

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 [27]:

```

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

```

In [28]:

```

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 [29]:

```

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

```

```

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

```

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

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

Out[29]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_
55660	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	
76127	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	

In [30]:

```
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[30]:

	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 [31]:

```
catogories = 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 catogories:
    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 catogory 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 placeing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_') # 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 [32]:

```
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 placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " #" abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_')
        sub_cat_list.append(temp.strip())

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

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

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

1.3 Text preprocessing

In [33]:

```
# 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 [34]:

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

Out[34]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
55660	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	Engineer STEAM the Print Classroom

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	04-27	Grades 3-5
						00:31:25	Tools
							Fo

In [35]:

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

In [36]:

```
# 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)
```

I have been fortunate enough to use the Fairy Tale STEM kits in my classroom as well as the STEM journals, which my students really enjoyed. I would love to implement more of the Lakeshore STEM kits in my classroom for the next school year as they provide excellent and engaging STEM lessons. My students come from a variety of backgrounds, including language and socioeconomic status. Many of them don't have a lot of experience in science and engineering and these kits give me the materials to provide these exciting opportunities for my students. Each month I try to do several science or STEM/STEAM projects. I would use the kits and robot to help guide my science instruction in engaging and meaningful ways. I can adapt the kits to my current language arts pacing guide where we already teach some of the material in the kits like tall tales (Paul Bunyan) or Johnny Appleseed. The following units will be taught in the next school year where I will implement these kits: magnets, motion, sink vs. float, robots. I often get to these units and don't know if I am teaching the right way or using the right materials. The kits will give me additional ideas, strategies, and lessons to prepare my students in science. It is challenging to develop high quality science activities. These kits give me the materials I need to provide my students with science activities that will go along with the curriculum in my classroom. Although I have some things (like magnets) in my classroom, I don't know how to use them effectively. The kits will provide me with the right amount of materials and show me how to use them in an appropriate way.

I teach high school English to students with learning and behavioral disabilities. My students all vary in their ability level. However, the ultimate goal is to increase all students literacy levels. This includes their reading, writing, and communication levels. I teach a really dynamic group of students. However, my students face a lot of challenges. My students all live in poverty and in a dangerous neighborhood. Despite these challenges, I have students who have the desire to defeat these challenges. My students all have learning disabilities and currently all are performing below grade level. My students are visual learners and will benefit from a classroom that fulfills their preferred learning style. The materials I am requesting will allow my students to be prepared for the classroom with the necessary supplies. Too often I am challenged with students who come to school unprepared for class due to economic challenges. I want my students to be able to focus on learning and not how they will be able to get school supplies. The supplies will last all year. Students will be able to complete written assignments and maintain a classroom journal. The chart paper will be used to make learning more visual in class and to create posters to aid students in their learning. The students have access to a classroom printer. The toner will be used to print student work that is completed on the classroom Chromebooks. I want to try and remove all barriers for the students learning and create opportunities for learning. One of the biggest barriers is the students not having the resources to get pens, paper, and folders. My students will be able to increase their literacy skills because of this project.

"Life moves pretty fast. If you don't stop and look around once in awhile, you could miss it." from the movie, Ferris Bueller's Day Off. Think back...what do you remember about your grandparents? How amazing would it be to be able to flip through a book to see a day in their lives? My second graders are voracious readers! They love to read both fiction and nonfiction books. Their favorite characters include Pete the Cat, Fly Guy, Piggie and Elephant, and Mercy Watson. They also love to read about insects, space and plants. My students are hungry bookworms! My students are eager to learn and read about the world around them. My kids love to be at school and are like little sponges absorbing everything around them. Their parents work long hours and usually do not see their children. My students are usually cared for by their grandparents or a family friend. Most of my students do not have someone who speaks English at home. Thus it is difficult for my students to acquire language. Now think forward... wouldn't it mean a lot to your kids, nieces or nephews or grandchildren, to be able to see a day in your life today 30 years from now? Memories are so precious to us and being able to share these memories with future generations will

Memories are so precious to us and being able to share these memories with future generations will be a rewarding experience. As part of our social studies curriculum, students will be learning about changes over time. Students will be studying photos to learn about how their community has changed over time. In particular, we will look at photos to study how the land, buildings, clothing, and schools have changed over time. As a culminating activity, my students will capture a slice of their history and preserve it through scrap booking. Key important events in their young lives will be documented with the date, location, and names. Students will be using photos from home and from school to create their second grade memories. Their scrap books will preserve their unique stories for future generations to enjoy. Your donation to this project will provide my second graders with an opportunity to learn about social studies in a fun and creative manner. Through their scrapbooks, children will share their story with others and have a historical document for the rest of their lives.

=====

"A person's a person, no matter how small." (Dr. Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. \r\nStudents in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans.\r\nOur school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, "\"Can we try cooking with REAL food?\" I will take their idea and create \"Common Core Cooking Lessons\" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it's healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. \r\nStudents will gain math and literature skills as well as a life long enjoyment for healthy cooking.annan

=====

My classroom consists of twenty-two amazing sixth graders from different cultures and backgrounds. They are a social bunch who enjoy working in partners and working with groups. They are hard-working and eager to head to middle school next year. My job is to get them ready to make this transition and make it as smooth as possible. In order to do this, my students need to come to school every day and feel safe and ready to learn. Because they are getting ready to head to middle school, I give them lots of choice- choice on where to sit and work, the order to complete assignments, choice of projects, etc. Part of the students feeling safe is the ability for them to come into a welcoming, encouraging environment. My room is colorful and the atmosphere is casual. I want them to take ownership of the classroom because we ALL share it together. Because my time with them is limited, I want to ensure they get the most of this time and enjoy it to the best of their abilities. Currently, we have twenty-two desks of differing sizes, yet the desks are similar to the ones the students will use in middle school. We also have a kidney table with crates for seating. I allow my students to choose their own spots while they are working independently or in groups. More often than not, most of them move out of their desks and onto the crates. Believe it or not, this has proven to be more successful than making them stay at their desks! It is because of this that I am looking toward the "Flexible Seating" option for my classroom.\r\nThe students look forward to their work time so they can move around the room. I would like to get rid of the constricting desks and move toward more "fun" seating options. I am requesting various seating so my students have more options to sit. Currently, I have a stool and a papasan chair I inherited from the previous sixth-grade teacher as well as five milk crate seats I made, but I would like to give them more options and reduce the competition for the "good seats". I am also requesting two rugs as not only more seating options but to make the classroom more welcoming and appealing. In order for my students to be able to write and complete work without desks, I am requesting a class set of clipboards. Finally, due to curriculum that requires groups to work together, I am requesting tables that we can fold up when we are not using them to leave more room for our flexible seating options.\r\nI know that with more seating options, they will be that much more excited about coming to school! Thank you for your support in making my classroom one students will remember forever!annan

=====

In [37]:

