

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>	<ul style="list-style-type: none">••	Title of the project. Examples: <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care & Hunger</code> <code>Health & Sports</code> <code>History & Civics</code> <code>Literacy & Language</code> <code>Math & Science</code> <code>Music & The Arts</code> <code>Special Needs</code> <code>Warmth</code> Examples: <ul style="list-style-type: none">• <code>Music & The Arts</code>• <code>Literacy & Language, Math & Science</code>
<code>school_state</code>		State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: <code>Literacy</code> <code>Literature & Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

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

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

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

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

Notes on the Essay Data

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

- __project_essay_1__: "Introduce us to your classroom"
- __project_essay_2__: "Tell us more about your students"
- __project_essay_3__: "Describe how your students will use the materials you're requesting"
- __project_essay_3__: "Close by sharing why your project will make a difference"

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

- __project_essay_1__: "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2__: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with `project_submitted_datetime` of 2016-05-17 and later, the values of `project_essay_3` and `project_essay_4` will be NaN.

In [1]:

```
%matplotlib inline

import 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 TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

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

```

```

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

```

In [2]:

```

import warnings
warnings.filterwarnings("ignore")

```

1.1 Reading Data

In [0]:

```

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

```

In [0]:

```

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 [0]:

```

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 [0]:

categories:

	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 [0]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [0]:

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

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

```

temp = temp.replace('&', '_')
sub_cat_list.append(temp.strip())

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

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

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

```

1.3 Text preprocessing

In [0]:

```

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

```

In [0]:

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

Out[0]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade

In [0]:

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

In [0]:

```

# 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 o

f 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 alongside 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 families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills.\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\nnnnnnn

=====

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

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\n\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs 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!nnnnn

=====

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\n\r\nThey also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nnnnn

=====

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As

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

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.

In [0]:

```
# 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 [0]:

```
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.

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.

In [0]:

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

es can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

In [0]:

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

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

In [0]:

```
# 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', \
            'himselves', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', \
            'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", \
            'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', \
            'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', \
            'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', \
            'before', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', \
            'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', \
            'few', 'more', \
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', \
            'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', \
            "mightn't", 'mustn', \
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', \
            "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

In [0]:

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


In [0]:

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

Out[0]:

'my kindergarten students varied disabilities ranging speech language delays cognitive delays gross fine motor delays autism they eager beavers 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 limitations students love coming school come eager learn explore have ever felt like ants pants needed groove move meeting this kids feel time the want able move learn say wobble chairs answer i love develop core enhances gross motor turn fine motor skills they also want learn games kids not want sit worksheets they want learn count jumping playing physical engagement key success the number toss color shape mats make happen my students forget work fun 6 year old de serves nannan'

1.4 Preprocessing of `project_title`

In [0]:

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

1.5 Preparing data for models

In [0]:

```
project_data.columns
```

Out[0]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category', 'project_title',
      'project_essay_1', 'project_essay_2', 'project_essay_3',
      'project_essay_4', 'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay'],
      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 (optional)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

1.5.1 Vectorizing Categorical data

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

In [0]:

```
# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (109248, 9)
```

In [0]:

```
# we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL',
'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encodig (109248, 30)
```

In [0]:

```
# you can do the similar thing with state, teacher_prefix and project_grade_category also
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [0]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

```
Shape of matrix after one hot encodig (109248, 16623)
```

In [0]:

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

1.5.2.2 TFIDF vectorizer

In [0]:

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

```
Shape of matrix after one hot encodig (109248, 16623)
```

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

```
'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

words = []
for i in preprocod_texts:
    words.extend(i.split(' '))

for i in preprocod_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(" , np.round(len(inter_words)/len(words)*100,3), "%) ")

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

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

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

'''
```

Out[0]:

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

```
with open('glove_vectors', 'wb') as f: pickle.dump(words_corpus, f)
```

In [0]:

```
# 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 [0]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors.append(vector)

print(len(avg_w2v_vectors))
print(len(avg_w2v_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248
[00:27<00:00, 3968.04it/s]
```

```
109248
300
```

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

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

In [0]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

print(len(tfidf_w2v_vectors))
print(len(tfidf_w2v_vectors[0]))
```

109248
300

```
# Similarly you can vectorize for title also
```

In [0]:

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

Out[0]:

(109248, 16663)

Computing Sentiment Scores

In [239]:

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

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

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students with the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelligences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of different backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a caring community of successful \
learners which can be seen through collaborative student project based learning in and out of the classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role play in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food i will take their idea \
and create common core cooking lessons where we learn important math and writing 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 apple sauce 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)

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

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

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

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')

```
and max_features=3000 )
```

- **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](#).

2. Logistic Regression

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

In [3]:

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

In [4]:

```
# Loading Data Sets
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

In [5]:

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

In [6]:

```
# summary table for storing AUC scores of every model

# http://zetcode.com/python/prettymtable/
from prettytable import PrettyTable
summary_table = PrettyTable(field_names = ['Vectorizer', 'Regularizer', 'Hyper parameter "C"',
                                           'AUC on Train Data', 'AUC on Test Data'])
```

In [7]:

```
data = project_data
y = list(data.project_is_approved.values)
data.drop(columns = 'project_is_approved', axis = 1, inplace = True)
```

In [142]:

```
# splitting the data into train, test, and cross validation
from sklearn.model_selection import train_test_split

X_data, X_test, y_data, y_test = train_test_split(data, y, test_size = 0.33, random_state = 0, stratify = y)
X_train, X_cv, y_train, y_cv = train_test_split(X_data, y_data, test_size = 0.3, random_state = 0, stratify = y_data)
```

In [9]:

```
c = Counter()
for i in y_train:
    c.update(str(i))
dict(c)
```

Out[9]:

```
{'1': 43479, '0': 7758}
```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [10]:

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

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


2.2.1 Encoding Numerical data

In [11]:

```
from sklearn.preprocessing import Normalizer
```

2.2.1.1 quantity

In [12]:

```
# train data
norm = Normalizer()
quantity_standardized = norm.fit_transform(X_train.quantity.values.reshape(-1, 1))
```

In [13]:

```
# cross validation data
cv_quantity_standardized = norm.transform(X_cv.quantity.values.reshape(-1, 1))
```

In [14]:

```
# test data
test_quantity_standardized = norm.transform(X_test.quantity.values.reshape(-1, 1))
```

2.2.1.2 teacher_number_of_previously_posted_projects

In [15]:

```
# train data
norm = Normalizer()
prev_posted_project_standardized =
norm.fit_transform(X_train.teacher_number_of_previously_posted_projects.values.reshape(-1, 1))
```

In [16]:

```
# cross validation data
cv_prev_posted_project_standardized =
norm.transform(X_cv.teacher_number_of_previously_posted_projects.values.reshape(-1, 1))
```

In [17]:

```
# test data
test_prev_posted_project_standardized =
norm.transform(X_test.teacher_number_of_previously_posted_projects.values.reshape(-1, 1))
```

2.2.1.3 price

In [18]:

```
# train data
norm = Normalizer()
price_standardized = norm.fit_transform(X_train.price.values.reshape(-1, 1))
```

In [19]:

```
# cross validation data
cv_price_standardized = norm.transform(X_cv.price.values.reshape(-1, 1))
```

In [20]:

```
# test data
test_price_standardized = norm.transform(X_test.price.values.reshape(-1, 1))
```

2.2.1.4 number of words in the title

In [215]:

```
# train data
norm = Normalizer()
word_count_title = norm.fit_transform(X_train.project_title.map(lambda x : len(x)).values.reshape(-1, 1))
```

In [216]:

```
# cross validation data
cv_word_count_title = norm.transform(X_cv.project_title.map(lambda x : len(x)).values.reshape(-1, 1))
```

In [217]:

```
# cross validation data
test_word_count_title = norm.transform(X_test.project_title.map(lambda x : len(x)).values.reshape(-1, 1))
```

2.2.1.5 number of words in the combine essays

In [218]:

```
# train data
norm = Normalizer()

# combine all essay (train)
essay = X_train.project_essay_1.map(str) + \
        X_train.project_essay_2.map(str) + \
        X_train.project_essay_3.map(str) + \
        X_train.project_essay_4.map(str)

word_count_essay = norm.fit_transform(essay.map(len).values.reshape(-1, 1))
```

In [219]:

```
# cross validation data

# combine all essay (cross validation)
essay = X_cv.project_essay_1.map(str) + \
        X_cv.project_essay_2.map(str) + \
        X_cv.project_essay_3.map(str) + \
        X_cv.project_essay_4.map(str)
cv_word_count_essay = norm.transform(essay.map(lambda x : len(x)).values.reshape(-1, 1))
```

In [220]:

```
# cross validation data
# combine all essay (test)
essay = X_test.project_essay_1.map(str) + \
        X_test.project_essay_2.map(str) + \
        X_test.project_essay_3.map(str) + \
        X_test.project_essay_4.map(str)
test_word_count_essay = norm.transform(essay.map(lambda x : len(x)).values.reshape(-1, 1))
```

2.2.1.6 sentiment score's of each of the essay

In [175]:

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

nltk.download('vader_lexicon')

ss = SentimentIntensityAnalyzer()
```

```
[nltk_data] Downloading package vader_lexicon to  
[nltk_data] C:\Users\Kamlesh\AppData\Roaming\nltk_data...
```

Out[175]:

True

In [211]:

```
def convert_to_array(ps):  
    """ Return array representation of polarity scores """  
    sentiment = list()  
    for p in ps:  
        temp = list()  
        for k in p.keys():  
            temp.append(p[k])  
        sentiment.append(np.array(temp))  
    return np.array(sentiment)
```

2.2.1.7.1 train data

In [212]:

```
# Finding polarity of each essay  
# essay1  
ps1 = convert_to_array(X_train.project_essay_1.map(lambda x : ss.polarity_scores(str(x))))  
# essay2  
ps2 = convert_to_array(X_train.project_essay_2.map(lambda x : ss.polarity_scores(str(x))))  
# essay3  
ps3 = convert_to_array(X_train.project_essay_3.map(lambda x : ss.polarity_scores(str(x))))  
# essay4  
ps4 = convert_to_array(X_train.project_essay_4.map(lambda x : ss.polarity_scores(str(x))))
```

2.2.1.7.2 cross validation data

In [213]:

```
# Finding polarity of each essay  
# essay1  
cv_ps1 = convert_to_array(X_cv.project_essay_1.map(lambda x : ss.polarity_scores(str(x))))  
# essay2  
cv_ps2 = convert_to_array(X_cv.project_essay_2.map(lambda x : ss.polarity_scores(str(x))))  
# essay3  
cv_ps3 = convert_to_array(X_cv.project_essay_3.map(lambda x : ss.polarity_scores(str(x))))  
# essay4  
cv_ps4 = convert_to_array(X_cv.project_essay_4.map(lambda x : ss.polarity_scores(str(x))))
```

2.2.1.7.1 test data

In [214]:

```
# Finding polarity of each essay  
# essay1  
test_ps1 = convert_to_array(X_test.project_essay_1.map(lambda x : ss.polarity_scores(str(x))))  
# essay2  
test_ps2 = convert_to_array(X_test.project_essay_2.map(lambda x : ss.polarity_scores(str(x))))  
# essay3  
test_ps3 = convert_to_array(X_test.project_essay_3.map(lambda x : ss.polarity_scores(str(x))))  
# essay4  
test_ps4 = convert_to_array(X_test.project_essay_4.map(lambda x : ss.polarity_scores(str(x))))
```

2.2.2 Preprocessing and Encoding Categorical Data

In [21]:

```
from sklearn.feature_extraction.text import CountVectorizer
```

2.2.2.1 project_subject_category

In [22]:

```
# replace & with _, remove spaces
def preprocess_subj_cat(subj_cat):
    subj_cat_list = list()
    for s in subj_cat:
        temp = ''
        for j in s.split(','):
            if 'The' in j.split(' '):
                j = j.replace('The', '')
            j = j.replace(' ', '')
            j = j.replace('&', '_')
            temp += j.strip() + ' '
        subj_cat_list.append(temp.strip())
    return subj_cat_list
```

In [23]:

```
# preprocessing train data
X_train['clean_category'] = preprocess_subj_cat(list(X_train.project_subject_categories.values))
X_train.drop(columns = 'project_subject_categories', axis = 1, inplace = True)
```

In [24]:

```
# preprocessing cross validation data
X_cv['clean_category'] = preprocess_subj_cat(list(X_cv.project_subject_categories.values))
X_cv.drop(columns = 'project_subject_categories', axis = 1, inplace = True)
```

In [25]:

```
# preprocessing test data
X_test['clean_category'] = preprocess_subj_cat(list(X_test.project_subject_categories.values))
X_test.drop(columns = 'project_subject_categories', axis = 1, inplace = True)
```

In [26]:

```
# encoding train data
cat_vectorizer = CountVectorizer(lowercase = False, binary = True)
category_one_hot = cat_vectorizer.fit_transform(X_train.clean_category.values)

print('Feature : ', cat_vectorizer.get_feature_names())
print('Shape of matrix after one hot encoding : ', category_one_hot.shape)
```

```
Feature :  ['AppliedLearning', 'Care_Hunger', 'Health_Sports', 'History_Civics',
'Literacy_Language', 'Math_Science', 'Music_Arts', 'SpecialNeeds', 'Warmth']
Shape of matrix after one hot encoding :  (51237, 9)
```

In [27]:

```
# encoding cross validation data
cv_category_one_hot = cat_vectorizer.transform(X_cv.clean_category.values)
print('Shape of matrix after one hot encoding : ', cv_category_one_hot.shape)
```

```
Shape of matrix after one hot encoding :  (21959, 9)
```

In [28]:

```
# encoding test data
test_category_one_hot = cat_vectorizer.transform(X_test.clean_category.values)
print('Shape of matrix after one hot encoding : ', test_category_one_hot.shape)
```

```
Shape of matrix after one hot encoding :  (36052, 9)
```

2.2.2.2 project_subject_subcategory

In [29]:

```
# replace & with _; remove spaces
def preprocess_subj_subcat(subj_subcat):
    subj_subcat_list = list()
    for s in subj_subcat:
        temp = ''
        for j in s.split(','):
            if 'The' in j.split(' '):
                j = j.replace('The', '')
            j = j.replace(' ', '')
            j = j.replace('&', '_')
            temp += j.strip() + ' '
        subj_subcat_list.append(temp.strip())
    return subj_subcat_list
```

In [30]:

```
# preprocessing train data
X_train['clean_subcategory'] = preprocess_subj_subcat(list(X_train.project_subject_subcategories.v
alues))
X_train.drop(columns = 'project_subject_subcategories', axis = 1, inplace = True)
```

In [31]:

```
# preprocessing cross validation data
X_cv['clean_subcategory'] = preprocess_subj_subcat(list(X_cv.project_subject_subcategories.values)
)
X_cv.drop(columns = 'project_subject_subcategories', axis = 1, inplace = True)
```

In [32]:

```
# preprocessing test data
X_test['clean_subcategory'] =
preprocess_subj_subcat(list(X_test.project_subject_subcategories.values))
X_test.drop(columns = 'project_subject_subcategories', axis = 1, inplace = True)
```

In [33]:

```
# encoding train data
subcat_vectorizer = CountVectorizer(lowercase = False, binary = True)
subcategory_one_hot = subcat_vectorizer.fit_transform(X_train.clean_subcategory.values)

print('Feature : ', subcat_vectorizer.get_feature_names())
print('Shape of matrix after one hot encoding : ', subcategory_one_hot.shape)
```

```
Feature :  ['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government',
'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics',
'EnvironmentalScience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness',
'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writing', 'M
athematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'Socia
lSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix after one hot encoding :  (51237, 30)
```

In [34]:

```
# encoding cross validation data
cv_subcategory_one_hot = subcat_vectorizer.transform(X_cv.clean_subcategory.values)
print('Shape of matrix after one hot encoding : ', cv_subcategory_one_hot.shape)
```

```
Shape of matrix after one hot encoding :  (21959, 30)
```

In [35]:

```
# encoding test data
test_subcategory_one_hot = subcat_vectorizer.transform(X_test.clean_subcategory.values)
print('Shape of matrix after one hot encoding : ', test_subcategory_one_hot.shape)
```

Shape of matrix after one hot encoding : (36052, 30)

2.2.2.3 project_grade_category

In [36]:

```
# replacing space by '_' and '-' by '_' in 'project_grade_category'
def preprocess_project_grade_category(grade_cat):
    project_grade_category = list()
    for p in grade_cat:
        p = p.strip()
        p = p.replace(' ', '_')
        p = p.replace('-', '_')
        project_grade_category.append(p.strip())
    return project_grade_category
```

In [37]:

```
# preprocessing train data
X_train['clean_grade_category'] =
preprocess_project_grade_category(list(X_train.project_grade_category.values))
X_train.drop(columns = 'project_grade_category', axis = 1, inplace = True)
```

In [38]:

```
# preprocessing cross validation data
X_cv['clean_grade_category'] = preprocess_project_grade_category(list(X_cv.project_grade_category.
values))
X_cv.drop(columns = 'project_grade_category', axis = 1, inplace = True)
```

In [39]:

```
# preprocessing test data
X_test['clean_grade_category'] =
preprocess_project_grade_category(list(X_test.project_grade_category.values))
X_test.drop(columns = 'project_grade_category', axis = 1, inplace = True)
```

In [40]:

```
# encoding train data
grade_vectorizer = CountVectorizer(lowercase = False, binary = True)
project_grade_category_one_hot =
grade_vectorizer.fit_transform(X_train.clean_grade_category.values)

print('Features : ', grade_vectorizer.get_feature_names())
print('Shape of matrix after one hot encoding : ', project_grade_category_one_hot.shape)
```

Features : ['Grades_3_5', 'Grades_6_8', 'Grades_9_12', 'Grades_PreK_2']
Shape of matrix after one hot encoding : (51237, 4)

In [41]:

```
# encoding cross validation data
cv_project_grade_category_one_hot = grade_vectorizer.transform(X_cv.clean_grade_category.values)
print('Shape of matrix after one hot encoding : ', cv_project_grade_category_one_hot.shape)
```

Shape of matrix after one hot encoding : (21959, 4)

