

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	Descr
<code>project_id</code>	A unique identifier for the proposed project. Example: p0:

Feature		Description	Example
project_title	•	Title of the project.	Art Will Make You Happy
	•		First Grade
project_grade_category		Grade level of students for which the project is targeted. One of the following enumerated values:	
	•		Grades Pre-K
	•		Grades K-1
	•		Grades 1-2
project_subject_categories		One or more (comma-separated) subject categories for the project from the following enumerated list of values:	
	•		Applied Learning
	•		Care & Health
	•		Health & Safety
	•		History & Civics
	•		Literacy & Language
	•		Math & Science
	•		Music & The Arts
	•		Special Interest
	•		World Languages
school_state	•		Music & The Arts
	•		Literacy & Language, Math & Science
school_state		State where school is located (Two-letter U.S. postal code) (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)	Example: CA

Feature	Description
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Example: Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to make sensory needs!<
<code>project_essay_1</code>	First application
<code>project_essay_2</code>	Second application
<code>project_essay_3</code>	Third application
<code>project_essay_4</code>	Fourth application
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016-01-12T12:43:56Z
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c0
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> • • • • • •
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. Example:

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



Notes on the Essay Data

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

- **project_essay_1:** "Introduce us to your classroom"
- **project_essay_2:** "Tell us more about your students"
- **project_essay_3:** "Describe how your students will use the materials you're requesting"
- **project_essay_3:** "Close by sharing why your project will make a difference"

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

- **project_essay_1:** "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- **project_essay_2:** "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

```
In [2]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

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

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

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

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

from tqdm import tqdm
import os

from plotly import plotly
```

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

```
C:\Users\narayana\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")
```

1.1 Reading Data

```
In [3]: project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

```
In [4]: 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 [5]: 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[5]:

	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 [6]: categories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/473019

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "
        if 'The' in j.split(): # this will split each of the category based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math
        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 [7]: sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python: https://stackoverflow.com/a/473019

# 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", "
        if 'The' in j.split(): # this will split each of the category based on space "Math
            j=j.replace('The', '') # if we have the words "The" we are going to replace it w
        j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math
        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]))

```

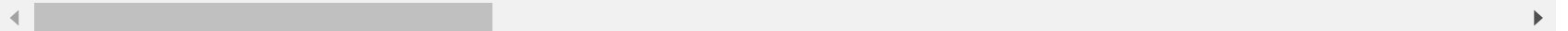
1.3 Text preprocessing

```
In [8]: # 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 [9]: project_data.head(2)
```

Out[9]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datet
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:40
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	2016-10-25 09:20



```
In [10]: ##### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

```
In [11]: # printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect. \r\n\r\n "The limits of your language are the limits of your world." -Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English along side of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills. \r\n\r\n By providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills. \r\n\r\n Parents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students. \r\n\r\n

=====

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

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

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day. \r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas. \r\n\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the su

ccess in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

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

=====
The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive

ve the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.nannan
=====

In [12]: [# https://stackoverflow.com/a/47091490/4084039](https://stackoverflow.com/a/47091490/4084039)

```
import re

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

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

```
In [13]: sent = decontracted(project_data['essay'].values[20000])  
print(sent)  
print("="*50)
```

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


```
In [14]: # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python,
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

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

```
In [15]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039  
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)  
print(sent)
```

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

```
In [16]: # https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "
you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they',
'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'l
'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'u
'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'd
'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over',
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', '
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'v
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now',
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'do
"hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
"mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn'
'won', "won't", 'wouldn', "wouldn't", "nan", "nannan"]
```

```
In [17]: # 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('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e.lower() for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109  
248 [01:25<00:00, 1274.44it/s]
```

```
In [18]: # after preprocesing  
preprocessed_essays[20000]
```

```
Out[18]: 'kindergarten students varied disabilities ranging speech language delays cognitive delay  
s gross fine motor delays autism eager beavers always strive work hardest working past li  
mitations materials ones seek students teach title school students receive free reduced p  
rice lunch despite disabilities limitations students love coming school come eager learn  
explore ever felt like ants pants needed groove move meeting kids feel time want able mov  
e learn say wobble chairs answer love develop core enhances gross motor turn fine motor s  
kills also want learn games kids not want sit worksheets want learn count jumping playing  
physical engagement key success number toss color shape mats make happen students forget  
work fun 6 year old deserves'
```

1.4 Preprocessing of project_title

Following Code blocks provided by me.


```
In [21]: project_data.columns
```

```
Out[21]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
               'project_submitted_datetime', 'project_grade_category', 'project_title',  
               'project_essay_1', 'project_essay_2', 'project_essay_3',  
               'project_essay_4', 'project_resource_summary',  
               'teacher_number_of_previously_posted_projects', 'project_is_approved',  
               'clean_categories', 'clean_subcategories', 'essay'],  
              dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (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/ \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/)

```
In [22]: # 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, binarize=True)
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", categories_one_hot.shape)
```

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

```
In [23]: # we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binarize=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", sub_categories_one_hot.shape)
```

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

```
In [24]: # you can do the similar thing with state, teacher_prefix and project_grade_category also
```

Following Code blocks provided by me.

```
In [25]: # Code took from original code provided.
states = project_data['school_state'].unique()
vectorizer = CountVectorizer(vocabulary=list(states), lowercase=False, binary=True)
vectorizer.fit(project_data['school_state'].values)
print(vectorizer.get_feature_names())

school_state_one_hot = vectorizer.transform(project_data['school_state'].values)
print("Shape of matrix after one hot encoding", school_state_one_hot.shape)

['IN', 'FL', 'AZ', 'KY', 'TX', 'CT', 'GA', 'SC', 'NC', 'CA', 'NY', 'OK', 'MA', 'NV', 'O
H', 'PA', 'AL', 'LA', 'VA', 'AR', 'WA', 'WV', 'ID', 'TN', 'MS', 'CO', 'UT', 'IL', 'MI',
'HI', 'IA', 'RI', 'NJ', 'MO', 'DE', 'MN', 'ME', 'WY', 'ND', 'OR', 'AK', 'MD', 'WI', 'SD',
'NE', 'NM', 'DC', 'KS', 'MT', 'NH', 'VT']
Shape of matrix after one hot encoding (109248, 51)
```

There are some NaN's in teacher_prefix column. replacing them with 'Mrs.' as that has high occurrence in that column.

```
In [26]: print("Number of NaN's before replacement in column: ", sum(project_data['teacher_prefix'].
project_data['teacher_prefix'] = project_data['teacher_prefix'].replace(np.nan, 'Mrs.', reg
print("Number of NaN's after replacement in column: ", sum(project_data['teacher_prefix'].i

# Output may show both zeros as I re-run this several times. But there are 3 zeros in origi

Number of NaN's before replacement in column: 3
Number of NaN's after replacement in column: 0
```



```
In [27]: # Code took from original code provided.
prefixes = project_data['teacher_prefix'].unique()
vectorizer = CountVectorizer(vocabulary=list(prefixes), lowercase=False, binary=True)
vectorizer.fit(project_data['teacher_prefix'].values)
print(vectorizer.get_feature_names())

teacher_prefix_one_hot = vectorizer.transform(project_data['teacher_prefix'].values)
print("Shape of matrix after one hot encoding", teacher_prefix_one_hot.shape)

['Mrs.', 'Mr.', 'Ms.', 'Teacher', 'Dr.']
Shape of matrix after one hot encoding (109248, 5)
```

```
In [28]: grades = project_data['project_grade_category'].unique()
vectorizer = CountVectorizer(vocabulary=list(grades), lowercase=False, binary=True)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())

project_grade_category_one_hot = vectorizer.transform(project_data['project_grade_category'].values)
print("Shape of matrix after one hot encoding", project_grade_category_one_hot.shape)

['Grades PreK-2', 'Grades 6-8', 'Grades 3-5', 'Grades 9-12']
Shape of matrix after one hot encoding (109248, 4)
```

Following Code blocks present in original notebook.

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

```
In [29]: # We are considering only the words which appeared in at least 10 documents(rows or project  
vectorizer = CountVectorizer(min_df=10)  
text_bow = vectorizer.fit_transform(preprocessed_essays)  
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

Shape of matrix after one hot encodig (109248, 16511)

```
In [30]: # you can vectorize the title also  
# before you vectorize the title make sure you preprocess it
```

Following Code blocks provided by me.

```
In [31]: # Code took from original code provided.  
# We are considering only the words which appeared in at least 5 documents(rows or projects  
# Reduced number as title has less words  
vectorizer = CountVectorizer(min_df=10)  
titles_bow = vectorizer.fit_transform(preprocessed_titles)  
print("Shape of matrix after one hot encodig ", titles_bow.shape)
```

Shape of matrix after one hot encodig (109248, 3222)

Following Code blocks present in original notebook.

1.5.2.2 TFIDF vectorizer

```
In [32]: from sklearn.feature_extraction.text import TfidfVectorizer  
vectorizer = TfidfVectorizer(min_df=10)  
text_tfidf = vectorizer.fit_transform(preprocessed_essays)  
print("Shape of matrix after one hot encoding ",text_tfidf.shape)
```

Shape of matrix after one hot encoding (109248, 16511)

1.5.2.3 Using Pretrained Models: Avg W2V

```
In [33]: '''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

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

# =====

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

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

inter_words = set(model.keys()).intersection(words)
```

```

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

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

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickl

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

...

```

Out[33]: `'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef (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 model[word] = embedding\n print ("Done.",len(model)," words loaded!")\n return model\nmodel = loadGloveModel(\'glove.42B.300d.txt\')\n\n# =====\nOutput:\n\nLoading Glove 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(\' \'))\n\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 that 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 wor`

```
ds_courpus[i] = model[i]\nprint("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 (http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport) pickle\nwith open('glove_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n\n
```

```
In [34]: # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickl  
# 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 [35]: # 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]))
```

```
In [37]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))) # getting
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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



```
In [39]: # Code took from original code provided.
# tfidf of project titles
vectorizer = TfidfVectorizer(min_df=10)
titles_tfidf = vectorizer.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ", titles_tfidf.shape)
```


1.5.3 Vectorizing Numerical features

```
In [43]: 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 [44]: # check this one: https://www.youtube.com/watch?v=0H0q0cLn3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScaler.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.
# 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 standar
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 [45]: price_standardized
```

```
Out[45]: array([[ -0.3905327 ],
 [  0.00239637],
 [  0.59519138],
 ...,
 [-0.15825829],
 [-0.61243967],
 [-0.51216657]])
```

Following Code blocks provided by me.

```
In [46]: warnings.filterwarnings("ignore")
# Code took from original code provided
scalar = StandardScaler()
scalar.fit(project_data['teacher_number_of_previously_posted_projects'].values.reshape(-1,
print(f"Mean : {scalar.mean_[0]}, Standard deviation : {np.sqrt(scalar.var_[0])}")

# Now standardize the data with above mean and variance.
previously_posted_projects_standardized = \
    scalar.transform(project_data['teacher_number_of_previously_posted_projects
print(previously_posted_projects_standardized)
```

Mean : 11.153165275336848, Standard deviation : 27.77702641477403

```
[[-0.40152481]
 [-0.14951799]
 [-0.36552384]
 ...
 [-0.29352189]
 [-0.40152481]
 [-0.40152481]]
```

Following Code blocks present in original notebook.

1.5.4 Merging all the above features

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

```
In [47]: print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16511)
(109248, 1)
```

```
In [48]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

```
Out[48]: (109248, 16551)
```

```
In [49]: # please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Assignment 4: Naive Bayes

1. Apply Multinomial NaiveBayes on these feature sets

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

2. The hyper paramter tuning(find best Alpha)

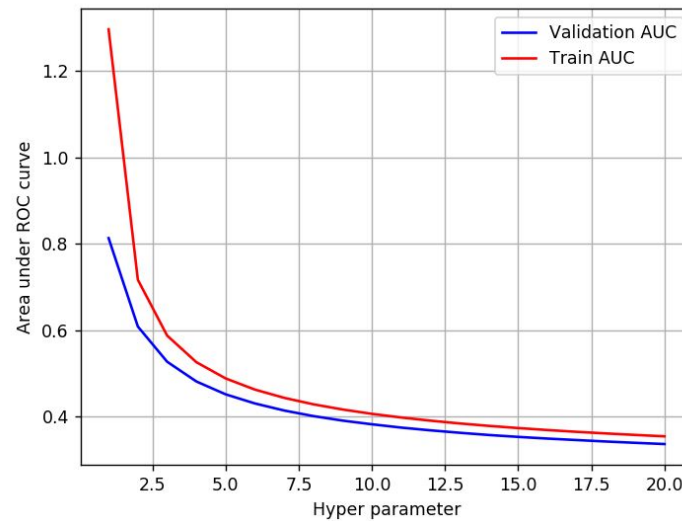
- Find the best hyper parameter which will give the maximum [AUC](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/>) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Feature importance

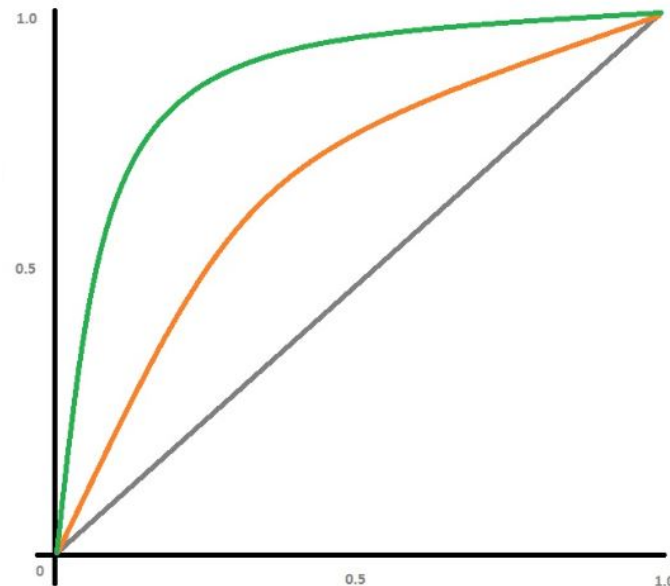
- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using values of `feature_log_prob_` parameter of [MultinomialNB](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html) (https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html) and print their corresponding feature names

4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure. Here on X-axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.



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



- Along with plotting ROC curve, you need to print the [confusion matrix](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tp-r-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tp-r-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

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 \(http://zetcode.com/python/prettytable/\)](http://zetcode.com/python/prettytable/)

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	7	0.78
TFIDF	Brute	12	0.79
W2V	Brute	10	0.78
TFIDFW2V	Brute	6	0.78

2. Naive Bayes

Some code blocks are taken from previous assignments. And some used the code present in original file ('4_DonorsChoose_NB.ipynb') which is mentioned in

comments.

Following Code blocks provided by me.

Adding a column `summary_numeric_bool` instead of `project_resource_summary` column which tells if resource summary has a number in it

```
In [50]: # ref: https://stackoverflow.com/questions/4138202/using-isdigit-for-floats
def nums_in_str(text):
    """
    Returns list of numbers present in the given string. Numbers := floats ints etc.
    """
    result = []
    for s in text.split():
        try:
            x = float(s)
            result.append(x)
        except:
            continue
    return result
```

```
In [51]: print(nums_in_str('HE44Llo 56 are -89 I 820.353 in -78.39 what .293 about 00'))

[56.0, -89.0, 820.353, -78.39, 0.293, 0.0]
```

```
In [52]: numbers_in_summary = np.array([len(nums_in_str(s)) for s in project_data['project_resource_
project_data['summary_numeric_bool'] = list(map(int, numbers_in_summary>0))
```

Taking Relevant columns as X (input data to model) and y (output class label)

In [53]: `project_data.columns`

Out[53]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
 'project_submitted_datetime', 'project_grade_category', 'project_title',
 'project_essay_1', 'project_essay_2', 'project_essay_3',
 'project_essay_4', 'project_resource_summary',
 'teacher_number_of_previously_posted_projects', 'project_is_approved',
 'clean_categories', 'clean_subcategories', 'essay', 'price', 'quantity',
 'summary_numeric_bool'],
 dtype='object')

In [54]: `project_data.head(2)`

Out[54]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datet
--	---------------	----	------------	----------------	--------------	-------------------------

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:40
---	--------	---------	----------------------------------	------	----	------------------

1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:20
---	--------	---------	---------------------------------	-----	----	------------------

2 rows × 21 columns



```
In [55]: # Categorical and numerical columns are listed below.
X_columns = ['teacher_prefix', 'school_state', 'project_grade_category', 'summary_numeric_b
            'teacher_number_of_previously_posted_projects', 'clean_categories', 'clean_sub
            'price', 'quantity']
X = project_data[X_columns]
y = project_data['project_is_approved']
```

Adding preprocessed_essays and preprocessed_titles as columns to X before splitting

```
In [56]: X['essay'] = preprocessed_essays
X['project_title'] = preprocessed_titles
X_columns.append('essay')
X_columns.append('project_title')
print('final columns used in input data are: ', X_columns)
```

final columns used in input data are: ['teacher_prefix', 'school_state', 'project_grade_category', 'summary_numeric_bool', 'teacher_number_of_previously_posted_projects', 'clean_categories', 'clean_subcategories', 'price', 'quantity', 'essay', 'project_title']

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

```
In [57]: # 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 cod
# 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
```

Not creating CV data as I am using K-fold validation

```
In [58]: # Code took from SAMPLE_SOLUTION notebook
# splitting into 80-20 ratio for train-test data
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, stratify=y)
```

```
In [59]: print(X_train.shape)
print(X_test.shape)
print('='*30)
print(y_train.shape)
print(y_test.shape)
```

```
(87398, 11)
```

```
(21850, 11)
```

```
=====
```

```
(87398,)
```

```
(21850,)
```

2.2 Make Data Model Ready: encoding numerical, categorical features

```
In [60]: # 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 cod
# 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
```

numerical columns

- teacher_number_of_previously_posted_projects
- price
- quantity

Leaving summary_numeric_bool as it is because it only has 0's and 1's in it.

categorical columns

- teacher_prefix
- school_state
- project_grade_category
- clean_categories
- clean_subcategories

Normalizing teacher_number_of_previously_posted_projects column

```
In [61]: warnings.filterwarnings("ignore")
# Code took from original Code provided.
scaler = StandardScaler()
scaler.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
print(f"Mean : {scaler.mean_[0]}, Standard deviation : {np.sqrt(scaler.var_[0])}")
```

Mean : 11.180107096272225, Standard deviation : 27.699906340212106

```
In [62]: warnings.filterwarnings("ignore")
X_train_tnppp_norm = scaler.transform(X_train['teacher_number_of_previously_posted_projects'])
X_test_tnppp_norm = scaler.transform(X_test['teacher_number_of_previously_posted_projects'])
```

Normalizing price column

```
In [63]: # Code took from original Code provided.
scaler = StandardScaler()
scaler.fit(X_train['price'].values.reshape(-1,1))
print(f"Mean : {scaler.mean_[0]}, Standard deviation : {np.sqrt(scaler.var_[0])}")
```

Mean : 297.4589361312616, Standard deviation : 366.1637483541025

```
In [64]: X_train_price_norm = scaler.transform(X_train['price'].values.reshape(-1,1))
X_test_price_norm = scaler.transform(X_test['price'].values.reshape(-1,1))
```

Normalizing quantity column

```
In [65]: warnings.filterwarnings("ignore")
# Code took from original Code provided.
scaler = StandardScaler()
scaler.fit(X_train['quantity'].values.reshape(-1,1))
print(f"Mean : {scaler.mean_[0]}, Standard deviation : {np.sqrt(scaler.var_[0])}")
```

Mean : 16.974873566900843, Standard deviation : 26.096669192815668

```
In [66]: warnings.filterwarnings("ignore")
X_train_quant_norm = scaler.transform(X_train['quantity'].values.reshape(-1,1))
X_test_quant_norm = scaler.transform(X_test['quantity'].values.reshape(-1,1))
```

Using a array to store column names data to use at last when interpreting the model

```
In [67]: # when combining the input matrix the order of columns is same as cat_num_columns
cat_num_columns = ['previously_posted_projects', 'price', 'quantity', 'summary_numeric_bool
```

Encoding teacher_prefix column

```
In [68]: # Code took from SAMPLE_SOLUTION notebook.
vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values)
print(vectorizer.get_feature_names())
```

['dr', 'mr', 'mrs', 'ms', 'teacher']

```
In [69]: # Code took from SAMPLE_SOLUTION notebook.  
X_train_prefix_ohe = vectorizer.transform(X_train['teacher_prefix'].values)  
X_test_prefix_ohe = vectorizer.transform(X_test['teacher_prefix'].values)  
  
print(X_train_prefix_ohe.shape, y_train.shape)  
print(X_test_prefix_ohe.shape, y_test.shape)  
  
(87398, 5) (87398,)  
(21850, 5) (21850,)
```

```
In [70]: cat_num_columns.extend(['prefix_'+i for i in vectorizer.get_feature_names()])
```

Encoding school_state column

```
In [71]: # Code took from SAMPLE_SOLUTION notebook.  
vectorizer = CountVectorizer()  
vectorizer.fit(X_train['school_state'].values)  
print(vectorizer.get_feature_names())  
  
['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']
```



```
In [72]: # Code took from SAMPLE_SOLUTION notebook.
X_train_school_ohe = vectorizer.transform(X_train['school_state'].values)
X_test_school_ohe = vectorizer.transform(X_test['school_state'].values)

print(X_train_school_ohe.shape, y_train.shape)
print(X_test_school_ohe.shape, y_test.shape)

(87398, 51) (87398,)
(21850, 51) (21850,)
```

```
In [73]: cat_num_columns.extend(['state_'+i for i in vectorizer.get_feature_names()])
print(len(cat_num_columns))

60
```

Encoding project_grade_category column

```
In [74]: # Code took from original Code provided.
grades = X_train['project_grade_category'].unique()
vectorizer = CountVectorizer(vocabulary=list(grades), lowercase=False, binary=True)
vectorizer.fit(X_train['project_grade_category'].values)
print(vectorizer.get_feature_names())

['Grades PreK-2', 'Grades 6-8', 'Grades 3-5', 'Grades 9-12']
```

```
In [75]: # Code took from SAMPLE_SOLUTION notebook.
X_train_grade_ohe = vectorizer.transform(X_train['project_grade_category'].values)
X_test_grade_ohe = vectorizer.transform(X_test['project_grade_category'].values)

print(X_train_grade_ohe.shape, y_train.shape)
print(X_test_grade_ohe.shape, y_test.shape)

(87398, 4) (87398,)
(21850, 4) (21850,)
```

```
In [76]: cat_num_columns.extend(vectorizer.get_feature_names())
print(len(cat_num_columns))
```

64

Encoding clean_categories column

```
In [77]: # Code took from original Code provided.
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binarize=False)
vectorizer.fit(X_train['clean_categories'].values)
print(vectorizer.get_feature_names())

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

```
In [78]: # Code took from SAMPLE_SOLUTION notebook.
X_train_categ_ohe = vectorizer.transform(X_train['clean_categories'].values)
X_test_categ_ohe = vectorizer.transform(X_test['clean_categories'].values)

print(X_train_categ_ohe.shape, y_train.shape)
print(X_test_categ_ohe.shape, y_test.shape)

(87398, 9) (87398,)
(21850, 9) (21850,)
```

```
In [79]: cat_num_columns.extend(['categ_'+i for i in vectorizer.get_feature_names()])
print(len(cat_num_columns))

73
```

Encoding clean_subcategories column

```
In [80]: # Code took from original Code provided.
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False,
vectorizer.fit(X_train['clean_subcategories'].values)
print(vectorizer.get_feature_names())

['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricu
lar', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hung
er', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'Co
llege_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopmen
t', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'Appli
edSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
```

```
In [81]: # Code took from SAMPLE_SOLUTION notebook.
X_train_subcat_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
X_test_subcat_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print(X_train_subcat_ohe.shape, y_train.shape)
print(X_test_subcat_ohe.shape, y_test.shape)

(87398, 30) (87398,)
(21850, 30) (21850,)
```

```
In [82]: cat_num_columns.extend(['subcateg_'+i for i in vectorizer.get_feature_names()])
print(len(cat_num_columns))

103
```

MultinomialNB is not accepting negative input values. So using MinMaxScaler to transform numerical columns to (0, 1) range before combining

```
In [83]: from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
scaler.fit(X_train_tnppp_norm)

X_train_tnppp_norm = scaler.transform(X_train_tnppp_norm)
X_test_tnppp_norm = scaler.transform(X_test_tnppp_norm)
```

```
In [84]: scaler = MinMaxScaler()
scaler.fit(X_train_price_norm)

X_train_price_norm = scaler.transform(X_train_price_norm)
X_test_price_norm = scaler.transform(X_test_price_norm)
```

```
In [85]: scaler = MinMaxScaler()
scaler.fit(X_train_quant_norm)

X_train_quant_norm = scaler.transform(X_train_quant_norm)
X_test_quant_norm = scaler.transform(X_test_quant_norm)
```

Combining categorical and numerical data for further use.