```
# 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"\'s", " is", 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 [38]:

```

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

```

\nA person is a person, no matter how small.\n (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. \n\nStudents in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans.\n\nOur school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum.Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, \n"Can we try cooking with REAL food?\n" I will take their idea and create \n"Common Core Cooking Lessons\n" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. \n\nStudents will gain math and literature skills as well as a life long enjoyment for healthy cooking.nannan

In [39]:

```

# \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)

```

A person is a person, no matter how small. (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum.Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, Can we try cooking with REAL food? I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.nannan

In [40]:

```

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

```

A person is a person no matter how small Dr Seuss I teach the smallest students with the biggest

A person is a person no matter how small Dr Seuss I teach the smallest students with the biggest enthusiasm for learning My students learn in many different ways using all of our senses and multiple intelligences I use a wide range of techniques to help all my students succeed Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures including Native Americans Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom Kindergarteners in my class love to work with hands on materials and have many different opportunities to practice a skill before it is mastered Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum Montana is the perfect place to learn about agriculture and nutrition My students love to role play in our pretend kitchen in the early childhood classroom I have had several kids ask me Can we try cooking with REAL food I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce make our own bread and mix up healthy plants from our classroom garden in the spring We will also create our own cookbooks to be printed and shared with families Students will gain math and literature skills as well as a life long enjoyment for healthy cooking nannan

In [41]:

```
# 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"]
```

In [42]:

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

```
100% |██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:52<00:00, 2099.72it/s]
```

7. 1401.

In [43]:

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

Out[43]:

```
'person person no matter small dr seuss teach smallest students biggest enthusiasm learning
students learn many different ways using senses multiple intelligences use wide range techniques h
elp students succeed students class come variety different backgrounds makes wonderful sharing exp
eriences cultures including native americans school caring community successful learners seen coll
aborative student project based learning classroom kindergarteners class love work hands materials
many different opportunities practice skill mastered social skills work cooperatively friends cruc
ial aspect kindergarten curriculum montana perfect place learn agriculture nutrition students love
role play pretend kitchen early childhood classroom several kids ask try cooking real food take id
ea create common core cooking lessons learn important math writing concepts cooking delicious heal
thy food snack time students grounded appreciation work went making food knowledge ingredients cam
e well healthy bodies project would expand learning nutrition agricultural cooking recipes us peel
apples make homemade applesauce make bread mix healthy plants classroom garden spring also create
cookbooks printed shared families students gain math literature skills well life long enjoyment he
althy cooking nannan'
```

1.4 Preprocessing of `project_title`

In [44]:

```
import re

# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
def remove_escape_sequences(phrase):
    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

#remove spacial character: https://stackoverflow.com/a/5843547/4084039
def remove_special_characters(phrase):
    return re.sub(r'^A-Za-z0-9+', ' ', phrase)
```

In [45]:

```
preprocessed_title = list()
for title in tqdm(project_data.project_title.values):
    title = decontracted(title)
    title = remove_escape_sequences(title)
    title = remove_special_characters(title)
    # https://gist.github.com/sebleier/554280
    title = ' '.join([w.lower() for w in title.split(' ') if w.lower() not in set(stopwords)])
    preprocessed_title.append(title.strip())
```

100% | 109248/109248
[00:03<00:00, 30159.73it/s]

1.5 Preparing data for models

In [46]:

```
project_data.columns
```

Out[46]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_grade_category', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', '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 (optinal)
- quantity : numerical (optinal)
- 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 [47]:

```
# 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 [48]:

```
# 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 [49]:

```
# school_states
vectorizer = CountVectorizer(vocabulary = set(project_data.school_state.values), lowercase = False,
binary = True)
school_states_one_hot = vectorizer.fit_transform(project_data.school_state.values)

print('Feature : ', 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 :  (109248, 51)
```

In [52]:

```
project_grade_category = list()
for p in project_data.project_grade_category.values:
    p = p.strip()
    p = p.replace(' ', '_')
    p = p.replace('-', '_')
    project_grade_category.append(p.strip())
project_data['clean_grade_category'] = project_grade_category
```

In [53]:

```
# project_grade_category
vectorizer = CountVectorizer(vocabulary = set(project_data.clean_grade_category.values), lowercase
= False, binary = True)
project_grade_category_one_hot = vectorizer.transform(project_data.clean_grade_category.values)

print('Features : ', 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 :  (109248, 4)
```

In [54]:

```
# replacing nan value by 'nan'
project_data.teacher_prefix.replace(to_replace = np.nan, value = 'nan', inplace = True)

# removing '.'
clean_teacher_prefix = list()
for tp in project_data.teacher_prefix.values:
    tp = tp.replace('.', '')
    clean_teacher_prefix.append(tp)

project_data.teacher_prefix = clean_teacher_prefix
```

In [55]:

```
# teacher_prefix
vectorizer = CountVectorizer(vocabulary = set(project_data.teacher_prefix.values), lowercase = False,
binary = True)
teacher_prefix_one_hot = vectorizer.transform(project_data.teacher_prefix.values)

print('Features : ', 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 :  (109248, 6)
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [56]:

```
# 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, 16512)
```

In [57]:

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

In [58]:

```
# preprocessed_title

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

Shape of matrix after BOW (109248, 2094)

1.5.2.2 TFIDF vectorizer

In [59]:

```
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, 16512)

In [60]:

```
# preprocessed_title

# Considering the words which appeared in at least min_df documents (rows or projects).
# using n_gram = 1
vectorizer = TfidfVectorizer(min_df = 20, ngram_range = (1, 1))
title_tfidf = vectorizer.fit_transform(preprocessed_title)
print("Shape of matrix after TFIDF ", title_tfidf.shape)
```

Shape of matrix after TFIDF (109248, 2094)

1.5.2.3 Using Pretrained Models: Avg W2V

In [61]:

```
'''
# 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-sa
ve-and-load-variables-in-python/

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

'''

```

Out[61]:

```

'\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# =====\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 [62]:

```

# 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 [63]:

```

# 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]))

```

109248
300

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:01<00:00, 86730.12it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[02:59<00:00, 610.23it/s]
```

In [67]:

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

In [68]:

```
# preprocessed_title

vectorizer = TfidfVectorizer()
tfidf_vec = vectorizer.fit(preprocessed_title)

tfidf_model = dict(zip(tfidf_vec.get_feature_names(), tfidf_vec.idf_))
tfidf_words = set(tfidf_vec.get_feature_names())
```

