# 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 | Desc |
|---|---|
| project_id | A unique identifier for the proposed project. **Example:** `p0` |
| project_title | Title of the project. **Exa** <br> • Art Will Make You H <br> • First Grad |
| project_grade_category | Grade level of students for which the project is targeted. One of the fo <br> enumerated v <br> • Grades P <br> • Grade <br> • Grade <br> • Grades |

| Feature | Desc |
|---|---|
| **project_subject_categories** | One or more (comma-separated) subject categories for the project fr following enumerated list of v<br><br>• Applied Lea<br>• Care & H<br>• Health & S<br>• History & C<br>• Literacy & Lan<br>• Math & Sc<br>• Music & The<br>• Special<br>• W<br><br>**Exan**<br><br>• Music & The<br>• Literacy & Language, Math & Sc |
| **school_state** | State where school is located ([Two-letter U.S. post:](https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_c)<br>**Exampl** |
| **project_subject_subcategories** | One or more (comma-separated) subject subcategories for the <br>**Exan**<br><br>• Lit<br>• Literature & Writing, Social Sci |
| **project_resource_summary** | An explanation of the resources needed for the project. **Exa**<br><br>• My students need hands on literacy materials to ma<br>sensory needs!< |
| **project_essay_1** | First application |
| **project_essay_2** | Second application |
| **project_essay_3** | Third application |
| **project_essay_4** | Fourth application |
| **project_submitted_datetime** | Datetime when project application was submitted. **Example:** 2016-<br>12:43:5 |
| **teacher_id** | A unique identifier for the teacher of the proposed project. **Ex:**<br>bdf8baa8fedef6bfeec7ae4ff1c |
| **teacher_prefix** | Teacher's title. One of the following enumerated v<br><br>•<br>•<br>•<br>•<br>•<br>Tea |
| **teacher_number_of_previously_posted_projects** | Number of project applications previously submitted by the same te<br>**Examp** |

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

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

| Feature | Description |
|---|---|
| **id** | A `project_id` value from the `train.csv` file. **Example:** `p036502` |
| **description** | Desciption of the resource. **Example:** `Tenor Saxophone Reeds, Box of 25` |

| Feature | Description |
|---|---|
| quantity | Quantity of the resource required. **Example:** 3 |
| price | Price of the resource required. **Example:** 9.95 |

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

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

| Label | Description |
|---|---|
| project_is_approved | A binary flag indicating whether DonorsChoose approved the project. A value of `0` indicates the project was not approved, and a value of `1` indicates the project was approved. |

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:
- **project_essay_1:** "Introduce us to your classroom"
- **project_essay_2:** "Tell us more about your students"
- **project_essay_3:** "Describe how your students will use the materials you're requesting"
- **project_essay_3:** "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:
- **project_essay_1:** "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- **project_essay_2:** "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [1]:

```python
%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
```

```
C:\Users\MONIKA KUMARI\Anaconda3\lib\site-packages\gensim\utils.py:1197: Use
rWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

## 1.1 Reading Data

In [2]:

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

```
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' 's
chool_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']
```

```
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']
```

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

## 1.2 preprocessing of `project_subject_categories`

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

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

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

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

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

## 1.3 preprocessing of `project_subject_subcategories`

In [6]:

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

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

sub_cat_list = []
for i in sub_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", "
        if 'The' in j.split(): # this will split each of the catogory based on space "Math
            j=j.replace('The','') # if we have the words "The" we are going to replace it w
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math
        temp +=j.strip()+" "+"#" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

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

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

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

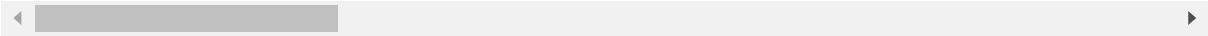## 1.3 Text preprocessing

In [7]:

```python
# 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)
```

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

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project |
|---|---|---|---|---|---|---|
| **0** | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | |
| **1** | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | |

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

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

==================================================

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.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more.With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade.  This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

==================================================

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch.  Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say.Wobble chairs are the answer and I love then 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 mess
age. Due to the volume of my speaker my students can't hear videos or book
s clearly and it isn't making the lessons as meaningful. But with the blue
tooth speaker my students will be able to hear and I can stop, pause and r
eplay it at any time.\r\nThe cart will allow me to have more room for stor
age of things that are needed for the day and has an extra part to it I ca
n 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 [10]:

```python
# 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 [11]:

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

My kindergarten students have varied disabilities ranging from speech and la
nguage delays, cognitive delays, gross/fine motor delays, to autism. They ar
e eager beavers and always strive to work their hardest working past their l
imitations. \r\n\r\nThe materials we have are the ones I seek out for my stu
dents. I teach in a Title I school where most of the students receive free o
r reduced price lunch.  Despite their disabilities and limitations, my stude
nts 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 yo
u were in a meeting? This is how my kids feel all the time. The want to be a
ble to move as they learn or so they say.Wobble chairs are the answer and I
love then 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 d
o 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
==================================================

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

My kindergarten students have varied disabilities ranging from speech and la
nguage delays, cognitive delays, gross/fine motor delays, to autism. They ar
e eager beavers and always strive to work their hardest working past their l
imitations.      The materials we have are the ones I seek out for my student
s. I teach in a Title I school where most of the students receive free or re
duced 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 wer
e in a meeting? This is how my kids feel all the time. The want to be able t
o move as they learn or so they say.Wobble chairs are the answer and I love
then 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 w
ant to sit and do worksheets. They want to learn to count by jumping and pla
ying. Physical engagement is the key to our success. The number toss and col
or and shape mats can make that happen. My students will forget they are doi
ng work and just have the fun a 6 year old deserves.nannan

```python
#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 la
nguage delays cognitive delays gross fine motor delays to autism They are ea
ger beavers and always strive to work their hardest working past their limit
ations 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 a
nts 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 le
arn or so they say Wobble chairs are the answer and I love then 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 work
sheets They want to learn to count by jumping and playing Physical engagemen
t is the key to our success The number toss and color and shape mats can mak
e that happen My students will forget they are doing work and just have the
fun a 6 year old deserves nannan

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

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

```
100%|████████████████████████████| 109248/109248 [02:10<00:00, 835.57i
t/s]
```

```python
# after preprocesing
preprocessed_essays[20000]
```

Out[16]:

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

Out[17]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project |
|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | |

# 1.4 Preprocessing of `project_title`

In [18]:

```
# printing some random project title.
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 [19]:

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

```
We Need To Move It While We Input It!
==================================================
```

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

```
100%|████████████████████████████| 109248/109248 [00:05<00:00, 19655.86i
t/s]
```

In [21]:

```python
#after preprocessing
preprocessed_titles[20000]
```

Out[21]:

```
'we need to move it while we input it'
```

In [22]:

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

Out[22]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_ |
|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | |

## 1.6 Preprocessing teacher_prefix

```
x = project_data['teacher_prefix'].replace(to_replace= np.nan, value= "mrs")
teacher_prefix_list = list(x.values)
preprocessed_teacher_prefix=[]
for l in tqdm (teacher_prefix_list):
    n = ""
    for e in l:
        e = e.replace('.', '')
        e = e.replace(',', '')
        n+= e
    preprocessed_teacher_prefix.append(n.lower().strip())

print(len(preprocessed_teacher_prefix))
```

```
100%|████████████████████████| 109248/109248 [00:00<00:00, 187241.90i
t/s]

109248
```

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

Out[24]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| **0** | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| **1** | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

## 1.7 Preprocessing project_grade_category

```python
project_grade_category_list = list(project_data['project_grade_category'].values)
preprocessed_project_grade_category=[]
for l in tqdm (project_grade_category_list):
    n = ""
    for e in l:
        e = e.replace(' ', '_')
        e = e.replace('-', '_')
        n+= e
    preprocessed_project_grade_category.append(n.lower().strip())

print(len(preprocessed_project_grade_category))
```

```
100%|████████████████████████| 109248/109248 [00:01<00:00, 71031.06i
t/s]

109248
```

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

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

## 2. Computing Sentiment Scores for preprocessed_essays

```python
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer
#nltk.download('vader_lexicon')
sid = SentimentIntensityAnalyzer()
neg= []
pos= []
neu= []
compound= []
for k in tqdm(project_data['preprocessed_essays'], position= 0, leave= True):
    a= sid.polarity_scores(k)['neg']
    b= sid.polarity_scores(k)['pos']
    c= sid.polarity_scores(k)['neu']
    d= sid.polarity_scores(k)['compound']
    neg.append(a)
    pos.append(b)
    neu.append(c)
    compound.append(d)
```

C:\Users\MONIKA KUMARI\Anaconda3\lib\site-packages\nltk\twitter\__init__.py:
20: UserWarning:

The twython library has not been installed. Some functionality from the twit
ter package will not be available.

100%|████████████████████████████████| 109248/109248 [28:33<00:00, 63.74i
t/s]

```python
project_data['neg'] = neg
project_data['pos'] = pos
project_data['neu'] = neu
project_data['compound'] = compound
```

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

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

2 rows × 22 columns

## 2.1 Number of words in combined essays

```python
word_count_essay = []
for w in project_data['preprocessed_essays']:
    s = len(w.split())
    word_count_essay.append(s)
```
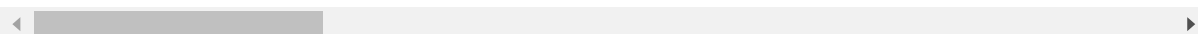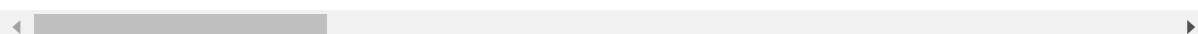
```python
project_data['word_count_essay'] = word_count_essay
project_data.head(2)
```

Out[31]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

2 rows × 23 columns

### 2.2 Number of words in preprocessed_titles

```python
word_count_title= []
for w in project_data['preprocessed_title']:
    s = len(w.split())
    word_count_title.append(s)
#print(word_count_title)
```

```python
project_data['word_count_title'] = word_count_title
project_data.head(2)
```

Out[33]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

2 rows × 24 columns

In [34]:

```python
# merging resource_data with project_data
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index
project_data = pd.merge(project_data, price_data, on='id', how='left')
```
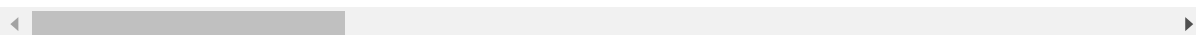
In [35]:

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

Out[35]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

2 rows × 26 columns

## 3. Preparing data for models

In [36]:

```python
project_data.columns
```

Out[36]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
       'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
       'project_essay_3', 'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approve
d',
       'clean_categories', 'clean_subcategories', 'preprocessed_essays',
       'preprocessed_title', 'preprocessed_teacher_prefix',
       'preprocessed_project_grade_category', 'neg', 'pos', 'neu', 'compoun
d',
       'word_count_essay', 'word_count_title', 'price', 'quantity'],
      dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data

- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

# 3.1 Loading Data

In [37]:

```
project_data.to_csv('data.csv',index=False)
```

In [38]:

```
data = pd.read_csv('data.csv')
data.head(2)
```

Out[38]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_date |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:4 |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | FL | 2016-10-25 09:2 |

2 rows × 26 columns

# 3.2 Splitting data into Train and cross validation(or test): Stratified Sampling
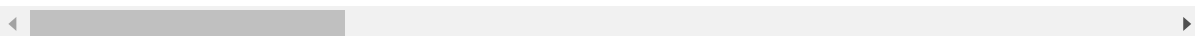
In [39]:

```python
y = data['project_is_approved'].values
X = data.drop(['project_is_approved'], axis=1)
X.head(1)
```