```
In [86]: from scipy.sparse import hstack
cat_num_train = hstack((X_train_tnppp_norm, X_train_price_norm, X_train_quant_norm,\
                        np.array(X_train['summary_numeric_bool']).reshape(-1, 1),\
                        X_train_prefix_ohe, X_train_grade_ohe, X_train_school_ohe, X_train_
cat_num_test = hstack((X_test_tnppp_norm, X_test_price_norm, X_test_quant_norm,\
                        np.array(X_test['summary_numeric_bool']).reshape(-1, 1),\
                        X_test_prefix_ohe, X_test_grade_ohe, X_test_school_ohe, X_test_categ
```

```
In [87]: print(cat_num_train.shape, y_train.shape)
print(cat_num_test.shape, y_test.shape)
```

```
(87398, 103) (87398,)
(21850, 103) (21850,)
```

```
In [88]: print(len(cat_num_columns))
```

```
103
```

2.3 Make Data Model Ready: encoding eassay, and project_title

```
In [89]: # 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 cod
# 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
```

Converting essay column to vector using Bag of Words (BoW).

```
In [90]: # Code took from original Code provided.
vectorizer = CountVectorizer(ngram_range=(1,3), min_df=15, max_features=50000)
vectorizer.fit(X_train['essay'].values)
print(len(vectorizer.get_feature_names()))
```

50000

```
In [91]: # Code took from SAMPLE_SOLUTION notebook.
X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
X_test_essay_bow = vectorizer.transform(X_test['essay'].values)

print(X_train_essay_bow.shape, y_train.shape)
print(X_test_essay_bow.shape, y_test.shape)
```

(87398, 50000) (87398,)
(21850, 50000) (21850,)

```
In [92]: essay_bow_columns = ['essay_'+i for i in vectorizer.get_feature_names()]  
print(len(essay_bow_columns))
```

50000

```
In [93]: import random  
random.sample(essay_bow_columns, 10)
```

```
Out[93]: ['essay_good school',  
          'essay_children express',  
          'essay_safe successful',  
          'essay_kids extremely',  
          'essay_murals',  
          'essay_need change',  
          'essay_students college bound',  
          'essay_love reading math',  
          'essay_population qualifies',  
          'essay_formula']
```

Converting essay column to vector using TFIDF Vectorizer.

```
In [94]: # Code took from original Code provided.  
vectorizer = TfidfVectorizer(ngram_range=(1,3), min_df=15, max_features=50000)  
vectorizer.fit(X_train['essay'].values)  
print(len(vectorizer.get_feature_names()))
```

50000

```
In [95]: # Code took from SAMPLE_SOLUTION notebook.
X_train_essay_tfidf = vectorizer.transform(X_train['essay'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)

print(X_train_essay_tfidf.shape, y_train.shape)
print(X_test_essay_tfidf.shape, y_test.shape)

(87398, 50000) (87398,)
(21850, 50000) (21850,)
```

```
In [96]: essay_tfidf_columns = ['essay_'+i for i in vectorizer.get_feature_names()]
print(len(essay_tfidf_columns))

50000
```

Converting project_title column to vector using Bag of Words (BoW).

```
In [97]: # Code took from original Code provided.
vectorizer = CountVectorizer(ngram_range=(1,3), min_df=10, max_features=10000)
vectorizer.fit(X_train['project_title'].values)
print(len(vectorizer.get_feature_names()))

5664
```



```
In [98]: # Code took from SAMPLE_SOLUTION notebook.
X_train_title_bow = vectorizer.transform(X_train['project_title'].values)
X_test_title_bow = vectorizer.transform(X_test['project_title'].values)

print(X_train_title_bow.shape, y_train.shape)
print(X_test_title_bow.shape, y_test.shape)

(87398, 5664) (87398,)
(21850, 5664) (21850,)
```

```
In [99]: title_bow_columns = ['title_'+i for i in vectorizer.get_feature_names()]
print(len(title_bow_columns))

5664
```

Converting project_title column to vector using TFIDF Vectorizer.

```
In [100]: # Code took from original Code provided.
vectorizer = TfidfVectorizer(ngram_range=(1,3), min_df=10, max_features=10000)
vectorizer.fit(X_train['project_title'].values)
print(len(vectorizer.get_feature_names()))

5664
```

```
In [101]: # Code took from SAMPLE_SOLUTION notebook.  
X_train_title_tfidf = vectorizer.transform(X_train['project_title'].values)  
X_test_title_tfidf = vectorizer.transform(X_test['project_title'].values)  
  
print(X_train_title_tfidf.shape, y_train.shape)  
print(X_test_title_tfidf.shape, y_test.shape)  
  
(87398, 5664) (87398,)  
(21850, 5664) (21850,)
```

```
In [102]: title_tfidf_columns = ['title_'+i for i in vectorizer.get_feature_names()]  
print(len(title_tfidf_columns))  
  
5664
```

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

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

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

Joining processed essay and project_title arrays with categorical and numerical data to form two types of matrices (BoW, TFIDF)

```
In [103]: bow_train = hstack((cat_num_train, X_train_essay_bow, X_train_title_bow)).tocsr()
bow_test = hstack((cat_num_test, X_test_essay_bow, X_test_title_bow)).tocsr()

tfidf_train = hstack((cat_num_train, X_train_essay_tfidf, X_train_title_tfidf)).tocsr()
tfidf_test = hstack((cat_num_test, X_test_essay_tfidf, X_test_title_tfidf)).tocsr()

print('='*30)
print(bow_train.shape)
print(bow_test.shape)
print('='*30)
print(tfidf_train.shape)
print(tfidf_test.shape)
print('='*30)
```

```
=====
(87398, 55767)
(21850, 55767)
=====
(87398, 55767)
(21850, 55767)
=====
```

```
In [104]: bow_columns = cat_num_columns + essay_bow_columns + title_bow_columns
tfidf_columns = cat_num_columns + essay_tfidf_columns + title_tfidf_columns

print(len(bow_columns))
print(len(tfidf_columns))
```

```
55767
55767
```

Want to normalize all data to same range as frequency of words in essays seems to be effecting results due to which we get same important features for both models and for both approved and rejected classes.

```
In [105]: from sklearn.preprocessing import MaxAbsScaler
scaler = MaxAbsScaler()
scaler.fit(bow_train)

bow_train = scaler.transform(bow_train)
bow_test = scaler.transform(bow_test)
```

```
In [106]: scaler = MaxAbsScaler()
scaler.fit(tfidf_train)

tfidf_train = scaler.transform(tfidf_train)
tfidf_test = scaler.transform(tfidf_test)
```

2.4.1 Applying Naive Bayes on BOW, SET 1

```
In [107]: # Please write all the code with proper documentation
```

Writing several functions to reuse them later

Function to plot AUC values with respect to hyper-parameter alpha given train data using K-fold validation

```
In [108]: from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import GridSearchCV
import math

# Code inside function took from SAMPLE_SOLUTION notebook
def auc_vs_K_plot(X_train, y_train, alphas, logplot=True):
    """
    Plots the AUC results for different alpha values on train and CV data
    Parameters:
    X_train, y_train - data which is used for K-fold validation and used to train Multinomi
    alphas - list of alpha values on which we have to train the data and plot the results
    """
    nb_model = MultinomialNB()
    parameters = {'alpha': alphas}
    clf = GridSearchCV(nb_model, parameters, cv=4, scoring='roc_auc')
    clf.fit(X_train, y_train)

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

    plt.figure(figsize=(12, 6))
    if logplot:
        # taking logs of alphas to plot a log-plot
        x_axis_ticks = [math.log10(i) for i in alphas]
    else:
        x_axis_ticks = alphas
    plt.plot(x_axis_ticks, train_auc, label='Train AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
    plt.gca().fill_between(x_axis_ticks, train_auc - train_auc_std, train_auc + train_auc_s

    plt.plot(x_axis_ticks, cv_auc, label='CV AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
```

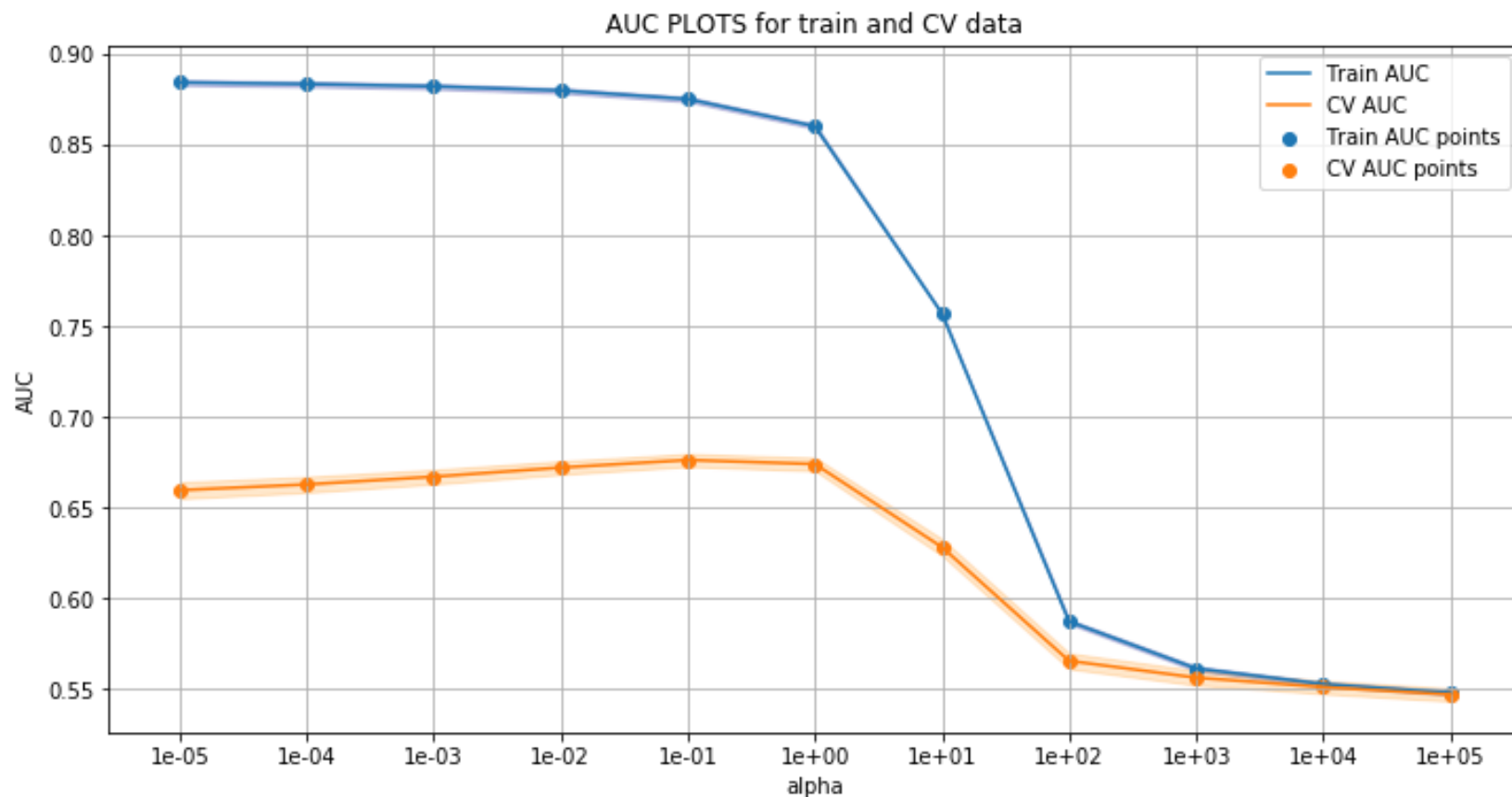
```
plt.gca().fill_between(x_axis_ticks, cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.