In [69]:

```
title_tfidf_w2v = list()

for sent in tqdm(preprocessed_title):
    temp_w2v = np.zeros(300)
    tfidf_weight = 0
    for word in sent.split():
        if word in glove_words and word in tfidf_words:
            vec = model[word]
            tfidf = tfidf_model[word] * (sent.count(word) / len(sent.split()))
            temp_w2v += vec * tfidf
            tfidf_weight += tfidf
    if tfidf_weight != 0:
        temp_w2v /= tfidf_weight
    title_tfidf_w2v.append(temp_w2v)
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248
[00:02<00:00, 42940.07it/s]
```

1.5.3 Vectorizing Numerical features

In [70]:

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

In [71]:

```
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler

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

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

# Now standardize the data with above maen and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

```
Mean : 298.1193425966608, Standard deviation : 367.49634838483496
```

In [72]:

```
price_standardized
```

Out[72]:

```
array([[ 1.16172762]
```



```
array([[ 1.19112702],
       [-0.23153793],
       [ 0.08402983],
       ...,
       [ 0.27450792],
       [-0.0282706 ],
       [-0.79625102]])
```

In [73]:

```
# quantity
scaler = StandardScaler()
quantity_standardized = scaler.fit_transform(project_data.quantity.values.reshape(-1, 1))
print(f'Before Standardization : \nmean : {scaler.mean_[0]} & std : {np.sqrt(scaler.var_[0])}')
```

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

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

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

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

Before Standardization :
mean : 16.965610354422964 & std : 26.18282191909318

In [74]:

```
# teacher_number_of_previously_posted_projects
scaler = StandardScaler()
prev_posted_project_standardized =
scaler.fit_transform(project_data.teacher_number_of_previously_posted_projects.values.reshape(-1, 1))
print(f'Before Standardization : \nmean : {scaler.mean_[0]} & std : {np.sqrt(scaler.var_[0])}')
```

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

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

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

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

Before Standardization :
mean : 11.153165275336848 & std : 27.77702641477403

1.5.4 Merging all the above features

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

In [75]:

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16512)
(109248, 1)
```

In [76]:

```
# 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[76]:

(109248, 16552)

Assignment 3: Apply KNN

1. [Task-1] Apply KNN(brute force version) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper parameter tuning to find best K

- Find the best hyper parameter which results in the maximum [AUC](#) value
- Find the best hyper parameter using k-fold cross validation (or) simple cross validation data
- Use gridsearch-cv or randomsearch-cv or write your own for loops to do this task

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, as shown in the figure
- Once you find the best hyper parameter, you need to train your model-M using the best hyper-param. Now, find the AUC on test data and plot the ROC curve on both train and test using model-M.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points

4. [Task-2]

- Select top 2000 features from feature **Set 2** using `SelectKBest` and then apply KNN on top of these features

```
from sklearn.datasets import load_digits
from sklearn.feature_selection import SelectKBest, chi2
X, y = load_digits(return_X_y=True)
X.shape
X_new = SelectKBest(chi2, k=20).fit_transform(X, y)
X_new.shape
=====
output:
(1797, 64)
(1797, 20)
```

- Repeat the steps 2 and 3 on the data matrix after feature selection

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](#)

Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on your train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](#).

4. For more details please go through this [link](#).

2. K Nearest Neighbor

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 Sets
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]:

```
# sorting data by time
project_data.sort_values(by = ['project_submitted_datetime'], inplace = True)
```

In [6]:

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

# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
summary_table = PrettyTable(field_names = ['Vectorizer', 'Model', 'Hyper parameter "K"', 'AUC'])
```

Considering 50000 data points

In [7]:

```
data = project_data.iloc[:50000]
y = list(data.project_is_approved.values)
data.drop(columns = 'project_is_approved', axis = 1, inplace = True)
```

In [8]:

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

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

In [9]:

```
c = Counter()
```

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

Out[9]:

```
{'1': 19695, '0': 3755}
```

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 Encoding Numerical data

In [10]:

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

2.2.1.1 quantity

In [11]:

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

In [12]:

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

In [13]:

```
# 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 [14]:

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

In [15]:

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

In [16]:

```
# 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 [17]:

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

In [18]:

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

In [19]:

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

2.2.2 Preprocessing and Encoding Categorical Data

In [20]:

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

2.2.2.1 project_subject_category

In [21]:

```
# 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 [22]:

```
# 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 [23]:

```
# 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 [24]:

```
# 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 [25]:

```
# 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 :  (23450, 9)
```

In [26]:

```
# 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 :  (10050, 9)
```

In [27]:

```
# 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 :  (16500, 9)
```

2.2.2.2 project_subject_subcategory

In [28]:

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

In [29]:

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

In [30]:

```
# 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 [31]:

```
# 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 [32]:

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

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

```
Feature :  ['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government',
'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics',
'EnvironmentalScience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness',
'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writing', 'M
athematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'Socia
lSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix after one hot encoding :  (23450, 30)
```

In [33]:

```
# 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 :  (10050, 30)
```

In [34]:

```
# 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 :  (16500, 30)
```

2.2.2.3 project_grade_category

In [35]:

```
# 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 [36]:

```
# 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 [37]:

```
# 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 [38]:

```
# 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 [39]:

```
# 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 :  (23450, 4)
```

In [40]:

```
# 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 :  (10050, 4)
```

In [41]:

```
# 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 :  (16500, 4)
```

2.2.2.4 teacher_prefix

In [42]:

```
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 [43]:

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

In [44]:

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

In [45]:

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

In [46]:


```
# 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']
Shape of matrix after one hot encoding :  (23450, 5)
```

In [47]:

```
# 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 :  (10050, 5)
```

In [48]:

```
# 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 :  (16500, 5)
```

2.2.2.5 school_sates

In [49]:

```
# 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 :  (23450, 51)
```

In [50]:

```
# 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 :  (10050, 51)
```

In [51]:

```
# 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 :  (16500, 51)
```

2.3 Make Data Model Ready: encoding eassay, and project_title

In [52]:

```
# 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
```

2.3.2 preprocessing essay and title

In [53]:

```
# 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 [54]:

```
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 [55]:

```
# 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', "do
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"]
```

2.3.2.1 preprocessing essay

In [56]:

```
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 [57]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 23450/23450
[00:22<00:00, 1024.04it/s]
```

In [58]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 10050/10050
[00:09<00:00, 1035.71it/s]
```

In [59]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:15<00:00, 1042.14it/s]
```

2.3.2.2 preprocessing title

In [60]:

```
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 = 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 [61]:

```

# project_title train data
preprocess_title(X_train)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 23450/23450
[00:00<00:00, 31639.10it/s]

```

In [62]:

```

# project_title cross validation data
preprocess_title(X_cv)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 10050/10050
[00:00<00:00, 31685.61it/s]