In [42]:

```
# encoding test data
test_project_grade_category_one_hot =
grade_vectorizer.transform(X_test.clean_grade_category.values)
print('Shape of matrix after one hot encoding : ', test_project_grade_category_one_hot.shape)
```

Shape of matrix after one hot encoding : (36052, 4)

2.2.2.4 teacher_prefix

In [43]:

```
def preprocess_teacher_prefix(data):  
    # replacing nan value by 'nan'  
    data.teacher_prefix.replace(to_replace = np.nan, value = 'nan', inplace = True)  
  
    # removing '.'  
    clean_teacher_prefix = list()  
    for tp in data.teacher_prefix.values:  
        tp = tp.replace('.', '')  
        clean_teacher_prefix.append(tp)  
  
    data.teacher_prefix = clean_teacher_prefix
```

In [44]:

```
# preprocessing train data  
preprocess_teacher_prefix(X_train)
```

In [45]:

```
# preprocessing cross validation data  
preprocess_teacher_prefix(X_cv)
```

In [46]:

```
# preprocessing test data  
preprocess_teacher_prefix(X_test)
```

In [47]:

```
# encoding train data  
teacher_vectorizer = CountVectorizer(lowercase = False, binary = True)  
teacher_prefix_one_hot = teacher_vectorizer.fit_transform(X_train.teacher_prefix.values)  
  
print('Features : ', teacher_vectorizer.get_feature_names())  
print('Shape of matrix after one hot encoding : ', teacher_prefix_one_hot.shape)
```

Features : ['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher', 'nan']
Shape of matrix after one hot encoding : (51237, 6)

In [48]:

```
# encoding cross validation data  
cv_teacher_prefix_one_hot = teacher_vectorizer.transform(X_cv.teacher_prefix.values)  
print('Shape of matrix after one hot encoding : ', cv_teacher_prefix_one_hot.shape)
```

Shape of matrix after one hot encoding : (21959, 6)

In [49]:

```
# encoding test data  
test_teacher_prefix_one_hot = teacher_vectorizer.transform(X_test.teacher_prefix.values)  
print('Shape of matrix after one hot encoding : ', test_teacher_prefix_one_hot.shape)
```

Shape of matrix after one hot encoding : (36052, 6)

2.2.2.5 school_sates

In [50]:

```
# encoding train data
state_vectorizer = CountVectorizer(lowercase = False, binary = True)
school_states_one_hot = state_vectorizer.fit_transform(X_train.school_state.values)

print('Feature : ', state_vectorizer.get_feature_names())
print('Shape of matrix after one hot encoding : ', school_states_one_hot.shape)
```

```
Feature :  ['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 matrix after one hot encoding :  (51237, 51)
```

In [51]:

```
# encoding cross validation data
cv_school_states_one_hot = state_vectorizer.transform(X_cv.school_state.values)
print('Shape of matrix after one hot encoding : ', cv_school_states_one_hot.shape)
```

```
Shape of matrix after one hot encoding :  (21959, 51)
```

In [52]:

```
# encoding test data
test_school_states_one_hot = state_vectorizer.transform(X_test.school_state.values)
print('Shape of matrix after one hot encoding : ', test_school_states_one_hot.shape)
```

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

2.3 Make Data Model Ready: encoding eassay, and project_title

In [53]:

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

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

2.3.1 preprocessing essay and title

In [54]:

```
# https://stackoverflow.com/a/47091490
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 [55]:

```
import re

def remove_escape_sequences(phrase):
    # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
    phrase = re.sub(r'\\', ' ', phrase)
    phrase = re.sub(r'\\n', ' ', phrase)
    phrase = re.sub(r'\\r', ' ', phrase)
    phrase = re.sub(r'\\t', ' ', phrase)
    return phrase

def remove_special_characters(phrase):
    return re.sub(r'[^A-Za-z0-9]+', ' ', phrase)
```

In [56]:

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

2.3.1.1 preprocessing essay

In [57]:

```
def preprocess_essay(dataframe):
    # essay train data
    dataframe_essay = dataframe.project_essay_1.map(str) + \
        dataframe.project_essay_2.map(str) + \
        dataframe.project_essay_3.map(str) + \
        dataframe.project_essay_4.map(str)

    # removing stop words, escape sequences, special character
    preprocessed_essay = list()
    for e in tqdm(dataframe_essay):
        e = decontracted(e)
        e = remove_escape_sequences(e)
        e = remove_special_characters(e)

        temp = ' '.join([word.lower() for word in e.split() if word.lower() not in set(stopwords)])
        preprocessed_essay.append(temp)

    # adding new column of preprocessed essay in dataframe
    dataframe['preprocessed_essay'] = preprocessed_essay
    # removing project_essay columns from dataframe
    dataframe.drop(columns = ['project_essay_1', 'project_essay_2', 'project_essay_3',
'project_essay_4'], axis = 1, inplace = True)
```

```
# essay train data
preprocess_essay(X_train)
```

In [59]:

```
# essay cross validation data
preprocess_essay(X_cv)
```

In [60]:

```
# essay test data
preprocess_essay(X_test)
```

2.3.1.2 preprocessing title

```
def preprocess_title(dataframe):
    # removing stop words, escape sequences, special character
    preprocessed_title = list()
    for e in tqdm(dataframe.project_title):
        e = decontracted(e)
        e = remove_escape_sequences(e)
        e = remove_special_characters(e)

        temp = ' '.join([word.lower() for word in e.split() if word.lower() not in set(stopwords)])
        preprocessed_title.append(temp)
    # adding new column of preprocessed_title in dataframe
    dataframe['preprocessed_title'] = preprocessed_title
    # removing project_title column from dataframe
    dataframe.drop(columns = ['project title'], axis = 1, inplace = True)
```

In [62]:

```
# project_title train data
preprocess title(X train)
```

In [63]:

```
# project_title cross validation data
preprocess title(X cv)
```

In [64]:

```
# project_title test data
preprocess title(X test)
```

2.3.2 Vectorizing text data (train)

2.3.2.1 Bag Of Words (BOW)

In [65]:

```
# preprocessed_essay train data

# Considering the words which appeared in at least min_df documents (rows)
# using n_gram = (1, 2) means uni-gram and bi-gram
essay_bow_vectorizer = CountVectorizer(min_df = 10, ngram_range = (1, 2), max_features = 5000)
essay_bow = essay_bow_vectorizer.fit_transform(X_train.preprocessed_essay)
print("Shape of matrix after BOW ", essay_bow.shape)
```

```
Shape of matrix after BOW (51237, 5000)
```

In [66]:

```
# preprocessed_essay cross validation data
cv_essay_bow = essay_bow_vectorizer.transform(X_cv.preprocessed_essay)
print("Shape of matrix after BOW ", cv_essay_bow.shape)
```

```
Shape of matrix after BOW (21959, 5000)
```

In [67]:

```
# preprocessed_essay test data
test_essay_bow = essay_bow_vectorizer.transform(X_test.preprocessed_essay)
print("Shape of matrix after BOW ", test_essay_bow.shape)
```

```
Shape of matrix after BOW (36052, 5000)
```

In [68]:

```
# preprocessed_title train data

# Considering the words which appeared in at least min_df documents (rows or projects).
# using n_gram = (1, 2) means uni-gram and bi-gram
title_bow_vectorizer = CountVectorizer(min_df = 10, ngram_range = (1, 2), max_features = 5000)
title_bow = title_bow_vectorizer.fit_transform(X_train.preprocessed_title)
print("Shape of matrix after BOW ", title_bow.shape)
```

```
Shape of matrix after BOW (51237, 3366)
```

In [69]:

```
# cv_preprocessed_title cross validation data
cv_title_bow = title_bow_vectorizer.transform(X_cv.preprocessed_title)
print("Shape of matrix after BOW ", cv_title_bow.shape)
```

```
Shape of matrix after BOW (21959, 3366)
```

In [70]:

```
# cv_preprocessed_title test data
test_title_bow = title_bow_vectorizer.transform(X_test.preprocessed_title)
print("Shape of matrix after BOW ", test_title_bow.shape)
```

```
Shape of matrix after BOW (36052, 3366)
```

2.3.2.2 Term Frequency Inverse Document Frequency (TFIDF)

In [71]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

In [72]:

```
# preprocessed_essay train data

# Considering the words which appeared in at least min_df documents(rows or projects).
# using n_gram = 1
essay_tfidf_vectorizer = TfidfVectorizer(min_df = 10, ngram_range = (1, 2), max_features = 5000)
essay_tfidf = essay_tfidf_vectorizer.fit_transform(X_train.preprocessed_essay)
print("Shape of matrix after TFIDF ", essay_tfidf.shape)
```

Shape of matrix after TFIDF (51237, 5000)

In [73]:

```
# preprocessed_essay cross validation data
cv_essay_tfidf = essay_tfidf_vectorizer.transform(X_cv.preprocessed_essay)
print("Shape of matrix after TFIDF ", cv_essay_tfidf.shape)
```

Shape of matrix after TFIDF (21959, 5000)

In [74]:

```
# preprocessed_essay test data
test_essay_tfidf = essay_tfidf_vectorizer.transform(X_test.preprocessed_essay)
print("Shape of matrix after TFIDF ", test_essay_tfidf.shape)
```

Shape of matrix after TFIDF (36052, 5000)

In [75]:

```
# preprocessed_title train data

# Considering the words which appeared in at least min_df documents(rows or projects).
# using n_gram = 1
title_tfidf_vectorizer = TfidfVectorizer(min_df = 10, ngram_range = (1, 2), max_features = 5000)
title_tfidf = title_tfidf_vectorizer.fit_transform(X_train.preprocessed_title)
print("Shape of matrix after TFIDF ", title_tfidf.shape)
```

Shape of matrix after TFIDF (51237, 3366)

In [76]:

```
# preprocessed_title cross validation data
cv_title_tfidf = title_tfidf_vectorizer.transform(X_cv.preprocessed_title)
print("Shape of matrix after TFIDF ", cv_title_tfidf.shape)
```

Shape of matrix after TFIDF (21959, 3366)

In [77]:

```
# preprocessed_title test data
test_title_tfidf = title_tfidf_vectorizer.transform(X_test.preprocessed_title)
print("Shape of matrix after TFIDF ", test_title_tfidf.shape)
```

Shape of matrix after TFIDF (36052, 3366)

2.3.2.3 Average Word 2 Vec (W2V)

In [78]:

```
# importing glove w2v model
import pickle
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f) # model is a dictionary
    glove_words = set(model.keys())
```

In [79]:

```
def avg_w2v(data, model, glove_words):
    ''' Return the avg_w2v representaion of data'''
    avg_w2v_list = list()
    for sent in tqdm(data):
        temp_w2v = np.zeros(300)
        count_words = 0
        for word in sent.split():
            if word in glove_words:
                temp_w2v += model[word]
                count_words += 1
        if count_words != 0:
            temp_w2v /= count_words
        avg_w2v_list.append(temp_w2v)
    return avg_w2v_list
```

In [80]:

```
# preprocessed_essay train data
essay_avg_w2v = avg_w2v(X_train.preprocessed_essay, model, glove words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 51237/51237  
[00:10<00:00, 4886.47it/s]
```

In [81]:

```
# preprocessed_essay cross validation data
cv_essay_avg_w2v = avg_w2v(X_cv.preprocessed_essay, model, glove_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 21959/21959  
[00:04<00:00, 4837.69it/s]
```

In [82]:

```
# preprocessed_essay test data
test_essay_avg_w2v = avg_w2v(X_test.preprocessed_essay, model, glove_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 36052/36052  
[00:07<00:00, 4910.83it/s]
```

In [83]:

```
# preprocessed_title train data
title_avg_w2v = avg_w2v(X_train.preprocessed_title, model, glove_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 51237/51237  
[00:00<00:00, 95139.70it/s]
```

In [84]:

```
# preprocessed_title cross validation data
cv title avg w2v = avg w2v(X cv.preprocessed title, model, glove words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 21959/21959  
[00:00<00:00, 93309.87it/s]
```

In [85]:

```
# preprocessed_title test data
test_title_avg_w2v = avg_w2v(X.test.preprocessed_title, model, glove_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 36052/36052  
[00:00<00:00, 95046.95it/s]
```

2.3.2.4 TFIDF weighted Word 2 Vec (W2V)

In [86]:

```
def tfidf_w2v(data, model, glove_words, tfidf_model, tfidf_words):
    ''' Return the tfidf weighted w2v representaion of data'''
    tfidf_w2v_list = list()
    for sent in tqdm(data):
        temp_w2v = np.zeros(300)
        tfidf_weight = 0
        for word in sent.split():
            if word in glove_words and word in tfidf_words:
                vec = model[word]
                tfidf = tfidf_model[word] * (sent.count(word) / len(sent.split()))
                temp_w2v += vec * tfidf
                tfidf_weight += tfidf
        if tfidf_weight != 0:
            temp_w2v /= tfidf_weight
        tfidf_w2v_list.append(temp_w2v)
    return tfidf_w2v_list
```

In [87]:

```
# preprocessed_essay train data

vectorizer = TfidfVectorizer()
essay_tfidf_vec = vectorizer.fit(X_train.preprocessed_essay)

# Making a dictionary with key as word and value as idf value of that word
essay_tfidf_model = dict(zip(essay_tfidf_vec.get_feature_names(), essay_tfidf_vec.idf_))
essay_tfidf_words = set(essay_tfidf_vec.get_feature_names())

essay_tfidf_w2v = tfidf_w2v(X_train.preprocessed_essay, model, glove_words, essay_tfidf_model,
essay_tfidf_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 51237/51237 [01:  
19<00:00, 644.80it/s]
```

In [88]:

```
# preprocessed_essay cross validation data
cv_essay_tfidf_w2v = tfidf_w2v(X_cv.preprocessed_essay, model, glove_words, essay_tfidf_model,
essay_tfidf_words)
```

```
100%|███████████████████████████████████████████| 21959/21959 [00:  
34<00:00, 624.95it/s]
```

In [89]:

```
# preprocessed_essay cross validation data
test_essay_tfidf_w2v = tfidf_w2v(X_test.preprocessed_essay, model, glove_words, essay_tfidf_model,
essay_tfidf_words)
```

```
100%|██████████████████████████████████████████████████████████████████| 36052/36052 [00:  
55<00:00, 646.19it/s]
```

In [90]:

```
# preprocessed_title train data

vectorizer = TfidfVectorizer()
title tfidf vec = vectorizer.fit(X_train.preprocessed_title)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 51237/51237  
[00:01<00:00, 44545.88it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████| 21959/21959  
[00:00<00:00, 44212.31it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 36052/36052  
[00:00<00:00, 44408.36it/s]
```

```
# plot confusion matrix python; https://scikit-learn.org/stable/auto\_examples/model\_selection/plot\_confusion\_matrix.html
from sklearn.metrics import confusion_matrix
from sklearn.utils.multiclass import unique_labels

def plot_confusion_matrix(y_true, y_pred, classes, title = None, cmap = plt.cm.Blues):
    """
    This function prints and plots the confusion matrix.
    """
    if not title:
        title = 'Confusion matrix'

    # Compute confusion matrix
```

```

# Compute confusion matrix
cm = confusion_matrix(y_true, y_pred)
# Only use the labels that appear in the data
classes = classes[unique_labels(y_true, y_pred)]
print('Confusion matrix')
print(cm)

fig, ax = plt.subplots()
im = ax.imshow(cm, interpolation = 'nearest', cmap = cmap)
ax.figure.colorbar(im, ax = ax)
# We want to show all ticks...
ax.set(xticks=np.arange(cm.shape[1]),
      yticks=np.arange(cm.shape[0]),
      # ... and label them with the respective list entries
      xticklabels=classes, yticklabels=classes,
      title=title,
      ylabel='True label',
      xlabel='Predicted label')

# Rotate the tick labels and set their alignment.
plt.setp(ax.get_xticklabels(), rotation=45, ha = "right", rotation_mode = "anchor")

# Loop over data dimensions and create text annotations.
fmt = 'd'
thresh = cm.max() / 2.
for i in range(cm.shape[0]):
    for j in range(cm.shape[1]):
        ax.text(j, i, format(cm[i, j], fmt),
                ha="center", va="center",
                color="white" if cm[i, j] > thresh else "black")
fig.tight_layout()
plt.grid()
plt.show()
return ax

```

2.4.1 Applying Logistic Regression on BOW, SET 1

In [96]:

```
# Please write all the code with proper documentation
```

In [97]:

```

# train data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train = hstack((school_states_one_hot, category_one_hot, subcategory_one_hot, project_grade_category_one_hot,
                  teacher_prefix_one_hot, quantity_standardized, prev_posted_project_standardized, price_standardized,
                  essay_bow, title_bow)).tocsr()
print('X_train shape : ', X_train.shape)

y_train = np.array(y_train)
print('y_train shape : ', y_train.shape)

```

```

X_train shape : (51237, 8469)
y_train shape : (51237,)

```

In [98]:

```

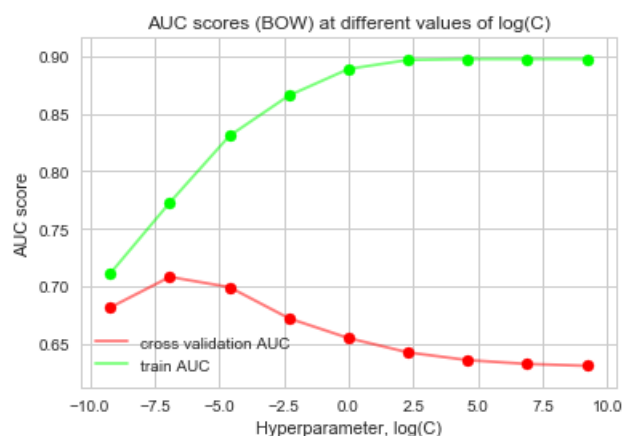
# cross validation data

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_cv = hstack((cv_school_states_one_hot, cv_category_one_hot, cv_subcategory_one_hot,
               cv_project_grade_category_one_hot, cv_teacher_prefix_one_hot,
               cv_quantity_standardized,
               cv_prev_posted_project_standardized, cv_price_standardized, cv_essay_bow, cv_title_bow)).tocsr()
print('X_cv shape : ', X_cv.shape)

```