Out[39]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_datetin |
|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:43: |

1 rows × 25 columns

In [40]:

```python
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_spli
from sklearn.model_selection import train_test_split

# X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.33, shuffle=Flase)#
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify= y) # th
#X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify


print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)
```

```
(73196, 25) (73196,)
(36052, 25) (36052,)
```

In [41]:

```python
#value count in numpy array #https://stackoverflow.com/a/28663910
unique, counts = np.unique(y_train, return_counts=True) #counts the majority and minority c
dict(zip(unique, counts))
```

Out[41]:

```
{0: 11083, 1: 62113}
```

## 3.3. Vectorizing Categorical data

- https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/ (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/)

## 3.3.1. encoding categorical features: clean_category

```
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
(73196, 9) (73196,)
(36052, 9) (36052,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'liter
acy_language', 'math_science', 'music_arts', 'specialneeds', 'warmth']
================================================================================
========================
```

### 3.3.2. encoding categorical features: clean_subcategory

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

# 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
(73196, 30) (73196,)
(36052, 30) (36052,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_governmen
t', 'college_careerprep', 'communityservice', 'earlydevelopment', 'economic
s', 'environmentalscience', 'esl', 'extracurricular', 'financialliteracy',
'foreignlanguages', 'gym_fitness', 'health_lifescience', 'health_wellness',
'history_geography', 'literacy', 'literature_writing', 'mathematics', 'musi
c', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 's
ocialsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
================================================================================
========================
```

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

### 3.3.3. encoding categorical features: School State

```
state_vectorizer = CountVectorizer()
state_vectorizer.fit(X_train['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 = state_vectorizer.transform(X_train['school_state'].values)
#X_cv_state_ohe = vectorizer.transform(X_cv['school_state'].values)
X_test_state_ohe = state_vectorizer.transform(X_test['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(state_vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(73196, 51) (73196,)
(36052, 51) (36052,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'i
a', 'id', 'il', 'in', 'ks', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo',
'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', 'ny', 'oh', 'ok', 'o
r', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv',
'wy']
================================================================================
=======================
```

### 3.3.4 encoding categorical features: teacher_prefix

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['preprocessed_teacher_prefix'].values) # fit has to happen only on t

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer.transform(X_train['preprocessed_teacher_prefix'].values)
#X_cv_teacher_ohe = vectorizer.transform(X_cv['preprocessed_teacher_prefix'].values)
X_test_teacher_ohe = vectorizer.transform(X_test['preprocessed_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
(73196, 5) (73196,)
(36052, 5) (36052,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
========================================================================
=======================
```

```
project_data['preprocessed_teacher_prefix'].value_counts()
```

```
mrs         57272
ms          38955
mr          10648
teacher      2360
dr             13
Name: preprocessed_teacher_prefix, dtype: int64
```

### 3.3.5. encoding categorical features: project_grade_category

```
grade_vectorizer = CountVectorizer()
grade_vectorizer.fit(X_train['preprocessed_project_grade_category'].values) # fit has to ha

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

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(grade_vectorizer.get_feature_names())
print("="*100)
```

```
After vectorizations
(73196, 4) (73196,)
(36052, 4) (36052,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
================================================================================
========================
```

## 3.4. Vectorizing Text data

### 3.4.1. Bag of words on preprocessed_essays

```python
# We are considering only the words which appeared in at least 10 documents(rows or project
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)


vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2), max_features = 5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train dat

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['preprocessed_essays'].values)
#X_cv_essay_bow = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_bow = vectorizer.transform(X_test['preprocessed_essays'].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("="*100)
```

```
(73196, 25) (73196,)
(36052, 25) (36052,)
========================================================================
======================
After vectorizations
(73196, 5000) (73196,)
(36052, 5000) (36052,)
========================================================================
======================
```

**3.4.2 Bag of words on preprocessed_titles**

```python
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)


vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2),max_features = 5000)
vectorizer.fit(X_train['preprocessed_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['preprocessed_title'].values)
#X_cv_title_bow = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_title_bow = vectorizer.transform(X_test['preprocessed_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("="*100)
```

```
(73196, 25) (73196,)
(36052, 25) (36052,)
================================================================================
=======================
After vectorizations
(73196, 5000) (73196,)
(36052, 5000) (36052,)
================================================================================
=======================
```

### 3.4.3 TFIDF vectorizer on preprocessed_essays

```python
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,2),max_features = 5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train dat

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['preprocessed_essays'].values)
#X_cv_essay_tfidf = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
#print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
(73196, 25) (73196,)
(36052, 25) (36052,)
================================================================================
========================
After vectorizations
(73196, 5000) (73196,)
(36052, 5000) (36052,)
================================================================================
========================
```

**3.4.4 TFIDF vectorization on preprocessed_titles**

```python
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=5,ngram_range=(1,2),max_features = 5000)
vectorizer.fit(X_train['preprocessed_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_titles_tfidf = vectorizer.transform(X_train['preprocessed_title'].values)
#X_cv_titles_tfidf = vectorizer.transform(X_cv['preprocessed_titles'].values)
X_test_titles_tfidf = vectorizer.transform(X_test['preprocessed_title'].values)

print("After vectorizations")
print(X_train_titles_tfidf.shape, y_train.shape)
#print(X_cv_titles_tfidf.shape, y_cv.shape)
print(X_test_titles_tfidf.shape, y_test.shape)
print("="*100)
```

```
(73196, 25) (73196,)
(36052, 25) (36052,)
================================================================================
======================
After vectorizations
(73196, 5000) (73196,)
(36052, 5000) (36052,)
================================================================================
======================
```

### 3.4.5. Using Pretrained Models: Avg W2V

In [53]:

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

# ==============================
Output:

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

# ==============================

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

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

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

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


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

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


'''
```

Out[53]:

'\n# Reading glove vectors in python: https://stackoverflow.com/a/3823034

In [54]:

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

In [55]:

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

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

```
100%|██████████████████████████████| 109248/109248 [01:06<00:00, 1649.36i
t/s]

109248
300
```

**3.4.5.1 Using Pretrained Models: Avg W2V on preprocessed_essays**

In [56]:

```python
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essays_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays'].values, position= 0, leave= True): # fc
    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_essays_train.append(vector)

print(len(avg_w2v_essays_train))
print(len(avg_w2v_essays_train[0]))
```

```
100%|████████████████████████████| 73196/73196 [00:43<00:00, 1668.03i
t/s]

73196
300
```

In [57]:

```python
avg_w2v_essays_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays'].values, position= 0, leave= True): # for
    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_essays_test.append(vector)
print(len(avg_w2v_essays_test))
print(len(avg_w2v_essays_test[0]))
```

```
100%|████████████████████████████| 36052/36052 [00:21<00:00, 1661.04i
t/s]

36052
300
```

**3.4.5.2 Using Pretrained Models: Avg W2V on project_title**

```python
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_titles_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_title'].values, position= 0, leave= True): # for
    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_titles_train.append(vector)

print(len(avg_w2v_titles_train))
print(len(avg_w2v_titles_train[0]))
#print(avg_w2v_titles_train[0])
```

```
100%|████████████████████████████| 73196/73196 [00:02<00:00, 29300.74i
t/s]

73196
300
```

```python
avg_w2v_titles_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_title'].values, position= 0, leave= True): # for
    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_titles_test.append(vector)
```

```
100%|████████████████████████████| 36052/36052 [00:01<00:00, 23840.23i
t/s]
```

**3.4.5.3 Using Pretrained Models: TFIDF weighted W2V on preprocessed_essays**

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

```python
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essay_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays'].values, position= 0, leave= True): # fo
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each words
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentenc
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettir
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay_train.append(vector)

print(len(tfidf_w2v_essay_train))
print(len(tfidf_w2v_essay_train[0]))
```

```
100%|████████████████████████| 73196/73196 [05:18<00:00, 229.86i
t/s]

73196
300
```

```python
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essay_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays'].values, position= 0, leave= True): # for
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each words
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentenc
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettir
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay_test.append(vector)

print(len(tfidf_w2v_essay_test))
print(len(tfidf_w2v_essay_test[0]))
```

```
100%|████████████████████████| 36052/36052 [02:35<00:00, 231.75i
t/s]

36052
300
```

### 3.4.5.4 Using Pretrained Models: TFIDF weighted W2V on Project_title

In [63]:

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

In [64]:

```python
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_title'].values, position= 0, leave= True): # for
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentenc
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettir
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_title_train.append(vector)

print(len(tfidf_w2v_title_train))
print(len(tfidf_w2v_title_train[0]))
```

```
100%|████████████████████████████| 73196/73196 [00:04<00:00, 14644.44i
t/s]

73196
300
```

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_title'].values, position= 0, leave= True): # for
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((senten
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # gettin
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_title_test.append(vector)

print(len(tfidf_w2v_title_test))
print(len(tfidf_w2v_title_test[0]))
```

```
100%|██████████████████████████████████| 36052/36052 [00:02<00:00, 14400.61i
t/s]

36052
300
```

## 3.5. Vectorizing Numerical features

In [66]:

```
data.head(1)
```

Out[66]:

| | Unnamed: 0 | id | teacher_id | school_state | project_submitted_datetin |
|---|---|---|---|---|---|
| **0** | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | IN | 2016-12-05 13:43: |

1 rows × 26 columns

### 3.5.1. encoding numerical features: price

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

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

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

# Now standardize the data with above maen and variance.
price_standardized_train = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
price_standardized_test = price_scalar.transform(X_test['price'].values.reshape(-1, 1))
print(price_standardized_train.shape,y_train.shape)
print(price_standardized_test.shape,y_test.shape)
price_scalar = StandardScaler()
price_scalar.fit(project_data['price'].values.reshape(-1,1)) # finding the mean and standar
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])

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

```
Mean : 297.89765683917153, Standard deviation : 371.09985227095217
(73196, 1) (73196,)
(36052, 1) (36052,)
Mean : 298.1193425966608, Standard deviation : 367.49634838483496
```

In [68]:

```python
price_standardized_train
```

Out[68]:

```
array([[-0.39888902],
       [-0.5967064 ],
       [-0.25986983],
       ...,
       [-0.5761459 ],
       [-0.65925021],
       [-0.21489542]])
```

**3.5.2 encoding numerical features: teacher_number_of_previously_posted_projects**

In [69]:

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

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

prev_projects_scalar = StandardScaler()
prev_projects_scalar.fit((X_train['teacher_number_of_previously_posted_projects'].values.as
print(f"Mean : {prev_projects_scalar.mean_[0]}, Standard deviation : {np.sqrt(prev_projects

# Now standardize the data with above maen and variance.
prev_projects_standardized_train = prev_projects_scalar.transform((X_train['teacher_number_
prev_projects_standardized_test = prev_projects_scalar.transform((X_test['teacher_number_of
print(prev_projects_standardized_train.shape,y_train.shape)
print(prev_projects_standardized_test.shape,y_test.shape)
```

```
Mean : 11.19103502923657, Standard deviation : 27.991769358648625
(73196, 1) (73196,)
(36052, 1) (36052,)
```

In [70]:

```
print(prev_projects_standardized_train)
```

```
[[-0.39979734]
 [ 1.70796509]
 [-0.39979734]
 ...
 [-0.32834777]
 [ 0.10034968]
 [-0.36407256]]
```

### 3.5.3 encoding numerical features: quantity

In [71]:

```
quantity_scalar = StandardScaler()
quantity_scalar.fit((X_train['quantity'].values.astype(float)).reshape(-1,1)) # finding the
print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(quantity_scalar.va

# Now standardize the data with above maen and variance.
quantity_standardized_train = quantity_scalar.transform((X_train['quantity'].values.astype(
quantity_standardized_test = quantity_scalar.transform((X_test['quantity'].values.astype(fl
print(quantity_standardized_train.shape,y_train.shape)
print(quantity_standardized_test.shape,y_test.shape)
```

```
Mean : 16.99729493414941, Standard deviation : 26.099569279799027
(73196, 1) (73196,)
(36052, 1) (36052,)
```

```
print(quantity_standardized_train)
```

```
[[ 1.11123309]
 [-0.11484078]
 [-0.45967406]
 ...
 [-0.1914704 ]
 [ 1.03460347]
 [-0.26810002]]
```

## 3.6. Merging all the above features

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

In [73]:

```
print(X_train_categories_ohe.shape)
print(X_test_categories_ohe.shape)
print(X_train_subcategories_ohe.shape)
print(X_test_subcategories_ohe.shape)
print(X_train_state_ohe.shape)
print(X_test_state_ohe.shape)
print(X_train_grade_ohe.shape)
print(X_test_grade_ohe.shape)

print(X_train_essay_bow.shape)
print(X_test_essay_bow.shape)
print(price_standardized_train.shape)
print(price_standardized_test.shape)
```

```
(73196, 9)
(36052, 9)
(73196, 30)
(36052, 30)
(73196, 51)
(36052, 51)
(73196, 4)
(36052, 4)
(73196, 5000)
(36052, 5000)
(73196, 1)
(36052, 1)
```

```python
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_1 = hstack((X_train_essay_bow, X_train_title_bow, X_train_categories_ohe, X_train_subc
#X_cr = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_resources_bow, X_cv_categories_ohe, X_
X_te_1 = hstack((X_test_essay_bow, X_test_title_bow, X_test_categories_ohe, X_test_subcateg

print("Final Data matrix")
print(X_tr_1.shape, y_train.shape)
#print(X_cr.shape, y_cv.shape)
print(X_te_1.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(73196, 10102) (73196,)
(36052, 10102) (36052,)
========================================================================
======================
```

```python
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_2 = hstack((X_train_essay_tfidf, X_train_titles_tfidf, X_train_categories_ohe, X_train
#X_cr = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_resources_bow, X_cv_categories_ohe, X_
X_te_2 = hstack((X_test_essay_tfidf, X_test_titles_tfidf, X_test_categories_ohe, X_test_sub

print("Final Data matrix")
print(X_tr_2.shape, y_train.shape)
#print(X_cr.shape, y_cv.shape)
print(X_te_2.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(73196, 10102) (73196,)
(36052, 10102) (36052,)
========================================================================
======================
```

```
from scipy.sparse import hstack
X_tr_3 = hstack((avg_w2v_essays_train, avg_w2v_titles_train, X_train_categories_ohe, X_trai
#X_cr = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_resources_bow, X_cv_categories_ohe, X_
X_te_3 = hstack((avg_w2v_essays_test, avg_w2v_titles_test, X_test_categories_ohe, X_test_su

print("Final Data matrix")
print(X_tr_3.shape, y_train.shape)
#print(X_cr.shape, y_cv.shape)
print(X_te_3.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(73196, 702) (73196,)
(36052, 702) (36052,)
========================================================================
=======================
```

```
from scipy.sparse import hstack
X_tr_4 = hstack((tfidf_w2v_essay_train, tfidf_w2v_title_train, X_train_categories_ohe, X_tr
#X_cr = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_resources_bow, X_cv_categories_ohe, X_
X_te_4 = hstack((tfidf_w2v_essay_test, tfidf_w2v_title_test, X_test_categories_ohe, X_test_

print("Final Data matrix")
print(X_tr_4.shape, y_train.shape)
#print(X_cr.shape, y_cv.shape)
print(X_te_4.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(73196, 702) (73196,)
(36052, 702) (36052,)
========================================================================
=======================
```

**Computing Sentiment Scores**

In [78]:

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

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

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest stu
for learning my students learn in many different ways using all of our senses and multiple
of techniques to help all my students succeed students in my class come from a variety of d
for wonderful sharing of experiences and cultures including native americans our school is
learners which can be seen through collaborative student project based learning in and out
in my class love to work with hands on materials and have many different opportunities to p
mastered having the social skills to work cooperatively with friends is a crucial aspect of
montana is the perfect place to learn about agriculture and nutrition my students love to r
in the early childhood classroom i have had several kids ask me can we try cooking with rea
and create common core cooking lessons where we learn important math and writing concepts w
food for snack time my students will have a grounded appreciation for the work that went in
of where the ingredients came from as well as how it is healthy for their bodies this proje
nutrition and agricultural cooking recipes by having us peel our own apples to make homemad
and mix up healthy plants from our classroom garden in the spring we will also create our o
shared with families students will gain math and literature skills as well as a life long e
nannan'
ss = sid.polarity_scores(for_sentiment)

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

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

Out[78]:

```
"import nltk\nfrom nltk.sentiment.vader import SentimentIntensityAnalyzer
\n\n# import nltk\n# nltk.download('vader_lexicon')\n\nsid = SentimentInte
nsityAnalyzer()\n\nfor_sentiment = 'a person is a person no matter how sma
ll dr seuss i teach the smallest students with the biggest enthusiasm for
learning my students learn in many different ways using all of our senses
and multiple intelligences i use a wide rangeof techniques to help all my
students succeed students in my class come from a variety of different bac
kgrounds which makesfor wonderful sharing of experiences and cultures incl
uding native americans our school is a caring community of successful lear
ners which can be seen through collaborative student project based learnin
g in and out of the classroom kindergarteners in my class love to work wit
h hands on materials and have many different opportunities to practice a s
kill before it ismastered having the social skills to work cooperatively w
ith friends is a crucial aspect of the kindergarten curriculummontana is t
he perfect place to learn about agriculture and nutrition my students love
to role play in our pretend kitchenin the early childhood classroom i have
had several kids ask me can we try cooking with real food i will take thei
r idea and create common core cooking lessons where we learn important mat
h and writing concepts while cooking delicious healthy food for snack time
my students will have a grounded appreciation for the work that went into
making the food and knowledge of where the ingredients came from as well a
s how it is healthy for their bodies this project would expand our learnin
g of nutrition and agricultural cooking recipes by having us peel our own
```

apples to make homemade applesauce make our own bread and mix up healthy p
lants from our classroom garden in the spring we will also create our own
cookbooks to be printed and shared with families students will gain math a
nd literature skills as well as a life long enjoyment for healthy cooking
nannan'\nss = sid.polarity_scores(for_sentiment)\n\nfor k in ss:\n    prin
t('{0}: {1}, '.format(k, ss[k]), end='')\n\n# we can use these 4 things as
features/attributes (neg, neu, pos, compound)\n# neg: 0.0, neu: 0.753, po
s: 0.247  compound: 0.93\n"

# 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 (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) 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 (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.

     

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

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 (http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html)](http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html) on [TfidfVectorizer (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html)](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html) of essay text, choose the number of components ( `n_components` ) using [elbow method (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/)** : 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 (http://zetcode.com/python/prettytable/)](http://zetcode.com/python/prettytable/)

**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. (https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf)](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf)

# 4. Support Vector Machines

## 4.1 Appling Support Vector Machines on different kind of featurization as mentioned in the instructions

Apply Support Vector Machines on different kind of featurization as mentioned in the instructions
For Every model that you work on make sure you do the step 2 and step 3 of instrucations

### 4.1.1 Applying SVM on BOW (SET 1)

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

svc = SGDClassifier(loss='hinge', penalty='l2', class_weight = 'balanced')
#clf= GridSearchCV(svc, parameters, cv= 10, scoring= 'roc_auc')

def hyper_param_tuning(X_train_sample, y_train_sample):
    #svc = SGDClassifier(loss='hinge', penalty='l2')
    parameters = [{'alpha': [10**-4, 10**-2, 10**-1, 10**0, 10**-1, 10**2, 10**4]}]
    clf= GridSearchCV(svc, parameters, cv= 10, scoring= 'roc_auc')
    clf.fit(X_train_sample, y_train_sample)

    results = pd.DataFrame.from_dict(clf.cv_results_)
    results = results.sort_values(['param_alpha'])

    train_auc= results['mean_train_score']
    train_auc_std= results['std_train_score']
    cv_auc = results['mean_test_score']
    cv_auc_std= results['std_test_score']
    a =  results['param_alpha']

    alpha = [10**-4, 10**-2, 10**-1, 10**0, 10**-1, 10**2, 10**4]
    log_a= []
    for a in alpha:
        log_a.append(np.log10(a))
    #print(log_a)
    plt.plot(log_a, train_auc, label='Train AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
    # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0
    plt.plot(log_a, cv_auc, label='Test AUC')
    # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
    # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='da
    plt.scatter(log_a, train_auc, label='Train AUC points')
    plt.scatter(log_a, cv_auc, label='Test AUC points')

    plt.legend()
    plt.xlabel("log_a: hyperparameter")
    #plt.xscale('log',basex=10)
    plt.ylabel("AUC")
    plt.title("Hyper parameter Vs AUC plot")
    plt.grid()
    plt.show()
    print('Best score: ',clf.best_score_)
    print('value of alpha with best score: ',clf.best_params_)
    print('='*75)
    print('Train AUC scores')
    print(results['mean_train_score'])
    print('CV AUC scores')
    print(results['mean_test_score'])
    return results.head()
```
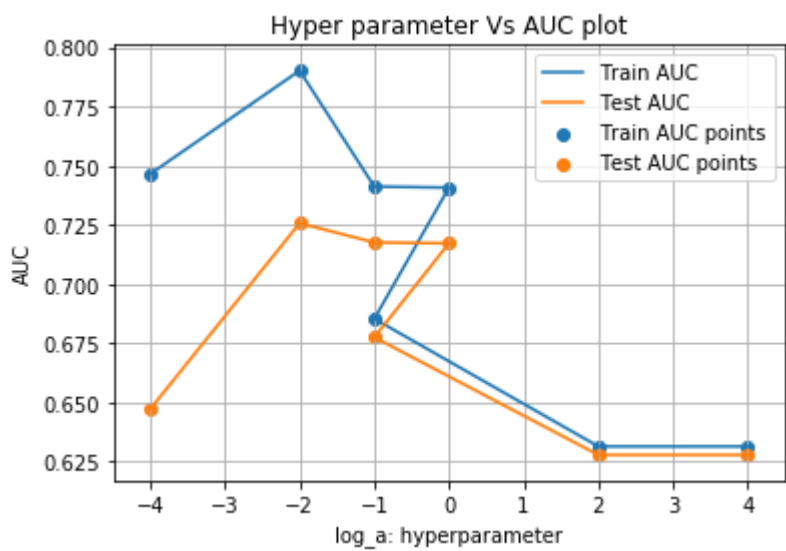
**4.1.1.1 Hyperparameter tuning Using L2 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l2', class_weight = 'balanced')
hyper_param_tuning(X_tr_1, y_train)
```



```
Best score:  0.725779239202959
value of alpha with best score:  {'alpha': 0.01}
================================================================================
Train AUC scores
0    0.746591
1    0.790243
2    0.741276
4    0.740822
3    0.685412
5    0.631464
6    0.631464
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.647265
1    0.725779
2    0.717650
4    0.717296
3    0.677577
5    0.627882
6    0.627881
Name: mean_test_score, dtype: float64
```

Out[80]:

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split |
|---|---|---|---|---|---|---|---|
| **0** | 0.658245 | 0.188894 | 0.007810 | 0.014398 | 0.0001 | {'alpha': 0.0001} | |
| **1** | 0.621843 | 0.013617 | 0.006246 | 0.007650 | 0.01 | {'alpha': 0.01} | |
| **2** | 0.648396 | 0.016010 | 0.009374 | 0.007653 | 0.1 | {'alpha': 0.1} | |

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split |
|---|---|---|---|---|---|---|---|
| **4** | 0.642149 | 0.012980 | 0.006246 | 0.007650 | 0.1 | {'alpha': 0.1} | |
| **3** | 0.656212 | 0.009883 | 0.009372 | 0.007653 | 1 | {'alpha': 1} | |

5 rows × 31 columns

## 4.1.1.2 Hyperparameter tuning Using L1 regulariser

```
svc = SGDClassifier(loss='hinge', penalty='l1', class_weight = 'balanced')
hyper_param_tuning(X_tr_1, y_train)
```



```
Best score:  0.6470730984440529
value of alpha with best score:  {'alpha': 0.0001}
===========================================================================
Train AUC scores
0    0.738338
1    0.624529
2    0.541795
4    0.536300
3    0.500000
5    0.500000
6    0.500000
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.647073
1    0.624630
2    0.541721
4    0.538469
3    0.500000
5    0.500000
6    0.500000
Name: mean_test_score, dtype: float64
```

Out[81]:

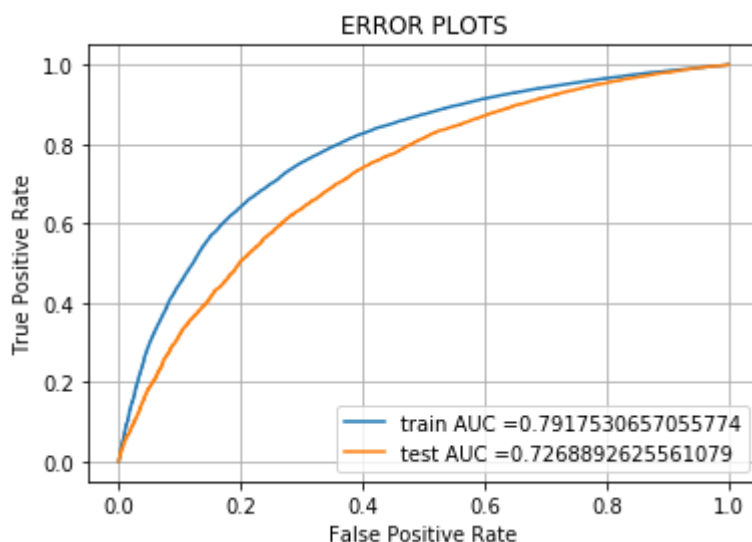| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| **0** | 1.514467 | 0.039808 | 0.006260 | 0.007667 | 0.0001 | {'alpha': 0.0001} | |
| **1** | 0.987989 | 0.017166 | 0.009372 | 0.007653 | 0.01 | {'alpha': 0.01} | |
| **2** | 0.959317 | 0.020005 | 0.007811 | 0.007811 | 0.1 | {'alpha': 0.1} | |
| **4** | 0.954625 | 0.033835 | 0.003123 | 0.006245 | 0.1 | {'alpha': 0.1} | |
| **3** | 0.976503 | 0.012599 | 0.007808 | 0.007808 | 1 | {'alpha': 1} | |

5 rows × 31 columns

### 4.1.1.3 Plotting Roc curve for Set 1

In [82]:

```python
def prob_predict(clf, data):
    y_data_pred = []
    y_data_pred.extend(clf.predict_proba(data)[:,1])
    return y_data_pred
```

In [83]:

```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.
# https://stackoverflow.com/questions/55250963/how-to-get-probabilities-for-sgdclassifier-l
from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV

svc_bow = SGDClassifier(alpha = 0.01, penalty= 'l2', class_weight = 'balanced')
model = CalibratedClassifierCV(svc_bow)
model.fit(X_tr_1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the p
# not the predicted outputs

y_train_pred = prob_predict(model, X_tr_1)
y_test_pred = prob_predict(model, X_te_1)

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

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("ERROR PLOTS")
plt.grid()
plt.show()
```

```python
# we are writing our own function for predict, with defined thresould
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshould, fpr, tpr):
    t = threshould[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
    return t

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

```python
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) 0.529106231378851 for threshold 0.83
```
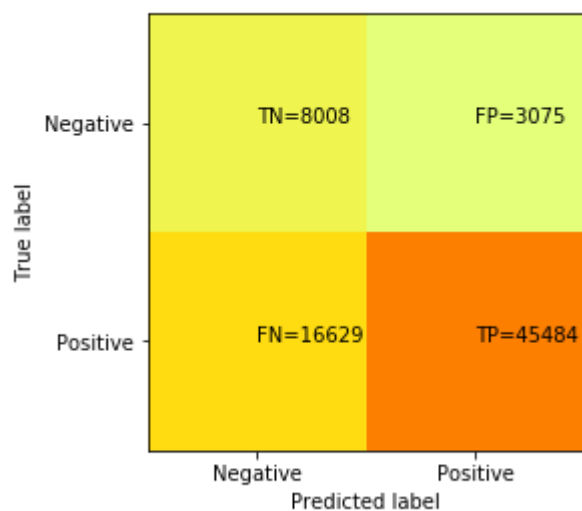
**4.1.1.4 Confusion Matrix**

```python
#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/
def get_confusion_matrix(y_train, y_train_pred):
    '''
    Generate matrix plot of confusion matrix with predicted and original labels.
    '''
    cm = confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
    plt.clf()
    plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Wistia)
    classNames = ['Negative','Positive']
    plt.ylabel('True label')
    plt.xlabel('Predicted label')
    tick_marks = np.arange(len(classNames))
    plt.xticks(tick_marks, classNames, rotation=0)
    plt.yticks(tick_marks, classNames)
    s = [['TN','FP'], ['FN', 'TP']]
    for i in range(2):
        for j in range(2):
            plt.text(j,i, str(s[i][j])+"="+str(cm[i][j]))
    plt.show()
```

```
print("TRAIN CONFUSION MATRIX")
get_confusion_matrix(y_train, y_train_pred)
```

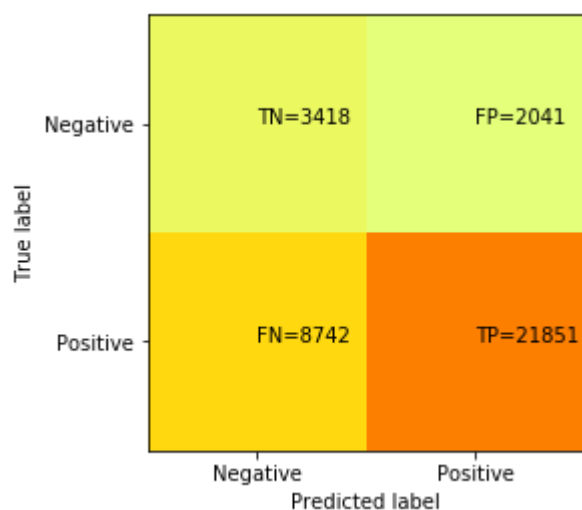TRAIN CONFUSION MATRIX

```
print("TEST CONFUSION MATRIX")
get_confusion_matrix(y_test, y_test_pred)
```
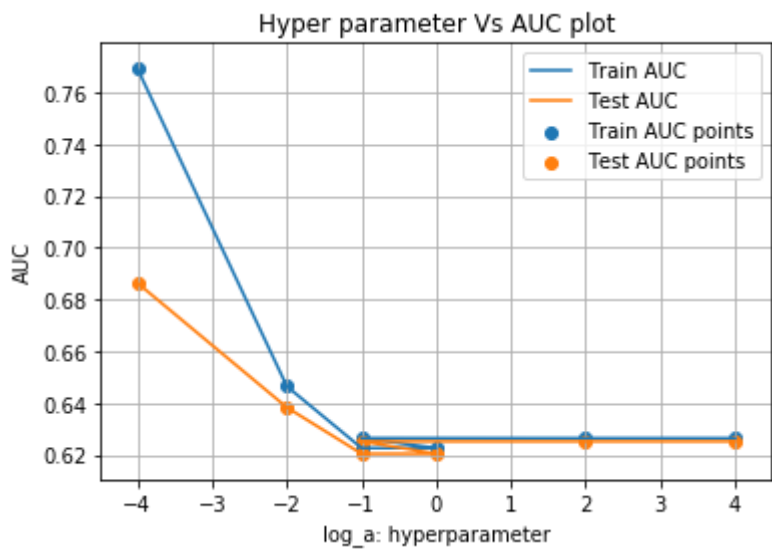
TEST CONFUSION MATRIX



## 4.1.2 Applying SVM on tfidf (SET 2)

**4.1.2.1 Hyperparameter tuning Using L2 regulariser**

```
# error plot for set 2
svc = SGDClassifier(loss='hinge', penalty='l2',class_weight = 'balanced')
hyper_param_tuning(X_tr_2, y_train)
```



```
Best score:  0.6863276312391463
value of alpha with best score:  {'alpha': 0.0001}
==============================================================================
Train AUC scores
0    0.769112
1    0.646413
2    0.622596
4    0.622582
3    0.626420
5    0.626386
6    0.626386
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.686328
1    0.638273
2    0.620387
4    0.620393
3    0.625160
5    0.625162
6    0.625162
Name: mean_test_score, dtype: float64
```
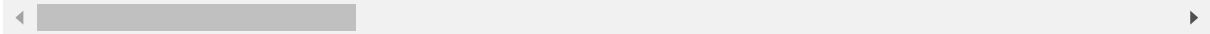
Out[98]:

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| 0 | 0.642389 | 0.021447 | 0.007883 | 0.007027 | 0.0001 | {'alpha': 0.0001} | |
| 1 | 0.662593 | 0.014622 | 0.009481 | 0.005867 | 0.01 | {'alpha': 0.01} | |
| 2 | 0.673175 | 0.013060 | 0.006331 | 0.006875 | 0.1 | {'alpha': 0.1} | |

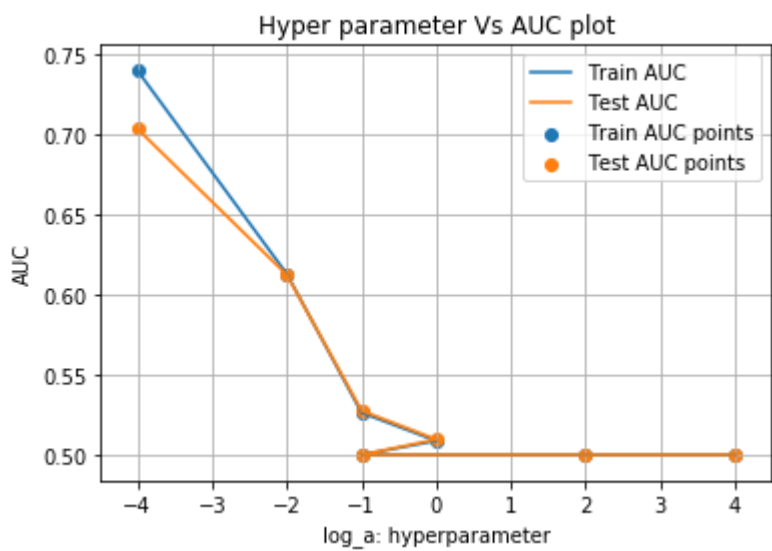| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| **4** | 0.664292 | 0.011695 | 0.007884 | 0.006048 | 0.1 | {'alpha': 0.1} | |
| **3** | 0.670337 | 0.019402 | 0.009433 | 0.007704 | 1 | {'alpha': 1} | |

5 rows × 31 columns

◀ ░░░░░░░░░ ▶

## 4.1.2.2 Hyperparameter tuning Using L1 regulariser

```
svc = SGDClassifier(loss='hinge', penalty='l1',class_weight = 'balanced')
hyper_param_tuning(X_tr_2, y_train)
```



Hyper parameter Vs AUC plot

```
Best score:  0.7031496325627694
value of alpha with best score:  {'alpha': 0.0001}
===========================================================================
Train AUC scores
0    0.739473
1    0.612683
2    0.526279
4    0.508714
3    0.500000
5    0.500000
6    0.500000
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.703150
1    0.612191
2    0.527543
4    0.509593
3    0.500000
5    0.500000
6    0.500000
Name: mean_test_score, dtype: float64
```

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| 0 | 1.313854 | 0.026492 | 0.008604 | 0.007360 | 0.0001 | {'alpha': 0.0001} | |
| 1 | 0.953073 | 0.026148 | 0.007807 | 0.007807 | 0.01 | {'alpha': 0.01} | |
| 2 | 1.000116 | 0.076159 | 0.004720 | 0.007210 | 0.1 | {'alpha': 0.1} | |
| 4 | 0.971824 | 0.015306 | 0.007804 | 0.007804 | 0.1 | {'alpha': 0.1} | |
| 3 | 0.976510 | 0.016008 | 0.006244 | 0.007648 | 1 | {'alpha': 1} | |

5 rows × 31 columns

## 4.1.2.3 Plotting Roc curve for Set 2

In [100]:

```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.
# https://stackoverflow.com/questions/55250963/how-to-get-probabilities-for-sgdclassifier-l
from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV

svc_bow = SGDClassifier(alpha = 0.0001, penalty= 'l1',class_weight = 'balanced')
model = CalibratedClassifierCV(svc_bow)
model.fit(X_tr_2, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the p
# not the predicted outputs

y_train_pred = prob_predict(model, X_tr_2)
y_test_pred = prob_predict(model, X_te_2)

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("ERROR PLOTS")
plt.grid()
plt.show()
```
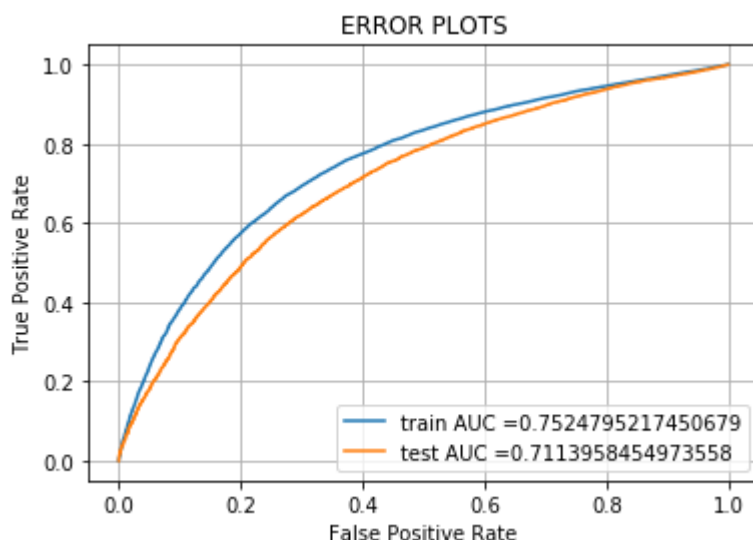


## 4.1.2.4 Confusion Matrix
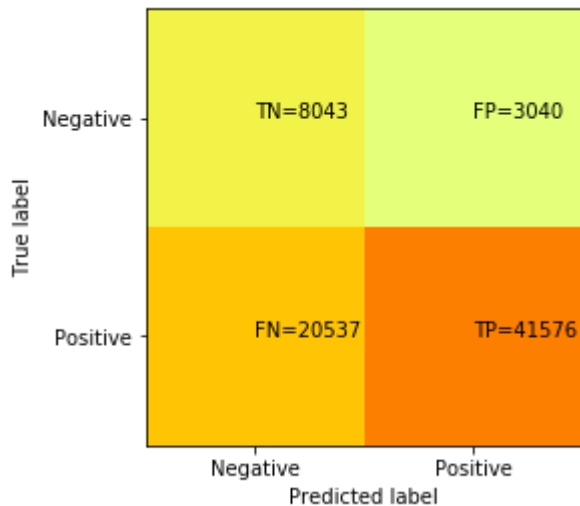
```
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

the maximum value of tpr*(1-fpr) 0.48575908689058667 for threshold 0.842

In [102]:

```
print("TRAIN CONFUSION MATRIX")
get_confusion_matrix(y_train, y_train_pred)
```
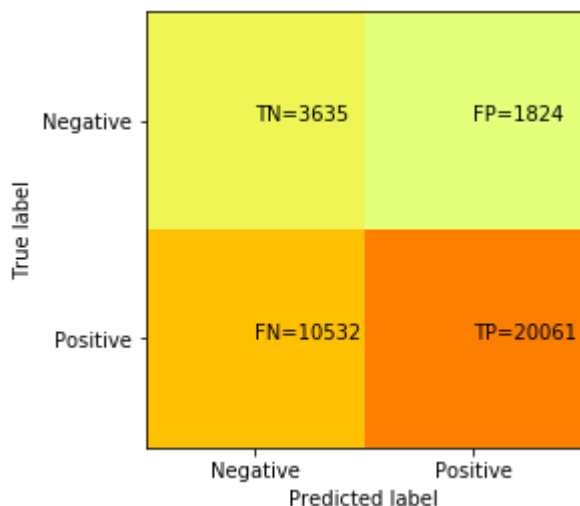
TRAIN CONFUSION MATRIX



In [103]:

```
print("TEST CONFUSION MATRIX")
get_confusion_matrix(y_test, y_test_pred)
```
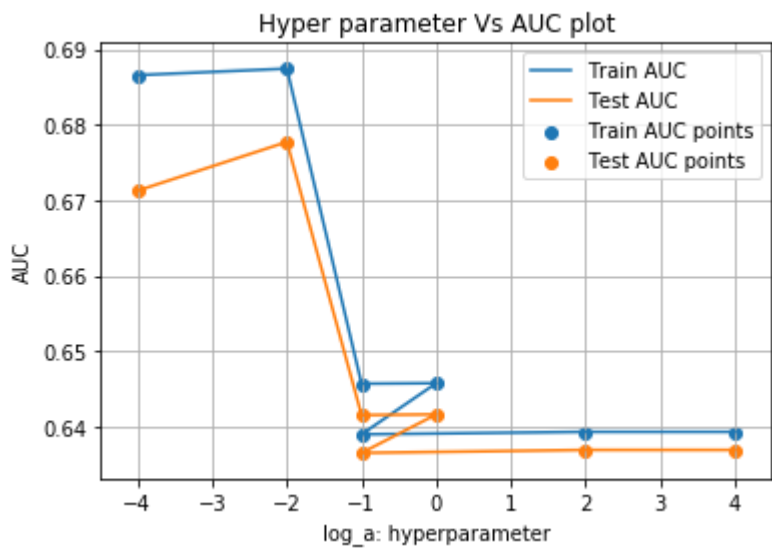
TEST CONFUSION MATRIX



## 4.1.3 Applying SVM on Avg W2V (SET 3)

**4.1.3.1 Hyperparameter tuning Using L2 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l2',class_weight = 'balanced')
hyper_param_tuning(X_tr_3, y_train)
```



Hyper parameter Vs AUC plot

```
Best score:  0.6777381043152693
value of alpha with best score:  {'alpha': 0.01}
=============================================================================
Train AUC scores
0     0.686574
1     0.687436
2     0.645749
4     0.645817
3     0.639049
5     0.639377
6     0.639377
Name: mean_train_score, dtype: float64
CV AUC scores
0     0.671285
1     0.677738
2     0.641598
4     0.641685
3     0.636589
5     0.636982
6     0.636981
Name: mean_test_score, dtype: float64
```

Out[105]:

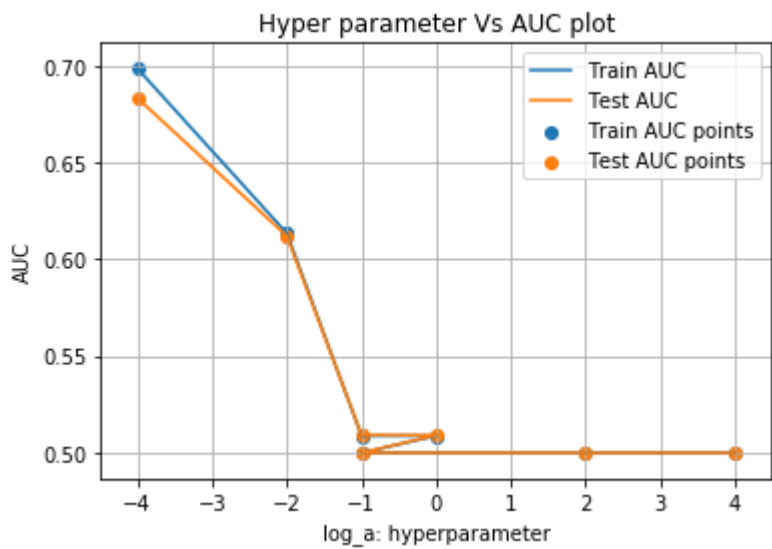| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| **0** | 1.809597 | 0.121142 | 0.021262 | 0.006171 | 0.0001 | {'alpha': 0.0001} | |
| **1** | 1.925756 | 0.064171 | 0.021871 | 0.007647 | 0.01 | {'alpha': 0.01} | |
| **2** | 1.941320 | 0.051871 | 0.023436 | 0.007812 | 0.1 | {'alpha': 0.1} | |
| **4** | 1.961229 | 0.042144 | 0.019278 | 0.006315 | 0.1 | {'alpha': 0.1} | |
| **3** | 2.171880 | 0.099263 | 0.024496 | 0.006506 | 1 | {'alpha': 1} | |

5 rows × 31 columns

**4.1.3.2 Hyperparameter tuning Using L1 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l1',class_weight = 'balanced')
hyper_param_tuning(X_tr_3, y_train)
```


Hyper parameter Vs AUC plot

```
Best score:  0.6832627743718859
value of alpha with best score:  {'alpha': 0.0001}
==========================================================================
Train AUC scores
0    0.698522
1    0.613417
2    0.508756
4    0.508771
3    0.500000
5    0.500000
6    0.500000
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.683263
1    0.611993
2    0.509226
4    0.509080
3    0.500000
5    0.500000
6    0.500000
Name: mean_test_score, dtype: float64
```

Out[106]:

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| 0 | 4.893264 | 0.226318 | 0.022692 | 0.007372 | 0.0001 | {'alpha': 0.0001} | |
| 1 | 4.947641 | 0.206390 | 0.025120 | 0.006791 | 0.01 | {'alpha': 0.01} | |
| 2 | 4.972842 | 0.066120 | 0.017276 | 0.004661 | 0.1 | {'alpha': 0.1} | |
| 4 | 4.895056 | 0.015959 | 0.021988 | 0.006836 | 0.1 | {'alpha': 0.1} | |
| 3 | 4.901596 | 0.014354 | 0.018085 | 0.005076 | 1 | {'alpha': 1} | |

5 rows × 31 columns

### 4.1.3.3 Plotting Roc curve for Set 3

In [107]:

```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.
# https://stackoverflow.com/questions/55250963/how-to-get-probabilities-for-sgdclassifier-L
from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV

svc_bow = SGDClassifier(alpha = 0.0001, penalty= 'l1',class_weight = 'balanced')
model = CalibratedClassifierCV(svc_bow)
model.fit(X_tr_3, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the p
# not the predicted outputs

y_train_pred = prob_predict(model, X_tr_3)
y_test_pred = prob_predict(model, X_te_3)

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("ERROR PLOTS")
plt.grid()
plt.show()
```
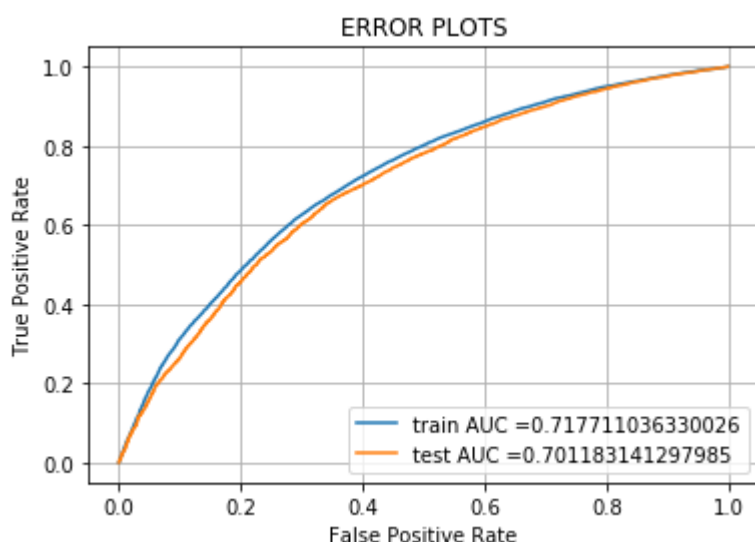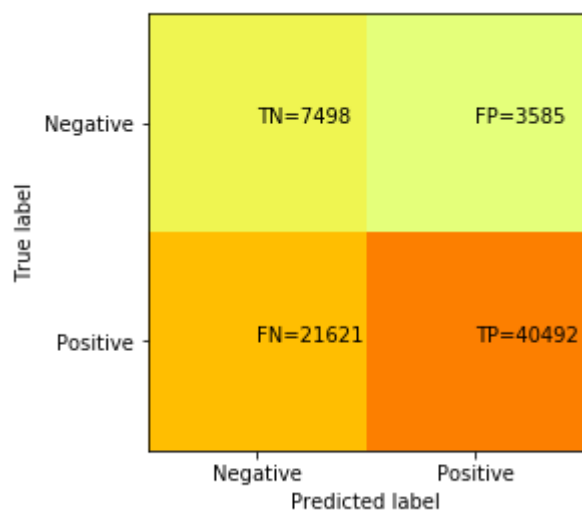


### 4.1.3.4 Confusion Matrix

```python
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

the maximum value of tpr*(1-fpr) 0.4410367968051243 for threshold 0.844

```python
print("TRAIN CONFUSION MATRIX")
get_confusion_matrix(y_train, y_train_pred)
```

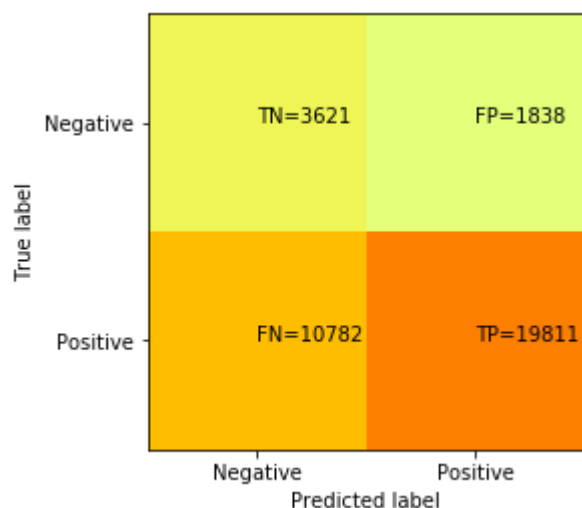TRAIN CONFUSION MATRIX

```python
print("TEST CONFUSION MATRIX")
get_confusion_matrix(y_test, y_test_pred)
```
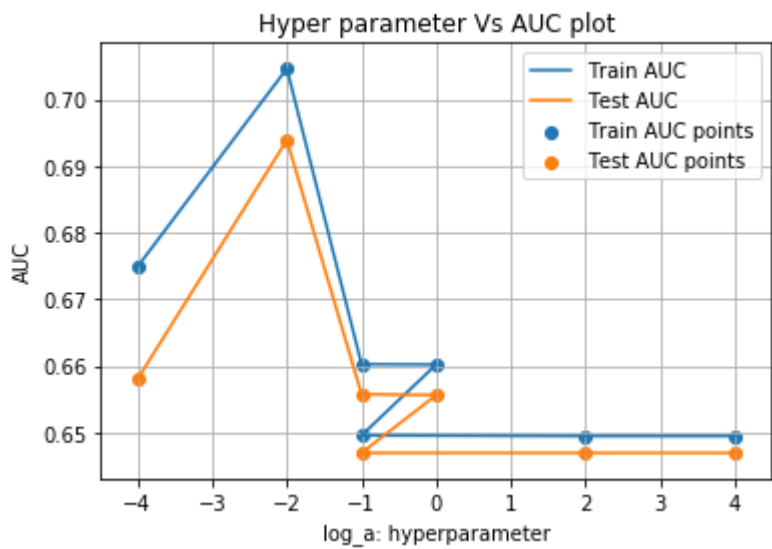
TEST CONFUSION MATRIX



## 4.1.4 Applying SVM on TFIDF W2V (SET 4)

**4.1.4.1 Hyperparameter tuning Using L2 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l2',class_weight = 'balanced')
hyper_param_tuning(X_tr_4, y_train)
```



Hyper parameter Vs AUC plot

```
Best score:  0.6939135664737353
value of alpha with best score:  {'alpha': 0.01}
===============================================================================
Train AUC scores
0    0.674902
1    0.704658
2    0.660267
4    0.660219
3    0.649583
5    0.649481
6    0.649481
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.658016
1    0.693914
2    0.655769
4    0.655636
3    0.646955
5    0.646951
6    0.646951
Name: mean_test_score, dtype: float64
```

Out[111]:

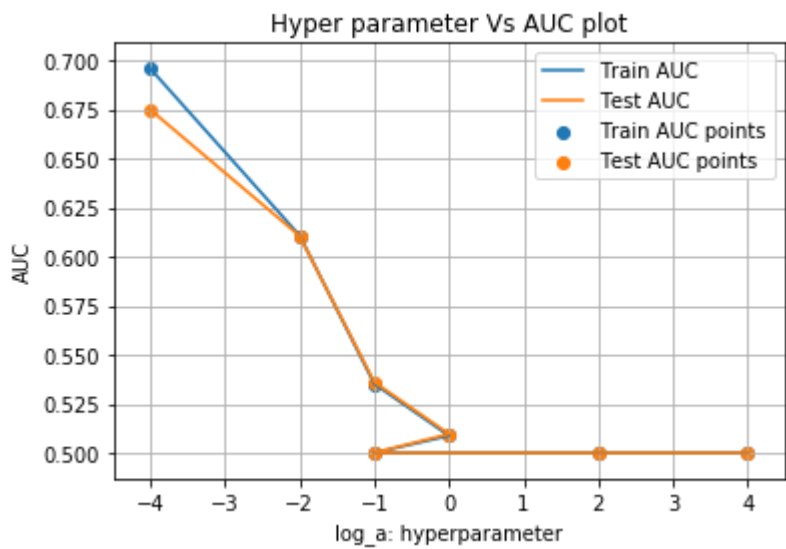| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| 0 | 2.017429 | 0.806404 | 0.021872 | 0.007652 | 0.0001 | {'alpha': 0.0001} | |
| 1 | 1.865801 | 0.026481 | 0.018750 | 0.006248 | 0.01 | {'alpha': 0.01} | |
| 2 | 2.119657 | 0.279874 | 0.021484 | 0.006647 | 0.1 | {'alpha': 0.1} | |
| 4 | 1.966674 | 0.080657 | 0.018975 | 0.002841 | 0.1 | {'alpha': 0.1} | |
| 3 | 1.951554 | 0.019369 | 0.021207 | 0.001001 | 1 | {'alpha': 1} | |

5 rows × 31 columns

**4.1.4.2 Hyperparameter tuning Using L1 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l1',class_weight = 'balanced')
hyper_param_tuning(X_tr_4, y_train)
```



```
Best score:  0.6750075999299053
value of alpha with best score:  {'alpha': 0.0001}
==========================================================================
Train AUC scores
0     0.695921
1     0.610815
2     0.535112
4     0.508714
3     0.500000
5     0.500000
6     0.500000
Name: mean_train_score, dtype: float64
CV AUC scores
0     0.675008
1     0.610616
2     0.536052
4     0.509593
3     0.500000
5     0.500000
6     0.500000
Name: mean_test_score, dtype: float64
```

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split0_t |
|---|---|---|---|---|---|---|---|
| **0** | 5.090693 | 0.457189 | 0.026819 | 0.006153 | 0.0001 | {'alpha': 0.0001} | |
| **1** | 5.078691 | 0.352764 | 0.023551 | 0.006988 | 0.01 | {'alpha': 0.01} | |
| **2** | 5.243652 | 0.267510 | 0.018024 | 0.005063 | 0.1 | {'alpha': 0.1} | |
| **4** | 5.125956 | 0.373796 | 0.019589 | 0.006338 | 0.1 | {'alpha': 0.1} | |
| **3** | 4.925173 | 0.034785 | 0.021875 | 0.007650 | 1 | {'alpha': 1} | |

5 rows × 31 columns

### 4.1.4.3 Plotting Roc curve for Set 4

In [113]:

```python
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.
# https://stackoverflow.com/questions/55250963/how-to-get-probabilities-for-sgdclassifier-l
from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV

svc_bow = SGDClassifier(alpha = 0.0001, penalty= 'l1',class_weight = 'balanced')
model = CalibratedClassifierCV(svc_bow)
model.fit(X_tr_4, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the p
# not the predicted outputs

y_train_pred = prob_predict(model, X_tr_4)
y_test_pred = prob_predict(model, X_te_4)

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("ERROR PLOTS")
plt.grid()
plt.show()
```
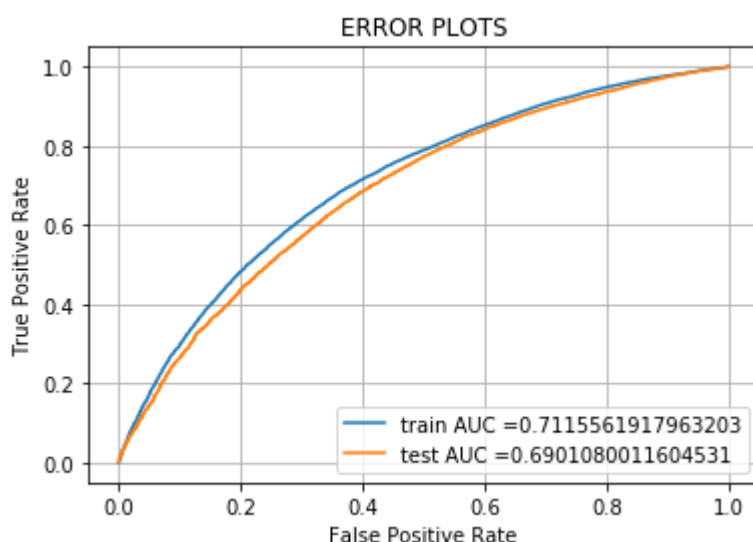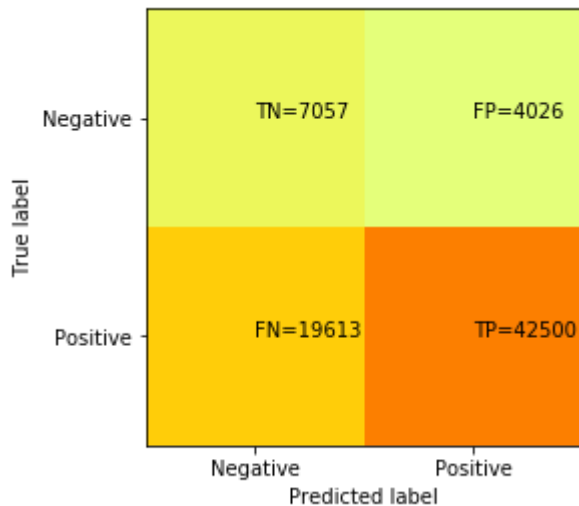


### 4.1.4.4 Confusion Matrix

In [114]:

```python
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

the maximum value of tpr*(1-fpr) 0.43568158954075625 for threshold 0.835

In [115]:

```python
print("TRAIN CONFUSION MATRIX")
get_confusion_matrix(y_train, y_train_pred)
```
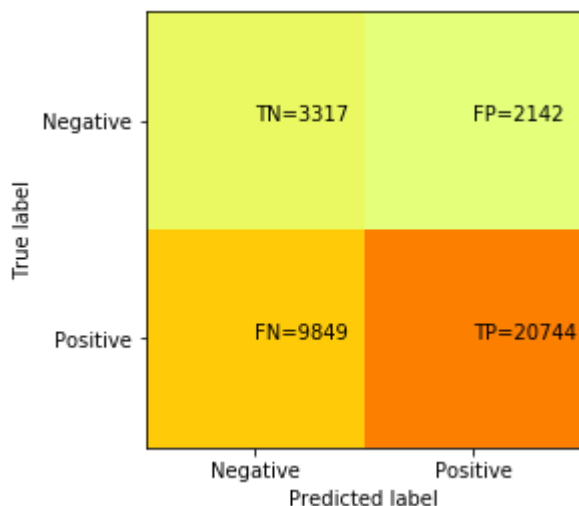
TRAIN CONFUSION MATRIX



In [116]:

```python
print("TEST CONFUSION MATRIX")
get_confusion_matrix(y_test, y_test_pred)
```

TEST CONFUSION MATRIX



## 4.2 Support Vector Machines with added Features `Set 5`

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

**4.2.1 Standardising word count essay**

In [117]:

```
word_count_essay_scalar = StandardScaler()
word_count_essay_scalar.fit((X_train['word_count_essay'].values.astype(float)).reshape(-1,1
print(f"Mean : {word_count_essay_scalar.mean_[0]}, Standard deviation : {np.sqrt(word_count

# Now standardize the data with above maen and variance.
word_count_essay_std_train = word_count_essay_scalar.transform((X_train['word_count_essay']
word_count_essay_std_test = word_count_essay_scalar.transform((X_test['word_count_essay'].v
print(word_count_essay_std_train.shape,y_train.shape)
print(word_count_essay_std_test.shape,y_test.shape)
```

```
Mean : 151.4318678616318, Standard deviation : 38.97341130425965
(73196, 1) (73196,)
(36052, 1) (36052,)
```

**4.2.2 Standardising word count title**

In [118]:

```
word_count_title_scalar = StandardScaler()
word_count_title_scalar.fit((X_train['word_count_title'].values.astype(float)).reshape(-1,1
print(f"Mean : {word_count_title_scalar.mean_[0]}, Standard deviation : {np.sqrt(word_count

# Now standardize the data with above maen and variance.
word_count_title_std_train = word_count_essay_scalar.transform((X_train['word_count_title']
word_count_title_std_test = word_count_essay_scalar.transform((X_test['word_count_title'].v
print(word_count_title_std_train.shape,y_train.shape)
print(word_count_title_std_test.shape,y_test.shape)
```

```
Mean : 4.334362533471775, Standard deviation : 1.7889369058844167
(73196, 1) (73196,)
(36052, 1) (36052,)
```

**4.2.3 Standardising sentiment scores of each essay**

In [119]:

```python
neg_scalar = StandardScaler()
neg_scalar.fit(X_train['neg'].values.reshape(-1,1)) # finding the mean and standard deviati
print(f"Mean : {neg_scalar.mean_[0]}, Standard deviation : {np.sqrt(neg_scalar.var_[0])}")

# Now standardize the data with above maen and variance.
neg_std_train = neg_scalar.transform(X_train['neg'].values.reshape(-1, 1))
neg_std_test =neg_scalar.transform(X_test['neg'].values.reshape(-1, 1))
print(neg_std_train.shape, y_train.shape)
print(neg_std_test.shape, y_test.shape)
```

Mean : 0.04515562325810154, Standard deviation : 0.033939307329286114
(73196, 1) (73196,)
(36052, 1) (36052,)

In [120]:

```python
pos_scalar = StandardScaler()
pos_scalar.fit(X_train['pos'].values.reshape(-1,1)) # finding the mean and standard deviati
print(f"Mean : {pos_scalar.mean_[0]}, Standard deviation : {np.sqrt(pos_scalar.var_[0])}")

# Now standardize the data with above maen and variance.
pos_std_train = pos_scalar.transform(X_train['pos'].values.reshape(-1, 1))
pos_std_test = pos_scalar.transform(X_test['pos'].values.reshape(-1, 1))
print(pos_std_train.shape, y_train.shape)
print(pos_std_test.shape, y_test.shape)
```

Mean : 0.2662285507404776, Standard deviation : 0.07392895289136912
(73196, 1) (73196,)
(36052, 1) (36052,)

In [121]:

```python
neu_scalar = StandardScaler()
neu_scalar.fit(X_train['neu'].values.reshape(-1,1)) # finding the mean and standard deviati
print(f"Mean : {neu_scalar.mean_[0]}, Standard deviation : {np.sqrt(neu_scalar.var_[0])}")

# Now standardize the data with above maen and variance.
neu_std_train = neu_scalar.transform(X_train['pos'].values.reshape(-1, 1))
neu_std_test = neu_scalar.transform(X_test['pos'].values.reshape(-1, 1))
print(neu_std_train.shape, y_train.shape)
print(neu_std_test.shape, y_test.shape)
```

Mean : 0.6886155664243948, Standard deviation : 0.07230568361813375
(73196, 1) (73196,)
(36052, 1) (36052,)

```python
compound_scalar = StandardScaler()
compound_scalar.fit(X_train['compound'].values.reshape(-1,1)) # finding the mean and standa
print(f"Mean : {compound_scalar.mean_[0]}, Standard deviation : {np.sqrt(compound_scalar.va

# Now standardize the data with above maen and variance.
compound_std_train = compound_scalar.transform(X_train['compound'].values.reshape(-1, 1))
compound_std_test = compound_scalar.transform(X_test['compound'].values.reshape(-1, 1))
print(compound_std_train.shape, y_train.shape)
print(compound_std_test.shape, y_test.shape)
```

```
Mean : 0.9586679026176292, Standard deviation : 0.15418529792107544
(73196, 1) (73196,)
(36052, 1) (36052,)
```

```python
print("************Categorical Features************")
print(X_train_categories_ohe.shape)
print(X_test_categories_ohe.shape)
print(X_train_subcategories_ohe.shape)
print(X_test_subcategories_ohe.shape)
print(X_train_state_ohe.shape)
print(X_test_state_ohe.shape)
print(X_train_grade_ohe.shape)
print(X_test_grade_ohe.shape)
print("***********Numerical Features************")
print(quantity_standardized_train.shape)
print(quantity_standardized_test.shape)
print(prev_projects_standardized_train.shape)
print(prev_projects_standardized_train.shape)
print(price_standardized_train.shape)
print(price_standardized_test.shape)
print(compound_std_train.shape)
print(compound_std_test.shape)
print(neu_std_train.shape)
print(neu_std_test.shape)
print(pos_std_train.shape)
print(pos_std_test.shape)
print(neg_std_train.shape)
print(neg_std_test.shape)
print(word_count_title_std_train.shape)
print(word_count_title_std_test.shape)
print(word_count_essay_std_train.shape)
print(word_count_essay_std_test.shape)
```

```
************Categorical Features************
(73196, 9)
(36052, 9)
(73196, 30)
(36052, 30)
(73196, 51)
(36052, 51)
(73196, 4)
(36052, 4)
***********Numerical Features************
(73196, 1)
(36052, 1)
(73196, 1)
(73196, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
(73196, 1)
(36052, 1)
```

```
print("Shape of X_train_essay_tfidf is : ", X_train_essay_tfidf.shape)
```

Shape of X_train_essay_tfidf is :  (73196, 5000)

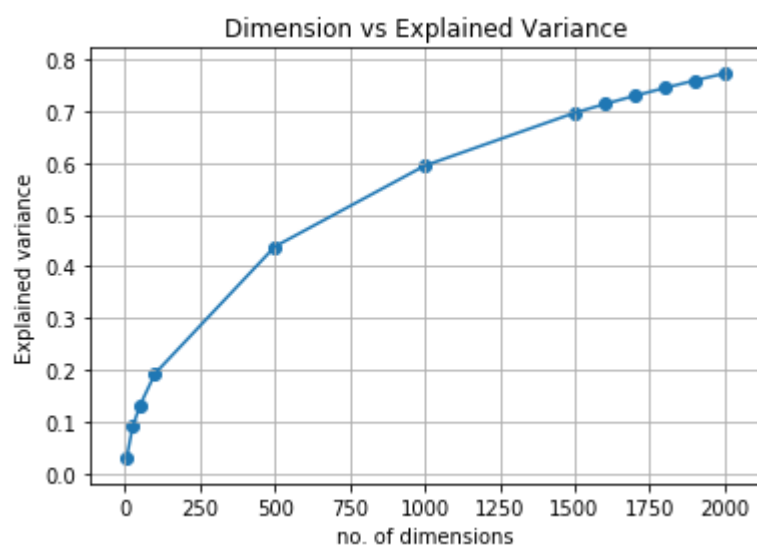**4.2.4 Selecting optimal dimension by elbow method**

In [127]:

```
from sklearn.decomposition import TruncatedSVD
dim = [5,25,50,100,500,1000,1500,1600,1700,1800,1900,2000]
explained_var_sum= []
for d in tqdm(dim , position=0, leave= True):
    svd = TruncatedSVD(n_components = d, n_iter = 5, random_state=42)
    svd.fit(X_train_essay_tfidf)
    explained_var_sum.append(svd.explained_variance_ratio_.sum())
print("explained_var_sum = ",explained_var_sum)
```

```
100%|████████████████████████████████████████████| 12/12 [50:05<00:00, 423.28s/
it]
```

explained_var_sum =  [0.0304031174741989, 0.09040373448772884, 0.13235548413
987375, 0.19332475462812745, 0.43824368846188566, 0.5942773035476773, 0.6962
678141540051, 0.7131064435970107, 0.7290412827072688, 0.7441782031952101, 0.
7586064587749346, 0.7723260251501587]

In [128]:

```
plt.plot(dim, explained_var_sum)
plt.xlabel('no. of dimensions')
plt.ylabel('Explained variance')
plt.scatter(dim, explained_var_sum)
plt.title("Dimension vs Explained Variance")
plt.grid()
plt.show()
#plt.legend()
```

```
optimal_dimension = 2000
svd = TruncatedSVD(n_components = optimal_dimension, n_iter = 5, random_state = 42)
svd.fit(X_train_essay_tfidf)
X_train_essay_tfidf_svd = svd.transform(X_train_essay_tfidf)
X_test_essay_tfidf_svd = svd.transform(X_test_essay_tfidf)
print(X_train_essay_tfidf_svd.shape)
print(X_test_essay_tfidf_svd.shape)
```

```
(73196, 2000)
(36052, 2000)
```

### 4.2.5 Merging all added features

```
# Merging all categorical and numerical features along with tfidf preprocessed essays
X_tr_5 = hstack((X_train_essay_tfidf_svd , X_train_categories_ohe, X_train_subcategories_oh
                X_train_teacher_ohe, X_train_grade_ohe, price_standardized_train, quantity
                prev_projects_standardized_train, compound_std_train, neu_std_train, pos_s
                word_count_title_std_train, word_count_essay_std_train)).tocsr()

#X_cr = hstack((X_cv_essay_bow, X_cv_title_bow, X_cv_resources_bow, X_cv_categories_ohe, X_

X_te_5 = hstack((X_test_essay_tfidf_svd, X_test_categories_ohe, X_test_subcategories_ohe, X
                X_test_grade_ohe, price_standardized_test, quantity_standardized_test, pre
                compound_std_test, neu_std_test, pos_std_test, neg_std_test, word_count_ti

print("Final Data matrix")
print(X_tr_5.shape, y_train.shape)
#print(X_cr.shape, y_cv.shape)
print(X_te_5.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(73196, 2108) (73196,)
(36052, 2108) (36052,)
================================================================================
======================
```
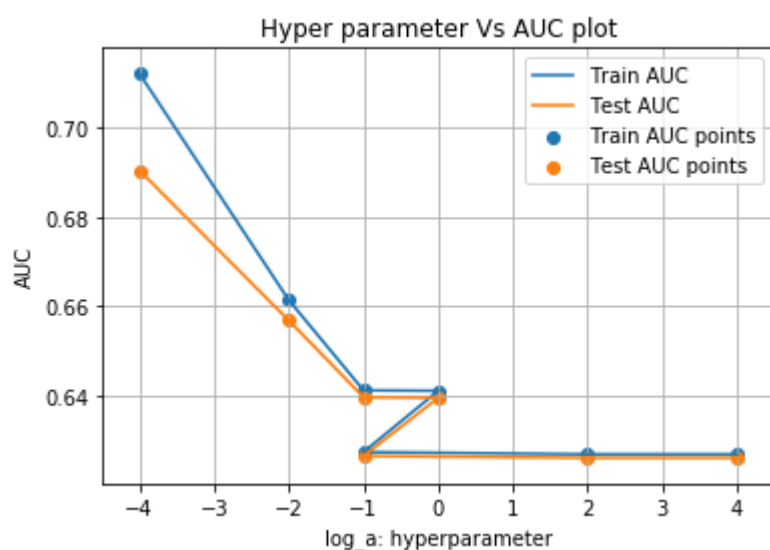
# 4.2.6 Applying SVM on added features (SET 5)

### 4.2.6.1 Hyperparameter tuning Using L2 regulariser

```
svc = SGDClassifier(loss='hinge', penalty='l2',class_weight = 'balanced')
hyper_param_tuning(X_tr_5, y_train)
```



```
Best score:  0.6903835150998499
value of alpha with best score:  {'alpha': 0.0001}
================================================================================
Train AUC scores
0    0.711999
1    0.661413
2    0.641293
4    0.641185
3    0.627509
5    0.627006
6    0.627006
Name: mean_train_score, dtype: float64
CV AUC scores
0    0.690384
1    0.656884
2    0.639755
4    0.639624
3    0.626645
5    0.626210
6    0.626210
Name: mean_test_score, dtype: float64
```

Out[131]:

| mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split |
| --- | --- | --- | --- | --- | --- | --- |

|   | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split |
|---|---|---|---|---|---|---|---|
| 0 | 7.038949 | 4.892151 | 0.226581 | 0.322936 | 0.0001 | {'alpha': 0.0001} | |
| 1 | 5.851091 | 0.130291 | 0.067183 | 0.018553 | 0.01 | {'alpha': 0.01} | |
| 2 | 6.777805 | 0.221484 | 0.092183 | 0.025629 | 0.1 | {'alpha': 0.1} | |
| 4 | 6.525335 | 0.117869 | 0.079680 | 0.008414 | 0.1 | {'alpha': 0.1} | |
| 3 | 6.704057 | 0.107950 | 0.074996 | 0.006246 | 1 | {'alpha': 1} | |

5 rows × 31 columns

**4.2.6.2 Hyperparameter tuning Using L1 regulariser**

```
svc = SGDClassifier(loss='hinge', penalty='l1',class_weight = 'balanced')
hyper_param_tuning(X_tr_5, y_train)
```



Hyper parameter Vs AUC plot

```
Best score:  0.7083341852098831
value of alpha with best score:  {'alpha': 0.0001}
===========================================================================
Train AUC scores
0     0.730139
1     0.636995
2     0.556542
4     0.523950
3     0.500000
5     0.500000
6     0.500000
Name: mean_train_score, dtype: float64
CV AUC scores
0     0.708334
1     0.637119
2     0.557302
4     0.526680
3     0.500000
5     0.500000
6     0.500000
Name: mean_test_score, dtype: float64
```

Out[132]:

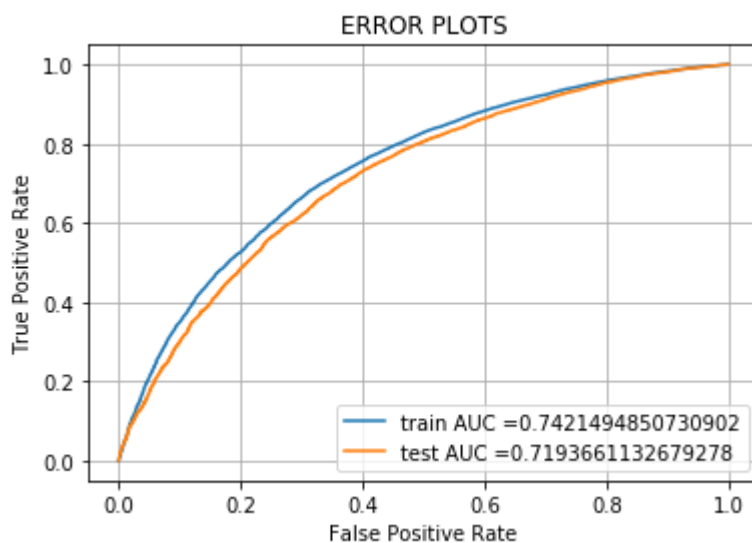| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_alpha | params | split |
|---|---|---|---|---|---|---|---|
| **0** | 16.392638 | 0.267898 | 0.081249 | 0.011702 | 0.0001 | {'alpha': 0.0001} | |
| **1** | 15.693627 | 0.059033 | 0.084376 | 0.017406 | 0.01 | {'alpha': 0.01} | |
| **2** | 16.650610 | 0.141382 | 0.116087 | 0.065646 | 0.1 | {'alpha': 0.1} | |
| **4** | 16.566165 | 0.252153 | 0.118751 | 0.103696 | 0.1 | {'alpha': 0.1} | |
| **3** | 16.828965 | 0.379355 | 0.076560 | 0.004688 | 1 | {'alpha': 1} | |

### 4.2.6.3 Plotting Roc curve for Set 5

In [133]:

```python
svc_bow = SGDClassifier(alpha = 0.0001, penalty= 'l1',class_weight = 'balanced')
model = CalibratedClassifierCV(svc_bow)
model.fit(X_tr_5, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the p
# not the predicted outputs

y_train_pred = prob_predict(model, X_tr_5)
y_test_pred = prob_predict(model, X_te_5)

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("ERROR PLOTS")
plt.grid()
plt.show()
```



### 4.2.6.4 Confusion Matrix

In [134]:

```python
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

the maximum value of tpr*(1-fpr) 0.4678768135216657 for threshold 0.837

```
print("TRAIN CONFUSION MATRIX")
get_confusion_matrix(y_train, y_train_pred)
```
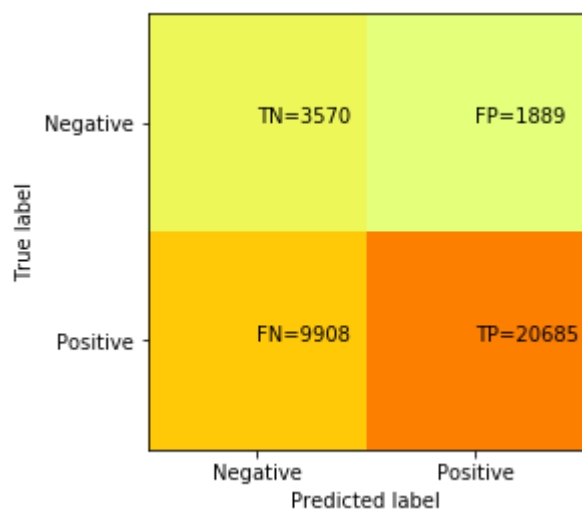
TRAIN CONFUSION MATRIX

```
print("TEST CONFUSION MATRIX")
get_confusion_matrix(y_test, y_test_pred)
```

TEST CONFUSION MATRIX



# 5. Conclusion

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

In [137]:

```python
from prettytable import PrettyTable

x = PrettyTable()
x.field_names = ["Vectorizer","Model","Hyperparameter","Regulariser", "Train AUC", "Test Au
x.add_row(["BOW","SVM", 0.01,"L2", 0.79, 0.72])
x.add_row(["-------------------","-----------","---------------","-----------","----------"
x.add_row(["TFIDF","SVM", 0.0001, "L1", 0.75, 0.71 ])
x.add_row(["-------------------","-----------","---------------","-----------","----------"
x.add_row(["AVG_W2V","SVM", 0.0001, "L1", 0.71, 0.70])
x.add_row(["-------------------","-----------","---------------","-----------","----------"
x.add_row(["TFIDF_W2V","SVM", 0.0001, "L1", 0.71, 0.69])
x.add_row(["-------------------","-----------","---------------","-----------","----------"
x.add_row(["TFIDF TRUNCATED SVD","SVM", 0.0001, "L1", 0.74, 0.71])
print(x)
```

```
+--------------------+-----------+---------------+-----------+-------
------+-----------+
|      Vectorizer    |    Model  | Hyperparameter | Regulariser |  Train
AUC  |   Test Auc  |
+--------------------+-----------+---------------+-----------+-------
------+-----------+
|        BOW         |    SVM    |      0.01     |     L2    |    0.
79   |    0.72     |
|  ----------------- | --------- | ------------- | --------- | ------
----- | ----------- |
|       TFIDF        |    SVM    |     0.0001    |     L1    |    0.
75   |    0.71     |
|  ----------------- | --------- | ------------- | --------- | ------
----- | ----------- |
|      AVG_W2V       |    SVM    |     0.0001    |     L1    |    0.
71   |    0.7      |
|  ----------------- | --------- | ------------- | --------- | ------
----- | ----------- |
|      TFIDF_W2V     |    SVM    |     0.0001    |     L1    |    0.
71   |    0.69     |
|  ----------------- | --------- | ------------- | --------- | ------
----- | ----------- |
| TFIDF TRUNCATED SVD |    SVM    |     0.0001    |     L1    |    0.
74   |    0.71     |
+--------------------+-----------+---------------+-----------+-------
------+-----------+
```