

# Support Vector Machine (SVM) 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. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Art Will Make You Happy!</li><li>• First Grade Fun</li></ul>
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none"><li>• Grades PreK-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none"><li>• Applied Learning</li><li>• Care &amp; Hunger</li><li>• Health &amp; Sports</li><li>• History &amp; Civics</li><li>• Literacy &amp; Language</li><li>• Math &amp; Science</li><li>• Music &amp; The Arts</li><li>• Special Needs</li><li>• Warmth</li></ul> <b>Examples:</b> <ul style="list-style-type: none"><li>• Music &amp; The Arts</li><li>• Literacy &amp; Language, Math &amp; Science</li></ul>
<code>school_state</code>	State where school is located ( <u>Two-letter U.S. postal code</u> ( <a href="https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes">https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes</a> )). <b>Example:</b> WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Literacy</li><li>• Literature &amp; Writing, Social Sciences</li></ul>
<code>project_resource_summary</code>	An explanation of the resources needed for the project. <b>Example:</b> <ul style="list-style-type: none"><li>• My students need hands on literacy materials to manage sensory needs!</li></ul>
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*
<code>project_essay_4</code>	Fourth application essay*

Feature	Description
project_submitted_datetime	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"><li>• nan</li><li>• Dr.</li><li>• Mr.</li><li>• Mrs.</li><li>• Ms.</li><li>• Teacher.</li></ul>
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. <b>Example:</b> 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 project_id value from the train.csv file. <b>Example:</b> p036502
description	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. <b>Example:</b> 3
price	Price of the resource required. <b>Example:</b> 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 [2]: import warnings
warnings.filterwarnings("ignore")
%matplotlib inline

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

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

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

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

from tqdm import tqdm
import os

from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
from sklearn.cross_validation import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.cross_validation import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn import cross_validation
```

## 1.1 Reading Data

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

```
In [4]: print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

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

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

```
In [5]: # how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]
```

```
#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
project_data.drop('project_submitted_datetime', axis=1, inplace=True)
project_data.sort_values(by=['Date'], inplace=True)
```

```
# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
project_data = project_data[cols]
```

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

Out[5]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject_categories	project_subject_subcategories	project_title
101880	5749	p096076	6eaa448903897a152320bd23a30147b2	Mrs.	CA	2016-01-05 00:00:00	Grades PreK-2	Math & Science	Mathematics	Math Madness
31477	47750	p185738	3afe10b996b7646d8641985a4b4b570d	Mrs.	UT	2016-01-05 01:05:00	Grades PreK-2	Math & Science	Mathematics	Math is Fun!

```
In [6]: 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[6]:

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

## 1.2 preprocessing of project\_subject\_categories

```
In [7]: 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 [8]: 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]))

```

## Preprocessing of teacher\_prefix

```

In [9]: #“Teacher prefix” data having the dots(.) and its has been observed the some rows are empty in this feature .
#the dot(.) and empty row available in the data consider as float datatype and it does not
# accepted by the .Split() - Pandas function , so removing the same.
# cleaning has been done for the same following references are used
# 1. Removing (.) from dataframe column - used ".str.replce" funtion (padas documentation)
# 2. for empty cell in datafram column - added the "Mrs." (in train data.csv) which has me mostly occured in data set.

project_data["teacher_prefix_clean"] = project_data["teacher_prefix"].str.replace(".", "")
project_data.head(2)
print(project_data.teacher_prefix_clean.shape)

(109248,)

```

## 1.4 Text preprocessing

```

In [10]: # 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 [11]: project_data.head(2)
```

Out[11]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title	project_essay_1	project_essay_2	project_essay_3	project_essay_4
101880	5749	p096076	6eaa448903897a152320bd23a30147b2	Mrs.	CA	2016-01-05 00:00:00	Grades PreK-2	Math Madness	A typical day in our classroom is full of encouragement and exploration. With common core and being more open to allowing students to make more mistakes has helped me improve and see students thinking in a different way. My students are math problem solvers who are enjoying math. I have 28 first graders who want to be heard and understood. Who want to enjoy math. By creating and finding different math games that continues to help them build fluency and number sense, my students are enjoying and doing math at their own pace. They are enjoying what they are learning and want to practice it in numerous ways. These materials will be used in math centers. Students will be able to explore and play games while practicing the skills they need. By playing these games they will be more engaged and will learn as they gain the skills they need to learn. They will practice learning their doubles, practice adding and subtracting and will be able to have fun. I want to create an environment in which my students are loving what they are learning. I will be able to use the donations for countless years. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.	I have 28 first graders who want to be heard and understood. Who want to enjoy math. By creating and finding different math games that continues to help them build fluency and number sense, my students are enjoying and doing math at their own pace. They are enjoying what they are learning and want to practice it in numerous ways. These materials will be used in math centers. Students will be able to explore and play games while practicing the skills they need. By playing these games they will be more engaged and will learn as they gain the skills they need to learn. They will practice learning their doubles, practice adding and subtracting and will be able to have fun. I want to create an environment in which my students are loving what they are learning. I will be able to use the donations for countless years. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.	These materials will be used in math centers. ...	I will be able to use the donations for countless years. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.
31477	47750	p185738	3afe10b996b7646d8641985a4b4b570d	Mrs.	UT	2016-01-05 01:05:00	Grades PreK-2	Math is Fun!	"The only way to learn mathematics is to do math. My students love coming to school and working ...	My students will be using these math manipulatives to help them learn. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.	My students will be using these math manipulatives to help them learn. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.	Learn to use the materials requested to create models of buildings that are imagined. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.

```
In [12]: # 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)
```

A typical day in our classroom is full of encouragement and exploration. With common core and being more open to allowing students to make more mistakes has helped me improve and see students thinking in a different way. My students are math problem solvers who are enjoying math. I have 28 first graders who want to be heard and understood. Who want to enjoy math. By creating and finding different math games that continues to help them build fluency and number sense, my students are enjoying and doing math at their own pace. They are enjoying what they are learning and want to practice it in numerous ways. These materials will be used in math centers. Students will be able to explore and play games while practicing the skills they need. By playing these games they will be more engaged and will learn as they gain the skills they need to learn. They will practice learning their doubles, practice adding and subtracting and will be able to have fun. I want to create an environment in which my students are loving what they are learning. I will be able to use the donations for countless years. I will be able to use the dice in numerous games. I will be able to provide some families with these games to try and continue to practice at home. I want to help my students in every way possible.

=====

My students have learned what amazing adventures they can have when reading picture books; now I want them to realize the endless possibilities for adventure that chapter books give! My students are all bright and inquisitive learners who love to read! My students live in the heart of the city, most in extreme poverty. All students at my school receive free breakfast and lunch; many families also utilize the food pantry that our runs out of its basement on the weekend. My classroom library is filled with picture books. Having these chapter books would enable my students to build their reading stamina within a book. They would be able to apply the many comprehension strategies and skills we've learned to deeper text. Having these chapter books in our classroom library will help my students become VORACIOUS readers! This love of reading will continue in their lives as they continue to second grade and beyond.

=====

"What are we working on today Mrs. Mistry?" is typically the question I get asked as excited students walk into my classroom ready to learn. The students at this school are so excited to walk into a room where they know they will have the chance to express themselves through art. \r\nThese K-5 suburban students are motivated to learn about anything that is handed to them. I enjoy supporting student learning with creative expression, and engagement in the art classroom! Architecture is such an important part of my students lives right now. As our city is growing, students are surrounded by tall buildings and construction. The materials for this project will allow my students to immerse themselves in something that connects art and what they are currently experiencing in their surroundings. \r\nThese growing minds will definitely enjoy seeing their drawing on paper come to life ! First the students will learn beginning steps in constructing a building. Students will learn the process by creating blueprints of buildings, and use the materials requested to create models of buildings that are imagined. nannan

=====

In [13]: `# 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 [14]: `sent = decontracted(project_data['essay'].values[2000])  
print(sent)  
print("="*50)`

My students attend a Title I school in downtown Oakland. Coming from diverse cultural/ethnic backgrounds and socioeconomically disadvantaged neighborhoods, these students know the meaning of perseverance & consistently give their best in all that they do. They are inquisitive, enthusiastic, curious, eager to explore new things and ask compelling questions. \r\nUnsatisfied by the cursory "textbook" explanation and uninterested in just memorizing the the correct formula to get a good grade, these students are always asking the how is and why is, thinking critically and analyzing the information that is presented to them. As a result of their hard work, they have consistently exceeded district norms on standardized testing. Between Math and ELA, we had a total of 14 perfect scores in our class, last year, on the SBAC.\r\nChallenges we face at our school include having limited or outdated technological equipment and software, no science lab, and little funding for resources beyond the basic educational supplies. These students will be contributing members of society one day. Sowing into them is sowing into tomorrow's leaders. During the time for the "curiosity project," students will be able to explore and research within a current unit/topic of study. Students will be led by curiosity, ask inquiring questions, do online research, read e-books (reference materials), and make presentations using the Amazon Kindle Fire. \r\nThis type of "inquiry-based" learning is aimed at sparking interest within students, encouraging them to initiate learning, giving them opportunity to pursue topics that fascinate them, teaching them vital research online research and organizational skills, and challenging them to present information they have learned to the entire class in a compelling and insightful way. It allows students to learn material beyond what is presented in class or in the textbook and is aligned with NGSS Science and Engineering Practices. Currently, our school's technological equipment is outdated and extremely limited. With this project funded, my class will have close to a 2:1 ratio of students to devices.

=====

In [15]: `# \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 students attend a Title I school in downtown Oakland. Coming from diverse cultural/ethnic backgrounds and socioeconomically disadvantaged neighborhoods, these students know the meaning of perseverance & consistently give their best in all that they do. They are inquisitive, enthusiastic, curious, eager to explore new things and ask compelling questions. Unsatisfied by the cursory textbook explanation and uninterested in just memorizing the the correct formula to get a good grade, these students are always asking the how is and why is, thinking critically and analyzing the information that is presented to them. As a result of their hard work, they have consistently exceeded district norms on standardized testing. Between Math and ELA, we had a total of 14 perfect scores in our class, last year, on the SBAC. Challenges we face at our school include having limited or outdated technological equipment and software, no science lab, and little funding for resources beyond the basic educational supplies. These students will be contributing members of society one day. Sowing into them is sowing into tomorrow's leaders. During the time for the "curiosity project," students will be able to explore and research within a current unit/topic of study. Students will be led by curiosity, ask inquiring questions, do online research, read e-books (reference materials), and make presentations using the Amazon Kindle Fire. This type of "inquiry-based" learning is aimed at sparking interest within students, encouraging them to initiate learning, giving them opportunity to pursue topics that fascinate them, teaching them vital research online research and organizational skills, and challenging them to present information they have learned to the entire class in a compelling and insightful way. It allows students to learn material beyond what is presented in class or in the textbook and is aligned with NGSS Science and Engineering Practices. Currently, our school's technological equipment is outdated and extremely limited. With this project funded, my class will have close to a 2:1 ratio of students to devices.

```
In [16]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

My students attend a Title I school in downtown Oakland. Coming from diverse cultural, ethnic backgrounds and socioeconomically disadvantaged neighborhoods, these students know the meaning of perseverance; consistently give their best in all that they do. They are inquisitive, enthusiastic, curious, eager to explore new things, and ask compelling questions. Unsatisfied by the cursory textbook explanation and uninterested in just memorizing the correct formula to get a good grade, these students are always asking the how, is, and why. Thinking critically and analyzing the information that is presented to them. As a result of their hard work, they have consistently exceeded district norms on standardized testing. Between Math and ELA, we had a total of 14 perfect scores in our class last year on the SBAC. Challenges we face at our school include having limited or outdated technological equipment and software, no science lab, and little funding for resources beyond the basic educational supplies. These students will be contributing members of society one day. Sowing into them is sowing into tomorrow's leaders. During the time for the curiosity project, students will be able to explore and research within a current unit topic of study. Students will be led by curiosity, ask inquiring questions, do online research, read e-books, reference materials, and make presentations using the Amazon Kindle Fire. This type of inquiry-based learning is aimed at sparking interest within students, encouraging them to initiate learning, giving them opportunity to pursue topics that fascinate them, teaching them vital research, online research, and organizational skills, and challenging them to present information they have learned to the entire class in a compelling and insightful way. It allows students to learn material beyond what is presented in class or in the textbook and is aligned with NGSS Science and Engineering Practices. Currently, our school's technological equipment is outdated and extremely limited. With this project funded, my class will have close to a 2:1 ratio of students to devices.

```
In [17]: # 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', 'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

### 1.4.1 Data Processing (Essay)

```
In [18]: # Combining all the above statements
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\r', ' ')
    sent = sent.replace('\\"', ' ')
    sent = sent.replace('\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109248 [01:00<00:00, 1803.63it/s]
```

```
In [19]: project data["preprocessed essays"] = preprocessed essays
```

```
In [20]: project_data.shape

Out[20]: (109248, 20)
```

1.4.2 Words in the Essay

```
In [21]: # https://stackoverflow.com/questions/49984905/count-number-of-words-per-row/49984998
project_data['essay_word_count'] = [len(x.split()) for x in project_data['preprocessed_essays'].tolist()]
```

```
In [22]: project_data.shape

Out[22]: (109248, 21)
```

```
In [23]: project_data.head(2)
```

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title	project_essay_1	project_essay_2	...	project_essay_4	prc
101880	5749	p096076	6eaa448903897a152320bd23a30147b2	Mrs.	CA	2016-01-05 00:00:00	Grades PreK-2	Math Madness	A typical day in our classroom is full of enco...	I have 28 first graders who want to be heard a...	...	I will be able to use these donations for coun...	My ma
31477	47750	p185738	3afe10b996b7646d8641985a4b4b570d	Mrs.	UT	2016-01-05 01:05:00	Grades PreK-2	Math is Fun!	"The only way to learn mathematics is to do m...	My students love coming to school and working ...	...	Learning about money is important so the stude...	My blo

2 rows × 21 columns

1.4.3 Sentiment Score for Essay

```
In [24]: import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer
```

```
In [25]: analyser = SentimentIntensityAnalyzer()
```

```
In [26]: positive = []
neutral= []
negative = []
compound = []
for a in (project_data["preprocessed_essays"]):
    P = analyser.polarity_scores(a)['pos']
    Neu = analyser.polarity_scores(a)['neu']
    Neg = analyser.polarity_scores(a)['neg']
    C = analyser.polarity_scores(a)['compound']
    positive.append(P)
    neutral.append(Neu)
    negative.append(Neg)
    compound.append(C)
```

```
In [27]: project_data["Positive_SC_Essay"] = positive

In [28]: project_data["Neutral_SC_Essay"] = neutral

In [29]: project_data["Negative_SC_Essay"] = negative

In [30]: project_data["Compound_SC_Essay"] = compound

In [31]: project_data.head(2)
```

Out[31]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title	project_essay_1	project_essay_2	...	clean_categories	cl
101880	5749	p096076	6eaa448903897a152320bd23a30147b2	Mrs.	CA	2016-01-05 00:00:00	Grades PreK-2	Math Madness	A typical day in our classroom is full of enco...	I have 28 first graders who want to be heard a...	...	Math_Science	Mi
31477	47750	p185738	3afe10b996b7646d8641985a4b4b570d	Mrs.	UT	2016-01-05 01:05:00	Grades PreK-2	Math is Fun!	"The only way to learn mathematics is to do m...	My students love coming to school and working ...	...	Math_Science	Mi

2 rows × 25 columns

1.5 Preprocessing of `project\_title`

```
In [32]: # Data processing for project titles
Title_clean = project_data.project_title
Title_clean.head(2)

Out[32]: 101880    Math Madness
31477    Math is Fun!
Name: project_title, dtype: object

In [33]: P = decontracted(project_data['project_title'].values[1])
print(P)

Math is Fun!