```
plt.legend()
plt.xlabel('Hyperparameter, log(C)')
plt.ylabel('AUC score')
plt.title('AUC scores (BOW) at different values of log(C)')
plt.show()
```



In [103]:

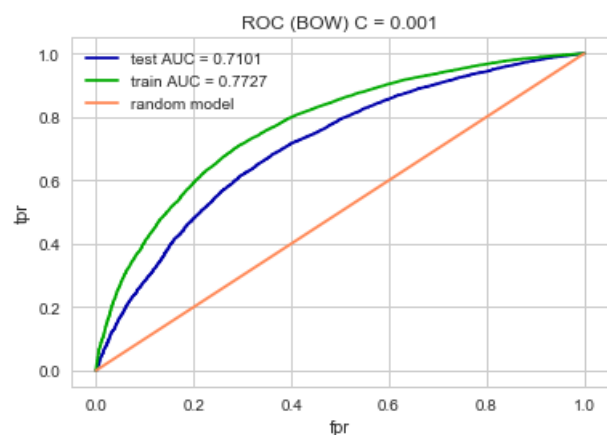
```
best_C = C_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best LR model (BOW)
best_LR_model_bow = LogisticRegression(penalty = 'l2', C = best_C, class_weight = 'balanced', n_jobs = -1)
# fitting the train data
best_LR_model_bow.fit(X_train, y_train)

# predicting the probability scores of train data and test data
predicted_test = batch_predict(best_LR_model_bow, X_test, 400)
predicted_train = batch_predict(best_LR_model_bow, X_train, 400)

# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)
```

In [104]:

```
# Plotting ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test, predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train, predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC (BOW) C = {best_C}')
plt.show()
```



In [105]:

```

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

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

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

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

```

In [106]:

```

# adding AUC scores to summary table
summary_table.add_row(['BOW', 'L2', best_C, '%.4f' % (roc_auc_score(y_train, predicted_train)),
                      '%.4f' % (roc_auc_score(y_test, predicted_test))])

# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

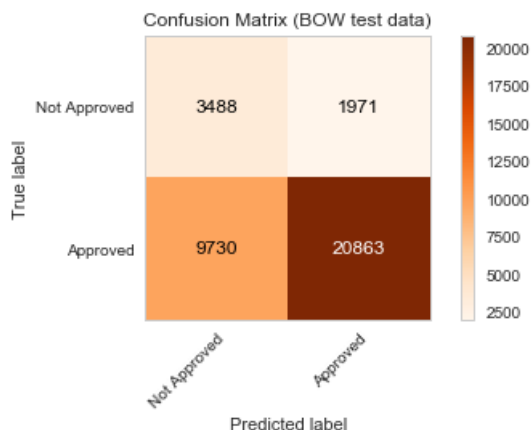
# Plotting confusion matrix for test data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (BOW test data)', cmap = plt.cm.Oranges)

```

the maximum value of $tpr*(1-fpr)$ 0.4357306131981015 for threshold 0.498

Confusion matrix

```
[[ 3488  1971]
 [ 9730 20863]]
```



Out[106]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c5038c88>

In [107]:

```

# predicting class labels for train data
pred_train = predict(predicted_train, train_th, train_fpr, train_tpr)

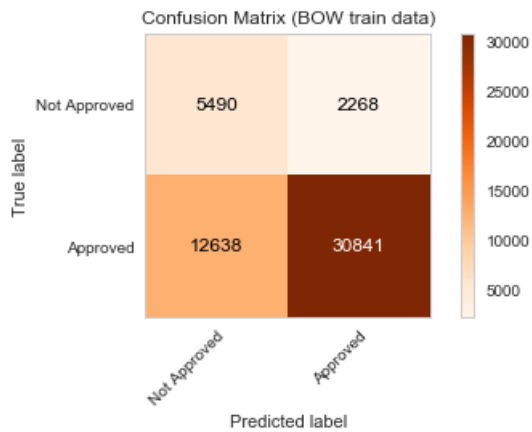
# Plotting confusion matrix for train data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_train, pred_train, classes = class_names,
                      title = 'Confusion Matrix (BOW train data)', cmap = plt.cm.Oranges)

```

the maximum value of $tpr*(1-fpr)$ 0.5019627311347308 for threshold 0.492

Confusion matrix

```
[[ 5490  2268]
 [12638 30841]]
```



Out[107]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c502d1d0>

2.4.2 Applying Logistic Regression on TFIDF, SET 2

In [108]:

```
# Please write all the code with proper documentation
```

In [109]:

```
# train data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train = hstack((school_states_one_hot, category_one_hot, subcategory_one_hot, project_grade_category_one_hot,
                  teacher_prefix_one_hot, quantity_standardized, prev_posted_project_standardized,
                  price_standardized, essay_tfidf, title_tfidf)).tocsr()
print('X_train shape : ', X_train.shape)

y_train = np.array(y_train)
print('y_train shape : ', y_train.shape)
```

```
X_train shape : (51237, 8469)
y_train shape : (51237,)
```

In [110]:

```
# cross validation data

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_cv = hstack((cv_school_states_one_hot, cv_category_one_hot, cv_subcategory_one_hot,
               cv_project_grade_category_one_hot, cv_teacher_prefix_one_hot,
               cv_quantity_standardized,
               cv_prev_posted_project_standardized, cv_price_standardized, cv_essay_tfidf, cv_title_tfidf)).tocsr()
print('X_cv shape : ', X_cv.shape)

y_cv = np.array(y_cv)
print('y_cv shape : ', y_cv.shape)
```

```
X_cv shape : (21959, 8469)
y_cv shape : (21959,)
```

In [111]:

```
# test data
```

```
X_test shape : (36052, 8469)
y_test shape : (36052,)
```

```
# training the LR_model_tfidf
C_range = list(map(lambda x : 10 ** x, range(-4, 5, 1)))
cv_auc_scores = dict()
train_auc_scores = dict()

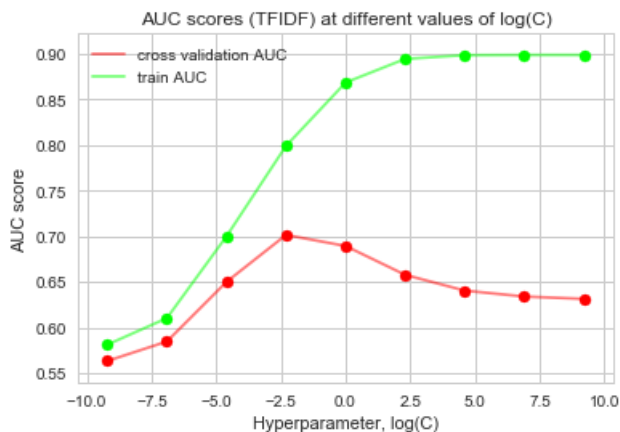
for C in tqdm(C_range):
    LR_model_tfidf = LogisticRegression(penalty = 'l2', C = C, class_weight = 'balanced', n_jobs =
-1)

    # fitting the train data
    LR_model_tfidf.fit(X_train, y_train)

    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(LR_model_tfidf, X_train, batch_size = 400)
    predicted_cv = batch_predict(LR_model_tfidf, X_cv, batch_size = 400)

    # calculating AUC score for cross validation and test data
    cv_auc_scores[C] = roc_auc_score(y_cv, predicted_cv)
    train_auc_scores[C] = roc_auc_score(y_train, predicted_train)
```

```
# plotting AUC scores
sns.set(style = 'whitegrid')
plt.plot(np.log(list(cv_auc_scores.keys())), list(cv_auc_scores.values()), color = '#FF000088', label = 'cross validation AUC')
plt.plot(np.log(list(train_auc_scores.keys())), list(train_auc_scores.values()), color = '#00FF0088', label = 'train AUC')
plt.scatter(np.log(list(cv_auc_scores.keys())), list(cv_auc_scores.values()), color = '#FF0000FF')
plt.scatter(np.log(list(train_auc_scores.keys())), list(train_auc_scores.values()), color = '#00FF00FF')
plt.legend()
plt.xlabel('Hyperparameter, log(C)')
plt.ylabel('AUC score')
plt.title('AUC scores (TFIDF) at different values of log(C)')
plt.show()
```



In [114]:

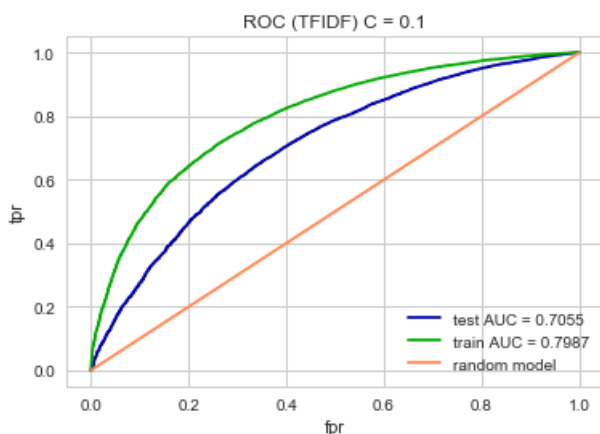
```
best_C = C_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best LR model (TFIDF)
best_LR_model_tfidf = LogisticRegression(penalty = 'l2', C = best_C, class_weight = 'balanced', n_jobs = -1)
# fitting the train data
best_LR_model_tfidf.fit(X_train, y_train)

# predicting the probability scores of train data and test data
predicted_test = batch_predict(best_LR_model_tfidf, X_test, 400)
predicted_train = batch_predict(best_LR_model_tfidf, X_train, 400)

# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)
```

In [115]:

```
# Plotting ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test, predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train, predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC (TFIDF) C = {best_C}')
plt.show()
```



In [116]:

```
# adding AUC scores to summary table
summary_table.add_row(['TFIDF', 'L2', best_C, '%.4f' % (roc_auc_score(y_train, predicted_train)), '%.4f' % (roc_auc_score(y_test, predicted_test))])

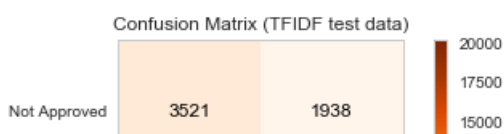
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

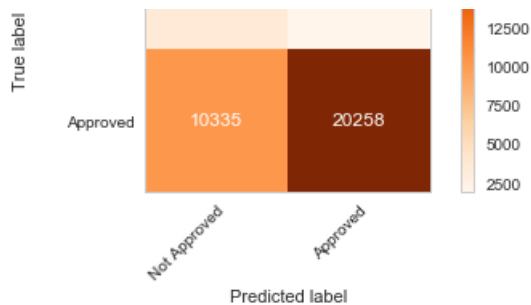
# Plotting confusion matrix for test data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (TFIDF test data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4270978948947868 for threshold 0.513

Confusion matrix

```
[[ 3521  1938]
 [10335 20258]]
```





Out[116]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c705d1d0>

In [117]:

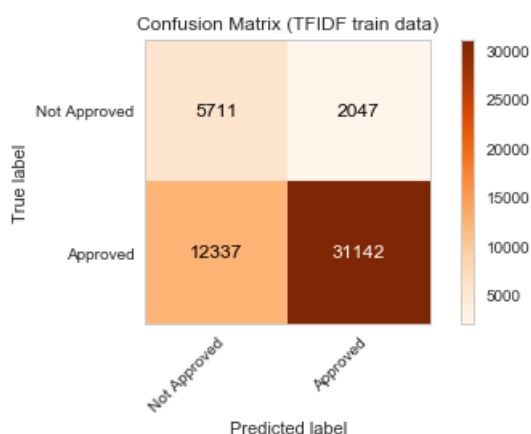
```
# predicting class labels for train data
pred_train = predict(predicted_train, train_th, train_fpr, train_tpr)

# Plotting confusion matrix for train data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_train, pred_train, classes = class_names,
                      title = 'Confusion Matrix (TFIDF train data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.5272654791267105 for threshold 0.499

Confusion matrix

```
[[ 5711  2047]
 [12337 31142]]
```



Out[117]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c4f6ae10>

2.4.3 Applying Logistic Regression on AVG W2V, SET 3

In [118]:

```
# Please write all the code with proper documentation
```

In [119]:

```
# train data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train = hstack((school_states_one_hot, category_one_hot, subcategory_one_hot, project_grade_category_one_hot,
                  teacher_prefix_one_hot, quantity_standardized, prev_posted_project_standardized, price_standardized,
                  essay_avg_w2v, title_avg_w2v))

# Logistic Regression
```

```
essay_avg_w2v, title_avg_w2v)).tocsr()
print('X_train shape : ', X_train.shape)
```

```
y_train = np.array(y_train)
print('y_train shape : ', y_train.shape)
```

```
X_train shape : (51237, 703)
y_train shape : (51237,)
```

In [120]:

```
# cross validation data
```

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

```
X_cv = hstack((cv_school_states_one_hot, cv_category_one_hot, cv_subcategory_one_hot,
               cv_teacher_prefix_one_hot, cv_quantity_standardized,
               cv_prev_posted_project_standardized,
               cv_price_standardized, cv_essay_avg_w2v, cv_title_avg_w2v)).tocsr()
print('X_cv shape : ', X_cv.shape)
```

```
y_cv = np.array(y_cv)
print('y_cv shape : ', y_cv.shape)
```

```
X_cv shape : (21959, 703)
y_cv shape : (21959,)
```

In [121]:

```
# test data
```

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

```
X_test = hstack((test_school_states_one_hot, test_category_one_hot, test_subcategory_one_hot,
                 test_project_grade_category_one_hot, test_teacher_prefix_one_hot,
                 test_quantity_standardized,
                 test_prev_posted_project_standardized, test_price_standardized,
                 test_essay_avg_w2v, test_title_avg_w2v)).tocsr()
print('X_test shape : ', X_test.shape)
```

```
y_test = np.array(y_test)
print('y_test shape : ', y_test.shape)
```

```
X_test shape : (36052, 703)
y_test shape : (36052,)
```

In [122]:

```
# training the LR_model_avg_w2v
```

```
C_range = list(map(lambda x : 10 ** x, range(-4, 5, 1)))
```

```
cv_auc_scores = dict()
```

```
train_auc_scores = dict()
```

```
for C in tqdm(C_range):
```

```
    LR_model_avg_w2v = LogisticRegression(penalty = 'l2', C = C, class_weight = 'balanced', n_jobs
    = -1)
```

```
    # fitting the train data
```

```
    LR_model_avg_w2v.fit(X_train, y_train)
```

```
    # predicting the probability scores of train data, cross validation data
```

```
    predicted_train = batch_predict(LR_model_avg_w2v, X_train, batch_size = 400)
```

```
    predicted_cv = batch_predict(LR_model_avg_w2v, X_cv, batch_size = 400)
```

```
    # calculating AUC score for cross validation and test data
```

```
    cv_auc_scores[C] = roc_auc_score(y_cv, predicted_cv)
```

```
    train_auc_scores[C] = roc_auc_score(y_train, predicted_train)
```

```
100% |████████████████████████████████████████████████████████████████████████████████| 9/9 [06
:37<00:00, 73.74s/it]
```

In [123]:

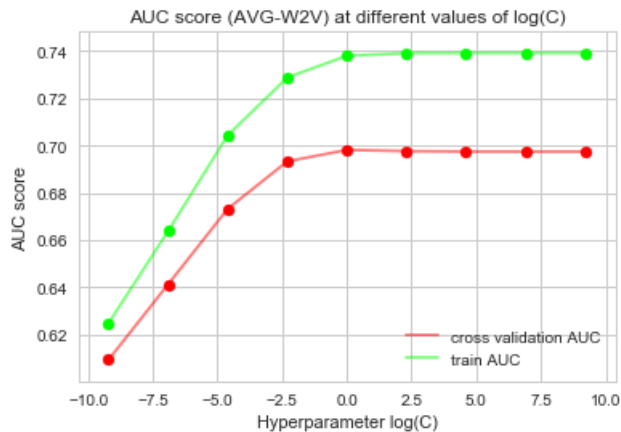
```
# plotting auc curve
```



```

sns.set(style = 'whitegrid')
plt.plot(np.log(list(cv_auc_scores.keys())) , list(cv_auc_scores.values()), color = '#FF000088', label = 'cross validation AUC')
plt.plot(np.log(list(train_auc_scores.keys())) , list(train_auc_scores.values()), color = '#00FF0088', label = 'train AUC')
plt.scatter(np.log(list(cv_auc_scores.keys())) , list(cv_auc_scores.values()), color = '#FF0000FF')
plt.scatter(np.log(list(train_auc_scores.keys())) , list(train_auc_scores.values()), color = '#00FF00FF')
plt.legend()
plt.xlabel('Hyperparameter log(C)')
plt.ylabel('AUC score')
plt.title('AUC score (AVG-W2V) at different values of log(C)')
plt.show()

```



In [124]:

```

best_C = C_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best LR model (AVG-W2V)
best_LR_model_avg_w2v = LogisticRegression(penalty = 'l2', C = best_C, class_weight = 'balanced', n_jobs = -1)
# fitting the train data
best_LR_model_avg_w2v.fit(X_train, y_train)

# predicting the probability scores of train data and test data
predicted_test = batch_predict(best_LR_model_avg_w2v, X_test, 400)
predicted_train = batch_predict(best_LR_model_avg_w2v, X_train, 400)

# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)

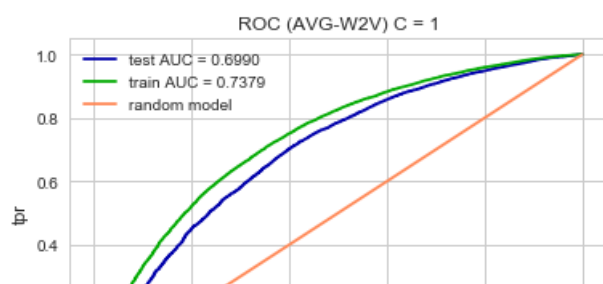
```

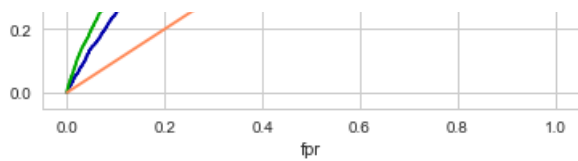
In [125]:

```

# Plots ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test, predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train, predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC (AVG-W2V) C = {best_C}')
plt.show()

```





In [126]:

```
# adding AUC scores to summary table
summary_table.add_row(['AVG-W2V', 'L2', best_C, '%.4f' % (roc_auc_score(y_train, predicted_train)),
                      '%.4f' % (roc_auc_score(y_test, predicted_test))])

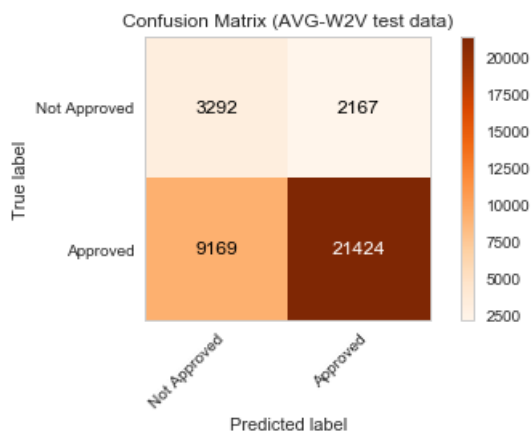
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

# Plotting confusion matrix for test data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (AVG-W2V test data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.42230402934695255 for threshold 0.477

Confusion matrix

```
[[ 3292  2167]
 [ 9169 21424]]
```



Out[126]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c7472128>

In [127]:

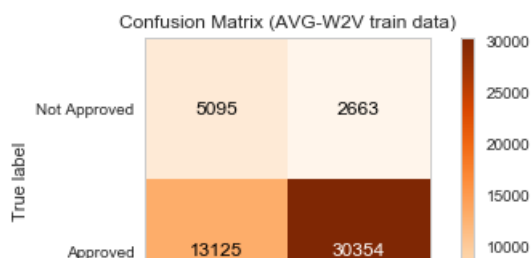
```
# predicting class labels for train data
pred_train = predict(predicted_train, train_th, train_fpr, train_tpr)

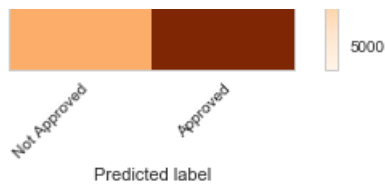
# Plotting confusion matrix for train data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_train, pred_train, classes = class_names,
                      title = 'Confusion Matrix (AVG-W2V train data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4584909798219432 for threshold 0.482

Confusion matrix

```
[[ 5095  2663]
 [13125 30354]]
```





Out[127]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c44023c8>

2.4.4 Applying Logistic Regression on TFIDF W2V, SET 4

In [128]:

```
# Please write all the code with proper documentation
```

In [129]:

```
# train data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train = hstack((school_states_one_hot, category_one_hot, subcategory_one_hot, project_grade_category_one_hot,
                  teacher_prefix_one_hot, quantity_standardized, prev_posted_project_standardized, price_standardized,
                  essay_tfidf_w2v, title_tfidf_w2v)).tocsr()
print('X_train shape : ', X_train.shape)

y_train = np.array(y_train)
print('y_train shape : ', y_train.shape)

X_train shape : (51237, 703)
y_train shape : (51237,)
```

In [130]:

```
# cross validation data

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_cv = hstack((cv_school_states_one_hot, cv_category_one_hot, cv_subcategory_one_hot,
               cv_project_grade_category_one_hot,
               cv_teacher_prefix_one_hot, cv_quantity_standardized,
               cv_prev_posted_project_standardized,
               cv_price_standardized, cv_essay_tfidf_w2v, cv_title_tfidf_w2v)).tocsr()
print('X_cv shape : ', X_cv.shape)

y_cv = np.array(y_cv)
print('y_cv shape : ', y_cv.shape)

X_cv shape : (21959, 703)
y_cv shape : (21959,)
```

In [131]:

```
# test data

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_test = hstack((test_school_states_one_hot, test_category_one_hot, test_subcategory_one_hot,
                 test_project_grade_category_one_hot, test_teacher_prefix_one_hot,
                 test_quantity_standardized,
                 test_prev_posted_project_standardized, test_price_standardized,
                 test_essay_tfidf_w2v, test_title_tfidf_w2v)).tocsr()
print('X_test shape : ', X_test.shape)

y_test = np.array(y_test)
```

```
X_test shape : (36052, 703)
y test shape : (36052,)
```

```
# training the LR_model_tfidf_w2v
C_range = list(map(lambda x : 10 ** x, range(-4, 5, 1)))
cv_auc_scores = dict()
train_auc_scores = dict()

for C in tqdm(C_range):
    LR_model_tfidf_w2v = LogisticRegression(penalty = 'l2', C = C, class_weight = 'balanced', n_jobs = -1)
    # fitting the train data
    LR_model_tfidf_w2v.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(LR_model_tfidf_w2v, X_train, batch_size = 400)
    predicted_cv = batch_predict(LR_model_tfidf_w2v, X_cv, batch_size = 400)

    # calculating AUC score for cross validation and test data
    cv_auc_scores[C] = roc_auc_score(y_cv, predicted_cv)
    train_auc_scores[C] = roc_auc_score(y_train, predicted_train)
```

In [133]:

AUC score (TFIDF-W2V) at different values of $\log(C)$

Hyperparameter $\log(C)$	cross validation AUC	train AUC
-10.0	0.630	0.643
-7.5	0.665	0.685
-5.0	0.687	0.716
-2.5	0.692	0.728
0.0	0.691	0.730
2.5	0.690	0.730
5.0	0.690	0.730
7.5	0.690	0.730
10.0	0.690	0.730

```
best_C = C_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best LR model (TFIDF-W2V)
best_LR_model_tfidf_w2v = LogisticRegression(penalty = 'l2', C = best_C, class_weight = 'balanced',
n_jobs = -1)
# fitting the train data
best_LR_model_tfidf_w2v.fit(X_train, y_train)

# predicting the probability scores of train data and test data
```

```

predicted_test = batch_predict(best_LR_model_tfidf_w2v, X_test, 400)
predicted_train = batch_predict(best_LR_model_tfidf_w2v, X_train, 400)

# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)

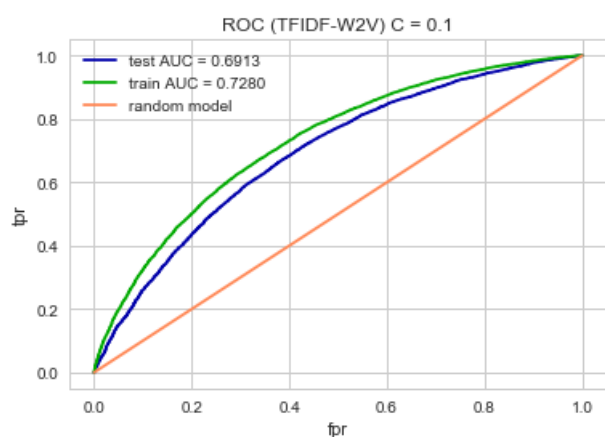
```

In [135]:

```

# Plots ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test,
predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train,
predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC (TFIDF-W2V) C = {best_C}')
plt.show()

```



In [136]:

```

# adding AUC scores to summary table
summary_table.add_row(['TFIDF-W2V', 'L2', best_C, '%.4f' % (roc_auc_score(y_train, predicted_train)
),
                        '%.4f' % (roc_auc_score(y_test, predicted_test))])

# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

# Plotting confusion matrix for test data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (TFIDF-W2V test data)', cmap = plt.cm.Oranges)

```

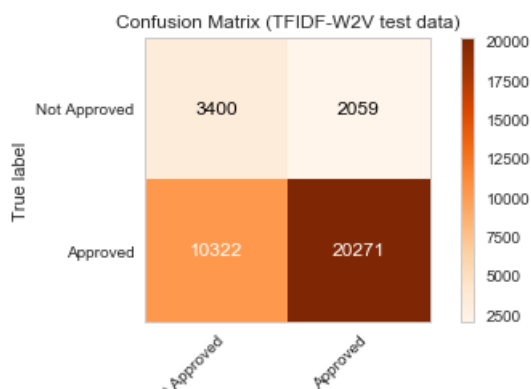
the maximum value of $tpr \cdot (1 - fpr)$ 0.41268523371991167 for threshold 0.486

Confusion matrix

```

[[ 3400  2059]
 [10322 20271]]

```



No

Predicted label

Out[136]:

<matplotlib.axes._subplots.AxesSubplot at 0x267ca9f20b8>

In [137]:

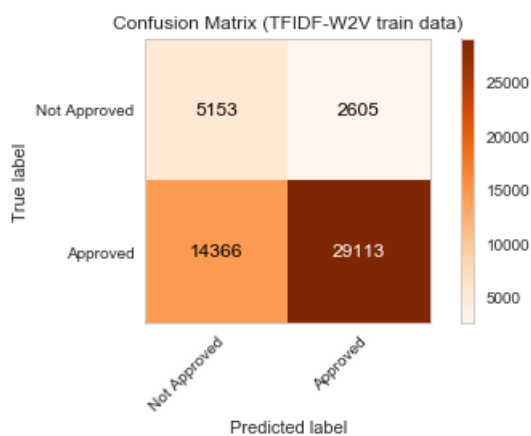
```
# predicting class labels for train data
pred_train = predict(predicted_train, train_th, train_fpr, train_tpr)

# Plotting confusion matrix for train data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_train, pred_train, classes = class_names,
                      title = 'Confusion Matrix (TFIDF-W2V train data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.44475186780809006 for threshold 0.487