```

In [63]:

```

# project_title test data
preprocess_title(X_test)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:00<00:00, 31097.97it/s]

```

2.3.2 Vectorizing text data (train)

2.3.2.1 Bag Of Words (BOW)

In [64]:

```

# preprocessed_essay train data

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

```

Shape of matrix after BOW (23450, 6486)

In [65]:

```

# 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 (10050, 6486)

In [66]:

```

# 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 (16500, 6486)

In [67]:

```
# 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 (23450, 680)

In [68]:

```
# 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 (10050, 680)

In [69]:

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

Shape of matrix after BOW (16500, 680)

2.3.2.2 Term Frequency Inverse Document Frequency (TFIDF)

In [70]:

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

In [71]:

```
# 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 (23450, 6486)

In [72]:

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

Shape of matrix after TFIDF (10050, 6486)

In [73]:

```
# 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 (16500, 6486)

In [74]:

```
# 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 (23450, 680)

In [75]:

```
# 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 (10050, 680)

In [76]:

```
# 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 (16500, 680)

2.3.2.3 Average Word 2 Vec (W2V)

In [77]:

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

In [78]:

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

In [79]:

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

```
100%|████████████████████████████████████████████████████████████████████████████████| 23450/23450
[00:05<00:00, 4248.43it/s]
```

In [80]:

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

In [81]:

In [82]:

In [83]:

In [84]:

2.3.2.4 TFIDF weighted Word 2 Vec (W2V)

In [85]:

In [86]:

```
# preprocessed_essay train data

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

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

$\frac{1}{x^2} = x^{-2}$

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

```
100%|███████████████████████████████████████████████████████████| 23450/23450 [00:  
40<00:00, 579.31it/s]
```

In [87]:

```
# preprocessed_essay cross validation data
```

```
cv_essay_tfidf_w2v = tfidf_w2v(X_cv.preprocessed_essay, model, glove_words, essay_tfidf_model,
essay_tfidf_words)
```

```
100%|███████████████████████████████████████████████████████████| 10050/10050 [00:  
17<00:00, 585.73it/s]
```

In [88]:

```
# preprocessed essay cross validation data
```

```
test_essay_tfidf_w2v = tfidf_w2v(X_test.preprocessed_essay, model, glove_words, essay_tfidf_model,
essay_tfidf_words)
```

```
100%|███████████████████████████████████████████████████████| 16500/16500 [00:  
28<00:00, 588.14it/s]
```

In [89]:

```
# preprocessed title train data
```

```
vectorizer = TfidfVectorizer()
title_tfidf_vec = vectorizer.fit(X_train.preprocessed_title)
```

```
# Making a dictionary with key as word and value as idf value of that word
```

```
title_tfidf_model = dict(zip(title_tfidf_vec.get_feature_names(), title_tfidf_vec.idf_))
title_tfidf_words = set(title_tfidf_vec.get_feature_names())
```

```
title_tfidf_w2v = tfidf_w2v(X_train.preprocessed_title, model, glove_words, title_tfidf_model,
title_tfidf_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 23450/23450  
[00:00<00:00, 38168.48it/s]
```

In [90]:

```
# preprocessed title cross validation data
```

```
cv_title_tfidf_w2v = TfidfVectorizer(X_cv.preprocessed_title, model, glove_words, title_tfidf_model,
                                     title_tfidf_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 10050/10050  
[00:00<00:00, 38315.18it/s]
```

In [91]:

```
# preprocessed title test data
```

```
test_title_tfidf_w2v = tfidf_w2v(X_test.preprocessed_title, model, glove_words, title_tfidf_model,
title_tfidf_words)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 35808.02it/s]
```

2.4 Applying KNN on different kind of featurization as mentioned in the instructions

Apply KNN 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 [92]:

```
# 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 [93]:

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score, roc_curve
```

In [94]:

```
# plot confusion matrix python; https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html
from sklearn.metrics import confusion_matrix
from sklearn.utils.multiclass import unique_labels

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

    # Compute confusion matrix
    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 KNN brute force on BOW, SET 1

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_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 : (23450, 7268)
y_train shape : (23450,)
```

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_bow, cv_title_bow)).tocsr()
print('X_cv shape : ', X_cv.shape)

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

```
X_cv shape : (10050, 7268)
y_cv shape : (10050,)
```

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_bow, test_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 : (16500, 7268)
y_test shape : (16500,)
```

In [99]:

```
def batch_predict(model, X, batch_size = 1000):
    predicted = list()
    for i in range(0, X.shape[0] - batch_size + 1, batch_size):
        pred = model.predict_proba(X[i : i + batch_size])[:, 1]
        predicted.extend(pred)
    if X.shape[0] % batch_size != 0:
        predicted.extend(model.predict_proba(X[X.shape[0] - X.shape[0] % batch_size :])[:, 1])
    return np.array(predicted)
```

In [103]:

```
# training the knn_model_bow
k_range = [1, 3, 5, 7, 9, 11, 15, 51, 81, 111, 151, 251, 351, 451, 651]
cv_auc_scores = dict()
train_auc_scores = dict()

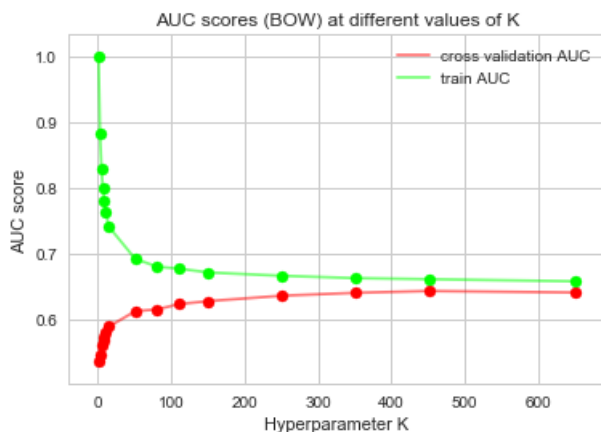
for k in tqdm(k_range):
    knn_model_bow = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute', n_jobs = -1)
    # fitting the train data
    knn_model_bow.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(knn_model_bow, X_train, batch_size = 400)
    predicted_cv = batch_predict(knn_model_bow, X_cv, batch_size = 400)

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

[illegible]

In [104]:

```
# plotting AUC scores
sns.set(style = 'whitegrid')
plt.plot(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF00088', label = 'cross validation AUC')
plt.plot(list(train_auc_scores.keys()), list(train_auc_scores.values()), color = '#00FF088', label = 'train AUC')
plt.scatter(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF000FF')
plt.scatter(list(train_auc_scores.keys()), list(train_auc_scores.values()), color= '#00FF0FF')
plt.legend()
plt.xlabel('Hyperparameter K')
plt.ylabel('AUC score')
plt.title('AUC scores (BOW) at different values of K')
plt.show()
```