plt.scatter(x_axis_ticks, train_auc, label='Train AUC points')
plt.scatter(x_axis_ticks, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("alpha")
# Setting x-ticks to match with actual alphas
if logplot:
    plt.xticks(x_axis_ticks, ["{: .0e}".format(i) for i in alphas])
plt.ylabel("AUC")
plt.title("AUC PLOTS for train and CV data")
plt.grid()
plt.show()
```

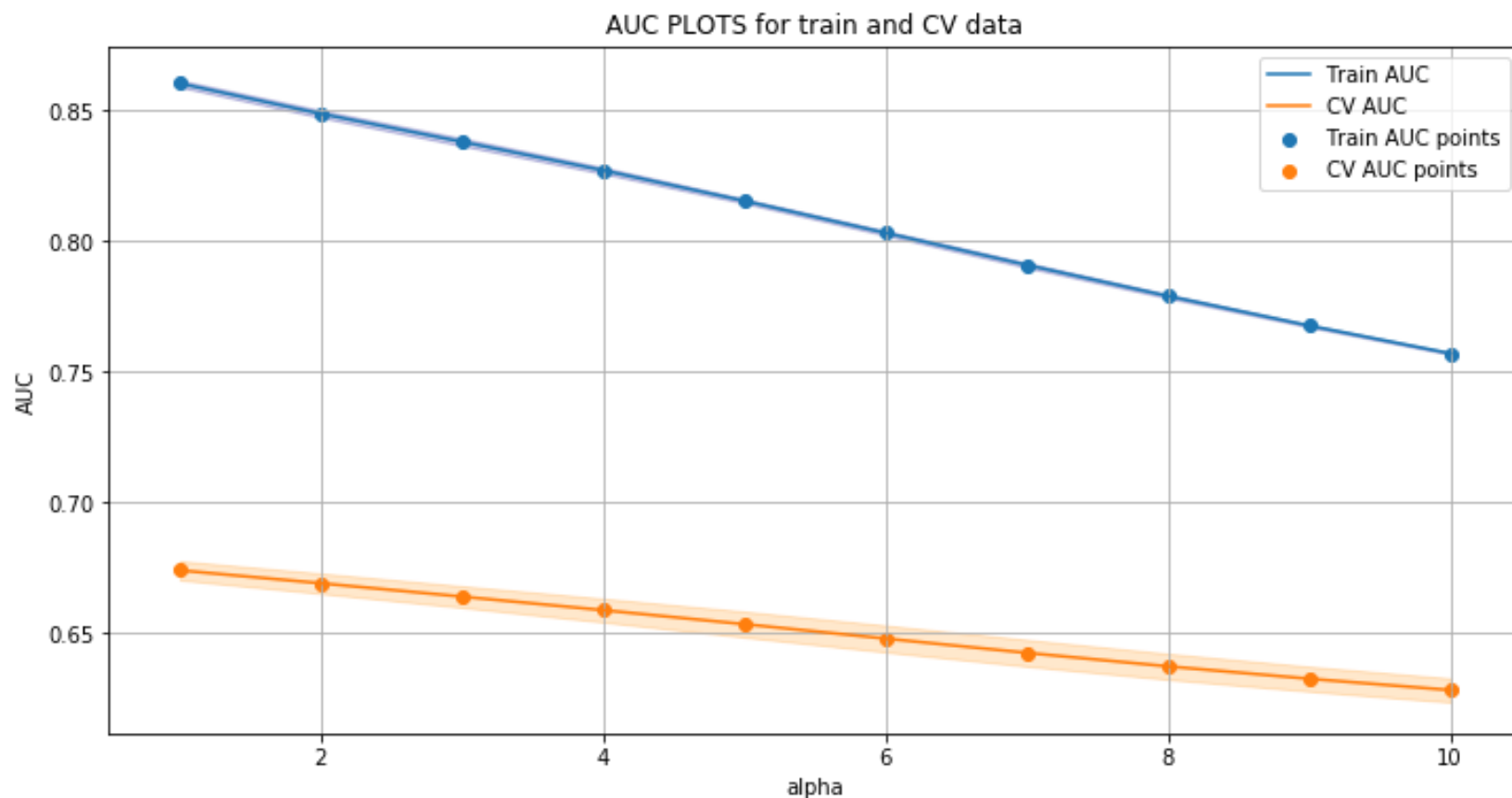
Function to plot ROC curves and print confusion matrices for train and test data. Function returns AUC Values for train, test data and also other values to interpret the model

```
In [109]: alphas = [10**i for i in range(-5, 6)]  
auc_vs_K_plot(bow_train, y_train, alphas, logplot=True)
```



We find maximum AUC for CV data at alpha = 1 (i.e. 1e+00) but the gap between CV and Train AUC values is less at alpha = 10 (i.e. 1e+01) even though AUC value for CV reduced a bit. We can check for other values in between 1 and 10 to see best alpha

```
In [110]: alphas = list(range(1, 11))  
auc_vs_K_plot(bow_train, y_train, alphas, logplot=False)
```



alpha = 6 seems to be fine as gap between AUC's is not decreasing much after 6. I will produce results for alpha = 1 and 10 before creating final table to see which alpha does good with test data. For now alpha = 6 seems to be good as its not that overfitting and have a good (or atleast decent) AUC value.


```

In [120]: from sklearn.metrics import roc_curve, auc, precision_recall_curve
from IPython.display import Markdown, display

# Code inside function took from SAMPLE_SOLUTION notebook
def ROC_conf_mat(X_train, y_train, X_test, y_test, best_alpha, plots = True):
    """
    Plots ROC Curve given a alpha value, Train data and Test data using MultinomialNB as mo
    And also plots confusion matrix for train data and test data taking a optimal threshold
    Returns Area Under ROC Curve for Train, Test data which can be taken as performance of
    and also returns feature_log_prob_, class_counts_ values of the model to interpret the
    """

    # Plotting ROC Curve code
    nb_model = MultinomialNB(alpha = best_alpha)
    nb_model.fit(X_train, y_train)

    y_train_pred = nb_model.predict_proba(X_train)[:, 1]
    y_test_pred = nb_model.predict_proba(X_test)[:, 1]

    train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
    test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

    result = {}

    result['train_auc'], result['test_auc'] = (auc(train_fpr, train_tpr), auc(test_fpr, tes
    result['feat_log_prob'] = nb_model.feature_log_prob_
    result['class_count'] = nb_model.class_count_

    if(plots):
        display(Markdown(f"**Analysis for aplha = {best_alpha}**"))

        plt.plot(train_fpr, train_tpr, label="train AUC =" + str(np.round(result['train_auc']
        plt.plot(test_fpr, test_tpr, label="test AUC =" + str(np.round(result['test_auc'], 3)
        plt.legend()
        plt.xlabel("False Positive rate")
        plt.ylabel("True Positive rate")

```

```
plt.title("ROC Curves for Train and Test data")
plt.grid()
plt.show()

# Printing confusion matrices code
# using precision_recall_curve to get f1_scores to get best threshold
thr_train = tr_thresholds[np.argmax(train_tpr*(1-train_fpr))]
thr_test = te_thresholds[np.argmax(test_tpr*(1-test_fpr))]

print(f"\nConfusion matrix for Train data with {thr_train} as threshold:")
predictions = []
for i in y_train_pred:
    if i >= thr_train:
        predictions.append(1)
    else:
        predictions.append(0)
ax = sns.heatmap(confusion_matrix(y_train, predictions), annot=True, fmt='g')
ax.set_yticklabels(['Rejected', 'Accepted'])
ax.set_xticklabels(['Rejected', 'Accepted'])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title('Confusion matrix for Train')
plt.show()

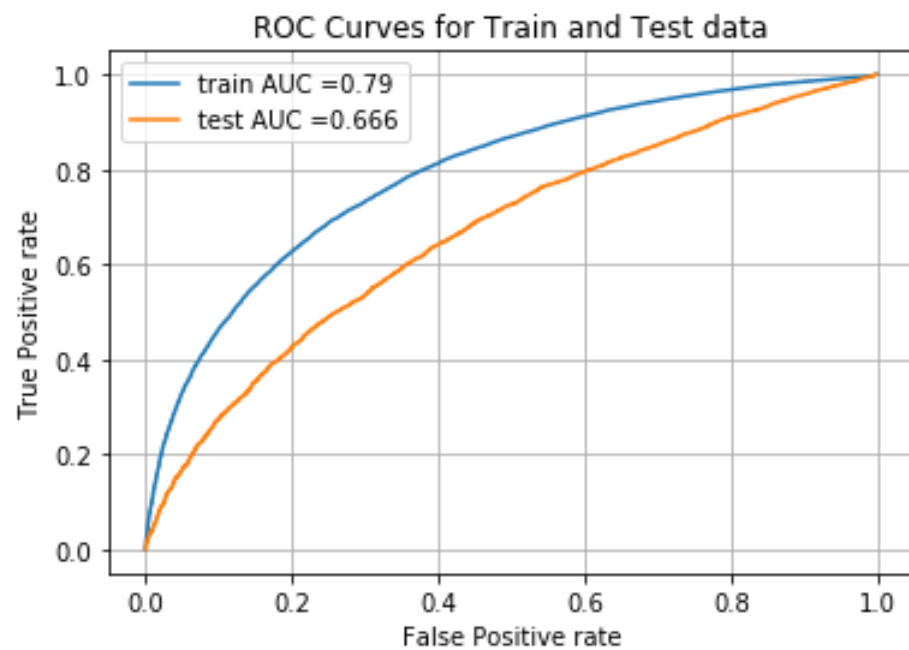
print(f"\nConfusion matrix for Test data with {thr_test} as threshold:")
predictions = []
for i in y_test_pred:
    if i >= thr_test:
        predictions.append(1)
    else:
        predictions.append(0)
ax = sns.heatmap(confusion_matrix(y_test, predictions), annot=True, fmt='g')
ax.set_yticklabels(['Rejected', 'Accepted'])
ax.set_xticklabels(['Rejected', 'Accepted'])
plt.xlabel("Predicted")
plt.ylabel("Actual")
```

```
plt.title('Confusion matrix for Test')  
plt.show()  
  
return result
```

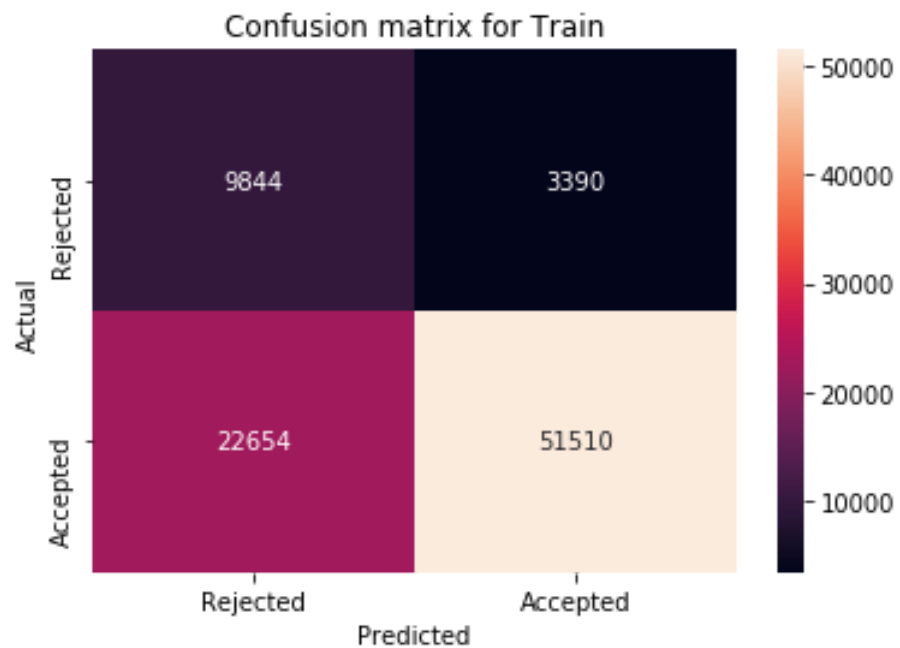
```
In [121]: bow_result = {}
```

```
In [122]: bow_result[6] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 6)
```

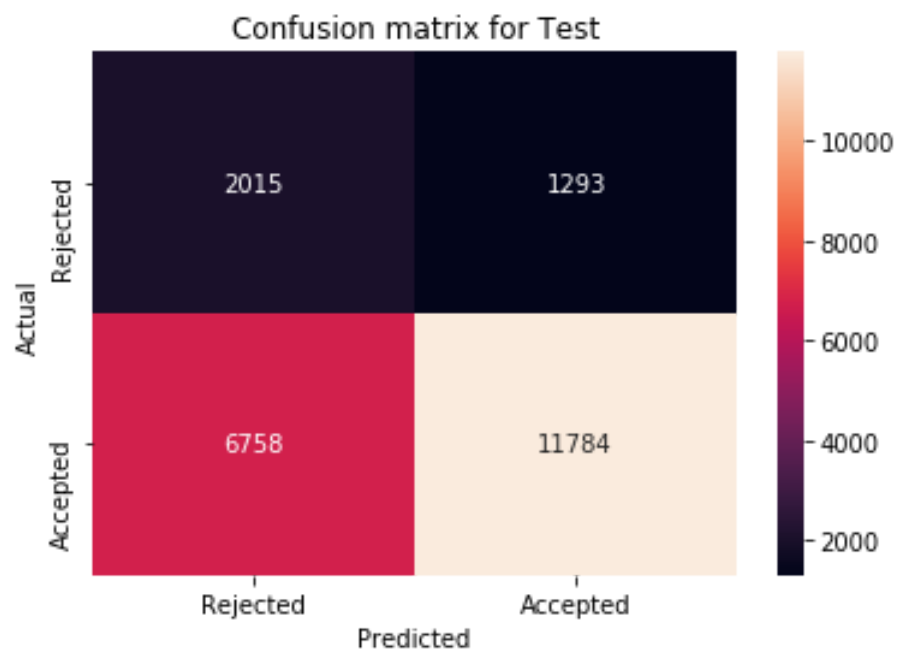
Analysis for alpha = 6



Confusion matrix for Train data with 0.9840810928255197 as threshold:



Confusion matrix for Test data with 0.9886345748476552 as threshold:



Thresholds seems to be so biased towards one class. To construct confusion matrix we consider only FPRs and TPRs so we see high difference in Positively predicted classes which is good as less number of predictions are in False Positive region and more number of predictions are in True Positive region. But if you see Negatively predicted classes there is lot of error. i.e. there are more points in False Negative region than in True Negative region. This is because we consider only FPRs and TPRs to get the threshold for confusion matrix.

2.4.1.1 Top 10 important features of positive class from SET 1

```
In [123]: # Please write all the code with proper documentation
```

```
In [124]: feat_log_prob = bow_result[6]['feat_log_prob']  
print(feat_log_prob.shape)
```

```
(2, 55767)
```

To determine which row corresponds to which class (0 or 1) we use class_count_ value in the result. Because we know positive classes have high number of data points

```
In [125]: print(bow_result[6]['class_count'])
```

```
[13234. 74164.]
```

So first row corresponds to probabilities of $P(f_i|y=0)$ and second row corresponding to $P(f_i|y=1)$ as we know there are more data points with $y=1$ than $y=0$

We find top 10 features which has high values in row 2 to get the feature names

```
In [126]: # Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-indices
indices = sorted(range(len(feat_log_prob[1])), key = lambda i: feat_log_prob[1][i], reverse=True)
print(indices)
```

```
[6, 72, 71, 7, 102, 40264, 101, 100, 17, 35924]
```

```
In [127]: best_feats_for_accepted = [bow_columns[i] for i in indices]
print(best_feats_for_accepted)
```

```
['prefix_mrs', 'categ_Literacy_Language', 'categ_Math_Science', 'prefix_ms', 'subcateg_Literacy', 'essay_students', 'subcateg_Mathematics', 'subcateg_Literature_Writing', 'state_de', 'essay_school']
```

Here we see some features for approving our results. When data is not normalized I got all columns from single words in essay which might be caused due to their high values in the data.

2.4.1.2 Top 10 important features of negative class from SET 1

```
In [128]: # Please write all the code with proper documentation
```

```
In [129]: # Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-indices
indices = sorted(range(len(feat_log_prob[0])), key = lambda i: feat_log_prob[0][i], reverse=True)
print(indices)
```

```
[6, 72, 71, 7, 101, 102, 40264, 100, 35924, 69]
```

```
In [130]: best_feats_for_rejected = [bow_columns[i] for i in indices]
print(best_feats_for_rejected)

['prefix_mrs', 'categ_Literacy_Language', 'categ_Math_Science', 'prefix_ms', 'subcateg_Ma
thematics', 'subcateg_Literacy', 'essay_students', 'subcateg_Literature_Writing', 'essay_
school', 'categ_SpecialNeeds']
```

All the top important words in accepted and rejected classes are same except for last ones. This may be due to high frequency of these features making them highly important. And most of the Negatively predicted data points are from Positive class (i.e. High False Negative Rate) So these features are same as the positive ones. So unbalance in data effecting our model a lot resulting in bad AUC result also

So Getting important words for classes in different way.

Getting Important words by taking difference between row 0 and row 1

```
In [131]: diff_bw_rows = feat_log_prob[0]-feat_log_prob[1]
```

If difference is less then we get highly important words in positive class and also less important in negative class. Similarly opposite is true.

```
In [132]: indices = sorted(range(len(diff_bw_rows)), key = lambda i: diff_bw_rows[i])
print(indices[:10])
print(indices[-10:])
```

```
[54341, 52451, 48580, 52287, 6879, 48582, 39736, 39737, 6184, 18691]
[55032, 55614, 47446, 23184, 51229, 50963, 51626, 52046, 52262, 53783]
```



```
In [133]: best_feats_for_accepted = [bow_columns[i] for i in indices[:10]]  
print('Best features for acceptance:', best_feats_for_accepted)
```

Best features for acceptance: ['title_rug', 'title_hokki stools', 'essay_wobble', 'title_headphones', 'essay_chromebooks', 'essay_wobble chairs', 'essay_stools', 'essay_stools al low', 'essay_chairs allow', 'essay_headphones']

```
In [134]: best_feats_for_rejected = [bow_columns[i] for i in indices[-10:]]  
print('Best features for rejection:', best_feats_for_rejected)
```

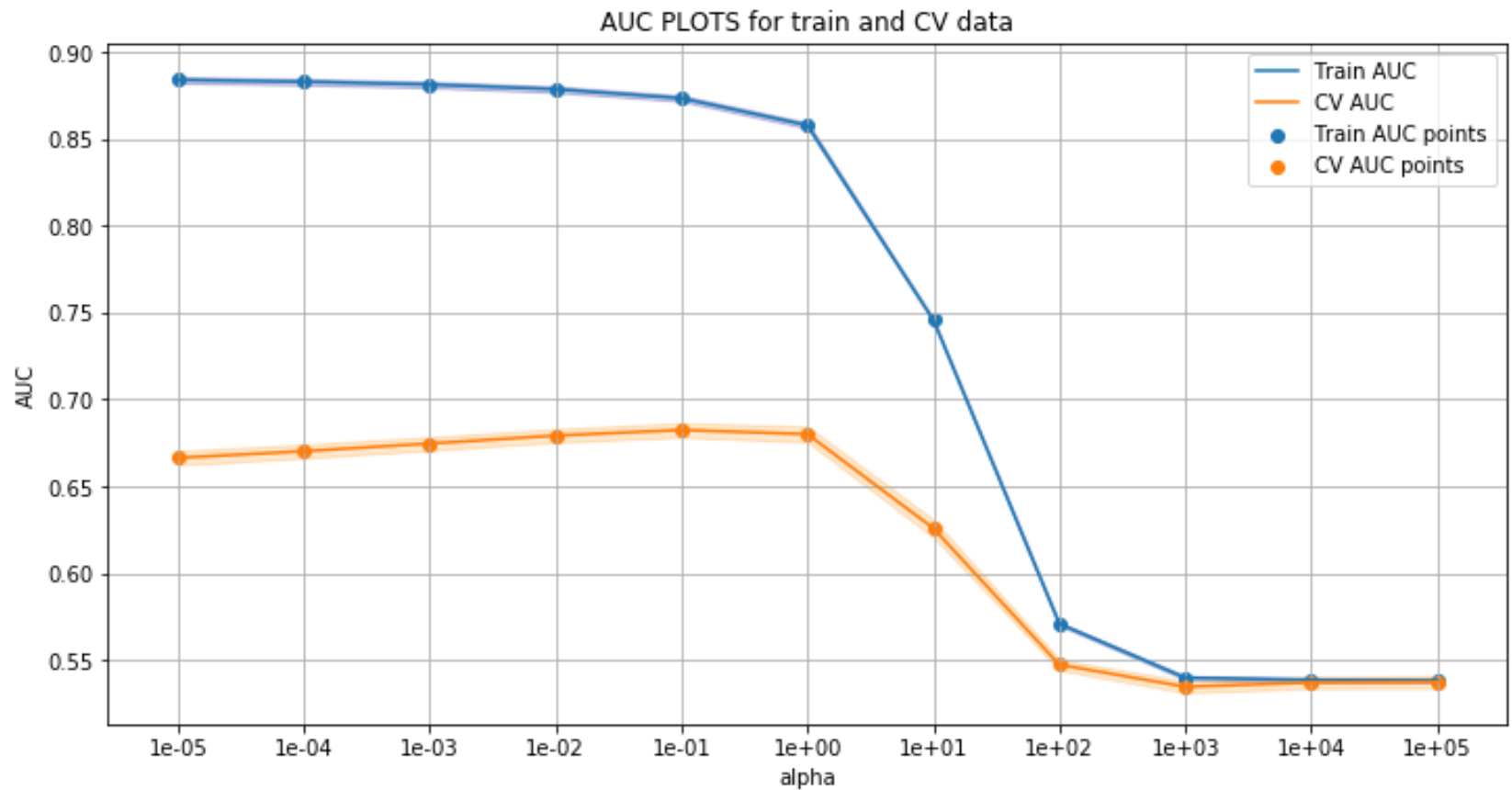
Best features for rejection: ['title_supplies learning', 'title_wish list', 'essay_visual hands learners', 'essay_learners visual', 'title_creating student', 'title_class supplies', 'title_entrepreneurship', 'title_germ', 'title_hands projects', 'title_outdoor fun']

These features may not be sensible as some may have low frequency and mostly present in one class. So these results also may not be good for interpreting and finding importance of words. lets see TFIDF results and hope they are good.

2.4.2 Applying Naive Bayes on TFIDF, SET 2

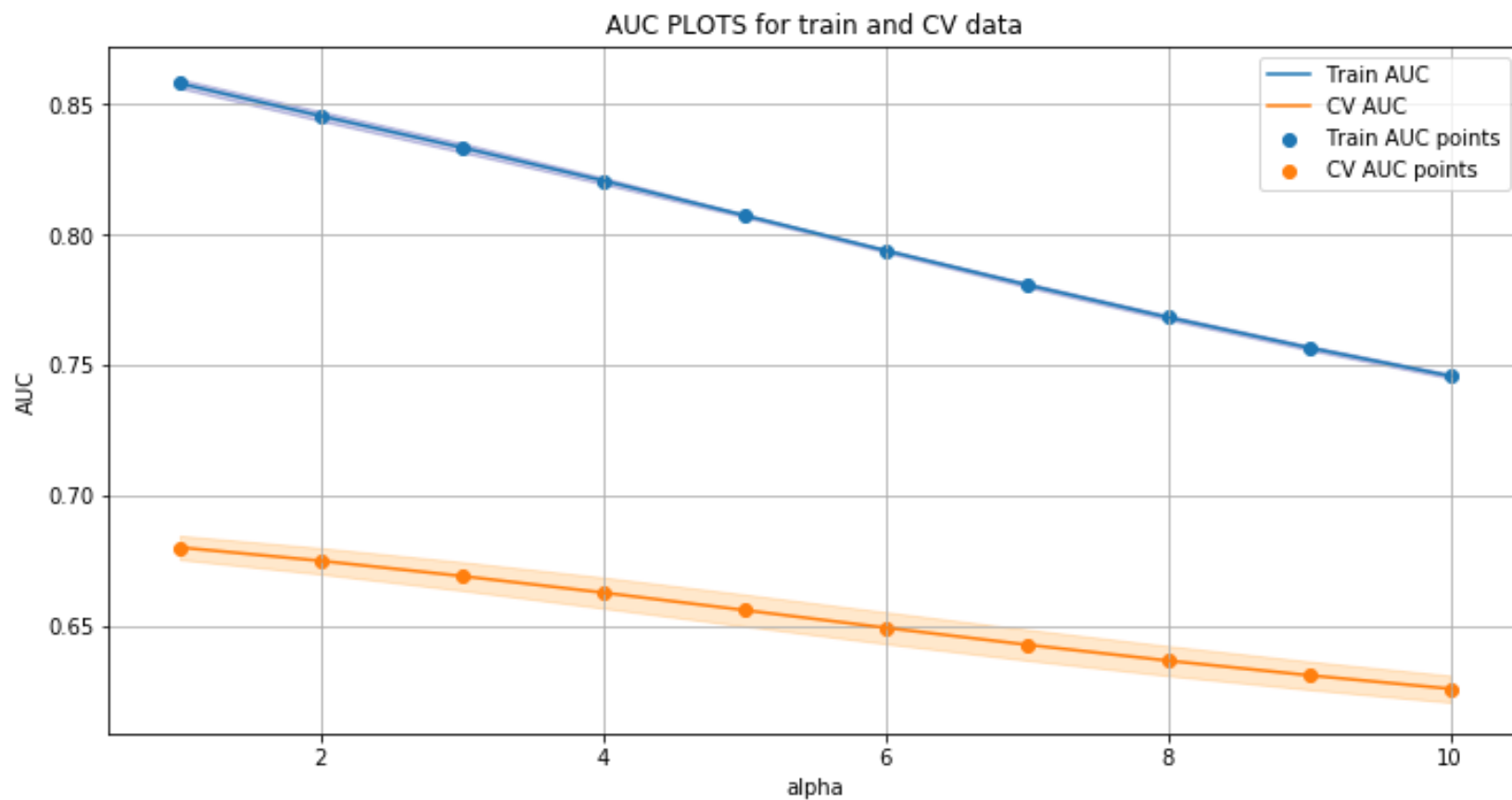
```
In [135]: # Please write all the code with proper documentation
```

```
In [136]: alphas = [10**i for i in range(-5, 6)]  
auc_vs_K_plot(tfidf_train, y_train, alphas, logplot=True)
```



Again searching for alpha in range [1, 10] to see a best alpha. And again the test data is run against other alphas (1,10) to compare results at last.

