

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

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

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

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

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

your neighborhood, and your school are all helpful.

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

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

In [1]:

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

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

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

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

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

from tqdm import tqdm
import os

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

1.1 Reading Data

In [2]:

```
project_data = pd.read_csv(r'C:\Users\utsav94\Desktop\train_data.csv')
resource_data = pd.read_csv(r'C:\Users\utsav94\Desktop\resources.csv')
```

In [3]:

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

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

```
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

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

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

Out[4]:

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

1.2 preprocessing of project_subject_categories

In [5]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [6]:

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

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

sub_cat_list = []
for i in sub_categories:
    temp = ""
```

```

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

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

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

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

```

1.3 Text preprocessing

In [7]:

```

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

```

In [8]:

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

Out[8]:

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

In [9]:

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

In [10]:

```

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

```

```

print("\n")
print(project_data['essay'].values[150])
print("\n")
print(project_data['essay'].values[1000])
print("\n")
print(project_data['essay'].values[20000])
print("\n")
print(project_data['essay'].values[99999])
print("\n")

```

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

=====

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

=====

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

=====

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 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 [11]:

```
# 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 [12]:

```
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. nannan

In [13]:

```
# \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', ' ')
```

```
sent = sent.replace('\\\\', '\\')
sent = sent.replace('\\\\n', '\\n')
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 [14]:

```
#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 [15]:

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

In [16]:

[illegible]

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

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

```
# similarly you can preprocess the titles also
# similarly you can preprocess the titles also
# similarly you can preprocess the titles also
print(project_data['project_title'].values[0])
print("="*50)
print(project_data['project_title'].values[150])
print("="*50)
print(project_data['project_title'].values[1000])
print("="*50)
print(project_data['project_title'].values[20000])
print("="*50)
print(project_data['project_title'].values[99999])
print("="*50)
```

```
sent_1 = decontracted(project_data['project_title'].values[500])
print(sent_1)
print("="*50)

preprocessed_titles = []
# tqdm is for printing the status bar
for sentence_1 in tqdm(project_data['project_title'].values):
```

```

for sentence_1 in tqdm(project_data['project_title'].values):
    sent_1 = decontracted(sentence_1)
    sent_1 = sent_1.replace('\\r', ' ')
    sent_1 = sent_1.replace('\\\"', ' ')
    sent_1 = sent_1.replace('\\n', ' ')
    sent_1 = re.sub('[^A-Za-z0-9]+', ' ', sent_1)
    # https://gist.github.com/sebleier/554280
    sent_1 = ' '.join(e for e in sent_1.split() if e not in stopwords)
    preprocessed_titles.append(sent_1.lower().strip())

```

Classroom Chromebooks for College Bound Seniors!

=====

100% |████████████████████████████████████████████████████████████████████████████████| 109248/109248 [00:04<00:00, 24034.91it/s]

1.5 Preparing data for models

In [20]:

```
project_data.columns
```

Out[20]:

```

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/>

In [21]:

```

# 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 [22]:

```
# 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 [23]:

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

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [24]:

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

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

In [25]:

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

1.5.2.2 TFIDF vectorizer

In [26]:

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

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

1.5.2.3 Using Pretrained Models: Avg W2V

In [27]:

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

```

    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

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

# =====

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

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

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

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

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

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

'''

```

Out[27]:

```

'\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        n
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\glove.42B.300d.txt')\n\n# =====\n\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split('\
'))\n\nfor i in preprocod_titles:\n    words.extend(i.split('\ '))\n\nprint("all the words in the
coupus", len(words))\n\nwords = set(words)\n\nprint("the unique words in the coupus",
len(words))\n\ninter_words = set(model.keys()).intersection(words)\n\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 = {}\n\nwords_glove =
set(model.keys())\n\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\n\nwith open(\glove_vectors', \wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'

```

In [28]:

```

# 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 [29]:

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

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

```
100%|████████████████████████████████████████| 109248/109248 [00:52<00:00, 2086.98it/s]
```

```
109248
300
```

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [30]:

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

In [31]:

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

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

```
100%|████████████████████████████████████████| 109248/109248 [07:27<00:00, 243.98it/s]
```

```
109248
300
```

In [32]:

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

1.5.3 Vectorizing Numerical features

In [33]:

```
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 [34]:

```
# 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 = 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. 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 mean and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

Mean : 298.1193425966608, Standard deviation : 367.49634838483496

In [35]:

```
price_standardized
```

Out[35]:

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

1.5.4 Merging all the above features

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

In [36]:

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

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

In [37]:

```
# 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[37]:

```
(109248, 16663)
```

In [38]:

```
# 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

ase = project_data['project_is_approved']
project_data['teacher_prefix'] = project_data['teacher_prefix'].replace(np.nan, 'teacher')
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '')
project_data['project_grade_category']
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace("-", "_")
")
```

In [39]:

```
project_data['preprocessed_essays'] = preprocessed_essays
project_data['preprocessed_titles'] = preprocessed_titles

#https://stackoverflow.com/questions/49984905/count-number-of-words-per-row/49984998
project_data['essay_count']=project_data['preprocessed_essays'].str.split().str.len()
project_data['title_count']=project_data['preprocessed_titles'].str.split().str.len()

#https://www.geeksforgeeks.org/python-textblob-sentiment-method/
#https://textblob.readthedocs.io/en/dev/quickstart.html#quickstart
from textblob import TextBlob
project_data['essay_sentiment'] = [ TextBlob(tb).sentiment.polarity for tb in project_data['essay']
]

project_data
```

Out[39]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetim
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	2016-10-25 09:22:10
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09
...

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetim
109243	38267	p048540	fadf72d6cd83ce6074f9be78a6fcd374	Mr.	MO	2016-06-17 12:02:31
109244	169142	p166281	1984d915cc8b91aa16b4d1e6e39296c6	Ms.	NJ	2017-01-11 12:49:39
109245	143653	p155633	cdbfd04aa041dc6739e9e576b1fb1478	Mrs.	NJ	2016-08-25 17:11:32
109246	164599	p206114	6d5675dbfafa1371f0e2f6f1b716fe2d	Mrs.	NY	2016-07-29 17:53:15
109247	128381	p191189	ca25d5573f2bd2660f7850a886395927	Ms.	VA	2016-06-29 09:17:01

109248 rows × 25 columns



Computing Sentiment Scores

In [40]:

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

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

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students w
ith the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelli
gences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of differen
t backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a carin
g 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 ki
ndergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role pla
y 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 co
oking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into maki
ng the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project woul
```



```

d expand our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make homemade apple
sauce make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create our own cook
books to be printed and \
shared with families students will gain math and literature skills as well as a life long enjoymen
t for healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

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

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

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

```

Assignment 7: SVM

1. [Task-1] Apply Support Vector Machines(SGDClassifier with hinge loss: Linear SVM) 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)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V)

2. The hyper paramter tuning (best alpha in range [10^{-4} to 10^4], and the best penalty among 'l1', 'l2')

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

4. [\[Task-2\] Apply the Support Vector Machines on these features by finding the best hyper paramter as suggested in step 2 and step 3](#)

- [Consider these set of features Set 5:](#)
 - [school_state](#) : categorical data
 - [clean_categories](#) : categorical data
 - [clean_subcategories](#) : categorical data
 - [project_grade_category](#) :categorical data
 - [teacher_prefix](#) : categorical data
 - [quantity](#) : numerical data
 - [teacher_number_of_previously_posted_projects](#) : numerical data
 - [price](#) : numerical data
 - [sentiment score's of each of the essay](#) : numerical data
 - [number of words in the title](#) : numerical data
 - [number of words in the combine essays](#) : numerical data
 - [Apply TruncatedSVD on TfidfVectorizer](#) of essay text, choose the number of components ('n_components') using [elbow method](#) : numerical data

• Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this [prettytable library link](#)

Note: Data Leakage

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

2. Support Vector Machines

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

In [41]:

```
#https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6

# ===== loading libraries =====
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn import model_selection
# =====

y13 = project_data['project_is_approved']
X_train, X_test, y_train, y_test = train_test_split(project_data, y13, stratify=y13, test_size=0.3)
```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [42]:

```
# 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

from sklearn.feature_extraction.text import CountVectorizer
vectorizer_train_tr_1 = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
binary=True)
categories_one_hot_train = vectorizer_train_tr_1.fit(X_train['clean_categories'].values)
#print(vectorizer_train_tr_1.get_feature_names())
#print("Shape of matrix after one hot encodig ",categories_one_hot_train_tr.shape)

categories_one_hot_train_tr =
categories_one_hot_train.transform(X_train['clean_categories'].values)
#print(vectorizer_train_tr_1.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_train_tr.shape)

#feature names encoding for X_test
#from sklearn.feature_extraction.text import CountVectorizer
#vectorizer_test_1 = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, bin
ary=True)
categories_one_hot_test = categories_one_hot_train.transform(X_test['clean_categories'].values)
#print(vectorizer_test_1.get_feature_names())
print("Shape of matrix after one hot encodig ",categories one hot test.shape)
```

```

# we use count vectorizer to convert the values into one for X_train
vectorizer_train_tr = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot_train = vectorizer_train_tr.fit(X_train['clean_subcategories'].values)
#print(vectorizer_train_tr.get_feature_names())
#print("Shape of matrix after one hot encoding ", sub_categories_one_hot_train_tr.shape)

sub_categories_one_hot_train_tr =
sub_categories_one_hot_train.transform(X_train['clean_subcategories'].values)
#print(vectorizer_train_tr.get_feature_names())
print("Shape of matrix after one hot encoding ", sub_categories_one_hot_train_tr.shape)

# we use count vectorizer to convert the values into one for X_test
vectorizer_test = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot_test = sub_categories_one_hot_train.transform(X_test['clean_subcategories'].values)
#print(vectorizer_test.get_feature_names())
print("Shape of matrix after one hot encoding ", sub_categories_one_hot_test.shape)

#school state for Xtrain
vectorizer_1_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_1_train_tr.fit(X_train['school_state'].values)
#print(vectorizer_1_train_tr.get_feature_names())
categories_state_1_train = vectorizer_1_train_tr.fit(X_train['school_state'].values)
#print("Shape of matrix after one hot encoding ", categories_state_1_train_tr.shape)

categories_state_1_train_tr = categories_state_1_train.transform(X_train['school_state'].values)
print("Shape of matrix after one hot encoding ", categories_state_1_train_tr.shape)

#school state for Xtest
vectorizer_1_test = CountVectorizer(lowercase=False, binary=True)
vectorizer_1_test.fit(X_test['school_state'].values)
#print(vectorizer_1_test.get_feature_names())
categories_state_1_test = categories_state_1_train.transform(X_test['school_state'].values)
print("Shape of matrix after one hot encoding ", categories_state_1_test.shape)

#Project grade category for X_train
#vectorizer_2_train = CountVectorizer(vocabulary=list(word_dict_134.keys()), lowercase=False, binary=True)
vectorizer_2_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_2_train_tr.fit(X_train['project_grade_category'].values)
#print(vectorizer_2_train_tr.get_feature_names())
categories_grade_train = vectorizer_2_train_tr.fit(X_train['project_grade_category'].values)
#print("Shape of matrix after one hot encoding ", categories_grade_train_tr.shape)

categories_grade_train_tr = categories_grade_train.transform(X_train['project_grade_category'].values)
print("Shape of matrix after one hot encoding ", categories_grade_train_tr.shape)

#Project grade category for X_test
#vectorizer_2_test = CountVectorizer(vocabulary=list(word_dict_134.keys()), lowercase=False, binary=True)
#vectorizer_2_test = CountVectorizer(lowercase=False, binary=True)
vectorizer_2_test.fit(X_test['project_grade_category'].values)
#print(vectorizer_2_test.get_feature_names())
categories_grade_test = categories_grade_train.transform(X_test['project_grade_category'].values)
print("Shape of matrix after one hot encoding ", categories_grade_test.shape)

from string import punctuation

#https://medium.com/@chaimgluck1/have-messy-text-data-clean-it-with-simple-lambda-functions-645918fcc2fc
#project_data.teacher_prefix = project_data.teacher_prefix.apply(lambda x:
x.translate(string.punctuation))

#https://stackoverflow.com/questions/50443494/error-in-removing-punctuation-float-object-has-no-attribute-translate
#project_data['teacher_prefix'] = project_data.fillna({'teacher_prefix':''})