In [105]:

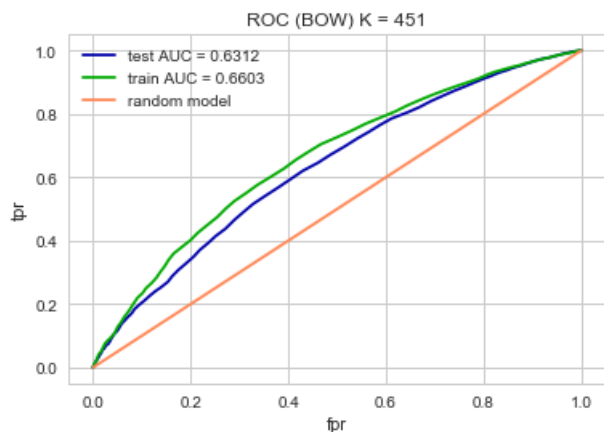
```
best_k = k_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best knn model (BOW)
best_knn_model_bow = KNeighborsClassifier(n_neighbors = best_k, algorithm = 'brute', n_jobs = -1)
# fitting the train data
best_knn_model_bow.fit(X_train, y_train)

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

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

In [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 (BOW) K = {best_k}')
plt.show()
```



In [110]:

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

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

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

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

In [108]:

```
# adding AUC score to summary table
summary_table.add_row(['BOW', 'Brute', 451, 0.6312])

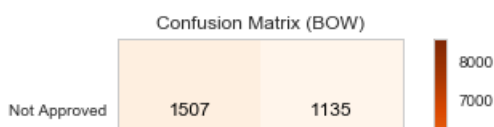
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

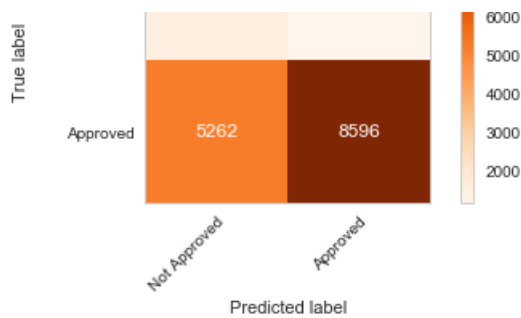
# Plotting confusion matrix
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (BOW)', cmap = plt.cm.Oranges)
```

the maximum value of $tpr*(1-fpr)$ 0.3538150390753669 for threshold 0.778

Confusion matrix

```
[[1507 1135]
 [5262 8596]]
```





Out[108]:

<matplotlib.axes._subplots.AxesSubplot at 0x16b874be940>

2.4.2 Applying KNN brute force on TFIDF, SET 2

In [109]:

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

In [110]:

```
# 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 : (23450, 7268)

y_train shape : (23450,)

In [111]:

```
# 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 : (10050, 7268)

y_cv shape : (10050,)

In [112]:

```
# 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_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 : (16500, 7268)
y_test shape : (16500,)

```

In [113]:

```

# training the knn_model_tfidf
k_range = [1, 3, 5, 7, 9, 11, 15, 51, 81, 111, 151, 251, 351, 451, 651]
cv_auc_scores = dict()
train_auc_scores = dict()

for k in tqdm(k_range):
    knn_model_tfidf = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute', n_jobs = -1)
    # fitting the train data
    knn_model_tfidf.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(knn_model_tfidf, X_train, batch_size = 400)
    predicted_cv = batch_predict(knn_model_tfidf, X_cv, batch_size = 400)

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

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 15/15
[16:55<00:00, 69.44s/it]

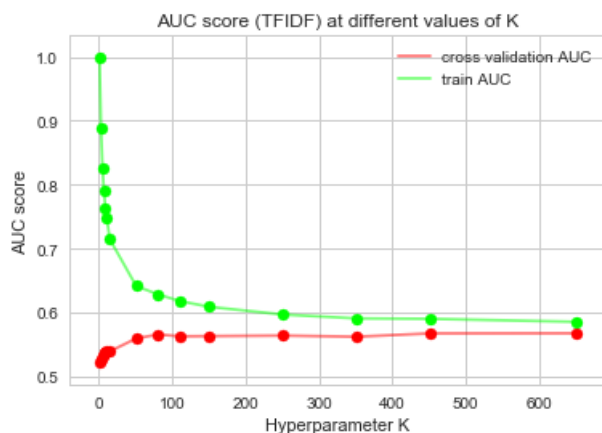
```

In [114]:

```

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

```



In [115]:

```

best_k = k_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best knn model (TFIDF)
best_knn_model_tfidf = KNeighborsClassifier(n_neighbors = best_k, algorithm = 'brute', n_jobs = -1)
# fitting the train data
best_knn_model_tfidf.fit(X_train, y_train)

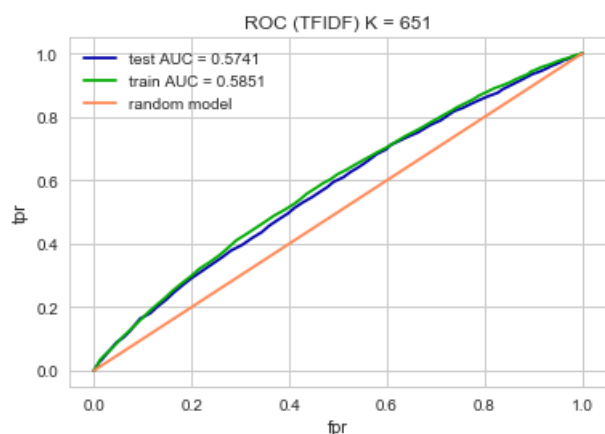
```

```
# predicting the probability scores of train_data, test_data
predicted_test = batch_predict(best_knn_model_tfidf, X_test, 400)
predicted_train = batch_predict(best_knn_model_tfidf, X_train, 400)

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

In [116]:

```
# 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) K = {best_k}')
plt.show()
```



In [117]:

```
# adding AUC score to summary table
summary_table.add_row(['TFIDF', 'Brute', 651, 0.5741])

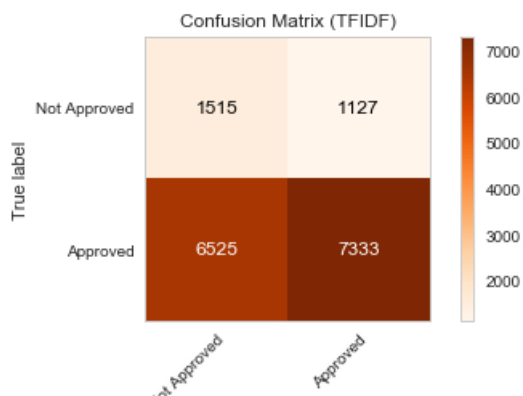
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.3034316981071884 for threshold 0.849