```
In [137]: alphas = list(range(1, 11))  
auc_vs_K_plot(tfidf_train, y_train, alphas, logplot=False)
```

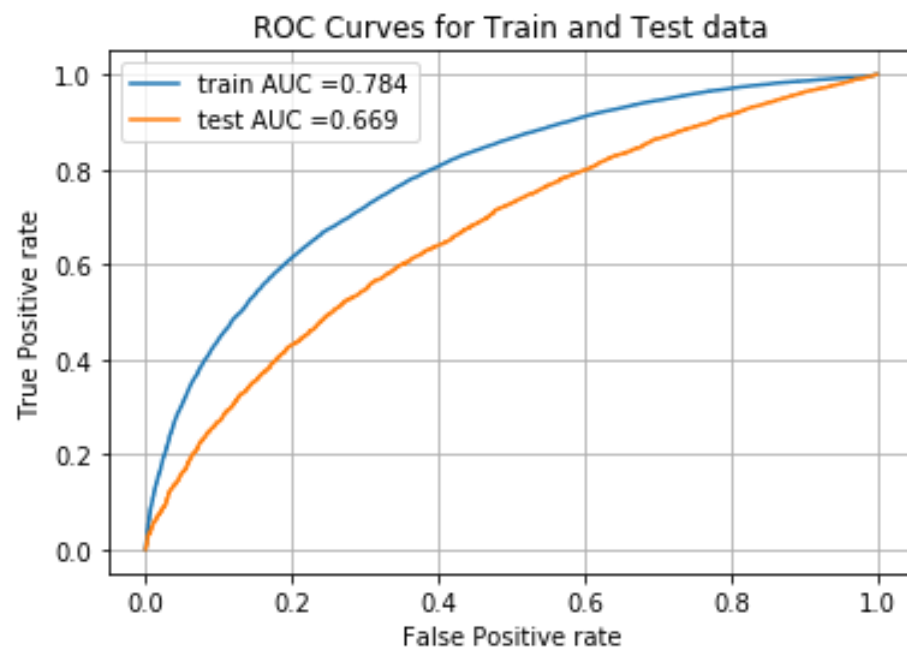


Again taking alpha = 6 as best alpha

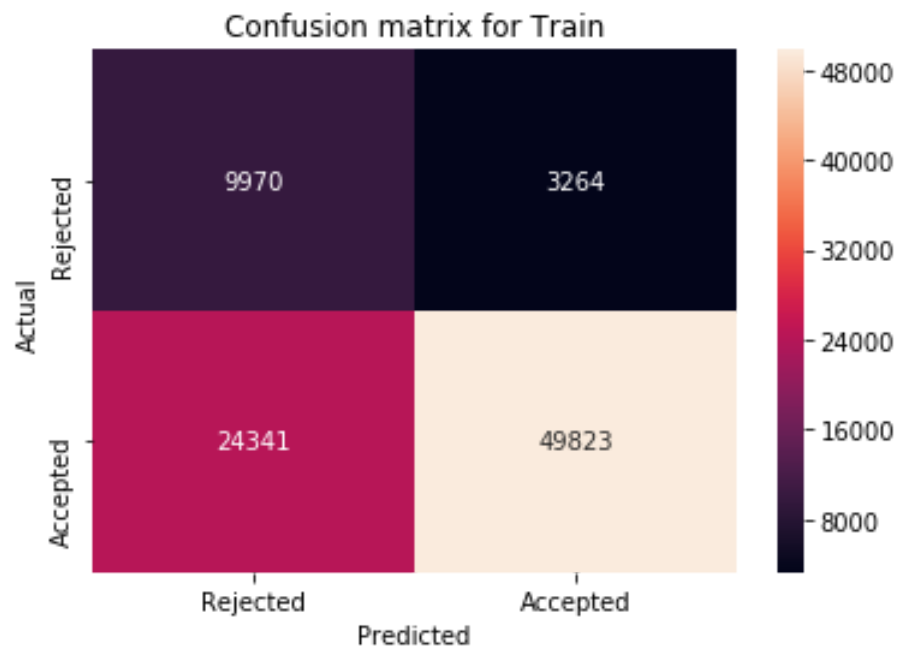
```
In [138]: tfidf_result = {}
```

```
In [139]: tfidf_result[6] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 6)
```

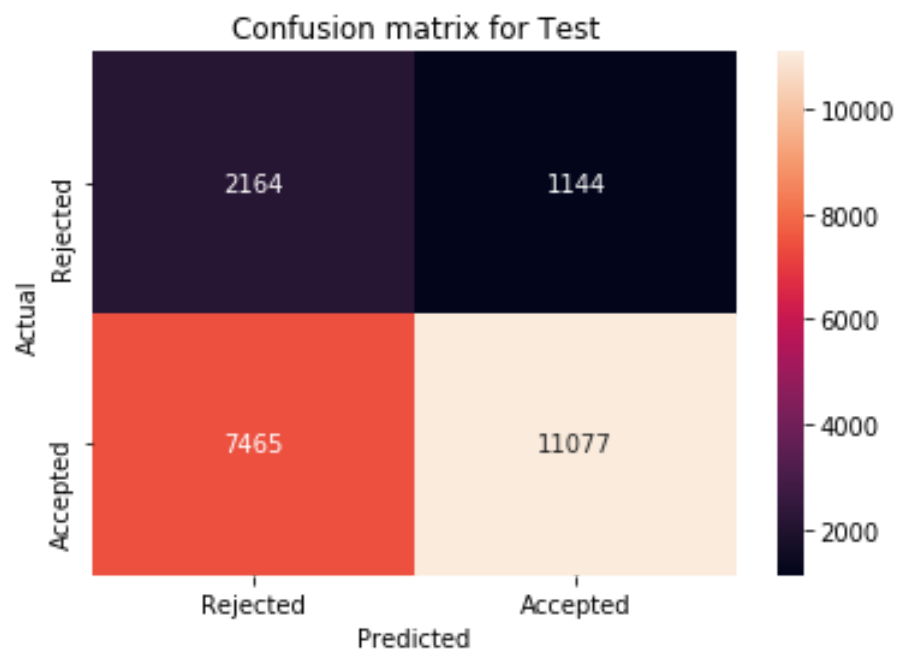
Analysis for alpha = 6



Confusion matrix for Train data with 0.9888149507766784 as threshold:



Confusion matrix for Test data with 0.9931608907738857 as threshold:



2.4.2.1 Top 10 important features of positive class from SET 2

```
In [140]: # Please write all the code with proper documentation
```

These results also may show same important features for both positive and negative classes as there is High False Negative Rate in these results also.

```
In [141]: feat_log_prob = tfidf_result[6]['feat_log_prob']  
print(feat_log_prob.shape)
```

```
(2, 55767)
```

To determine which row corresponds to which class (0 or 1) we use class_count_ value in the result. Because we know positive classes have high number of data points

```
In [142]: print(tfidf_result[6]['class_count'])
```

```
[13234. 74164.]
```

So first row corresponds to probabilities of $P(f_i|y=0)$ and second row corresponding to $P(f_i|y=1)$ as we know there are more data points with $y=1$ than $y=0$

We find top 10 features which has high values in row 2 to get the feature names

```
In [143]: # Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-indices
indices = sorted(range(len(feat_log_prob[1])), key = lambda i: feat_log_prob[1][i], reverse=True)
print(indices)
```

```
[6, 72, 71, 7, 40264, 102, 101, 100, 35924, 17]
```

```
In [144]: best_feats_for_accepted = [tfidf_columns[i] for i in indices]
print(best_feats_for_accepted)
```

```
['prefix_mrs', 'categ_Literacy_Language', 'categ_Math_Science', 'prefix_ms', 'essay_students', 'subcateg_Literacy', 'subcateg_Mathematics', 'subcateg_Literature_Writing', 'essay_school', 'state_de']
```

Here we see got same features as BOW model for approving our results

2.4.2.2 Top 10 important features of negative class from SET 2

```
In [145]: # Please write all the code with proper documentation
```

```
In [146]: # Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-indices
indices = sorted(range(len(feat_log_prob[0])), key = lambda i: feat_log_prob[0][i], reverse=True)
print(indices)
```

```
[6, 72, 71, 7, 40264, 101, 102, 100, 35924, 23195]
```

```
In [147]: best_feats_for_rejected = [tfidf_columns[i] for i in indices]
print(best_feats_for_rejected)
```

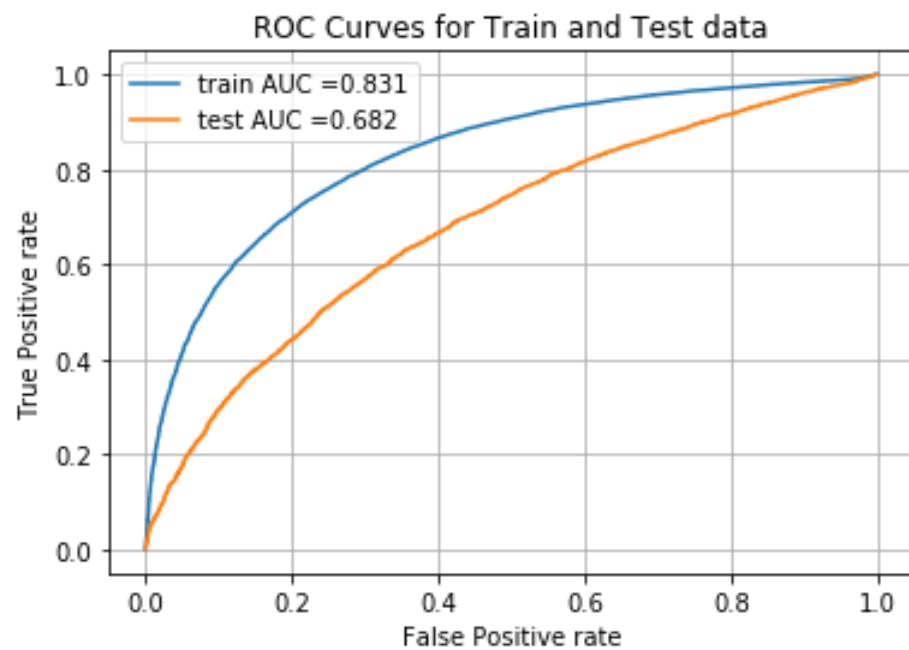
```
['prefix_mrs', 'categ_Literacy_Language', 'categ_Math_Science', 'prefix_ms', 'essay_students', 'subcateg_Mathematics', 'subcateg_Literacy', 'subcateg_Literature_Writing', 'essay_school', 'essay_learning']
```

These are same results for both models. The same reason of unbalance data effecting our model applies here.

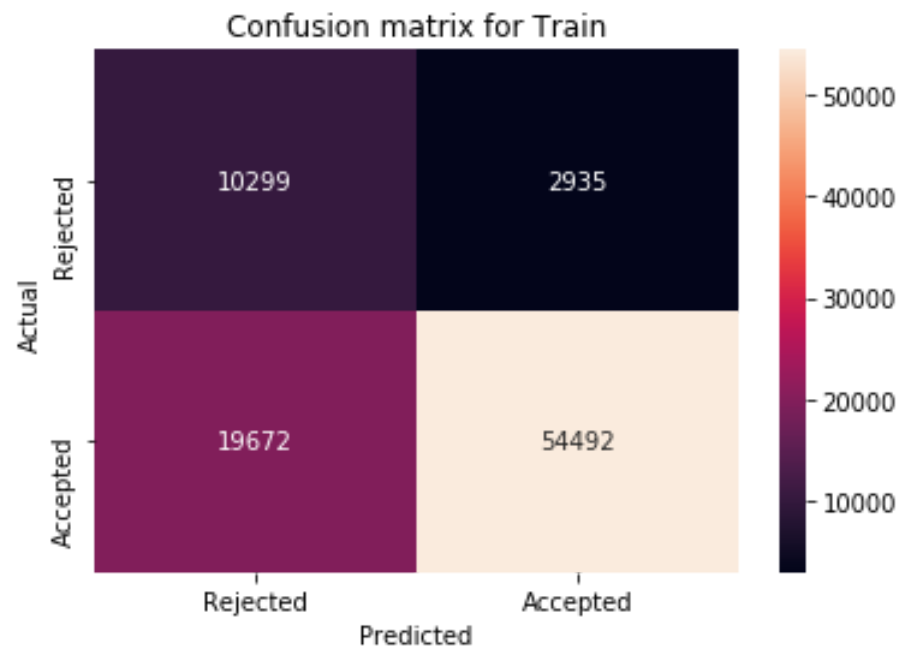
Doing analysis on other alphas for BOW data