```

```

project_data['teacher_prefix'] = project_data['teacher_prefix'].replace(np.nan, 'teacher')

#teacher_prefix for X_train
vectorizer_3_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_3_train_tr.fit(X_train['teacher_prefix'].values)
#print(vectorizer_3_train_tr.get_feature_names())
categories_teacher_prefix_train = vectorizer_3_train_tr.fit(X_train['teacher_prefix'].values)
#print("Shape of matrix after one hot encodig ",categories_teacher_prefix_train_tr.shape)

categories_teacher_prefix_train_tr =
categories_teacher_prefix_train.transform(X_train['teacher_prefix'].values)
print("Shape of matrix after one hot encodig ",categories_teacher_prefix_train_tr.shape)

#teacher_prefix for X_test
#vectorizer_3_test = CountVectorizer(lowercase=False, binary=True)
#vectorizer_3_test.fit(X_test['teacher_prefix'].values)
#print(vectorizer_3_test.get_feature_names())
categories_teacher_prefix_test = categories_teacher_prefix_train.transform(X_test['teacher_prefix']
].values)
print("Shape of matrix after one hot encodig ",categories_teacher_prefix_test.shape)

```

```

Shape of matrix after one hot encodig (76473, 9)
Shape of matrix after one hot encodig (32775, 9)
Shape of matrix after one hot encodig (76473, 30)
Shape of matrix after one hot encodig (32775, 30)
Shape of matrix after one hot encodig (76473, 51)
Shape of matrix after one hot encodig (32775, 51)
Shape of matrix after one hot encodig (76473, 4)
Shape of matrix after one hot encodig (32775, 4)
Shape of matrix after one hot encodig (76473, 6)
Shape of matrix after one hot encodig (32775, 6)

```

2.3 Make Data Model Ready: encoding eassay, and project_title

In [43]:

```

# 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

#Vectorizing essays for X_train.

#First fitting the vector
vectorizer_essays_train_tr = CountVectorizer(min_df=10, ngram_range=(2,2), max_features=5000)
text_bow_train = vectorizer_essays_train_tr.fit(X_train['preprocessed_essays'])
#print("Shape of matrix after one hot encodig ",text_bow_train_tr.shape)

#transforming train data
text_bow_train_tr = text_bow_train.transform(X_train['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_bow_train_tr.shape)

#transforming text data.
vectorizer_essays_test = CountVectorizer(min_df=10, ngram_range=2)
text_bow_test = text_bow_train.transform(X_test['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_bow_test.shape)

# Vectorizing Title for X_train
vectorizer_title_tr_tr = CountVectorizer(min_df=10)
title_bow_tr = vectorizer_title_tr_tr.fit(X_train['preprocessed_titles'])
#print("Shape of matrix after one hot encodig ",title_bow_tr_tr.shape)

title_bow_tr_tr = title_bow_tr.transform(X_train['preprocessed_titles'])

```

```

title_bow_tr_tr = title_bow_tr.transform(X_train['preprocessed_titles'])
print("Shape of matrix after one hot encodig ",title_bow_tr_tr.shape)

# Vectorizing Title for X_test
vectorizer_title_test = CountVectorizer(min_df=10)
title_bow_test = title_bow_tr.transform(X_test['preprocessed_titles'])
print("Shape of matrix after one hot encodig ",title_bow_test.shape)

#Vectorizing using tfidf for essays for x_train
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essays_train_tr = TfidfVectorizer(min_df=10, ngram_range=(2,2), max_features=5000)
text_tfidf_train = vectorizer_tfidf_essays_train_tr.fit(X_train['preprocessed_essays'])
#print("Shape of matrix after one hot encodig ",text_tfidf_train_tr.shape)

text_tfidf_train_tr = text_tfidf_train.transform(X_train['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_tfidf_train_tr.shape)

#Vectorizing using tfidf for essays for x_test
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essays_test = TfidfVectorizer(min_df=10)
text_tfidf_test = text_tfidf_train.transform(X_test['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_tfidf_test.shape)

#Vectorizing using tfidf for titles using X_train
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_tr = TfidfVectorizer(min_df=10)
title_tfidf_train = vectorizer_tfidf_tr.fit(X_train['preprocessed_titles'])
#print("Shape of matrix after one hot encodig ",title_tfidf_train_tr.shape)

title_tfidf_train_tr = title_tfidf_train.transform(X_train['preprocessed_titles'])
print("Shape of matrix after one hot encodig ",title_tfidf_train_tr.shape)

#Vectorizing using tfidf for titles using X_test
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
title_tfidf_test = title_tfidf_train.transform(X_test['preprocessed_titles'])
print("Shape of matrix after one hot encodig ",title_tfidf_test.shape)

```

```

Shape of matrix after one hot encodig (76473, 5000)
Shape of matrix after one hot encodig (32775, 5000)
Shape of matrix after one hot encodig (76473, 2696)
Shape of matrix after one hot encodig (32775, 2696)
Shape of matrix after one hot encodig (76473, 5000)
Shape of matrix after one hot encodig (32775, 5000)
Shape of matrix after one hot encodig (76473, 2696)
Shape of matrix after one hot encodig (32775, 2696)

```

In [44]:

```

avg_w2v_vectors_1_train_tr = []
for sentence in tqdm(X_train['preprocessed_titles']): # 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_1_train_tr.append(vector)

print(len(avg_w2v_vectors_1_train_tr))
print(len(avg_w2v_vectors_1_train_tr[0]))

avg_w2v_vectors_1_test = []
for sentence in tqdm(X_test['preprocessed_titles']): # 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_1_test.append(vector)

print(len(avg_w2v_vectors_1_test))
print(len(avg_w2v_vectors_1_test[0]))

avg_w2v_vectors_1_train_tr_essay = []
for sentence in tqdm(X_train['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_1_train_tr_essay.append(vector)

print(len(avg_w2v_vectors_1_train_tr_essay))
print(len(avg_w2v_vectors_1_train_tr_essay[0]))

avg_w2v_vectors_1_test_essay = []
for sentence in tqdm(X_test['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_1_test_essay.append(vector)

print(len(avg_w2v_vectors_1_test_essay))
print(len(avg_w2v_vectors_1_test_essay[0]))

```

[illegible]

76473
300

[illegible]

32775
300

```
100% |██████████| 76473/76473 [01:05<00:00, 1175.33it/s]
```

76473
300

```
100%|███████████| 32775/32775 [00:23<00:00, 1392.68it/s]
```

32775
300

In [45]:

```
# Similarly you can vectorize for title also, first for train
tfidf_model_1_train_tr = TfidfVectorizer()
tfidf_model_1_train_tr.fit(X_train['preprocessed_titles'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary_1 = dict(zip(tfidf_model_1_train_tr.get_feature_names(), list(tfidf_model_1_train_tr.idf_)))
tfidf_words_1_train_tr = set(tfidf_model_1_train_tr.get_feature_names())
```

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_1_train_tr = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_titles']): # 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_1_train_tr):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary_1[word]*(sentence.count(word)/len(sentence.split())) # getting the
            tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_1_train_tr.append(vector)

print(len(tfidf_w2v_vectors_1_train_tr))
print(len(tfidf_w2v_vectors_1_train_tr[0]))

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_1_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_titles']): # 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_1_train_tr):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary_1[word]*(sentence.count(word)/len(sentence.split())) # getting the
            tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_1_test.append(vector)

print(len(tfidf_w2v_vectors_1_test))
print(len(tfidf_w2v_vectors_1_test[0]))
```

```
100%|███████████| 76473/76473 [00:06<00:00, 12599.25it/s]
```

76473
300

```
100% |██████████| 32775/32775 [00:02<00:00, 11942.19it/s]
```

32775
300

In [46]:

```
# Similarly you can vectorize for title also, first for train
tfidf_model_1_train_tr_essay = TfidfVectorizer()
tfidf_model_1_train_tr_essay.fit(X_train['preprocessed_essays'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary_1 = dict(zip(tfidf_model_1_train_tr_essay.get_feature_names(),
list(tfidf_model_1_train_tr_essay.idf)))
tfidf_words_1_train_tr_essay = set(tfidf_model_1_train_tr_essay.get_feature_names())

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_1_train_tr_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(3000) # as word vectors are of zero length
```



```

#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar_train_tr.mean_[0]}, Standard
deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar_train_tr.var_[0])}")
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_train_tr =
teacher_number_of_previously_posted_projects_scalar_train_tr.transform(X_train['teacher_number_of_p
reviously_posted_projects'].values.reshape(-1, 1))

teacher_number_of_previously_posted_projects_scalar_test = StandardScaler()
teacher_number_of_previously_posted_projects_scalar_test.fit(X_test['teacher_number_of_previously_p
sted_projects'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar.mean_[0]}, Standard deviation
: {np.sqrt(teacher_number_of_previously_posted_projects_scalar.var_[0])}")
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_test =
teacher_number_of_previously_posted_projects_scalar_test.transform(X_test['teacher_number_of_previc
usly_posted_projects'].values.reshape(-1, 1))

# 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_X_train = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
#X.shape

#categorical, numerical features + project_title(BOW)
X1_train_tr = hstack((categories_one_hot_train_tr, sub_categories_one_hot_train_tr,
categories_state_1_train_tr, categories_grade_train_tr, categories_teacher_prefix_train_tr,
price_standardized_train_tr, teacher_number_of_previously_posted_projects_standardized_train_tr, t
itle_bow_train_tr, text_bow_train_tr))

X1_test = hstack((categories_one_hot_test, sub_categories_one_hot_test, categories_state_1_test, c
ategories_grade_test, categories_teacher_prefix_test, price_standardized_test,
teacher_number_of_previously_posted_projects_standardized_test, title_bow_test, text_bow_test))

#categorical, numerical features + project_title(TFIDF)
X2_train_tr = hstack((categories_one_hot_train_tr, sub_categories_one_hot_train_tr,
categories_state_1_train_tr, categories_grade_train_tr, categories_teacher_prefix_train_tr,
price_standardized_train_tr, teacher_number_of_previously_posted_projects_standardized_train_tr, t
itle_tfidf_train_tr, text_tfidf_train_tr))

X2_test = hstack((categories_one_hot_test, sub_categories_one_hot_test, categories_state_1_test, c
ategories_grade_test, categories_teacher_prefix_test, price_standardized_test,
teacher_number_of_previously_posted_projects_standardized_test, title_tfidf_test, text_tfidf_test)
)

#categorical, numerical features + project_title(AVG W2V)
X3_train_tr = hstack((categories_one_hot_train_tr, sub_categories_one_hot_train_tr,
categories_state_1_train_tr, categories_grade_train_tr, categories_teacher_prefix_train_tr,
price_standardized_train_tr, teacher_number_of_previously_posted_projects_standardized_train_tr, a
vg_w2v_vectors_1_train_tr_essay, avg_w2v_vectors_1_train_tr))

X3_test = hstack((categories_one_hot_test, sub_categories_one_hot_test, categories_state_1_test, c
ategories_grade_test, categories_teacher_prefix_test, price_standardized_test,
teacher_number_of_previously_posted_projects_standardized_test, avg_w2v_vectors_1_test_essay, avg_
w2v_vectors_1_test))

#categorical, numerical features + project_title(TFIDF W2V)
X4_train_tr = hstack((categories_one_hot_train_tr, sub_categories_one_hot_train_tr,
categories_state_1_train_tr, categories_grade_train_tr, categories_teacher_prefix_train_tr,
price_standardized_train_tr, teacher_number_of_previously_posted_projects_standardized_train_tr, t
fidf_w2v_vectors_1_train_tr_essay, tfidf_w2v_vectors_1_train_tr))

X4_test = hstack((categories_one_hot_test, sub_categories_one_hot_test, categories_state_1_test, c
ategories_grade_test, categories_teacher_prefix_test, price_standardized_test,
teacher_number_of_previously_posted_projects_standardized_test, tfidf_w2v_vectors_1_test_essay, tf
idf_w2v_vectors_1_test))

```

2.4 Applying Support Vector Machines on different kind of featurization as mentioned in the instructions

Apply Support Vector Machines 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 [48]:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets, neighbors
from matplotlib.colors import ListedColormap

#https://www.ritchieng.com/machine-learning-efficiently-search-tuning-param/

# imports
from sklearn.datasets import load_iris
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import cross_val_score
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.model_selection import GridSearchCV
import matplotlib.pyplot as plt
from sklearn.linear_model import SGDClassifier
from sklearn.metrics import roc_auc_score, roc_curve, f1_score, auc

parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}

model = SGDClassifier(loss = 'hinge',penalty='l2', class_weight= "balanced" )

grid = GridSearchCV(model,param_grid=parameters,cv=5,scoring='roc_auc', return_train_score=True)
```

