

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

| Feature | Description |
|--|---|
| <code>project_id</code> | A unique identifier for the proposed project. Example: p036502 |
| <code>project_title</code> | Title of the project. Examples: <ul style="list-style-type: none">• Art Will Make You Happy!• First Grade Fun |
| <code>project_grade_category</code> | Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none">• Grades PreK-2• Grades 3-5• Grades 6-8• Grades 9-12 |
| <code>project_subject_categories</code> | One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none">• Applied Learning• Care & Hunger• Health & Sports• History & Civics• Literacy & Language• Math & Science• Music & The Arts• Special Needs• Warmth Examples: <ul style="list-style-type: none">• Music & The Arts• Literacy & Language, Math & Science |
| <code>school_state</code> | State where school is located (Two-letter U.S. postal code). Example: WY |
| <code>project_subject_subcategories</code> | One or more (comma-separated) subject subcategories for the project. Examples: <ul style="list-style-type: none">• Literacy |

| Feature | Description |
|---|---|
| <code>project_resource_summary</code> | An explanation of the resources needed for the project. Example: <ul style="list-style-type: none"> • My students need hands on literacy materials to manage sensory needs! |
| <code>project_essay_1</code> | First application essay* |
| <code>project_essay_2</code> | Second application essay* |
| <code>project_essay_3</code> | Third application essay* |
| <code>project_essay_4</code> | Fourth application essay* |
| <code>project_submitted_datetime</code> | Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245 |
| <code>teacher_id</code> | A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56 |
| <code>teacher_prefix</code> | Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher. |
| <code>teacher_number_of_previously_posted_projects</code> | Number of project applications previously submitted by the same teacher. Example: 2 |

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

- `__project_essay_1__` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."

your neighborhood, and your school are all helpful.

- `__project_essay_2__` "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

```
C:\Users\Shashank\Anaconda3\lib\site-packages\gensim\utils.py:1209: UserWarning: detected Windows;
aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

1.1 Reading Data

In [2]:

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

In [3]:

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

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

```
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories']
```

```
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

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

```
Number of data points in train data (1541272, 4)
['id' 'description' 'quantity' 'price']
```

Out[4]:

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

1.2 preprocessing of project_subject_categories

In [5]:

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

```
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('&', '_')
        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 [7]:

```
# 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 [8]:

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

Out[8]:

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

In [9]:

```
y=project_data['project_is_approved']
```

In [10]:

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

In [11]:

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

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English 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 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 of picture will not be taken before the first day of school.

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!nannan

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

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is a makeup of 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 on smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the Bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\n\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letters, words and pictures for students to learn about different letters and it is more accessible.nannan

In [12]:

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

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

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

In [13]:

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

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

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 [14]:

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

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

In [15]:

```
#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 [16]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
            'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', \
            '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', 'ma', 'may', 'mayn't', 'mightn', "mightn't", 'shan', "shan't", 'shouldn', "shouldn't", 'won', "won't", 'wouldn', "wouldn't"]
```


◀ ▶

```

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

project_data['clean_project_title'] = project_title_list
project_data.drop(['project_title'], axis=1, inplace=True)

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

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

```

In [23]:

```
project_data['project_title_list'] = project_title_list
```

In [24]:

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

1.5 Preparing data for models

In [25]:

```
project_data.columns
```

Out[25]:

```

Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category',
      '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', 'preprocessed_essays',
      'project_title_list'],
      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 [26]:

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

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

In [27]:

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

['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 encoding (109248, 30)
```

In [28]:

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

In [29]:

```
#onehotencoding for school_state
one_hot_encoding_school_state=pd.get_dummies(project_data.school_state)

print("Shape of dataframe for school_state", one_hot_encoding_school_state.shape)
```

Shape of dataframe for school_state (109248, 51)

In [30]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix=pd.get_dummies(project_data.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix.shape)
```

Shape of dataframe for teacher_prefix (109248, 5)

In [31]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category=pd.get_dummies(project_data.project_grade_category)

print("Shape of dataframe for project_grade_category", one_hot_encoding_project_grade_category.shape)
```

Shape of dataframe for project_grade_category (109248, 4)

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [32]:

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

Shape of matrix after one hot encoding (109248, 16623)

In [33]:

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

In [34]:

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

Shape of matrix after one hot encoding (109248, 3349)

1.5.2.2 TFIDF vectorizer

In [35]:

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

Shape of matrix after one hot encoding (109248, 16623)

1.5.2.3 Using Pretrained Models: Avg W2V

In [36]:

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

# =====
Output:

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

# =====
```

```

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

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

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove 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-save-and-load-variables-in-python/

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

'''

```

Out[36]:

```

'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\r',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n
word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        m
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel('\glove.42B.300d.txt')\n\n# =====\n\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preproced_texts:\n    words.extend(i.split('\
'))\n\nfor i in preproced_titles:\n    words.extend(i.split('\ '))\nprint("all the words in the
coupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus",
len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha
t are present in both glove vectors and our coupus", len(inter_words),
(" , np.round(len(inter_words)/len(words)*100,3), "%) ") \n\nwords_courpus = {}\nwords_glove =
set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\r
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\nwith open('\glove_vectors', '\wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'

```

In [37]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-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 [38]:

```

# 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:

```

```

    cnt_words -= 1
    vector /= cnt_words
    avg_w2v_vectors.append(vector)

```

```

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

```

100%|

109248/109248 [01:03<00:00, 1719.44it/s]

109248
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [39]:

```

# 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 [40]:

```

# 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]))

```

100%|

109248/109248 [06:48<00:00, 267.54it/s]

109248
300

In [41]:

```

# Similarly you can vectorize for title also

```

In [42]:

```

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

```

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_project_title = []; # the avg-w2v for each sentence/review is stored in this list

for sentence in tqdm(project_title): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_project_title.append(vector)

print(len(tfidf_w2v_vectors_project_title))
print(len(tfidf_w2v_vectors_project_title[0]))
```

109248
300

In [44]:

In [45]:

```
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.    ... 399.    287.
73    5.5 ].
# Reshape your data either using array.reshape(-1, 1) or array.reshape(1, -1)
```

```
price normalized
```

```
array([[0.00098843, 0.00191166, 0.00330448, ..., 0.00153418, 0.00046704,
        0.00070265]])
```

1.5.4 Merging all the above features

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

In [47]:

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_normalized.T.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 1)
```

In [48]:

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

Out[48]:

```
(109248, 16663)
```

In [49]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# 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
```

Assignment 4: Naive Bayes

1. Apply Multinomial NaiveBayes on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)

2. The hyper paramter tuning(find best Alpha)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- 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. Feature importance

- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using absolute values of `coef_` parameter of [MultinomialNB](#) and print their corresponding feature names

4. 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. Here on X-axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.
- 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](#).

5. 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](#)

2. Naive Bayes

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

In [50]:

```
# 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 [51]:

```
from sklearn.model_selection import train_test_split
X1_train, X_test_bow, y1_train, y_test_bow = train_test_split(
    project_data, y, test_size=0.20, stratify=y, random_state=42)
X_cv_bow, X_train_bow, y_cv_bow, y_train_bow = train_test_split(X1_train, y1_train, test_size=0.70, stratify=y, random_state=42)
```

In [52]:

```
X_train_bow.head(2)
```

Out[52]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime |
|-------|------------|---------|----------------------------------|----------------|--------------|----------------------------|
| 69975 | 54429 | p152820 | b6f1a59555245d7c795957340b5bfa43 | Ms. | LA | 2016-12-01 21:21:45 |
| 8122 | 114259 | p105475 | 7a732b77e49f18e4f1e1b99e2860f50a | Mrs. | NJ | 2017-01-13 07:54:28 |

2.2 Make Data Model Ready: encoding numerical, categorical features

In [101]:

```
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1
# normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html
from sklearn.preprocessing import normalize

# price_normalized = normalize(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287. 73 5.5 ].
# Reshape your data either using array.reshape(-1, 1) or array.reshape(1, -1)
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
```

```
# the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1
# normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html

# price_normalized = normalize(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287. 73 5.5 ].
# Reshape your data either using array.reshape(-1, 1) or array.reshape(1, -1)

price_normalized_train_bow = normalize(X_train_bow['price'].values.reshape(-1,1))
```

In [56]:

```
# Now standardize the data with above mean and variance.
price_normalized_cv_bow = normalize(X_cv_bow['price'].values.reshape(-1, 1))
```

In [57]:

```
price_normalized_test_bow = normalize(X_test_bow['price'].values.reshape(-1, 1))
```

In [58]:

```
#onehotencoding for school_state
one_hot_encoding_school_state_train_bow=pd.get_dummies(X_train_bow.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_train_bow.shape)
```

Shape of dataframe for school_state (61179, 51)

In [59]:

```
#onehotencoding for school_state
one_hot_encoding_school_state_cv_bow=pd.get_dummies(X_cv_bow.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_cv_bow.shape)
```

Shape of dataframe for school_state (26219, 51)

In [60]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix_train_bow=pd.get_dummies(X_train_bow.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_train_bow.shape)
```

Shape of dataframe for teacher_prefix (61179, 5)

In [61]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix_cv_bow=pd.get_dummies(X_cv_bow.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_cv_bow.shape)
```

Shape of dataframe for teacher_prefix (26219, 5)

In [62]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category_train_bow=pd.get_dummies(X_train_bow.project_grade_category)

print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_train_bow.shape)
```

Shape of dataframe for project_grade_category (61179, 4)

In [63]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category_cv_bow=pd.get_dummies(X_cv_bow.project_grade_category)

print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_cv_bow.shape)
```

Shape of dataframe for project_grade_category (26219, 4)

In [64]:

```
# 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_train_bow = vectorizer.fit_transform(X_train_bow['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_train_bow.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (61179, 9)
```

In [65]:

```
# we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True
)
categories_one_hot_cv_bow = vectorizer.transform(X_cv_bow['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_cv_bow.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (26219, 9)
```

In [66]:

```
categories_one_hot_test_bow = vectorizer.transform(X_test_bow['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_test_bow.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig (21850, 9)
```

In [67]:

```
# 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_train_bow =
vectorizer.fit_transform(X_train_bow['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_train_bow.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 (61179, 30)
```

In [68]:

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

```
# 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_cv_bow = vectorizer.transform(X_cv_bow['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ", sub_categories_one_hot_cv_bow.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 (26219, 30)
```

In [69]:

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

```
sub_categories_one_hot_test_bow = vectorizer.transform(X_test_bow['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ", sub_categories_one_hot_test_bow.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 (21850, 30)
```

In [70]:

```
#onehotencoding for school_state
```

```
one_hot_encoding_school_state_test_bow=pd.get_dummies(X_test_bow.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_test_bow.shape)
```

Shape of dataframe for school_state (21850, 51)

In [71]:

```
#onehotencoding for teacher_prefix
```

```
one_hot_encoding_teacher_prefix_test_bow=pd.get_dummies(X_test_bow.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_test_bow.shape)
```

Shape of dataframe for teacher_prefix (21850, 5)

In [72]:

```
#onehotencoding for project_grade_category
```

```
one_hot_encoding_project_grade_category_test_bow=pd.get_dummies(X_test_bow.project_grade_category)

print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_test_bow.shape)
```

Shape of dataframe for project_grade_category (21850, 4)

2.3 Make Data Model Ready: encoding essay, and project_title

In [73]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
```

```
vectorizer = CountVectorizer(min_df=10)
text_essay_train_bow = vectorizer.fit_transform(X_train_bow['preprocessed_essays'])
```

```
print("Shape of matrix after one hot encodig ",text_essay_train_bow.shape)
```

Shape of matrix after one hot encodig (61179, 13308)

In [74]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
text_essay_cv_bow = vectorizer.transform(X_cv_bow['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_essay_cv_bow.shape)
```

Shape of matrix after one hot encodig (26219, 13308)

In [75]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
text_bow_essay_test = vectorizer.transform(X_test_bow['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_bow_essay_test.shape)
```

Shape of matrix after one hot encodig (21850, 13308)

In [76]:

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

Shape of matrix after one hot encodig (61179, 375)

In [77]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
text_bow_project_title_cv= vectorizer.transform(X_cv_bow['project_title_list'])
print("Shape of matrix after one hot encodig ",text_bow_project_title_cv.shape)
```

Shape of matrix after one hot encodig (26219, 375)

In [78]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
text_bow_project_title_test = vectorizer.transform(X_test_bow['project_title_list'])
print("Shape of matrix after one hot encodig ",text_bow_project_title_test.shape)
```

Shape of matrix after one hot encodig (21850, 375)

In [79]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatinating a sparse matrix and a dense matirx :)
Data_model_ready_num_cat = hstack((categories_one_hot,
sub_categories_one_hot,one_hot_encoding_school_state,one_hot_encoding_teacher_prefix,one_hot_encodi
ng_project_grade_category,
price_normalized.T))
Data_model_ready_num_cat.shape
```

Out[79]:

(109248, 100)

2.4 Applying NB() on different kind of featurization as mentioned in the instructions

Apply Naive Bayes on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instructions

2.4.1 Applying Naive Bayes on BOW, SET 1

In [80]:

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

In [103]:

```
# 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 :)
bow_data_matrix_train=
hstack((price_normalized_train_bow,text_bow_project_title_train,text_essay_train_bow,sub_categories
_one_hot_train_bow,categories_one_hot_train_bow,one_hot_encoding_project_grade_category_train_bow,
one_hot_encoding_teacher_prefix_train_bow,one_hot_encoding_school_state_train_bow))
bow_data_matrix_train.shape
```

Out[103]:

```
(61179, 13783)
```

In [104]:

```
text_bow_project_title_train.shape
```

Out[104]:

```
(61179, 375)
```

In [105]:

```
# 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 :)
bow_data_matrix_cv=
hstack((one_hot_encoding_school_state_cv_bow,one_hot_encoding_teacher_prefix_cv_bow,one_hot_encoding
_project_grade_category_cv_bow,categories_one_hot_cv_bow,sub_categories_one_hot_cv_bow,price_norma
lized_cv_bow,text_essay_cv_bow,
text_bow_project_title_cv))
bow_data_matrix_cv.shape
```

Out[105]:

```
(26219, 13783)
```

In [106]:

```
# 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 :)
bow_data_matrix_test=
hstack((one_hot_encoding_school_state_test_bow,one_hot_encoding_teacher_prefix_test_bow,one_hot_enc
oding_project_grade_category_test_bow,categories_one_hot_test_bow,sub_categories_one_hot_test_bow,
price_normalized_test_bow,text_bow_essay_test,
text_bow_project_title_test))
bow_data_matrix_test.shape
```

Out[106]:

```
(21850, 13783)
```

In [124]:

```
from scipy.sparse import coo_matrix
m = coo_matrix(bow_data_matrix_train)
m1 = m.tocsr()
```

In [125]:

```
new_bow_data_matrix_train=m1[:60001]
```

In [126]:

```
new_y_train_bow=y_train_bow[:60001]
```

In [127]:

```
from scipy.sparse import coo_matrix
m2 = coo_matrix(bow_data_matrix_test)
m3 = m2.tocsr()
```

In [128]:

```
new_bow_data_matrix_test=m3[:20001]
```

In [129]:

```
new_y_test_bow=y_test_bow[:20001]
```

In [130]:

```
from scipy.sparse import coo_matrix
m4 = coo_matrix(bow_data_matrix_cv)
m5 = m4.tocsr()
```

In [131]:

```
new_bow_data_matrix_cv=m5[:20001]
```

In [132]:

```
new_y_cv_bow=y_cv_bow[:20001]
```

In [148]:

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

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

    return y_data_pred
```

In [137]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.naive_bayes import MultinomialNB
gnb_bow = MultinomialNB()
param_grid = {'alpha':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001]}
```

```

clf = GridSearchCV(gnb_bow, param_grid, cv=10, scoring='roc_auc')
clf.fit(new_bow_data_matrix_train, new_y_train_bow)

train_auc_bow= clf.cv_results_['mean_train_score']
train_auc_std_bow= clf.cv_results_['std_train_score']
cv_auc_bow = clf.cv_results_['mean_test_score']
cv_auc_std_bow= clf.cv_results_['std_test_score']

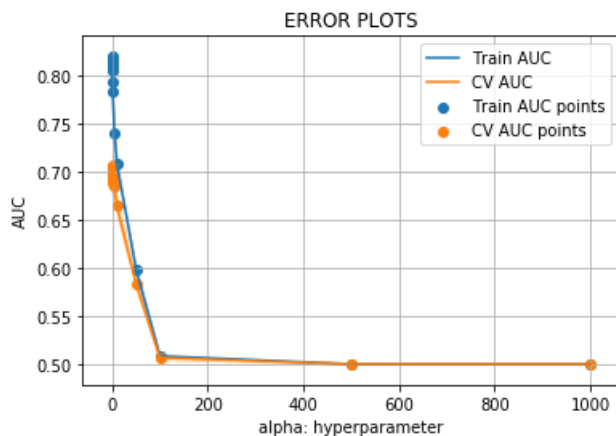
plt.plot(param_grid['alpha'], train_auc_bow, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(param_grid['alpha'], train_auc_bow - train_auc_std_bow, train_auc_bow +
train_auc_std_bow, alpha=0.2, color='darkblue')

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

plt.scatter(param_grid['alpha'], train_auc_bow, label='Train AUC points')
plt.scatter(param_grid['alpha'], cv_auc_bow, label='CV AUC points')

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

```



In [152]:

```

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

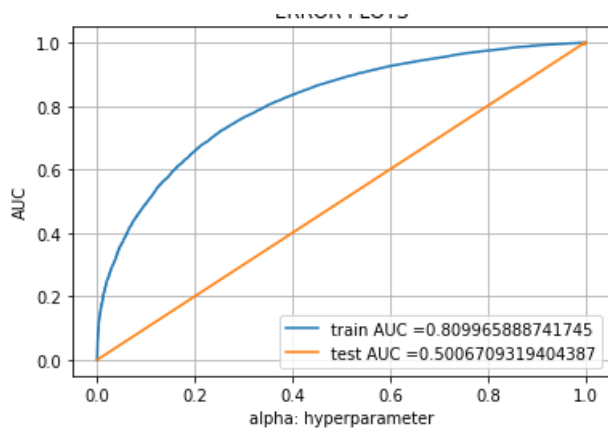
mnb = MultinomialNB(alpha=.0001)
mnb.fit(new_bow_data_matrix_train, new_y_train_bow)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

y_train_pred = batch_predict(mnb, new_bow_data_matrix_train)
y_test_pred = batch_predict(mnb, new_bow_data_matrix_test)

train_fpr_bow, train_tpr_bow, tr_thresholds_bow = roc_curve(new_y_train_bow, y_train_pred)
test_fpr_bow, test_tpr_bow, te_thresholds_bow = roc_curve(new_y_test_bow, y_test_pred)

plt.plot(train_fpr_bow, train_tpr_bow, label="train AUC =" + str(auc(train_fpr_bow, train_tpr_bow)))
plt.plot(test_fpr_bow, test_tpr_bow, label="test AUC =" + str(auc(test_fpr_bow, test_tpr_bow)))
plt.legend()
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```

In [156]:

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

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

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

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

In [158]:

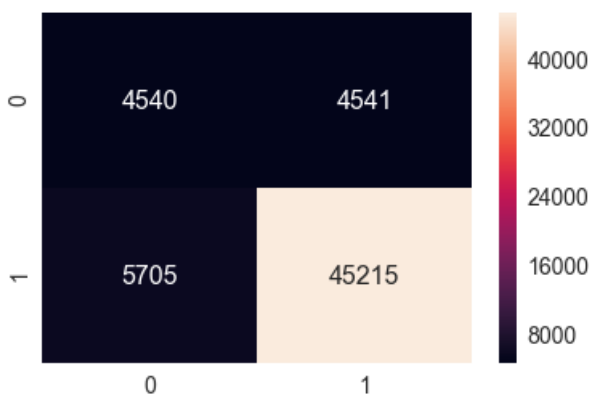
```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
df_cm=confusion_matrix(new_y_train_bow, predict(y_train_pred, tr_thresholds_bow, train_fpr_bow,
train_fpr_bow))

sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.2499999969683947 for threshold 0.154

Out[158]:

<matplotlib.axes._subplots.AxesSubplot at 0x2269e2d5080>



In [159]:

```
print("Test confusion matrix")

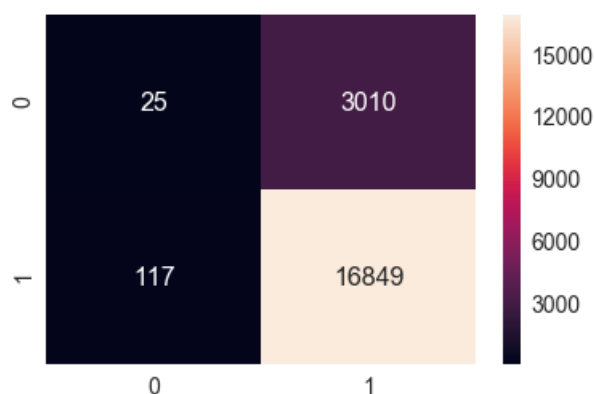
df_cm_test=confusion_matrix(new_y_test_bow, predict(y_test_pred, tr_thresholds_bow, test_fpr_bow, test_fpr_bow))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm_test, annot=True,annot_kws={"size": 16}, fmt='g')
```

Test confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.008169380294151943 for threshold 1.0

Out[159]:

<matplotlib.axes._subplots.AxesSubplot at 0x2269e2fe128>



TOP 10 important features from both Positive and Negative class from set 1

In [162]:

```
#Code Reference:https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-learn-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
    coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
    top = zip(coefs_with_fns[:n], coefs_with_fns[-(n + 1):-1])
    print("\t\t\tPositive\t\t\t\t\tNegative")

    print("_____")
    for (coef_1, fn_1), (coef_2, fn_2) in top:
        print("\t%.4f\t%-15s\t\t\t\t\t%.4f\t%-15s" % (coef_1, fn_1, coef_2, fn_2))

show_most_informative_features(vectorizer,mnb)
```

| | Positive | Negative |
|----------|----------------------------------|-----------------------------------|
| -14.2691 | build | -5.0406 16 |
| -14.2691 | chrome | -9.1124 life |
| -14.2691 | empoweringstudentsthroughart | -10.2402 jump |
| -14.2691 | welovetoread | -10.3064 flexibleseatingclassroom |
| -14.0868 | crazyforchromebooks | -10.4405 boom |
| -14.0868 | readingisfun | -10.6912 movement |
| -14.0868 | tfalldown | -10.7311 on |
| -13.9326 | creative | -10.7428 part3 |
| -13.9326 | flexibleseatingforactivelearners | -10.8351 fiction |
| -13.9326 | happy | -10.8547 wecan |

2.4.2 Applying Naive Bayes on TFIDF, SET 2

In [3]:

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

In [163]:

```
from sklearn.model_selection import train_test_split
X1_train, X_test_tfidf, y1_train, y_test_tfidf = train_test_split(
    project_data, y, test_size=0.20, stratify=y, random_state=42)
X_cv_tfidf, X_train_tfidf, y_cv_tfidf, y_train_tfidf = train_test_split(X1_train, y1_train, test_size=0.70,
    stratify=y1_train, random_state=42)
```

In [165]:

```
# check this one: https://www.youtube.com/watch?v=0H0qOcln3Z4&t=530s
# the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1
# normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html
from sklearn.preprocessing import normalize

# price_normalized = normalize(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5 ].
# Reshape your data either using array.reshape(-1, 1) or array.reshape(1, -1)

price_normalized_tfidf_train = normalize(X_train_tfidf['price'].values.reshape(-1, 1))
```

In [166]:

```
# Now standardize the data with above mean and variance.
price_normalized_cv_tfidf = normalize(X_cv_tfidf['price'].values.reshape(-1, 1))
```

In [167]:

```
# Now standardize the data with above mean and variance.
price_normalized_test_tfidf = normalize(X_test_tfidf['price'].values.reshape(-1, 1))
```

In [168]:

```
#onehotencoding for school_state
one_hot_encoding_school_state_train_tfidf=pd.get_dummies(X_train_tfidf.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_train_tfidf.shape)
```

Shape of dataframe for school_state (61179, 51)

In [169]:

```
#onehotencoding for school_state
one_hot_encoding_school_state_cv_tfidf=pd.get_dummies(X_cv_tfidf.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_cv_tfidf.shape)
```

Shape of dataframe for school_state (26219, 51)

In [170]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix_train_tfidf=pd.get_dummies(X_train_tfidf.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_train_tfidf.shape)
```

Shape of dataframe for teacher_prefix (61179, 5)

In [171]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix_cv_tfidf=pd.get_dummies(X_cv_tfidf.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_cv_tfidf.shape)
```

Shape of dataframe for teacher_prefix (26219, 5)

In [172]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category_train_tfidf=pd.get_dummies(X_train_tfidf.project_grade_category)
```

```
print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_train_tfidf.shape)
```

Shape of dataframe for project_grade_category (61179, 4)

In [173]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category_cv_tfidf=pd.get_dummies(X_cv_tfidf.project_grade_category)
```

```
print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_cv_tfidf.shape)
```

Shape of dataframe for project_grade_category (26219, 4)

In [174]:

```
# 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_train_tfidf = vectorizer.fit_transform(X_train_tfidf['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_train_tfidf.shape)
```