```
In [148]: bow_result[1] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 1)
```

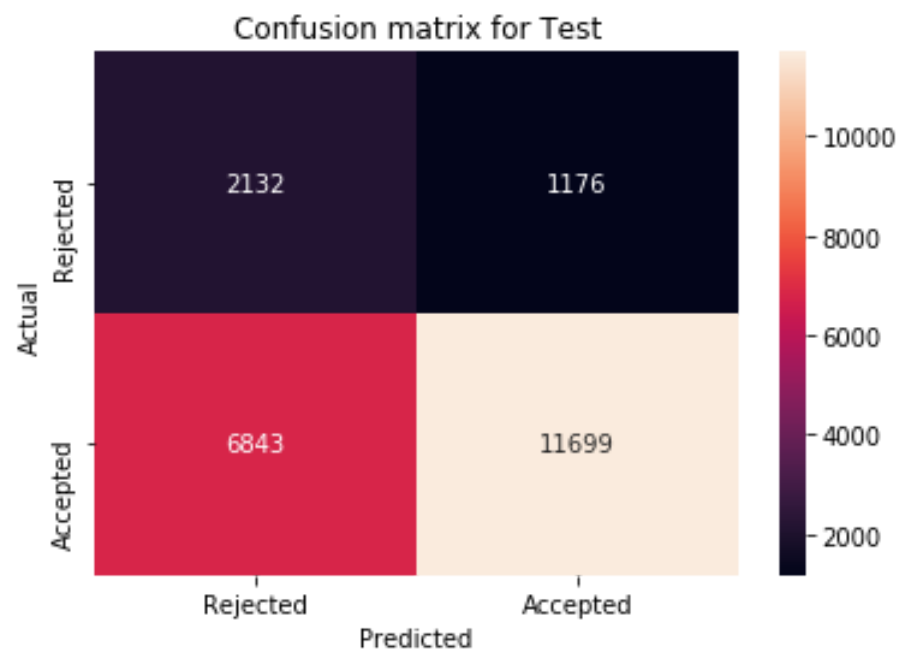
Analysis for alpha = 1



Confusion matrix for Train data with 0.8706226047720625 as threshold:

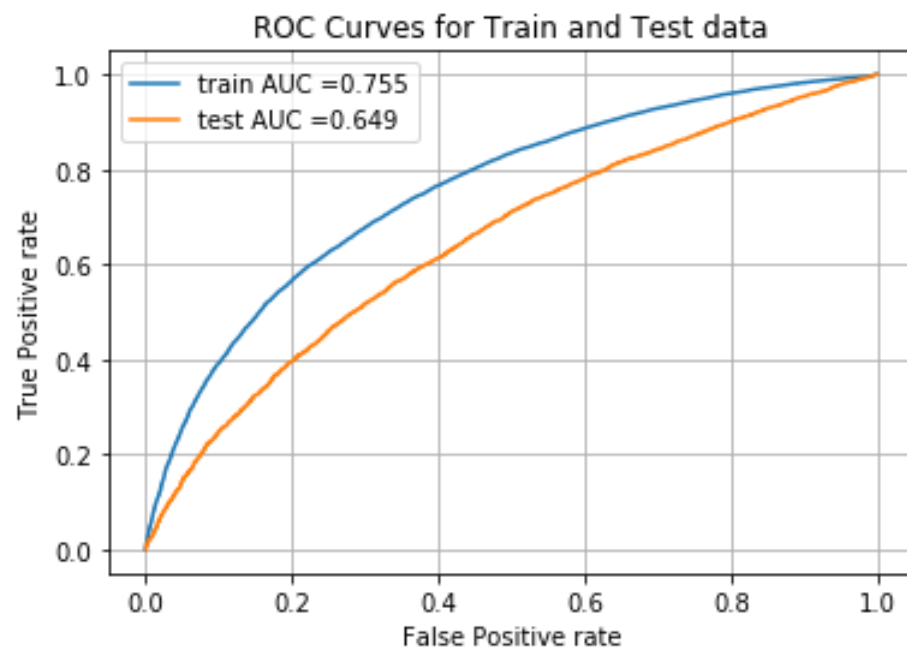


Confusion matrix for Test data with 0.9374107325757828 as threshold:

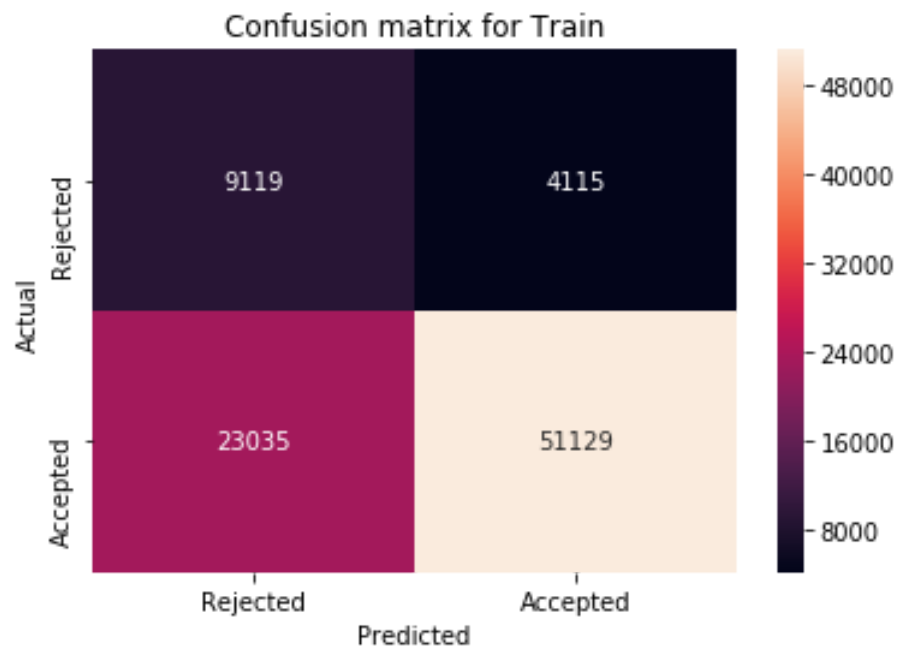



```
In [149]: bow_result[10] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 10)
```

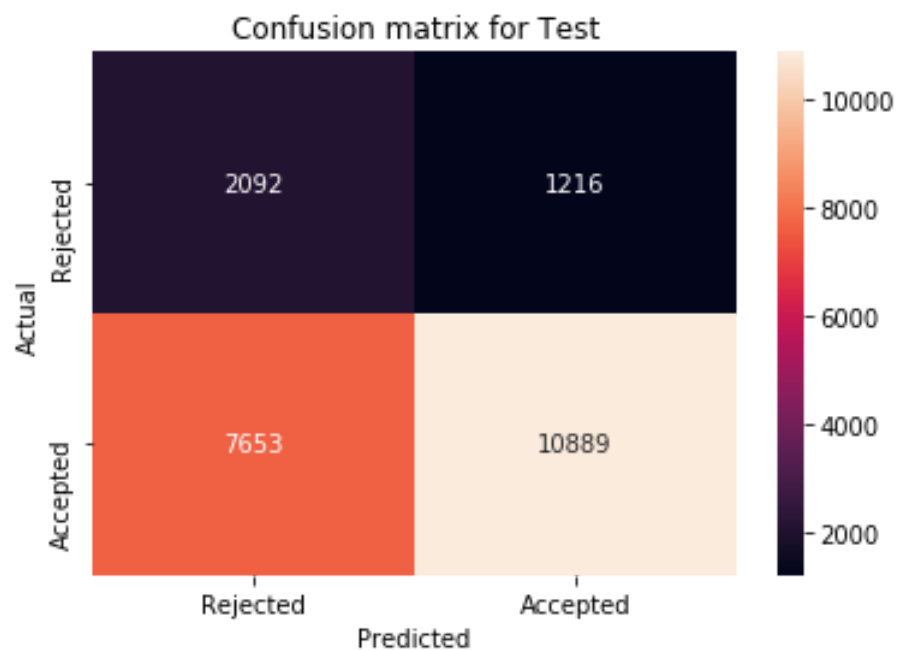
Analysis for alpha = 10



Confusion matrix for Train data with 0.9979513613956511 as threshold:



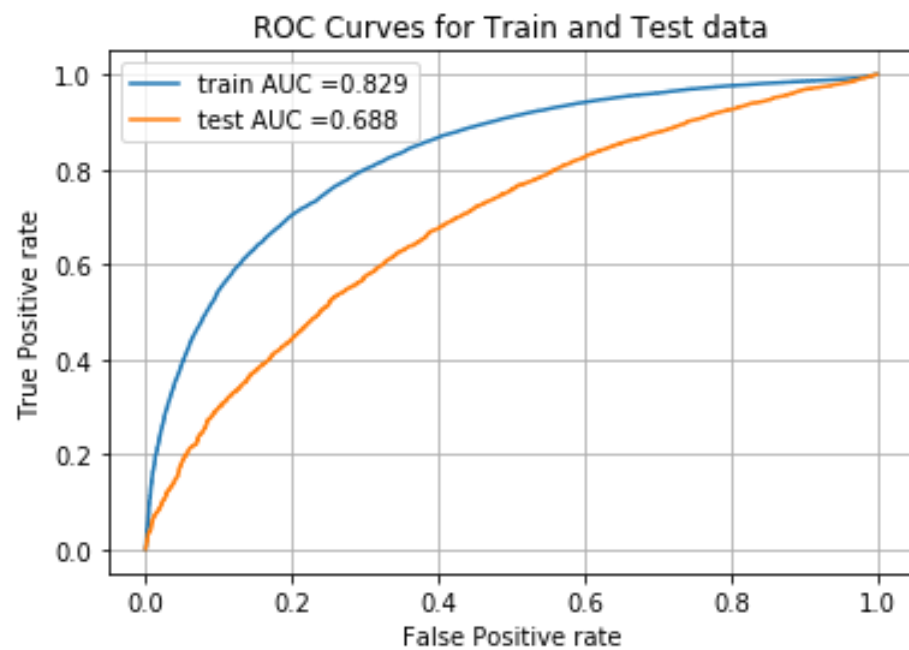
Confusion matrix for Test data with 0.9990785638498187 as threshold:



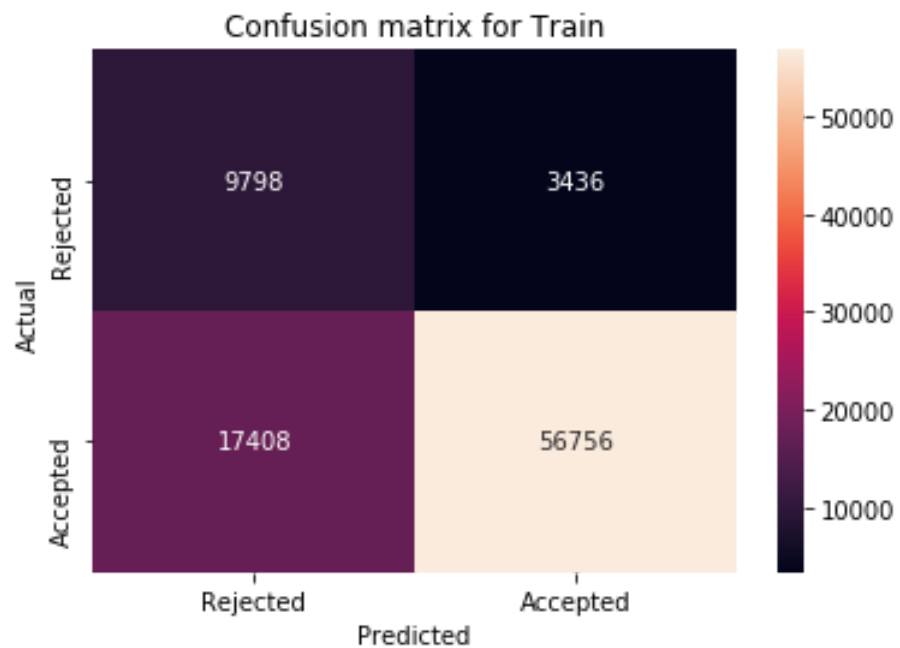
Doing analysis on other alphas for TFIDF data

