

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]:

cat_list = []

	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
    sub_cat_list.append(temp.strip())
```

```

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

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

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

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

```

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.

nannan

=====

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.

nannan

=====

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.

nannan

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

In [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())
```



```
[02:07<00:00, 859.30it/s]
```

```
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\nle\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
[01:00<00:00, 1806.88it/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)

In [0]:

```
# 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 [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 [3]:

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

In [4]:

```
# 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 [5]:

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

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

In [6]:

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

In [7]:

```
# 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(project_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 [8]:

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

Out[8]:

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

1. class 1 (Approved) is Majority
2. class 0 (Not Approved) is Minority

In [9]:

```
# upsample the minority class; https://elitedatascience.com/imbalanced-classes
from sklearn.utils import resample, shuffle

X_train['output'] = y_train

# splitting train data frame by classes
train_minority = X_train[X_train.output == 0]
train_majority = X_train[X_train.output == 1]
```

```

train_minority_upsampled = resample(train_minority, replace = True,
                                     n_samples = train_majority.shape[0],
                                     random_state = 0)
print('Upsampled minority shape : ', train_minority_upsampled.shape)

# concatenate upsampled minority data and majority data
X_train = pd.concat([train_majority, train_minority_upsampled])
# shuffle dataframe pandas; https://stackoverflow.com/a/39673666
X_train = shuffle(X_train)
y_train = X_train.output.values
X_train.drop(columns = 'output', axis = 1, inplace = True)

```

Upsampled minority shape : (43479, 19)

In [11]:

```

print('X_train shape : ', X_train.shape)
print('y_train shape : ', y_train.shape)

```

X_train shape : (86958, 18)
y_train shape : (86958,)

2.2 Make Data Model Ready: encoding numerical, categorical features

In [12]:

```

# 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 Endcoding Numerical features

In [13]:

```

from sklearn.preprocessing import Normalizer

```

2.2.1.1 quantity

In [14]:

```

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

```

In [15]:

```

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

```

In [16]:

```

# 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 [17]:

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

In [18]:

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

In [19]:

```
# 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 [20]:

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

In [21]:

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

In [22]:

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

2.2.2 Preprocessing and Encoding Categorical Data

In [23]:

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

2.2.2.1 project_subject_category

In [24]:

```
# 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 [25]:

```
# 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 [26]:

```
# 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 [27]:

```
# 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 [28]:

```
# 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 :  (86958, 9)
```

In [29]:

```
# 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 [30]:

```
# 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 [31]:

```
# replace & with _; remove spaces\n",
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 [32]:

```
# 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 [33]:

```
# 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 [34]:

```
# 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 [35]:

```
# encoding train data  
#count = Counter()  
#for s in X_train.clean_subcategory.values:  
#    count.update(s.split())  
  
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 : (86958, 30)

In [36]:

```
# 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 [37]:

```
# 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 [38]:

```
# 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 [39]:

```
# 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 [40]:

```
# 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 [41]:

```
# 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 [42]:

```
# 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 :  (86958, 4)
```

In [43]:

```
# 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 [44]:

```
# 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 [45]:

```
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 [46]:

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

In [47]:

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

In [48]:

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

In [49]:

```
# 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 :  (86958, 6)
```

In [50]:

```
# 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 [51]:

```
# 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_states

In [52]:

```
# 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 :  (86958, 51)
```

In [53]:

```
# 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 [54]:

```
# 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 [55]:

```
# 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 [56]:

```
# 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 [57]:

```
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 [58]:

```
# 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 [59]:

```
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)
```

In [60]:

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

```
100%|██████████████████████████████████████████████████████████████████████████████| 86958/86958  
[01:19<00:00, 1099.65it/s]
```

In [61]:

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

```
100%|██████████████████████████████████████████████████████████████████████████| 21959/21959  
[00:20<00:00, 1084.18it/s]
```

In [62]:

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

```
100%|██████████████████████████████████████████████████████████████████████████████| 36052/36052  
[00:33<00:00, 1082.6lit/s]
```

2.3.1.2 preprocessing title

In [63]:

```
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 [64]:

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

```
100%|██████████████████████████████████████████████████████████████████████████| 86958/86958  
[00:02<00:00, 32359.18it/s]
```

In [65]:

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

```
100%|██████████████████████████████████████████████████████████████████████████| 21959/21959  
[00:00<00:00, 32218.42it/s]
```

In [66]:

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

```
100%|██████████████████████████████████████████████████████████████████████████████| 36052/36052  
[00:01<00:00, 32294.11it/s]
```

2.3.2 Vectorizing text data (train)

2.3.2.1 Bag Of Words (BOW)

In [67]:

```
# preprocessed_essay train data

# Considering the words which appeared in at least min_df documents(rows or projects).
# using n_gram = 1
essayBowVectorizer = CountVectorizer(min_df = 20, ngram_range = (1, 1))
essayBow = essayBowVectorizer.fit_transform(X_train.preprocessed_essay)
print("Shape of matrix after BOW ", essayBow.shape)
```


Shape of matrix after BOW (86958, 11332)

In [68]:

```
# 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, 11332)

In [69]:

```
# 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, 11332)

In [70]:

```
# preprocessed_title train data

# Considering the words which appeared in at least min_df documents (rows or projects).
# using n_gram = 1
title_bow_vectorizer = CountVectorizer(min_df = 20, ngram_range = (1, 1))
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 (86958, 1826)

In [71]:

```
# 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, 1826)

In [72]:

```
# test_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, 1826)

2.3.2.2 Term Frequency Inverse Document Frequency (TFIDF)

In [73]:

```
# 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 = 20, ngram_range = (1, 1))
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 (86958, 11332)

In [74]:

```
# preprocessed_essay cross validation data
cv_essay_tfidf = essay_tfidf_vectorizer.transform(X_cv.preprocessed_essay)
```

```
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, 11332)

In [75]:

```
# 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, 11332)

In [76]:

```
# 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 = 20, ngram_range = (1, 1))
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 (86958, 1826)

In [77]:

```
# 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, 1826)

In [78]:

```
# 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, 1826)

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

In [79]:

```
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score, roc_curve
```

In [80]:

```
# print 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
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 Naive Bayes on BOW, SET 1

In [81]:

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

In [82]:

```

# 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 : (86958, 13261)
y_train shape : (86958,)

```

In [83]:

```

# 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))

```

```
ow)).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, 13261)
y_cv shape : (21959,)
```

In [84]:

```
# 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_bow, te
                 st_title_bow)).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, 13261)
y_test shape : (36052,)
```

In [85]:

```
# Getting features name (BOW)
features_bow = list() # stores the column names of train/test/cv data
features_bow.extend(state_vectorizer.get_feature_names())
features_bow.extend(cat_vectorizer.get_feature_names())
features_bow.extend(subcat_vectorizer.get_feature_names())
features_bow.extend(grade_vectorizer.get_feature_names())
features_bow.extend(teacher_vectorizer.get_feature_names())
features_bow.extend(['quantity', 'prev_posted_project', 'price'])
features_bow.extend(essay_bow_vectorizer.get_feature_names())
features_bow.extend(title_bow_vectorizer.get_feature_names())
print('No. of features : ', len(features_bow))
```

```
No. of features : 13261
```

In [86]:

```
def batch_predict(model, X, batch_size = 1000):
    ''' Predict the log probability score in the batches
        and return all the log probability scores in a single array'''
    predicted = list()
    for i in range(0, X.shape[0] - batch_size + 1, batch_size):
        pred = model.predict_log_proba(X[i : i + batch_size])[:, 1]
        predicted.extend(pred)
    if X.shape[0] % batch_size != 0:
        predicted.extend(model.predict_log_proba(X[X.shape[0] - X.shape[0] % batch_size :])[:, 1])
    return np.array(predicted)
```

In [87]:

```
# training the NaiveBayes model bow
alpha_range = list(map(lambda x : 10 ** x, [-4, -3, -2, -1, 0, 1, 2, 3, 4]))
cv_auc_scores = dict() # stores the AUC scores of cross validation data
train_auc_scores = dict() # stores the AUC scores of train data

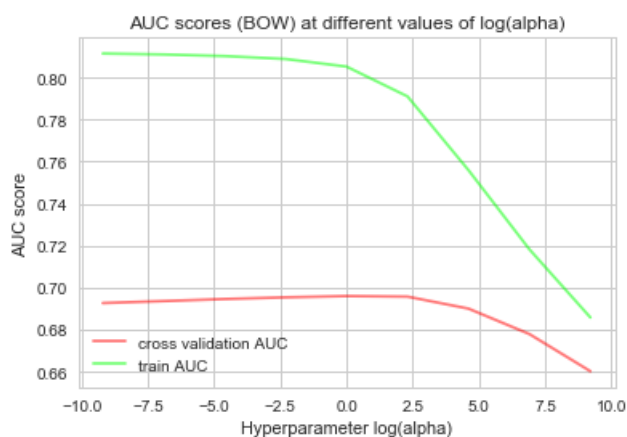
for alpha in tqdm(alpha_range):
    model = MultinomialNB(alpha = alpha)
    # fitting the train data
    model.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(model, X_train, batch_size = 400)
    predicted_cv = batch_predict(model, X_cv, batch_size = 400)
```

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

100% | 9/9 [00:02<00:00, 3.83it/s]

In [88]:

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



In [89]:

```
best_alpha = alpha_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best MultinomialNaiveBayes model (BOW)
best_model = MultinomialNB(alpha = best_alpha)
# fitting the train data
best_model.fit(X_train, y_train)

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

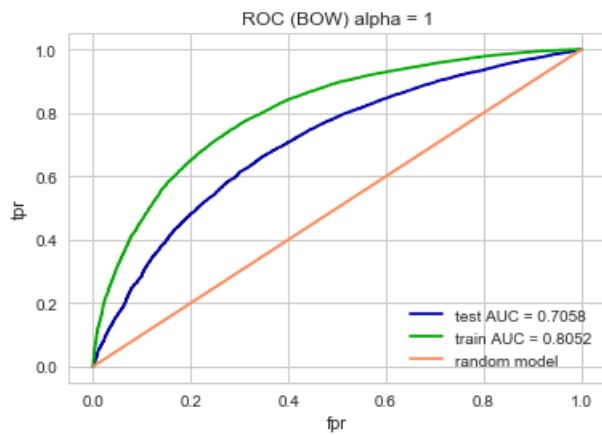
# calculates fpr and tpr
test_fpr, test_tpr, test_th = roc_curve(y_test, predicted_test)
train_fpr, train_tpr, train_th = roc_curve(y_train, predicted_train)

# adding AUC score to summary table
summary_table.add_row(['BOW', f'{best_alpha}', '%.4f' % (roc_auc_score(y_train, predicted_train)),
                      '%.4f' % (roc_auc_score(y_test, predicted_test))])
```

In [90]:

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

```
plt.ylabel('tpr')
plt.title(f'ROC (BOW) alpha = {best_alpha}')
plt.show()
```



In [91]:

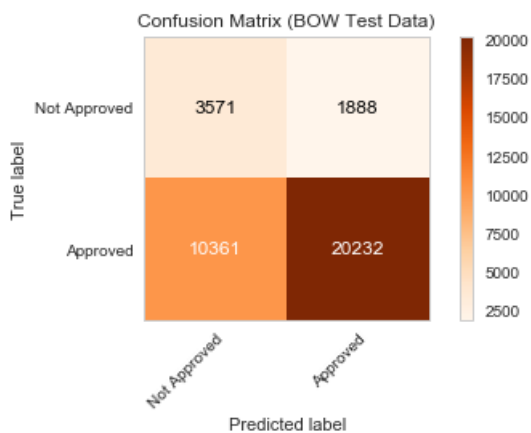
```
def predict_class(proba, tpr, fpr, threshold):
    """Predict the class label"""
    t = threshold[np.argmax(tpr * (1 - fpr))]
    prediction = list()
    for i in proba:
        if i >= t:
            prediction.append(1)
        else:
            prediction.append(0)
    return prediction
```

In [92]:

```
pred_test = predict_class(predicted_test, test_tpr, test_fpr, test_th)

# 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)
```

Confusion matrix
[[3571 1888]
[10361 20232]]



Out [92]:

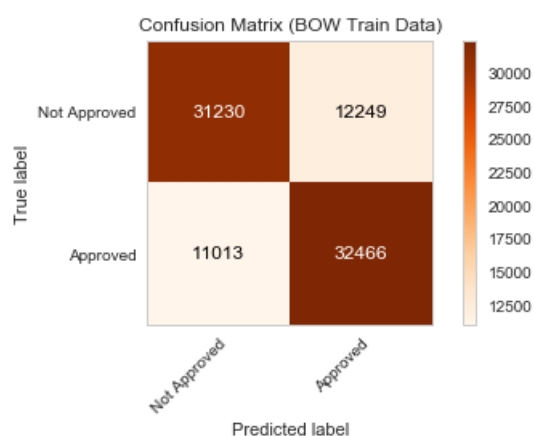
<matplotlib.axes._subplots.AxesSubplot at 0x27496b97b70>

In [93]:

```
pred_train = predict_class(predicted_train, train_tpr, train_fpr, train_th)
```

```
# 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)
```

```
Confusion matrix
[[31230 12249]
 [11013 32466]]
```



```
Out[93]:
<matplotlib.axes._subplots.AxesSubplot at 0x2749935b208>
```

2.4.1.1 Top 10 important features of positive class from SET 1

```
In [94]:
```

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

```
In [95]:
```

```
# top 10 important features using naive bayes; https://stackoverflow.com/a/50530697
indices = (best_model.feature_log_prob_[1].argsort()[::-1][:10])
print(np.array(features_bow)[indices]) # print the top 10 column(feature) names of positive class
```

```
['students' 'school' 'learning' 'classroom' 'not' 'learn' 'help'
 'quantity' 'price' 'many']
```

2.4.1.2 Top 10 important features of negative class from SET 1

```
In [96]:
```

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

```
In [97]:
```

```
# top 10 important features using naive bayes; https://stackoverflow.com/a/50530697
indices = best_model.feature_log_prob_[0].argsort()[::-1][:10]
print(np.array(features_bow)[indices]) # print the top 10 column(feature) names of negative class
```

```
['students' 'school' 'learning' 'classroom' 'not' 'learn' 'help'
 'quantity' 'price' 'nannan']
```

2.4.2 Applying Naive Bayes on TFIDF, SET 2

```
In [98]:
```

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

In [99]:

```
# 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 : (86958, 13261)
y_train shape : (86958,)
```

In [100]:

```
# 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, 13261)
y_cv shape : (21959,)
```

In [101]:

```
# 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, test_title_tfidf)).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, 13261)
y_test shape : (36052,)
```

In [102]:

```
# Getting features name (TFIDF)
features_tfidf = list() # stores the column names of train/test/cv data
features_tfidf.extend(state_vectorizer.get_feature_names())
features_tfidf.extend(cat_vectorizer.get_feature_names())
features_tfidf.extend(subcat_vectorizer.get_feature_names())
features_tfidf.extend(grade_vectorizer.get_feature_names())
features_tfidf.extend(teacher_vectorizer.get_feature_names())
features_tfidf.extend(['quantity', 'prev_posted_project', 'price'])
features_tfidf.extend(essay_tfidf_vectorizer.get_feature_names())
features_tfidf.extend(title_tfidf_vectorizer.get_feature_names())
print('No. of features : ', len(features_tfidf))
```


No. of features : 13261

In [103]:

```
# training the MultinomialNaiveBayes model tfidf
alpha_range = list(map(lambda x : 10 ** x, list(np.arange(-4, 4, 1.0))))
cv_auc_scores = dict()
train_auc_scores = dict()

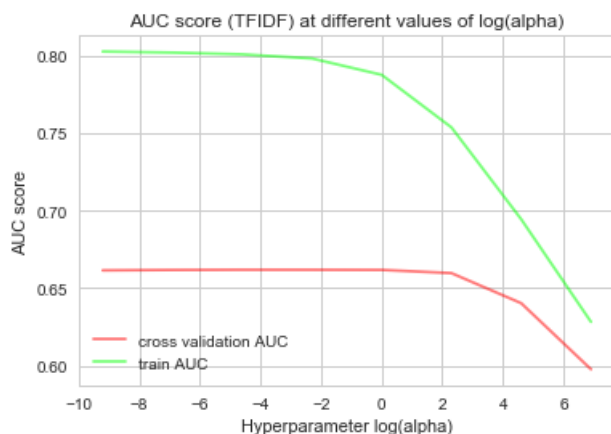
for alpha in tqdm(alpha_range):
    model = MultinomialNB(alpha = alpha)
    # fitting the train data
    model.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(model, X_train, batch_size = 400)
    predicted_cv = batch_predict(model, X_cv, batch_size = 400)

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

[illegible]

In [104]:

```
# 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.legend()
plt.xlabel('Hyperparameter log(alpha)')
plt.ylabel('AUC score')
plt.title('AUC score (TFIDF) at different values of log(alpha)')
plt.show()
```



In [105]:

```
best_alpha = alpha_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best MultinomialNaiveBayes model (TFIDF)
best_model = MultinomialNB(alpha = best_alpha)
# fitting the train data
best_model.fit(X_train, y_train)

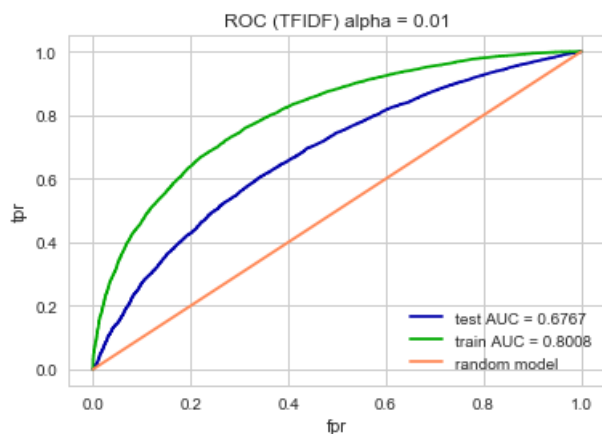
# predicting the probability scores of train data, test data
predicted_test = batch_predict(best_model, X_test, 400)
predicted_train = batch_predict(best_model, X_train, 400)

# 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)
```

```
# adding AUC score to summary table
summary_table.add_row(['TFIDF', f'{best_alpha}', '%.4f' % (roc_auc_score(y_train, predicted_train))
,
                        '%.4f' % (roc_auc_score(y_test, predicted_test))])
```

In [106]:

```
# 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) alpha = {best_alpha}')
plt.show()
```

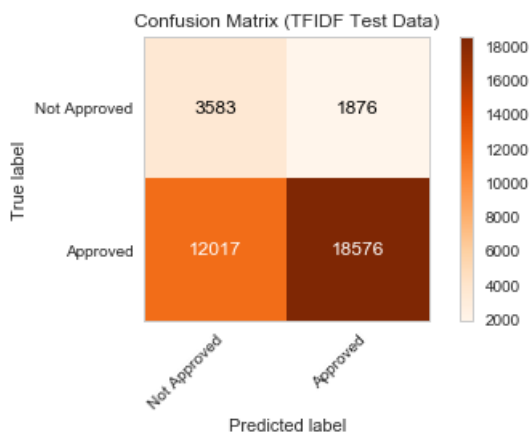


In [107]:

```
pred_test = predict_class(predicted_test, test_tpr, test_fpr, test_th)

# 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)
```

Confusion matrix
[[3583 1876]
[12017 18576]]



Out[107]:

<matplotlib.axes._subplots.AxesSubplot at 0x274993d9898>

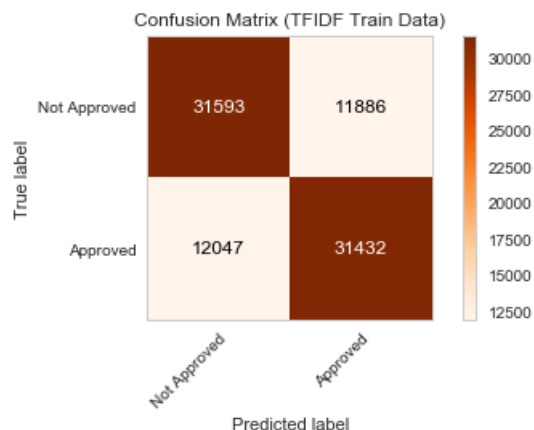
In [108]:

```
In [100]:
```

```
pred_train = predict_class(predicted_train, train_tpr, train_fpr, train_th)

# 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)
```

```
Confusion matrix
[[31593 11886]
 [12047 31432]]
```



```
Out[108]:
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x274933a6320>
```

2.4.2.1 Top 10 important features of positive class from SET 2

```
In [109]:
```

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

```
In [110]:
```

```
# top 10 important features using naive bayes; https://stackoverflow.com/a/50530697
indices = best_model.feature_log_prob_[1].argsort()[::-1][:10]
print(np.array(features_tfidf)[indices]) # print the top 10 column(feature) names of positive
class
```

```
['quantity' 'price' 'prev_posted_project' 'Mrs' 'Literacy_Language'
 'Grades_PreK_2' 'Math_Science' 'Ms' 'Grades_3_5' 'Literacy']
```

2.4.2.2 Top 10 important features of negative class from SET 2

```
In [111]:
```

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

```
In [112]:
```

```
# top 10 important features using naive bayes; https://stackoverflow.com/a/50530697
indices = best_model.feature_log_prob_[0].argsort()[::-1][:10]
print(np.array(features_tfidf)[indices]) # print the top 10 column(feature) names of negative
class
```

```
['price' 'quantity' 'prev_posted_project' 'Mrs' 'Literacy_Language'
 'Grades_PreK_2' 'Math_Science' 'Ms' 'Grades_3_5' 'Mathematics']
```

3. Conclusions

In [113]:

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

In [114]:

```
print(summary_table)
```

Vectorizer	Hyper parameter "alpha"	AUC on Train Data	AUC on Test Data
BOW	1	0.8052	0.7058
TFIDF	0.01	0.8008	0.6767

1. Multinomial Naive Bayes model by using Bag Of Word for text encoding gives better AUC score on both train and test data
2. Hence for Donors Choose dataset BOW word representation should be used