In [34]: # \r \n \t and -- remove from string python: http://texthandler.com/info/remove-line-breaks-python/
P = P.replace('\r', ' ')
P = P.replace('\\"', ' ')
P = P.replace('\n', ' ')
P = P.replace('--', ' ')
print(P)

Math is Fun!
```

1.5.1 Data Pracessing (Project Title)

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109248 [00:02<00:00, 38987.39it/s]
```

1 [27] [\[1, 2, 3, 4, 5, 6, 7\]](#)

\_\_\_\_\_

2 rows  $\times$  27 columns

◀ ▶

```
In [39]: #As recommended in the Lecture video, splitting the Data in Train, Test and Cross validation data set
#before applying Vectorization to avoid the data leakage issues.
# As suggested to use stratify sampling, Referred following site for code
# https://stackoverflow.com/questions/29438265/stratified-train-test-split-in-scikit-Learn

# split the data set into train and test
X_train, X_test, y_train, y_test = cross_validation.train_test_split(project_data, project_data['project_is_approved'], test_size=0.3, stratify = project_data['project_is_approved'])

# split the train data set into cross validation train and cross validation test
X_train, X_cv, y_train, y_cv = cross_validation.train_test_split(X_train, y_train, test_size=0.3, stratify=y_train)
```

```
In [40]: #Removing the class label from the data set, in our case the class label is "project is approved"
#From all Train, Test and Cross validation data set

#Train Data

X_train.drop(['project_is_approved'], axis = 1, inplace = True)

#Test Data

X_test.drop(['project_is_approved'], axis = 1, inplace = True)

#Cross Validation data

X_cv.drop(['project_is_approved'], axis = 1, inplace = True)
```

## 1.6 Preparing data for models

```
In [41]: project_data.columns
```

```
Out[41]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
               'Date', 'project_grade_category', 'project_title', 'project_essay_1',
               'project_essay_2', 'project_essay_3', 'project_essay_4',
               'project_resource_summary',
               'teacher_number_of_previously_posted_projects', 'project_is_approved',
               'clean_categories', 'clean_subcategories', 'teacher_prefix_clean',
               'essay', 'preprocessed_essays', 'essay_word_count', 'Positive_SC_Essay',
               'Neutral_SC_Essay', 'Negative_SC_Essay', 'Compound_SC_Essay',
               'preprocessed_Titles', 'title_word_count'],
              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

### 1.6.1 Vectorizing Categorical data

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

#### Project\_categories - Vectorization



```
In [42]: # 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)

vectorizer.fit(X_train['clean_categories'].values)

categories_one_hot_train = vectorizer.transform(X_train['clean_categories'].values)
categories_one_hot_cv = vectorizer.transform(X_cv['clean_categories'].values)
categories_one_hot_test = vectorizer.transform(X_test['clean_categories'].values)

print(vectorizer.get_feature_names())

print("Shape of matrix after one hot encodig ",categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig ",categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig ",categories_one_hot_test.shape)

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig  (53531, 9)
Shape of matrix after one hot encodig  (22942, 9)
Shape of matrix after one hot encodig  (32775, 9)
```

### Project\_sub\_categories - Vectorization

```
In [43]: # we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)

vectorizer.fit(X_train['clean_subcategories'].values)

sub_categories_one_hot_train = vectorizer.transform(X_train['clean_subcategories'].values)
sub_categories_one_hot_cv = vectorizer.transform(X_cv['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer.transform(X_test['clean_subcategories'].values)

print(vectorizer.get_feature_names())

print("Shape of matrix after one hot encodig ",sub_categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig ",sub_categories_one_hot_test.shape)

['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encodig  (53531, 30)
Shape of matrix after one hot encodig  (22942, 30)
Shape of matrix after one hot encodig  (32775, 30)
```

### School\_State - Vectorization

```
In [44]: # we use count vectorizer to convert the values into one hot encoded features
from collections import Counter
my_counter_state = Counter()
for word in project_data['school_state'].values:
    my_counter_state.update(word.split())

state_dict = dict(my_counter_state)
sorted_state_dict = dict(sorted(state_dict.items(), key=lambda kv: kv[1]))

vectorizer = CountVectorizer(vocabulary=list(sorted_state_dict.keys()), lowercase=False, binary=True)

vectorizer.fit(X_train['school_state'].values)

school_state_one_hot_train = vectorizer.transform(X_train['school_state'].values)
school_state_one_hot_cv = vectorizer.transform(X_cv['school_state'].values)
school_state_one_hot_test = vectorizer.transform(X_test['school_state'].values)

print(vectorizer.get_feature_names())

print("Shape of matrix after one hot encodig ",school_state_one_hot_train.shape)
print("Shape of matrix after one hot encodig ",school_state_one_hot_cv.shape)
print("Shape of matrix after one hot encodig ",school_state_one_hot_test.shape)
```

```
['VT', 'WY', 'ND', 'MT', 'RI', 'SD', 'NE', 'DE', 'AK', 'NH', 'WV', 'ME', 'HI', 'DC', 'NM', 'KS', 'IA', 'ID', 'AR', 'CO', 'MN', 'OR', 'KY', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'A
L', 'WI', 'VA', 'AZ', 'NJ', 'OK', 'WA', 'MA', 'LA', 'OH', 'MO', 'IN', 'PA', 'MI', 'SC', 'GA', 'IL', 'NC', 'FL', 'NY', 'TX', 'CA']
Shape of matrix after one hot encodig  (53531, 51)
Shape of matrix after one hot encodig  (22942, 51)
Shape of matrix after one hot encodig  (32775, 51)
```

### teacher\_prefix - Vectorization

```
In [45]: #“Teacher prefix” data having the dots(.) and its has been observed the some rows are empty in this feature .
#the dot(.) and empty row available in the data consider as float datatype and it does not
# accepted by the .Split() - Pandas function , so removing the same.
# cleaning has been done for the same following references are used
# 1. Removing (.) from dataframe column - used ".str.replce" funtion (padas documentation)
# 2. for empty cell in datafram column - added the "Mrs." (in train data.csv) which has me mostly occured in data set.

project_data["teacher_prefix_clean"] = project_data["teacher_prefix"].str.replace(".", "")
project_data.head(2)
print(project_data.teacher_prefix_clean.shape)

(109248,)
```

```
In [46]: from collections import Counter
my_counter_T = Counter()
for word in project_data["teacher_prefix_clean"].values:

    my_counter_T.update(word.split())

Teacher_dict = dict(my_counter_T)
sorted_Teacher_dict = dict(sorted(Teacher_dict.items(), key=lambda kv: kv[1]))

vectorizer = CountVectorizer(vocabulary=list(Teacher_dict.keys()), lowercase=False, binary=True)
#vectorizer.fit(project_data.teacher_prefix_clean.values)

vectorizer.fit(X_train["teacher_prefix_clean"].values)
print(vectorizer.get_feature_names())

Teacher_Prefix_one_hot_train = vectorizer.transform(X_train["teacher_prefix_clean"].values)
Teacher_Prefix_one_hot_cv = vectorizer.transform(X_cv["teacher_prefix_clean"].values)
Teacher_Prefix_one_hot_test = vectorizer.transform(X_test["teacher_prefix_clean"].values)

print("Shape of matrix after one hot encodig ",Teacher_Prefix_one_hot_train.shape)
print("Shape of matrix after one hot encodig ",Teacher_Prefix_one_hot_cv.shape)
print("Shape of matrix after one hot encodig ",Teacher_Prefix_one_hot_test.shape)

['Mrs', 'Ms', 'Mr', 'Teacher', 'Dr']
Shape of matrix after one hot encodig  (53531, 5)
Shape of matrix after one hot encodig  (22942, 5)
Shape of matrix after one hot encodig  (32775, 5)
```

### project\_grade\_category - Vectorization

```
In [47]: # Used this as reference to avoid the space between grades and category ,
# it has split the string with comma , now getting four project grade category as required.
# https://stackoverflow.com/questions/4071396/split-by-comma-and-strip-whitespace-in-python
from collections import Counter
my_counter_project_grade_category= Counter()
for word in project_data['project_grade_category'].values:
    my_counter_project_grade_category.update(word.split(','))

project_grade_category_dict = dict(my_counter_project_grade_category)
sorted_project_grade_category_prefix_dict = dict(sorted(project_grade_category_dict.items(), key=lambda kv: kv[1]))

vectorizer = CountVectorizer(vocabulary=list(project_grade_category_dict.keys()), lowercase=False, binary=True)

vectorizer.fit(X_train["project_grade_category"].values)
print(vectorizer.get_feature_names())

project_grade_category_one_hot_train = vectorizer.transform(X_train["project_grade_category"].values)
project_grade_category_one_hot_cv = vectorizer.transform(X_cv["project_grade_category"].values)
project_grade_category_one_hot_test = vectorizer.transform(X_test["project_grade_category"].values)

print("Shape of matrix after one hot encoding ",project_grade_category_one_hot_train.shape)
print("Shape of matrix after one hot encoding ",project_grade_category_one_hot_cv.shape)
print("Shape of matrix after one hot encoding ",project_grade_category_one_hot_test.shape)

['Grades PreK-2', 'Grades 9-12', 'Grades 6-8', 'Grades 3-5']
Shape of matrix after one hot encoding  (53531, 4)
Shape of matrix after one hot encoding  (22942, 4)
Shape of matrix after one hot encoding  (32775, 4)
```

## 1.6.2 Vectorizing Text data

### 1.6.2.1 Bag of words

#### Train Data Vectorization - BOW (essays)

```
In [48]: # We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train["preprocessed_essays"])
bow_essays_train = vectorizer.fit_transform(X_train["preprocessed_essays"])
print("Shape of matrix after one hot encoding ",bow_essays_train.shape)

Shape of matrix after one hot encoding  (53531, 12597)
```

#### CV Data Vectorization - BOW (essays)

```
In [49]: bow_essays_cv = vectorizer.transform(X_cv["preprocessed_essays"])
print("Shape of matrix after one hot encoding ",bow_essays_cv.shape)

Shape of matrix after one hot encoding  (22942, 12597)
```

**Test Data Vectorization - BOW (essays)**

```
In [50]: bow_essays_test = vectorizer.transform(X_test["preprocessed_essays"])
print("Shape of matrix after one hot encoding ",bow_essays_test.shape)
```

Shape of matrix after one hot encoding (32775, 12597)

**Train Data Vectorization - BOW (Project Titles)**

```
In [51]: # We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
bow_title_train = vectorizer.fit_transform(X_train["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_train.shape)
```

Shape of matrix after one hot encodig (53531, 2199)

**CV Data Vectorization - BOW (Project Titles)**

```
In [52]: bow_title_cv = vectorizer.transform(X_cv["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_cv.shape)
```

Shape of matrix after one hot encodig (22942, 2199)

**Test Data Vectorization - BOW (Project Titles)**

```
In [53]: bow_title_test = vectorizer.transform(X_test["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_test.shape)
```

Shape of matrix after one hot encodig (32775, 2199)

**1.6.2.2 TFIDF vectorizer****Train Data Vectorization - TFIDF (essays)**

```
In [54]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
tfidf_essays_train = vectorizer.fit_transform(X_train["preprocessed_essays"])
print("Shape of matrix after one hot encodig ",tfidf_essays_train.shape)
```

Shape of matrix after one hot encodig (53531, 12597)

**CV Data Vectorization - TFIDF (essays)**

```
In [55]: tfidf_essays_cv = vectorizer.transform(X_cv["preprocessed_essays"])
print("Shape of matrix after one hot encodig ",tfidf_essays_cv.shape)
```

Shape of matrix after one hot encodig (22942, 12597)

**Test Data Vectorization - TFIDF (essays)**

```
In [56]: tfidf_essays_test = vectorizer.transform(X_test["preprocessed_essays"])
print("Shape of matrix after one hot encodig ",tfidf_essays_test.shape)
```

Shape of matrix after one hot encodig (32775, 12597)

**Train Data Vectorization - TFIDF (Project Titles)**

```
In [57]: vectorizer = CountVectorizer(min_df=10)
tfidf_title_train = vectorizer.fit_transform(X_train["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_train.shape)
```

Shape of matrix after one hot encodig (53531, 2199)

**CV Data Vectorization - TFIDF (Project Titles)**

```
In [58]: tfidf_title_cv = vectorizer.transform(X_cv["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_cv.shape)
```

Shape of matrix after one hot encodig (22942, 2199)

**Test Data Vectorization - TFIDF (Project Titles)**

```
In [59]: tfidf_title_test = vectorizer.transform(X_test["preprocessed_Titles"])
print("Shape of matrix after one hot encodig ",bow_title_test.shape)
```

Shape of matrix after one hot encodig (32775, 2199)

**1.6.2.3 Using Pretrained Models: Avg W2V**

```
In [60]: # 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')
```

Loading Glove Model

1917495it [04:31, 7058.11it/s]

Done. 1917495 words loaded!

```
In [61]: words = []
for i in X_train["preprocessed_essays"]:
    words.extend(i.split(' '))

for i in X_train["preprocessed_essays"]:
    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))
```

all the words in the coupus 16222090  
the unique words in the coupus 42750  
The number of words that are present in both glove vectors and our coupus 39068 ( 91.387 %)  
word 2 vec length 39068

```
In [62]: words = []
for i in X_train["preprocessed_Titles"]:
    words.extend(i.split(' '))

for i in X_train["preprocessed_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))
```

all the words in the coupus 464328  
the unique words in the coupus 12267  
The number of words that are present in both glove vectors and our coupus 11731 ( 95.631 %)  
word 2 vec length 11731

```
In [63]: import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)
```

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

### Train Data Vectorization - AGV\_W2V (essays)

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essays_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train["preprocessed_essays"]): # 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_essays_train.append(vector)

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

[illegible]

53531  
300

### CV Data Vectorization - AGV\_W2V (essays)

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essays_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv["preprocessed_essays"]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essays_cv.append(vector)

print(len(avg_w2v_essays_cv))
print(len(avg_w2v_essays_cv[0]))
```

[illegible]

22942  
300

### Test Data Vectorization - AGV\_W2V (essays)



```
In [67]: # average Word2Vec
# compute average word2vec for each review.
avg_w2v_essays_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test["preprocessed_essays"]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essays_test.append(vector)

print(len(avg_w2v_essays_test))
print(len(avg_w2v_essays_test[0]))
```

[illegible]

32775

300

### Train Data Vectorization - AGV\_W2V (Project Titles)

```
In [68]: # average Word2Vec
# compute average word2vec for each review.
avg_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train["preprocessed_Titles"]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title_train.append(vector)

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

[illegible]

53531

300

### CV Data Vectorization - AGV\_W2V (Project Titles)

```
100%|██████████████████████████████████████████████████████████████████████████| 22942/22942 [00:00<00:00, 57336.81it/s]
22942
300
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 32775/32775 [00:00<00:00, 54953.25it/s]
32775
300
```

### Train Data Vectorization - TFIDF\_W2V (essays)

```
In [72]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essays_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train["preprocessed_essays"]): # for each review/sentence
    vector = np.zeros(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_essays_train.append(vector)

print(len(tfidf_w2v_essays_train))
print(len(tfidf_w2v_essays_train[0]))
```

100%|██| 53531/53531 [01:56<00:00, 457.99it/s]

53531  
300

### CV Data Vectorization - TFIDF\_W2V (essays)

```
In [73]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_essays_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv["preprocessed_essays"]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the 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_essays_cv.append(vector)

print(len(tfidf_w2v_essays_cv))
print(len(tfidf_w2v_essays_cv[0]))
```

100%|██| 22942/22942 [00:48<00:00, 477.07it/s]

22942  
300

### Test Data Vectorization - TFIDF\_W2V (essays)

```
100%|██████████████████████████████████████████████████████████████████████████████| 32775/32775 [01:07<00:00, 483.08it/s]
32775
300
```

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train["preprocessed_Titles"])
# 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_titles = set(tfidf_model.get_feature_names())
```

### Train Data Vectorization - TFIDF\_W2V (Project Titles)

```
100%|██████████████████████████████████████████████████████████████████████████████| 53531/53531 [00:01<00:00, 27190.77it/s]
53531
300
```

### CV Data Vectorization - TFIDF\_W2V (Project Titles)

```
In [77]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv["preprocessed_Titles"]): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words_titles ):
            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_cv.append(vector)

print(len(tfidf_w2v_title_cv))
print(len(tfidf_w2v_title_cv[0]))
```

---

```
100%|██████████████████████████████████████| 22942/22942 [00:00<00:00, 26547.55it/s]
```

```
22942  
300
```

### Test Data Vectorization - TFIDF\_W2V (Project Titles)

```
In [78]: # average Word2Vec  
# compute average word2vec for each review.  
tfidf_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list  
for sentence in tqdm(X_test["preprocessed_Titles"]): # for each review/sentence  
    vector = np.zeros(300) # as word vectors are of zero length  
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review  
    for word in sentence.split(): # for each word in a review/sentence  
        if (word in glove_words) and (word in tfidf_words_titles):  
            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_test.append(vector)  
  
print(len(tfidf_w2v_title_test))  
print(len(tfidf_w2v_title_test[0]))
```

100%|██| 32775/32775 [00:01<00:00, 27177.97it/s]

32775  
300

### 1.6.3 Vectorizing Numerical features

1.6.3.1 Vectorizing Numerical features - Price

```
In [79]: price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()

# Merging the project data train , Cv , test with price from resource data
X_train = pd.merge(X_train, price_data, on='id', how='left')
X_cv = pd.merge(X_cv, price_data, on='id', how='left')
X_test = pd.merge(X_test, price_data, on='id', how='left')
```

```
In [80]: X_train.head(2)
```

Out[80]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title	project_essay_1	project_essay_2	...	preprocessed_essays	essays
0	17618	p108429	7b7eec8da7fdc201396ae1e33f096fac	Mrs.	IL	2016-09-25 19:45:00	Grades 3-5	iTry iLearn iSucceed with an iPad	There are no students at our school, only cade...	These iPads will greatly assist my students wi...	...	there no students school cadets children walk ...	152
1	21505	p091469	f39b8fe67437286195a887a7a426b58b	Mrs.	UT	2016-07-09 23:16:00	Grades PreK-2	Engaging Supplies and Books for Learning	My 2nd graders are a great group of eager, mot...	These materials will help students to make con...	...	my 2nd graders great group eager motivated cur...	130

2 rows × 28 columns

```
In [81]: #https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html

from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['price'].values.reshape(-1,1))

price_train = normalizer.transform(X_train['price'].values.reshape(-1,1))
price_cv = normalizer.transform(X_cv['price'].values.reshape(-1,1))
price_test = normalizer.transform(X_test['price'].values.reshape(-1,1))

print(price_train.shape)
print(price_cv.shape)
print(price_test.shape)
```

(53531, 1)  
(22942, 1)  
(32775, 1)

1.6.3.2 Vectorizing Numerical features - teacher\_number\_of\_previously\_posted\_projects

```
In [82]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

prev_post_train = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
prev_post_cv = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
prev_post_test = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

print(prev_post_train.shape)
print(prev_post_cv.shape)
print(prev_post_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.3 Vectorizing Numerical features - Quantity

```
In [83]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['quantity'].values.reshape(-1,1))

Quantity_train = normalizer.transform(X_train['quantity'].values.reshape(-1,1))
Quantity_cv = normalizer.transform(X_cv['quantity'].values.reshape(-1,1))
Quantity_test = normalizer.transform(X_test['quantity'].values.reshape(-1,1))

print(Quantity_train.shape)
print(Quantity_cv.shape)
print(Quantity_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.4 Vectorizing Numerical features - Project Title word count

```
In [84]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['title_word_count'].values.reshape(-1,1))

title_word_count_train = normalizer.transform(X_train['title_word_count'].values.reshape(-1,1))
title_word_count_cv = normalizer.transform(X_cv['title_word_count'].values.reshape(-1,1))
title_word_count_test = normalizer.transform(X_test['title_word_count'].values.reshape(-1,1))

print(title_word_count_train.shape)
print(title_word_count_cv.shape)
print(title_word_count_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.4 Vectorizing Numerical features - Essay word count

```
In [85]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['essay_word_count'].values.reshape(-1,1))

essay_word_count_train = normalizer.transform(X_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_cv = normalizer.transform(X_cv['essay_word_count'].values.reshape(-1,1))
essay_word_count_test = normalizer.transform(X_test['essay_word_count'].values.reshape(-1,1))

print(essay_word_count_train.shape)
print(essay_word_count_cv.shape)
print(essay_word_count_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.5 Vectorizing Numerical features - Essay Sentiment score – Positive

```
In [86]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['Positive_SC_Essay'].values.reshape(-1,1))

Positive_SC_Essay_train = normalizer.transform(X_train['Positive_SC_Essay'].values.reshape(-1,1))
Positive_SC_Essay_cv = normalizer.transform(X_cv['Positive_SC_Essay'].values.reshape(-1,1))
Positive_SC_Essay_test = normalizer.transform(X_test['Positive_SC_Essay'].values.reshape(-1,1))

print(Positive_SC_Essay_train.shape)
print(Positive_SC_Essay_cv.shape)
print(Positive_SC_Essay_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.6 Vectorizing Numerical features - Essay Sentiment score – Neutral

```
In [87]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['Neutral_SC_Essay'].values.reshape(-1,1))

Neutral_SC_Essay_train = normalizer.transform(X_train['Neutral_SC_Essay'].values.reshape(-1,1))
Neutral_SC_Essay_cv = normalizer.transform(X_cv['Neutral_SC_Essay'].values.reshape(-1,1))
Neutral_SC_Essay_test = normalizer.transform(X_test['Neutral_SC_Essay'].values.reshape(-1,1))

print(Neutral_SC_Essay_train.shape)
print(Neutral_SC_Essay_cv.shape)
print(Neutral_SC_Essay_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.7 Vectorizing Numerical features - Essay Sentiment score – Negative



```
In [88]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['Negative_SC_Essay'].values.reshape(-1,1))

Negative_SC_Essay_train = normalizer.transform(X_train['Negative_SC_Essay'].values.reshape(-1,1))
Negative_SC_Essay_cv = normalizer.transform(X_cv['Negative_SC_Essay'].values.reshape(-1,1))
Negative_SC_Essay_test = normalizer.transform(X_test['Negative_SC_Essay'].values.reshape(-1,1))

print(Negative_SC_Essay_train.shape)
print(Negative_SC_Essay_cv.shape)
print(Negative_SC_Essay_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

### 1.6.3.8 Vectorizing Numerical features - Essay Sentiment score – Compound

```
In [89]: from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['Compound_SC_Essay'].values.reshape(-1,1))

Compound_SC_Essay_train = normalizer.transform(X_train['Compound_SC_Essay'].values.reshape(-1,1))
Compound_SC_Essay_cv = normalizer.transform(X_cv['Compound_SC_Essay'].values.reshape(-1,1))
Compound_SC_Essay_test = normalizer.transform(X_test['Compound_SC_Essay'].values.reshape(-1,1))

print(Compound_SC_Essay_train.shape)
print(Compound_SC_Essay_cv.shape)
print(Compound_SC_Essay_test.shape)

(53531, 1)
(22942, 1)
(32775, 1)
```

## Assignment 7:SVM

### 1. [Task-1] Apply Support Vector Machines(SGDClassifier with hinge loss: Linear SVM) on these feature sets

- Set 1: categorical, numerical features + project\_title(BOW) + preprocessed\_eassay (BOW)
- Set 2: categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (TFIDF)
- Set 3: categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)
- Set 4: categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_eassay (TFIDF W2V)

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

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

### 3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
- Along with plotting ROC curve, you need to print the confusion matrix (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tp-r-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.



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

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

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

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

### 4. [Task-2] Apply the Support Vector Machines on these features by finding the best hyper paramter as suggested in step 2 and step 3 (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

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

- Consider these set of features Set 5 : (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)
  - **school\_state** : categorical data
  - **clean\_categories** : categorical data
  - **clean\_subcategories** : categorical data
  - **project\_grade\_category** :categorical data
  - **teacher\_prefix** : categorical data
  - **quantity** : numerical data
  - **teacher\_number\_of\_previously\_posted\_projects** : numerical data
  - **price** : numerical data
  - **sentiment score's of each of the essay** : numerical data
  - **number of words in the title** : numerical data
  - **number of words in the combine essays** : numerical data
- Apply (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)**TruncatedSVD** (<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>) on **TfidfVectorizer** ([https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.TfidfVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html)) of essay text, choose the number of components (`n_components`) using **elbow method** (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/>) : numerical data

### • Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link (<http://zetcode.com/python/prettytable/>)



**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-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method `fit_transform()` on you train data, and apply the method `transform()` on cv/test data.
4. For more details please go through this [link](https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf). (<https://soundcloud.com/applied-ai-course/leakage-bow-and-tfidf>)

## 2. Support Vector Machines

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

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

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

#### 2.4.1 Applying SGD Classifier (with Loss = 'hinge') on BOW, SET 1

```
In [90]: from scipy.sparse import hstack
X_train_bow = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train ,Teacher_Prefix_one_hot_train,project_grade_category_one_hot_train,bow_essays_train,bow_title_train,price_train,prev_post_train)).tocsr()
X_train_bow.shape
```

```
Out[90]: (53531, 14897)
```

```
In [91]: X_cv_bow = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv ,Teacher_Prefix_one_hot_cv,project_grade_category_one_hot_cv,bow_essays_cv,bow_title_cv,price_cv,prev_post_cv)).tocsr()
X_cv_bow.shape
```

```
Out[91]: (22942, 14897)
```

```
In [92]: X_test_bow = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test ,Teacher_Prefix_one_hot_test,project_grade_category_one_hot_test,bow_essays_test,bow_title_test,price_test,prev_post_test)).tocsr()
X_test_bow.shape
```

```
Out[92]: (32775, 14897)
```

#### GridSearchCV - Finding the best hyper parameter That maximum AUC value

```
In [93]: from sklearn.linear_model import SGDClassifier
import math
```

#### With L1 Regularizer

```
In [95]: # https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV

sgd = SGDClassifier(loss="hinge",penalty='l1',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

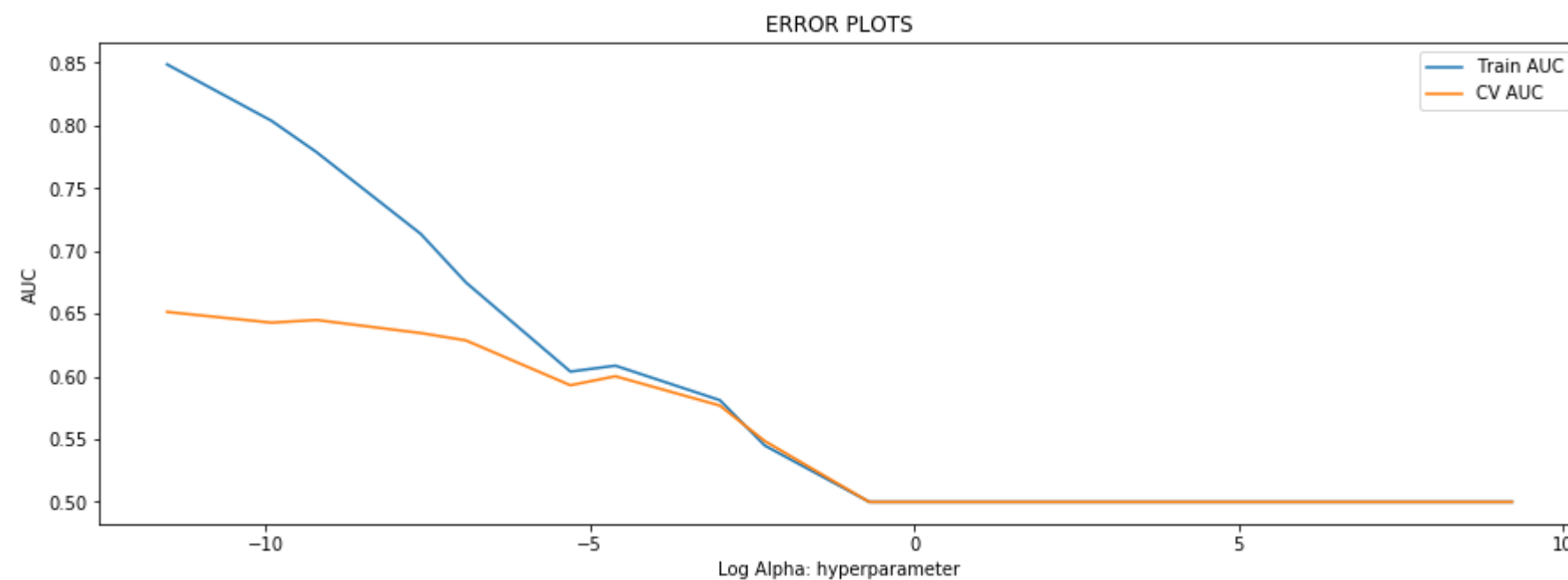
tuned_parameters = [{'alpha': Cs}]

clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_bow, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [16,9]
plt.show()
```



In [96]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.8487617527469937	0.6514610857592005
5e-05	-9.903487552536127	0.8038042781814289	0.6429078609012243
0.0001	-9.210340371976182	0.7788174830037868	0.6450251439928333
0.0005	-7.600902459542082	0.7135108659328492	0.6345847235930451
0.001	-6.907755278982137	0.6749227494937656	0.6287420823445666
0.005	-5.298317366548036	0.6039131297234986	0.5930110253982379
0.01	-4.605170185988091	0.6086202660276775	0.6001937712419987
0.05	-2.995732273553991	0.5811128058484601	0.5768303059459602
0.1	-2.3025850929940455	0.5449011560446215	0.5485012292627028
0.5	-0.6931471805599453	0.5	0.5
1	0.0	0.5	0.5
5	1.6094379124341003	0.5	0.5
10	2.302585092994046	0.5	0.5
50	3.912023005428146	0.5	0.5
100	4.605170185988092	0.5	0.5
500	6.214608098422191	0.5	0.5
1000	6.907755278982137	0.5	0.5
2500	7.824046010856292	0.5	0.5
5000	8.517193191416238	0.5	0.5
10000	9.210340371976184	0.5	0.5

With L2 Regularizer

```
In [97]: # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.GridSearchCV.html
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV

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

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

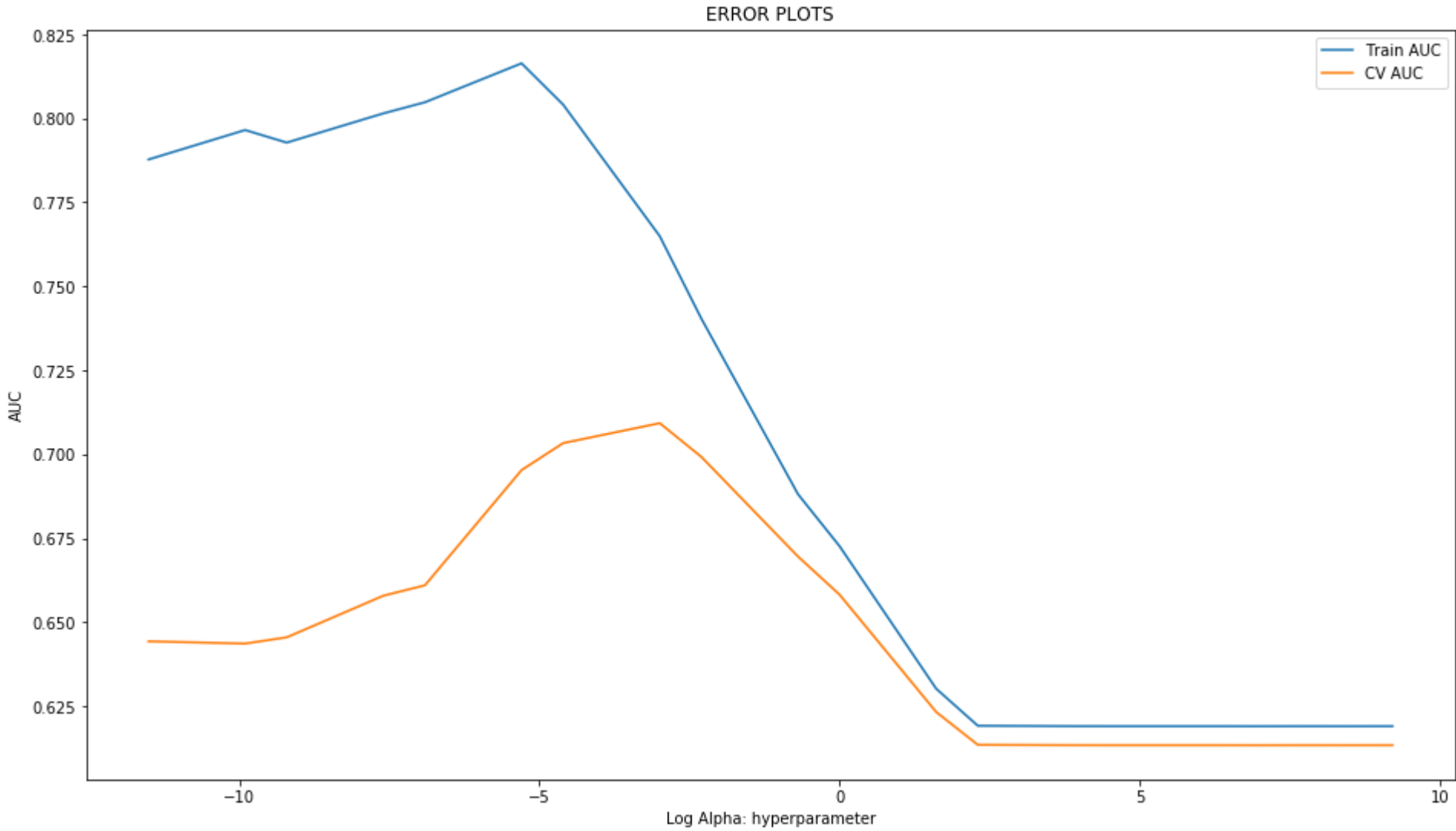
tuned_parameters = [{'alpha': Cs}]

clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_bow, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [15,5]
plt.show()
```



In [98]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.787762071955261	0.644345777583368
5e-05	-9.903487552536127	0.7965308552053161	0.6436757205633944
0.0001	-9.210340371976182	0.7927885757261879	0.6455340246661738
0.0005	-7.600902459542082	0.801480190076245	0.6579222363285815
0.001	-6.907755278982137	0.804794645874658	0.6610225268501833
0.005	-5.298317366548036	0.8163840565398622	0.695312431417753
0.01	-4.605170185988091	0.8040863242240556	0.7033456224761007
0.05	-2.995732273553991	0.7649413604962195	0.7092838388121385
0.1	-2.3025850929940455	0.7403964754729392	0.6992428535251365
0.5	-0.6931471805599453	0.6881561943476386	0.6696014760757534
1	0.0	0.6726627629299009	0.6582366305556617
5	1.6094379124341003	0.6302527966185135	0.6233164641270065
10	2.302585092994046	0.6192224448201072	0.6135454944134426
50	3.912023005428146	0.6190938664595825	0.6134460168163379
100	4.605170185988092	0.6190887083067025	0.6134420408921606
500	6.214608098422191	0.6190880095517571	0.6134414053960832
1000	6.907755278982137	0.6190746034495699	0.6134257653618401
2500	7.824046010856292	0.6190870802261856	0.6134408681088969
5000	8.517193191416238	0.6190780698485819	0.6134306474887857
10000	9.210340371976184	0.6190892886440303	0.6134428967797573

Using Best alpha Value – Training the Model

In [99]:

```
#Taking the Optimal hypermeter from L2 Regularizer for FPR, TPR plot and confusion matrix  
best_alpha_bow = 0.01
```



```

In [100]: #https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

SGD = SGDClassifier(loss="hinge",penalty='l2',alpha= best_alpha_bow,class_weight = 'balanced')
SGD.fit(X_train_bow, y_train)

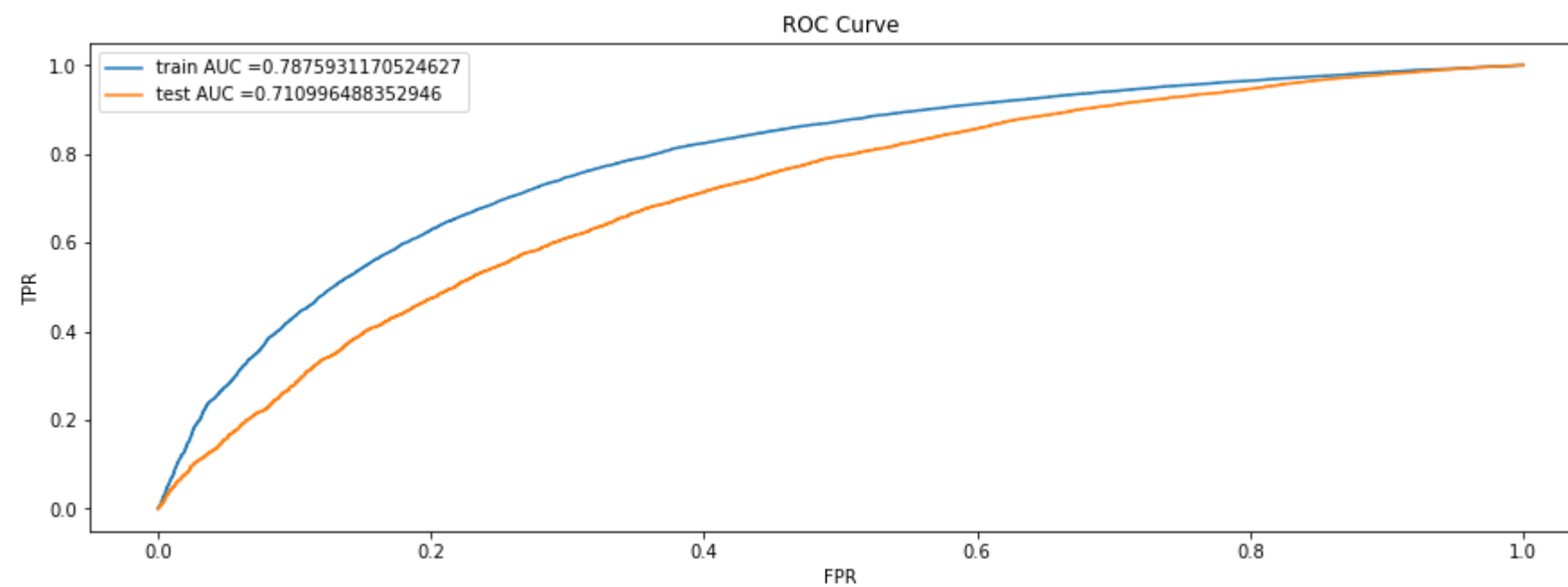
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_bow = SGD.decision_function(X_train_bow)
y_test_pred_bow = SGD.decision_function(X_test_bow)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_bow)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_bow)

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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")
plt.rcParams["figure.figsize"] = [5,5]
plt.show()

```



## Confusion Matrix

### Train confusion matrix

In [101]: *#https://stackoverflow.com/questions/28719067/roc-curve-and-cut-off-point-python*

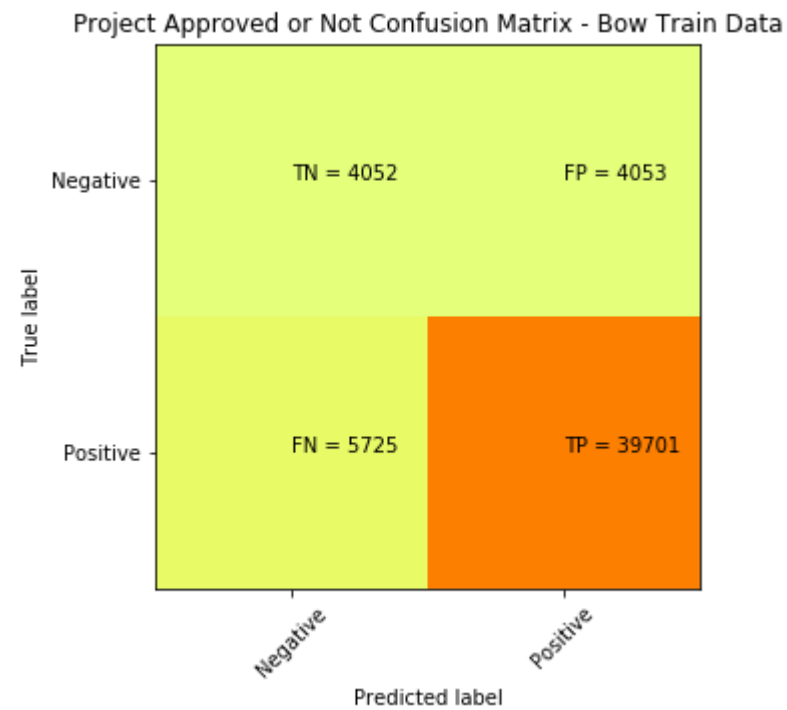
```
def predict(proba, threshold, fpr, tpr):
    t = threshold[np.argmax(fpr*(1-tpr))]
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

In [102]: **from sklearn.metrics import confusion\_matrix**  
print("Train confusion matrix")  
bow\_train\_confusion\_matrix = confusion\_matrix(y\_train, predict(y\_train\_pred\_bow, tr\_thresholds, train\_fpr, train\_fpr))  
print(bow\_train\_confusion\_matrix)

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999999619430507 for threshold -0.406
[[ 4052  4053]
 [ 5725 39701]]
```

In [103]: <http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/>

```
plt.clf()
plt.imshow(bow_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Bow Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(bow_train_confusion_matrix[i][j]))
plt.show()
```



### Test confusion matrix

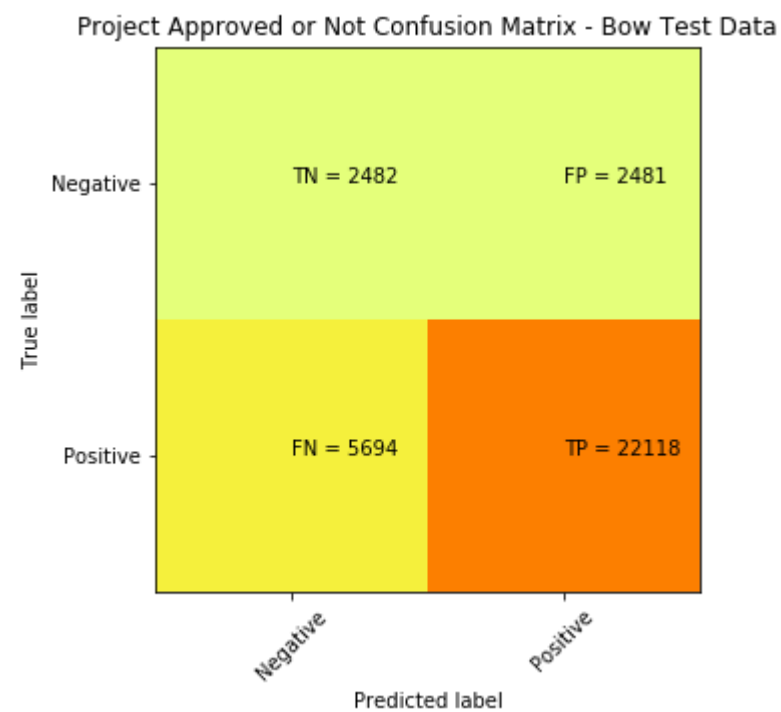
In [104]: 

```
print("Train confusion matrix")
bow_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_bow, te_thresholds, test_fpr, test_fpr))
print(bow_test_confusion_matrix)
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999998985034083 for threshold -0.178
[[ 2482  2481]
 [ 5694 22118]]
```

In [105]: <http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/>

```
plt.clf()
plt.imshow(bow_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Bow Test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(bow_test_confusion_matrix[i][j]))
plt.show()
```



## 2.4.2 Applying SGD Classifier (with Loss = 'hinge') on TFIDF, SET 2

In [106]: `X_train_tfidf = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train , Teacher_Prefix_one_hot_train, project_grade_category_one_hot_train, tfidf_essays_train, tfidf_title_train, price_train, prev_post_train)).tocsr()`  
`X_train_tfidf.shape`

Out[106]: (53531, 14897)

In [107]: `X_cv_tfidf = hstack((categories_one_hot_cv, sub_categories_one_hot_cv, school_state_one_hot_cv , Teacher_Prefix_one_hot_cv, project_grade_category_one_hot_cv, tfidf_essays_cv, tfidf_title_cv, price_cv, prev_post_cv)).tocsr()`  
`X_cv_tfidf.shape`

Out[107]: (22942, 14897)

```
In [108]: X_test_tfidf = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test ,Teacher_Prefix_one_hot_test,project_grade_category_one_hot_test,tfidf_essays_test,tfidf_title_test,price_test,prev_post_test)).tocsr()  
X_test_tfidf.shape
```

```
Out[108]: (32775, 14897)
```

### GridSearchCV - Finding the best hyper parameter That maximum AUC value

#### With L1 Regularizer

```

In [109]: from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV

sgd = SGDClassifier(loss="hinge",penalty='l1',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000 ]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)
tuned_parameters = [{'alpha': Cs}]

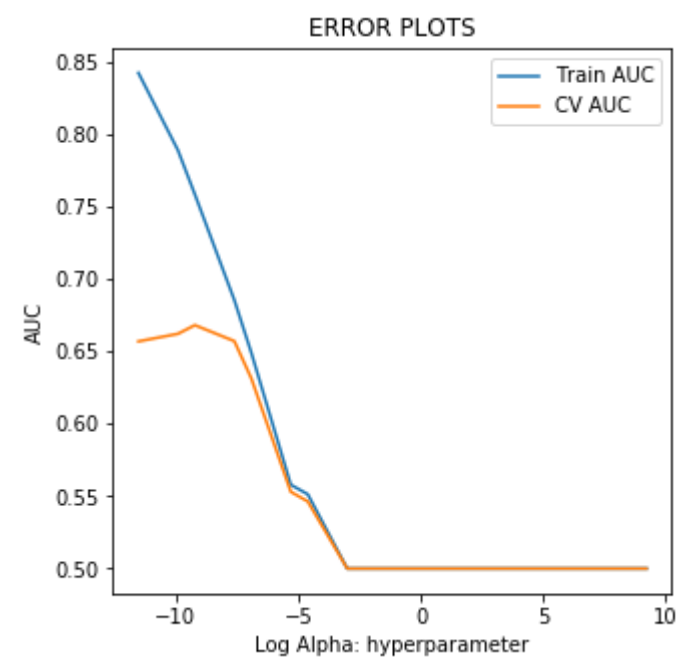
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [16,9]
plt.show()

```



```
In [110]: #http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']
x.add_column(column_names[0],Cs)
x.add_column(column_names[1],log_alphas)
x.add_column(column_names[2],train_auc)
x.add_column(column_names[3],cv_auc)
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.8421733129630317	0.6568444318362728
5e-05	-9.903487552536127	0.7888540045681257	0.6620571961337651
0.0001	-9.210340371976182	0.7582731797707817	0.6680349645066885
0.0005	-7.600902459542082	0.6853597151826659	0.6570192388739426
0.001	-6.907755278982137	0.6489923458409169	0.6312874557793773
0.005	-5.298317366548036	0.5578835606056799	0.5529866896312495
0.01	-4.605170185988091	0.5512930508361578	0.5463214642443495
0.05	-2.995732273553991	0.5	0.5
0.1	-2.3025850929940455	0.5	0.5
0.5	-0.6931471805599453	0.5	0.5
1	0.0	0.5	0.5
5	1.6094379124341003	0.5	0.5
10	2.302585092994046	0.5	0.5
50	3.912023005428146	0.5	0.5
100	4.605170185988092	0.5	0.5
500	6.214608098422191	0.5	0.5
1000	6.907755278982137	0.5	0.5
2500	7.824046010856292	0.5	0.5
5000	8.517193191416238	0.5	0.5
10000	9.210340371976184	0.5	0.5

With L2 Regularizer

```
In [111]: from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV

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

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000 ]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)
tuned_parameters = [{'alpha': Cs}]

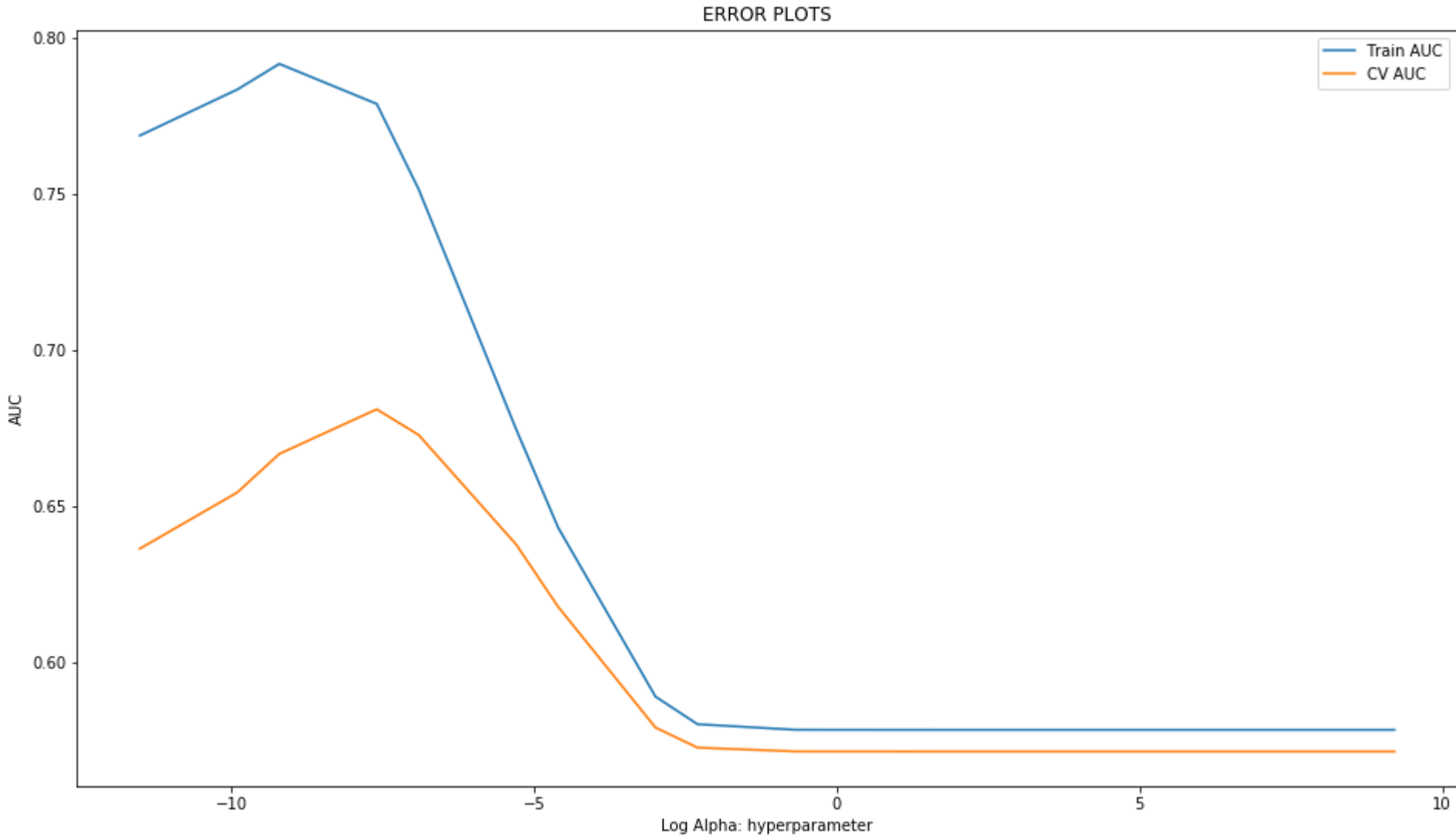
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [16,9]
plt.show()
```





In [112]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.768617298881629	0.63638274877399
5e-05	-9.903487552536127	0.7833537140392721	0.6543765060020961
0.0001	-9.210340371976182	0.7915541856025307	0.6667079468353564
0.0005	-7.600902459542082	0.778738430733302	0.6809868214749007
0.001	-6.907755278982137	0.7513742451237814	0.6727578966553305
0.005	-5.298317366548036	0.6746102924910219	0.6377275380849647
0.01	-4.605170185988091	0.643156482943018	0.6177848113912728
0.05	-2.995732273553991	0.5890017850262187	0.5791015995080276
0.1	-2.3025850929940455	0.5801744386417115	0.572677675313687
0.5	-0.6931471805599453	0.578398268363215	0.5714654335508755
1	0.0	0.5783795297678326	0.5714471577055937
5	1.6094379124341003	0.5783731946771192	0.5714380073135858
10	2.302585092994046	0.5783735614556755	0.5714380072214376
50	3.912023005428146	0.5783731946771192	0.5714380073135858
100	4.605170185988092	0.5783731946771192	0.5714380073135858
500	6.214608098422191	0.5783731946771192	0.5714380073135858
1000	6.907755278982137	0.5783731946771192	0.5714380073135858
2500	7.824046010856292	0.5783731946771192	0.5714380073135858
5000	8.517193191416238	0.5783731946771192	0.5714380073135858
10000	9.210340371976184	0.5783731946771192	0.5714380073135858

Using Best alpha Value – Training the Model

In [113]:

```
#Taking the Optimal hypermeter from L2 Regularizer for FPR, TPR plot and confusion matrix  
best_alpha_tfidf = 0.001
```

In [114]: [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.metrics.roc\\_curve](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve)  
**from sklearn.metrics import roc\_curve, auc**

```
SGD = SGDClassifier(loss="hinge",penalty='l2',alpha= best_alpha_tfidf,class_weight = 'balanced')
SGD.fit(X_train_tfidf, y_train)
```

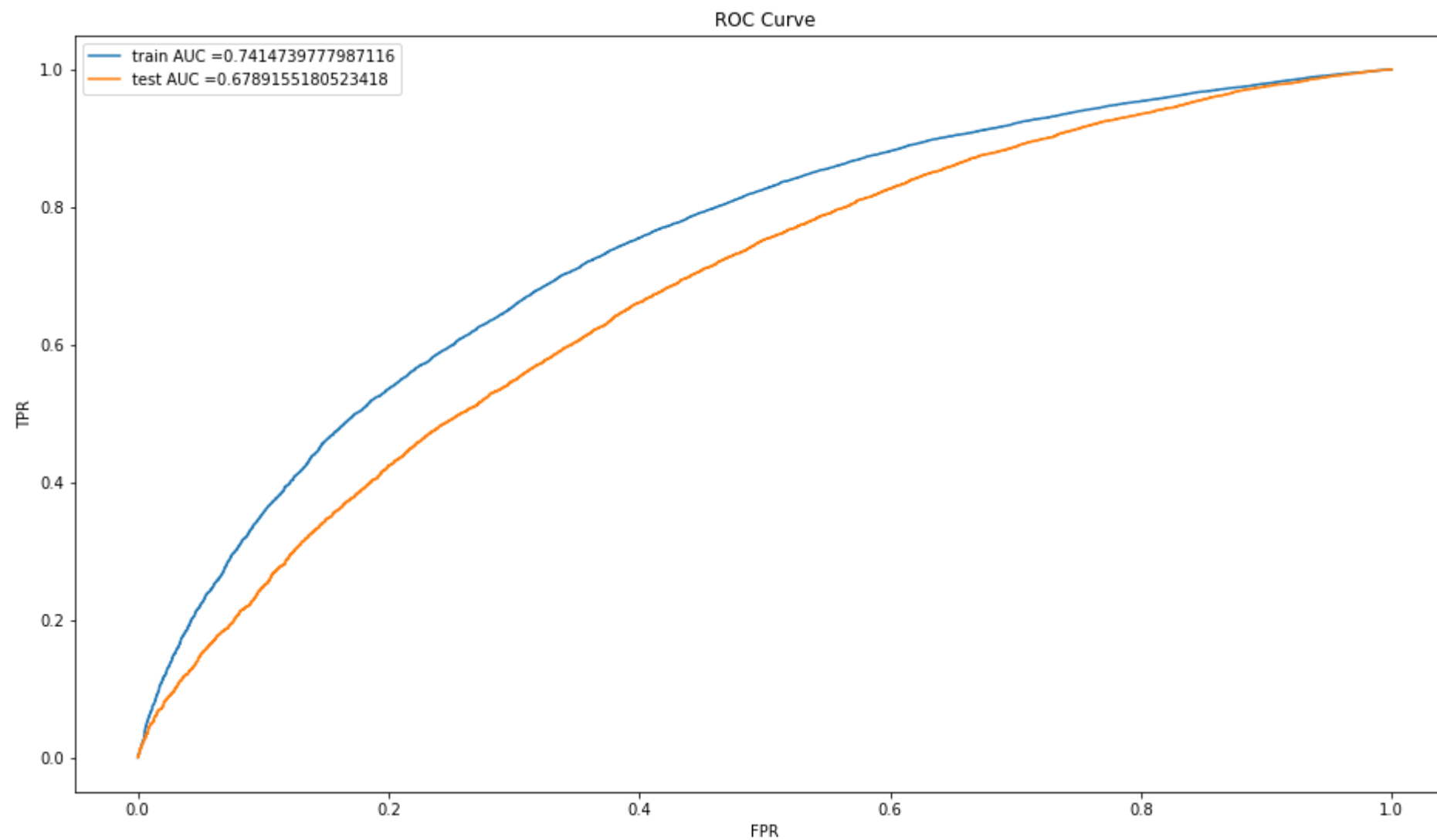
*# roc\_auc\_score(y\_true, y\_score) the 2nd parameter should be probability estimates of the positive class  
 # not the predicted outputs*

```
y_train_pred_tfidf = SGD.decision_function(X_train_tfidf)
y_test_pred_tfidf = SGD.decision_function(X_test_tfidf)
```

```
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_tfidf)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_tfidf)
```

```
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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")
```

```
plt.show()
```



## Confusion Matrix

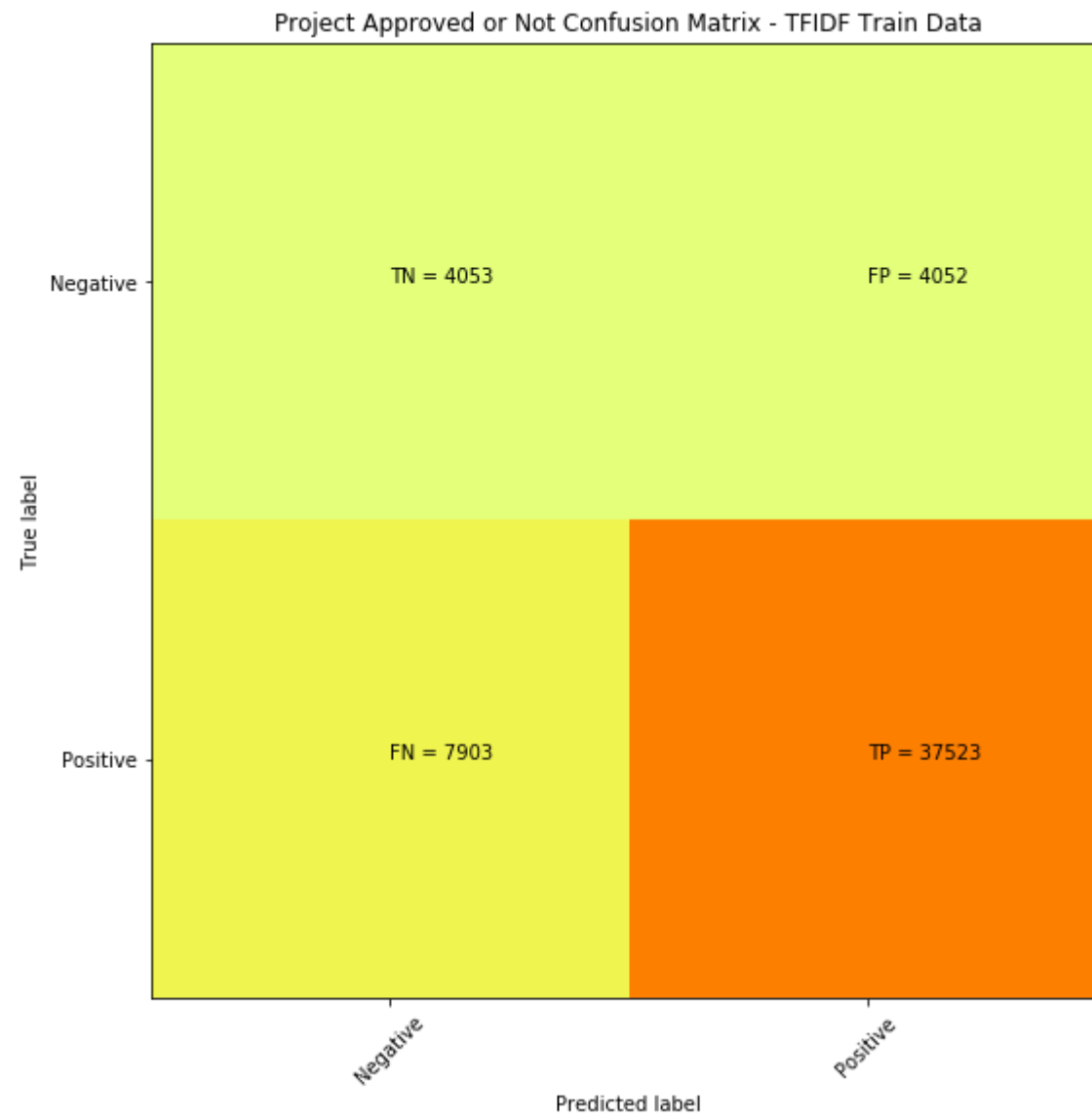
### Train confusion matrix

```
In [115]: from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
tfidf_train_confusion_matrix = confusion_matrix(y_train, predict(y_train_pred_tfidf, tr_thresholds, train_fpr, train_fpr))
print(tfidf_train_confusion_matrix)
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.2499999961943051 for threshold -0.356  
[[ 4053 4052]  
 [ 7903 37523]]

In [116]: <http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/>

```
plt.clf()
plt.imshow(tfidf_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - TFIDF Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(tfidf_train_confusion_matrix[i][j]))
plt.show()
```



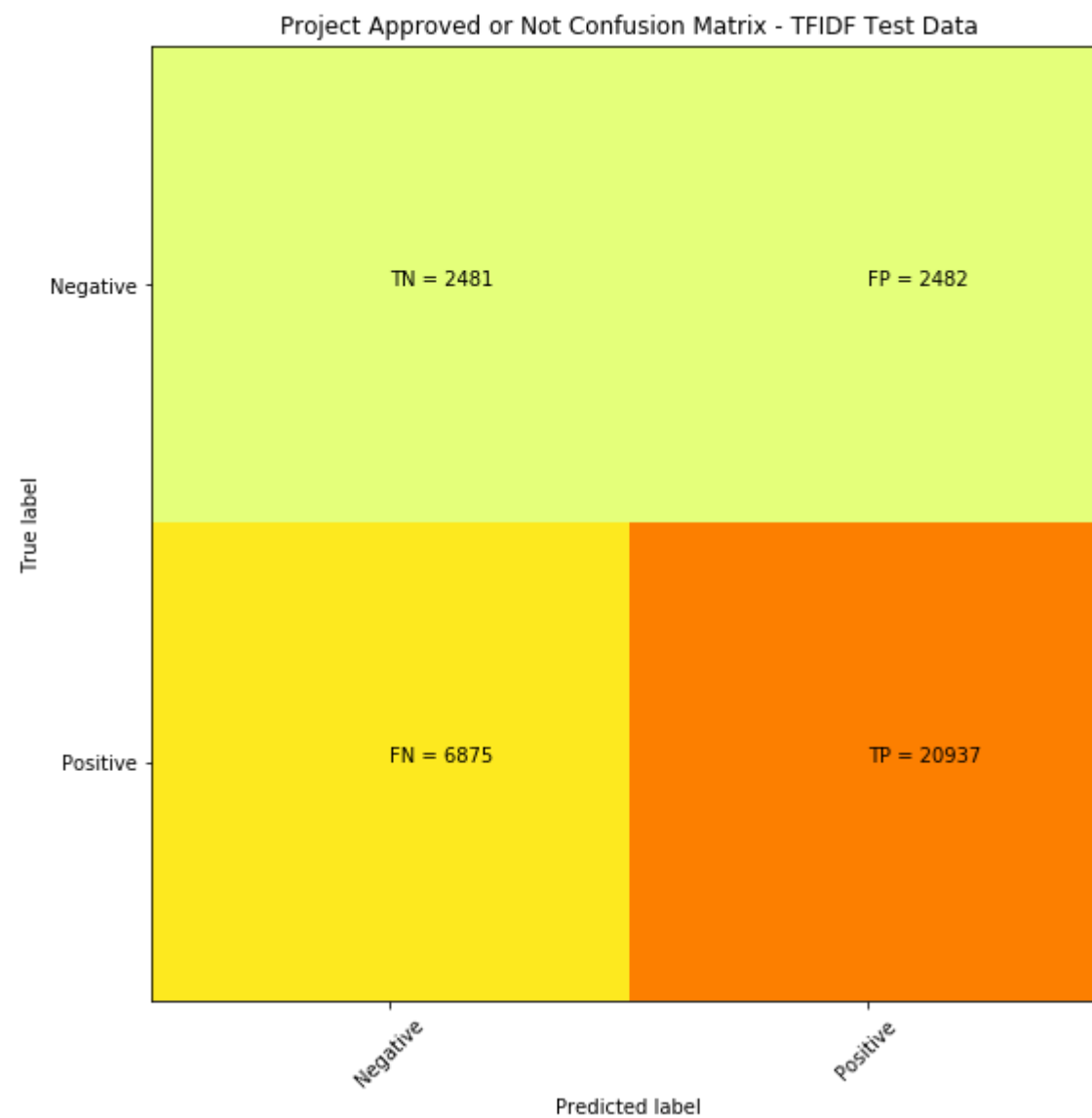
**Test confusion matrix**

```
In [117]: print("Test confusion matrix")
tfidf_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_tfidf, te_thresholds, test_fpr, test_fpr))
print(tfidf_test_confusion_matrix)
```

```
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24999998985034083 for threshold -0.214
[[ 2481  2482]
 [ 6875 20937]]
```

```
In [118]: ##http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/
```

```
plt.clf()
plt.imshow(tfidf_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - TFIDF Test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(tfidf_test_confusion_matrix[i][j]))
plt.show()
```



### 2.4.3 Applying SGD Classifier (with Loss = 'hinge') on AVG W2V, SET 3

```
In [119]: X_train_avg_w2v = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train ,Teacher_Prefix_one_hot_train,project_grade_category_one_hot_train,avg_
w2v_essays_train,avg_w2v_title_train,price_train,prev_post_train)).tocsr()
X_train_avg_w2v.shape
```

```
Out[119]: (53531, 701)
```

```
In [120]: X_cv_avg_w2v = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv ,Teacher_Prefix_one_hot_cv,project_grade_category_one_hot_cv,avg_w2v_essays_cv,avg_
w2v_title_cv,price_cv,prev_post_cv)).tocsr()
X_cv_avg_w2v.shape
```

```
Out[120]: (22942, 701)
```

```
In [121]: X_test_avg_w2v = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test ,Teacher_Prefix_one_hot_test,project_grade_category_one_hot_test,avg_w2v_es
says_test,avg_w2v_title_test,price_test,prev_post_test)).tocsr()
X_test_avg_w2v.shape
```

```
Out[121]: (32775, 701)
```

#### GridSearchCV - Finding the best hyper parameter That maximum AUC value

#### With L1 Regularizer

```
In [122]: sgd = SGDClassifier(loss="hinge",penalty='l1',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]

log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

tuned_parameters = [{'alpha': Cs}]

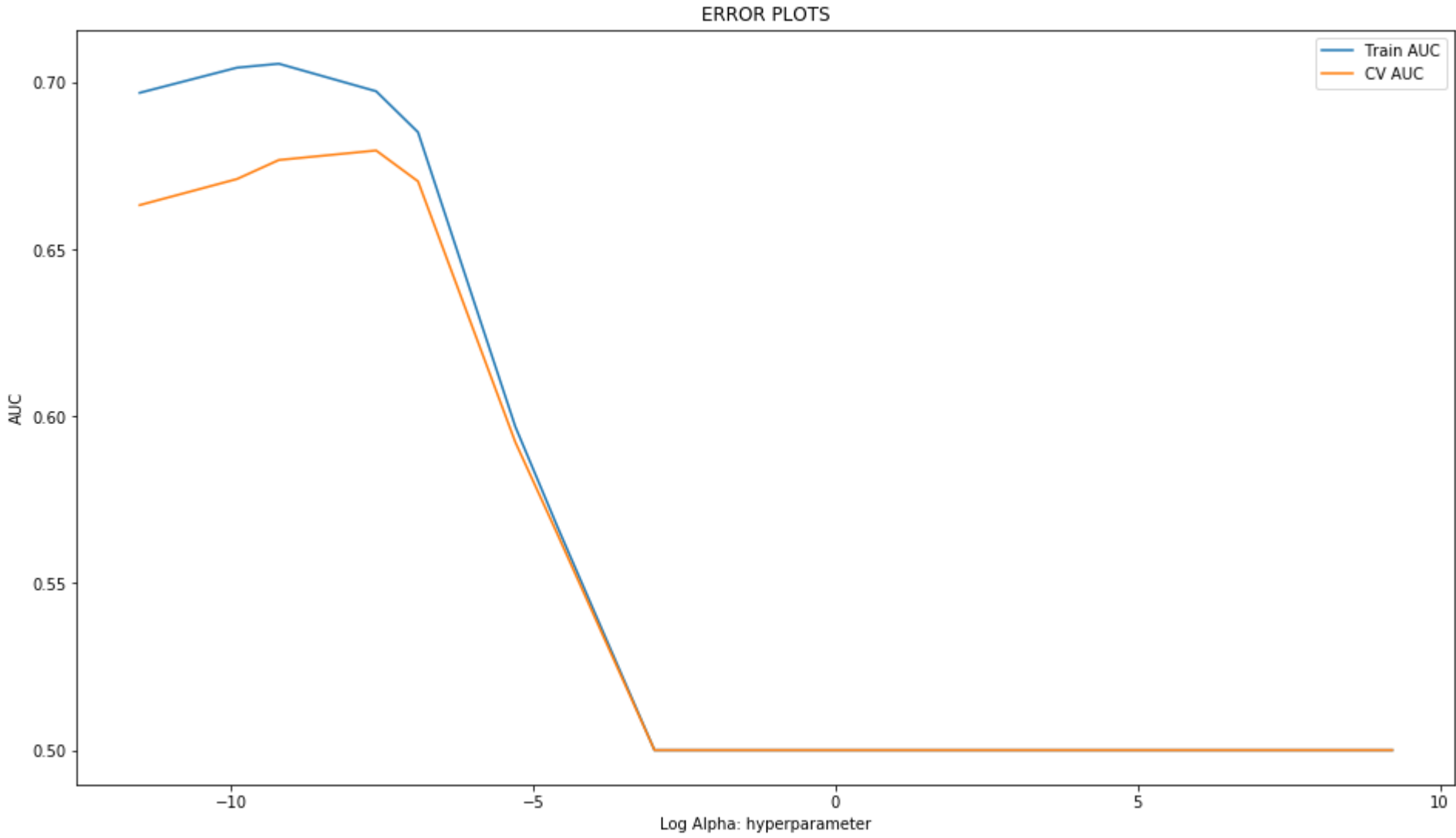
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_avg_w2v, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [16,9]
plt.show()
```





In [123]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.6967933843372721	0.6631751388663257
5e-05	-9.903487552536127	0.7043231461033933	0.6709812700285693
0.0001	-9.210340371976182	0.7054965151751547	0.6766773400894286
0.0005	-7.600902459542082	0.6972606460186701	0.6795229946510464
0.001	-6.907755278982137	0.6849825938825539	0.6702847681608023
0.005	-5.298317366548036	0.5968558517543366	0.5922072609006755
0.01	-4.605170185988091	0.5670282986476558	0.5648060308747387
0.05	-2.995732273553991	0.5	0.5
0.1	-2.3025850929940455	0.5	0.5
0.5	-0.6931471805599453	0.5	0.5
1	0.0	0.5	0.5
5	1.6094379124341003	0.5	0.5
10	2.302585092994046	0.5	0.5
50	3.912023005428146	0.5	0.5
100	4.605170185988092	0.5	0.5
500	6.214608098422191	0.5	0.5
1000	6.907755278982137	0.5	0.5
2500	7.824046010856292	0.5	0.5
5000	8.517193191416238	0.5	0.5
10000	9.210340371976184	0.5	0.5

With L2 Regularizer

```
In [124]: sgd = SGDClassifier(loss="hinge",penalty='l2',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]

log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

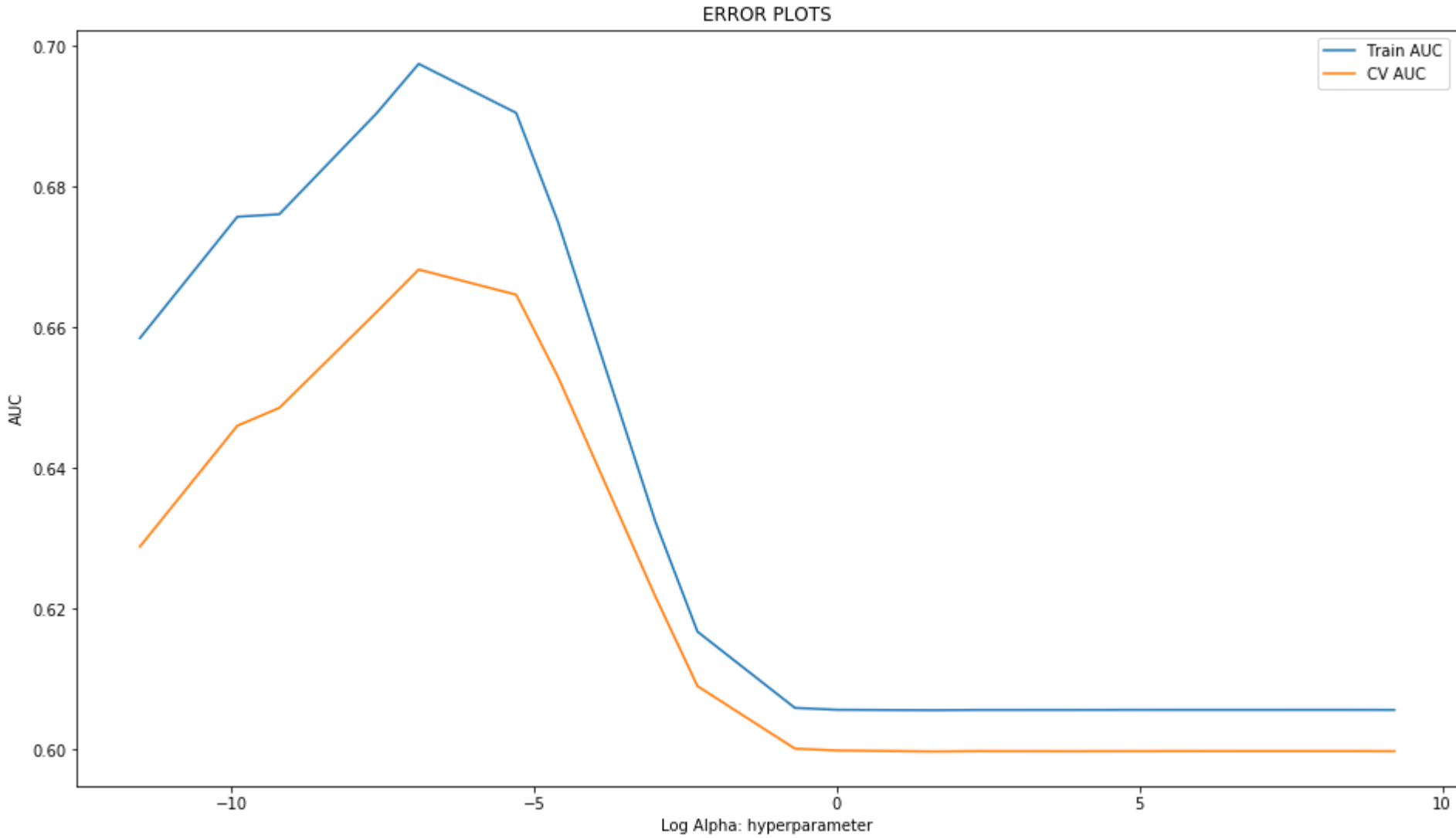
tuned_parameters = [{'alpha': Cs}]

clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_avg_w2v, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [15,5]
plt.show()
```



In [125]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.6584782624632616	0.6288530184778335
5e-05	-9.903487552536127	0.6757325017170218	0.646036990772455
0.0001	-9.210340371976182	0.676087381484403	0.6485680201339994
0.0005	-7.600902459542082	0.6904546610057037	0.6622151669522631
0.001	-6.907755278982137	0.6974620255438543	0.6682338579674761
0.005	-5.298317366548036	0.6904966862044253	0.664659125844879
0.01	-4.605170185988091	0.674956635859825	0.6529368893428131
0.05	-2.995732273553991	0.6323338158628165	0.6217104346017245
0.1	-2.3025850929940455	0.6167654505040124	0.6090230089277895
0.5	-0.6931471805599453	0.60592871435106	0.6001412839158565
1	0.0	0.6056623545992886	0.5998772466578578
5	1.6094379124341003	0.6056007416664545	0.5997322300923305
10	2.302585092994046	0.6056408734872613	0.599781778833833
50	3.912023005428146	0.6056421678750114	0.5997689009143903
100	4.605170185988092	0.6056480246934931	0.599776038026744
500	6.214608098422191	0.6056523821816464	0.5997864911890326
1000	6.907755278982137	0.6056513046316162	0.5997846091301583
2500	7.824046010856292	0.6056524820877934	0.5997857170952513
5000	8.517193191416238	0.605652294567221	0.5997861897971972
10000	9.210340371976184	0.6056448265111912	0.5997711900635405

Using Best alpha Value – Training the Model

In [126]:

```
#Taking the Optimal hypermeter from L2 Regularizer for FPR, TPR plot and confusion matrix  
best_alpha_avg_w2v = 0.001
```

In [127]: [#https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.metrics.roc\\_curve](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve)  
**from sklearn.metrics import roc\_curve, auc**

```
SGD = SGDClassifier(loss="hinge",penalty='l2',alpha= best_alpha_avg_w2v,class_weight = 'balanced')
SGD.fit(X_train_avg_w2v, y_train)
```

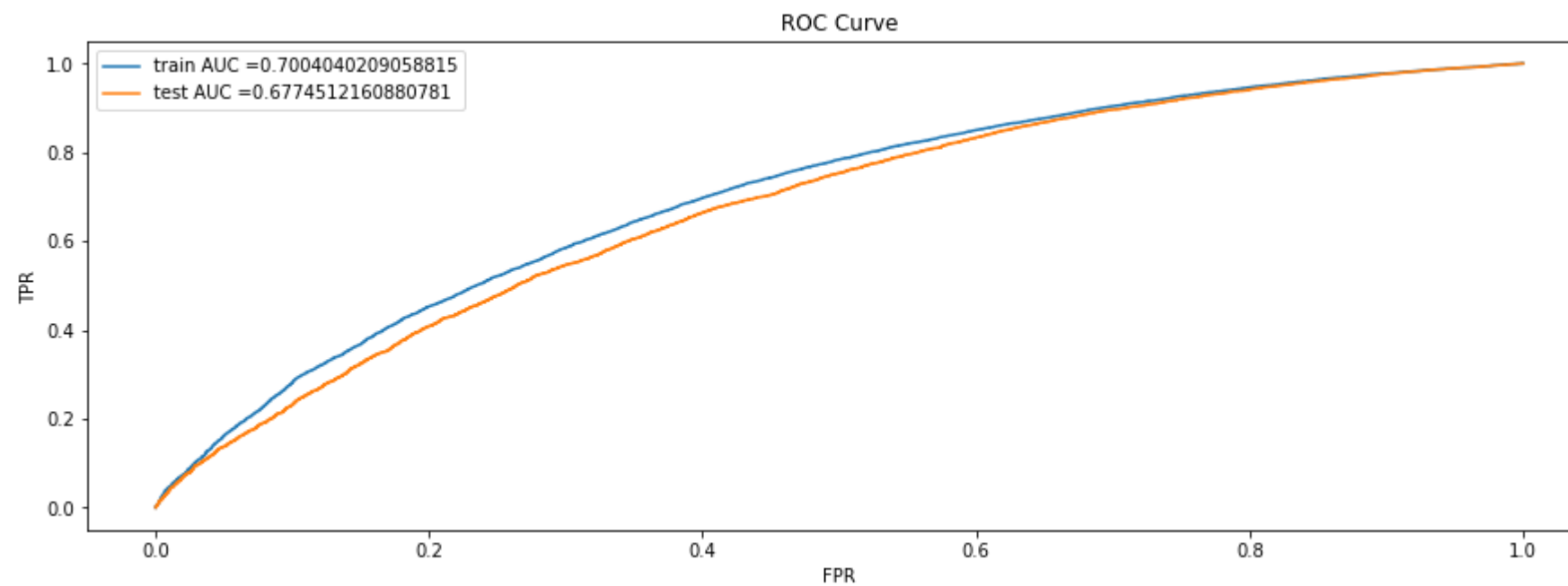
```
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
```

```
y_train_pred_avg_w2v = SGD.decision_function(X_train_avg_w2v)
y_test_pred_avg_w2v = SGD.decision_function(X_test_avg_w2v)
```

```
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_avg_w2v)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_avg_w2v)
```

```
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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")
```

```
plt.show()
```



## Confusion Matrix

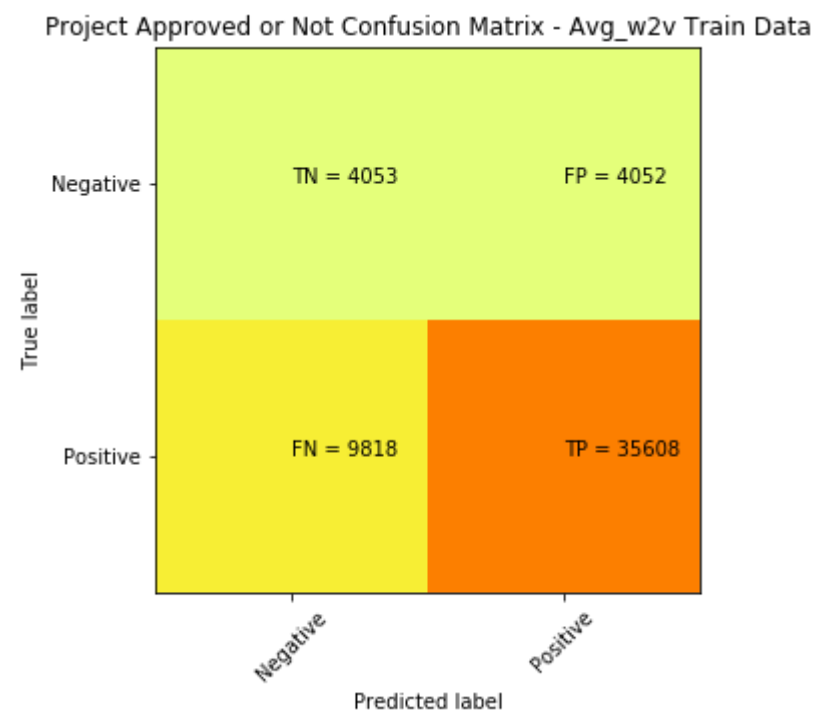
### Train confusion matrix

```
In [128]: from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
avg_w2v_train_confusion_matrix = confusion_matrix(y_train, predict(y_train_pred_avg_w2v, tr_thresholds, train_fpr, train_fpr))
print(avg_w2v_train_confusion_matrix)

#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(avg_w2v_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Avg_w2v Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(avg_w2v_train_confusion_matrix[i][j]))
plt.show()
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.2499999961943051 for threshold -1.009  
[[ 4053 4052]  
[ 9818 35608]]



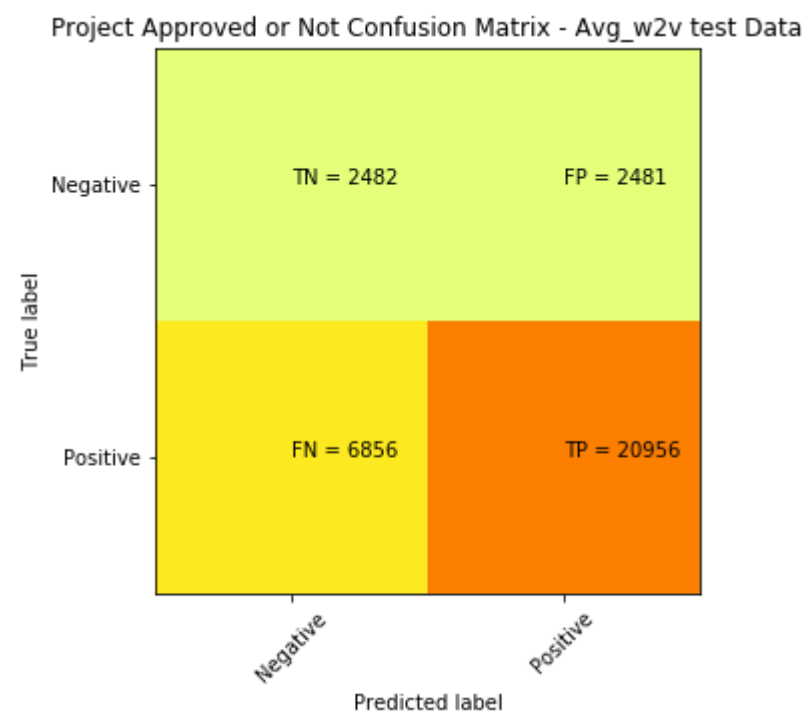
**Test confusion matrix**

```
In [129]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
avg_w2v_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_avg_w2v, te_thresholds, test_fpr, test_fpr))
print(avg_w2v_test_confusion_matrix)

#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(avg_w2v_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Avg_w2v test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(avg_w2v_test_confusion_matrix[i][j]))
plt.show()
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999998985034083 for threshold -0.943  
[[ 2482 2481]  
[ 6856 20956]]



#### 2.4.4 Applying SGD Classifier (with Loss = 'hinge') on TFIDF W2V, SET 4

```
In [130]: X_train_tfidf_w2v = hstack((categories_one_hot_train, sub_categories_one_hot_train, school_state_one_hot_train ,Teacher_Prefix_one_hot_train,project_grade_category_one_hot_train,tfidf_w2v_essays_train,tfidf_w2v_title_train,price_train,prev_post_train)).tocsr()
X_train_tfidf_w2v.shape
```

Out[130]: (53531, 701)



```
In [131]: X_cv_tfidf_w2v = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv ,Teacher_Prefix_one_hot_cv,project_grade_category_one_hot_cv,tfidf_w2v_essays_cv,
tfidf_w2v_title_cv,price_cv,prev_post_cv)).tocsr()
X_cv_tfidf_w2v.shape
```

Out[131]: (22942, 701)

```
In [132]: X_test_tfidf_w2v = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test ,Teacher_Prefix_one_hot_test,project_grade_category_one_hot_test,tfidf_w2
v_essays_test,tfidf_w2v_title_test,price_test,prev_post_test)).tocsr()
X_test_tfidf_w2v.shape
```

Out[132]: (32775, 701)

### GridSearchCV - Finding the best hyper parameter That maximum AUC value

#### With L1 Regularizer

```

In [133]: sgd = SGDClassifier(loss="hinge",penalty='l1',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)
tuned_parameters = [{'alpha': Cs}]

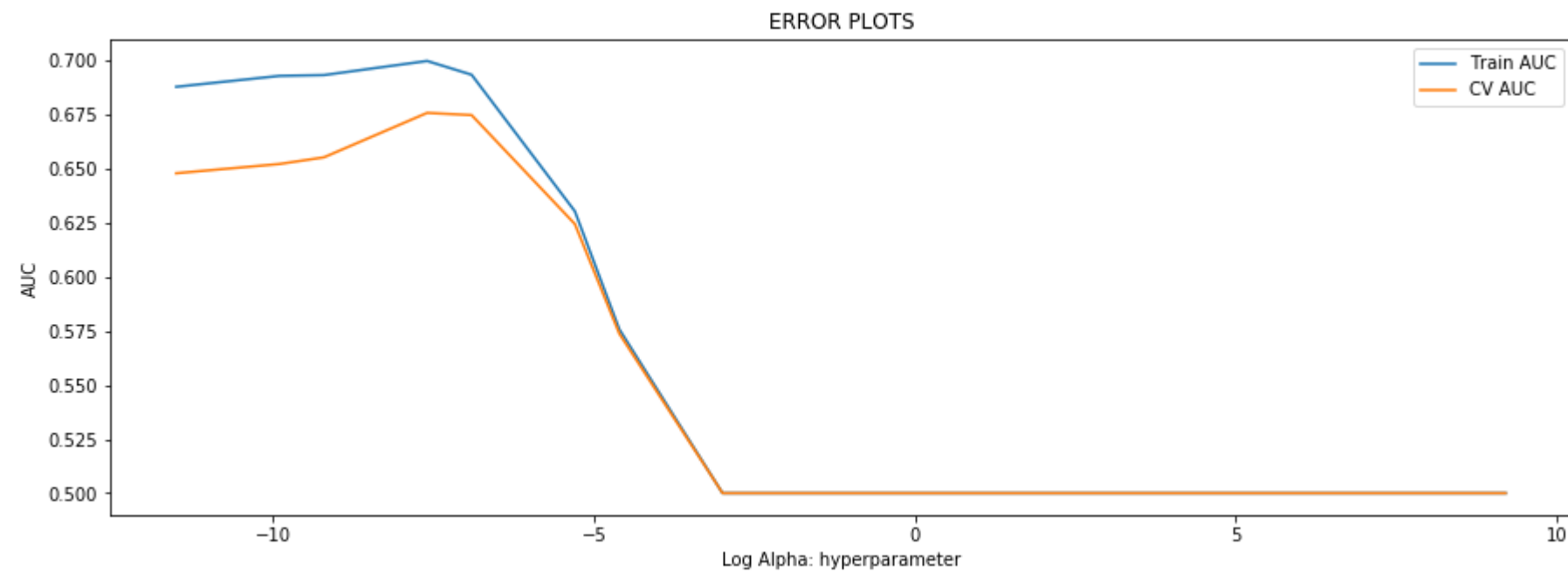
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf_w2v, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [16,9]
plt.show()

```



In [134]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.6878117960348722	0.6478490274529382
5e-05	-9.903487552536127	0.6928442224720367	0.6521286221586523
0.0001	-9.210340371976182	0.6932191165105279	0.655203001815302
0.0005	-7.600902459542082	0.6997417091665676	0.6757589895909295
0.001	-6.907755278982137	0.6933737059251537	0.6747284417297719
0.005	-5.298317366548036	0.630283029401252	0.6243882084210974
0.01	-4.605170185988091	0.5757207931905136	0.5737150166816652
0.05	-2.995732273553991	0.5	0.5
0.1	-2.3025850929940455	0.5	0.5
0.5	-0.6931471805599453	0.5	0.5
1	0.0	0.5	0.5
5	1.6094379124341003	0.5	0.5
10	2.302585092994046	0.5	0.5
50	3.912023005428146	0.5	0.5
100	4.605170185988092	0.5	0.5
500	6.214608098422191	0.5	0.5
1000	6.907755278982137	0.5	0.5
2500	7.824046010856292	0.5	0.5
5000	8.517193191416238	0.5	0.5
10000	9.210340371976184	0.5	0.5

With L2 Regularizer

```

In [135]: sgd = SGDClassifier(loss="hinge",penalty='l2',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)
tuned_parameters = [{'alpha': Cs}]

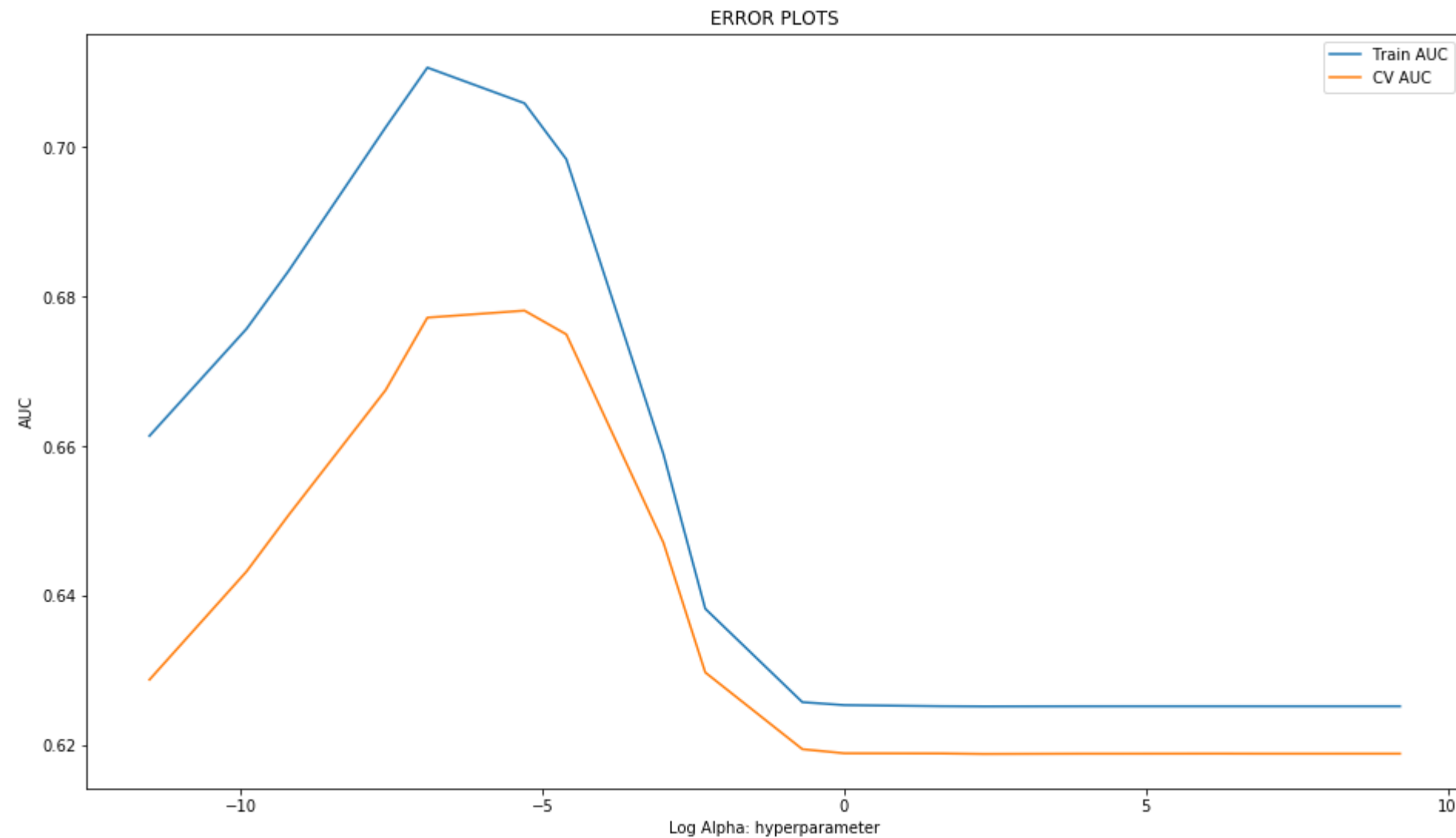
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf_w2v, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [15,5]
plt.show()

```



In [136]:

```
#http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']
x.add_column(column_names[0],Cs)
x.add_column(column_names[1],log_alphas)
x.add_column(column_names[2],train_auc)
x.add_column(column_names[3],cv_auc)
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.6613439819388263	0.6287363602650604
5e-05	-9.903487552536127	0.6756722492690646	0.6432222513280771
0.0001	-9.210340371976182	0.6833772733262767	0.6507196495568632
0.0005	-7.600902459542082	0.7025884470884612	0.6674549044897007
0.001	-6.907755278982137	0.710596490010524	0.6771601588355193
0.005	-5.298317366548036	0.7058292767019272	0.6780927830101462
0.01	-4.605170185988091	0.6983419066341211	0.6749212160215339
0.05	-2.995732273553991	0.6588956982541859	0.6470522311767265
0.1	-2.3025850929940455	0.6382276964329522	0.6297085311247125
0.5	-0.6931471805599453	0.625728486993104	0.6194405268551364
1	0.0	0.6253383970889832	0.6188959366593468
5	1.6094379124341003	0.6251876245314478	0.6188720573348495
10	2.302585092994046	0.6251605905219231	0.6188032213613509
50	3.912023005428146	0.6251772562834935	0.6188389154621348
100	4.605170185988092	0.6251799486418008	0.618839959284811
500	6.214608098422191	0.6251818247394532	0.6188505184322237
1000	6.907755278982137	0.6251803377181875	0.6188401874628014
2500	7.824046010856292	0.625179368118623	0.6188397474271884
5000	8.517193191416238	0.6251804926739527	0.6188423382830138
10000	9.210340371976184	0.6251805883489067	0.6188421183419975

Using Best alpha Value – Training the Model

In [137]:

```
#Taking the Optimal hypermeter from L2 Regularize for FPR, TPR plot and confusion matrix
best_alpha_tfidf_w2v = 0.001
```

```
In [138]: #https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

SGD = SGDClassifier(loss="hinge",penalty='l2',alpha= best_alpha_tfidf_w2v,class_weight = 'balanced')
SGD.fit(X_train_tfidf_w2v, y_train)

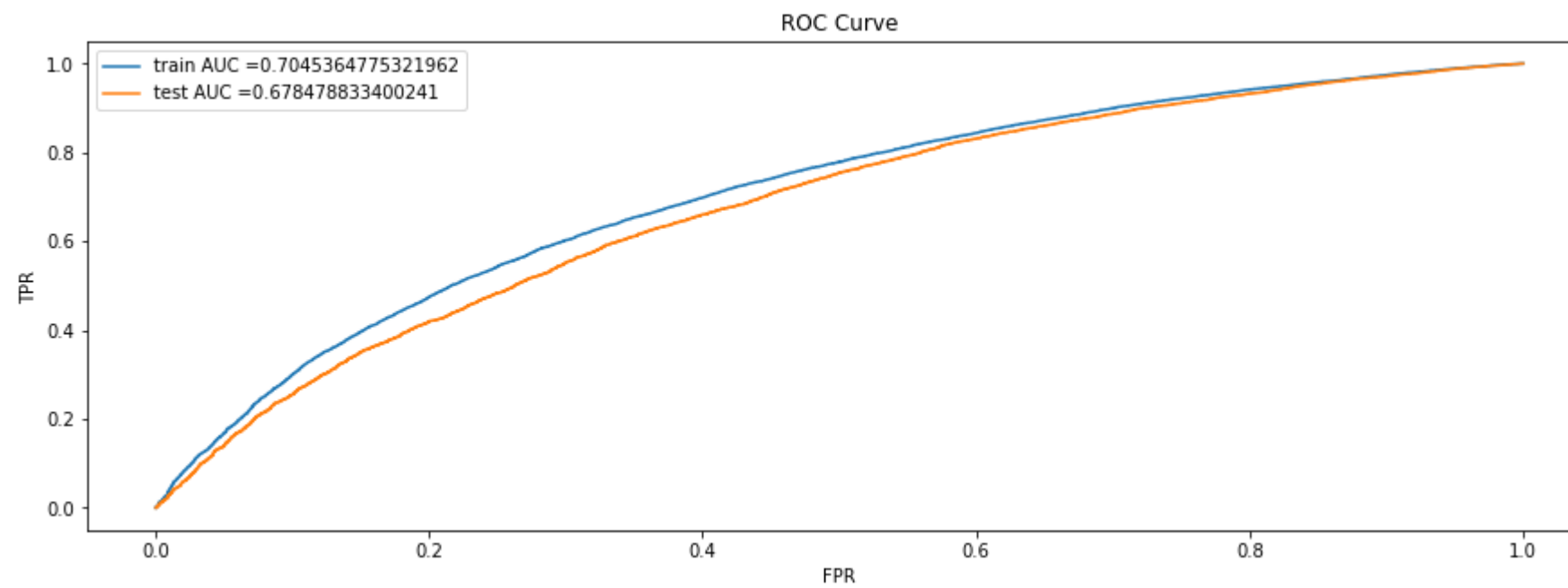
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_tfidf_w2v = SGD.decision_function(X_train_tfidf_w2v)
y_test_pred_tfidf_w2v = SGD.decision_function(X_test_tfidf_w2v)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_tfidf_w2v)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_tfidf_w2v)

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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")

plt.show()
```



## Confusion Matrix

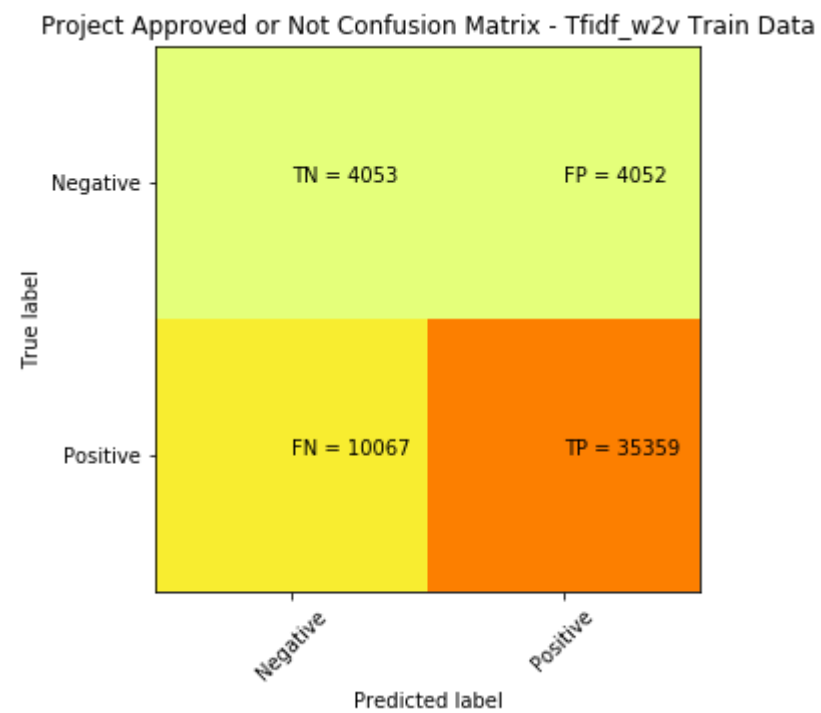
### Train confusion matrix

```
In [139]: from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
tfidf_w2v_train_confusion_matrix = confusion_matrix(y_train, predict(y_train_pred_tfidf_w2v, tr_thresholds, train_fpr, train_fpr))
print(tfidf_w2v_train_confusion_matrix)

#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(tfidf_w2v_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Tfidf_w2v Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(tfidf_w2v_train_confusion_matrix[i][j]))
plt.show()
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.2499999961943051 for threshold -0.495  
[[ 4053 4052]  
[10067 35359]]



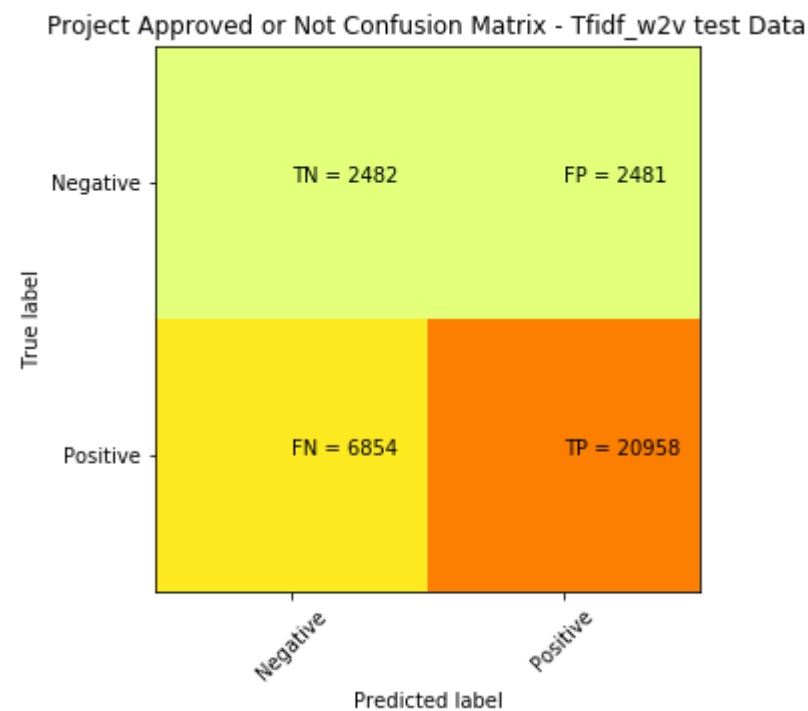
**Test confusion matrix**

```
In [140]: from sklearn.metrics import confusion_matrix
print("Test confusion matrix")
tfidf_w2v_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_tfidf_w2v, te_thresholds, test_fpr, test_fpr))
print(tfidf_w2v_test_confusion_matrix)

#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(avg_w2v_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Tfidf_w2v test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(tfidf_w2v_test_confusion_matrix[i][j]))
plt.show()
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999998985034083 for threshold -0.429  
[[ 2482 2481]  
[ 6854 20958]]



## 2.5 Categorical , Numerical and Truncated Text TFIDF SVD



```
In [141]: from sklearn.decomposition import TruncatedSVD
## https://medium.com/@jonathan_hui/machine-learning-singular-value-decomposition-svd-principal-component-analysis-pca-1d45e885e491
## https://galaxydatatech.com/2018/07/15/singular-value-decomposition/

svd = TruncatedSVD(n_components= 3000, n_iter=5, random_state=42)
svd.fit(tfidf_essays_train)
```

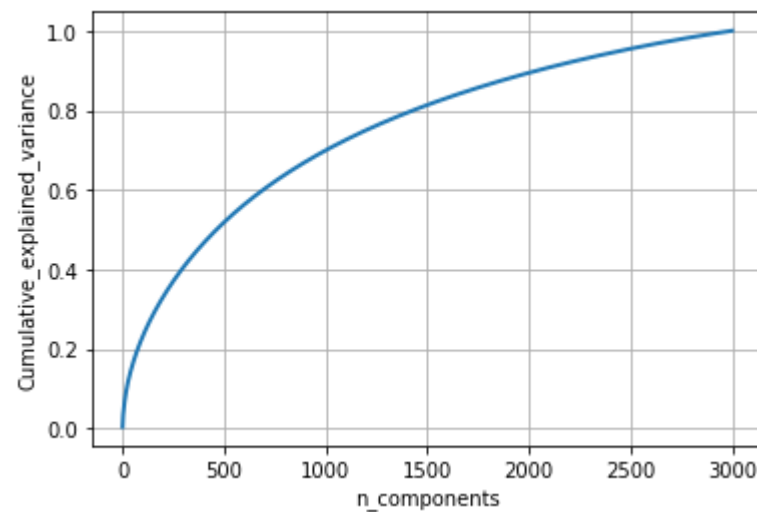
```
Out[141]: TruncatedSVD(algorithm='randomized', n_components=3000, n_iter=5,
      random_state=42, tol=0.0)
```

```
In [142]: percentage_var_explained = svd.explained_variance_ / np.sum(svd.explained_variance_);
```

```
In [143]: cum_var_explained = np.cumsum(percentage_var_explained)
```

```
In [144]: # Plot the Truncated SVD spectrum
plt.figure(1, figsize=(6, 4))

plt.clf()
plt.plot(cum_var_explained, linewidth=2)
plt.axis('tight')
plt.grid()
plt.xlabel('n_components')
plt.ylabel('Cumulative_explained_variance')
plt.show()
```



**100% Variance explained with 3000 no's of component**

**Train data - Truncated Text TFIDF SVD**

```
In [145]: svd.fit(tfidf_essays_train)
svd_tfidf_train = svd.transform(tfidf_essays_train)
svd_tfidf_train.shape
```

```
Out[145]: (53531, 3000)
```

**CV data - Truncated Text TFIDF SVD**

```
In [146]: svd_tfidf_cv = svd.transform(tfidf_essays_cv)
          svd_tfidf_cv.shape
```

```
Out[146]: (22942, 3000)
```

### Test data - Truncated Text TFIDF SVD

```
In [147]: svd_tfidf_test = svd.transform(tfidf_essays_test)
          svd_tfidf_test.shape
```

```
Out[147]: (32775, 3000)
```

```
In [149]: X_train_cat_num_freatures = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train ,Teacher_Prefix_one_hot_train,project_grade_category_one_hot_train,price_train,prev_post_train,Quantity_train,title_word_count_train,essay_word_count_train,Positive_SC_Essay_train,Neutral_SC_Essay_train,Negative_SC_Essay_train,Compound_SC_Essay_train,svd_tfidf_train)).tocsr()
          X_train_cat_num_freatures.shape
```

```
Out[149]: (53531, 3108)
```

```
In [150]: X_cv_cat_num_freatures = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv ,Teacher_Prefix_one_hot_cv,project_grade_category_one_hot_cv,price_cv,prev_post_cv,Quantity_cv,title_word_count_cv,essay_word_count_cv,Positive_SC_Essay_cv,Neutral_SC_Essay_cv,Negative_SC_Essay_cv,Compound_SC_Essay_cv,svd_tfidf_cv)).tocsr()
          X_cv_cat_num_freatures.shape
```

```
Out[150]: (22942, 3108)
```

```
In [151]: X_test_cat_num_freatures = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test ,Teacher_Prefix_one_hot_test,project_grade_category_one_hot_test,price_test,prev_post_test,Quantity_test,title_word_count_test,essay_word_count_test,Positive_SC_Essay_test,Neutral_SC_Essay_test,Negative_SC_Essay_test,Compound_SC_Essay_test,svd_tfidf_test)).tocsr()
          X_test_cat_num_freatures.shape
```

```
Out[151]: (32775, 3108)
```

### GridSearchCV - Finding the best hyper parameter That maximum AUC value

#### With L1 Regularizer

```

In [152]: sgd = SGDClassifier(loss="hinge",penalty='l1',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

tuned_parameters = [{'alpha': Cs}]

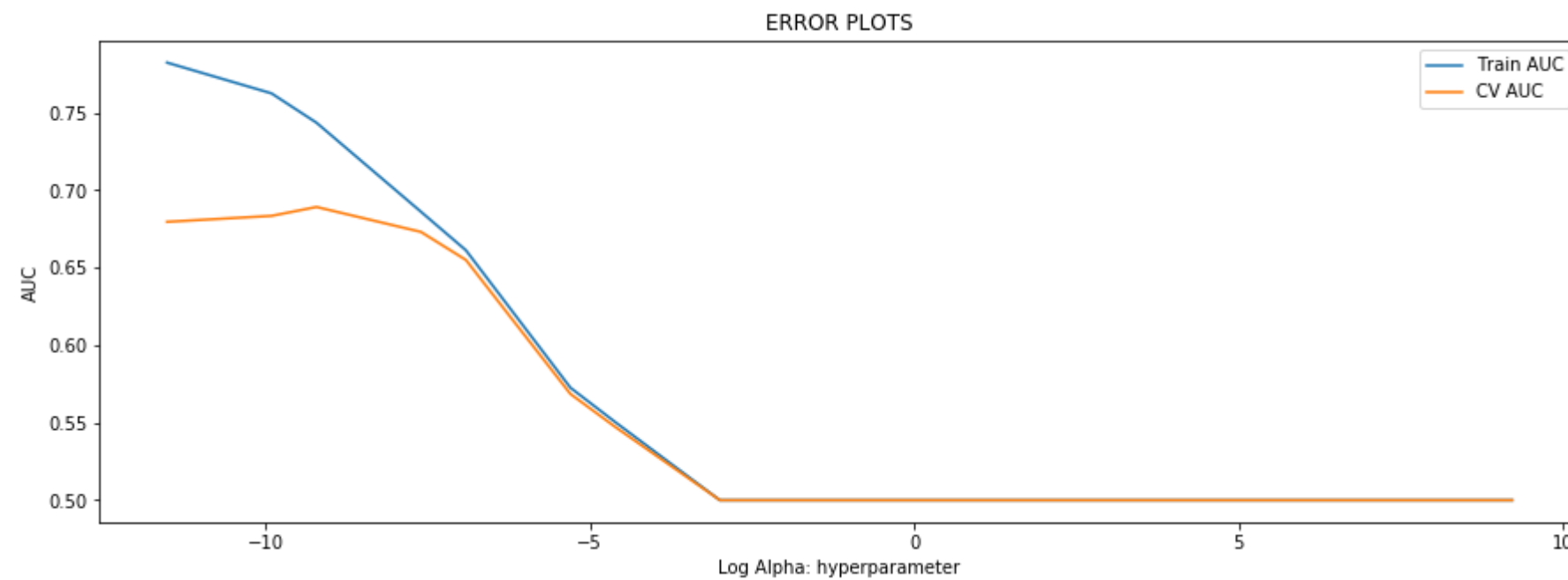
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_cat_num_freatures, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [15,5]
plt.show()

```



In [153]:

```
#http://zetcode.com/python/prettytable/  
from prettytable import PrettyTable  
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable  
x = PrettyTable()  
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']  
x.add_column(column_names[0],Cs)  
x.add_column(column_names[1],log_alphas)  
x.add_column(column_names[2],train_auc)  
x.add_column(column_names[3],cv_auc)  
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.7823811972458339	0.6796173945730339
5e-05	-9.903487552536127	0.7625481963929746	0.6834591769441515
0.0001	-9.210340371976182	0.743604661604138	0.6892040067316438
0.0005	-7.600902459542082	0.6861367669467402	0.6731130913235798
0.001	-6.907755278982137	0.6612584152346356	0.6550069318885746
0.005	-5.298317366548036	0.572594312845497	0.568568097906158
0.01	-4.605170185988091	0.5503298401486485	0.5471876430002948
0.05	-2.995732273553991	0.5	0.5
0.1	-2.3025850929940455	0.5	0.5
0.5	-0.6931471805599453	0.5	0.5
1	0.0	0.5	0.5
5	1.6094379124341003	0.5	0.5
10	2.302585092994046	0.5	0.5
50	3.912023005428146	0.5	0.5
100	4.605170185988092	0.5	0.5
500	6.214608098422191	0.5	0.5
1000	6.907755278982137	0.5	0.5
2500	7.824046010856292	0.5	0.5
5000	8.517193191416238	0.5	0.5
10000	9.210340371976184	0.5	0.5

With L2 Regularizer

```

In [154]: sgd = SGDClassifier(loss="hinge",penalty='l2',class_weight = 'balanced')

Cs = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 500, 1000, 2500, 5000, 10000]
log_alphas = []
for P in Cs :
    T = math.log(P)
    log_alphas.append(T)

tuned_parameters = [{'alpha': Cs}]

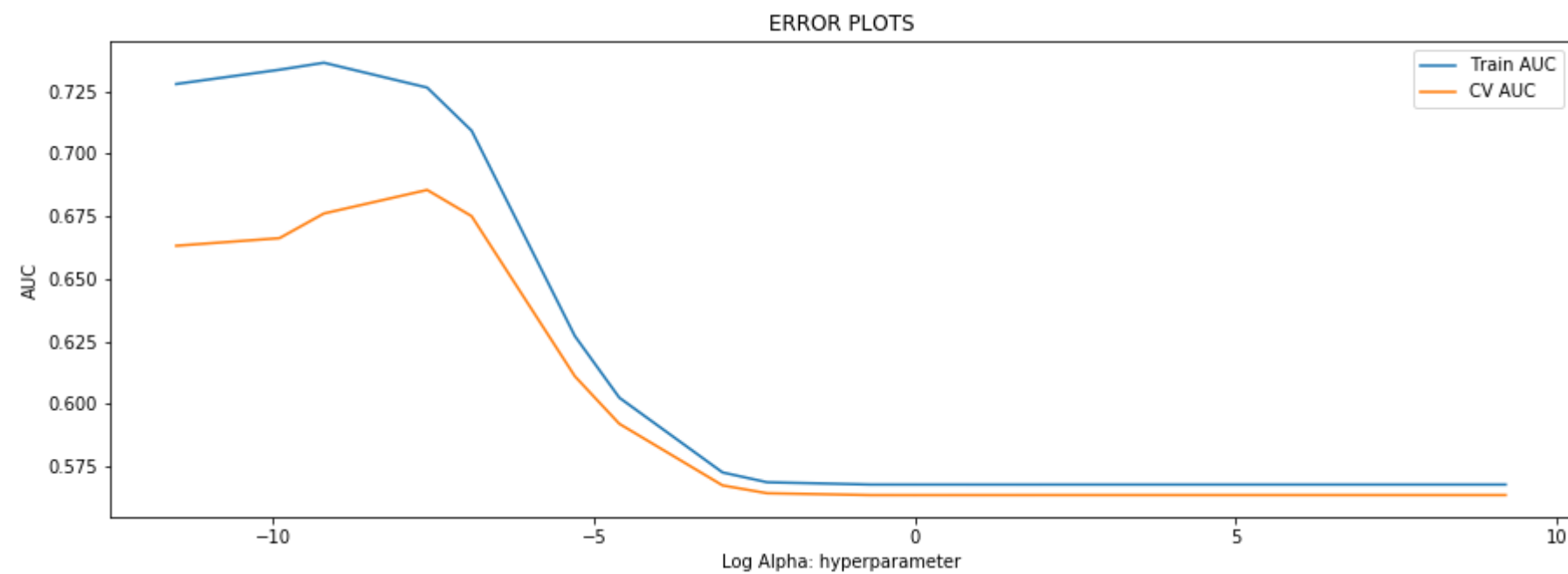
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_cat_num_freatures, y_train)

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

plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')

plt.legend()
plt.xlabel("Log Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.rcParams["figure.figsize"] = [15,5]
plt.show()

```



```
In [155]: #http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
column_names = ['alphas', 'log_alphas', 'train_auc', 'cv_auc']
x.add_column(column_names[0],Cs)
x.add_column(column_names[1],log_alphas)
x.add_column(column_names[2],train_auc)
x.add_column(column_names[3],cv_auc)
print(x)
```

alphas	log_alphas	train_auc	cv_auc
1e-05	-11.512925464970229	0.7278302427851875	0.6631775773542985
5e-05	-9.903487552536127	0.7335325702983807	0.6662057479036522
0.0001	-9.210340371976182	0.7363271711264558	0.6760734794638164
0.0005	-7.600902459542082	0.7264016032422164	0.685488786397706
0.001	-6.907755278982137	0.7091715503729019	0.6750091302555048
0.005	-5.298317366548036	0.6270893826987011	0.6110151754515842
0.01	-4.605170185988091	0.6024555145028612	0.5920271969426875
0.05	-2.995732273553991	0.5725772251790371	0.5673952608730214
0.1	-2.3025850929940455	0.5686549393745495	0.5642635128401746
0.5	-0.6931471805599453	0.5677353229634879	0.5635209916911124
1	0.0	0.5677303029614793	0.5635153777328795
5	1.6094379124341003	0.567726523585357	0.5635236809503644
10	2.302585092994046	0.5677318364056148	0.5635310540937316
50	3.912023005428146	0.5677305979496095	0.5635283816386261
100	4.605170185988092	0.5677311948215388	0.5635293593914507
500	6.214608098422191	0.5677307710496692	0.5635289764080142
1000	6.907755278982137	0.5677310623027084	0.5635284956943455
2500	7.824046010856292	0.5677311132968048	0.5635290578818704
5000	8.517193191416238	0.5677311132968048	0.5635290578818704
10000	9.210340371976184	0.5677311132968048	0.5635290578818704

Using Best aplha Value – Training the Model

```
In [156]: #Taking the Optimal hypermeter from L2 Regularize for FPR, TPR plot and confusion matrix
best_alpha_cat_num_freatures = 0.001
```

```

In [158]: SGD = SGDClassifier(loss="log",alpha= best_alpha_cat_num_freatures,class_weight = 'balanced')
SGD.fit(X_train_cat_num_freatures, y_train)

# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

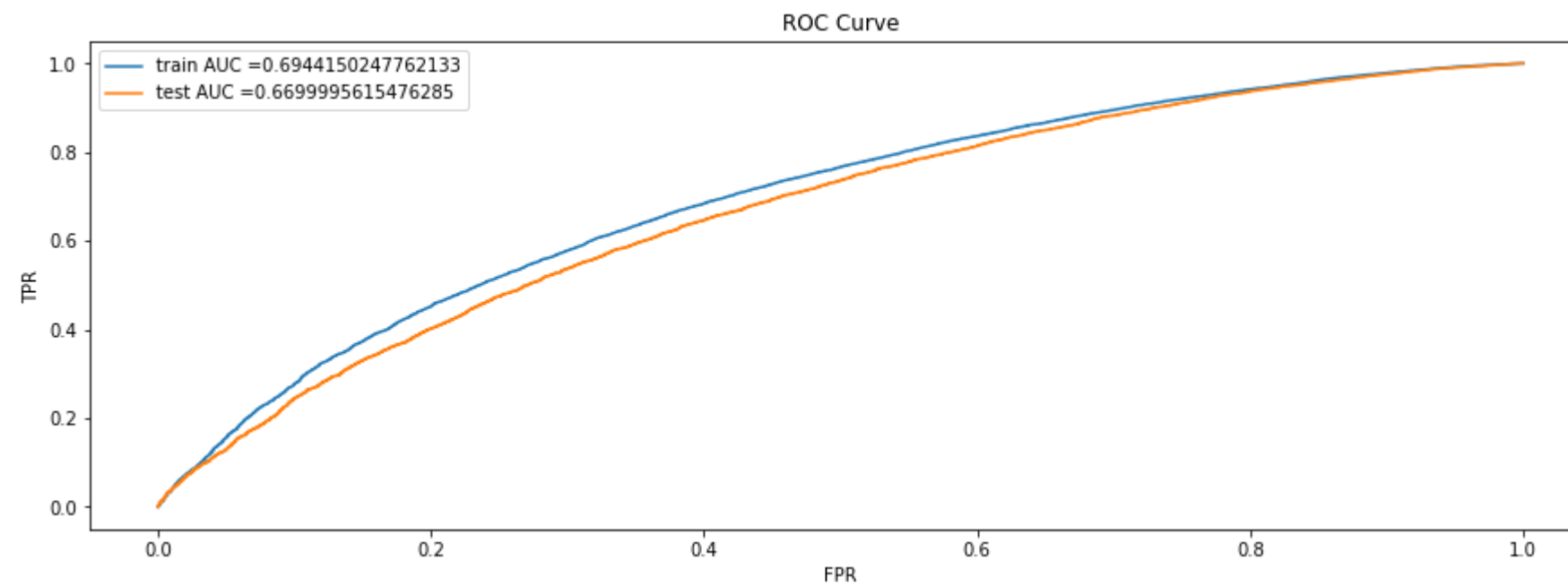
y_train_pred_cat_num_freatures = SGD.decision_function(X_train_cat_num_freatures)
y_test_pred_cat_num_freatures = SGD.decision_function(X_test_cat_num_freatures)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_cat_num_freatures)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_cat_num_freatures)

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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")

plt.show()

```



## Confusion Matrix

### Train confusion matrix

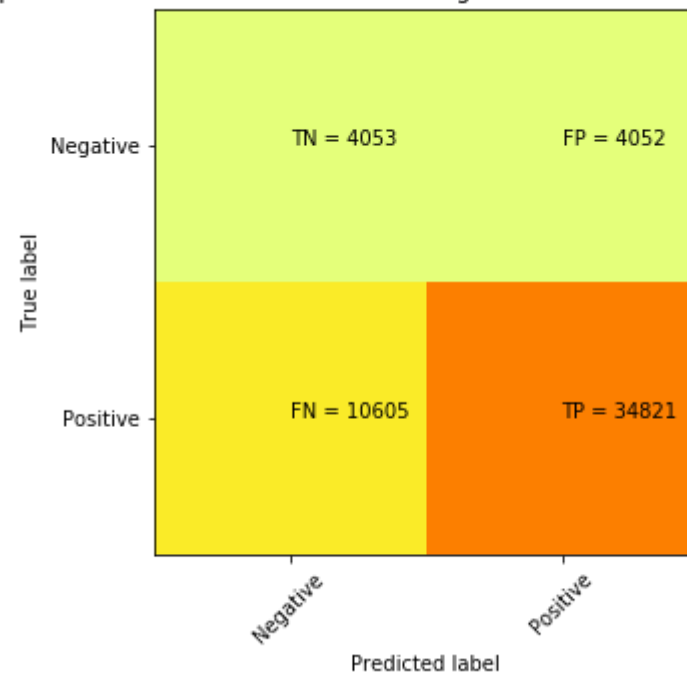
```
In [159]: from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
cat_num_freatures_train_confusion_matrix = confusion_matrix(y_train, predict(y_train_pred_cat_num_freatures, tr_thresholds, train_fpr, train_fpr))
print(cat_num_freatures_train_confusion_matrix)

#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(cat_num_freatures_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+" = "+str(cat_num_freatures_train_confusion_matrix[i][j]))
plt.show()
```

Train confusion matrix  
the maximum value of tpr\*(1-fpr) 0.2499999961943051 for threshold -0.363  
[[ 4053 4052]  
[10605 34821]]

Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Train Data



Test confusion matrix



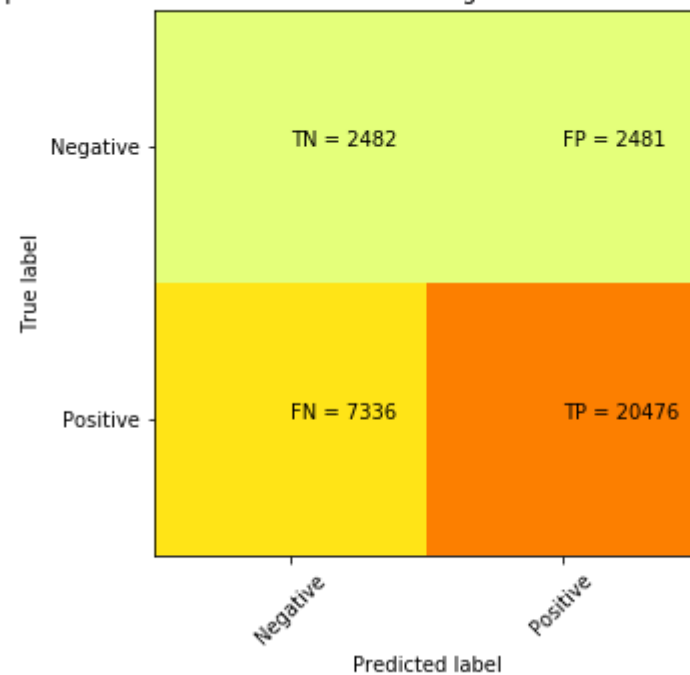
```
In [160]: print("Train confusion matrix")
cat_num_freatures_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_cat_num_freatures, te_thresholds, test_fpr, test_fpr))
print(cat_num_freatures_test_confusion_matrix)

##http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

plt.clf()
plt.imshow(cat_num_freatures_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
plt.xticks(tick_marks, classNames, rotation=45)
plt.yticks(tick_marks, classNames)
s = [['TN', 'FP'], ['FN', 'TP']]
for i in range(2):
    for j in range(2):
        plt.text(j,i, str(s[i][j])+ " = "+str(cat_num_freatures_test_confusion_matrix[i][j]))
plt.show()
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999998985034083 for threshold -0.337  
[[ 2482 2481]  
[ 7336 20476]]

Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Test Data



### 3. Conclusions

```
In [162]: ## L1 Reg. optimal hyperparameter has been taken from Hyperparameter vs AUC Pretty table.
## L2 Reg. optimal hyperparameter has been taken from FPR vs TPR plot
```

```
In [161]: from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]
x.add_row(["BOW", "SGD Classifier (with Loss = 'hinge')L1-Reg ", 0.0001, 0.64])
x.add_row(["BOW", "SGD Classifier (with Loss = 'hinge')L2-Reg ", 0.01, 0.71])
x.add_row(["TFIDF", "SGD Classifier (with Loss = 'hinge')L1-Reg ", 0.0001, 0.67])
x.add_row(["TFIDF", "SGD Classifier (with Loss = 'hinge')L2-Reg ", 0.001, 0.68])
x.add_row(["AVG W2V", "SGD Classifier (with Loss = 'hinge')L1-Reg ", 0.0005, 0.68])
x.add_row(["AVG W2V", "SGD Classifier (with Loss = 'hinge')L2-Reg ", 0.001, 0.68])
x.add_row(["TFIDF W2V", "SGD Classifier (with Loss = 'hinge')L1-Reg ", 0.001, 0.67])
x.add_row(["TFIDF W2V", "SGD Classifier (with Loss = 'hinge')L2-Reg ", 0.001, 0.68])
x.add_row(["Cat, Num and Truncated Text SVD freatures", "SGD Classifier (with Loss = 'hinge')L1-Reg ", 0.0001, 0.69])
x.add_row(["Cat, Num and Truncated Text SVD freatures", "SGD Classifier (with Loss = 'hinge')L2-Reg ", 0.001, 0.68])
print(x)
```

Vectorizer	Model	Hyper Parameter	AUC
BOW	SGD Classifier (with Loss = 'hinge')L1-Reg	0.0001	0.64
BOW	SGD Classifier (with Loss = 'hinge')L2-Reg	0.01	0.71
TFIDF	SGD Classifier (with Loss = 'hinge')L1-Reg	0.0001	0.67
TFIDF	SGD Classifier (with Loss = 'hinge')L2-Reg	0.001	0.68
AVG W2V	SGD Classifier (with Loss = 'hinge')L1-Reg	0.0005	0.68
AVG W2V	SGD Classifier (with Loss = 'hinge')L2-Reg	0.001	0.68
TFIDF W2V	SGD Classifier (with Loss = 'hinge')L1-Reg	0.001	0.67
TFIDF W2V	SGD Classifier (with Loss = 'hinge')L2-Reg	0.001	0.68
Cat, Num and Truncated Text SVD freatures	SGD Classifier (with Loss = 'hinge')L1-Reg	0.0001	0.69
Cat, Num and Truncated Text SVD freatures	SGD Classifier (with Loss = 'hinge')L2-Reg	0.001	0.68