```
In [150]: tfidf_result[1] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 1)
```

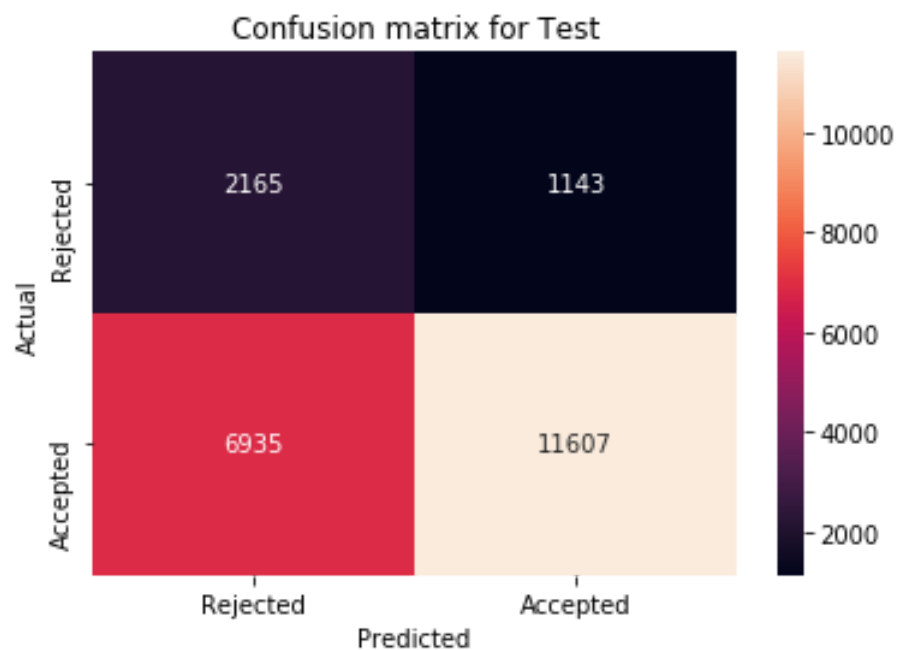
Analysis for alpha = 1



Confusion matrix for Train data with 0.8080632063246167 as threshold:

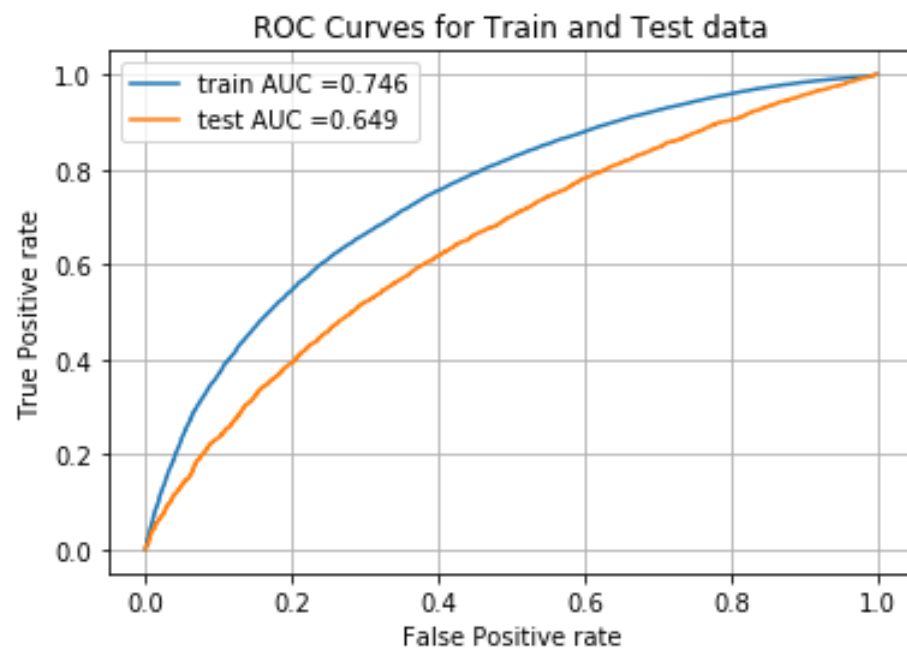


Confusion matrix for Test data with 0.9378109902165341 as threshold:

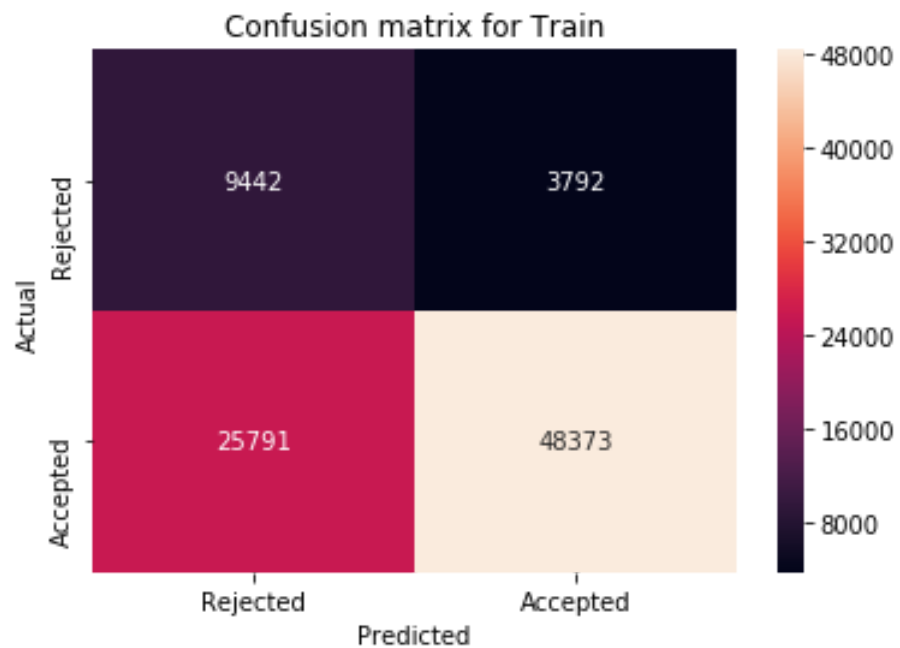



```
In [151]: tfidf_result[10] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 10)
```

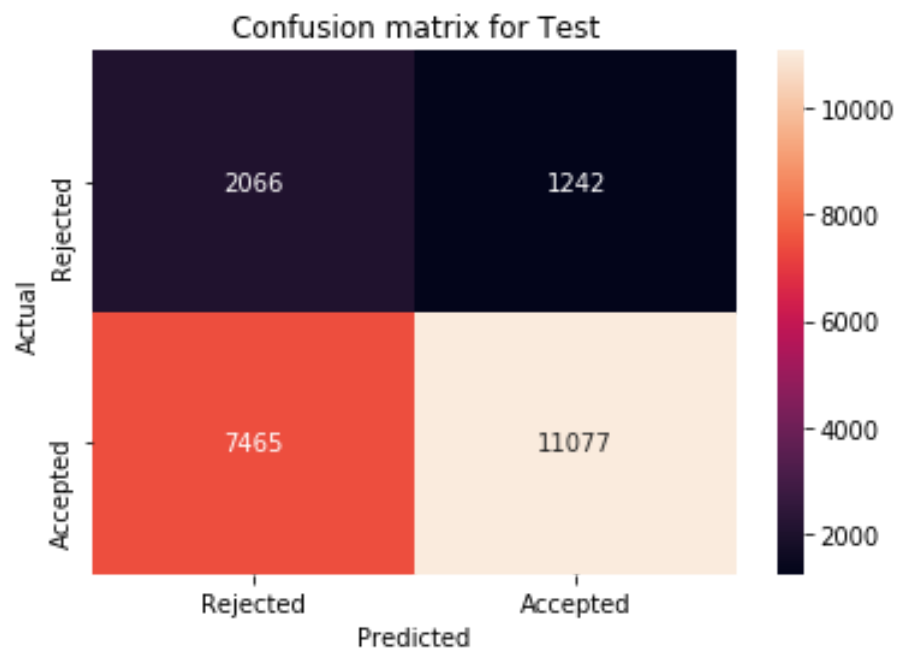
Analysis for alpha = 10



Confusion matrix for Train data with 0.9989699530834106 as threshold:



Confusion matrix for Test data with 0.9992929614556644 as threshold:



Seeing important features with alpha = 1 model as it did better than other alphas.

```
In [152]: feat_log_prob = bow_result[1]['feat_log_prob']
# Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-
indices = sorted(range(len(feat_log_prob[1])), key = lambda i: feat_log_prob[1][i], reverse
best_feats = [bow_columns[i] for i in indices]
print("important features for acceptance acc to BOW:", best_feats)
print("="*125)
indices = sorted(range(len(feat_log_prob[0])), key = lambda i: feat_log_prob[0][i], reverse
best_feats = [bow_columns[i] for i in indices]
print("important features for rejection acc to BOW:", best_feats)
```

```
important features for acceptance acc to BOW: ['prefix_mrs', 'categ_Literacy_Language',
'categ_Math_Science', 'prefix_ms', 'subcateg_Literacy', 'essay_students', 'subcateg_Mathe
matics', 'subcateg_Literature_Writing', 'state_de', 'essay_school']
```

```
=====
important features for rejection acc to BOW: ['prefix_mrs', 'categ_Literacy_Language', 'c
ateg_Math_Science', 'prefix_ms', 'subcateg_Mathematics', 'subcateg_Literacy', 'essay_stud
ents', 'subcateg_Literature_Writing', 'essay_school', 'categ_SpecialNeeds']
```

```
In [153]: feat_log_prob = tfidf_result[1]['feat_log_prob']
# Code took from here: https://stackoverflow.com/questions/13070461/get-index-of-the-top-n-
indices = sorted(range(len(feat_log_prob[1])), key = lambda i: feat_log_prob[1][i], reverse
best_feats = [tfidf_columns[i] for i in indices]
print("important features for acceptance acc to TFIDF:", best_feats)
print("="*125)
indices = sorted(range(len(feat_log_prob[0])), key = lambda i: feat_log_prob[0][i], reverse
best_feats = [tfidf_columns[i] for i in indices]
print("important features for rejection acc to TFIDF:", best_feats)
```

```
important features for acceptance acc to TFIDF: ['prefix_mrs', 'categ_Literacy_Language',
'categ_Math_Science', 'prefix_ms', 'essay_students', 'subcateg_Literacy', 'subcateg_Mathe
matics', 'subcateg_Literature_Writing', 'essay_school', 'state_de']
=====
=====
```

```
important features for rejection acc to TFIDF: ['prefix_mrs', 'categ_Literacy_Language',
'categ_Math_Science', 'prefix_ms', 'essay_students', 'subcateg_Mathematics', 'subcateg_Li
teracy', 'subcateg_Literature_Writing', 'essay_school', 'essay_learning']
```

Above results didnt change much from alpha = 6.

```
In [154]: bow_result[0.1] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 0.1, plots=False)
bow_result[3] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 3, plots=False)
bow_result[8] = ROC_conf_mat(bow_train, y_train, bow_test, y_test, 8, plots=False)
tfidf_result[0.1] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 0.1, plots=False)
tfidf_result[3] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 3, plots=False)
tfidf_result[8] = ROC_conf_mat(tfidf_train, y_train, tfidf_test, y_test, 8, plots=False)
```

3. Conclusions

```
In [155]: # Please compare all your models using Prettytable library
```

```
In [156]: # Sorting results in the result dictionary by the alpha values
bow_result = dict(sorted(bow_result.items()))
tfidf_result = dict(sorted(tfidf_result.items()))
```

```
In [157]: from prettytable import PrettyTable
table = PrettyTable()
table.field_names = ['Vectorizer', 'alpha-HyperParameter', 'Train AUC', 'Test AUC']

for hp, res in bow_result.items():
    table.add_row(['BOW', hp, np.round(res['train_auc'], 3), np.round(res['test_auc'], 3)])
for hp, res in tfidf_result.items():
    table.add_row(['TFIDF', hp, np.round(res['train_auc'], 3), np.round(res['test_auc'], 3)])
print(table)
```

Vectorizer	alpha-HyperParameter	Train AUC	Test AUC
BOW	0.1	0.842	0.684
BOW	1	0.831	0.682
BOW	3	0.814	0.676
BOW	6	0.79	0.666
BOW	8	0.773	0.657
BOW	10	0.755	0.649
TFIDF	0.1	0.841	0.689
TFIDF	1	0.829	0.688
TFIDF	3	0.811	0.682
TFIDF	6	0.784	0.669
TFIDF	8	0.765	0.659
TFIDF	10	0.746	0.649

SUMMARY:

- alpha at 0.1 and 1 have good results (i.e. Test AUC) than other alphas but the Train AUCs are high which may tell us that those are overfitting. And alpha = 6 seems to fine as both Train AUCs is not that high and Test AUC is good enough.
- We didnt get good feature interpretation for negative class, as our model is highly effected by unbalance in data. Which leads to high false negative rate. which means we have so many true positive points predicted as negative which leads to the same important features for both positive and negative classes.
- But we can tell that the important features observed might be important for approval of our project. Some of the features we got are Mrs. and Mr. prefix, Mathematics and Literacy in subject categories and subcategories and 'Students' word in essay

In []: