

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• Literature & Writing, Social Sciences

Feature	Description
<code>project_resource_summary</code>	Description 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."
- `__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 chart_studio 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('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

In [3]:

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

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

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

In [4]:

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

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

Out[4]:

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

2. Preprocessing

2.1 preprocessing of project_subject_categories

In [5]:

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

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

2.2 preprocessing of project_subject_subcategories

In [6]:

```
sub_catogories = 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_catogories:
```

```

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

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

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

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

```

2.3 Text preprocessing of essay

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
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10

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(project_data["essay"].values[0])
print("="*50)
print(project_data["essay"].values[150])
print("="*50)
print(project_data["essay"].values[1000])
print("="*50)
print(project_data["essay"].values[20000])
print("="*50)
print(project_data["essay"].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. 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 along side of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills. By providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills. Parents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.

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 group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still. nannan

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.

My class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.

They 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.

Your generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.

It 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!

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in

a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\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. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\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', ' ')
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive de

lays, 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", \
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself'
, \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 't
heir', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these',
'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'd
o', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'whil
e', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'bef
ore', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'a
gain', '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]:

```
# Combining all the above students
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\n', ' ')
```



```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248 [01:31<00:00, 1196.48it/s]
```

```
# 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 deserves nannan'

```
project_data['clean_essay'] = preprocessed_essays
project_data.drop(['essay'], axis=1, inplace=True)
project_data.head(2)
```

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

```
# similarly you can preprocess the titles also
```

```
# printing some random reviews
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)
```

Educational Support for English Learners at Home

More Movement with Hokki Stools

Sailing Into a Super 4th Grade Year

We Need To Move It While We Input It!

Inspiring Minds by Enhancing the Educational Experience

In [21]:

```
# 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 [22]:

```
# Combining all the above students
from tqdm import tqdm
preprocessed_titles = []
# tqdm is for printing the status bar

# https://gist.github.com/sebleier/554280

for sentence in tqdm(project_data['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_titles.append(sent.lower().strip())
```

100% |██| 109248/109248 [00:03<00:00, 27344.39it/s]

In [23]:

```
# after preprocesing
preprocessed_titles[20000]
```

Out[23]:

'need move input'

In [24]:

```
project_data['clean_project_title'] = preprocessed_titles
project_data.drop(['project_title'], axis=1, inplace=True)
project_data.head(2)
```

Out[24]:

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

2.5 Cleaning data of project_grade_category

In [25]:

```
#cleaning project_grade_category

grades = list(project_data['project_grade_category'].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
grade_list = []
for i in grades:
    i = i.replace('-', '_')
    i = i.replace(' ', '')

    grade_list.append(i)
```

In [26]:

```
project_data['clean_grade_category'] = grade_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data.head(2)
```

Out[26]:

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

2.6 Dropping unnecessary columns

In [27]:

```
#project_data.drop(['id'], axis=1, inplace=True)
project_data.drop(['teacher_id'], axis=1, inplace=True)
project_data.drop(['project_essay_1'], axis=1, inplace=True)
project_data.drop(['project_essay_2'], axis=1, inplace=True)
project_data.drop(['project_essay_3'], axis=1, inplace=True)
project_data.drop(['project_essay_4'], axis=1, inplace=True)
```

```
project_data.drop(['project_id'], axis=1, inplace=True)
project_data.drop(['project_resource_summary'], axis=1, inplace=True)
project_data.drop(['Unnamed: 0'], axis=1, inplace=True)
project_data.head(2)
```

Out[27]:

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
0	p253737	Mrs.	IN	2016-12-05 13:43:57	0
1	p258326	Mr.	FL	2016-10-25 09:22:10	7

2.7 Adding price column in our dataframe

In [28]:

```
resource_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1541272 entries, 0 to 1541271
Data columns (total 4 columns):
id                1541272 non-null object
description       1540980 non-null object
quantity         1541272 non-null int64
price            1541272 non-null float64
dtypes: float64(1), int64(1), object(2)
memory usage: 47.0+ MB
```

In [29]:

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

Out[29]:

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
0	p253737	Mrs.	IN	2016-12-05 13:43:57	0
1	p258326	Mr.	FL	2016-10-25 09:22:10	7

In [30]:

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

In [31]:

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

Out[31]:

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
0	p253737	Mrs.	IN	2016-12-05 13:43:57	0
1	p258326	Mr.	FL	2016-10-25 09:22:10	7

2.8 Adding quantity column in our dataframe

In [32]:

```
resource_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1541272 entries, 0 to 1541271
Data columns (total 4 columns):
id            1541272 non-null object
description   1540980 non-null object
quantity      1541272 non-null int64
price         1541272 non-null float64
dtypes: float64(1), int64(1), object(2)
memory usage: 47.0+ MB
```

In [33]:

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

Out[33]:

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
0	p253737	Mrs.	IN	2016-12-05 13:43:57	0
1	p258326	Mr.	FL	2016-10-25 09:22:10	7

In [34]:

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

In [35]:

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

Out[35]:

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
0	p253737	Mrs.	IN	2016-12-05 13:43:57	0

	id	teacher_prefix	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projec
1	p258326	Mr.	FL	2016-10-25 09:22:10	7

2.9 Preprocessing of teacher_prefix

In [36]:

```
import re
prefix = list(project_data['teacher_prefix'].values)

prefix_list = []

for i in prefix:

    j=str(i)
    j=j.lower()
    j = re.sub(r"\.", "", j)

    prefix_list.append(j)

#print(prefix_list)
```

In [37]:

```
project_data['clean_teacher_prefix'] = prefix_list
project_data.drop(['teacher_prefix'], axis=1, inplace=True)
project_data.head(2)
```

Out[37]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

2.10 Preprocessing of school_state

In [38]:

```
state = list(project_data['school_state'].values)

state_list = []

for i in state:

    j=str(i)
    j=j.lower()

    state_list.append(j)

#print(state_list)
```

In [39]:

```
project_data['clean_school_state'] = state_list
#project_data.drop(['school_state'], axis=1, inplace=True)
project_data.head(2)
```

Out[39]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

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](#).

3. Support Vector Machines

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

In [60]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
```

```
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
```

```
X = project_data.drop(['project_is_approved','id'], axis=1)
X.head(2)
```

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

```
# split the data set into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,shuffle=False)
```

```
print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)
```

```
(87398, 12) (87398,)
(21850, 12) (21850,)
```

3.2 Make Data Model Ready: encoding numerical, categorical features

In [61]:

```
# 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
```

3.2.1 encoding categorical features: School State

In [62]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_oh = vectorizer.transform(X_train['clean_school_state'].values)
X_cv_state_oh = vectorizer.transform(X_cv['clean_school_state'].values)
X_test_state_oh = vectorizer.transform(X_test['clean_school_state'].values)

print("After vectorizations")
print(X_train_state_oh.shape, y_train.shape)
#print(X_cv_state_oh.shape, y_cv.shape)
print(X_test_state_oh.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(87398, 51) (87398,)
(21850, 51) (21850,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks',
'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', '
ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
=====
```

3.2.2 encoding categorical features: teacher prefix

In [63]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_teacher_prefix'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_oh = vectorizer.transform(X_train['clean_teacher_prefix'].values)
X_cv_teacher_oh = vectorizer.transform(X_cv['clean_teacher_prefix'].values)
X_test_teacher_oh = vectorizer.transform(X_test['clean_teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_oh.shape, y_train.shape)
#print(X_cv_teacher_oh.shape, y_cv.shape)
print(X_test_teacher_oh.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(87398, 6) (87398,)
(21850, 6) (21850,)
['dr', 'mr', 'mrs', 'ms', 'nan', 'teacher']
=====
```

3.2.3 encoding categorical features: project_grade_category

In [64]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_grade_category'].values) # fit has to happen only on train data
```

```

vectorizer.fit(X_train['clean_grade_category'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer.transform(X_train['clean_grade_category'].values)
#X_cv_grade_ohe = vectorizer.transform(X_cv['clean_grade_category'].values)
X_test_grade_ohe = vectorizer.transform(X_test['clean_grade_category'].values)

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
#print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

After vectorizations
(87398, 4) (87398,)
(21850, 4) (21850,)
['grades3_5', 'grades6_8', 'grades9_12', 'gradesprek_2']
=====

```

3.2.4 encoding categorical features: project_subject_categories

In [65]:

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_categories_ohe = vectorizer.transform(X_train['clean_categories'].values)
#X_cv_categories_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_categories_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_categories_ohe.shape, y_train.shape)
#print(X_cv_categories_ohe.shape, y_cv.shape)
print(X_test_categories_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