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

Shape of matrix after one hot encodig (61179, 9)

In [175]:

```
# we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
categories_one_hot_cv_tfidf = vectorizer.transform(X_cv_tfidf['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_cv_tfidf.shape)
```

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

Shape of matrix after one hot encodig (26219, 9)

In [176]:

```
categories_one_hot_test_tfidf = vectorizer.transform(X_test_tfidf['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_test_tfidf.shape)
```

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

Shape of matrix after one hot encodig (21850, 9)

In [177]:

```
# 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_train_tfidf = vectorizer.fit_transform(X_train_tfidf['clean_subcategories'])
```

```
.values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_train_tfidf.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 (61179, 30)
```

In [178]:

```
# 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_cv_tfidf = vectorizer.transform(X_cv_tfidf['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_cv_tfidf.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 (26219, 30)
```

In [179]:

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

sub_categories_one_hot_test_tfidf =
vectorizer.transform(X_test_tfidf['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_test_tfidf.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 (21850, 30)
```

In [180]:

```
#onehotencoding for school_state
one_hot_encoding_school_state_test_tfidf=pd.get_dummies(X_test_tfidf.school_state)
print("Shape of dataframe for school_state", one_hot_encoding_school_state_test_tfidf.shape)
```

Shape of dataframe for school_state (21850, 51)

In [181]:

```
#onehotencoding for teacher_prefix
one_hot_encoding_teacher_prefix_test_tfidf=pd.get_dummies(X_test_tfidf.teacher_prefix)

print("Shape of dataframe for teacher_prefix", one_hot_encoding_teacher_prefix_test_tfidf.shape)
```

Shape of dataframe for teacher_prefix (21850, 5)

In [182]:

```
#onehotencoding for project_grade_category
one_hot_encoding_project_grade_category_test_tfidf=pd.get_dummies(X_test_tfidf.project_grade_category)
```

```
print("Shape of dataframe for project_grade_category",
one_hot_encoding_project_grade_category_test_tfidf.shape)
```

Shape of dataframe for project_grade_category (21850, 4)

In [183]:

```
vectorizer = TfidfVectorizer(min_df=10)
tfidf_essay_train = vectorizer.fit_transform(X_train_tfidf['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",tfidf_essay_train.shape)
```

Shape of matrix after one hot encodig (61179, 13308)

In [184]:

```
tfidf_essay_cv = vectorizer.transform(X_cv_tfidf['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",tfidf_essay_cv.shape)
```

Shape of matrix after one hot encodig (26219, 13308)

In [185]:

```
tfidf_essay_test = vectorizer.transform(X_test_tfidf['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",tfidf_essay_test.shape)
```

Shape of matrix after one hot encodig (21850, 13308)

In [186]:

```
vectorizer = TfidfVectorizer(min_df=10)
tfidf_project_title_train = vectorizer.fit_transform(X_train_tfidf['project_title_list'])
print("Shape of matrix after one hot encodig ",tfidf_project_title_train.shape)
```

Shape of matrix after one hot encodig (61179, 375)

In [187]:

```
tfidf_project_title_cv = vectorizer.transform(X_cv_tfidf['project_title_list'])
print("Shape of matrix after one hot encodig ",tfidf_project_title_cv.shape)
```

Shape of matrix after one hot encodig (26219, 375)

In [188]:

```
tfidf_project_title_test = vectorizer.transform(X_test_tfidf['project_title_list'])
print("Shape of matrix after one hot encodig ",tfidf_project_title_test.shape)
```

Shape of matrix after one hot encodig (21850, 375)

In [191]:

```
# 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 :)
tfidf_data_matrix_train=
hstack((one_hot_encoding_school_state_train_tfidf,one_hot_encoding_teacher_prefix_train_tfidf,one_
hot_encoding_project_grade_category_train_tfidf,categories_one_hot_train_tfidf,sub_categories_one_h
ot_train_tfidf,price_normalized_tfidf_train,tfidf_essay_train,tfidf_project_title_train))
tfidf_data_matrix_train.shape
```

Out[191]:

(61179, 13783)

In [192]:

In [192]:

```
# 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 :)
tfidf_data_matrix_test=
hstack((one_hot_encoding_school_state_test_tfidf,one_hot_encoding_teacher_prefix_test_tfidf,one_hot_encoding_project_grade_category_test_tfidf,categories_one_hot_test_tfidf,sub_categories_one_hot_test_tfidf,price_normalized_test_tfidf,tfidf_essay_test,tfidf_project_title_test))
tfidf_data_matrix_test.shape
```

Out[192]:

(21850, 13783)

In [193]:

```
# 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 :)
tfidf_data_matrix_cv=
hstack((one_hot_encoding_school_state_cv_tfidf,one_hot_encoding_teacher_prefix_cv_tfidf,one_hot_encoding_project_grade_category_cv_tfidf,categories_one_hot_cv_tfidf,sub_categories_one_hot_cv_tfidf,price_normalized_cv_tfidf,tfidf_essay_cv,tfidf_project_title_cv))
tfidf_data_matrix_cv.shape
```

Out[193]:

(26219, 13783)

In [194]:

```
from scipy.sparse import coo_matrix
m = coo_matrix(tfidf_data_matrix_train)
m1 = m.tocsr()
```

In [195]:

```
new_tfidf_data_matrix_train=m1[:60001]
```

In [196]:

```
new_y_train_tfidf=y_train_tfidf[:60001]
```

In [197]:

```
from scipy.sparse import coo_matrix
m2 = coo_matrix(tfidf_data_matrix_test)
m3 = m2.tocsr()
```

In [198]:

```
new_tfidf_data_matrix_test=m3[:20001]
```

In [199]:

```
new_y_test_tfidf=y_test_tfidf[:20001]
```

In [200]:

```
new_y_test_tfidf.shape
```

Out[200]:

(20001,)

In [201]:

```
from scipy.sparse import coo_matrix
m4 = coo_matrix(tfidf_data_matrix_cv)
m5 = m4.tocsr()
```

In [211]:

```
new_tfidf_data_matrix_cv=m5[:20001]
```

In [212]:

```
new_y_cv_tfidf=y_cv_tfidf[:20001]
```

In [213]:

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

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

    return y_data_pred
```

In [205]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.naive_bayes import MultinomialNB
gnb_tfidf = MultinomialNB()
param_grid = {'alpha':[1000,500,100,50,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001,0.0005,0.0001]}

clf = GridSearchCV(gnb_tfidf, param_grid, cv=10, scoring='roc_auc')
clf.fit(new_tfidf_data_matrix_train,new_y_train_tfidf)

train_auc_tfidf= clf.cv_results_['mean_train_score']
train_auc_std_tfidf= clf.cv_results_['std_train_score']
cv_auc_tfidf = clf.cv_results_['mean_test_score']
cv_auc_std_tfidf= clf.cv_results_['std_test_score']

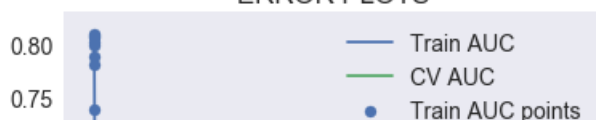
plt.plot(param_grid['alpha'], train_auc_tfidf, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(param_grid['alpha'],train_auc_tfidf - train_auc_std_tfidf,train_auc_tfidf +
train_auc_std_tfidf,alpha=0.2,color='darkblue')

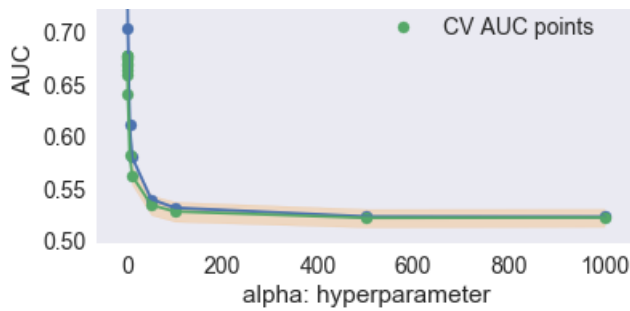
plt.plot(param_grid['alpha'], cv_auc_tfidf, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(param_grid['alpha'],cv_auc_tfidf - cv_auc_std_tfidf,cv_auc_tfidf +
cv_auc_std_tfidf,alpha=0.2,color='darkorange')

plt.scatter(param_grid['alpha'], train_auc_tfidf, label='Train AUC points')
plt.scatter(param_grid['alpha'], cv_auc_tfidf, label='CV AUC points')

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

ERROR PLOTS





In [206]:

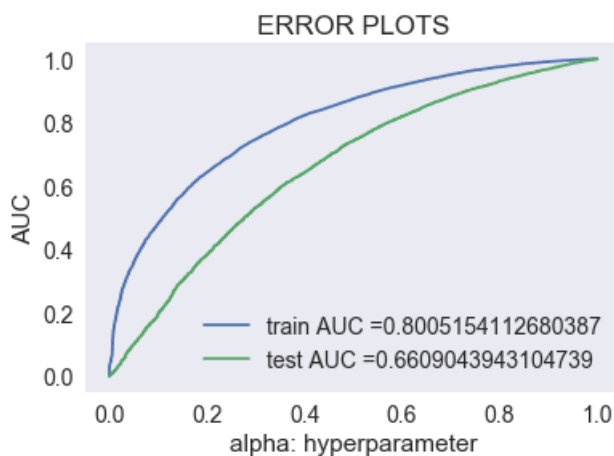
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

mnb = MultinomialNB(alpha=.0001)
mnb.fit(new_tfidf_data_matrix_train, new_y_train_tfidf)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(mnb, new_tfidf_data_matrix_train)
y_test_pred = batch_predict(mnb, new_tfidf_data_matrix_test)

train_fpr_tfidf, train_tpr_tfidf, tr_thresholds_tfidf = roc_curve(new_y_train_tfidf, y_train_pred)
test_fpr_tfidf, test_tpr_tfidf, te_thresholds_tfidf = roc_curve(new_y_test_tfidf, y_test_pred)

plt.plot(train_fpr_tfidf, train_tpr_tfidf, label="train AUC =" + str(auc(train_fpr_tfidf, train_tpr_tfidf)))
plt.plot(test_fpr_tfidf, test_tpr_tfidf, label="test AUC =" + str(auc(test_fpr_tfidf, test_tpr_tfidf)))
plt.legend()
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [207]:

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

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

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

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

```

        predictions.append(1)
    else:
        predictions.append(0)
    return predictions

```

In [215]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
df_cm=confusion_matrix(new_y_train_tfidf, predict(y_train_pred, tr_thresholds_tfidf,
train_fpr_tfidf, train_fpr_tfidf))

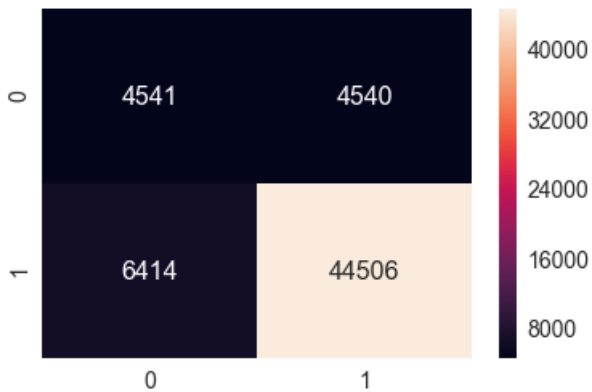
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

```

Train confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.24999999696839467 for threshold 0.785

Out[215]:

<matplotlib.axes._subplots.AxesSubplot at 0x2268f9baa90>



In [217]:

```

print("Test confusion matrix")

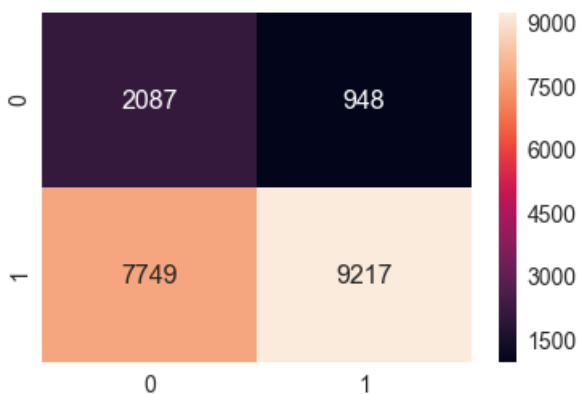
df_cm_test=confusion_matrix(new_y_test_tfidf, predict(y_test_pred, tr_thresholds_tfidf, test_fpr_tfidf,
test_fpr_tfidf))
sns.set(font_scale=1.4)#for label size
sns.heatmap(df_cm_test, annot=True,annot_kws={"size": 16}, fmt='g')

```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.2499999728592017 for threshold 0.886

Out[217]:

<matplotlib.axes._subplots.AxesSubplot at 0x2268fa0d5f8>



TOP 10 important features from both Positive and Negative class from set 2

In [218]:

```
#Code Reference:https://stackoverflow.com/questions/11116697/how-to-get-most-informative-features-for-scikit-learn-classifiers
def show_most_informative_features(vectorizer, clf, n=10):
    feature_names = vectorizer.get_feature_names()
    coefs_with_fns = sorted(zip(clf.coef_[0], feature_names))
    top = zip(coefs_with_fns[:n], coefs_with_fns[-(n + 1):-1])
    print("\t\t\tPositive\t\t\t\t\tNegative")

print("_____")
for (coef_1, fn_1), (coef_2, fn_2) in top:
    print("\t%.4f\t%-15s\t\t\t\t\t%.4f\t%-15s" % (coef_1, fn_1, coef_2, fn_2))

show_most_informative_features(vectorizer,mnb)
```

| Positive | Negative |
|-----------------------------|----------------------------------|
| -13.8906 one | -2.7759 flexibleseatingclassroom |
| -13.5598 readingforsuccess | -3.4089 chromebooks |
| -13.5426 moreyouread | -3.4994 collaborate |
| -13.5056 kid | -3.6779 class |
| -13.4754 sims | -3.7561 coding |
| -13.4166 go | -3.8192 chromebooksforall |
| -13.3838 part3 | -3.8518 chromebooksforlearning |
| -13.3835 timeforkids | -3.9350 flexibleseating |
| -13.3646 wigglewhilewelearn | -4.1400 flexiblelearning |
| -13.3255 ms | -4.3588 flexibleclassroomseating |

3. Conclusions

In [220]:

```
# comparing all models using Prettytable library
```

In [219]:

```
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Featurization", "train_auc", "test_auc", "threshold_for train", "tpr*(1-fpr) for train", "threshold_for_test", "tpr*(1-fpr) for test" ,]
x.add_row(["BOW",0.8005,0.6609,0.785,0.2499,0.886,0.2499 ])
x.add_row(["TFIDF",0.8099,0.5006,0.154,0.2499,1,0.0081 ])

print(x)
```

| Featurization | train_auc | test_auc | threshold_for train | tpr*(1-fpr) for train | threshold_for_test | tpr*(1-fpr) for test |
|---------------|-----------|----------|---------------------|-----------------------|--------------------|----------------------|
| BOW | 0.8005 | 0.6609 | 0.785 | 0.2499 | 0.886 | 0.2499 |
| TFIDF | 0.8099 | 0.5006 | 0.154 | 0.2499 | 1 | 0.0081 |