Confusion matrix

```
[[ 5153  2605]
 [14366 29113]]
```



Out[137]:

<matplotlib.axes._subplots.AxesSubplot at 0x267c410fef0>

2.5 Logistic Regression with added Features `Set 5`

In [0]:

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

2.5.1 Encoding data for set 5

In [223]:

```
# train data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matirx :)
X_train = hstack((school_states_one_hot, category_one_hot, subcategory_one_hot, project_grade_category_one_hot,
```

```

        teacher_prefix_one_hot, quantity_standardized, prev_posted_project_standardized, price_standardized,
        word_count_essay, word_count_title, ps1, ps2, ps3, ps4)).tocsr()
print('X_train shape : ', X_train.shape)

y_train = np.array(y_train)
print('y_train shape : ', y_train.shape)

```

X_train shape : (51237, 121)
y_train shape : (51237,)

In [226]:

```

# cross validation data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_cv = hstack((cv_school_states_one_hot, cv_category_one_hot, cv_subcategory_one_hot,
               cv_project_grade_category_one_hot,
               cv_teacher_prefix_one_hot, cv_quantity_standardized,
               cv_prev_posted_project_standardized, cv_price_standardized,
               cv_word_count_essay, cv_word_count_title, cv_ps1, cv_ps2, cv_ps3, cv_ps4)).tocsr()
print('X_cv shape : ', X_cv.shape)

y_cv = np.array(y_cv)
print('y_cv shape : ', y_cv.shape)

```

X_cv shape : (21959, 121)
y_cv shape : (21959,)

In [227]:

```

# cross validation data

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

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_test = hstack((test_school_states_one_hot, test_category_one_hot, test_subcategory_one_hot,
                 test_project_grade_category_one_hot, test_teacher_prefix_one_hot,
                 test_quantity_standardized,
                 test_prev_posted_project_standardized, test_price_standardized,
                 test_word_count_essay,
                 test_word_count_title, test_ps1, test_ps2, test_ps3, test_ps4)).tocsr()
print('X_test shape : ', X_test.shape)

y_test = np.array(y_test)
print('y_test shape : ', y_test.shape)

```

X_test shape : (36052, 121)
y_test shape : (36052,)

In [228]:

```

# training the LR_model
C_range = list(map(lambda x : 10 ** x, range(-4, 5, 1)))
cv_auc_scores = dict()
train_auc_scores = dict()

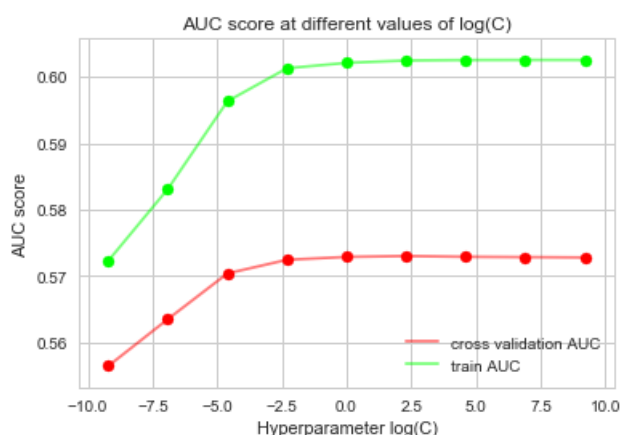
for C in tqdm(C_range):
    LR_model = LogisticRegression(penalty = 'l2', C = C, class_weight = 'balanced', n_jobs = -1)
    # fitting the train data
    LR_model.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(LR_model, X_train, batch_size = 400)
    predicted_cv = batch_predict(LR_model, X_cv, batch_size = 400)

    # calculating AUC score for cross validation and test data
    cv_auc_scores[C] = roc_auc_score(y_cv, predicted_cv)
    train_auc_scores[C] = roc_auc_score(y_train, predicted_train)

```

In [229]:

```
# plotting auc curve
sns.set(style = 'whitegrid')
plt.plot(np.log(list(cv_auc_scores.keys())), list(cv_auc_scores.values()), color = '#FF000088', label = 'cross validation AUC')
plt.plot(np.log(list(train_auc_scores.keys())), list(train_auc_scores.values()), color = '#00FF0088', label = 'train AUC')
plt.scatter(np.log(list(cv_auc_scores.keys())), list(cv_auc_scores.values()), color = '#FF0000FF')
plt.scatter(np.log(list(train_auc_scores.keys())), list(train_auc_scores.values()), color = '#00FF00FF')
plt.legend()
plt.xlabel('Hyperparameter log(C)')
plt.ylabel('AUC score')
plt.title('AUC score at different values of log(C)')
plt.show()
```



In [230]:

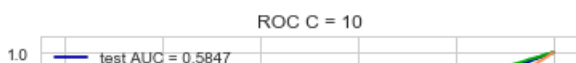
```
best_C = C_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best LR model
best_LR_model = LogisticRegression(penalty = 'l2', C = best_C, class_weight = 'balanced', n_jobs = -1)
# fitting the train data
best_LR_model.fit(X_train, y_train)

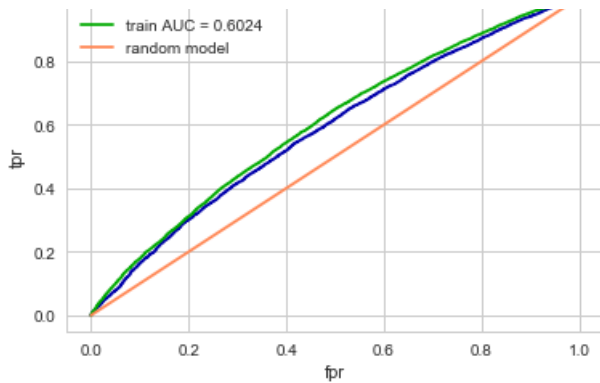
# predicting the probability scores of train data and test data
predicted_test = batch_predict(best_LR_model, X_test, 400)
predicted_train = batch_predict(best_LR_model, X_train, 400)

# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)
```

In [231]:

```
# Plots ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test, predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train, predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC C = {best_C}')
plt.show()
```





In [232]:

```
# adding AUC scores to summary table
summary_table.add_row(['Sentiment Score', 'L2', best_C, '%.4f' % (roc_auc_score(y_train, predicted_train)),
                      '%.4f' % (roc_auc_score(y_test, predicted_test))])

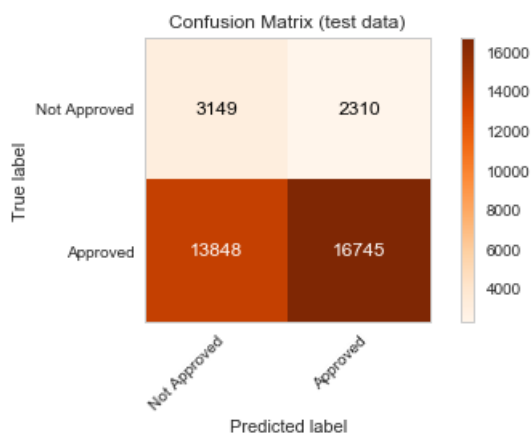
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

# Plotting confusion matrix for test data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (test data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.3157349449877268 for threshold 0.508

Confusion matrix

```
[[ 3149  2310]
 [13848 16745]]
```



Out[232]:

<matplotlib.axes._subplots.AxesSubplot at 0x267cbcb5ef0>

In [233]:

```
# predicting class labels for train data
pred_train = predict(predicted_train, train_th, train_fpr, train_tpr)

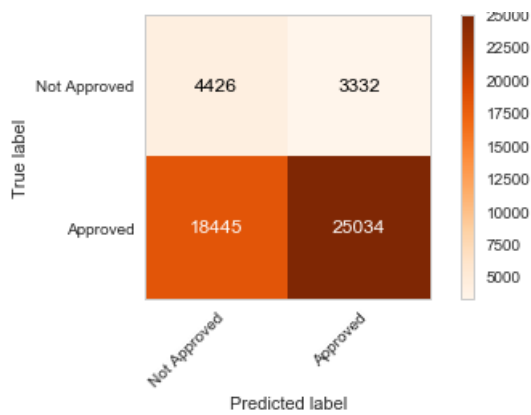
# Plotting confusion matrix for train data
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_train, pred_train, classes = class_names,
                      title = 'Confusion Matrix (train data)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.32848257408445913 for threshold 0.501

Confusion matrix

```
[[ 4426  3332]
 [18445 25034]]
```

Confusion Matrix (train data)



Out[233]:

<matplotlib.axes._subplots.AxesSubplot at 0x267ca01e320>

3. Conclusion

In [234]:

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

In [235]:

```
print(summary_table)
```

Vectorizer	Regularizer	Hyper parameter "C"	AUC on Train Data	AUC on Test Data
BOW	L2	0.001	0.7727	0.7101
TFIDF	L2	0.1	0.7987	0.7055
AVG-W2V	L2	1	0.7379	0.6990
TFIDF-W2V	L2	0.1	0.7280	0.6913
Sentiment Score	L2	10	0.6024	0.5847

1. By using BOW vectorizer on text data, both Train and Test AUC scores are high compare to others
2. Every model (except sentiment score) gives almost similar performance on test data but model(using TFIDF) gives highest AUC score on train data