Confusion matrix

```
[[1515 1127]
 [6525 7333]]
```



Predicted label

Out[117]:

<matplotlib.axes._subplots.AxesSubplot at 0x16b9a2769b0>

2.4.3 Applying KNN brute force on AVG W2V, SET 3

In [118]:

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

In [95]:

```
# train data

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

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

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

X_train shape : (23450, 702)

y_train shape : (23450,)

In [96]:

```
# 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_avg_w2v, cv_title_avg_w2v)).tocsr()
print('X_cv shape : ', X_cv.shape)

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

X_cv shape : (10050, 702)

y_cv shape : (10050,)

In [97]:

```
# test data

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_test = hstack((test_school_states_one_hot, test_category_one_hot, test_subcategory_one_hot,
                 test_project_grade_category_one_hot, test_teacher_prefix_one_hot,
                 test_quantity_standardized,
                 test_prev_posted_project_standardized, test_price_standardized,
                 test_essay_avg_w2v, test_title_avg_w2v)).tocsr()
print('X_test shape : ', X_test.shape)

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

X_test shape : (16500, 702)

y_test shape : (16500,)

In [100]:

```
# training the knn_model_avg_w2v
k_range = [1, 3, 5, 11, 31, 51, 111, 151, 251, 351, 451, 651]
cv_auc_scores = dict()
train_auc_scores = dict()

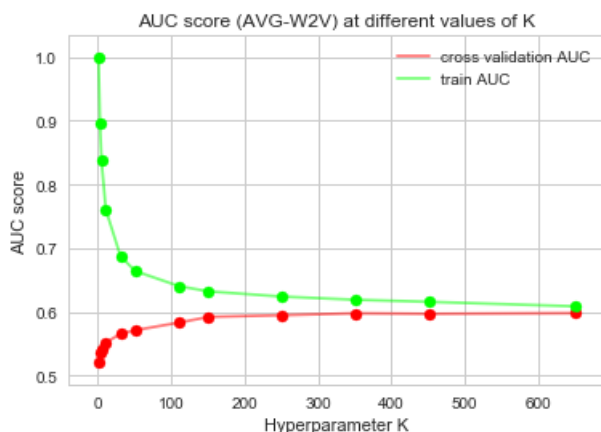
for k in tqdm(k_range):
    knn_model_avg_w2v = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute', n_jobs = -1)
    # fitting the train data
    knn_model_avg_w2v.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(knn_model_avg_w2v, X_train, batch_size = 400)
    predicted_cv = batch_predict(knn_model_avg_w2v, X_cv, batch_size = 400)

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

[illegible]

In [101]:

```
# plotting auc curve
sns.set(style = 'whitegrid')
plt.plot(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF00088', label = 'cross validation AUC')
plt.plot(list(train_auc_scores.keys()), list(train_auc_scores.values()), color = '#00FF088', label = 'train AUC')
plt.scatter(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF000FF')
plt.scatter(list(train_auc_scores.keys()), list(train_auc_scores.values()), color = '#00FF0FF')
plt.legend()
plt.xlabel('Hyperparameter K')
plt.ylabel('AUC score')
plt.title('AUC score (AVG-W2V) at different values of K')
plt.show()
```



In [102]:

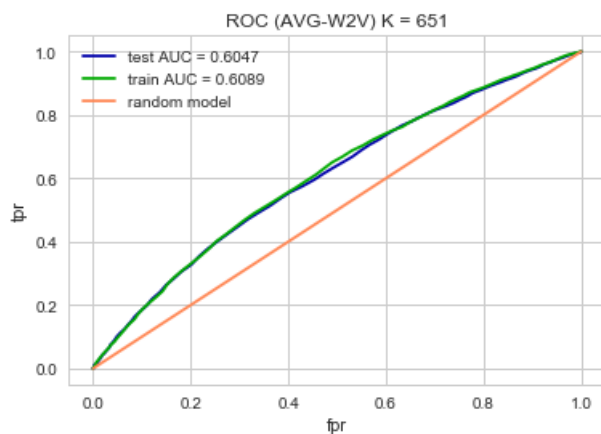
```
best_k = k_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best knn model (AVG-W2V)
best_knn_model_avg_w2v = KNeighborsClassifier(n_neighbors = best_k, algorithm = 'brute', n_jobs = -1)
# fitting the train data
best_knn_model_avg_w2v.fit(X_train, y_train)

# predicting the probability scores of train data, test data
predicted_test = batch_predict(best_knn_model_avg_w2v, X_test, 400)
predicted_train = batch_predict(best_knn_model_avg_w2v, X_train, 400)

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

In [103]:

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



In [118]:

```
# adding AUC score to summary table
summary_table.add_row(['AVG-W2V', 'Brute', 651, 0.6047])

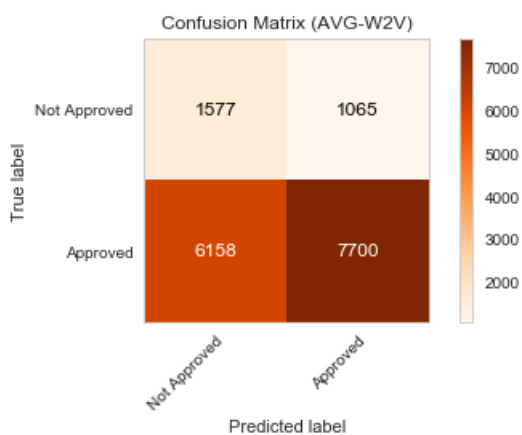
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.33165690852246466 for threshold 0.846

Confusion matrix

```
[[1577 1065]
 [6158 7700]]
```



Out[118]:

<matplotlib.axes._subplots.AxesSubplot at 0x1cb397cb518>

2.4.4 Applying KNN brute force on TFIDF W2V, SET 4

In []:

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

In [119]:

```
# train data

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

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

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

X_train shape : (23450, 702)
y_train shape : (23450,)
```

In [120]:

```
# cross validation data

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

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

X_cv shape : (10050, 702)
y_cv shape : (10050,)
```

In [121]:

```
# test data

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

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

X_test shape : (16500, 702)
y_test shape : (16500,)
```