After vectorizations
(87398, 9) (87398,)
(21850, 9) (21850,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language', 'math_scienc
e', 'music_arts', 'specialneeds', 'warmth']
=====

```

3.2.5 encoding categorical features: project_subject_subcategories

In [66]:

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_subcategories_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
#X_cv_subcategories_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_subcategories_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_subcategories_ohe.shape, y_train.shape)
#print(X_cv_subcategories_ohe.shape, y_cv.shape)
print(X_test_subcategories_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

After vectorizations
(87398, 30) (87398,)
(21850, 30) (21850,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government', 'college_careerprep', 'co
mmunityservice', 'earlydevelopment', 'economics', 'environmentalscience', 'esl', 'extracurricular', 'fi
nancialliteracy', 'foreignlanguages', 'gym_fitness', 'health_lifescience', 'health_wellness', 'history_

```

```
geography', 'literacy', 'literature_writing', 'mathematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialsciences', 'specialneeds', 'teamsports', 'visualarts', 'war  
mth']
```

3.2.6 encoding numerical feature: price

In [67]:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['price'].values.reshape(-1,1))

X_train_price_scaler = scaler.transform(X_train['price'].values.reshape(-1,1))
#X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
X_test_price_scaler = scaler.transform(X_test['price'].values.reshape(-1,1))

# X_train_price_scaler = X_train_price_scaler.reshape(-1,1)
# #X_cv_price_norm = X_cv_price_norm.reshape(-1,1)
# X_test_price_scaler = X_test_price_scaler.reshape(-1,1)

print("After vectorizations")
print(X_train_price_scaler.shape, y_train.shape)
#print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_scaler.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)
```

In [68]:

```
print(X_train_price_scaler)
```

```
[[-3.87742480e-01]
 [ 5.98715095e-04]
 [ 5.86472187e-01]
 ...
 [-4.08584892e-01]
 [-3.19460050e-01]
 [-7.65891064e-01]]
```

3.2.7 encoding numerical feature: teacher_number_of_previously_posted_projects

In [69]:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_train_posted_project_scaler = scaler.transform(X_train['teacher_number_of_previously_posted_projects']
].values.reshape(-1,1))
#X_cv_posted_project_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].v
alues.reshape(1,-1))
```

```
X_test_posted_project_scaler = scaler.transform(X_test['teacher_number_of_previously_posted_projects'].
values.reshape(-1,1))

# X_train_posted_project_scaler = X_train_posted_project_scaler.reshape(-1,1)
# #X_cv_posted_project_norm = X_cv_posted_project_norm.reshape(-1,1)
# X_test_posted_project_scaler = X_test_posted_project_scaler.reshape(-1,1)

print("After vectorizations")
print(X_train_posted_project_scaler.shape, y_train.shape)
#print(X_cv_posted_project_norm.shape, y_cv.shape)
print(X_test_posted_project_scaler.shape, y_test.shape)
print("=="*100)
```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)

3.3 Make Data Model Ready: encoding eassay, and project_title

In [70]:

```
# 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
```

3.3.1 encoding essay

3.3.1.1 encoding essay : BOW

In [71]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(project_data['clean_essay'].values)

X_train_essay_bow = vectorizer.transform(X_train['clean_essay'].values)
#X_cv_essay_bow = vectorizer.transform(X_cv['clean_essay'].values)
X_test_essay_bow = vectorizer.transform(X_test['clean_essay'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
#print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("=="*100)
```

After vectorizations
(87398, 16623) (87398,)
(21850, 16623) (21850,)

3.3.1.2 encoding essay : TFIDF

In [72]:

```
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(project_data['clean_essay'].values)
```

```
After vectorizations
(87398, 16623) (87398,)
(21850, 16623) (21850,)
```

```
87398
300
[-1.27397677e-02  7.68853460e-02  3.02681217e-02 -1.45521483e-02
-5.59825432e-02 -5.82986609e-02 -3.06017739e+00  1.07062693e-01
 8.74551988e-03  2.24763975e-03 -6.28171696e-02  6.73061976e-02
 8.94857317e-02 -1.11388755e-01 -3.24216280e-02 -1.02995041e-01
 9.46111870e-02 -1.19286409e-01  5.89181522e-02 -2.35565801e-02
 1.25929906e-01  6.36091658e-02 -5.46851255e-02 -1.72957752e-02
 4.75857494e-02 -3.92132745e-02  1.08420205e-01 -1.24281491e-01
-1.10755396e-01 -8.02830404e-02 -2.44689957e-01 -7.33243634e-02
 1.54909975e-02  1.51013340e-01 -5.79608099e-02 -2.26476386e-02
-2.09521919e-02 -1.49075714e-02 -1.08899453e-02 -2.14418609e-02
-3.53230944e-02  1.60913349e-01 -1.64618584e-01 -6.79025839e-02
 1.77766901e-02 -7.86097832e-02  7.24294863e-02  4.72571323e-02
-2.36638609e-02 -1.21562659e-01 -4.97333540e-03 -8.04816335e-03
 6.82426640e-02  2.47284584e-02 -9.49800745e-03 -9.26380627e-02
-8.41076273e-03  3.44850199e-02 -1.48079828e-01 -1.67029234e-02
-4.54254963e-02 -7.35387839e-02  1.96613133e-01 -2.61459937e-02
-1.03438194e-01  1.68149717e-01  7.19528578e-02 -1.82409149e-02
 1.50301033e-01 -3.29286460e-03 -1.62415158e-01  4.93090981e-02
```

4.62819478e-02	-5.99997515e-02	-8.09824267e-02	-8.50705155e-02
-9.21186712e-02	-7.17673851e-03	1.32652209e-01	-4.55777038e-02
4.37017170e-02	-4.42119586e-01	1.24375191e-01	-1.05845279e-01
-5.37251273e-02	-3.13178380e-02	7.83762528e-02	-8.54408524e-02
1.45877899e-01	3.83897224e-02	-3.75069512e-02	-3.03952522e-02
2.75739491e-02	-4.16052764e-02	-4.43663888e-02	-1.00978147e-01
-2.26430846e+00	6.76389938e-02	7.98292789e-02	1.08467383e-01
-1.54405525e-01	-3.60443025e-02	4.99527394e-02	-9.92731478e-02
6.28913317e-02	8.31479503e-02	5.72756957e-02	-1.37570841e-01
1.41705969e-02	-3.77781068e-02	-1.48204745e-02	9.78382174e-03
3.78192120e-02	2.59109519e-01	2.51970925e-02	4.17840149e-02
-2.10480143e-01	8.42777112e-02	9.78638149e-02	5.35765673e-02
-6.66809248e-02	6.01530006e-02	4.98731041e-02	-7.82657720e-02
1.98550634e-02	9.24674894e-02	3.80812075e-02	-1.87524658e-03
1.62647466e-02	1.20124207e-01	5.56539534e-02	2.49189050e-02
-1.27551429e-02	-4.35180311e-03	1.83787006e-02	-1.03308411e-01
1.61365922e-01	-2.55553156e-02	9.53843168e-02	2.50750973e-01
3.91210702e-02	3.08078261e-02	7.70642602e-02	-8.10137168e-02
-7.28019375e-02	-6.05851465e-02	1.48389460e-02	-1.22267579e-01
1.93090911e-01	-3.76394571e-02	-4.78297025e-02	-8.99626435e-02
1.77333559e-02	-3.82686342e-02	3.28577938e-02	1.38311398e-02
-1.35536335e-03	-3.60368460e-02	-3.46428640e-02	-8.68583012e-02
7.52905528e-03	-1.83548609e-02	-4.23556311e-02	-8.82063211e-02
-3.20427945e-02	7.84589923e-02	-1.19621107e-01	1.82180292e-02
5.46003486e-02	-9.01364758e-02	4.80001851e-02	-3.19122522e-02
-3.65354489e-02	-1.85317088e-01	-7.56150801e-03	6.71040590e-02
-3.94159375e-02	4.31732335e-02	-8.40638609e-02	2.37002857e-02
-7.03248553e-03	2.21203011e-01	-3.24874986e-02	-3.47626637e-02
6.58347478e-02	-8.42859557e-02	-6.91604951e-02	-7.07387217e-02
6.49915466e-03	4.38022981e-03	-7.84354472e-03	-1.11407840e-01
-2.30519317e-02	-1.16810683e-03	2.05156727e-02	-7.76687944e-02
3.09946590e-02	-5.35495627e-02	1.21696166e-01	-3.93369644e-02
5.53015925e-02	-6.95430596e-02	-1.97074056e-02	1.07098339e-01
-5.70889155e-02	5.13375994e-02	5.72898267e-02	-3.18427062e-02
1.40936877e-01	6.01723975e-02	1.50338677e-02	2.12080610e-02
2.58220863e-02	-9.06774708e-02	-2.15105894e-02	-2.37361832e-02
-5.55374820e-02	-7.55825826e-02	1.25404497e-02	5.94258590e-02
-1.72682746e-01	-6.53030658e-02	-1.37772133e-01	5.03286814e-02
-2.16628925e+00	1.28998291e-01	-1.37406853e-01	-7.31322261e-02
-7.24280516e-02	-1.47837163e-01	4.68618137e-03	6.26378969e-02
2.05788447e-02	-1.80698882e-02	-7.68966783e-02	-9.15185978e-02
8.26376925e-02	-4.54822814e-02	-1.38261037e-02	9.22719944e-02
-1.41306781e-01	1.49135329e-02	-2.40908088e-01	4.16931534e-02
-4.88187633e-02	-2.07904484e-02	-3.91970273e-02	-1.08958678e-01
8.31100994e-03	5.49559565e-04	2.17592404e-02	1.02302736e-01
3.48695832e-02	1.41747758e-02	9.75991180e-02	1.60644534e-02
4.03670281e-02	-8.51861435e-02	1.14762005e-01	-8.95438360e-02
9.12223006e-02	6.12844596e-03	-5.10804764e-02	5.89501000e-02
3.09595783e-02	-1.0626299e-01	-1.02069508e-01	1.13595123e-02
-4.04622174e-03	9.9207298e-02	-3.20287950e-02	-6.59488944e-03
-1.46744683e-01	1.42886306e-01	-5.09363708e-02	4.60311876e-02
9.23073627e-			

```
avg_w2v_essay_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['clean_essay'].values): # 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_essay_test.append(vector)
```


3.3.2 encoding titles

3.3.2.1 encoding titles : BOW

In [79]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(project_data['clean_project_title'].values)

X_train_title_bow = vectorizer.transform(X_train['clean_project_title'].values)
#X_cv_title_bow = vectorizer.transform(X_cv['clean_project_title'].values)
X_test_title_bow = vectorizer.transform(X_test['clean_project_title'].values)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
#print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations
(87398, 3222) (87398,)
(21850, 3222) (21850,)

3.3.2.2 encoding titles : TFIDF

In [80]:

```
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(project_data['clean_project_title'].values)

X_train_title_tfidf = vectorizer.transform(X_train['clean_project_title'].values)
#X_cv_title_tfidf= vectorizer.transform(X_cv['clean_project_title'].values)
X_test_title_tfidf = vectorizer.transform(X_test['clean_project_title'].values)

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
#print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations
(87398, 3222) (87398,)
(21850, 3222) (21850,)

3.3.2.3 encoding titles : AVG W2V

In [81]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

In [82]:

```
# average Word2Vec
```



```
# compute average wordvec for each review.
avg_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['clean_project_title'].values): # 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_title_train.append(vector)

print(len(avg_w2v_title_train))
print(len(avg_w2v_title_train[0]))
print(avg_w2v_title_train[0])
```

87398

-1.1758000e-02	-1.8464800e-02	-2.0872020e-01	-3.9510000e-03
-5.7714400e-01	-1.8090080e-01	-2.8288200e-01	-2.4662120e-01
-1.8806540e+00	4.4765400e-01	-2.9412700e-01	-1.7280000e-02
-3.1931600e-01	-1.9190500e-01	-1.1642000e-02	1.7475600e-01
1.3068840e-01	1.1943000e-01	-1.7219524e-01	1.9224000e-02
2.2620000e-01	-1.0821980e-01	1.3789060e-01	2.6989320e-01
-2.4364960e-01	-1.3650800e-01	-3.0984180e-01	-3.9546200e-02
-1.1410800e-01	-6.6744640e-02	1.6330620e-01	-4.0601000e-01
9.3793000e-02	-8.3026800e-02	9.0567600e-02	3.1595600e-01
1.6786620e-01	1.0099860e-01	3.5043600e-02	6.6221200e-02
-3.5907800e-02	-2.4589760e-01	2.6006800e-01	-8.0637000e-02
1.5359624e-01	-1.1078680e-01	-5.6956400e-02	2.2253080e-01
3.5808000e-02	-1.8873860e-01	-2.5032660e-01	3.1674000e-02
-2.2424700e-01	2.7863640e-01	2.2622600e-02	1.3753300e-01
-2.3369620e-01	2.8058040e-01	5.0818000e-02	-3.4805800e-02
1.7916600e-01	-7.5374000e-02	7.1228900e-02	1.7556000e-01
-5.8004120e-01	-2.0522500e-01	-1.3367960e-01	1.3656000e-02
-2.9052200e-02	1.3698600e-02	1.1746340e-01	-2.3288400e-02
2.7706200e-01	1.6106000e-01	-2.0183340e-01	5.7781800e-02
-2.0954400e-01	-1.4111260e-02	-3.1186860e-01	-2.9536360e-02
-1.7226500e-01	3.5709400e-01	2.9448200e-01	8.5600000e-05

```
avg_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['clean_project_title'].values): # 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_title_test.append(vector)
```

3.3.2.4 encoding titles : TFIDF W2V

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

In [85]:

```
100%|██████████████████████████████████████████████████████████████████████████████| 87398/87398 [00:03<00:00  
0, 22176.27it/s]
```

In [86]:

```
100%|██████████████████████████████████████████████████████████████████████████| 21850/21850 [00:01<00:00:00, 21309.05it/s]
```

In [87]:

In [88]:

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

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

```

X_tr_bow = hstack((X_train_state_ohc, X_train_teacher_ohc, X_train_grade_ohc, X_train_categories_ohc, X_train_subcategories_ohc, X_train_price_scaler, X_train_posted_project_scaler, X_train_essay_bow, X_train_title_bow)).tocsr()

X_te_bow = hstack((X_test_state_ohc, X_test_teacher_ohc, X_test_grade_ohc, X_test_categories_ohc, X_test_subcategories_ohc, X_test_price_scaler, X_test_posted_project_scaler, X_test_essay_bow, X_test_title_bow)).tocsr()

y_train_bow = y_train
y_test_bow = y_test

print("Final Data matrix")
print(X_tr_bow.shape, y_train_bow.shape)

print(X_te_bow.shape, y_test_bow.shape)
print("=="*100)

```

```

Final Data matrix
(87398, 19947) (87398,)
(21850, 19947) (21850,)
=====

```

3.4.1.1 Hyperparameter Tuning

In [89]:

```

from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier
import matplotlib.pyplot as plt

c= [10**i for i in range(-4,4)]
tuned_parameters = [{'alpha':c}]

clf_bow = SGDClassifier(loss='hinge',penalty='l2',class_weight='balanced', n_jobs=-1)

#Using GridSearchCV
model_bow = GridSearchCV(clf_bow, tuned_parameters, scoring = 'roc_auc',verbose=5,n_jobs=-1,return_train_score=True)
model_bow.fit(X_tr_bow, y_train_bow)

print(model_bow.best_estimator_)
print(model_bow.score(X_te_bow, y_test_bow))

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks      | elapsed: 14.5s
[Parallel(n_jobs=-1)]: Done 22 out of 24 | elapsed: 18.1s remaining: 1.6s
[Parallel(n_jobs=-1)]: Done 24 out of 24 | elapsed: 18.6s finished

```

```

SGDClassifier(alpha=0.01, 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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.7381240694829668

```

In [90]:

```

train_auc= model_bow.cv_results_['mean_train_score']
cv_auc = model_bow.cv_results_['mean_test_score']

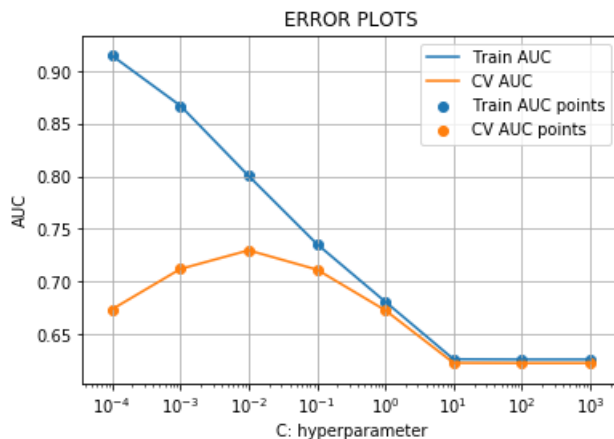
# https://stackoverflow.com/a/48803361/4084039

```

```
plt.plot(c, train_auc, label='Train AUC')
plt.plot(c, cv_auc, label='CV AUC')
plt.xscale('log')

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

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



In [91]:

```
print(model_bow.best_estimator_)
print(model_bow.score(X_te_bow, y_test_bow))
```

```
SGDClassifier(alpha=0.01, 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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.7381240694829668
```

3.4.1.2 Testing the performance of the model on test data, plotting ROC Curves

In [92]:

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

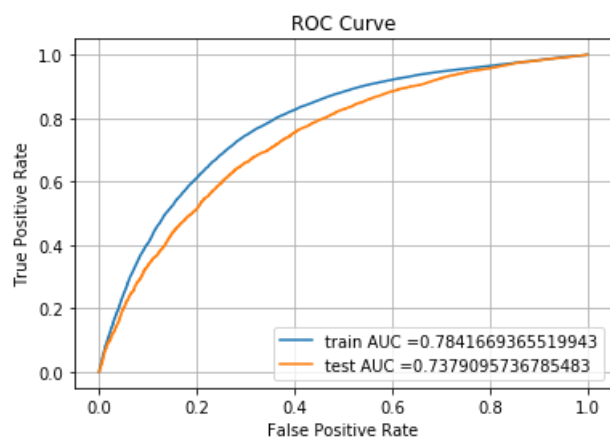
clf_bow = SGDClassifier(alpha=0.01, loss='hinge', penalty='l2', class_weight='balanced', n_jobs=-1)
clf_bow.fit(X_tr_bow, y_train_bow)

y_train_pred = clf_bow.decision_function(X_tr_bow)
y_test_pred = clf_bow.decision_function(X_te_bow)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train_bow, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test_bow, y_test_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
```

```
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



In [93]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [94]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

=====

the maximum value of tpr*(1-fpr) 0.5224341885261056 for threshold -0.154

In [95]:

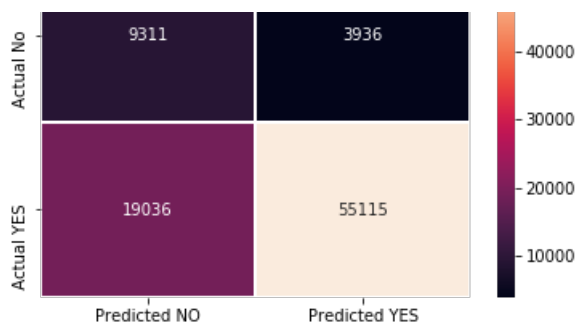
```
def get_confusion_matrix(y, y_pred):
    df = pd.DataFrame(confusion_matrix(y, y_pred), range(2), range(2))
    df.columns = ['Predicted NO', 'Predicted YES']
    df = df.rename({0: 'Actual No', 1: 'Actual YES'})
    sns.heatmap(df, annot=True, fmt='g', linewidth=0.5)
```

In [96]:

```
print("Train confusion matrix")
get_confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
```

Train confusion matrix

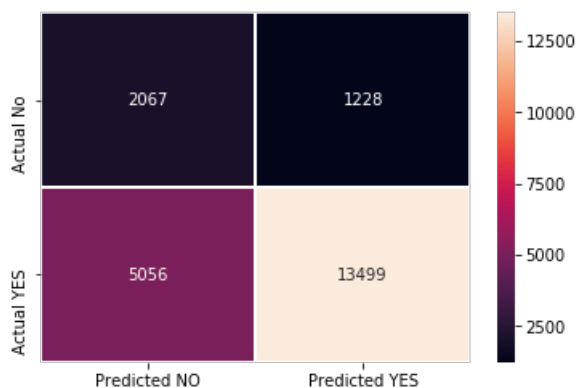




In [97]:

```
print("Test confusion matrix")
get_confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
```

Test confusion matrix



3.4.2 Applying Support Vector Machines on TFIDF, SET 2

In [98]:

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

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

X_tr_tfidf = hstack((X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_categories_ohe,
X_train_subcategories_ohe, X_train_price_scaler, X_train_posted_project_scaler, X_train_essay_tfidf, X_train_title_tfidf)).tocsr()
X_te_tfidf = hstack((X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_categories_ohe, X_test_subcategories_ohe,
X_test_price_scaler, X_test_posted_project_scaler, X_test_essay_tfidf, X_test_title_tfidf)).tocsr()

y_train_tfidf = y_train
y_test_tfidf = y_test

print("Final Data matrix")
print(X_tr_tfidf.shape, y_train_tfidf.shape)
print(X_te_tfidf.shape, y_test_tfidf.shape)
print("="*100)
```

Final Data matrix
(87398, 19947) (87398,)
(21850, 19947) (21850,)

=====

In [99]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

c= [10**i for i in range(-4,4)]
```

```

tuned_parameters = [{'alpha':c}]

clf_tfidf = SGDClassifier(loss='hinge',penalty='l2',class_weight='balanced', n_jobs=-1)

#Using GridSearchCV
model_tfidf = GridSearchCV(clf_tfidf, tuned_parameters, scoring = 'roc_auc',verbose=5,n_jobs=-1,return_
train_score=True)
model_tfidf.fit(X_tr_tfidf, y_train_tfidf)

print(model_tfidf.best_estimator_)

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks      | elapsed:    7.8s
[Parallel(n_jobs=-1)]: Done 22 out of  24 | elapsed:  10.8s remaining:    0.9s
[Parallel(n_jobs=-1)]: Done 24 out of  24 | elapsed:  11.1s finished

```

```

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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)

```

In [100]:

```

import matplotlib.pyplot as plt

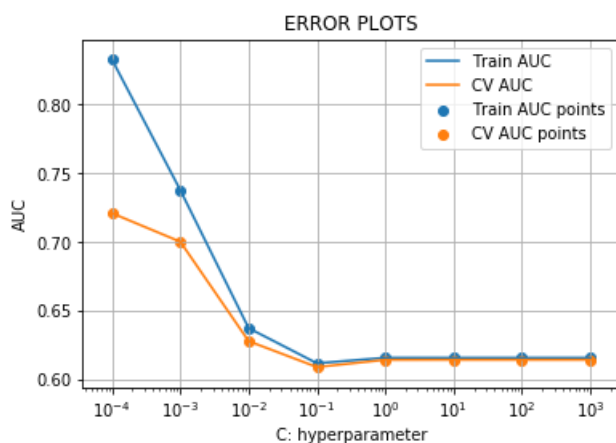
train_auc= model_tfidf.cv_results_['mean_train_score']
cv_auc = model_tfidf.cv_results_['mean_test_score']

# https://stackoverflow.com/a/48803361/4084039
plt.plot(c, train_auc, label='Train AUC')
plt.plot(c, cv_auc, label='CV AUC')
plt.xscale('log')

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

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

```



In [101]:

```

print(model_tfidf.best_estimator_)
print(model_tfidf.score(X_test_tfidf, y_test_tfidf))

```



```
print(model_clf.score(X_te_clf, y_test_clf))
```

```
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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.7336865219220715
```

3.4.2.2 Testing the performance of the model on test data, plotting ROC Curves

In [102]:

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

clf_tfidf = SGDClassifier(alpha=0.0001, loss='hinge', penalty='l2', class_weight='balanced', n_jobs=-1)

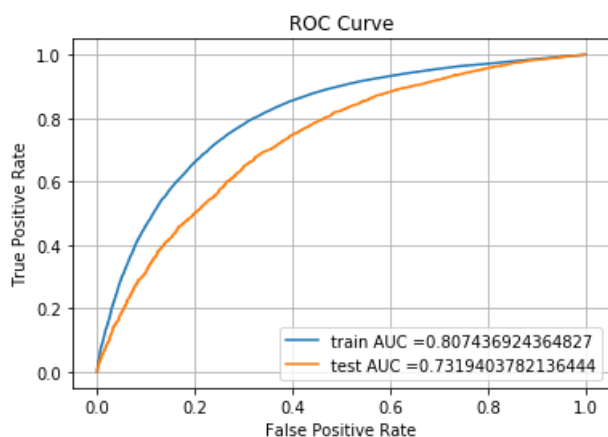
clf_tfidf.fit(X_tr_tfidf, y_train_tfidf)

#print(clf.predict_proba(X_te_bow)[:,:1])

y_train_pred = clf_tfidf.decision_function(X_tr_tfidf)
y_test_pred = clf_tfidf.decision_function(X_te_tfidf)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train_tfidf, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test_tfidf, y_test_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



In [103]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
```

```

else:
    predictions.append(0)
return predictions

```

In [104]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)

```

=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.5489126204322521 for threshold -0.101

In [105]:

```

def get_confusion_matrix(y, y_pred):

    df = pd.DataFrame(confusion_matrix(y, y_pred), range(2), range(2))
    df.columns = ['Predicted NO', 'Predicted YES']
    df = df.rename({0: 'Actual No', 1: 'Actual YES'})
    sns.heatmap(df, annot=True, fmt='g', linewidth=0.5)

```

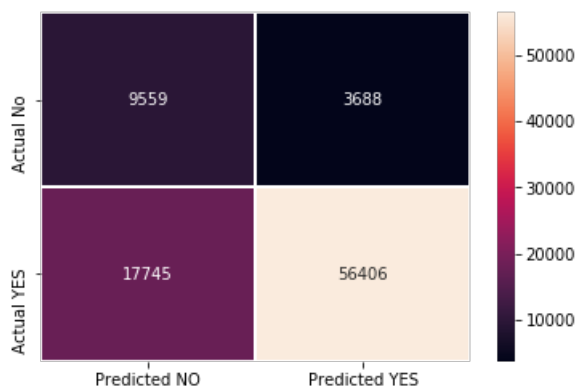
In [106]:

```

print("Train confusion matrix")
get_confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))

```

Train confusion matrix



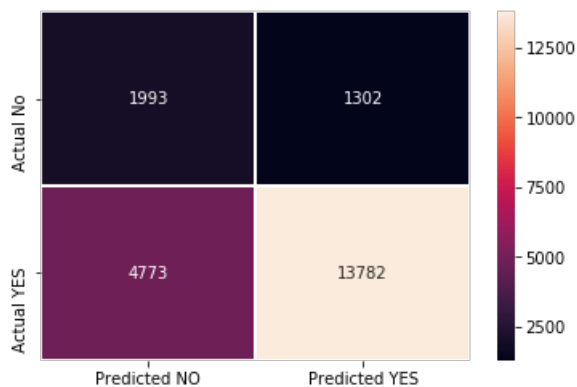
In [107]:

```

print("Test confusion matrix")
get_confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))

```

Test confusion matrix



3.4.3 Applying Support Vector Machines on AVG W2V, SET 3

In [108]:

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

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

X_tr_avgw2v = hstack((X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_categories_ohe,
                     X_train_subcategories_ohe, X_train_price_scaler, X_train_posted_project_scaler, avg_w2v_essay_train,
                     avg_w2v_title_train)).tocsr()
X_te_avgw2v = hstack((X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_categories_ohe, X_test_subcategories_ohe,
                     X_test_price_scaler, X_test_posted_project_scaler, avg_w2v_essay_test, avg_w2v_title_test)).tocsr()

y_train_avgw2v = y_train
y_test_avgw2v = y_test

print("Final Data matrix")
print(X_tr_avgw2v.shape, y_train_avgw2v.shape)
print(X_te_avgw2v.shape, y_test_avgw2v.shape)
print("=="*100)
```

```
Final Data matrix
(87398, 702) (87398,)
(21850, 702) (21850,)
```

3.4.3.1 Hyperparameter Tuning

In [109]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import SGDClassifier

c= [10**i for i in range(-4,4)]
tuned_parameters = [{'alpha':c}]

clf_avgw2v = SGDClassifier(loss='hinge',penalty='l2',class_weight='balanced', n_jobs=-1)

#Using GridSearchCV
model_avgw2v = GridSearchCV(clf_avgw2v, tuned_parameters, scoring = 'roc_auc',verbose=5,n_jobs=-1,return_train_score=True)
model_avgw2v.fit(X_tr_avgw2v, y_train_avgw2v)

print(model_avgw2v.best_estimator_)
```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks | elapsed: 24.3s
[Parallel(n_jobs=-1)]: Done 22 out of 24 | elapsed: 33.2s remaining: 2.9s
[Parallel(n_jobs=-1)]: Done 24 out of 24 | elapsed: 34.8s finished
```

```
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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

In [110]:

```
import matplotlib.pyplot as plt
```

```
train_auc= model_avgw2v.cv_results_['mean_train_score']
```

```

train_auc = model_avgw2v.cv_results_['mean_train_score']
cv_auc = model_avgw2v.cv_results_['mean_test_score']

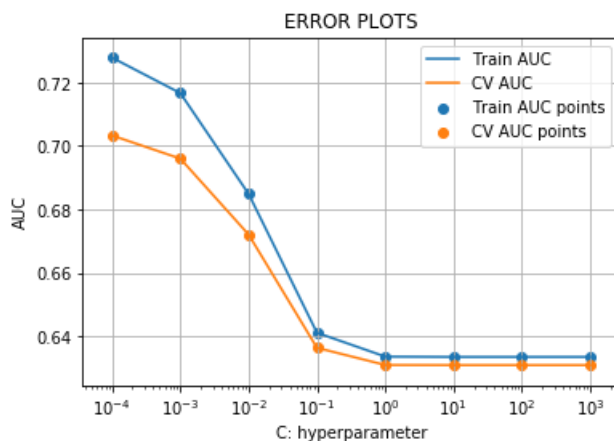
# https://stackoverflow.com/a/48803361/4084039
plt.plot(c, train_auc, label='Train AUC')
plt.plot(c, cv_auc, label='CV AUC')

plt.xscale('log')

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

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

```



In [111]:

```

print(model_avgw2v.best_estimator_)
print(model_avgw2v.score(X_te_avgw2v, y_test_avgw2v))

```

```

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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.7142788551773036

```

3.4.3.2 Testing the performance of the model on test data, plotting ROC Curves

In [112]:

```

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

clf_avgw2v = SGDClassifier(alpha=0.0001, loss='hinge', penalty='l2', class_weight='balanced', n_jobs=-1)
clf_avgw2v.fit(X_tr_avgw2v, y_train_avgw2v)

#print(clf.predict_proba(X_te_avgw2v)[:1])

y_train_pred = clf_avgw2v.decision_function(X_tr_avgw2v)
y_test_pred = clf_avgw2v.decision_function(X_te_avgw2v)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train_avgw2v, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test_avgw2v, y_test_pred)

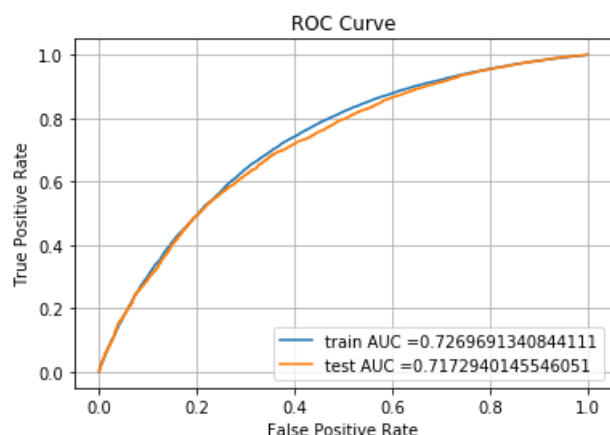
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()

```

```

plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.grid()
plt.show()

```



In [113]:

```

# we are writing our own function for predict, with defined threshould
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

```

In [114]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)

```

=====

the maximum value of tpr*(1-fpr) 0.4516837981201981 for threshold -0.391

In [115]:

```

def get_confusion_matrix(y, y_pred):

    df = pd.DataFrame(confusion_matrix(y, y_pred), range(2), range(2))
    df.columns = ['Predicted NO', 'Predicted YES']
    df = df.rename({0: 'Actual No', 1: 'Actual YES'})
    sns.heatmap(df, annot=True, fmt='g', linewidth=0.5)

```

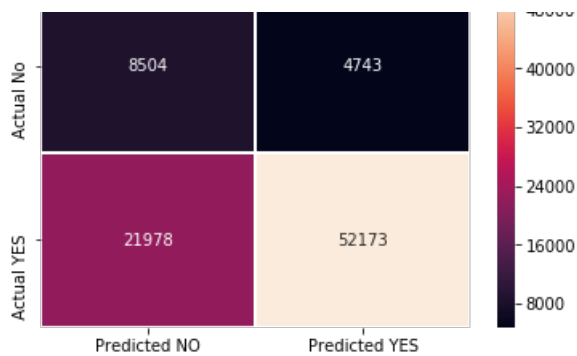
In [116]:

```

print("Train confusion matrix")
get_confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))

```

Train confusion matrix



In [117]:

```
print("Test confusion matrix")
get_confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
```

Test confusion matrix



3.4.4 Applying Support Vector Machines on TFIDF W2V, SET 4

In [118]:

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

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

X_tr_tfidfw2v = hstack((X_train_state_oh, X_train_teacher_oh, X_train_grade_oh, X_train_categories_oh,
                        X_train_subcategories_oh, X_train_price_scaler, X_train_posted_project_scaler, essay_tfidf_w2v_train,
                        title_tfidf_w2v_train)).tocsr()
X_te_tfidfw2v = hstack((X_test_state_oh, X_test_teacher_oh, X_test_grade_oh, X_test_categories_oh,
                        X_test_subcategories_oh, X_test_price_scaler, X_test_posted_project_scaler, essay_tfidf_w2v_test,
                        title_tfidf_w2v_test)).tocsr()

y_train_tfidfw2v = y_train
y_test_tfidfw2v = y_test

print("Final Data matrix")
print(X_tr_tfidfw2v.shape, y_train_tfidfw2v.shape)
print(X_te_tfidfw2v.shape, y_test_tfidfw2v.shape)
print("=="*100)
```

Final Data matrix
(87398, 702) (87398,)
(21850, 702) (21850,)

3.4.4.1 Hyperparameter Tuning

In [119]:

```

from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

c= [10**i for i in range(-4,4)]
tuned_parameters = [{'alpha':c}]

clf_tfidf2v = SGDClassifier(loss='hinge',penalty='l2',class_weight='balanced', n_jobs=-1)

#Using GridSearchCV
model_tfidf2v = GridSearchCV(clf_tfidf2v, tuned_parameters, scoring = 'roc_auc',verbose=5,n_jobs=-1,return_train_score=True)
model_tfidf2v.fit(X_train_tfidf2v, y_train_tfidf2v)

print(model_tfidf2v.best_estimator_)

```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks      | elapsed: 23.6s
[Parallel(n_jobs=-1)]: Done 22 out of 24 | elapsed: 33.0s remaining: 2.9s
[Parallel(n_jobs=-1)]: Done 24 out of 24 | elapsed: 33.4s finished

```

```

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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)

```

In [120]:

```

import matplotlib.pyplot as plt

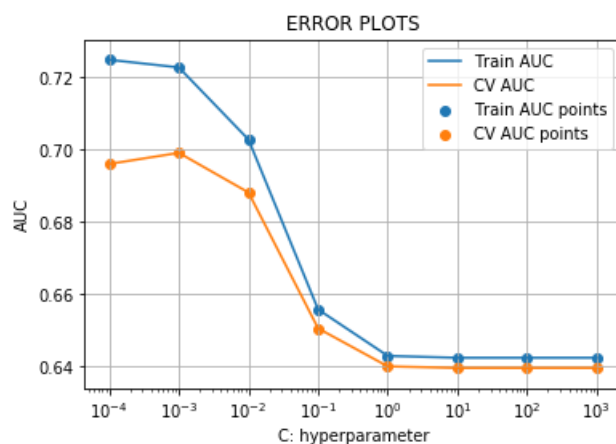
train_auc= model_tfidf2v.cv_results_['mean_train_score']
cv_auc = model_tfidf2v.cv_results_['mean_test_score']

# https://stackoverflow.com/a/48803361/4084039
plt.plot(c, train_auc, label='Train AUC')
plt.plot(c, cv_auc, label='CV AUC')

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

plt.legend()
plt.xscale('log')
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [121]:

```
print(model_tfidfw2v.best_estimator_)
print(model_tfidfw2v.score(X_te_tfidfw2v, y_test_tfidfw2v))
```

```
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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
0.7122319855378076
```

3.4.4.2 Testing the performance of the model on test data, plotting ROC Curves

In [122]:

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

clf = SGDClassifier(alpha=0.001, loss='hinge', penalty='l2', class_weight='balanced', n_jobs=-1)

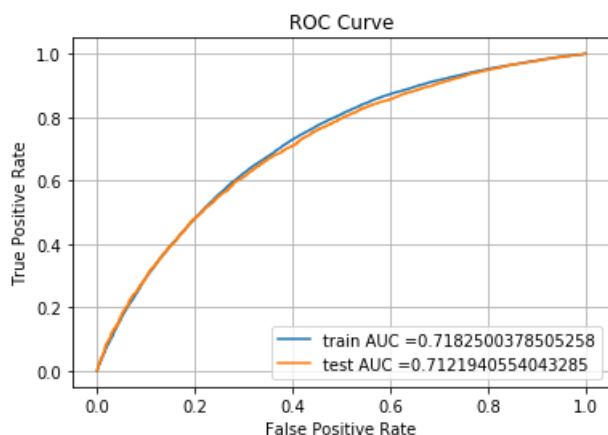
clf.fit(X_tr_tfidfw2v, y_train_tfidfw2v)

#print(clf.predict_proba(X_te_bow)[: ,1])

y_train_pred = clf.decision_function(X_tr_tfidfw2v)
y_test_pred = clf.decision_function(X_te_tfidfw2v)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train_tfidfw2v, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test_tfidfw2v, y_test_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



In [123]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
```



```

    if i>=threshold:
        predictions.append(1)
    else:
        predictions.append(0)
return predictions

```

In [124]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)

```

=====

the maximum value of tpr*(1-fpr) 0.44026265908631795 for threshold -0.212

In [125]:

```

def get_confusion_matrix(y, y_pred):

    df = pd.DataFrame(confusion_matrix(y, y_pred), range(2), range(2))
    df.columns = ['Predicted NO', 'Predicted YES']
    df = df.rename({0: 'Actual No', 1: 'Actual YES'})
    sns.heatmap(df, annot=True, fmt='g', linewidth=0.5)

```

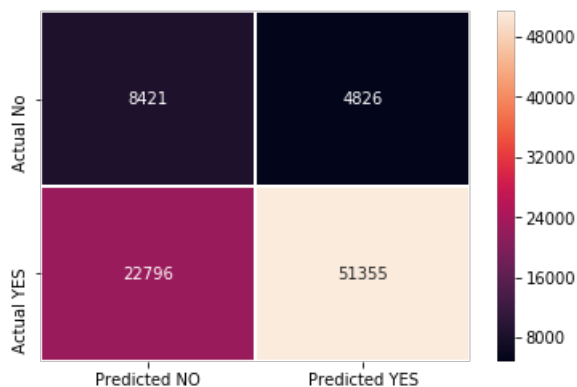
In [126]:

```

print("Train confusion matrix")
get_confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))

```

Train confusion matrix



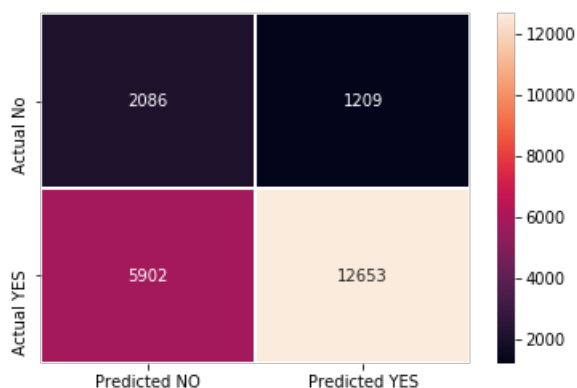
In [127]:

```

print("Test confusion matrix")
get_confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))

```

Test confusion matrix



3.5 Support Vector Machines with added Features `Set 5`

In [128]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

3.5.1.1 Adding quantity column in our dataframe

In [129]:

```
resource_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1541272 entries, 0 to 1541271
Data columns (total 4 columns):
id                1541272 non-null object
description       1540980 non-null object
quantity         1541272 non-null int64
price            1541272 non-null float64
dtypes: float64(1), int64(1), object(2)
memory usage: 47.0+ MB
```

In [130]:

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

Out[130]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

In [131]:

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

In [132]:

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

In [133]:

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

Out[133]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

3.5.1.2 Adding no_of_words_title in our dataframe

In [134]:

```
words = []
for title in project_data['clean_project_title']:
    words.append(len(title.split()))

project_data['no_of_words_title'] = words
project_data.head(2)
```

Out[134]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

3.5.1.3 Adding no_of_words_essay in our dataframe

In [135]:

```
words = []
for essay in project_data['clean_essay']:
    words.append(len(essay.split()))

project_data['no_of_words_essay'] = words
project_data.head(2)
```

Out[135]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
1	p258326	FL	2016-10-25 09:22:10	7	1

3.5.1.4 Adding sentiment_score_essay in our dataframe

In [136]:

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

sid = SentimentIntensityAnalyzer()

score = []
for essay in project_data['clean_essay']:
    compound = sid.polarity_scores(essay) ["compound"]

    if compound >= 0.5:
        score.append('positive')
    elif compound <= -0.5:
        score.append('negative')
    else:
        score.append('neutral')

project_data['sentiment_score_essay'] = score

project_data.head(2)
```

Out[136]:

	id	school_state	project_submitted_datetime	teacher_number_of_previously_posted_projects	project_is_a
0	p253737	IN	2016-12-05 13:43:57	0	0
1	p258326	FL	2016-10-25 09:22:10	7	1

In [137]:

```
print(project_data['clean_essay'][5])
```

i moving 2nd grade 3rd grade beginning next school year i takings current students move i teach inclusi on classroom includes students adhd sld well autistic students my students work hard achieving goals no matter struggles may the school i teach houses great deal autistic students well ell students my studen t love read work challenge they also love move around they work better able move room different areas r ather usual set these flexible seating options allow students different seating options instead sitting traditional desk chair able use flexible seating tools reduce stress anxiety these tools beneficial stu dents special needs also students it proven fact students moving oxygen going brain means learning taki ng place these flexible seating options allow students move traditional seat allows reduce stress class room this project significantly help students reduce stress anxiety standardized testing the students 3 rd grade required take state mandated test this puts great deal stress students perform well test if st udents able work throughout year less stressful classroom assistance flexible seating obtain skills nee ded successful standardized test nannan

In [138]:

```
print(score[5])
```

positive

3.5.2 Splitting data into train and test

In [139]:

```
from sklearn.model_selection import train_test_split

X = project_data.drop(['project_is_approved', 'id', 'clean_project_title', 'project_submitted_datetime'],
axis=1)
X.head(2)

y = project_data['project_is_approved'].values

# split the data set into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=False)

print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)

(87398, 13) (87398,)
(21850, 13) (21850,)
```

In [140]:

```
X_train.head(2)
```

Out[140]:

	school_state	teacher_number_of_previously_posted_projects	clean_categories	clean_subcategories	clean_ess:
0	IN	0	Literacy_Language	ESL Literacy	my students english learners working english s...
1	FL	7	History_Civics Health_Sports	Civics_Government TeamSports	our student arrive schoo eager learn they po...

3.5.3 Make data model ready

3.5.3.1 encoding numerical feature: quantity

In [141]:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['quantity'].values.reshape(-1,1))

X_train_quantity_scaler = scaler.transform(X_train['quantity'].values.reshape(1,-1))
```

```

X_test_quantity_scaler = scaler.transform(X_test['quantity'].values.reshape(1,-1))

X_train_quantity_scaler = X_train_quantity_scaler.reshape(-1,1)

X_test_quantity_scaler = X_test_quantity_scaler.reshape(-1,1)

print("After vectorizations")
print(X_train_quantity_scaler.shape, y_train.shape)

print(X_test_quantity_scaler.shape, y_test.shape)
print("=="*100)

```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)

3.5.3.2 encoding numerical feature: price

In [142]:

```

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['price'].values.reshape(-1,1))

X_train_price_scaler = scaler.transform(X_train['price'].values.reshape(1,-1))

X_test_price_scaler = scaler.transform(X_test['price'].values.reshape(1,-1))

X_train_price_scaler = X_train_price_scaler.reshape(-1,1)

X_test_price_scaler = X_test_price_scaler.reshape(-1,1)

print("After vectorizations")
print(X_train_price_scaler.shape, y_train.shape)

print(X_test_price_scaler.shape, y_test.shape)
print("=="*100)

```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)

3.5.3.3 encoding numerical feature: teacher_number_of_previously_posted_projects

In [143]:

```

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_train_posted_project_scaler = scaler.transform(X_train['teacher_number_of_previously_posted_projects']
].values.reshape(1,-1))

```

```

X_test_posted_project_scaler = scaler.transform(X_test['teacher_number_of_previously_posted_projects'].
values.reshape(1,-1))

X_train_posted_project_scaler = X_train_posted_project_scaler.reshape(-1,1)

X_test_posted_project_scaler = X_test_posted_project_scaler.reshape(-1,1)


print("After vectorizations")
print(X_train_posted_project_scaler.shape, y_train.shape)

print(X_test_posted_project_scaler.shape, y_test.shape)
print("=="*100)

```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)

3.5.3.4 encoding numerical feature: no_of_words_title

In [144]:

```

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['no_of_words_title'].values.reshape(-1,1))

X_train_no_of_words_title_scaler = scaler.transform(X_train['no_of_words_title'].values.reshape(1,-1))

X_test_no_of_words_title_scaler = scaler.transform(X_test['no_of_words_title'].values.reshape(1,-1))


X_train_no_of_words_title_scaler = X_train_no_of_words_title_scaler.reshape(-1,1)

X_test_no_of_words_title_scaler = X_test_no_of_words_title_scaler.reshape(-1,1)


print("After vectorizations")
print(X_train_no_of_words_title_scaler.shape, y_train.shape)

print(X_test_no_of_words_title_scaler.shape, y_test.shape)
print("=="*100)

```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)

3.5.3.5 encoding numerical feature: no_of_words_essay

In [145]:

```

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
scaler.fit(X_train['no_of_words_essay'].values.reshape(-1,1))

X_train_no_of_words_essay_scaler = scaler.transform(X_train['no_of_words_essay'].values.reshape(1,-1))

```

```

X_test_no_or_words_essay_scaler = scaler.transform(X_test['no_or_words_essay'].values.reshape(1,-1))

X_train_no_of_words_essay_scaler = X_train_no_of_words_essay_scaler.reshape(-1,1)

X_test_no_of_words_essay_scaler = X_test_no_of_words_essay_scaler.reshape(-1,1)

print("After vectorizations")
print(X_train_no_of_words_essay_scaler.shape, y_train.shape)

print(X_test_no_of_words_essay_scaler.shape, y_test.shape)
print("="*100)

```

```

After vectorizations
(87398, 1) (87398,)
(21850, 1) (21850,)
=====

```

3.5.3.6 encoding categorical features: School State

In [146]:

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_ohe = vectorizer.transform(X_train['clean_school_state'].values)
#X_cv_state_ohe = vectorizer.transform(X_cv['clean_school_state'].values)
X_test_state_ohe = vectorizer.transform(X_test['clean_school_state'].values)

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
#print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

After vectorizations
(87398, 51) (87398,)
(21850, 51) (21850,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks',
'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', '
ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
=====

```

3.5.3.7 encoding categorical features: clean_teacher_prefix

In [147]:

```

vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_teacher_prefix'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer.transform(X_train['clean_teacher_prefix'].values)
#X_cv_teacher_ohe = vectorizer.transform(X_cv['clean_teacher_prefix'].values)
X_test_teacher_ohe = vectorizer.transform(X_test['clean_teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
#print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

```

```

After vectorizations
(87398, 6) (87398,)
(21850, 6) (21850,)
['dr', 'mr', 'mrs', 'ms', 'nan', 'teacher']
=====

```


3.5.3.8 encoding categorical features: clean_grade_category

In [148]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_grade_category'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohc = vectorizer.transform(X_train['clean_grade_category'].values)
#X_cv_grade_ohc = vectorizer.transform(X_cv['clean_grade_category'].values)
X_test_grade_ohc = vectorizer.transform(X_test['clean_grade_category'].values)

print("After vectorizations")
print(X_train_grade_ohc.shape, y_train.shape)
#print(X_cv_grade_ohc.shape, y_cv.shape)
print(X_test_grade_ohc.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations
(87398, 4) (87398,)
(21850, 4) (21850,)
['grades3_5', 'grades6_8', 'grades9_12', 'gradesprek_2']
=====

3.5.3.9 encoding categorical features: clean_categories

In [149]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_categories_ohc = vectorizer.transform(X_train['clean_categories'].values)
#X_cv_categories_ohc = vectorizer.transform(X_cv['clean_categories'].values)
X_test_categories_ohc = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_categories_ohc.shape, y_train.shape)
#print(X_cv_categories_ohc.shape, y_cv.shape)
print(X_test_categories_ohc.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations
(87398, 9) (87398,)
(21850, 9) (21850,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language', 'math_scienc
e', 'music_arts', 'specialneeds', 'warmth']
=====

3.5.3.10 encoding categorical features: clean_subcategories

In [150]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_subcategories_ohc = vectorizer.transform(X_train['clean_subcategories'].values)
#X_cv_subcategories_ohc = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_subcategories_ohc = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_subcategories_ohc.shape, y_train.shape)
#print(X_cv_subcategories_ohc.shape, y_cv.shape)
print(X_test_subcategories_ohc.shape, y_test.shape)
print(vectorizer.get_feature_names())
```

```
print("="*100)
```

After vectorizations

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

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

```
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government', 'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience', 'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness', 'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'mathematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialsciences', 'specialneeds', 'teamsports', 'visualarts', 'war_mth']
```

3.5.3.11 encoding categorical features: sentiment_score_essay

In [151]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['sentiment_score_essay'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_sentiment_ohe = vectorizer.transform(X_train['sentiment_score_essay'].values)
X_test_sentiment_ohe = vectorizer.transform(X_test['sentiment_score_essay'].values)

print("After vectorizations")
print(X_train_sentiment_ohe.shape, y_train.shape)
print(X_test_sentiment_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
```

After vectorizations

```
(87398, 3) (87398,)
```

```
(21850, 3) (21850,)
```

```
['negative', 'neutral', 'positive']
```

3.5.3.12 encoding essay: TFIDF

In [152]:

```
X_train.head(2)
```

Out[152]:

	school_state	teacher_number_of_previously_posted_projects	clean_categories	clean_subcategories	clean_essay
0	IN	0	Literacy_Language	ESL Literacy	my students english learners working english s...
1	FL	7	History_Civics Health_Sports	Civics_Government TeamSports	our student arrive school eager learn they po...

In [153]:

```
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(project_data['clean_essay'].values)

X_train_essay_tfidf = vectorizer.transform(X_train['clean_essay'].values)
#X_cv_essay_tfidf = vectorizer.transform(X_cv['clean_essay'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['clean_essay'].values)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
```

```
After vectorizations
(87398, 16623) (87398,)
(21850, 16623) (21850,)
```

In [154]:

(87398, 5000) (87398,)
(21850, 5000) (21850,)

```
from sklearn.decomposition import TruncatedSVD
component = [10,50,100,500,1000,1500,2000,2500]
variance = []

for i in tqdm(component):
    svd = TruncatedSVD(n_components = i)
    svd.fit(X_train_essay_tf)
    variance.append(svd.explained_variance_ratio_.sum())
```

In [156]:

A line graph showing the relationship between the number of components (X-axis) and the variance explained (Y-axis). The X-axis is labeled 'Component' and ranges from 0 to 2500 with major ticks at 0, 500, 1000, 1500, 2000, and 2500. The Y-axis is labeled 'Variance' and ranges from 0.0 to 1.0 with major ticks at 0.2, 0.4, 0.6, and 0.8. The curve starts at (0,0) and rises steeply, reaching a variance of approximately 0.65 at component 500, 0.8 at component 1000, and continues to rise more gradually, reaching approximately 0.95 at component 2500.

Component	Variance
0	0.00
100	0.20
200	0.30
500	0.65
1000	0.80
1500	0.88
2000	0.92
2500	0.95

```
svd = TruncatedSVD(n_components = 2000)
svd.fit(X_train_essay_tf)
X_train_essay_tfidf_svd = svd.transform(X_train_essay_tf)
X_test_essay_tfidf_svd = svd.transform(X_test_essay_tf)
```

In [159]:

```
X_train_essay_tf.shape
```

Out[159]:

```
(87398, 5000)
```

In [160]:

```
X_train_essay_tfidf_svd.shape
```

Out[160]:

```
(87398, 2000)
```

3.5.4 Applying Support Vector Machines on set 5

In [161]:

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

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

X_tr = hstack((X_train_state_ohc, X_train_teacher_ohc, X_train_grade_ohc, X_train_categories_ohc, X_train_subcategories_ohc, X_train_price_scaler, X_train_quantity_scaler, X_train_posted_project_scaler, X_train_no_of_words_title_scaler, X_train_no_of_words_essay_scaler, X_train_sentiment_ohc, X_train_essay_tfidf_svd).to_csr())
X_te = hstack((X_test_state_ohc, X_test_teacher_ohc, X_test_grade_ohc, X_test_categories_ohc, X_test_subcategories_ohc, X_test_price_scaler, X_test_quantity_scaler, X_test_posted_project_scaler, X_test_no_of_words_title_scaler, X_test_no_of_words_essay_scaler, X_test_sentiment_ohc, X_test_essay_tfidf_svd).to_csr())

y_train = y_train
y_test = y_test

print("Final Data matrix")
print(X_tr.shape, y_train.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(87398, 2108) (87398,)
(21850, 2108) (21850,)
```

=====

3.5.4.1 Hyperparameter Tuning

In [162]:

```
from sklearn.model_selection import GridSearchCV

c= [10**i for i in range(-4,4)]
tuned_parameters = [{'alpha':c}]

clf_set5 = SGDClassifier(loss='hinge',penalty='l2',class_weight='balanced', n_jobs=-1)

#Using GridSearchCV
model_set5 = GridSearchCV(clf_set5, tuned_parameters, scoring = 'roc_auc',verbose=5,n_jobs=-1,return_train_score=True)
model_set5.fit(X_tr, y_train)

# y_train_pred = batch_predict(neigh, X_tr_bow)
# y_cv_pred = batch_predict(neigh, X_cv_bow)
```

```
print(model_set5.best_estimator_)
```

Fitting 3 folds for each of 8 candidates, totalling 24 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks      | elapsed: 3.3min
[Parallel(n_jobs=-1)]: Done 22 out of 24 | elapsed: 3.9min remaining: 21.0s
[Parallel(n_jobs=-1)]: Done 24 out of 24 | elapsed: 3.9min finished
```

```
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=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
```

In [163]:

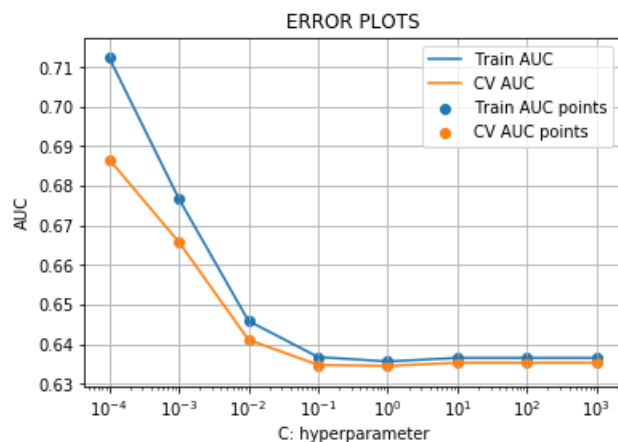
```
import matplotlib.pyplot as plt

train_auc= model_set5.cv_results_['mean_train_score']
cv_auc = model_set5.cv_results_['mean_test_score']

# https://stackoverflow.com/a/48803361/4084039
plt.plot(c, train_auc, label='Train AUC')
plt.plot(c, cv_auc, label='CV AUC')

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

plt.legend()
plt.xscale('log')
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [164]:

```
print(model_set5.best_estimator_)
print(model_set5.score(X_te, y_test))
```

```
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=-1, penalty='l2',
power_t=0.5, random_state=None, shuffle=True, tol=0.001,
validation_fraction=0.1, verbose=0, warm_start=False)
0.6912440601271943
```

3.5.4.2 Testing the performance of the model on test data, plotting ROC Curves

In [165]:

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

clf_set5 = SGDClassifier(alpha=0.0001, loss='hinge', penalty='l2', class_weight='balanced', n_jobs=-1)

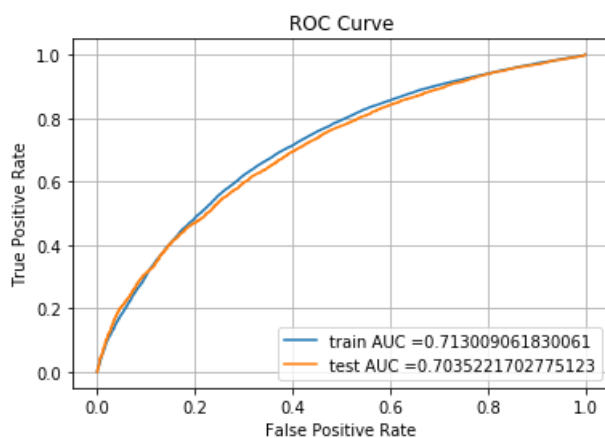
clf_set5.fit(X_tr, y_train)

#print(clf.predict_proba(X_teBow)[:,:])

y_train_pred = clf_set5.decision_function(X_tr)
y_test_pred = clf_set5.decision_function(X_te)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



In [166]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i >= threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [167]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.43671995737884056 for threshold -0.184

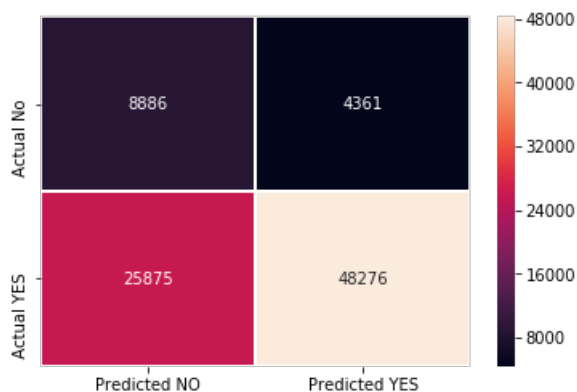
In [168]:

```
def get_confusion_matrix(y, y_pred):
    df = pd.DataFrame(confusion_matrix(y, y_pred), range(2), range(2))
    df.columns = ['Predicted NO', 'Predicted YES']
    df = df.rename({0: 'Actual No', 1: 'Actual YES'})
    sns.heatmap(df, annot=True, fmt='g', linewidth=0.5)
```

In [169]:

```
print("Train confusion matrix")
get_confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
```

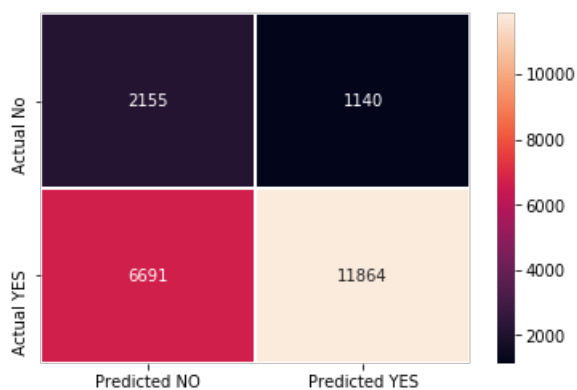
Train confusion matrix



In [170]:

```
print("Test confusion matrix")
get_confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
```

Test confusion matrix



4. Conclusion

In [171]:

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

```

from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyperparameter (alpha)", "AUC"]

x.add_row(["BOW", "SVM", 0.01, 0.73790])
x.add_row(["TFIDF", "SVM", 0.0001, 0.73194])
x.add_row(["ACG W2V", "SVM", 0.0001, 0.71729])
x.add_row(["TFIDF W2V", "SVM", 0.001, 0.71219])
x.add_row(["Set 5", "SVM", 0.0001, 0.70352])

print(x)

```

Vectorizer	Model	Hyperparameter (alpha)	AUC
BOW	SVM	0.01	0.7379
TFIDF	SVM	0.0001	0.73194
ACG W2V	SVM	0.0001	0.71729
TFIDF W2V	SVM	0.001	0.71219
Set 5	SVM	0.0001	0.70352