In [49]:

```
grid.fit(X1_train_tr, y_train)
```

Out[49]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight='balanced',
                                     early_stopping=False, epsilon=0.1,
                                     eta0=0.0, fit_intercept=True,
                                     l1_ratio=0.15, learning_rate='optimal',
                                     loss='hinge', max_iter=1000,
                                     n_iter_no_change=5, n_jobs=None,
                                     penalty='l2', power_t=0.5,
                                     random_state=None, shuffle=True, tol=0.001,
                                     validation_fraction=0.1, verbose=0,
                                     warm_start=False),
             iid='warn', n_jobs=None,
             param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0,
                                   1000.0, 10000.0]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring='roc_auc', verbose=0)
```

In [50]:

```
grid.cv_results_
```

Out[50]:

```
{'mean_fit_time': array([3.10011096, 1.00862155, 0.62584052, 0.44814382, 0.32181988,
                        0.29902964, 0.27804508, 0.28044038, 0.29623165]),
 'std_fit_time': array([1.44251986, 0.21098851, 0.10803691, 0.09300722, 0.05021732,
                        0.02841128, 0.01539656, 0.01431126, 0.03759548]),
 'mean_score_time': array([0.17802348, 0.09135113, 0.058567 , 0.01738977, 0.01198821,
                        0.01199617, 0.01239252, 0.0117928 , 0.01259623]),
 'std_score_time': array([0.19998306, 0.0685404 , 0.09265748, 0.01282381, 0.00141316,
                        0.00141731, 0.00162367, 0.00074811, 0.00119428]),
 'param_alpha': masked_array(data=[0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0,
                                   10000.0],
                              mask=[False, False, False, False, False, False, False, False,
                                    False],
                              fill_value='?').
```

```

dtype=object),
'params': [{'alpha': 0.0001},
{'alpha': 0.001},
{'alpha': 0.01},
{'alpha': 0.1},
{'alpha': 1},
{'alpha': 10.0},
{'alpha': 100.0},
{'alpha': 1000.0},
{'alpha': 10000.0}],
'split0_test_score': array([0.65219429, 0.67691361, 0.69784841, 0.66523498, 0.63494585,
0.63444983, 0.6344496 , 0.6344496 , 0.6344496 ]),
'split1_test_score': array([0.64908852, 0.67226193, 0.68675385, 0.65979736, 0.63760358,
0.63733438, 0.63733438, 0.63733438, 0.63733438]),
'split2_test_score': array([0.65910902, 0.68222245, 0.69908285, 0.66867602, 0.65238456,
0.652127 , 0.652127 , 0.652127 , 0.652127 ]),
'split3_test_score': array([0.64907271, 0.66857193, 0.68147021, 0.6497208 , 0.62995977,
0.62962516, 0.62962516, 0.62962516, 0.62962516]),
'split4_test_score': array([0.64953836, 0.67333873, 0.69102475, 0.6680945 , 0.64501377,
0.64457615, 0.64457615, 0.64457615, 0.64457615]),
'mean_test_score': array([0.65180064, 0.67466176, 0.69123602, 0.66230458, 0.63998138,
0.63962238, 0.63962233, 0.63962233, 0.63962233]),
'std_test_score': array([0.00383363, 0.0046224 , 0.00664542, 0.00703343, 0.00788178,
0.00791006, 0.00791009, 0.00791009, 0.00791009]),
'rank_test_score': array([4, 2, 1, 3, 5, 6, 7, 7, 7]),
'split0_train_score': array([0.83851196, 0.81703047, 0.75557785, 0.67900774, 0.64570199,
0.64506862, 0.64506846, 0.64506846, 0.64506846]),
'split1_train_score': array([0.84025425, 0.81937055, 0.75740987, 0.68045076, 0.65009564,
0.64949409, 0.64949409, 0.64949409, 0.64949409]),
'split2_train_score': array([0.83730589, 0.81743827, 0.75370203, 0.67858303, 0.64799976,
0.64738611, 0.64738611, 0.64738611, 0.64738611]),
'split3_train_score': array([0.84041059, 0.81989374, 0.75760927, 0.68254252, 0.65095786,
0.6503803 , 0.6503803 , 0.6503803 , 0.6503803 ]),
'split4_train_score': array([0.83912093, 0.81908953, 0.75719194, 0.68079737, 0.64820588,
0.64754327, 0.64754327, 0.64754327, 0.64754327]),
'mean_train_score': array([0.83912072, 0.81856451, 0.75629819, 0.68027628, 0.64859223,
0.64797448, 0.64797445, 0.64797445, 0.64797445]),
'std_train_score': array([0.00115003, 0.00112373, 0.00148408, 0.00140813, 0.00182813,
0.00184788, 0.00184793, 0.00184793, 0.00184793])}]

```

In [51]:

```
pd.DataFrame(grid.cv_results_)
```

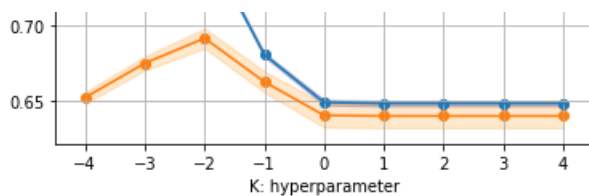
Out[51]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_sco
0	3.100111	1.442520	0.178023	0.199983	0.0001	{'alpha': 0.0001}	0.652194	0.649089
1	1.008622	0.210989	0.091351	0.068540	0.001	{'alpha': 0.001}	0.676914	0.672262
2	0.625841	0.108037	0.058567	0.092657	0.01	{'alpha': 0.01}	0.697848	0.686754
3	0.448144	0.093007	0.017390	0.012824	0.1	{'alpha': 0.1}	0.665235	0.659797
4	0.321820	0.050217	0.011988	0.001413	1	{'alpha': 1}	0.634946	0.637604
5	0.299030	0.028411	0.011996	0.001417	10	{'alpha': 10.0}	0.634450	0.637334
6	0.278045	0.015397	0.012393	0.001624	100	{'alpha': 100.0}	0.634450	0.637334
7	0.280440	0.014311	0.011793	0.000748	1000	{'alpha': 1000.0}	0.634450	0.637334
8	0.296232	0.037595	0.012596	0.001194	10000	{'alpha': 10000.0}	0.634450	0.637334

◀ ▶

```
train_auc1 = grid_cv_results['mean_train_score']
```

[illegible]



In [54]:

```
#https://stackoverflow.com/questions/39200265/attributeerror-probability-estimates-are-not-available-for-loss-hinge
from sklearn.calibration import CalibratedClassifierCV

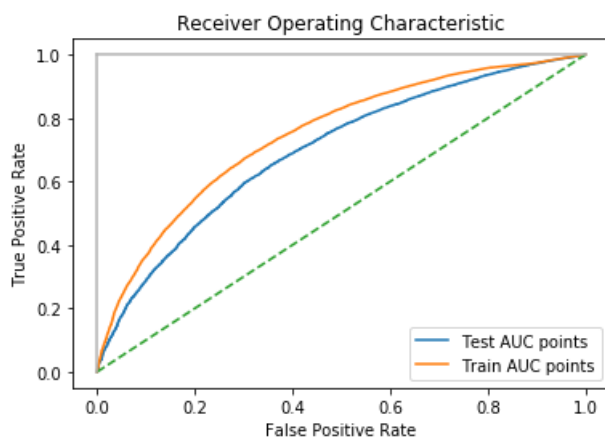
meowl = SGDClassifier(loss = 'hinge',penalty='l2', alpha = 0.01, class_weight= "balanced")
fit_meowl = meowl.fit(X1_train_tr, y_train)
calibrator = CalibratedClassifierCV(fit_meowl, cv='prefit')
Trained_model_BOW = calibrator.fit(X1_train_tr, y_train)
```

In [55]:

```
# Get predicted probabilities
y_score_test = Trained_model_BOW.predict_proba(X1_test)[:,-1]
y_score_train = Trained_model_BOW.predict_proba(X1_train_tr)[:,-1]

# Create true and false positive rates
false_positive_rate, true_positive_rate, threshold = roc_curve(y_test, y_score_test)
false_positive_rate1, true_positive_rate1, threshold1 = roc_curve(y_train, y_score_train)

# Plot ROC curve
plt.title('Receiver Operating Characteristic')
plt.plot(false_positive_rate, true_positive_rate, label='Test AUC points')
plt.plot(false_positive_rate1, true_positive_rate1, label='Train AUC points')
plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0] , c=".7"), plt.plot([1, 1] , c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()
```



In [56]:

```
import numpy as np
from sklearn.metrics import roc_auc_score

roc_auc_score(y_test, y_score_test)
```

Out[56]:

0.6966059338167592

In [57]:

```
#https://chrisalbon.com/machine_learning/model_evaluation/generate_text_reports_on_performance/
```

```

from sklearn.metrics import classification_report

# Create list of target class names
#class_names = project_data['project_is_approved'].target_names

# Train model and make predictions
y_hat = Trained_model_BOW.predict(X1_test)

print(classification_report(y_test, y_hat))

```

	precision	recall	f1-score	support
0	0.41	0.04	0.07	4963
1	0.85	0.99	0.92	27812
accuracy			0.85	32775
macro avg	0.63	0.52	0.50	32775
weighted avg	0.79	0.85	0.79	32775

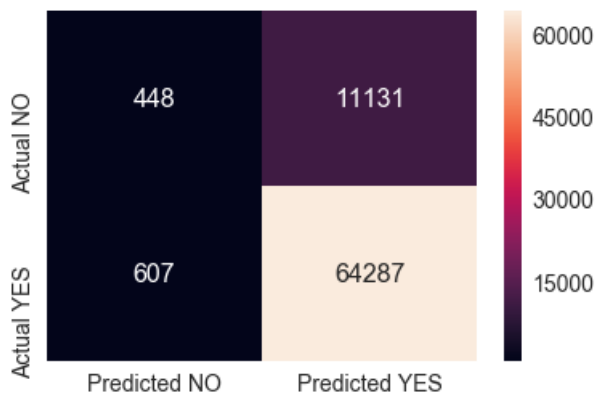
In [58]:

```

#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), range(2),range(2))
    df_cm.columns = ['Predicted NO','Predicted YES']
    df_cm = df_cm.rename({0: 'Actual NO', 1: 'Actual YES'})
    sns.set(font_scale=1.4)#for label size
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

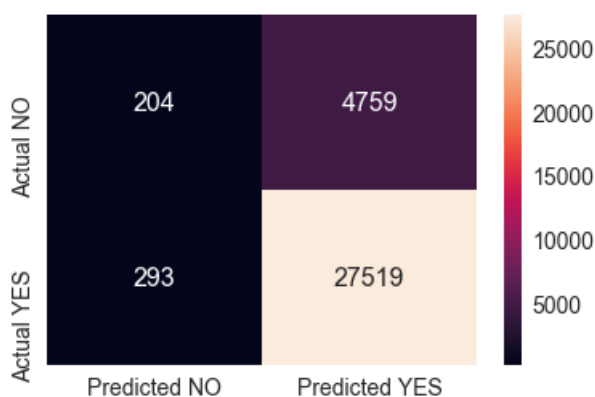
get_confusion_matrix(Trained_model_BOW, X1_train_tr, y_train)

```



In [59]:

```
get_confusion_matrix(Trained_model_BOW, X1_test, y_test)
```



In [60]:

```
parameters={'alpha': [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
```

```
parameters = {'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0, 10000.0]}
```

```
model = SGDClassifier(loss = 'hinge',penalty='l2', class_weight= "balanced")
```

```
grid2 = GridSearchCV(model,param_grid=parameters,cv=5,scoring='roc_auc', return_train_score=True)
```

In [61]:

```
grid2.fit(X2_train_tr, y_train)
```

Out[61]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight='balanced',
                                     early_stopping=False, epsilon=0.1,
                                     eta0=0.0, fit_intercept=True,
                                     l1_ratio=0.15, learning_rate='optimal',
                                     loss='hinge', max_iter=1000,
                                     n_iter_no_change=5, n_jobs=None,
                                     penalty='l2', power_t=0.5,
                                     random_state=None, shuffle=True, tol=0.001,
                                     validation_fraction=0.1, verbose=0,
                                     warm_start=False),
             iid='warn', n_jobs=None,
             param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0,
                                   1000.0, 10000.0]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring='roc_auc', verbose=0)
```

In [62]:

```
grid2.cv_results_
```

Out[62]:

```
{'mean_fit_time': array([1.19071698, 0.75317149, 0.37798443, 0.35320225, 0.27224488,
                        0.28263831, 0.26764674, 0.27903986, 0.30002842]),
 'std_fit_time': array([0.22724884, 0.11799134, 0.06187858, 0.0840365 , 0.02683703,
                        0.01615825, 0.01740781, 0.02480051, 0.0299315 ]),
 'mean_score_time': array([0.01379156, 0.01238875, 0.02198777, 0.01139331, 0.01119676,
                        0.01179304, 0.01199279, 0.01219196, 0.01439118]),
 'std_score_time': array([0.00330843, 0.0022456 , 0.01858878, 0.00135585, 0.00075195,
                        0.00222623, 0.00089447, 0.00213477, 0.0027254 ]),
 'param_alpha': masked_array(data=[0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0,
                                   10000.0],
                             mask=[False, False, False, False, False, False, False, False,
                                   False],
                             fill_value='?',
                             dtype=object),
 'params': [{'alpha': 0.0001},
            {'alpha': 0.001},
            {'alpha': 0.01},
            {'alpha': 0.1},
            {'alpha': 1},
            {'alpha': 10.0},
            {'alpha': 100.0},
            {'alpha': 1000.0},
            {'alpha': 10000.0}],
 'split0_test_score': array([0.68843468, 0.67283341, 0.61750091, 0.60814761, 0.61296754,
                            0.61307277, 0.61307277, 0.61307277, 0.61307277]),
 'split1_test_score': array([0.6819294 , 0.66994927, 0.61596257, 0.60514497, 0.61211056,
                            0.61233491, 0.61233491, 0.61233491, 0.61233491]),
 'split2_test_score': array([0.69000196, 0.67890282, 0.62367737, 0.61526345, 0.62196 ,
                            0.62209892, 0.62209892, 0.62209892, 0.62209892]),
 'split3_test_score': array([0.67897182, 0.66083175, 0.61015572, 0.59909265, 0.60474343,
                            0.60478455, 0.60478419, 0.60478419, 0.60478419]),
 'split4_test_score': array([0.68125497, 0.67064386, 0.62737482, 0.61792926, 0.61987617,
                            0.61939504, 0.61939504, 0.61939504, 0.61939504]),
 'mean_test_score': array([0.68411864, 0.67063222, 0.61893406, 0.60911536, 0.61433139,
                          0.61433711, 0.61433703, 0.61433703, 0.61433703]),
 'std_test_score': array([0.00430642, 0.00582698, 0.006028 , 0.00682107, 0.00612637,
                          0.00604464, 0.00604476, 0.00604476, 0.00604476]),
 'rank_test_score': array([1, 2, 3, 9, 8, 4, 5, 5, 5]),
 'split0_train_score': array([0.81326027, 0.7101342 , 0.62667422, 0.61175962, 0.61584179,
                              0.61584179, 0.61584179, 0.61584179, 0.61584179])}
```

```

0.61579534, 0.61579534, 0.61579534, 0.61579534)),
'split1_train_score': array([0.81361215, 0.71617423, 0.62631769, 0.61106227, 0.61591583,
0.61595535, 0.61595535, 0.61595535, 0.61595535]),
'split2_train_score': array([0.81027291, 0.71290068, 0.6264094 , 0.60957464, 0.61372217,
0.61367843, 0.61367843, 0.61367843, 0.61367843]),
'split3_train_score': array([0.80985675, 0.71235318, 0.62964373, 0.61300675, 0.61782454,
0.61784953, 0.61784936, 0.61784936, 0.61784936]),
'split4_train_score': array([0.81416676, 0.70949992, 0.62821313, 0.61080246, 0.6145523 ,
0.61448348, 0.61448348, 0.61448348, 0.61448348]),
'mean_train_score': array([0.81223377, 0.71221244, 0.62745163, 0.61124115, 0.61557132,
0.61555243, 0.61555239, 0.61555239, 0.61555239]),
'std_train_score': array([0.00179918, 0.00236077, 0.00129318, 0.0011303 , 0.00139475,
0.00142496, 0.0014249 , 0.0014249 , 0.0014249 ])}

```

In [63]:

```

train_auc2 = grid2.cv_results_['mean_train_score']
train_auc_std2 = grid2.cv_results_['std_train_score']
cv_auc2 = grid2.cv_results_['mean_test_score']
cv_auc_std2 = grid2.cv_results_['std_test_score']

log_k_range=[]
for a in tqdm(k_range):

    b = np.log10(a)
    log_k_range.append(b)

plt.plot(log_k_range, train_auc2, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, train_auc2 - train_auc_std2, train_auc2 + train_auc_std2,alpha=
0.2,color='darkblue')

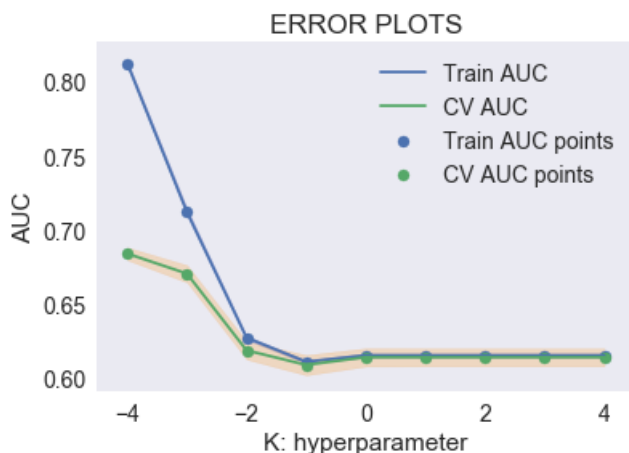
plt.plot(log_k_range, cv_auc2, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, cv_auc2 - cv_auc_std2, cv_auc2 + cv_auc_std2,alpha=0.2,color='d
arkorange')

plt.scatter(log_k_range, train_auc2, label='Train AUC points')
plt.scatter(log_k_range, cv_auc2, label='CV AUC points')

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

```

100% |████████████████████████████████████████████████████████████████████████████████| 9/9 [00:00<00:00, 9000.65it/s]



In [64]:

```

# examine the best model

# Single best score achieved across all params (k)
print(grid2.best_score_)

```



```
# Dictionary containing the parameters (k) used to generate that score
print(grid2.best_params_)

# Actual model object fit with those best parameters
# Shows default parameters that we did not specify
print(grid2.best_estimator_)
```

```
0.6841186385013454
{'alpha': 0.0001}
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

In [65]:

```
#https://stackoverflow.com/questions/39200265/attributeerror-probability-estimates-are-not-available-for-loss-hinge
from sklearn.calibration import CalibratedClassifierCV

meow2 = SGDClassifier(loss = 'hinge',penalty='l2', alpha = 0.001, class_weight= "balanced")
fit_meow2 = meow2.fit(X2_train_tr, y_train)
calibrator = CalibratedClassifierCV(fit_meow2, cv='prefit')

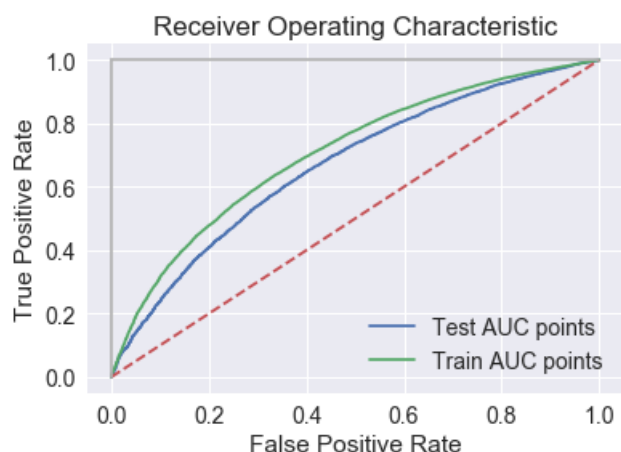
Trained_model_TFIDF = calibrator.fit(X2_train_tr, y_train)
```

In [66]:

```
# Get predicted probabilities
y_score_test2 = Trained_model_TFIDF.predict_proba(X2_test)[: ,1]
y_score_train2 = Trained_model_TFIDF.predict_proba(X2_train_tr)[: ,1]

# Create true and false positive rates
false_positive_rate2, true_positive_rate2, threshold2 = roc_curve(y_test, y_score_test2)
false_positive_rate21, true_positive_rate21, threshold21 = roc_curve(y_train, y_score_train2)

# Plot ROC curve
plt.title('Receiver Operating Characteristic')
plt.plot(false_positive_rate2, true_positive_rate2, label='Test AUC points')
plt.plot(false_positive_rate21, true_positive_rate21, label='Train AUC points')
plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0], c=".7"), plt.plot([1, 1], c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()
```



In [67]:

```
import numpy as np
```

```
from sklearn.metrics import roc_auc_score

roc_auc_score(y_test, y_score_test2)
```

Out[67]:
0.6665373816580681

In [68]:

```
#https://chrisalbon.com/machine_learning/model_evaluation/generate_text_reports_on_performance/
from sklearn.metrics import classification_report

# Create list of target class names
#class_names = project_data['project_is_approved'].target_names

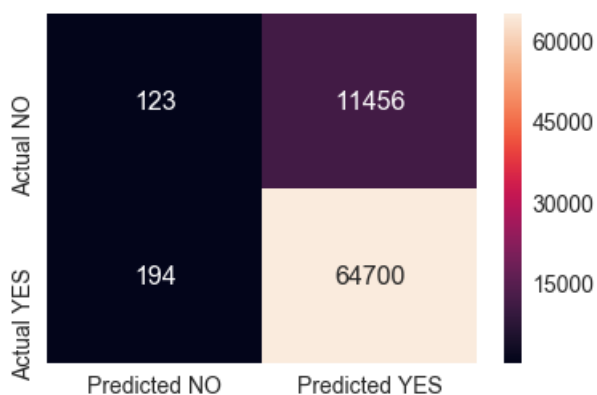
# Train model and make predictions
y_hat_2 = Trained_model_TFIDF.predict(X2_test)

print(classification_report(y_test, y_hat_2))
```

	precision	recall	f1-score	support
0	0.37	0.01	0.02	4963
1	0.85	1.00	0.92	27812
accuracy			0.85	32775
macro avg	0.61	0.50	0.47	32775
weighted avg	0.78	0.85	0.78	32775

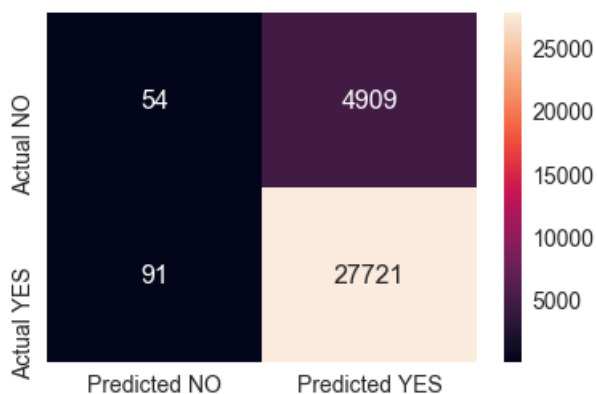
In [69]:

```
get_confusion_matrix(Trained_model_TFIDF, X2_train_tr, y_train)
```



In [70]:

```
get_confusion_matrix(Trained_model_TFIDF, X2_test, y_test)
```



In [71]:

```
## [1:]
```

```
grid3 = GridSearchCV(model,param_grid=parameters,cv=5,scoring='roc_auc', return_train_score=True)
```

```
In [72]:
```

```
grid3.fit(X3_train_tr, y_train)
```

```
Out[72]:
```

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight='balanced',
                                     early_stopping=False, epsilon=0.1,
                                     eta0=0.0, fit_intercept=True,
                                     l1_ratio=0.15, learning_rate='optimal',
                                     loss='hinge', max_iter=1000,
                                     n_iter_no_change=5, n_jobs=None,
                                     penalty='l2', power_t=0.5,
                                     random_state=None, shuffle=True, tol=0.001,
                                     validation_fraction=0.1, verbose=0,
                                     warm_start=False),
             iid='warn', n_jobs=None,
             param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0,
                                    1000.0, 10000.0]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring='roc_auc', verbose=0)
```

```
In [73]:
```

```
# view the complete results (list of named tuples)
grid3.cv_results_
```

```
Out[73]:
```

```
{'mean_fit_time': array([9.1277566 , 5.14123425, 2.68061409, 1.95706596, 1.5711586 ,
                        1.5323462 , 1.40959511, 1.37022471, 1.35522742]),
 'std_fit_time': array([1.2119242 , 0.24298151, 0.2982299 , 0.19145088, 0.26515753,
                        0.18043059, 0.08060969, 0.04540232, 0.0344082 ]),
 'mean_score_time': array([0.18319316, 0.07996273, 0.07141833, 0.03657784, 0.03678098,
                        0.03657103, 0.03657851, 0.04552269, 0.05008402]),
 'std_score_time': array([0.25302672, 0.09048309, 0.06889564, 0.00149288, 0.0023155 ,
                        0.00080661, 0.00135545, 0.01007746, 0.02530443]),
 'param_alpha': masked_array(data=[0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0,
                                   10000.0],
                             mask=[False, False, False, False, False, False, False, False,
                                   False],
                             fill_value='?',
                             dtype=object),
 'params': [{'alpha': 0.0001},
             {'alpha': 0.001},
             {'alpha': 0.01},
             {'alpha': 0.1},
             {'alpha': 1},
             {'alpha': 10.0},
             {'alpha': 100.0},
             {'alpha': 1000.0},
             {'alpha': 10000.0}],
 'split0_test_score': array([0.7049816 , 0.69856666, 0.67359334, 0.63546369, 0.63207593,
                        0.63212131, 0.63210724, 0.6321073 , 0.6321084 ]),
 'split1_test_score': array([0.70293816, 0.69310745, 0.67347426, 0.63303196, 0.63099156,
                        0.63123927, 0.63126006, 0.63125986, 0.63125843]),
 'split2_test_score': array([0.70998338, 0.70297658, 0.67673471, 0.6384014 , 0.63448591,
                        0.63480884, 0.63479603, 0.63479387, 0.6347942 ]),
 'split3_test_score': array([0.69703151, 0.68802843, 0.66331077, 0.62745529, 0.62393862,
                        0.62374793, 0.6237467 , 0.62374636, 0.62374636]),
 'split4_test_score': array([0.70188606, 0.69973988, 0.67559658, 0.63951825, 0.63196181,
                        0.63150445, 0.63150572, 0.63150598, 0.63150435]),
 'mean_test_score': array([0.70336418, 0.69648371, 0.67254185, 0.63477399, 0.63069073,
                        0.63068434, 0.63068313, 0.63068266, 0.63068233]),
 'std_test_score': array([0.00421666, 0.00529135, 0.00477705, 0.00430672, 0.00356709,
                        0.00369249, 0.0036897 , 0.00368935, 0.00368939]),
 'rank_test_score': array([1, 2, 3, 4, 5, 6, 7, 8, 9]),
 'split0_train_score': array([0.72567063, 0.71579144, 0.6833411 , 0.63889628, 0.63248786,
                        0.6324593 , 0.63244046, 0.6324397 , 0.63244103]),
 'split1_train_score': array([0.7301668 , 0.71795369, 0.68526338, 0.63825235, 0.6333629
```

```

split1_train_score': array([0.7501000 , 0.71793399, 0.6826330 , 0.6825233 , 0.6333029 ,
0.63335997, 0.63337739, 0.6333761 , 0.63337623]),
'split2_train_score': array([0.72496608, 0.71457002, 0.6820006 , 0.63791973, 0.63188189,
0.63180734, 0.63179626, 0.63179589, 0.63179539]),
'split3_train_score': array([0.73185594, 0.71808259, 0.68651657, 0.64045541, 0.63461542,
0.6344364 , 0.63443559, 0.63443534, 0.63443534]),
'split4_train_score': array([0.72408593, 0.71555978, 0.68377008, 0.63877762, 0.63251977,
0.63249727, 0.63250031, 0.6324997 , 0.63249991]),
'mean_train_score': array([0.72734907, 0.7163915 , 0.68417835, 0.63886028, 0.63297357,
0.63291206, 0.63291 , 0.63290934, 0.63290958]),
'std_train_score': array([0.00307881, 0.00139067, 0.00156556, 0.00087256, 0.00094666,
0.00090793, 0.00091374, 0.00091375, 0.00091373])}

```

In [74]:

```

train_auc3 = grid3.cv_results_['mean_train_score']
train_auc_std3 = grid3.cv_results_['std_train_score']
cv_auc3 = grid3.cv_results_['mean_test_score']
cv_auc_std3 = grid3.cv_results_['std_test_score']

log_k_range = []
for a in tqdm(k_range):

    b = np.log10(a)
    log_k_range.append(b)

plt.plot(log_k_range, train_auc3, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, train_auc3 - train_auc_std3, train_auc3 + train_auc_std3, alpha
=0.2,color='darkblue')

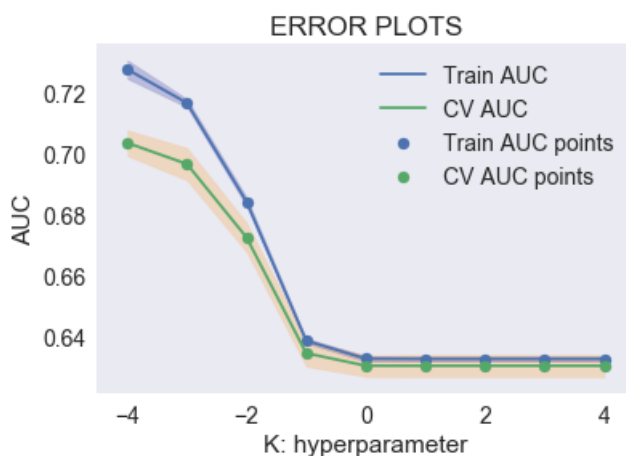
plt.plot(log_k_range, cv_auc3, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, cv_auc3 - cv_auc_std3, cv_auc3 + cv_auc_std3, alpha=0.2,color='
darkorange')

plt.scatter(log_k_range, train_auc3, label='Train AUC points')
plt.scatter(log_k_range, cv_auc3, label='CV AUC points')

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

```

100% |████████████████████████████████████████████████████████████████████████████████| 9/9 [00:00<00:00, 9007.10it/s]



In [75]:

```

# examine the best model

# Single best score achieved across all params (k)
print(grid3.best_score_)

# Dictionary containing the parameters (k) used to generate that score

```

```
# Dictionary containing the parameters (k) used to generate that score
print(grid3.best_params_)
```

```
# Actual model object fit with those best parameters
# Shows default parameters that we did not specify
print(grid3.best_estimator_)
```

```
0.7033641826577617
{'alpha': 0.0001}
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

In [76]:

```
#https://stackoverflow.com/questions/39200265/attributeerror-probability-estimates-are-not-available-for-loss-hinge
```

```
from sklearn.calibration import CalibratedClassifierCV
```

```
meow3 = SGDClassifier(loss = 'hinge',penalty='l2', alpha = 0.0001, class_weight= "balanced")
fit_meow3 = meow3.fit(X3_train_tr, y_train)
calibrator3 = CalibratedClassifierCV(fit_meow3, cv='prefit')

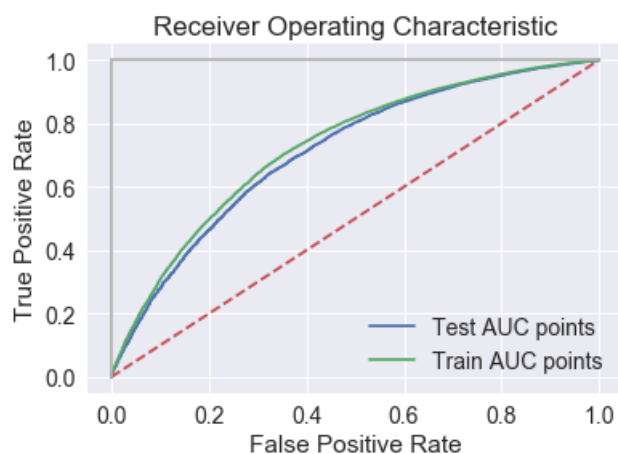
Trained_model_AVGW2V = calibrator3.fit(X3_train_tr, y_train)
```

In [77]:

```
# Get predicted probabilities
y_score_test3 = Trained_model_AVGW2V.predict_proba(X3_test)[: ,1]
y_score_train3 = Trained_model_AVGW2V.predict_proba(X3_train_tr)[: ,1]

# Create true and false positive rates
false_positive_rate3, true_positive_rate3, threshold3 = roc_curve(y_test, y_score_test3)
false_positive_rate31, true_positive_rate31, threshold31 = roc_curve(y_train, y_score_train3)

# Plot ROC curve
plt.title('Receiver Operating Characteristic')
plt.plot(false_positive_rate3, true_positive_rate3, label='Test AUC points')
plt.plot(false_positive_rate31, true_positive_rate31, label='Train AUC points')
plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0], c=".7"), plt.plot([1, 1], c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()
```



In [78]:

```
roc_auc_score(y_test, y_score_test3)
```

Out[78]:

```
0.7112676014561521
```

0.7113076214584791

In [79]:

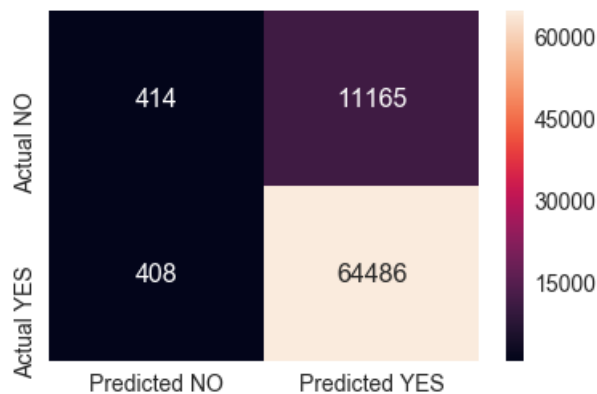
```
# Train model and make predictions
y_hat_3 = Trained_model_AVGW2V.predict(X3_test)

print(classification_report(y_test, y_hat_3))
```

	precision	recall	f1-score	support
0	0.48	0.04	0.07	4963
1	0.85	0.99	0.92	27812
accuracy			0.85	32775
macro avg	0.67	0.51	0.49	32775
weighted avg	0.80	0.85	0.79	32775

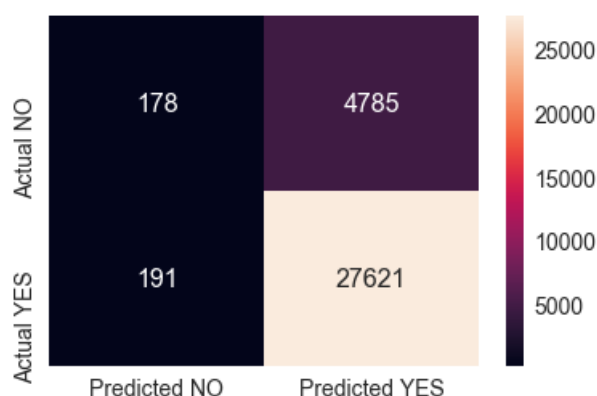
In [80]:

```
get_confusion_matrix(Trained_model_AVGW2V, X3_train_tr, y_train)
```



In [81]:

```
get_confusion_matrix(Trained_model_AVGW2V, X3_test, y_test)
```



In [82]:

```
grid4 = GridSearchCV(model,param_grid=parameters,cv=5,scoring='roc_auc', return_train_score=True)
```

In [83]:

```
grid4.fit(X4_train_tr, y_train)
```

Out[83]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
```

```

        class_weight='balanced',
        early_stopping=False, epsilon=0.1,
        eta0=0.0, fit_intercept=True,
        l1_ratio=0.15, learning_rate='optimal',
        loss='hinge', max_iter=1000,
        n_iter_no_change=5, n_jobs=None,
        penalty='l2', power_t=0.5,
        random_state=None, shuffle=True, tol=0.001,
        validation_fraction=0.1, verbose=0,
        warm_start=False),
    iid='warn', n_jobs=None,
    param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0,
                          1000.0, 10000.0]},
    pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
    scoring='roc_auc', verbose=0)

```

In [84]:

```

# view the complete results (list of named tuples)
grid4.cv_results_

```

Out[84]:

```

{'mean_fit_time': array([8.92487769, 4.57996874, 2.2271214 , 2.01405959, 1.49735017,
        1.38780403, 1.31624551, 1.2348968 , 1.34103799]),
 'std_fit_time': array([0.30316102, 0.45837322, 0.38588888, 0.52726612, 0.22039492,
        0.19939297, 0.06390014, 0.00508728, 0.12034699]),
 'mean_score_time': array([0.05434499, 0.03358841, 0.03357983, 0.03456759, 0.03358374,
        0.03378396, 0.03397865, 0.03378024, 0.03637538]),
 'std_score_time': array([0.0368064 , 0.00101741, 0.00101918, 0.00185819, 0.00101343,
        0.0011595 , 0.00167303, 0.00074273, 0.00242101]),
 'param_alpha': masked_array(data=[0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0,
        10000.0],
        mask=[False, False, False, False, False, False, False, False,
        False],
        fill_value='?',
        dtype=object),
 'params': [{'alpha': 0.0001},
        {'alpha': 0.001},
        {'alpha': 0.01},
        {'alpha': 0.1},
        {'alpha': 1},
        {'alpha': 10.0},
        {'alpha': 100.0},
        {'alpha': 1000.0},
        {'alpha': 10000.0}],
 'split0_test_score': array([0.69779076, 0.7037028 , 0.69224334, 0.65279929, 0.644538 ,
        0.6441 , 0.64410237, 0.64410356, 0.6441026 ]),
 'split1_test_score': array([0.69274503, 0.69657309, 0.68620833, 0.64750515, 0.640926 ,
        0.6407242 , 0.64072191, 0.64072121, 0.64072121]),
 'split2_test_score': array([0.70611234, 0.70948667, 0.69794467, 0.6538139 , 0.64627984,
        0.64622718, 0.64622435, 0.64622435, 0.64622528]),
 'split3_test_score': array([0.69194857, 0.69209921, 0.6801414 , 0.64164664, 0.63477883,
        0.63442127, 0.63441861, 0.63441861, 0.63441941]),
 'split4_test_score': array([0.69449532, 0.70061027, 0.69207964, 0.65539446, 0.64156744,
        0.64036397, 0.64035808, 0.64035848, 0.64036014]),
 'mean_test_score': array([0.69661846, 0.7004944 , 0.68972341, 0.65023175, 0.64161803,
        0.64116735, 0.64116508, 0.64116526, 0.64116575]),
 'std_test_score': array([0.00515422, 0.00595027, 0.00606079, 0.00504412, 0.00393743,
        0.00401387, 0.00401468, 0.00401486, 0.00401462]),
 'rank_test_score': array([2, 1, 3, 4, 5, 6, 9, 8, 7]),
 'split0_train_score': array([0.72540826, 0.72156531, 0.70181442, 0.65545016, 0.64356107,
        0.6430791 , 0.64308432, 0.64308506, 0.64308453]),
 'split1_train_score': array([0.72238667, 0.72398102, 0.70349772, 0.65525741, 0.64490374,
        0.64442173, 0.64442218, 0.64442206, 0.64442206]),
 'split2_train_score': array([0.72261761, 0.71999585, 0.69962756, 0.65237719, 0.64260876,
        0.64211907, 0.64211907, 0.64211907, 0.64212004]),
 'split3_train_score': array([0.72741192, 0.72401599, 0.70457976, 0.65647642, 0.64610578,
        0.64564186, 0.6456391 , 0.64563909, 0.64564029]),
 'split4_train_score': array([0.72336057, 0.72238214, 0.70180287, 0.65458818, 0.64367912,
        0.64307975, 0.64307914, 0.64307978, 0.6430796 ]),
 'mean_train_score': array([0.72423701, 0.72238806, 0.70226447, 0.65482987, 0.64417169,
        0.6436683 , 0.64366876, 0.64366901, 0.6436693 ]),
 'std_train_score': array([0.00191135, 0.00152232, 0.00168822, 0.00136768, 0.00121118,
        0.00122933, 0.00122806, 0.00122791, 0.00122812])}

```

In [85]:

```
train_auc4 = grid4.cv_results_['mean_train_score']
train_auc_std4 = grid4.cv_results_['std_train_score']
cv_auc4 = grid4.cv_results_['mean_test_score']
cv_auc_std4 = grid4.cv_results_['std_test_score']

log_k_range = []
for a in tqdm(k_range):

    b = np.log10(a)
    log_k_range.append(b)

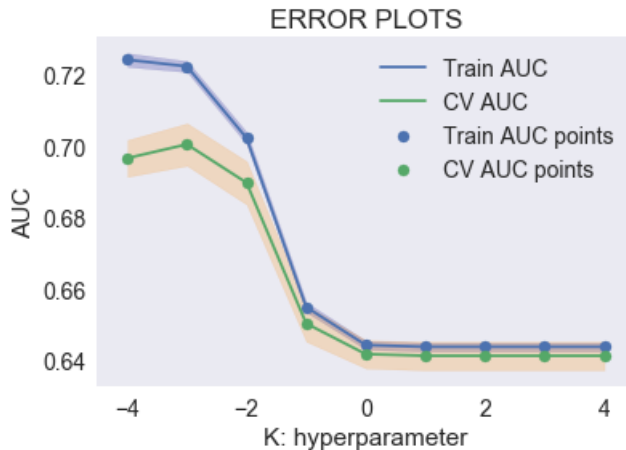
plt.plot(log_k_range, train_auc4, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, train_auc4 - train_auc_std4, train_auc4 + train_auc_std4, alpha=0.2, color='darkblue')

plt.plot(log_k_range, cv_auc4, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, cv_auc4 - cv_auc_std4, cv_auc4 + cv_auc_std4, alpha=0.2, color='darkorange')

plt.scatter(log_k_range, train_auc4, label='Train AUC points')
plt.scatter(log_k_range, cv_auc4, label='CV AUC points')

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

```
100%|██████████| 9/9 [00:00<00:00, 9002.80it/s]
```



In [86]:

```
# examine the best model

# Single best score achieved across all params (k)
print(grid4.best_score_)

# Dictionary containing the parameters (k) used to generate that score
print(grid4.best_params_)

# Actual model object fit with those best parameters
# Shows default parameters that we did not specify
print(grid4.best_estimator_)

0.7004944028625657
{'alpha': 0.001}
SGDClassifier(alpha=0.001, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
```



```
max_iter=1000, n_iter_no_change=5, n_jobs=None, penalty='l2',
power_t=0.5, random_state=None, shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0, warm_start=False)
```

In [87]:

```
#https://stackoverflow.com/questions/39200265/attributeerror-probability-estimates-are-not-available-for-loss-hinge
from sklearn.calibration import CalibratedClassifierCV

meow4 = SGDClassifier(loss = 'hinge',penalty='l2', alpha = 0.001, class_weight= "balanced")
fit_meow4 = meow4.fit(X4_train_tr, y_train)
calibrator4 = CalibratedClassifierCV(fit_meow4, cv='prefit')

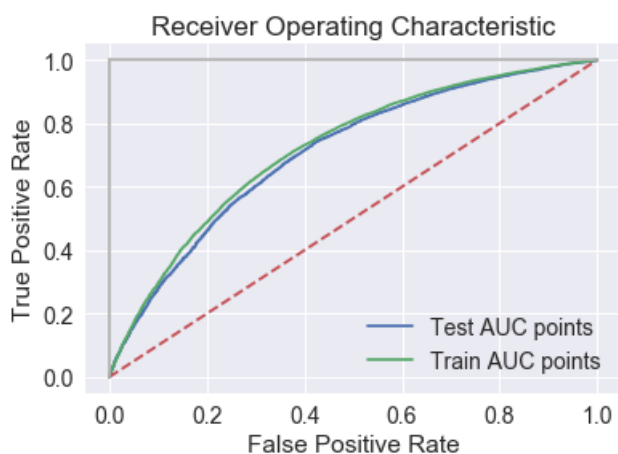
Trained_model_TFIDFW2V = calibrator4.fit(X4_train_tr, y_train)
```

In [88]:

```
# Get predicted probabilities
y_score_test4 = Trained_model_TFIDFW2V.predict_proba(X4_test)[:,:1]
y_score_train4 = Trained_model_TFIDFW2V.predict_proba(X4_train_tr)[:,:1]

# Create true and false positive rates
false_positive_rate4, true_positive_rate4, threshold4 = roc_curve(y_test, y_score_test4)
false_positive_rate41, true_positive_rate41, threshold41 = roc_curve(y_train, y_score_train4)

# Plot ROC curve
plt.title('Receiver Operating Characteristic')
plt.plot(false_positive_rate4, true_positive_rate4, label='Test AUC points')
plt.plot(false_positive_rate41, true_positive_rate41, label='Train AUC points')
plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0], c=".7"), plt.plot([1, 1], c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()
```



In [89]:

```
roc_auc_score(y_test, y_score_test4)
```

Out[89]:

```
0.7073230442597238
```

In [90]:

```
y_hat_4 = Trained_model_AVGW2V.predict(X4_test)

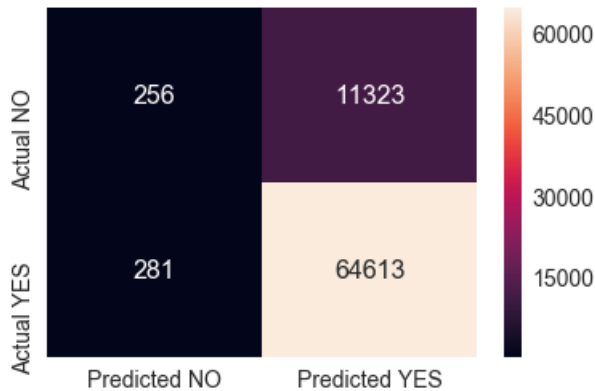
print(classification_report(y_test, y_hat_4))
```

```
precision    recall  f1-score   support
```

0	0.45	0.05	0.10	4963
1	0.85	0.99	0.92	27812
accuracy			0.85	32775
macro avg	0.65	0.52	0.51	32775
weighted avg	0.79	0.85	0.79	32775

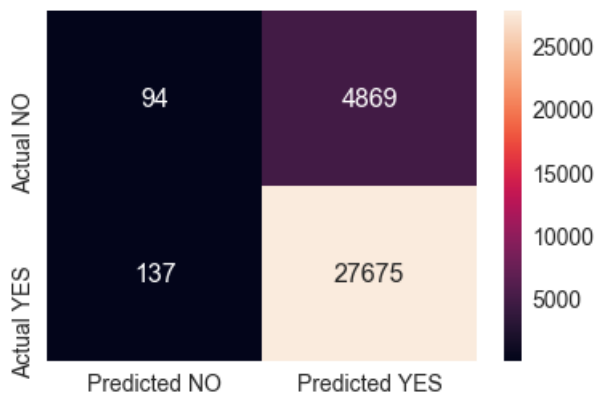
In [91]:

```
get_confusion_matrix(Trained_model_TFIDFW2V, X4_train_tr, y_train)
```



In [92]:

```
get_confusion_matrix(Trained_model_TFIDFW2V, X4_test, y_test)
```



2.5 Support Vector Machines with added Features `Set 5`

In [93]:

```
#randomly removing 50,000 features from the dataset
remove_n = 60000
drop_indices = np.random.choice(project_data.index, remove_n, replace=False)
project_data = project_data.drop(drop_indices)
```

In [94]:

```
print(project_data.shape)
```

(49248, 25)

In [95]:

```
y13 = project_data['project_is_approved']
X_train_new, X_test_new, y_train_new, y_test_new = train_test_split(project_data, y13, stratify=y13,
, test_size=0.3)
```

In [96]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_train_tr_1 = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
binary=True)
categories_one_hot_train = vectorizer_train_tr_1.fit(X_train_new['clean_categories'].values)
#print(vectorizer_train_tr_1.get_feature_names())
#print("Shape of matrix after one hot encodig ",categories_one_hot_train_tr.shape)

categories_one_hot_train_tr = categories_one_hot_train.transform(X_train_new['clean_categories'].v
alues)
#print(vectorizer_train_tr_1.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_train_tr.shape)

#feature names encoding for X_test
#from sklearn.feature_extraction.text import CountVectorizer
#vectorizer_test_1 = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, bin
ary=True)
categories_one_hot_test = categories_one_hot_train.transform(X_test_new['clean_categories'].values
)
#print(vectorizer_test_1.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot_test.shape)

# we use count vectorizer to convert the values into one for X_train
vectorizer_train_tr = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False
, binary=True)
sub_categories_one_hot_train = vectorizer_train_tr.fit(X_train_new['clean_subcategories'].values)
#print(vectorizer_train_tr.get_feature_names())
#print("Shape of matrix after one hot encodig ",sub_categories_one_hot_train_tr.shape)

sub_categories_one_hot_train_tr =
sub_categories_one_hot_train.transform(X_train_new['clean_subcategories'].values)
#print(vectorizer_train_tr.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_train_tr.shape)

# we use count vectorizer to convert the values into one for X_test
vectorizer_test = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, bi
nary=True)
sub_categories_one_hot_test =
sub_categories_one_hot_train.transform(X_test_new['clean_subcategories'].values)
#print(vectorizer_test.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_test.shape)

#school state for Xtrain
vectorizer_1_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_1_train_tr.fit(X_train_new['school_state'].values)
#print(vectorizer_1_train_tr.get_feature_names())
categories_state_1_train = vectorizer_1_train_tr.fit(X_train_new['school_state'].values)
#print("Shape of matrix after one hot encodig ",categories_state_1_train_tr.shape)

categories_state_1_train_tr =
categories_state_1_train.transform(X_train_new['school_state'].values)
print("Shape of matrix after one hot encodig ",categories_state_1_train_tr.shape)

#school state for Xtest
vectorizer_1_test = CountVectorizer(lowercase=False, binary=True)
vectorizer_1_test.fit(X_test_new['school_state'].values)
#print(vectorizer_1_test.get_feature_names())
categories_state_1_test = categories_state_1_train.transform(X_test_new['school_state'].values)
print("Shape of matrix after one hot encodig ",categories_state_1_test.shape)

#Project grade category for X_train
#vectorizer_2_train = CountVectorizer(vocabulary=list(word_dict_134.keys()), lowercase=False, bina
ry=True)
vectorizer_2_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_2_train_tr.fit(X_train_new['project_grade_category'].values)
#print(vectorizer_2_train_tr.get_feature_names())
categories_grade_train = vectorizer_2_train_tr.fit(X_train_new['project_grade_category'].values)
#print("Shape of matrix after one hot encodig ",categories_grade_train_tr.shape)

categories_grade_train_tr = categories_grade_train.transform(X_train_new['project_grade_category']
```

```

.values)
print("Shape of matrix after one hot encodig ",categories_grade_train_tr.shape)

#Project grade category for X_test
#vectorizer_2_test = CountVectorizer(vocabulary=list(word_dict_134.keys()), lowercase=False, binary=True)
#vectorizer_2_test = CountVectorizer(lowercase=False, binary=True)
#vectorizer_2_test.fit(X_test['project_grade_category'].values)
#print(vectorizer_2_test.get_feature_names())
categories_grade_test =
categories_grade_train.transform(X_test_new['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",categories_grade_test.shape)

from string import punctuation

#https://medium.com/@chaimgluck1/have-messy-text-data-clean-it-with-simple-lambda-functions-645918fcc2fc
#project_data.teacher_prefix = project_data.teacher_prefix.apply(lambda x:
x.translate(string.punctuation))

#https://stackoverflow.com/questions/50443494/error-in-removing-punctuation-float-object-has-no-attribute-translate
#project_data['teacher_prefix'] = project_data.fillna({'teacher_prefix':''})

project_data['teacher_prefix'] = project_data['teacher_prefix'].replace(np.nan, 'teacher')

#teacher_prefix for X_train
vectorizer_3_train_tr = CountVectorizer(lowercase=False, binary=True)
vectorizer_3_train_tr.fit(X_train_new['teacher_prefix'].values)
#print(vectorizer_3_train_tr.get_feature_names())
categories_teacher_prefix_train = vectorizer_3_train_tr.fit(X_train_new['teacher_prefix'].values)
#print("Shape of matrix after one hot encodig ",categories_teacher_prefix_train_tr.shape)

categories_teacher_prefix_train_tr =
categories_teacher_prefix_train.transform(X_train_new['teacher_prefix'].values)
print("Shape of matrix after one hot encodig ",categories_teacher_prefix_train_tr.shape)

#teacher_prefix for X_test
#vectorizer_3_test = CountVectorizer(lowercase=False, binary=True)
#vectorizer_3_test.fit(X_test['teacher_prefix'].values)
#print(vectorizer_3_test.get_feature_names())
categories_teacher_prefix_test =
categories_teacher_prefix_train.transform(X_test_new['teacher_prefix'].values)
print("Shape of matrix after one hot encodig ",categories_teacher_prefix_test.shape)

```

```

Shape of matrix after one hot encodig (34473, 9)
Shape of matrix after one hot encodig (14775, 9)
Shape of matrix after one hot encodig (34473, 30)
Shape of matrix after one hot encodig (14775, 30)
Shape of matrix after one hot encodig (34473, 51)
Shape of matrix after one hot encodig (14775, 51)
Shape of matrix after one hot encodig (34473, 4)
Shape of matrix after one hot encodig (14775, 4)
Shape of matrix after one hot encodig (34473, 6)
Shape of matrix after one hot encodig (14775, 6)

```

In [97]:

```

from sklearn.preprocessing import StandardScaler

#price scalar and standardized for train
price_scalar_train_tr = StandardScaler()
price_scalar_train_tr.fit(X_train_new['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {price_scalar_train_tr.mean_[0]}, Standard deviation :
{np.sqrt(price_scalar_train_tr.var_[0])}")
price_standardized_train_tr = price_scalar_train_tr.transform(X_train_new['price'].values.reshape(-1, 1))

#price scalar and standardized for test
price_scalar_test = StandardScaler()
price_scalar_test.fit(X_test_new['price'].values.reshape(-1,1)) # finding the mean and standard deviation of this data

```

```

deviation of this data
#print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")
price_standardized_test = price_scalar_test.transform(X_test_new['price'].values.reshape(-1, 1))

teacher_number_of_previously_posted_projects_scalar_train_tr = StandardScaler()
teacher_number_of_previously_posted_projects_scalar_train_tr.fit(X_train_new['teacher_number_of_previously_posted_projects'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar_train_tr.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar_train_tr.var_[0])}")
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_train_tr =
teacher_number_of_previously_posted_projects_scalar_train_tr.transform(X_train_new['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

teacher_number_of_previously_posted_projects_scalar_test = StandardScaler()
teacher_number_of_previously_posted_projects_scalar_test.fit(X_test_new['teacher_number_of_previously_posted_projects'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar.var_[0])}")
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_test =
teacher_number_of_previously_posted_projects_scalar_test.transform(X_test_new['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

```

In [98]:

```

essay_sentiment_score = StandardScaler()
essay_sentiment_score.fit(X_train_new['essay_sentiment'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar_train_tr.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar_train_tr.var_[0])}")
# Now standardize the data with above mean and variance.
essay_sentiment_score_train =
essay_sentiment_score.transform(X_train_new['essay_sentiment'].values.reshape(-1, 1))
essay_sentiment_score_test = essay_sentiment_score.transform(X_test_new['essay_sentiment'].values.reshape(-1, 1))
print(essay_sentiment_score_train.shape)
print(essay_sentiment_score_test.shape)

title_count = StandardScaler()
title_count.fit(X_train_new['title_count'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar_train_tr.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar_train_tr.var_[0])}")
# Now standardize the data with above mean and variance.
title_count_train = title_count.transform(X_train_new['title_count'].values.reshape(-1, 1))
title_count_test = title_count.transform(X_test_new['title_count'].values.reshape(-1, 1))
print(title_count_train.shape)
print(title_count_test.shape)

essay_count = StandardScaler()
essay_count.fit(X_train_new['essay_count'].values.reshape(-1,1)) # finding the mean and standard deviation of this data
#print(f"Mean : {teacher_number_of_previously_posted_projects_scalar_train_tr.mean_[0]}, Standard deviation : {np.sqrt(teacher_number_of_previously_posted_projects_scalar_train_tr.var_[0])}")
# Now standardize the data with above mean and variance.
essay_count_train = essay_count.transform(X_train_new['essay_count'].values.reshape(-1, 1))
essay_count_test = essay_count.transform(X_test_new['essay_count'].values.reshape(-1, 1))
print(essay_count_train.shape)
print(essay_count_test.shape)

```

```

(34473, 1)
(14775, 1)
(34473, 1)
(14775, 1)
(34473, 1)
(14775, 1)

```

In [99]:

```

#Vectorizing using tfidf for essays for x_train
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10, ngram_range=(2,2), max_features=5000)

```

```

text_tfidf_train = vectorizer_tfidf_essays_train_tr.fit(X_train_new['preprocessed_essays'])
#print("Shape of matrix after one hot encodig ",text_tfidf_train_tr.shape)

text_tfidf_train_tr = text_tfidf_train.transform(X_train_new['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_tfidf_train_tr.shape)

text_tfidf_test_tr = text_tfidf_train.transform(X_test_new['preprocessed_essays'])
print("Shape of matrix after one hot encodig ",text_tfidf_test_tr.shape)

```

```

Shape of matrix after one hot encodig (34473, 5000)
Shape of matrix after one hot encodig (14775, 5000)

```

In [100]:

```

#Tfidf vectorization of essays
text_tfidf_train_tr

```

Out[100]:

```

<34473x5000 sparse matrix of type '<class 'numpy.float64'>'
  with 1458431 stored elements in Compressed Sparse Row format>

```

In [101]:

```

from sklearn.decomposition import TruncatedSVD
svd = TruncatedSVD()
dim_range = [2,5,10,50,100,200,500,700,1000,2000,2500,4000,4500]
var_range = []
for i in dim_range:
    svd = TruncatedSVD(n_components = i, random_state = 42, n_iter = 5)
    new_text_tfidf_train_tr = svd.fit(text_tfidf_train_tr)
    var_range.append(new_text_tfidf_train_tr.explained_variance_ratio_.sum())

```

In [102]:

```

var_range

```

Out[102]:

```

[0.005871345012842383,
 0.014556153312884356,
 0.025152786025342783,
 0.07857288168614358,
 0.12291813458431945,
 0.19063259438460384,
 0.33031854945247874,
 0.40086104491839203,
 0.48891197683760557,
 0.7002624219092254,
 0.7774592454700449,
 0.9399950950408562,
 0.9759679645766187]

```

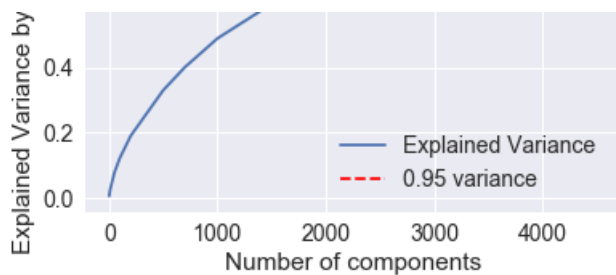
In [103]:

```

plt.plot(dim_range,var_range,label = 'Explained Variance ')
plt.axhline(0.95,linestyle = '--',color = 'r',label = '0.95 variance')
plt.xlabel("Number of components")
plt.ylabel("Explained Variance by the Components")
plt.legend()
plt.title("Components vs Explained Variance")
sns.despine()
plt.show()

```





In [104]:

```
from sklearn.decomposition import TruncatedSVD
svd_model = TruncatedSVD(n_components = 4100, n_iter = 3)
svd_model.fit(text_tfidf_train_tr)
svd_train = svd_model.transform(text_tfidf_train_tr)
```

In [105]:

```
svd_test = svd_model.transform(text_tfidf_test_tr)
```

In [106]:

```
#categorical, numerical features + project title(BOW)
X5_train_tr = hstack((categories_one_hot_train_tr, sub_categories_one_hot_train_tr,
categories_state_1_train_tr, categories_grade_train_tr, categories_teacher_prefix_train_tr,
price_standardized_train_tr, teacher_number_of_previously_posted_projects_standardized_train_tr, e
ssay_sentiment_score_train, title_count_train, essay_count_train, svd_train))

X5_test = hstack((categories_one_hot_test, sub_categories_one_hot_test, categories_state_1_test, c
ategories_grade_test, categories_teacher_prefix_test, price_standardized_test,
teacher_number_of_previously_posted_projects_standardized_test, essay_sentiment_score_test,
title_count_test, essay_count_test, svd_test))
```

In [107]:

```
grid5 = GridSearchCV(model, param_grid=parameters, cv=5, scoring='roc_auc', return_train_score=True)
```

In [108]:

```
grid5.fit(X5_train_tr, y_train_new)
```

Out[108]:

```
GridSearchCV(cv=5, error_score='raise-deprecating',
             estimator=SGDClassifier(alpha=0.0001, average=False,
                                     class_weight='balanced',
                                     early_stopping=False, epsilon=0.1,
                                     eta0=0.0, fit_intercept=True,
                                     l1_ratio=0.15, learning_rate='optimal',
                                     loss='hinge', max_iter=1000,
                                     n_iter_no_change=5, n_jobs=None,
                                     penalty='l2', power_t=0.5,
                                     random_state=None, shuffle=True, tol=0.001,
                                     validation_fraction=0.1, verbose=0,
                                     warm_start=False),
             iid='warn', n_jobs=None,
             param_grid={'alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0,
                                     1000.0, 10000.0]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring='roc_auc', verbose=0)
```

In [109]:

```
grid5.cv_results_
```

Out[109]:

```
{'mean_fit_time': array([175.26609268, 8.6209487, 4.88827763, 3.79143233,
4.11223736, 3.79460063, 3.82012134, 3.82371998,
3.85928459]),
'std_fit_time': array([2.40155958e+02, 8.02372131e-01, 9.34102227e-01, 1.34289739e-01,
2.25603080e-01, 1.56599681e-01, 2.43072095e-01, 9.03390501e-02,
1.55325280e-01]),
'mean_score_time': array([13.66900005, 0.55801072, 0.83631744, 0.20223174, 0.1425518,
0.1085331, 0.23648705, 0.17545238, 0.14491625]),
'std_score_time': array([2.43148483e+01, 3.54282205e-01, 6.52147423e-01, 8.70420178e-02,
6.64972278e-02, 3.93283392e-03, 2.61109585e-01, 1.35048990e-01,
6.96478263e-02]),
'param_alpha': masked_array(data=[0.0001, 0.001, 0.01, 0.1, 1, 10.0, 100.0, 1000.0,
10000.0],
mask=[False, False, False, False, False, False, False, False,
False],
fill_value='?',
dtype=object),
'params': [{'alpha': 0.0001},
{'alpha': 0.001},
{'alpha': 0.01},
{'alpha': 0.1},
{'alpha': 1},
{'alpha': 10.0},
{'alpha': 100.0},
{'alpha': 1000.0},
{'alpha': 10000.0}],
'split0_test_score': array([0.66968817, 0.65584912, 0.63775104, 0.63050196, 0.6276125,
0.6283094, 0.62830711, 0.62830711, 0.62830711]),
'split1_test_score': array([0.66429896, 0.6555879, 0.62974008, 0.61857832, 0.61301061,
0.61460464, 0.61460464, 0.61460464, 0.61460464]),
'split2_test_score': array([0.66662933, 0.65071684, 0.63172921, 0.62198464, 0.62249975,
0.62378868, 0.62378868, 0.62378868, 0.62378868]),
'split3_test_score': array([0.66662606, 0.66654223, 0.64073976, 0.63559878, 0.63745391,
0.63798932, 0.63798932, 0.63798932, 0.63798932]),
'split4_test_score': array([0.65214232, 0.64415987, 0.61022628, 0.59957314, 0.59608,
0.59705308, 0.59706061, 0.59706061, 0.59706061]),
'mean_test_score': array([0.66387732, 0.6545713, 0.63003771, 0.62124783, 0.61933165,
0.62034932, 0.62035037, 0.62035037, 0.62035037]),
'std_test_score': array([0.0061119, 0.00733759, 0.01067183, 0.01240169, 0.0140539,
0.01387096, 0.01386816, 0.01386816, 0.01386816]),
'rank_test_score': array([1, 2, 3, 4, 9, 8, 5, 5, 5]),
'split0_train_score': array([0.81644632, 0.70358128, 0.64440142, 0.62559396, 0.6222035,
0.62349644, 0.62349741, 0.62349741, 0.62349741]),
'split1_train_score': array([0.82169643, 0.70693274, 0.64284618, 0.62549747, 0.62190695,
0.62293458, 0.62293458, 0.62293458, 0.62293458]),
'split2_train_score': array([0.80810108, 0.70087706, 0.64192321, 0.6240468, 0.62060157,
0.62154294, 0.62154294, 0.62154294, 0.62154294]),
'split3_train_score': array([0.81815808, 0.71221338, 0.64089753, 0.62383753, 0.6198298,
0.62068942, 0.62068942, 0.62068942, 0.62068942]),
'split4_train_score': array([0.80705007, 0.71334407, 0.64666341, 0.63022362, 0.62575233,
0.62657926, 0.62658907, 0.62658907, 0.62658907]),
'mean_train_score': array([0.8142904, 0.70738971, 0.64334635, 0.62583988, 0.62205883,
0.62304853, 0.62305068, 0.62305068, 0.62305068]),
'std_train_score': array([0.0057478, 0.00481352, 0.00201932, 0.00230735, 0.00203854,
0.00202543, 0.0020289, 0.0020289, 0.0020289])}
```

In [110]:

```
pd.DataFrame(grid5.cv_results_)
```

Out[110]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_sco
0	175.266093	240.155958	13.669000	24.314848	0.0001	{'alpha': 0.0001}	0.669688	0.664299
1	8.620949	0.802372	0.558011	0.354282	0.001	{'alpha': 0.001}	0.655849	0.655588
2	4.888278	0.934102	0.836317	0.652147	0.01	{'alpha': 0.01}	0.637751	0.629740
3	3.791432	0.134290	0.202232	0.087042	0.1	{'alpha': 0.1}	0.630502	0.618578

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_sco
4	4.112237	0.225603	0.142552	0.066497	1	{'alpha': 1}	0.627613	0.613011
5	3.794601	0.156600	0.108533	0.003933	10	{'alpha': 10.0}	0.628309	0.614605
6	3.820121	0.243072	0.236487	0.261110	100	{'alpha': 100.0}	0.628307	0.614605
7	3.823720	0.090339	0.175452	0.135049	1000	{'alpha': 1000.0}	0.628307	0.614605
8	3.859285	0.155325	0.144916	0.069648	10000	{'alpha': 10000.0}	0.628307	0.614605

9 rows × 21 columns



In [111]:

```
# examine the best model

# Single best score achieved across all params (k)
print(grid5.best_score_)

# Dictionary containing the parameters (k) used to generate that score
print(grid5.best_params_)

# Actual model object fit with those best parameters
# Shows default parameters that we did not specify

print(grid5.best_estimator_)
```

```
0.6638773173916696
{'alpha': 0.0001}
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=None, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

In [112]:

```
train_auc2 = grid5.cv_results_['mean_train_score']
train_auc_std2 = grid5.cv_results_['std_train_score']
cv_auc2 = grid5.cv_results_['mean_test_score']
cv_auc_std2 = grid5.cv_results_['std_test_score']

log_k_range = []
for a in tqdm(k_range):

    b = np.log10(a)
    log_k_range.append(b)

plt.plot(log_k_range, train_auc2, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, train_auc2 - train_auc_std2, train_auc2 + train_auc_std2, alpha=
0.2, color='darkblue')

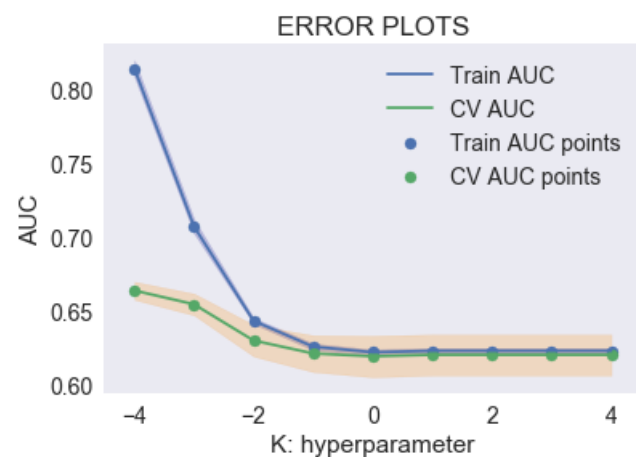
plt.plot(log_k_range, cv_auc2, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(log_k_range, cv_auc2 - cv_auc_std2, cv_auc2 + cv_auc_std2, alpha=0.2, color='darkorange')

plt.scatter(log_k_range, train_auc2, label='Train AUC points')
plt.scatter(log_k_range, cv_auc2, label='CV AUC points')

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

```
plt.show()
```

100% | 9/9 [00:00<00:00, 12.49it/s]



In [113]:

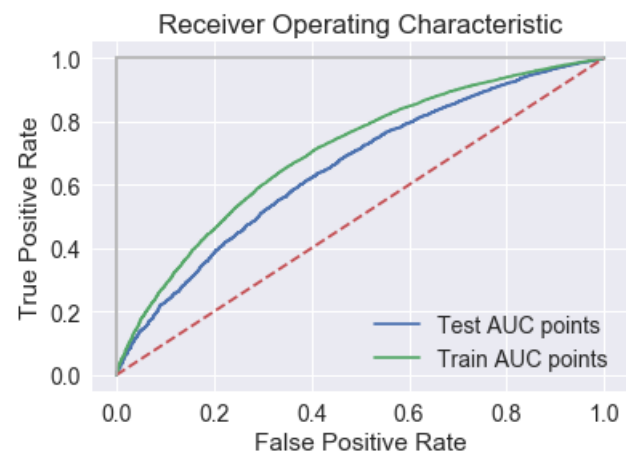
```
meowl = SGDClassifier(loss = 'hinge',penalty='l2', alpha = 0.001, class_weight= "balanced")
fit_meowl = meowl.fit(X5_train_tr, y_train_new)
calibrator = CalibratedClassifierCV(fit_meowl, cv='prefit')
Trained_model_FINAL = calibrator.fit(X5_train_tr, y_train_new)
```

In [114]:

```
# Get predicted probabilities
y_score_test5 = Trained_model_FINAL.predict_proba(X5_test)[: ,1]
y_score_train5 = Trained_model_FINAL.predict_proba(X5_train_tr)[: ,1]

# Create true and false positive rates
false_positive_rate5, true_positive_rate5, threshold5 = roc_curve(y_test_new, y_score_test5)
false_positive_rate51, true_positive_rate51, threshold51 = roc_curve(y_train_new, y_score_train5)

# Plot ROC curve
plt.title('Receiver Operating Characteristic')
plt.plot(false_positive_rate5, true_positive_rate5, label='Test AUC points')
plt.plot(false_positive_rate51, true_positive_rate51, label='Train AUC points')
plt.plot([0, 1], ls="--")
plt.plot([0, 0], [1, 0] , c=".7"), plt.plot([1, 1] , c=".7")
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.legend()
plt.show()
```



In [115]:

```
roc_auc_score(y_test_new, y_score_test5)
```

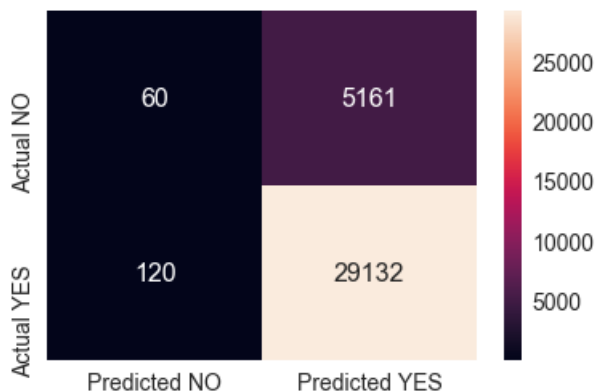
Out[115]:

Out[115]:

0.6523915531027964

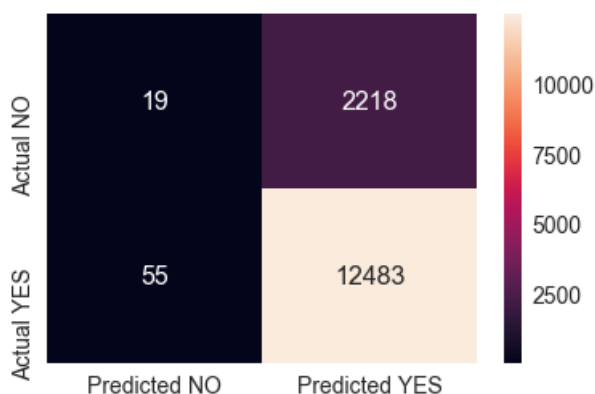
In [116]:

```
get_confusion_matrix(Trained_model_FINAL, X5_train_tr, y_train_new)
```



In [117]:

```
get_confusion_matrix(Trained_model_FINAL, X5_test, y_test_new)
```



3. Conclusion

In [118]:

```
# Please compare all your models using Prettytable library
from prettytable import PrettyTable

DBZ = PrettyTable()
DBZ.field_names = ["Vectorizer", "Model", "Hyperparameter", "Train-AUC", "Test-AUC"]

DBZ.add_row(["BOW", "SVM", grid.best_params_, grid.best_score_, roc_auc_score(y_test, y_score_test
)])
DBZ.add_row(["TF-IDF", "SVM", grid2.best_params_, grid2.best_score_, roc_auc_score(y_test, y_score
_test2)])
DBZ.add_row(["AVGW2V", "SVM", grid3.best_params_, grid3.best_score_, roc_auc_score(y_test, y_score
_test3)])
DBZ.add_row(["TFIDFW2V", "SVM", grid4.best_params_, grid4.best_score_, roc_auc_score(y_test, y_sco
re_test4)])
DBZ.add_row(["No text data", "SVM", grid5.best_params_, grid5.best_score_, roc_auc_score(y_test_ne
w, y_score_test5)])

print(DBZ)
```

Vectorizer	Model	Hyperparameter	Train-AUC	Test-AUC
BOW	SVM	{'alpha': 0.01}	0.6912360212476011	0.6966059338167592
TF-IDF	SVM	{'alpha': 0.0001}	0.6841186385013454	0.6665373816580681

TFIDF	SVM	{'alpha': 0.0001}	0.6611188888888889	0.6666678888888889
AVGW2V	SVM	{'alpha': 0.0001}	0.7033641826577617	0.7113076214584791
TFIDFW2V	SVM	{'alpha': 0.001}	0.7004944028625657	0.7073230442597238
No text data	SVM	{'alpha': 0.0001}	0.6638773173916696	0.6523915531027964
+-----+-----+-----+-----+-----+				