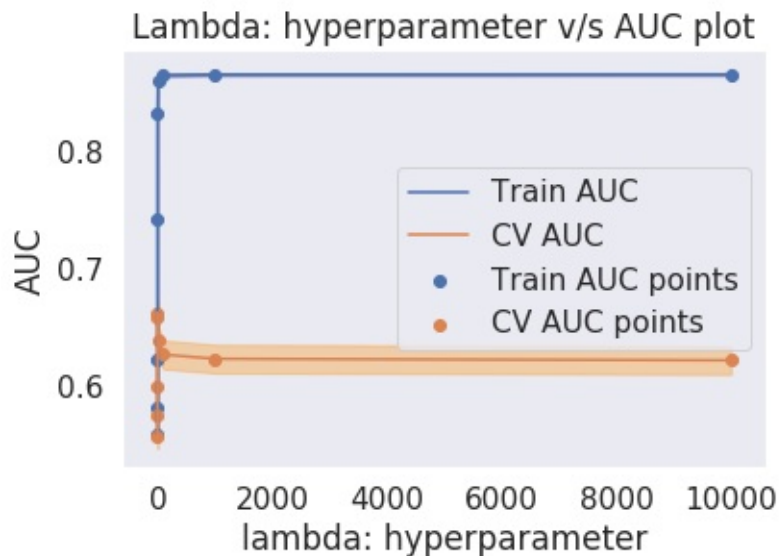
plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_
std,train_auc + train_auc_std,alpha=0.3,color='darkblue')
```



```
plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')

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

plt.legend()
plt.xlabel("lambda: hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC plot")
plt.grid()
plt.show()
```



In [95]:

```
#https://datascience.stackexchange.com/questions/21877/how-to-use-the-output-of-gridsearch
#choosing the best hyperparameter
```

```
clf.best_params_
```

Out[95]:

```
{'C': 1}
```

Train the model using the best hyperparameter value

In [96]:

```
best_c2 = clf.best_params_['C']
```

In [97]:

```
#Train the model using the best hyperparameter found from GridSearch/RandomSearch/SimpleCV
```

```
from sklearn.metrics import roc_curve, auc
```

```
model = LogisticRegression(C = best_c2)
```

```
model.fit(X_train, y_train)
```

```
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
```

```
# not the predicted outputs
```

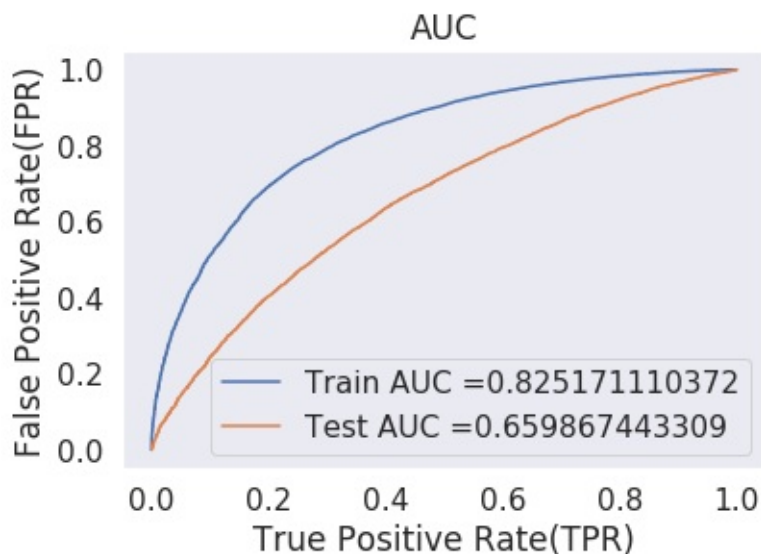
```
y_train_pred = batch_predict(model, X_train)
```

```
y_test_pred = batch_predict(model, X_test)
```

```
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("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Confusion Matrix of Train and Test Data

In [98]:

```
#our objective here is to make auc the maximum
#so we find the best threshold that will give the least fpr
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_p
```

```
red, best_t)))
```

the maximum value of $\text{tpr} \times (1 - \text{fpr})$ 0.56260740964

9 for threshold 0.828

Train confusion matrix

```
[[ 5586  1840]
 [10490 31125]]
```

In [99]:

```
#plotting confusion matrix using seaborn's heatmap
# https://stackoverflow.com/questions/35572000/how-can-i-plot
-a-confusion-matrix
```

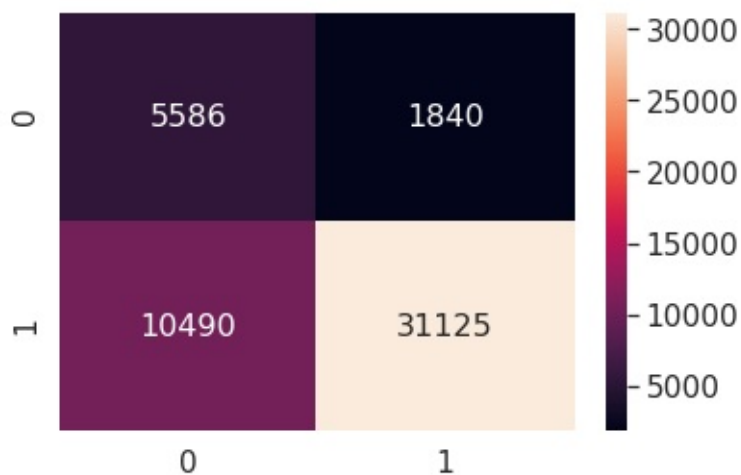
```
print("Train data confusion matrix")
```

```
confusion_matrix_df_train = pd.DataFrame(confusion_matrix(y_train,
predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_train, annot=True, annot_kws={
"size": 16}, fmt='g')
```

Train data confusion matrix

Out[99]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05cffddfd0>
```



In [100]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

Test confusion matrix

```
[[ 2812  2647]
 [ 8979 21614]]
```

In [101]:

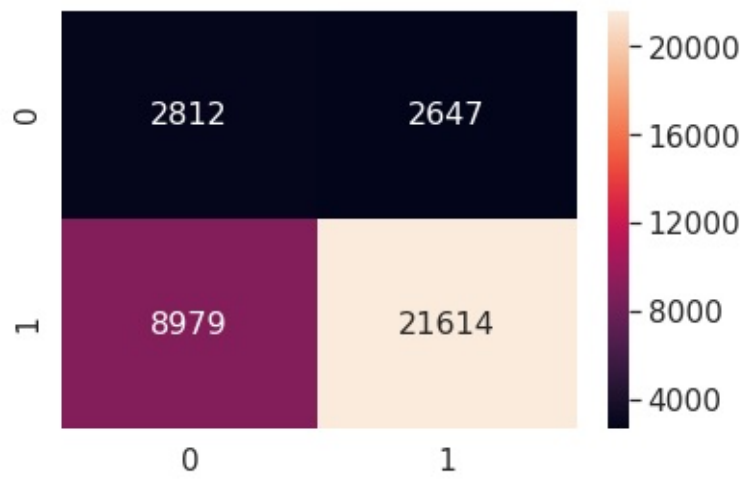
```
print("Test data confusion matrix")

confusion_matrix_df_test = pd.DataFrame(confusion_matrix(y_test,
predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_test, annot=True, annot_kws={"
size": 16}, fmt='g')
```

Test data confusion matrix

Out[101]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05cfef1d30>
```



Set 3 : Categorical, Numerical features + Project_title(AVG W2V) + Preprocessed_essay (AVG W2V)

In [102]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train = hstack((train_categories_one_hot, train_subcategories_one_hot, train_school_state_category_one_hot , train_project_grade_category_one_hot, train_teacher_prefix_categories_one_hot, price_normalized_train, quantity_normalized_train, previously_posted_projects_normalized_train, title_word_count_normalized_train, essay_word_count_normalized_train, train_avg_w2v_essays, train_avg_w2v_titles)).tocsr()
X_test = hstack((test_categories_one_hot, test_subcategories_one_hot, test_school_state_category_one_hot , test_project_grade_category_one_hot, test_teacher_prefix_categories_one_hot, price_normalized_test, quantity_normalized_test, previously_posted_projects_normalized_test, title_word_count_normalized_test, essay_word_count_normalized_test, test_avg_w2v_essays, test_avg_w2v_titles)).tocsr()
X_cv = hstack((cv_categories_one_hot, cv_subcategories_one_hot, cv_school_state_category_one_hot , cv_project_grade_category_one_hot, cv_teacher_prefix_categories_one_hot, price_normalized_cv, quantity_normalized_cv, previously_posted_projects_normalized_cv, title_word_count_normalized_cv, essay_word_count_normalized_cv, cv_avg_w2v_essays, cv_avg_w2v_titles)).tocsr()
```

In [103]:

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

```
(49041, 704)
(24155, 704)
(36052, 704)
```

Using GridSearchCV (K fold Cross Validation) to determine the best hyperparameter

In [104]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

parameters = {'C':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10*
*1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(lr, parameters, cv= 10, scoring='roc_auc')

clf.fit(X_train, y_train)

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

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



```

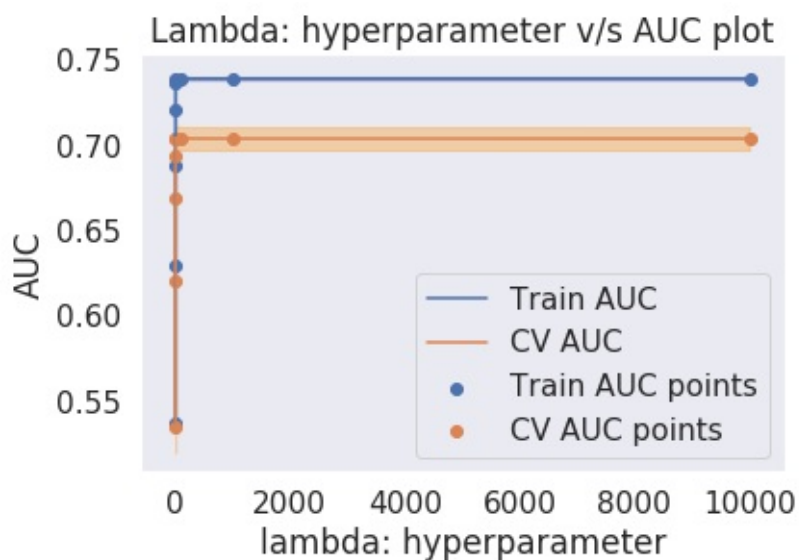
std,train_auc + train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("lambda: hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC plot")
plt.grid()
plt.show()

```



Train the model using the best hyper parameter value

In [105]:

```
#https://datascience.stackexchange.com/questions/21877/how-to  
-use-the-output-of-gridsearch  
#choosing the best hyperparameter  
clf.best_params_
```

Out[105]:

```
{'C': 1}
```

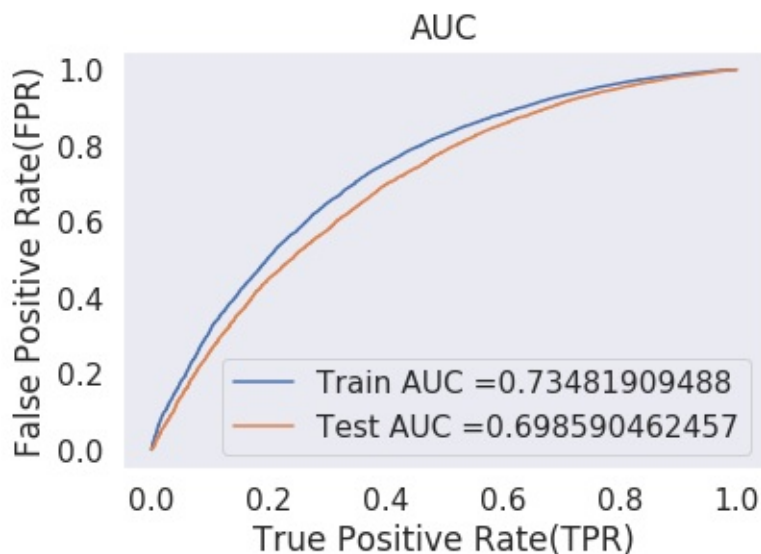
In [107]:

```
best_c3 = clf.best_params_['C']
```

In [108]:

```
#Train the model using the best hyperparameter found from Gri  
dSearch/RandomSearch/SimpleCV  
  
from sklearn.metrics import roc_curve, auc  
  
model = LogisticRegression(C = best_c3)  
  
model.fit(X_train, 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_train)  
y_test_pred = batch_predict(model, X_test)  
  
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("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Confusion Matrix of Train and Test Data

In [109]:

```
#our objective here is to make auc the maximum
#so we find the best threshold that will give the least fpr
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_p
```

```
red, best_t)))
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.45894408878

5 for threshold 0.834

Train confusion matrix

```
[[ 4832  2594]
 [12263 29352]]
```

In [110]:

```
#plotting confusion matrix using seaborn's heatmap
# https://stackoverflow.com/questions/35572000/how-can-i-plot
-a-confusion-matrix
```

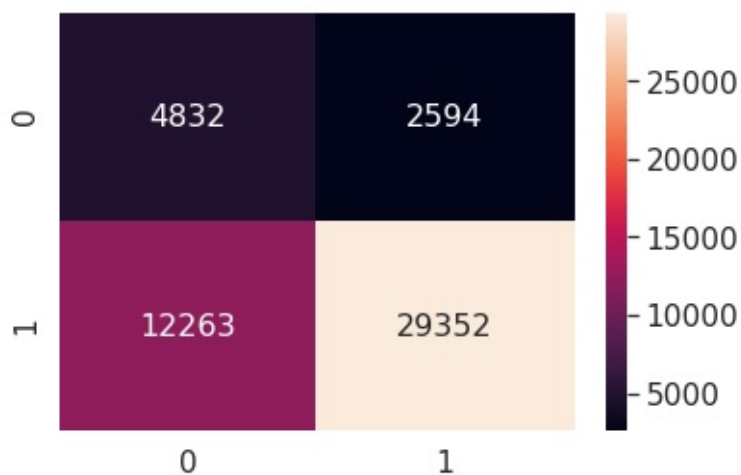
```
print("Train data confusion matrix")
```

```
confusion_matrix_df_train = pd.DataFrame(confusion_matrix(y_train,
predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_train, annot=True, annot_kws={
"size": 16}, fmt='g')
```

Train data confusion matrix

Out[110]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05d54f9630>
```



In [111]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

Test confusion matrix

```
[[ 3266  2193]
 [ 9201 21392]]
```

In [112]:

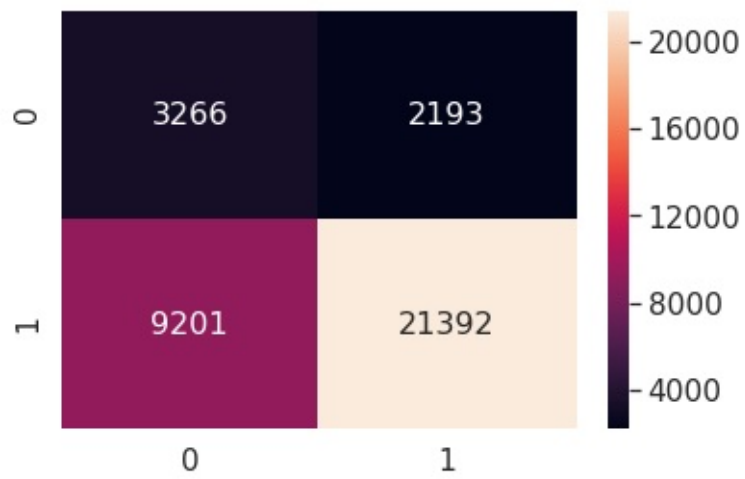
```
print("Test data confusion matrix")

confusion_matrix_df_test = pd.DataFrame(confusion_matrix(y_test,
predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_test, annot=True, annot_kws={"
size": 16}, fmt='g')
```

Test data confusion matrix

Out[112]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05e0195860>
```



Set 4 : Categorical, Numerical features + Project_title(TFIDF W2V) + Preprocessed_essay (TFIDF W2V)

In [113]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
```

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

```
X_train = hstack((train_categories_one_hot, train_subcategories_one_hot, train_school_state_category_one_hot , train_project_grade_category_one_hot, train_teacher_prefix_categories_one_hot, price_normalized_train, quantity_normalized_train, previously_posted_projects_normalized_train, title_word_count_normalized_train, essay_word_count_normalized_train, train_tfidf_w2v_essays, train_tfidf_w2v_titles)).tocsr()
```

```
X_test = hstack((test_categories_one_hot, test_subcategories_one_hot, test_school_state_category_one_hot , test_project_grade_category_one_hot, test_teacher_prefix_categories_one_hot, price_normalized_test, quantity_normalized_test, previously_posted_projects_normalized_test, title_word_count_normalized_test, essay_word_count_normalized_test, test_tfidf_w2v_essays , test_tfidf_w2v_titles)).tocsr()
```

```
X_cv = hstack((cv_categories_one_hot, cv_subcategories_one_hot, cv_school_state_category_one_hot , cv_project_grade_category_one_hot, cv_teacher_prefix_categories_one_hot, price_normalized_cv, quantity_normalized_cv, previously_posted_projects_normalized_cv, title_word_count_normalized_cv, essay_word_count_normalized_cv, cv_tfidf_w2v_essays, cv_tfidf_w2v_titles)).tocsr()
```

In [114]:

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

```
(49041, 704)
(24155, 704)
(36052, 704)
```

Using GridSearchCV (K fold Cross Validation) to determine the best hyperparameter

In [115]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

parameters = {'C':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10*
*1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(lr, parameters, cv= 10, scoring='roc_auc')

clf.fit(X_train, y_train)

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

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



```

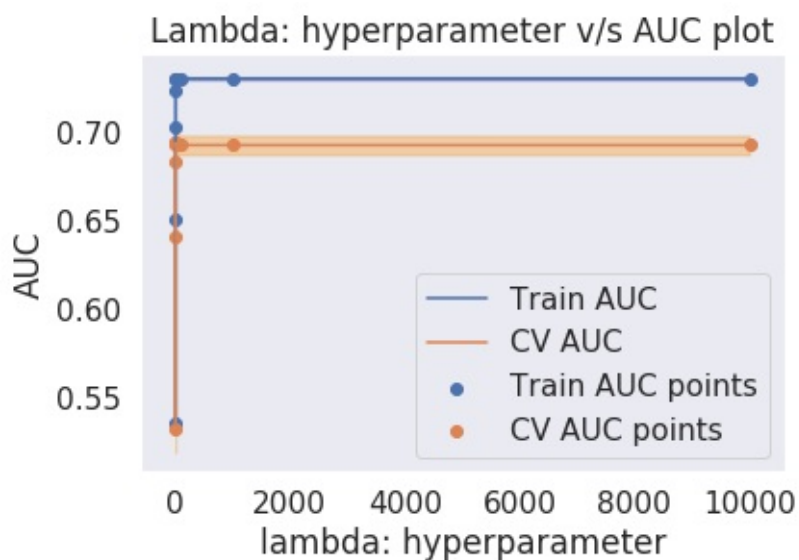
std,train_auc + train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/
48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv
_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("lambda: hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC plot")
plt.grid()
plt.show()

```



In [116]:

```

#https://datascience.stackexchange.com/questions/21877/how-to
-use-the-output-of-gridsearch

```

```
#choosing the best hyperparameter  
clf.best_params_
```

Out[116]:

```
{'C': 0.1}
```

Train the model using the best hyperparameter value

In [117]:

```
best_c4 = clf.best_params_['C']
```

In [118]:

```
#Train the model using the best hyperparameter found from GridSearch/RandomSearch/SimpleCV
```

```
from sklearn.metrics import roc_curve, auc
```

```
model = LogisticRegression(C = best_c4)
```

```
model.fit(X_train, 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_train)
```

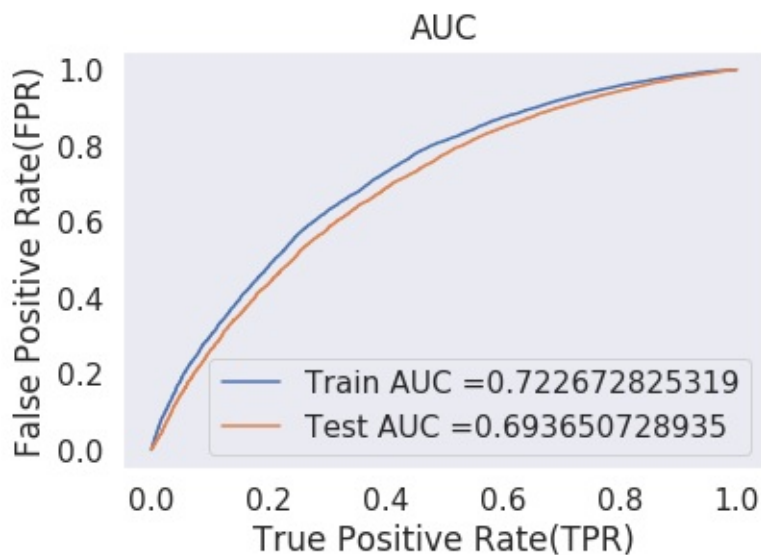
```
y_test_pred = batch_predict(model, X_test)
```

```
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)
```

```
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("True Positive Rate(TPR)")  
plt.ylabel("False Positive Rate(FPR)")  
plt.title("AUC")  
plt.grid()  
plt.show()
```



Confusion Matrix of Train and Test Data

In [119]:

```
#our objective here is to make auc the maximum  
#so we find the best threshold that will give the least fpr  
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)  
print("Train confusion matrix")
```

```
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
```

the maximum value of $\text{tpr} \times (1 - \text{fpr})$ 0.44151784571

6 for threshold 0.839

Train confusion matrix

```
[[ 4916  2510]
 [13860 27755]]
```

In [120]:

```
#plotting confusion matrix using seaborn's heatmap
# https://stackoverflow.com/questions/35572000/how-can-i-plot
-a-confusion-matrix
```

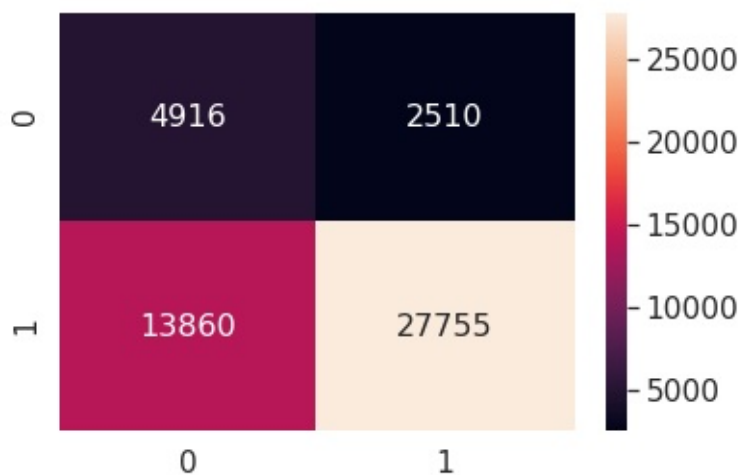
```
print("Train data confusion matrix")
```

```
confusion_matrix_df_train = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_train, annot=True, annot_kws={"size": 16}, fmt='g')
```

Train data confusion matrix

Out[120]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05d2450320>
```



In [121]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

Test confusion matrix
[[3422 2037]
 [10436 20157]]

In [122]:

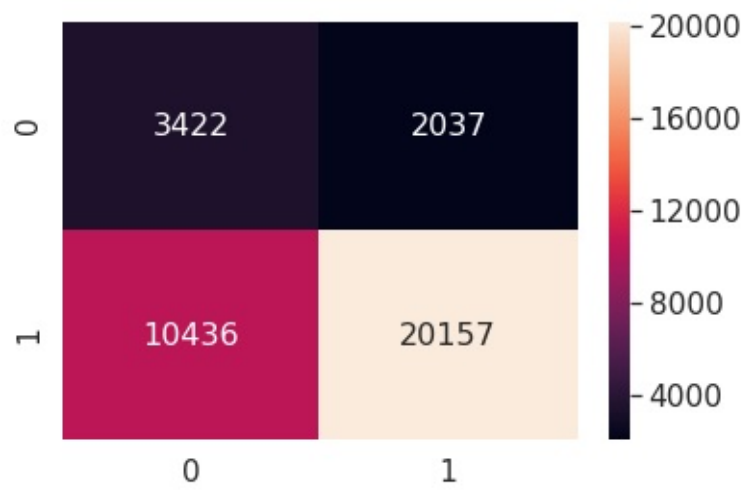
```
print("Test data confusion matrix")

confusion_matrix_df_test = pd.DataFrame(confusion_matrix(y_test,
predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_test, annot=True, annot_kws={"
size": 16}, fmt='g')
```

Test data confusion matrix

Out[122]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f05d6f64a90>



Logistic Regression with added Features

Set5

Set 5 : Categorical features + Numerical features + Essay Sentiment Scores

In [123]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train = hstack((train_categories_one_hot, train_subcategories_one_hot, train_school_state_category_one_hot , train_project_grade_category_one_hot, train_teacher_prefix_categories_one_hot, price_normalized_train, quantity_normalized_train, previously_posted_projects_normalized_train, title_word_count_normalized_train, essay_word_count_normalized_train, sent_neg_train, sent_pos_train, sent_neu_train, sent_compound_train)).tocsr()

X_test = hstack((test_categories_one_hot, test_subcategories_one_hot, test_school_state_category_one_hot , test_project_grade_category_one_hot, test_teacher_prefix_categories_one_hot, price_normalized_test, quantity_normalized_test, previously_posted_projects_normalized_test, title_word_count_normalized_test, essay_word_count_normalized_test, sent_neg_test, sent_pos_test, sent_neu_test, sent_compound_test)).tocsr()

X_cv = hstack((cv_categories_one_hot, cv_subcategories_one_hot, cv_school_state_category_one_hot , cv_project_grade_category_one_hot, cv_teacher_prefix_categories_one_hot, price_normalized_cv, quantity_normalized_cv, previously_posted_projects_normalized_cv, title_word_count_normalized_cv, essay_word_count_normalized_cv, sent_neg_cv, sent_pos_cv, sent_neu_cv, sent_compound_cv)).tocsr()
```

In [124]:


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

```
(49041, 108)
(24155, 108)
(36052, 108)
```

Using GridSearchCV (K fold Cross Validation) to determine the best hyperparameter

In [125]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()

parameters = {'C':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10*
*1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(lr, parameters, cv= 10, scoring='roc_auc')

clf.fit(X_train, y_train)

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

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

```

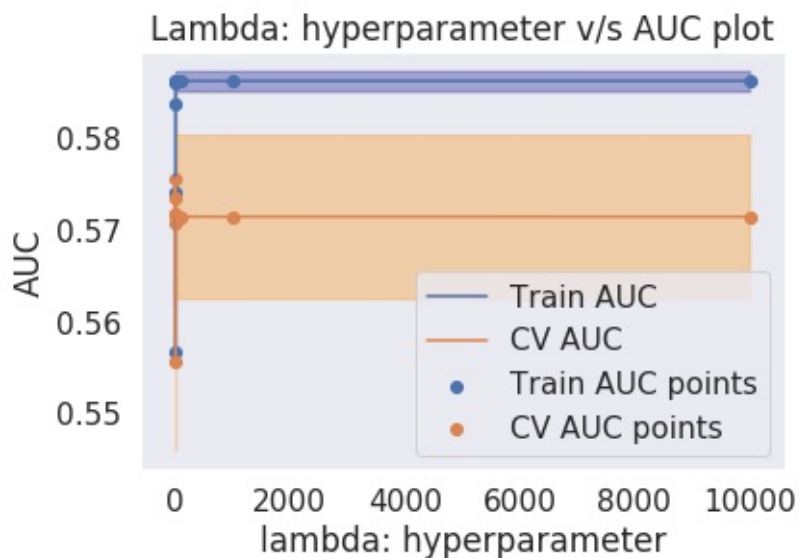
std,train_auc + train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("lambda: hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC plot")
plt.grid()
plt.show()

```



Train the model using the best hyper parameter value

In [126]:

```
#https://datascience.stackexchange.com/questions/21877/how-to  
-use-the-output-of-gridsearch  
#choosing the best hyperparameter  
clf.best_params_
```

Out[126]:

```
{'C': 0.01}
```

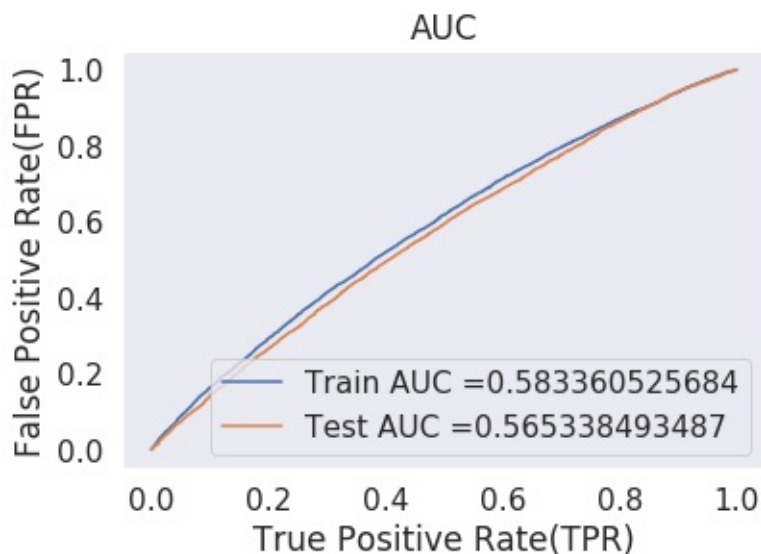
In [127]:

```
best_c5 = clf.best_params_['C']
```

In [129]:

```
#Train the model using the best hyperparameter found from Gri  
dSearch/RandomSearch/SimpleCV  
  
from sklearn.metrics import roc_curve, auc  
  
model = LogisticRegression(C = best_c5)  
  
model.fit(X_train, 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_train)  
y_test_pred = batch_predict(model, X_test)  
  
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("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()
```



Confusion Matrix of Train and Test Data

In [130]:

```
#our objective here is to make auc the maximum
#so we find the best threshold that will give the least fpr
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_p
```

```
red, best_t)))
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.31380898201

2 for threshold 0.848

Train confusion matrix

```
[[ 4188  3238]
 [18459 23156]]
```

In [131]:

```
#plotting confusion matrix using seaborn's heatmap
# https://stackoverflow.com/questions/35572000/how-can-i-plot
-a-confusion-matrix
```

```
print("Train data confusion matrix")
```

```
confusion_matrix_df_train = pd.DataFrame(confusion_matrix(y_train,
predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
```

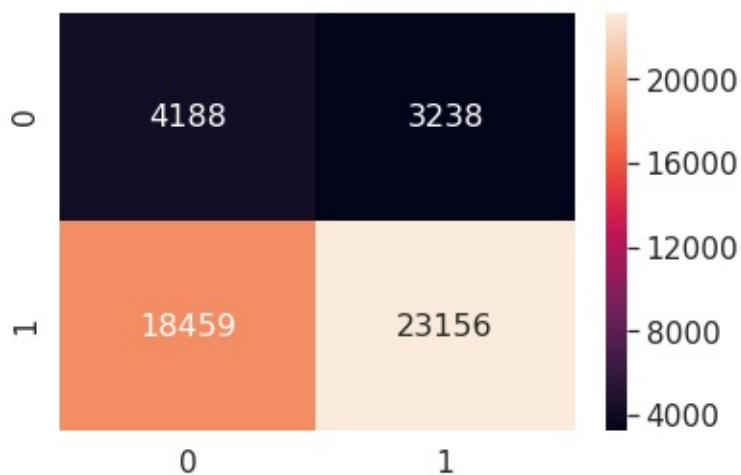
```
sns.set(font_scale=1.4)#for label size
```

```
sns.heatmap(confusion_matrix_df_train, annot=True, annot_kws={
"size": 16}, fmt='g')
```

Train data confusion matrix

Out[131]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05e15b3d68>
```



In [132]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

Test confusion matrix

```
[[ 2913  2546]
 [13503 17090]]
```

In [133]:

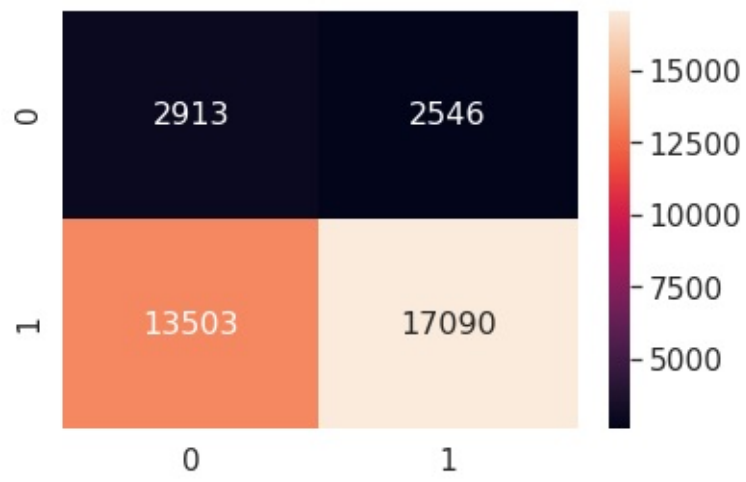
```
print("Test data confusion matrix")

confusion_matrix_df_test = pd.DataFrame(confusion_matrix(y_test,
predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
sns.set(font_scale=1.4)#for label size
sns.heatmap(confusion_matrix_df_test, annot=True, annot_kws={"
size": 16}, fmt='g')
```

Test data confusion matrix

Out[133]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f05d7e02898>
```



Conclusion

In [134]:

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

from prettytable import PrettyTable

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

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

x.add_row(["BOW", "Logistic Regression", 0.01, 0.67])
x.add_row(["TFIDF", "Logistic Regression", 1, 0.66])
x.add_row(["AVG W2V", "Logistic Regression", 1.0, 0.69])
x.add_row(["TFIDF W2V", "Logistic Regression", 0.1, 0.69])
x.add_row(["WITHOUT TEXT", "Logistic Regression", 0.01, 0.56]
)

print(x)
```

```
+-----+-----+-----+
-----+-----+
| Vectorizer | Model | Alpha:H
yper Parameter | AUC |
+-----+-----+-----+
-----+-----+
| BOW | Logistic Regression |
0.01 | 0.67 |
```


	TFIDF		Logistic Regression	
	1		0.66	
	AVG W2V		Logistic Regression	
	1.0		0.69	
	TFIDF W2V		Logistic Regression	
	0.1		0.69	
	WITHOUT TEXT		Logistic Regression	
	0.01		0.56	
+-----+-----+-----				
-----+-----+				