In [124]:

```
# training the knn_model_tfidf_w2v
k_range = [1, 3, 5, 11, 31, 51, 111, 151, 251, 351, 451, 651]
cv_auc_scores = dict()
train_auc_scores = dict()

for k in tqdm(k_range):
```

```

knn_model_tfidf_w2v = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute', n_jobs = -1)
# fitting the train data
knn_model_tfidf_w2v.fit(X_train, y_train)
# predicting the probability scores of train data, cross validation data
predicted_train = batch_predict(knn_model_tfidf_w2v, X_train, batch_size = 400)
predicted_cv = batch_predict(knn_model_tfidf_w2v, X_cv, batch_size = 400)

# calculating AUC score for cross validation and test data
cv_auc_scores[k] = roc_auc_score(y_cv[:predicted_cv.shape[0]], predicted_cv)
train_auc_scores[k] = roc_auc_score(y_train[:predicted_train.shape[0]], predicted_train)

```

```

0%|          | 0
[00:00<?, ?it/s]

8%|          | 1/12 [08:55<
8:07, 535.21s/it]

17%|         | 2/12
[17:40<1:28:43, 532.33s/it]

25%|         | 3/12 [26:31<
19:46, 531.79s/it]

33%|         | 4/12 [35:18<
10:42, 530.34s/it]

42%|         | 5/12
[44:05<1:01:45, 529.32s/it]

50%|         | 6/12 [52:5
<52:53, 528.94s/it]

58%|         | 7/12 [1:01:4
<44:04, 528.83s/it]

67%|         | 8/12
[1:10:32<35:17, 529.29s/it]

75%|         | 9/12
[1:19:22<26:28, 529.47s/it]

83%|         | 10/12 [1:28:1
2<17:39, 529.62s/it]

92%|         | 11/12
[1:37:10<08:52, 532.17s/it]

100%|        | 12/12
[1:46:01<00:00, 532.00s/it]

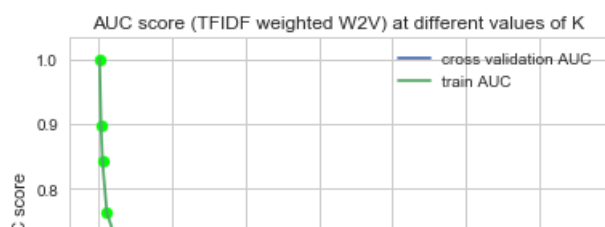
```

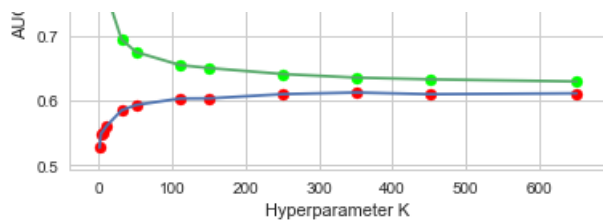
In [125]:

```

# plotting AUC score
sns.set(style = 'whitegrid')
plt.plot(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), label = 'cross validation AUC')
plt.plot(list(train_auc_scores.keys()), list(train_auc_scores.values()), label = 'train AUC')
plt.scatter(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF0000FF')
plt.scatter(list(train_auc_scores.keys()), list(train_auc_scores.values()), color = '#00FF00FF')
plt.legend()
plt.xlabel('Hyperparameter K')
plt.ylabel('AUC score')
plt.title('AUC score (TFIDF weighted W2V) at different values of K')
plt.show()

```





In [126]:

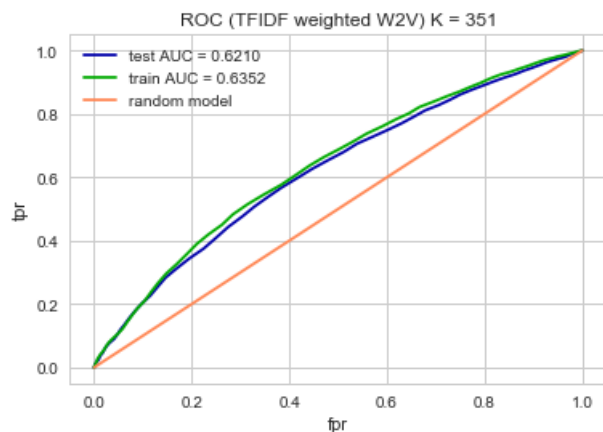
```
best_k = k_range[np.argmax(np.array(list(cv_auc_scores.values())))]
# best knn model (TFIDF-W2V)
best_knn_model_tfidf_w2v = KNeighborsClassifier(n_neighbors = best_k, algorithm = 'brute', n_jobs =
-1)
# fitting the train data
best_knn_model_tfidf_w2v.fit(X_train, y_train)

# predicting the probability scores of train_data, test_data
predicted_test = batch_predict(best_knn_model_tfidf_w2v, X_test, 400)
predicted_train = batch_predict(best_knn_model_tfidf_w2v, X_train, 400)

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

In [127]:

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



In [128]:

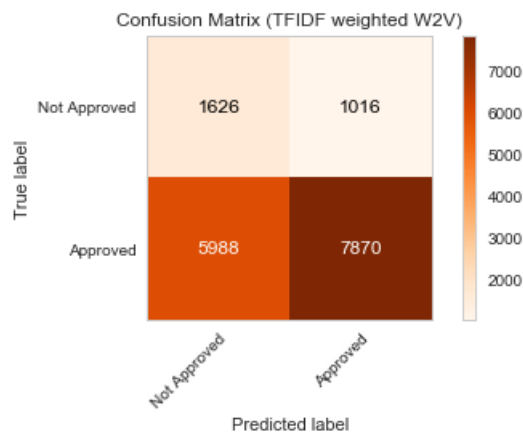
```
# adding AUC score to summary table
summary_table.add_row(['TFIDF-W2V', 'Brute', 351, 0.621])

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

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.3495118487953241 for threshold 0.84
Confusion matrix

```
[[1626 1016]
 [5988 7870]]
```



Out[128]:

<matplotlib.axes._subplots.AxesSubplot at 0x1cb39803080>

2.5 Feature selection with `SelectKBest`

In [129]:

```
# 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 [130]:

```
from sklearn.feature_selection import SelectKBest, chi2
```

In [131]:

```
# 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)

# cross validation data
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)
```

```
# test data
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_train shape : (23450, 7268)
y_train shape : (23450,)
X_cv shape : (10050, 7268)
y_cv shape : (10050,)
X_test shape : (16500, 7268)
y_test shape : (16500,)
```

In [132]:

```
# selecting top 2000 features
select_model = SelectKBest(chi2, k = 2000)
X_train = select_model.fit_transform(X_train, y_train)
X_cv = select_model.transform(X_cv)
X_test = select_model.transform(X_test)
print('X_train shape : ', X_train.shape)
print('X_cv shape : ', X_cv.shape)
print('X_test shape : ', X_test.shape)
```

```
X_train shape : (23450, 2000)
X_cv shape : (10050, 2000)
X_test shape : (16500, 2000)
```

In [133]:

```
# training the knn_model_tfidf
k_range = [1, 3, 5, 7, 11, 15, 31, 51, 111, 151, 251, 351, 551]
cv_auc_scores = dict()
train_auc_scores = dict()

for k in tqdm(k_range):
    knn_model_tfidf = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute', n_jobs = -1)
    # fitting the train data
    knn_model_tfidf.fit(X_train, y_train)
    # predicting the probability scores of train data, cross validation data
    predicted_train = batch_predict(knn_model_tfidf, X_train, batch_size = 400)
    predicted_cv = batch_predict(knn_model_tfidf, X_cv, batch_size = 400)

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

```
0%|          | C
[00:00<?, ?it/s]

8%|          | 1/13 [00:
10:00, 50.07s/it]

15%|         | 2/13
[01:39<09:07, 49.74s/it]

23%|         | 3/13 [02:
<08:25, 50.56s/it]

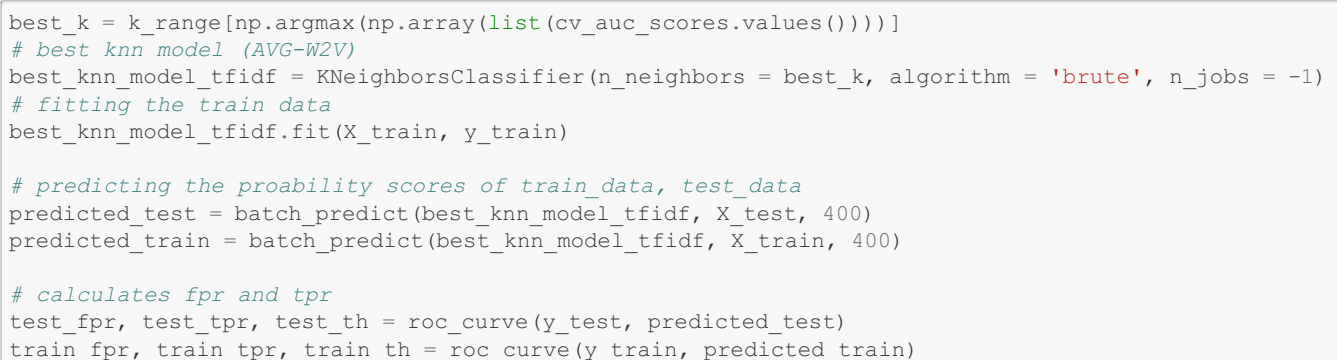
31%|         | 4/13 [03:
<07:40, 51.19s/it]

38%|         | 5/13
[04:16<06:52, 51.59s/it]

46%|         | 6/13
```

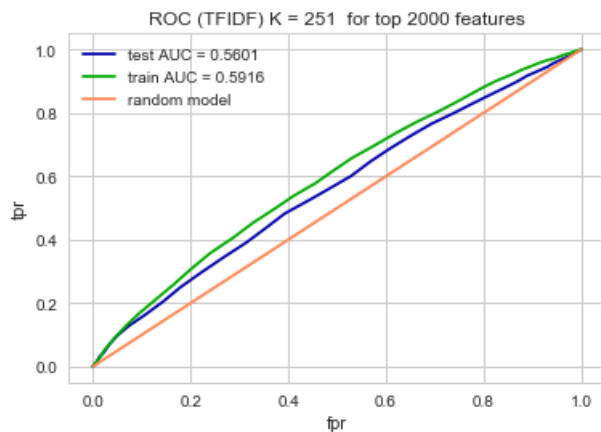
[illegible]

```
# plotting auc curve
sns.set(style = 'whitegrid')
plt.plot(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF00088', label = 'cross validation AUC')
plt.plot(list(train_auc_scores.keys()), list(train_auc_scores.values()), color = '#00FF088', label = 'train AUC')
plt.scatter(list(cv_auc_scores.keys()), list(cv_auc_scores.values()), color = '#FF000FF')
plt.scatter(list(train_auc_scores.keys()), list(train_auc_scores.values()), color= '#00FF0FF')
plt.legend()
plt.xlabel('Hyperparameter K')
plt.ylabel('AUC score')
plt.title('AUC score (TFIDF) at different values of K for top 2000 features')
plt.show()
```



In [136]:

```
# Plots ROC curve
sns.set(style = 'whitegrid')
plt.plot(test_fpr, test_tpr, color = '#0000AA', label = 'test AUC = %.4f' % (roc_auc_score(y_test,
predicted_test)))
plt.plot(train_fpr, train_tpr, color = '#00AA00', label = 'train AUC = %.4f' % (roc_auc_score(y_train,
predicted_train)))
plt.plot([0, 1], [0, 1], color = '#FF8855', label = 'random model')
plt.legend()
plt.xlabel('fpr')
plt.ylabel('tpr')
plt.title(f'ROC (TFIDF) K = {best_k} for top 2000 features')
plt.show()
```



In [137]:

```
# adding AUC score to summary table
summary_table.add_row(['TFIDF (Top 2000 features)', 'Brute', 251, 0.5601])

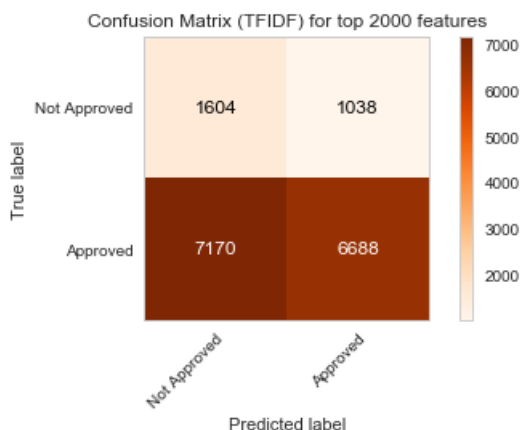
# predicting class labels for test data
pred_test = predict(predicted_test, test_th, test_fpr, test_tpr)

# Plotting confusion matrix
class_names = np.array(['Not Approved', 'Approved'])
plot_confusion_matrix(y_test, pred_test, classes = class_names,
                      title = 'Confusion Matrix (TFIDF) for top 2000 features', cmap =
plt.cm.Oranges)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.29299975560483765 for threshold 0.849

Confusion matrix

```
[[1604 1038]
 [7170 6688]]
```



Out[137]:

<matplotlib.axes._subplots.AxesSubplot at 0x1cb3cbf6cc0>

3. Conclusions

In [138]:

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

In [139]:

```
print(summary_table)
```

Vectorizer	Model	Hyper parameter "K"	AUC
BOW	Brute	451	0.6312
TFIDF	Brute	651	0.5741
AVG-W2V	Brute	651	0.6047
TFIDF-W2V	Brute	351	0.621
TFIDF (Top 2000 features)	Brute	251	0.5601