

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*

Feature	Description
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

- __project_essay_1__: "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2__: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

```

1.1 Reading Data

In [2]:

```

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

```

In [3]:

```

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

```

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

```

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

```

In [4]:

```

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

```

Number of data points in train data (1541272, 4)

```
['id' 'description' 'quantity' 'price']
```

Out[4]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00

In [5]:

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

1.2 preprocessing of project_subject_categories

In [6]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [7]:

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

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

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

```

        temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
        sub_cat_list.append(temp.strip())

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

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

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

```

1.3 Text preprocessing

In [8]:

```

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

```

In [9]:

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

Out[9]:

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

In [10]:

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

In [11]:

```

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

```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our s

chool. \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 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\n\nannan

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\n\r\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\n\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs a lot of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager learners and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\n\r\nThey also want to learn through games, my kids don't want to sit it 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 d

disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time. The cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible. nannan

In [12]:

```
# https://stackoverflow.com/a/47091490/4084039
import re

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

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

In [13]:

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

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

In [14]:

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

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

In [15]:

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

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

In [16]:

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

In [17]:

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



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

In [19]:

```
# Updating dataframe for clean project title and remove old project title
project_data['clean_essay'] = preprocessed_essays
project_data.drop(['essay'], axis=1, inplace=True)
project_data.head(2)
```

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

In [20]:

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

In [21]:

```
# after preprocessing
preprocessed_title[20000]
```

Out[21]:

'need move input'

In [22]:

```
# Updating dataframe for clean project title and remove old project title
project_data['clean_project_title'] = preprocessed_title
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_submitted_datetime	project_grade_c
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Gra

Preprocessing project_grade

In [23]:

```
# similarly you can preprocess the project_grade also
# Combining all the above students
from tqdm import tqdm
preprocessed_grade = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_grade_category'].values):
    sent = decontracted(sentence)
    sent = sent.replace(' ', '_')
    sent = sent.replace('-', '_')
    sent = sent.lower()
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_grade.append(sent.strip())
```

100% | 109248/109248 [00:00<00:00, 135395.76it/s]

In [24]:

```
preprocessed_grade[:10]
```

Out[24]:

```
['grades_prek_2',
'grades_6_8',
'grades_6_8',
'grades_prek_2',
'grades_prek_2',
'grades_3_5',
'grades_6_8',
'grades_3_5',
'grades_prek_2',
'grades_prek_2']
```

In [25]:

```
# Updating dataframe for clean project title and remove old project title
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data['project_grade_category'] = preprocessed_grade
project_data.head(2)
```

Out[25]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_essay_1
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	My students are English learners that are work...
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Our students arrive to our school eager to lea...

In [26]:

```
# remove unnecessary column: https://cmdlinetips.com/2018/04/how-to-drop-one-or-more-columns-in-pandas-dataframe/
project_data = project_data.drop(['Unnamed: 0', 'id', 'teacher_id', 'project_submitted_datetime', \
                                  'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4', \
                                  'project_resource_summary'], axis=1)
```

In [27]:

```
project_data.head()
```

Out[27]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_category
0	Mrs.	IN	0	0	154.60	23	Literacy_Language_Arts
1	Mr.	FL	7	1	299.00	1	History_Civics_Health_Sports
2	Ms.	AZ	1	0	516.85	22	Health_Sports
3	Mrs.	KY	4	1	232.90	4	Literacy_Language_Math_Science
4	Mrs.	TX	1	1	67.98	4	Math_Science

Check whether each column contain NaN or Not

In [28]:

```
project_data['teacher_prefix'].isnull().values.any()
```

Out[28]:

True

In [29]:

```
project_data['school_state'].isnull().values.any()
```

Out[29]:

False

In [30]:

```
project_data['teacher_number_of_previously_posted_projects'].isnull().values.any()
```

Out[30]:

False

In [31]:

```
project_data['project_is_approved'].isnull().values.any()
```

Out[31]:

False

In [32]:

```
project_data['price'].isnull().values.any()
```

Out[32]:

False

In [33]:

```
project_data['quantity'].isnull().values.any()
```

Out[33]:

False

In [34]:

```
project_data['clean_categories'].isnull().values.any()
```

Out[34]:

False

In [35]:

```
project_data['clean_subcategories'].isnull().values.any()
```

Out[35]:

False

In [36]:

```
project_data['clean_essay'].isnull().values.any()
```

Out[36]:

False

In [37]:

```
project_data['clean_project_title'].isnull().values.any()
```

Out[37]:

False

In [38]:

```
project_data['project_grade_category'].isnull().values.any()
```

Out[38]:

False

Since we got 'teacher prefix' attributes which contain NaN. Let check how many NaN are contain in this attributes

In [39]:

```
project_data['teacher_prefix'].isnull().sum().sum()
```

Out[39]:

3

1.5 Preparing data for models

In [40]:

```
project_data.columns
```

Out[40]:

```
Index(['teacher_prefix', 'school_state',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'price', 'quantity', 'clean_categories', 'clean_subcategories',  
      'clean_essay', 'clean_project_title', 'project_grade_category'],  
      dtype='object')
```

we are going to consider

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

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

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical

- price : numerical

1.5.1 Vectorizing Categorical data

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

In [0]:

```
# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", categories_one_hot.shape)
```

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

In [0]:

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

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

In [0]:

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

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [0]:

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

```
Shape of matrix after one hot encoding (109248, 16623)
```

In [0]:

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

1.5.2.2 TFIDF vectorizer

In [0]:

```

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

```

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

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

```

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

# =====
Output:

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

# =====

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

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

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

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

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

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

'''

```

Out[0]:

```

'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\nndef loadGloveModel(\ngloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\'r\', encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n        word = splitLine[0]\n

```

In [0]:

In [0]:

In [0]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
```



```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248 [03:22<00  
:00, 540.56it/s]
```

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

In [0]:

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

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

In [0]:

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

Out[0]:

```
(109248, 16663)
```

Computing Sentiment Scores

In [0]:

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

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

sid = SentimentIntensityAnalyzer()

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

D:\installed\Anaconda3\lib\site-packages\nltk\twitter__init__.py:20: UserWarning:

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

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

Assignment 5: Logistic Regression

1. [Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay ('BOW with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay ('TFIDF with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

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

3. Representation of results

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

4. [Task-2] Apply Logistic Regression on the below feature set **Set 5** by finding the best hyper parameter as suggested in step 2 and step 3.

5. 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

And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3.

6. Conclusion

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

Note: Data Leakage

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

2. Logistic Regression

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

In [41]:

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

In [41]:

```
# Combine the train.csv and resource.csv
# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
from sklearn.model_selection import train_test_split

# https://www.geeksforgeeks.org/python-pandas-dataframe-sample/
# Take 50k dataset
project_data = project_data.sample(n=50000)
# Remove that row which contain NaN. We observed that only 3 rows that contain NaN
project_data = project_data[pd.notnull(project_data['teacher_prefix'])]
project_data.shape
```

Out[41]:

(49999, 11)

In [42]:

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

Out[42]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_category
19774	Mr.	MA	7	1	499.99	1	Math_Sci AppliedLea

9183	Ms.	CA	6	1	651.95	4	Math_Sci
------	-----	----	---	---	--------	---	----------

In [127]:

```
# Split train and test
tr_X, ts_X, tr_y, ts_y, = train_test_split(project_data, project_data['project_is_approved'].values, te
```

```

st_size=0.33, random_state=1, stratify=project_data['project_is_approved'].values)
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)

# After train data, We are going to perform KFold Cross validation at the time of training model

# Reset index of df
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)
tr_X.drop(['project_is_approved'], axis=1, inplace=True)
ts_X.drop(['project_is_approved'], axis=1, inplace=True)

print('Shape of train data:', tr_X.shape)
print('Shape of test data:', ts_X.shape)

```

Shape of train data: (33499, 10)
Shape of test data: (16500, 10)

In [128]:

```
tr_X.head(2)
```

Out[128]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	price	quantity	clean_categories	clean_subcategories
0	Ms.	MS	0	336.48	7	Math_Science	AppliedSciences

1	Ms.	CA	25	119.95	12	Music_Arts	VisualArts
---	-----	----	----	--------	----	------------	------------

In [129]:

```
ts_X.head(2)
```

Out[129]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	price	quantity	clean_categories	clean_subcategorie
0	Ms.	CA	13	329.28	37	Literacy_Language	Literac

1	Mrs.	NY	0	369.61	31	Literacy_Language SpecialNeeds	Literature_Writin SpecialNeed
---	------	----	---	--------	----	-----------------------------------	----------------------------------

In [130]:

```

print('Shape of Train Data', [tr_X.shape, tr_y.shape])
print('Shape of Test Data', [ts_X.shape, ts_y.shape])

```

Shape of Train Data [(33499, 10), (33499,)]
Shape of Test Data [(16500, 10), (16500,)]

2.2 Make Data Model Ready: encoding numerical, categorical features

In [131]:

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

# # when you plot any graph make sure you use
#     # a. Title, that describes your plot, this will be very helpful to the reader
#     # b. Legends if needed
#     # c. X-axis label
#     # d. Y-axis label

# # For Numerical with train data
# ### 1) quantity

from sklearn.preprocessing import Normalizer
# # normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.Normalizer.html

quantity_scaler = Normalizer()
quantity_scaler.fit(tr_X['quantity'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
quantity_normalized = quantity_scaler.transform(tr_X['quantity'].values.reshape(1, -1))

# ### 2) price

# # the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1

price_scaler = Normalizer()
price_scaler.fit(tr_X['price'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
price_normalized = price_scaler.transform(tr_X['price'].values.reshape(1, -1))

# ### 3) For teacher_number_of_previously_projects

# # We are going to represent the teacher_number_of_previously_posted_projects, as numerical values within the range 0-1

teacher_number_of_previously_posted_projects_scaler = Normalizer()
teacher_number_of_previously_posted_projects_scaler.fit(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
teacher_number_of_previously_posted_projects_normalized = teacher_number_of_previously_posted_projects_scaler.transform(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
```

In [132]:

```
print('Shape of quantity:', quantity_normalized.T.shape)
print('Shape of price:', price_normalized.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', teacher_number_of_previously_posted_projects_normalized.T.shape)
```

Shape of quantity: (33499, 1)

Shape of price: (33499, 1)

Shape of teacher_number_of_previously_posted_projects: (33499, 1)

In [133]:

```
quantity_normalized.T
```

Out[133]:

```
array([[0.001206 ],
       [0.00206743],
       [0.00103372],
       ...,
       [0.00068914],
       [0.00051686],
       [0.00051686]])
```

In [134]:

```
price_normalized.T
```

Out[134]:

```
array([[0.00389795],
       [0.00138956],
       [0.00068916],
       ...,
       [0.00124846],
       [0.00115266],
       [0.00083408]])
```

In [135]:

```
teacher_number_of_previously_posted_projects_normalized.T
```

Out[135]:

```
array([[0.          ],
       [0.00456236],
       [0.00054748],
       ...,
       [0.00857724],
       [0.00036499],
       [0.          ]])
```

In [136]:

```
# # Transform numerical attributes for test data
ts_price = price_scalar.transform(ts_X['price'].values.reshape(1,-1))
ts_quantity = quantity_scalar.transform(ts_X['quantity'].values.reshape(1,-1))
ts_teacher_number_of_previously_posted_projects = \
teacher_number_of_previously_posted_projects_scalar.transform(ts_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
```

In [137]:

```
print('-----Test data-----')
print('Shape of quantity:', ts_quantity.T.shape)
print('Shape of price:', ts_price.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', ts_teacher_number_of_previously_posted_projects.T.shape)
```

```
-----Test data-----
Shape of quantity: (16500, 1)
Shape of price: (16500, 1)
Shape of teacher_number_of_previously_posted_projects: (16500, 1)
```

In [138]:

```
# For categorical with train data
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# One hot encoding for school state

### 1) school_state
print('=====\\n')
# Count Vectorize with vocabulary contains unique code of school state and we are doing boolean BoW
vectorizer_school_state = CountVectorizer(vocabulary=tr_X['school_state'].unique(), lowercase=False, binary=True)
vectorizer_school_state.fit(tr_X['school_state'].values)
print('List of feature in school_state', vectorizer_school_state.get_feature_names())

# Transform train data
school_state_one_hot = vectorizer_school_state.transform(tr_X['school_state'].values)
print("\\nShape of school_state matrix after one hot encoding ", school_state_one_hot.shape)
```

```

### 2) project_subject_categories
print('=====\\n')
vectorizer_categories = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_categories.fit(tr_X['clean_categories'].values)
print('List of features in project_subject_categories',vectorizer_categories.get_feature_names())

# Transform train data
categories_one_hot = vectorizer_categories.transform(tr_X['clean_categories'].values)
print("\\nShape of project_subject_categories matrix after one hot encoding ",categories_one_hot.shape)

### 3) project_subject_subcategories
print('=====\\n')
vectorizer_subcategories = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_subcategories.fit(tr_X['clean_subcategories'].values)
print('List of features in project_subject_subcategories',vectorizer_subcategories.get_feature_names())

# Transform train data
subcategories_one_hot = vectorizer_subcategories.transform(tr_X['clean_subcategories'].values)
print("\\nShape of project_subject_subcategories matrix after one hot encoding ",subcategories_one_hot.shape)

### 4) project_grade_category
print('=====\\n')
# One hot encoding for project_grade_category

# Count Vectorize with vocabulary contains unique code of project_grade_category and we are doing boolean BoW
vectorizer_grade_category = CountVectorizer(vocabulary=tr_X['project_grade_category'].unique(), lowercase=False, binary=True)
vectorizer_grade_category.fit(tr_X['project_grade_category'].values)
print('List of features in project_grade_category',vectorizer_grade_category.get_feature_names())

# Transform train data
project_grade_category_one_hot = vectorizer_grade_category.transform(tr_X['project_grade_category'].values)
print("\\nShape of project_grade_category matrix after one hot encoding ",project_grade_category_one_hot.shape)

### 5) teacher_prefix
print('=====\\n')
# One hot encoding for teacher_prefix

# Count Vectorize with vocabulary contains unique code of teacher_prefix and we are doing boolean BoW
# Since some of the data is filled with nan. So we update the nan to 'None' as a string
# tr_X['teacher_prefix'] = tr_X['teacher_prefix'].fillna('None')
vectorizer_teacher_prefix = CountVectorizer(vocabulary=tr_X['teacher_prefix'].unique(), lowercase=False, binary=True)
vectorizer_teacher_prefix.fit(tr_X['teacher_prefix'].values)
print('List of features in teacher_prefix',vectorizer_teacher_prefix.get_feature_names())

# Transform train data
teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(tr_X['teacher_prefix'].values)
print("\\nShape of teacher_prefix matrix after one hot encoding ",teacher_prefix_one_hot.shape)

```

List of feature in school_state ['MS', 'CA', 'NC', 'SC', 'FL', 'MN', 'OR', 'MI', 'GA', 'DC', 'NY', 'CO', 'OH', 'WA', 'MA', 'MD', 'TX', 'PA', 'AZ', 'LA', 'CT', 'OK', 'IL', 'NE', 'RI', 'AR', 'VA', 'TN', 'KS', 'NM', 'ID', 'WI', 'MO', 'IN', 'UT', 'AL', 'NV', 'NH', 'HI', 'NJ', 'SD', 'KY', 'AK', 'WV', 'IA', 'ME', 'DE', 'WY', 'ND', 'MT', 'VT']

Shape of school_state matrix after one hot encoding (33499, 51)

List of features in project_subject_categories ['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']

Shape of project_subject_categories matrix after one hot encoding (33499, 9)

List of features in project_subject_subcategories ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other']


```
, 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
```

```
Shape of project_subject_subcategories matrix after one hot encoding (33499, 30)
```

```
List of features in project_grade_category ['grades_3_5', 'grades_prek_2', 'grades_6_8', 'grades_9_12']
```

```
Shape of project_grade_category matrix after one hot encoding (33499, 4)
```

```
List of features in teacher_prefix ['Ms.', 'Mrs.', 'Mr.', 'Teacher', 'Dr.']
```

```
Shape of teacher_prefix matrix after one hot encoding (33499, 5)
```

```
In [139]:
```

```
vectorizer_school_state.get_feature_names()[0], len(vectorizer_school_state.get_feature_names())
```

```
Out[139]:
```

```
('MS', 51)
```

```
In [140]:
```

```
# Transform categorical for test data
ts_school_state = vectorizer_school_state.transform(ts_X['school_state'].values)
ts_project_subject_category = vectorizer_categories.transform(ts_X['clean_categories'].values)
ts_project_subject_subcategory = vectorizer_subcategories.transform(ts_X['clean_subcategories'].values)
ts_project_grade_category = vectorizer_grade_category.transform(ts_X['project_grade_category'].values)
ts_teacher_prefix = vectorizer_teacher_prefix.transform(ts_X['teacher_prefix'].values)
```

```
In [141]:
```

```
print('-----Test data-----')
print('Shape of school_state:', ts_school_state.shape)
print('Shape of project_subject_categories:', ts_project_subject_category.shape)
print('Shape of project_subject_subcategories:', ts_project_subject_subcategory.shape)
print('Shape of project_grade_category:', ts_project_grade_category.shape)
print('Shape of teacher_prefix:', ts_teacher_prefix.shape)
```

```
-----Test data-----
Shape of school_state: (16500, 51)
Shape of project_subject_categories: (16500, 9)
Shape of project_subject_subcategories: (16500, 30)
Shape of project_grade_category: (16500, 4)
Shape of teacher_prefix: (16500, 5)
```

2.3 Make Data Model Ready: encoding eassay, and project_title

```
In [142]:
```

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

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Note:

We already have preprocessed both essay and project_title in Text processing section (1.3 and 1.4) above

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

Apply Logistic Regression on different kind of featurization as mentioned in the instructions

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

In [60]:

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

BoW

In [62]:

```
### BoW in Essay and Title on Train

# # We are considering only the bigram words which appeared in at least 10 documents with max feature =
5000(rows or projects).
vectorizer_bow = CountVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)
tr_essay = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on train",tr_essay.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = CountVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)
tr_title = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig ",tr_title.shape)

### BoW in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on test",ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig on test",ts_title.shape)
```

Shape of essay matrix after one hot encodig on train (33499, 5000)

Shape of title matrix after one hot encodig (33499, 2296)

Shape of essay matrix after one hot encodig on test (16500, 5000)

Shape of title matrix after one hot encodig on test (16500, 2296)

In [63]:

```
print('Shape of normalized essay in train data', tr_essay.shape)
print('Shape of normalized title in train data', tr_title.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

Shape of normalized essay in train data (33499, 5000)

Shape of normalized title in train data (33499, 2296)

Shape of normalized essay in test data (16500, 5000)

```
Shape of normalized essay in test data (16500, 5000),
Shape of normalized title in test data (16500, 2296)
```

TFIDF

In [102]:

```
### BoW in Essay and Title on Train

# # We are considering only the bigram words which appeared in at least 10 documents with max feature =
5000(rows or projects).
vectorizer_bow = TfidfVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)
tr_essay = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on train",tr_essay.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = TfidfVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)
tr_title = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding ",tr_title.shape)

### BoW in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on test",ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding on test",ts_title.shape)
```

```
Shape of essay matrix after one hot encoding on train (33499, 5000)
Shape of title matrix after one hot encoding (33499, 2296)
=====
```

```
Shape of essay matrix after one hot encoding on test (16500, 5000)
Shape of title matrix after one hot encoding on test (16500, 2296)
```

In [103]:

```
print('Shape of normalized essay in train data', tr_essay.shape)
print('Shape of normalized title in train data', tr_title.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

```
Shape of normalized essay in train data (33499, 5000)
Shape of normalized title in train data (33499, 2296)
=====
```

```
Shape of normalized essay in test data (16500, 5000)
Shape of normalized title in test data (16500, 2296)
```

AVG W2V

In [59]:

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

In [60]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_essav'].values): # for each review/sentence
```

```

vector = np.zeros(300) # as word vectors are of zero length
cnt_words = 0; # num of words with a valid vector in the sentence/review
for word in sentence.split(): # for each word in a review/sentence
    if word in glove_words:
        vector += model[word]
        cnt_words += 1
if cnt_words != 0:
    vector /= cnt_words
avg_w2v_essay.append(vector)

avg_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title.append(vector)

tr_essay = np.array(avg_w2v_essay)
tr_title = np.array(avg_w2v_title)
print('===== Train Essay =====')
print(len(avg_w2v_essay))
print(len(avg_w2v_essay[0]))
print('===== Train Title =====')
print(len(avg_w2v_title))
print(len(avg_w2v_title[0]))
# print(avg_w2v_essay[0])

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33499/33499 [00:06<00:
00, 4919.23it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 33499/33499 [00:00<00:0
0, 92276.50it/s]

```

```

===== Train Essay =====
33499
300
===== Train Title =====
33499
300

```

In [61]:

```

# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay.append(vector)

avg_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title.append(vector)

ts_essay = np.array(avg_w2v_essay)

```

```

es_essay = np.array(avg_w2v_essay)
ts_title = np.array(avg_w2v_title)
print('===== Test Essay =====')
print(len(avg_w2v_essay))
print(len(avg_w2v_essay[0]))
print('===== Test Title =====')
print(len(avg_w2v_title))
print(len(avg_w2v_title[0]))
# print(avg_w2v_essay[0])

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:03<00:00, 4701.90it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:00<00:00, 83119.67it/s]

```

```

===== Test Essay =====
16500
300
===== Test Title =====
16500
300

```

TFIDF W2V

In [101]:

```

# Tfidf weighted w2v on essay in train
tfidf_model = TfidfVectorizer()
tfidf_model.fit(tr_X['clean_essay'].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())

# tfidf Word2Vec
# compute average word2vec for each essay
tfidf_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay.append(vector)

tr_essay = np.array(tfidf_w2v_essay)

# compute average word2vec for each essay for test data
tfidf_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay.append(vector)

```

```
ts_essay = np.array(tfidf_w2v_essay)
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 33499/33499 [00:53<00:00, 625.53it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:26<00:00, 625.76it/s]
```

In [102]:

```
# tfidf Word2Vec on title
# compute average word2vec for each title for train data
tfidf_model = TfidfVectorizer()
tfidf_model.fit(tr_X['clean_project_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())

tfidf_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_title.append(vector)

tr_title = np.array(tfidf_w2v_title)

# compute average word2vec for each title for test data
tfidf_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_title.append(vector)

ts_title = np.array(tfidf_w2v_title)
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 33499/33499 [00:00<00:00, 44983.02it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:00<00:00, 45681.29it/s]
```

In [103]:

```
print('Train essay and title shape:',tr_essay.shape,tr_title.shape)
print('Test essay and title shape:',ts_essay.shape,ts_title.shape)
```

```
Train essay and title shape: (33499, 300) (33499, 300)
Test essay and title shape: (16500, 300) (16500, 300)
```

Merge Them

In [104]:

```
quantity_normalized.T.shape, price_normalized.T.shape, teacher_number_of_previously_posted_projects_normalized.T.shape
```

Out[104]:

```
((33499, 1), (33499, 1), (33499, 1))
```

In [105]:

```
school_state_one_hot.shape, categories_one_hot.shape, subcategories_one_hot.shape, project_grade_category_one_hot.shape, \
    teacher_prefix_one_hot.shape
```

Out[105]:

```
((33499, 51), (33499, 9), (33499, 30), (33499, 4), (33499, 5))
```

In [106]:

```
tr_essay.shape, tr_title.shape
```

Out[106]:

```
((33499, 300), (33499, 300))
```

In [107]:

```
# for train data
from scipy.sparse import hstack
tr_X = hstack((quantity_normalized.T, price_normalized.T, teacher_number_of_previously_posted_projects_normalized.T, \
    school_state_one_hot, categories_one_hot, subcategories_one_hot, project_grade_category_one_hot, \
    teacher_prefix_one_hot, tr_essay, tr_title))
tr_X.shape
```

Out[107]:

```
(33499, 702)
```

In [108]:

```
tr_X = tr_X.toarray()
```

In [109]:

```
tr_X
```

Out[109]:

```
array([[ 1.20600340e-03,  3.89795435e-03,  0.00000000e+00, ...,
         2.38950319e-01,  2.69916929e-01,  2.71074947e-01],
       [ 2.06743441e-03,  1.38956141e-03,  4.56235943e-03, ...,
         6.92896427e-02,  2.94690910e-01,  3.29211798e-01],
       [ 1.03371720e-03,  6.89162221e-04,  5.47483132e-04, ...,
         6.22374811e-02,  2.65066111e-01,  1.79199912e-01],
       ...,
       [ 6.89144802e-04,  1.24846214e-03,  8.57723573e-03, ...,
        -2.73954265e-01, -3.48057225e-01,  9.30665295e-02],
       [ 5.16858602e-04,  1.15265828e-03,  3.64988754e-04, ...,
         2.96568898e-01,  2.92006406e-01,  3.65948406e-01],
       [ 5.16858602e-04,  8.34084383e-04,  0.00000000e+00, ...,
         1.92752096e-01, -7.95984087e-03,  3.20504472e-02]])
```

In [110]:

```
tr_X.shape, tr_y.shape
```

Out[110]:

```
((33499, 702), (33499,))
```

In [111]:

```
# for test data
ts_X = hstack((ts_quantity.T, ts_price.T, ts_teacher_number_of_previously_posted_projects.T, ts_school_
state, \
               ts_project_subject_category, ts_project_subject_subcategory, ts_project_grade_category, \
               ts_teacher_prefix, ts_essay, ts_title))
ts_X.shape
```

Out[111]:

```
(16500, 702)
```

In [112]:

```
ts_X = ts_X.toarray()
```

In [113]:

```
ts_X.shape, ts_y.shape
```

Out[113]:

```
((16500, 702), (16500,))
```

In [114]:

```
# check whether data still contain NaN or infinity or not
```

In [115]:

```
np.any(np.isnan(tr_X)), np.any(np.isnan(ts_X))
```

Out[115]:

```
(False, False)
```

In [116]:

```
np.all(np.isfinite(tr_X)), np.all(np.isfinite(ts_X))
```

Out[116]:

```
(True, True)
```

In [117]:

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

In [118]:

```
# we are writing our own function for predict, with defined threshold
```



```

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

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

def plot_cm(feature_names, tr_thresholds, train_fpr, train_tpr, y_train, y_train_pred, y_test, y_test_pred):
    """
    Parameters:
    feature_name - (string) Write feature to print the plot title
    tr_thresholds - train threshold value
    train_fpr = FPR for train data
    train_tpr - TPR for train data
    y_true - test class data
    y_pred - test prediction value

    Return:
    Plot the confusion matrix for Train and Test Data
    """
    best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
    print("Train confusion matrix")
    cm = metrics.confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
    plt.figure(figsize = (10,7))
    sns.heatmap(cm, annot=True, fmt="d")
    plt.xlabel('Predicted Class')
    plt.ylabel('True Class')
    plt.title('Confusion matrix for Train Data when Logistic Regr with {0} features'.format(feature_names))

    print("Test confusion matrix")
    cm = metrics.confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
    plt.figure(figsize = (10,7))
    sns.heatmap(cm, annot=True, fmt="d")
    plt.xlabel('Predicted Class')
    plt.ylabel('True Class')
    plt.title('Confusion matrix for Test Data when Logistic Regr with {0} features'.format(feature_names))

```

In [157]:

```
clf = LogisticRegression(class_weight='balanced', penalty='l1')
```

In [158]:

```
parameters = {'C':[10**-4,10**-3,10**-2,10**-1,1,10,100,10**3,10**4]}
```

-- LR on BoW [SET 1] --

In [81]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.3s
```

```
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 1.7s remaining: 0.0s
```

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.2s
```

```
[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 3.4s remaining: 0.0s
```

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5883290902799887, total= 1.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5880566543664572, total= 1.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5752570200511509, total= 1.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6792375648698786, total= 1.5s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6740486650891329, total= 1.4s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6709267358268193, total= 1.4s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6832449778887497, total= 3.5s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6858548386124801, total= 3.2s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6809568340994321, total= 3.3s
[CV] C=1 .....
[CV] ..... C=1, score=0.6282123583960982, total= 8.6s
[CV] C=1 .....
[CV] ..... C=1, score=0.6236528232395617, total= 5.8s
[CV] C=1 .....
[CV] ..... C=1, score=0.6193442302891617, total= 8.1s
[CV] C=10 .....
[CV] ..... C=10, score=0.6006990931524928, total= 19.4s
[CV] C=10 .....
[CV] ..... C=10, score=0.5896756706945333, total= 29.6s
[CV] C=10 .....
[CV] ..... C=10, score=0.5884590455818532, total= 47.0s
[CV] C=100 .....
[CV] ..... C=100, score=0.5921668899720484, total=11.5min
[CV] C=100 .....
[CV] ..... C=100, score=0.5882192057024227, total= 1.6min
[CV] C=100 .....
[CV] ..... C=100, score=0.5834921458923842, total= 1.3min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.588635078496696, total= 9.9min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.5878175196630935, total= 1.2min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.582240637178308, total= 56.4s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.5878234831220301, total= 8.4min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.5854510936104864, total= 37.4s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.5804486485678351, total= 51.8s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 38.8min finished
```

Out[81]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,
             fit_intercept=True, intercept_scaling=1, max_iter=100,
             multi_class='warn', n_jobs=None, penalty='l1', random_state=None,
             solver='warn', tol=0.0001, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=None,
```

```

best_params=None, max_iter=10000),
param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]},
pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
scoring='roc_auc', verbose=3)

```

In [82]:

```

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

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']
K = parameters['C']

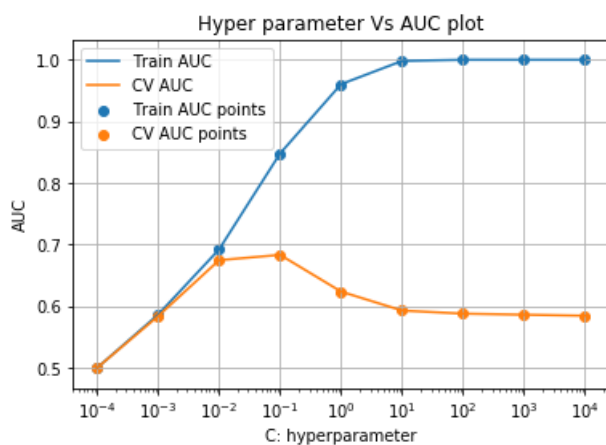
plt.plot(K, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')

plt.plot(K, cv_auc, label='CV 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='darkorange')

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

plt.xscale('log')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

```



In [83]:

```

best_C = clf.best_params_['C']
best_C

```

Out[83]:

0.1

In [84]:

```

lr = LogisticRegression(C=best_C, class_weight='balanced', penalty='l1')
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = np.argmax(lr.predict_proba(tr_X),axis=1)
y_test_pred = np.argmax(lr.predict_proba(ts_X),axis=1)

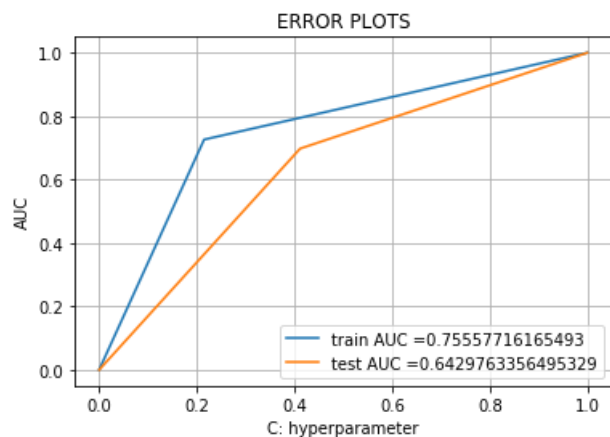
```

```

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [85]:

```

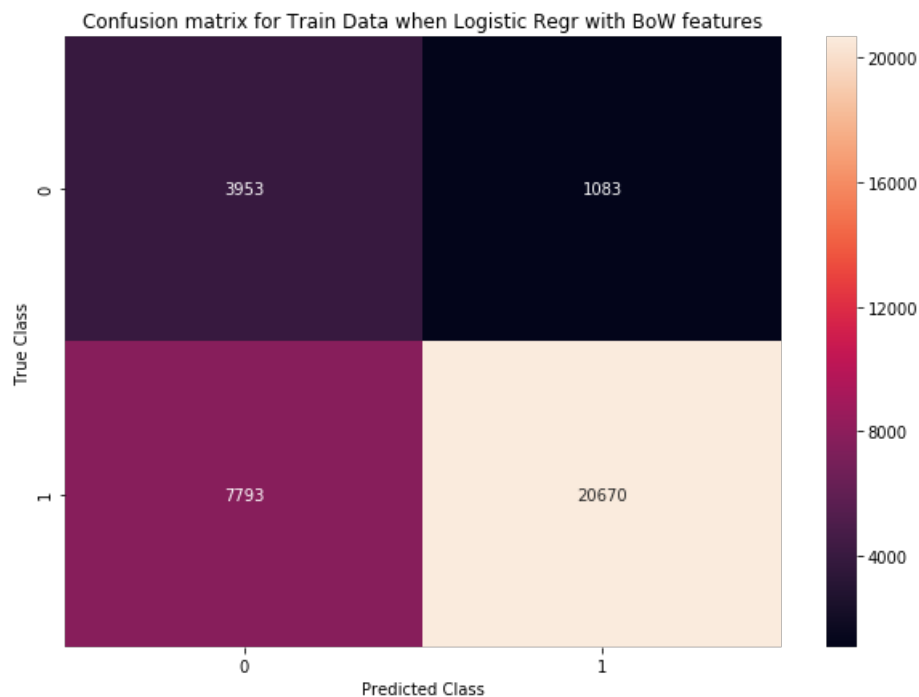
plot_cm('BoW', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)

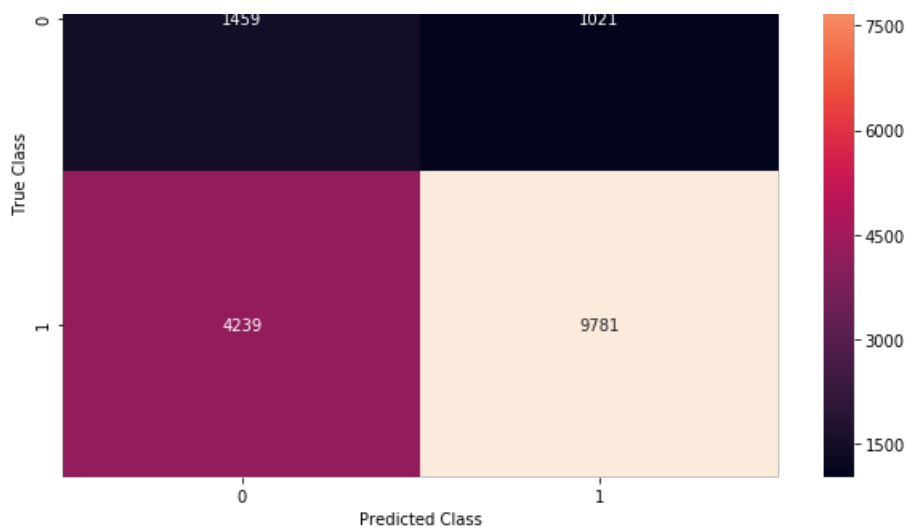
```

the maximum value of $tpr * (1 - fpr)$ 0.5700341792336229 for threshold 1

Train confusion matrix

Test confusion matrix





-- LR on TFIDF [SET 2] --

In [120]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.3s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 1.7s remaining: 0.0s

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.2s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 3.5s remaining: 0.0s

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 1.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 1.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 1.4s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 1.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.549267210165852, total= 1.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.5407977174641213, total= 1.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.5352152113807829, total= 1.3s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6613266925928567, total= 1.5s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6577688302179386, total= 1.4s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6513978478547835, total= 1.5s
[CV] C=1 .....
[CV] ..... C=1, score=0.6861283416712951, total= 6.5s
[CV] C=1 .....
[CV] ..... C=1, score=0.6812253740532537, total= 4.3s
[CV] C=1 .....
[CV] ..... C=1, score=0.6798992109269909, total= 5.6s
```

```
[CV] C=10 .....
[CV] ..... C=10, score=0.6188139471117775, total= 48.0s
[CV] C=10 .....
[CV] ..... C=10, score=0.6057293649255209, total= 43.4s
[CV] C=10 .....
[CV] ..... C=10, score=0.6078100978278662, total= 6.2s
[CV] C=100 .....
[CV] ..... C=100, score=0.5871981359859468, total= 30.5s
[CV] C=100 .....
[CV] ..... C=100, score=0.5818222974608471, total= 1.6min
[CV] C=100 .....
[CV] ..... C=100, score=0.5813352516893765, total= 45.8s
[CV] C=1000 .....
[CV] ..... C=1000, score=0.5838612354579484, total= 57.8s
[CV] C=1000 .....
[CV] ..... C=1000, score=0.5807882336812143, total= 1.2min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.5800975627773933, total= 1.3min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.5817457140934488, total=10.3min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.580383220659531, total= 32.1s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.5785917697048077, total= 21.0s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 20.0min finished
```

Out[120]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,
             fit_intercept=True, intercept_scaling=1, max_iter=100,
             multi_class='warn', n_jobs=None, penalty='l1', random_state=None,
             solver='warn', tol=0.0001, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

In [121]:

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

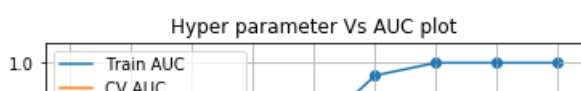
train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
K = parameters['C']

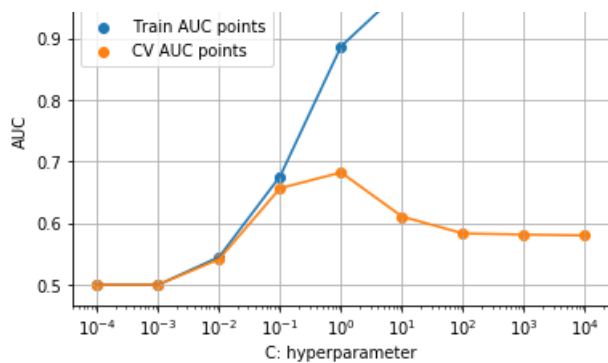
plt.plot(K, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# plt.gca().fill_between(K, train_auc - train_auc_std, train_auc + train_auc_std, alpha=0.2, color='darkblue')

plt.plot(K, cv_auc, label='CV 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='darkorange')

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

plt.xscale('log')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()
```





In [122]:

```
best_C = clf.best_params_['C']
best_C
```

Out[122]:

1

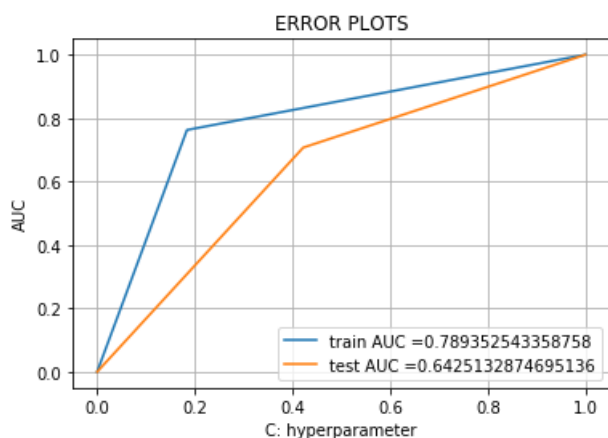
In [123]:

```
lr = LogisticRegression(C=best_C, class_weight='balanced', penalty='l1')
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = np.argmax(lr.predict_proba(tr_X), axis=1)
y_test_pred = np.argmax(lr.predict_proba(ts_X), axis=1)

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, y_test_pred)

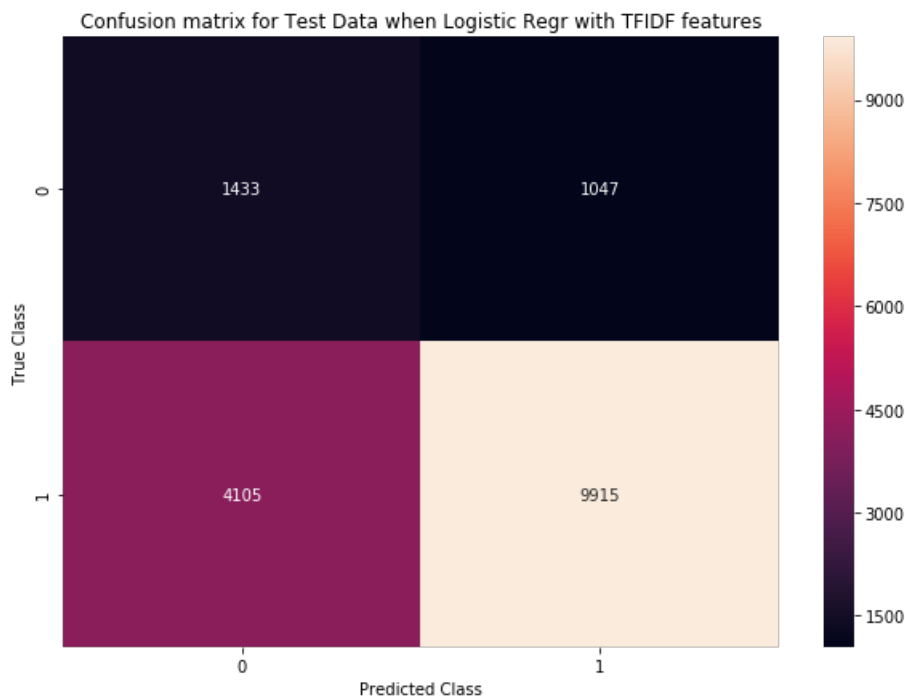
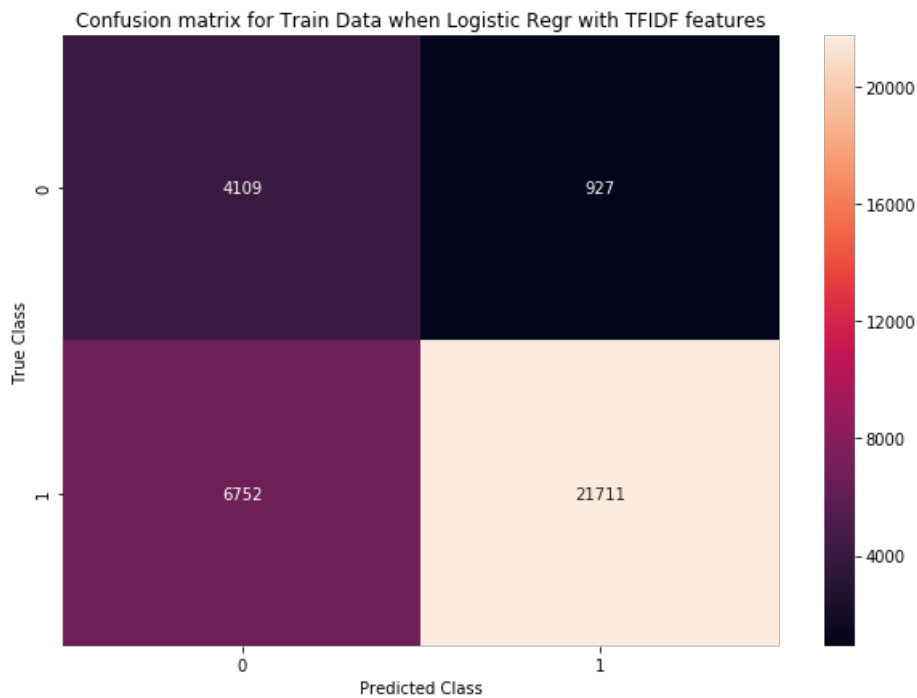
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [124]:

```
plot_cm('TFIDF', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.6223713243147737 for threshold 1
 Train confusion matrix
 Test confusion matrix



-- LR on AVG W2V [SET 3] --

In [79]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 0.3s
```



```
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.4s remaining: 0.0s
```

```
[CV] C=0.0001 .....  
[CV] ..... C=0.0001, score=0.5, total= 0.3s
```

```
[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.8s remaining: 0.0s
```

```
[CV] C=0.0001 .....  
[CV] ..... C=0.0001, score=0.5, total= 0.3s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, score=0.5, total= 0.3s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, score=0.5, total= 0.3s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, score=0.5, total= 0.3s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, score=0.5791924068085009, total= 1.8s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, score=0.5756955456340351, total= 0.9s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, score=0.5768074838288566, total= 0.8s  
[CV] C=0.1 .....  
[CV] ..... C=0.1, score=0.6832546807019675, total= 1.7min  
[CV] C=0.1 .....  
[CV] ..... C=0.1, score=0.6720099814462083, total= 1.6min  
[CV] C=0.1 .....  
[CV] ..... C=0.1, score=0.6719858883841425, total= 1.5min  
[CV] C=1 .....  
[CV] ..... C=1, score=0.69215290451349, total= 8.1min  
[CV] C=1 .....  
[CV] ..... C=1, score=0.694662635600618, total=14.1min  
[CV] C=1 .....  
[CV] ..... C=1, score=0.6883716504803199, total= 7.6min  
[CV] C=10 .....  
[CV] ..... C=10, score=0.6896474146547735, total=29.8min  
[CV] C=10 .....  
[CV] ..... C=10, score=0.6951233469495756, total=29.3min  
[CV] C=10 .....  
[CV] ..... C=10, score=0.688644548028955, total=27.7min  
[CV] C=100 .....  
[CV] ..... C=100, score=0.6890092680336082, total=34.9min  
[CV] C=100 .....  
[CV] ..... C=100, score=0.6947606902694192, total=29.3min  
[CV] C=100 .....  
[CV] ..... C=100, score=0.6882282335698953, total=10.0min  
[CV] C=1000 .....  
[CV] ..... C=1000, score=0.6889698942414879, total=24.8min  
[CV] C=1000 .....  
[CV] ..... C=1000, score=0.6947185944948855, total=31.9min  
[CV] C=1000 .....  
[CV] ..... C=1000, score=0.6881732486413401, total= 6.2min  
[CV] C=10000 .....  
[CV] ..... C=10000, score=0.6889507770626772, total=38.2min  
[CV] C=10000 .....  
[CV] ..... C=10000, score=0.6947127074166425, total=35.7min  
[CV] C=10000 .....  
[CV] ..... C=10000, score=0.6881648868780115, total= 5.3min
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 337.8min finished
```

Out[79]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',  
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,  
             fit_intercept=True, intercept_scaling=1, max_iter=100,  
             multi_class='warn', n_jobs=None, penalty='l1', random_state=None,  
             solver='warn', tol=0.0001, verbose=0, warm_start=False),  
             fit_params=None, iid='warn', n_jobs=None,  
             param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]},  
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',  
             scoring='roc_auc', verbose=3)
```

In [80]:

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

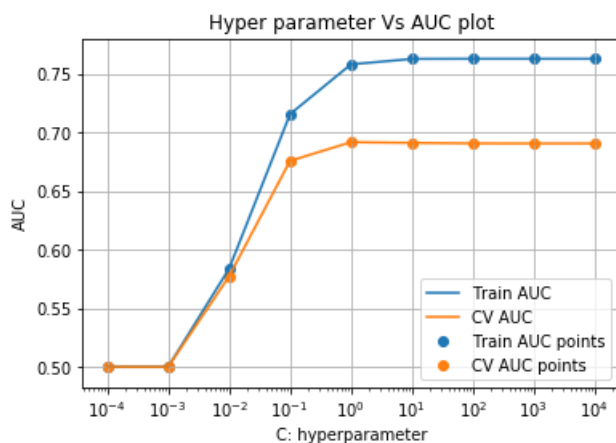
train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
K = parameters['C']

plt.plot(K, 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=0.2,color='darkblue')

plt.plot(K, cv_auc, label='CV 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='darkorange')

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

plt.xscale('log')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()
```



In [81]:

```
best_C = clf.best_params_['C']
best_C
```

Out[81]:

1

In [82]:

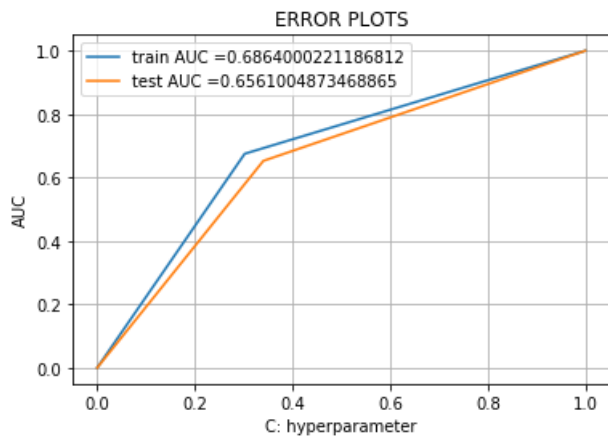
```
lr = LogisticRegression(C=best_C, class_weight='balanced', penalty='l1')
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = np.argmax(lr.predict_proba(tr_X),axis=1)
y_test_pred = np.argmax(lr.predict_proba(ts_X),axis=1)

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, y_test_pred)

plt.plot(train_fpr, train_tpr, label="train AUC = "+str(auc(train_fpr, train_tpr)))
```

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



In [83]:

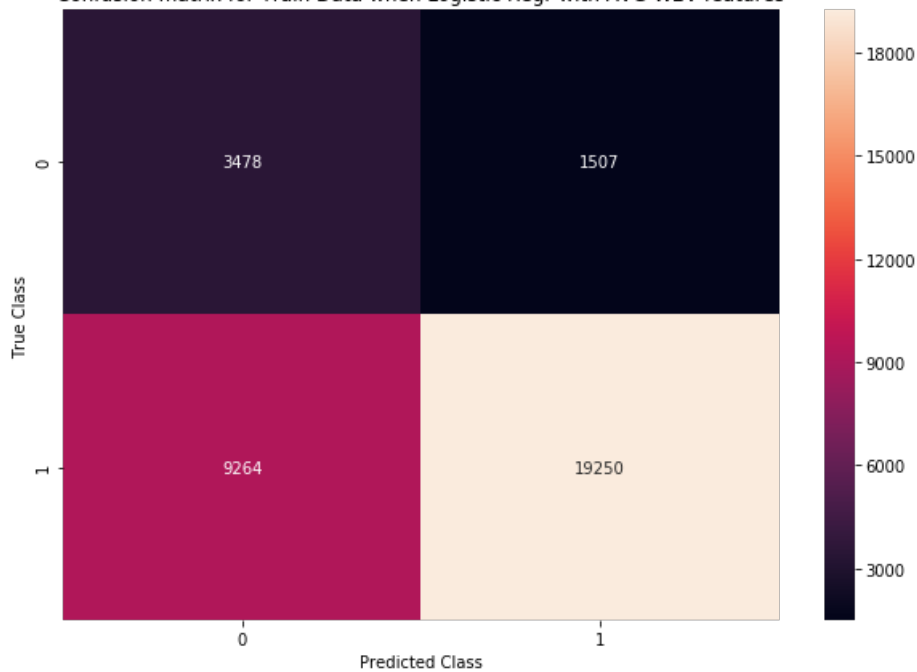
```
plot_cm('AVG W2V', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr * (1 - fpr)$ 0.47101745722543237 for threshold 1

Train confusion matrix

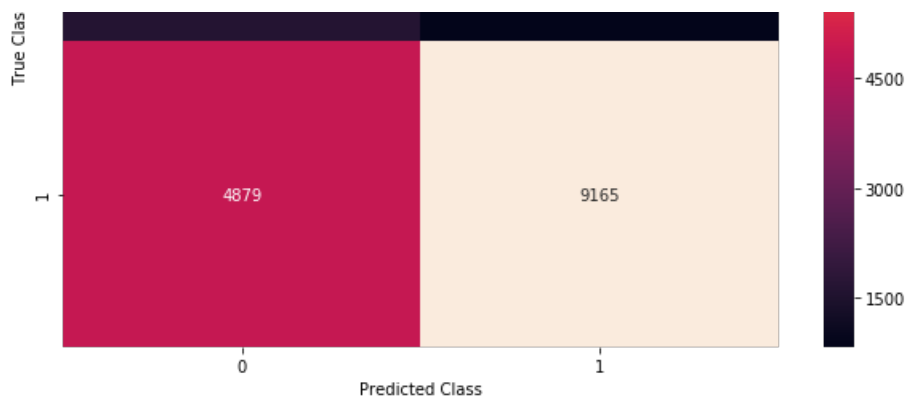
Test confusion matrix

Confusion matrix for Train Data when Logistic Regr with AVG W2V features



Confusion matrix for Test Data when Logistic Regr with AVG W2V features





-- SET 4 --

In [121]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 0.3s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.4s remaining: 0.0s

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 0.3s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.8s remaining: 0.0s

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5, total= 0.3s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 0.3s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 0.3s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5, total= 0.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6037702621522273, total= 1.4s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6137518982662238, total= 1.5s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.6105503661945564, total= 1.9s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6865800569843853, total= 42.1s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.6853857397240416, total= 48.7s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.675688882604897, total= 42.8s
[CV] C=1 .....
[CV] ..... C=1, score=0.6847149293139149, total= 5.7min
[CV] C=1 .....
[CV] ..... C=1, score=0.6904104559573749, total= 5.9min
[CV] C=1 .....
[CV] ..... C=1, score=0.6762809207872422, total=11.9min
[CV] C=10 .....
[CV] ..... C=10, score=0.6828909478892293, total=21.4min
[CV] C=10 .....
[CV] ..... C=10, score=0.6901128103455588, total=15.2min
[CV] C=10 .....
```

```
[CV] ..... C=10, score=0.6759443598132641, total=23.6min
[CV] C=100 .....
[CV] ..... C=100, score=0.6825186060158344, total=21.2min
[CV] C=100 .....
[CV] ..... C=100, score=0.6898111767129973, total=15.0min
[CV] C=100 .....
[CV] ..... C=100, score=0.6756626570744572, total=27.3min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6824758772221347, total=24.4min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6897852862291112, total=17.2min
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6756618335674627, total=19.3min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6824751809010521, total=24.5min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6897793358489515, total=17.7min
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6756351646101797, total=30.8min
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 283.5min finished
```

Out[121]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,
             fit_intercept=True, intercept_scaling=1, max_iter=100,
             multi_class='warn', n_jobs=None, penalty='l1', random_state=None,
             solver='warn', tol=0.0001, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

In [122]:

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

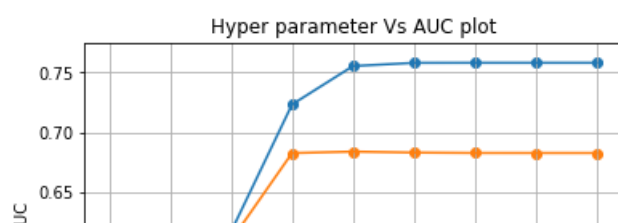
train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
K = parameters['C']

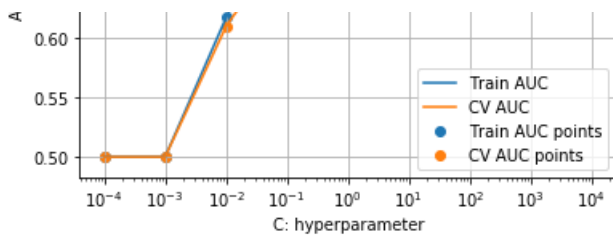
plt.plot(K, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')

plt.plot(K, cv_auc, label='CV 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='darkorange')

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

plt.xscale('log')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()
```





In [123]:

```
best_C = clf.best_params_['C']
best_C
```

Out[123]:

1

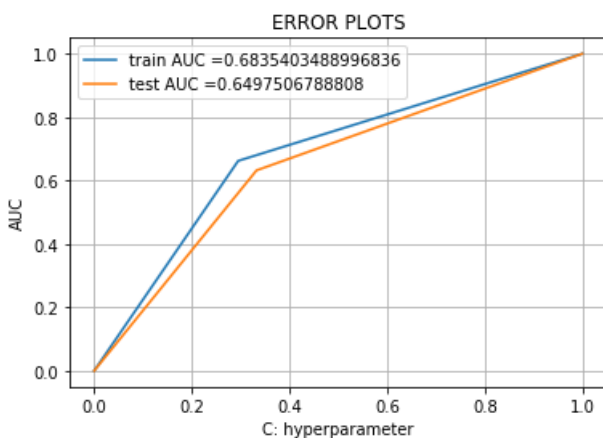
In [124]:

```
lr = LogisticRegression(C=best_C, class_weight='balanced', penalty='l1')
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = np.argmax(lr.predict_proba(tr_X), axis=1)
y_test_pred = np.argmax(lr.predict_proba(ts_X), axis=1)

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [125]:

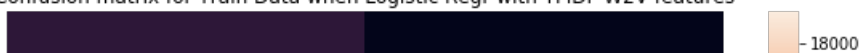
```
plot_cm('TFIDF W2V', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

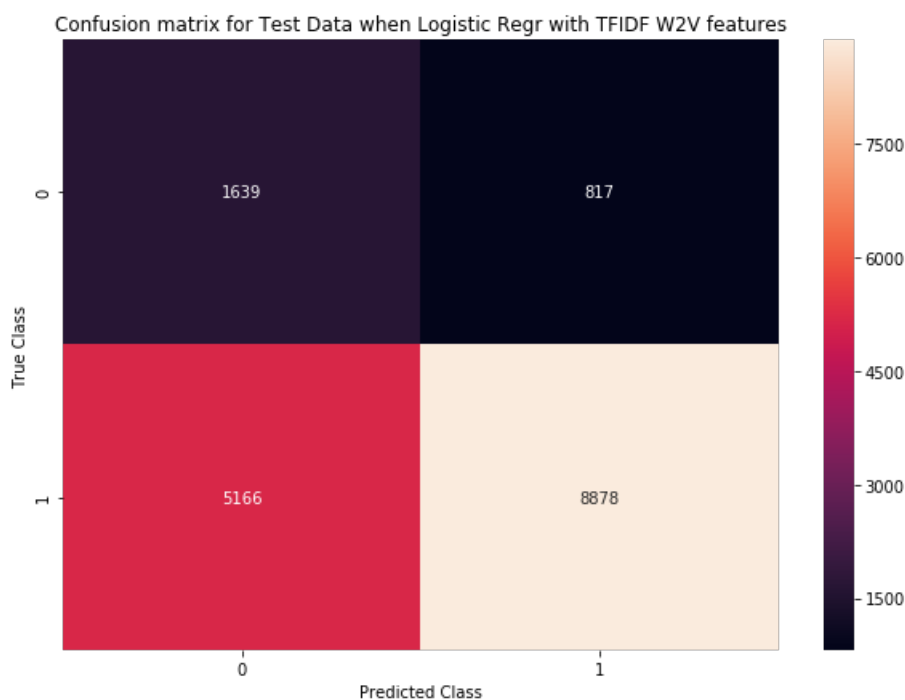
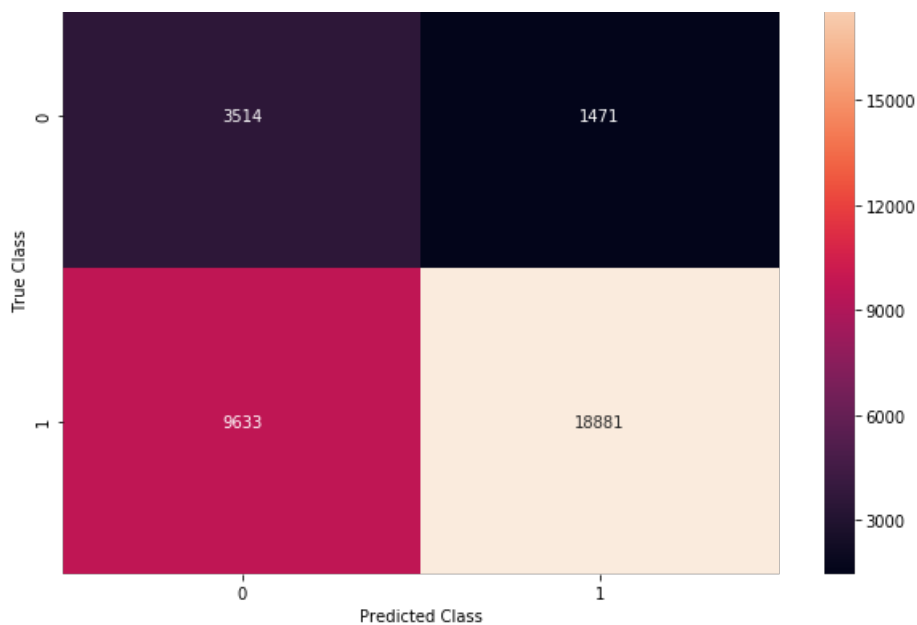
the maximum value of $tpr \cdot (1 - fpr)$ 0.46677054379804916 for threshold 1

Train confusion matrix

Test confusion matrix

Confusion matrix for Train Data when Logistic Regr with TFIDF W2V features





2.5 Logistic Regression with added Features `Set 5`

In [0]:

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

Count number of words in essay and title

In [143]:

```
tr_essay = []
# For each row in the training data, get the length of each essay
```

```
# To calculate number of words, just take the length of each essay
for i in range(tr_X.shape[0]):
    tr_essay.append(len(tr_X['clean_essay'][i]))
```

In [144]:

```
tr_essay = np.array(tr_essay).reshape(-1,1)
tr_essay.shape
```

Out[144]:

```
(33499, 1)
```

In [145]:

```
tr_title = []
# To calculate number of words, just take the length of each title
for i in range(tr_X.shape[0]):
    tr_title.append(len(tr_X['clean_project_title'][i]))
```

In [146]:

```
tr_title = np.array(tr_title).reshape(-1,1)
tr_title.shape
```

Out[146]:

```
(33499, 1)
```

In [147]:

```
ts_essay = []
# To calculate number of words, just take the length of each essay
for i in range(ts_X.shape[0]):
    ts_essay.append(len(ts_X['clean_essay'][i]))

ts_title = []
# To calculate number of words, just take the length of each title
for i in range(ts_X.shape[0]):
    ts_title.append(len(ts_X['clean_project_title'][i]))

ts_essay = np.array(ts_essay).reshape(-1,1)
ts_title = np.array(ts_title).reshape(-1,1)
```

In [148]:

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

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

sid = SentimentIntensityAnalyzer()
tr_sen_essay = []

for i in tqdm(range(tr_X.shape[0])):
    ss = sid.polarity_scores(tr_X['clean_essay'][i])
    tr_sen_essay.append([ss['neg'],ss['neu'],ss['pos'],ss['compound']])

tr_sen_title = []

for i in tqdm(range(tr_X.shape[0])):
    ss = sid.polarity_scores(tr_X['clean_project_title'][i])
    tr_sen_title.append([ss['neg'],ss['neu'],ss['pos'],ss['compound']])

tr_sen_essay = np.array(tr_sen_essay)
tr_sen_title = np.array(tr_sen_title)
# 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
```



```
100%|██████████████████████████████████████████████████████████| 33499/33499 [00:50<00  
:00, 659.72it/s]  
100%|██████████████████████████████████████████████████████████| 33499/33499 [00:02<00:0  
0, 14243.58it/s]
```

In [149]:

| tr_sen_essay |

Out[149]:

```
array([[ 0.202 ,  0.596 ,  0.202 , -0.25  ],
       [ 0.066 ,  0.707 ,  0.226 ,  0.9774],
       [ 0.062 ,  0.592 ,  0.346 ,  0.9923],
       ...,
       [ 0.066 ,  0.704 ,  0.23  ,  0.9621],
       [ 0.134 ,  0.663 ,  0.203 ,  0.8316],
       [ 0.012 ,  0.809 ,  0.179 ,  0.9812]])
```

In [150]:

| sen | title |

Out[150]:

```
array([[ 0.      ,  1.      ,  0.      ,  0.      ],
       [ 0.      ,  1.      ,  0.      ,  0.      ],
       [ 0.      ,  1.      ,  0.      ,  0.      ],
       ...,
       [ 0.444 ,  0.556 ,  0.      , -0.34  ],
       [ 0.      ,  1.      ,  0.      ,  0.      ],
       [ 0.      ,  0.645 ,  0.355 ,  0.5106]])
```

In [151]:

```
ts_sen_essay = []

for i in tqdm(range(ts_X.shape[0])):
    ss = sid.polarity_scores(ts_X['clean_essay'][i])
    ts_sen_essay.append([ss['neg'], ss['neu'], ss['pos'], ss['compound']])

ts_sen_title = []

for i in tqdm(range(ts_X.shape[0])):
    ss = sid.polarity_scores(ts_X['clean_project_title'][i])
    ts_sen_title.append([ss['neg'], ss['neu'], ss['pos'], ss['compound']])

ts_sen_essay = np.array(ts_sen_essay)
ts_sen_title = np.array(ts_sen_title)
```

```
100%|██████████████████████████████████████████████████████████████| 16500/16500 [00:25<00  
:00, 651.32it/s]  
100%|██████████████████████████████████████████████████████████████| 16500/16500 [00:01<00:0  
0, 13851.99it/s]
```

In [152]:

ts sen essay

Out[152]:

```
array([[0.024 , 0.714 , 0.261 , 0.9897],
       [0.      , 0.61  , 0.39  , 0.9937],
       [0.054 , 0.77  , 0.176 , 0.9296],
       ...,
       [0.      , 0.78  , 0.22  , 0.9657],
       [0.046 , 0.612 , 0.342 , 0.9899],
       [0.023 , 0.654 , 0.324 , 0.9858]])
```

```
[0.023 , 0.634 , 0.324 , 0.9998]])
```

In [153]:

```
ts_sen_title
```

Out[153]:

```
array([[0.      , 1.      , 0.      , 0.      ],
       [0.      , 0.612 , 0.388 , 0.2263],
       [0.      , 0.597 , 0.403 , 0.4019],
       ...,
       [0.      , 1.      , 0.      , 0.      ],
       [0.      , 1.      , 0.      , 0.      ],
       [0.      , 1.      , 0.      , 0.      ]])
```

Merge them

In [154]:

```
# for train data
from scipy.sparse import hstack
tr_X = hstack((quantity_normalized.T, price_normalized.T, teacher_number_of_previously_posted_projects_
normalized.T, \
               school_state_one_hot, categories_one_hot, subcategories_one_hot, project_grade_category_o
ne_hot, \
               teacher_prefix_one_hot, tr_essay, tr_title, tr_sen_essay, tr_sen_title))
tr_X.shape
```

Out[154]:

```
(33499, 112)
```

In [155]:

```
# for test data
ts_X = hstack((ts_quantity.T, ts_price.T, ts_teacher_number_of_previously_posted_projects.T, ts_school_
state, \
               ts_project_subject_category, ts_project_subject_subcategory, ts_project_grade_category, \
               ts_teacher_prefix, ts_essay, ts_title, ts_sen_essay, ts_sen_title))
ts_X.shape
```

Out[155]:

```
(16500, 112)
```

In [159]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
```

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5601513802033384, total= 0.0s
```

```
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.0s remaining: 0.0s
```

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.565986455921926, total= 0.0s
```

```
[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.0s remaining: 0.0s
```

```
[CV] C=0.0001 .....
[CV] ..... C=0.0001, score=0.5694816289525801, total= 0.0s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5254810470896627, total= 0.0s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5395141957713054, total= 0.0s
[CV] C=0.001 .....
[CV] ..... C=0.001, score=0.5456391693880405, total= 0.0s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.573157265382524, total= 0.0s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.5822482751810276, total= 0.0s
[CV] C=0.01 .....
[CV] ..... C=0.01, score=0.5846514512980496, total= 0.1s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.5922153202032499, total= 0.4s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.5969023207115641, total= 0.2s
[CV] C=0.1 .....
[CV] ..... C=0.1, score=0.5984552655797388, total= 0.2s
[CV] C=1 .....
[CV] ..... C=1, score=0.6234328502764078, total= 0.9s
[CV] C=1 .....
[CV] ..... C=1, score=0.6255459948560863, total= 4.3s
[CV] C=1 .....
[CV] ..... C=1, score=0.634189134471344, total= 0.8s
[CV] C=10 .....
[CV] ..... C=10, score=0.6273669377887755, total= 1.2s
[CV] C=10 .....
[CV] ..... C=10, score=0.6291576224053336, total= 2.4s
[CV] C=10 .....
[CV] ..... C=10, score=0.6397706114932185, total= 2.7s
[CV] C=100 .....
[CV] ..... C=100, score=0.6278210657384073, total= 2.9s
[CV] C=100 .....
[CV] ..... C=100, score=0.6294515332040709, total= 2.7s
[CV] C=100 .....
[CV] ..... C=100, score=0.6401839486577596, total= 3.0s
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6277213652197747, total= 1.3s
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6294913501096073, total= 4.3s
[CV] C=1000 .....
[CV] ..... C=1000, score=0.6400957067159656, total= 1.3s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6279072196468892, total= 5.5s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6294979968108494, total= 2.9s
[CV] C=10000 .....
[CV] ..... C=10000, score=0.6402110610418859, total= 2.3s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 42.1s finished
```

Out[159]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,
             fit_intercept=True, intercept_scaling=1, max_iter=100,
             multi_class='warn', n_jobs=None, penalty='l1', random_state=None,
             solver='warn', tol=0.0001, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]}},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

In [160]:

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

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']
K = parameters['C']

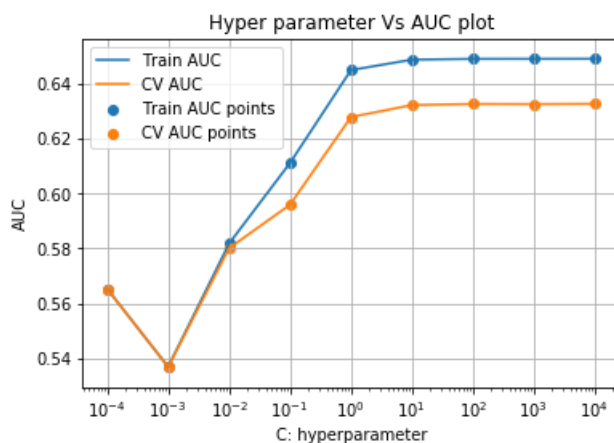
plt.plot(K, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# plt.gca().fill_between(K, train_auc - train_auc_std, train_auc + train_auc_std, alpha=0.2, color='darkblue')

plt.plot(K, cv_auc, label='CV 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='darkorange')

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

plt.xscale('log')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

```



In [161]:

```

best_C = clf.best_params_['C']
best_C

```

Out[161]:

10000

In [162]:

```

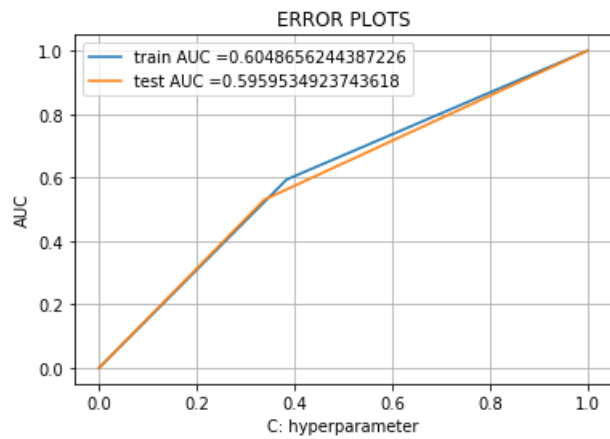
lr = LogisticRegression(C=best_C, class_weight='balanced', penalty='l1')
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = np.argmax(lr.predict_proba(tr_X), axis=1)
y_test_pred = np.argmax(lr.predict_proba(ts_X), axis=1)

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [163]:

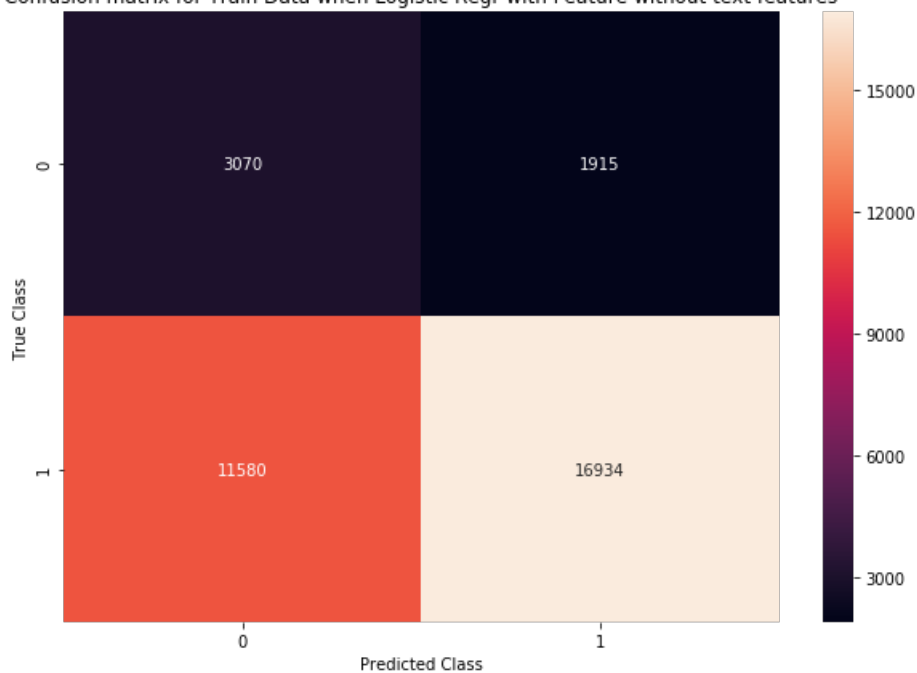
```
plot_cm('Feature without text', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.36574182110053244 for threshold 1

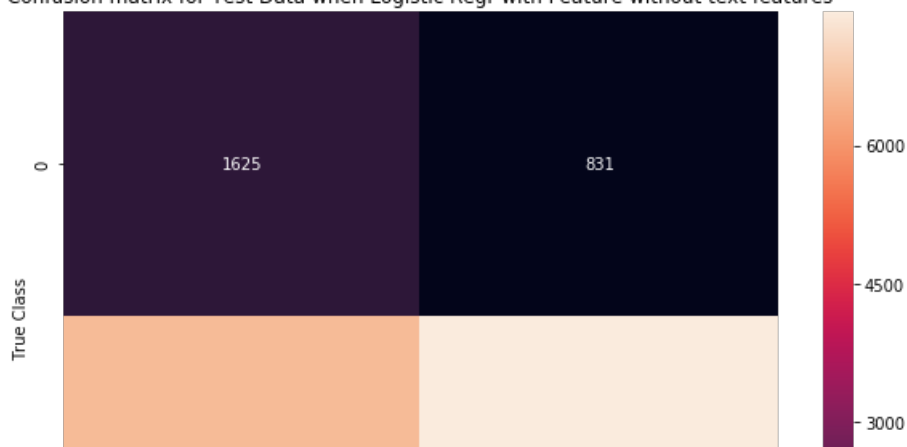
Train confusion matrix

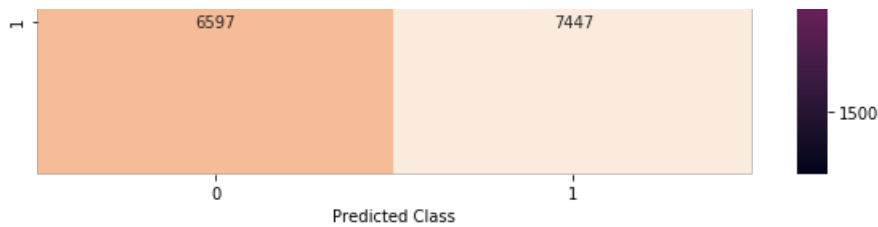
Test confusion matrix

Confusion matrix for Train Data when Logistic Regr with Feature without text features



Confusion matrix for Test Data when Logistic Regr with Feature without text features





In []:

3. Conclusion

In [164]:

```
# Please compare all your models using Prettytable library
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Features", "Model", "Alpha", "Maximum AUC score",]

x.add_row(["BoW", "Logistic Regr", 0.1, 0.64297])
x.add_row(["TFIDF", "MultinomialNB", 1, 0.64251])
x.add_row(["AVG W2V", "Logistic Regr", 1, 0.65610])
x.add_row(["TFIDF W2V", "Logistic Regr", 1, 0.64975])
x.add_row(["Feature WITHOUT Text", "Logistic Regr", 10000, 0.59595])
print(x)
```

Features	Model	Alpha	Maximum AUC score
BoW	Logistic Regr	0.1	0.64297
TFIDF	MultinomialNB	1	0.64251
AVG W2V	Logistic Regr	1	0.6561
TFIDF W2V	Logistic Regr	1	0.64975
Feature WITHOUT Text	Logistic Regr	10000	0.59595

In []:

In []: