

# Logistic Regression on DonorsChoose data set

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*                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| <code>project_submitted_datetime</code>                   | Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| <code>teacher_id</code>                                   | A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| <code>teacher_prefix</code>                               | 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>                                                                                                                                                                                                                                                                                                                                                                            |
| <code>teacher_number_of_previously_posted_projects</code> | 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                                                                                   |
|-----------------|-----------------------------------------------------------------------------------------------|
| <code>id</code> | A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502 |

| Feature     | Description                                                                   |
|-------------|-------------------------------------------------------------------------------|
| 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 |
|--------|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|----------------------------|
| 101880 | 5749       | p096076 | 6eaa448903897a152320bd23a30147b2 | Mrs.           | CA           | 2016-01-05 00:00:00 | Grades PreK-2          | Math                       |
| 31477  | 47750      | p185738 | 3afe10b996b7646d8641985a4b4b570d | Mrs.           | UT           | 2016-01-05 01:05:00 | Grades PreK-2          | Math                       |

```
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                         |
|--------|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|---------------------------------------|
| 101880 | 5749       | p096076 | 6eaa448903897a152320bd23a30147b2 | Mrs.           | CA           | 2016-01-05 00:00:00 | Grades PreK-2          | Math & Science, Warmth, Care & Hunger |
| 31477  | 47750      | p185738 | 3afe10b996b7646d8641985a4b4b570d | Mrs.           | UT           | 2016-01-05 01:05:00 | Grades PreK-2          | Math & Science, Warmth, Care & Hunger |

```
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 in 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 these 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. nannan

=====

```
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 is 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. nannan

```
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 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 is 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 nannan

```
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', "were \
            n't", \
            'won', "won't", 'wouldn', "wouldn't"]
```

### 1.4.1 Data Processing (Essay)



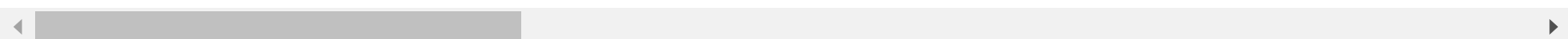
```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248 [00:54<00:00, 2002.78it/
s]
```

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

Out[23]:

|        | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | Date                | project_grade_category | pr       |
|--------|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|----------|
| 101880 | 5749       | p096076 | 6eaa448903897a152320bd23a30147b2 | Mrs.           | CA           | 2016-01-05 00:00:00 | Grades PreK-2          | Ma<br>Ma |
| 31477  | 47750      | p185738 | 3afe10b996b7646d8641985a4b4b570d | Mrs.           | UT           | 2016-01-05 01:05:00 | Grades PreK-2          | Ma       |

2 rows x 21 columns



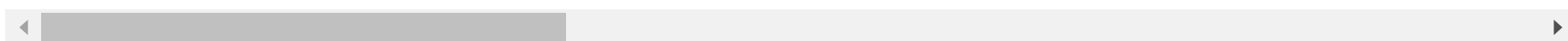
```
In [28]: project_data["Neutral SC Essay"] = neutral
```

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

Out[31]:

|        | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | Date                | project_grade_category | pr       |
|--------|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|----------|
| 101880 | 5749       | p096076 | 6eaa448903897a152320bd23a30147b2 | Mrs.           | CA           | 2016-01-05 00:00:00 | Grades PreK-2          | Ma<br>Ma |
| 31477  | 47750      | p185738 | 3afe10b996b7646d8641985a4b4b570d | Mrs.           | UT           | 2016-01-05 01:05:00 | Grades PreK-2          | Ma       |

2 rows × 25 columns



```
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!

```
In [35]: # Combining all the above statements
from tqdm import tqdm
preprocessed_Titles = []
# tqdm is for printing the status bar
for Pance in tqdm(project_data['project_title'].values):
    P = decontracted(Pance)
    P = P.replace('\\r', ' ')
    P = P.replace('\\\"', ' ')
    P = P.replace('\\n', ' ')
    P = re.sub('[^A-Za-z0-9]+', ' ', P)
    # https://gist.github.com/sebleier/554280
    P = ' '.join(e for e in P.split() if e not in stopwords)
    preprocessed_Titles.append(P.lower().strip())
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109248 [00:02<00:00, 44254.57it/s]
```

```
In [36]: project data["preprocessed Titles"] = preprocessed Titles
```

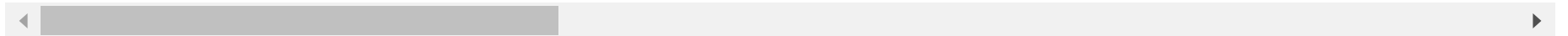
```
In [37]: project_data['title_word_count'] = [len(x.split()) for x in project_data['preprocessed_Titles'].tolist()]
```

In [38]: `project_data.head(2)`

Out[38]:

|        | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | Date                | project_grade_category | project_title |
|--------|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|---------------|
| 101880 | 5749       | p096076 | 6eaa448903897a152320bd23a30147b2 | Mrs.           | CA           | 2016-01-05 00:00:00 | Grades PreK-2          | Me            |
| 31477  | 47750      | p185738 | 3afe10b996b7646d8641985a4b4b570d | Mrs.           | UT           | 2016-01-05 01:05:00 | Grades PreK-2          | Me            |

2 rows × 27 columns



## Train , Cross Validation and Test Data Split

```
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_S
cience', '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', 'CharacterEduc
ation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelo
pment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNee
ds', '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', 'C
O', 'MN', 'OR', 'KY', 'MS', 'NV', 'MD', 'CT', 'TN', 'UT', 'AL', '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 avoide 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 encodig ",project_grade_category_one_hot_train.shape)
print("Shape of matrix after one hot encodig ",project_grade_category_one_hot_cv.shape)
print("Shape of matrix after one hot encodig ",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 encodig  (53531, 4)
Shape of matrix after one hot encodig  (22942, 4)
Shape of matrix after one hot encodig  (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(ngram_range=(2,2), min_df=10, max_features = 5000)
vectorizer.fit(X_train["preprocessed_essays"])
bow_essays_train = vectorizer.fit_transform(X_train["preprocessed_essays"])
print("Shape of matrix after one hot encodig ",bow_essays_train.shape)

Shape of matrix after one hot encodig  (53531, 5000)
```

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

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

Shape of matrix after one hot encodig (22942, 5000)

### 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, 5000)

### 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(ngram_range=(2,2), min_df=10, max_features = 5000)
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, 1864)

### 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, 1864)

### 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, 1864)

## 1.6.2.2 TFIDF vectorizer

### Train Data Vectorization - TFIDF (essays)

```
In [54]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(ngram_range=(2,2), min_df=10, max_features = 5000)
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, 5000)

### 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, 5000)

### 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, 5000)

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

```
In [57]: vectorizer = CountVectorizer(ngram_range=(2,2), min_df=10, max_features = 5000)
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, 1864)

### 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, 1864)

### 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, 1864)

### 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 [03:58, 8027.35it/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 16229882

the unique words in the coupus 42668

The number of words that are present in both glove vectors and our coupus 38939 ( 91.26 %)

word 2 vec length 38939





```
In [66]: # 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]))  
  
100%|██████████████████████████████████████████████████████████████████████████| 22942/22942 [00:06<00:00, 3760.28it/  
s]  
  
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]))  
  
100%|██████████████████████████████████████████████████████████| 32775/32775 [00:07<00:00, 4281.19it/  
s]  
  
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]))  
  
100%|██████████████████████████████████████████████████████████████████████████| 53531/53531 [00:00<00:00, 80127.14it/  
s]  
  
53531  
300
```

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

```
In [69]: # average Word2Vec
# compute average word2vec for each review.
avg_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
    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_cv.append(vector)

print(len(avg_w2v_title_cv))
print(len(avg_w2v_title_cv[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 22942/22942 [00:00<00:00, 73261.59it/s]
```

```
22942
300
```

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

```
In [70]: # average Word2Vec
# compute average word2vec for each review.
avg_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
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title_test.append(vector)

print(len(avg_w2v_title_test))
print(len(avg_w2v_title_test[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 32775/32775 [00:00<00:00, 72109.76it/s]
```

```
32775
300
```

### 1.6.2.3 Using Pretrained Models: TFIDF weighted W2V

```
In [71]: # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train["preprocessed_essays"])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
```

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

```
100%|██████████████████████████████████████████████████████████████████████████████| 53531/53531 [01:35<00:00, 562.56it/
s]

53531
300
```

## CV Data Vectorization - TFIDF W2V (essays)

```
100% |██████████████████████████████████████████████████████████████████████████████| 22942/22942 [00:40<00:00, 561.42it/
s]

22942
300
```

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

```
In [75]: # 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())
```

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

```
100%|██████████████████████████████████████████████████████████████| 22942/22942 [00:00<00:00, 36467.67it/s]
22942
300
```

### Test Data Vectorization - TFIDF W2V (Project Titles)

```
100%|██████████████████████████████████████████████████████████████████████████████| 32775/32775 [00:00<00:00, 35582.75it/s]
32775
300
```

### 1.6.3 Vectorizing Numerical features

### 1.6.3.1 Vectorizing Numerical features - Price

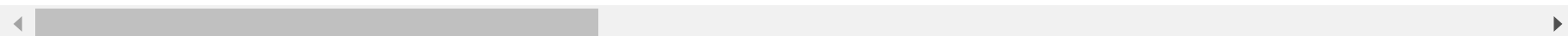
```
# Merging the project data train , Cv , test with price from resource data
X_train = pd.merge(X_train, price_data, on='id', how='left')
X_cv = pd.merge(X_cv, price_data, on='id', how='left')
X_test = pd.merge(X_test, price_data, on='id', how='left')
```

In [80]: X\_train.head(2)

Out[80]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | Date                | project_grade_category | project_title                                |
|---|------------|---------|----------------------------------|----------------|--------------|---------------------|------------------------|----------------------------------------------|
| 0 | 123439     | p242877 | 0262e6419ee4373bfe75364968b2c8a4 | Mrs.           | FL           | 2016-05-26 19:20:00 | Grades PreK-2          | Let's Wiggle and Wobble Our Way to Learning! |
| 1 | 23580      | p119560 | f80a46ee3a6b18fd856c887f0ae3454c | Mrs.           | OH           | 2016-09-19 20:56:00 | Grades PreK-2          | Technology Comes Alive!                      |

2 rows × 28 columns



In [81]: [#https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html](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 5: Logistic Regression



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


- Set 1: categorical, numerical features + project\_title(BOW) + preprocessed\_eassay (`BOW with bi-grams` with `min\_df=10` and `max\_features=5000`)
- Set 2: categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (`TFIDF with bi-grams` with `min\_df=10` and `max\_features=5000`)
- Set 3: categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)
- Set 4: categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_essay (TFIDF W2V)

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

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

### 3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.  

- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.  

- Along with plotting ROC curve, you need to print the confusion matrix (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.

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

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

#### 5. Consider these set of features Set 5 :

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category :categorical data
- teacher\_prefix : categorical data
- quantity : numerical data
- teacher\_number\_of\_previously\_posted\_projects : numerical data
- price : numerical data
- sentiment score's of each of the essay : numerical data
- number of words in the title : numerical data
- number of words in the combine essays : numerical data

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

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

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

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

## 2. Logistic Regression

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

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

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

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

**BOW with bi-grams with min\_df=10 and max\_features=5000**

```
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, 6965)
```

```
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, 6965)
```

```
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, 6965)
```

```
In [93]: from sklearn.metrics import roc_auc_score
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs

    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate unti the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000]))[:,1])
    # we will be predicting for the last data points
    y_data_pred.extend(clf.predict_proba(data[tr_loop:]))[:,1])

    return y_data_pred
```

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

```
In [102]: from sklearn.linear_model import SGDClassifier
import math
```

```

In [104]: # 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="log",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}]

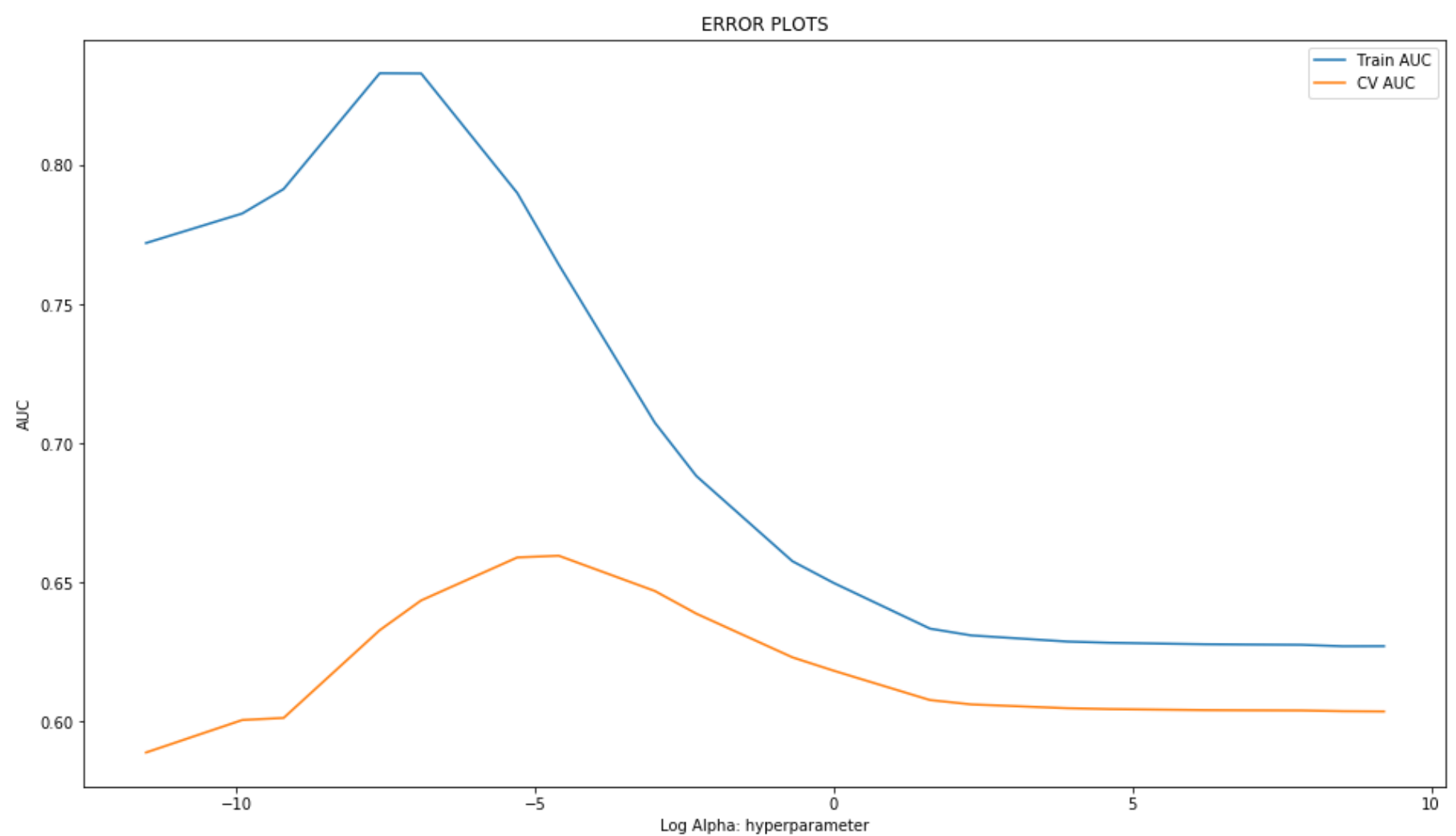
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_bow, y_train)

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 [105]: #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 | 0.7718653002377307 | 0.5889895758370888 |
| 5e-05  | -9.903487552536127  | 0.7824532103413321 | 0.6006447616705843 |
| 0.0001 | -9.210340371976182  | 0.7912356142351258 | 0.6013963105899292 |
| 0.0005 | -7.600902459542082  | 0.8328145665089149 | 0.632903732238323  |
| 0.001  | -6.907755278982137  | 0.8327320334722476 | 0.6435876114938036 |
| 0.005  | -5.298317366548036  | 0.7898088630068872 | 0.6590088945306684 |
| 0.01   | -4.605170185988091  | 0.7641518565889228 | 0.6596229974693796 |
| 0.05   | -2.995732273553991  | 0.7073631172099284 | 0.646962844459953  |
| 0.1    | -2.3025850929940455 | 0.6882706597650561 | 0.638841651390341  |
| 0.5    | -0.6931471805599453 | 0.6576530233801318 | 0.6231527288255172 |
| 1      | 0.0                 | 0.6498274514248273 | 0.6183384794083445 |
| 5      | 1.6094379124341003  | 0.6334609897484842 | 0.607819070132629  |
| 10     | 2.302585092994046   | 0.631014087758406  | 0.606248726058912  |
| 50     | 3.912023005428146   | 0.628802543078828  | 0.6048809366723331 |
| 100    | 4.605170185988092   | 0.6283857306238282 | 0.6045772896590605 |
| 500    | 6.214608098422191   | 0.6278448254340567 | 0.604216578285192  |
| 1000   | 6.907755278982137   | 0.6277438360524511 | 0.6041509213150295 |
| 2500   | 7.824046010856292   | 0.6276462444230452 | 0.6040854924435074 |
| 5000   | 8.517193191416238   | 0.6271468603537826 | 0.6038204142548022 |
| 10000  | 9.210340371976184   | 0.6271703061629293 | 0.6037343839552527 |

### Using Best alpha Value – Training the Model

```
In [106]: best_alpha_bow = 0.001
```

```
In [107]: #https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

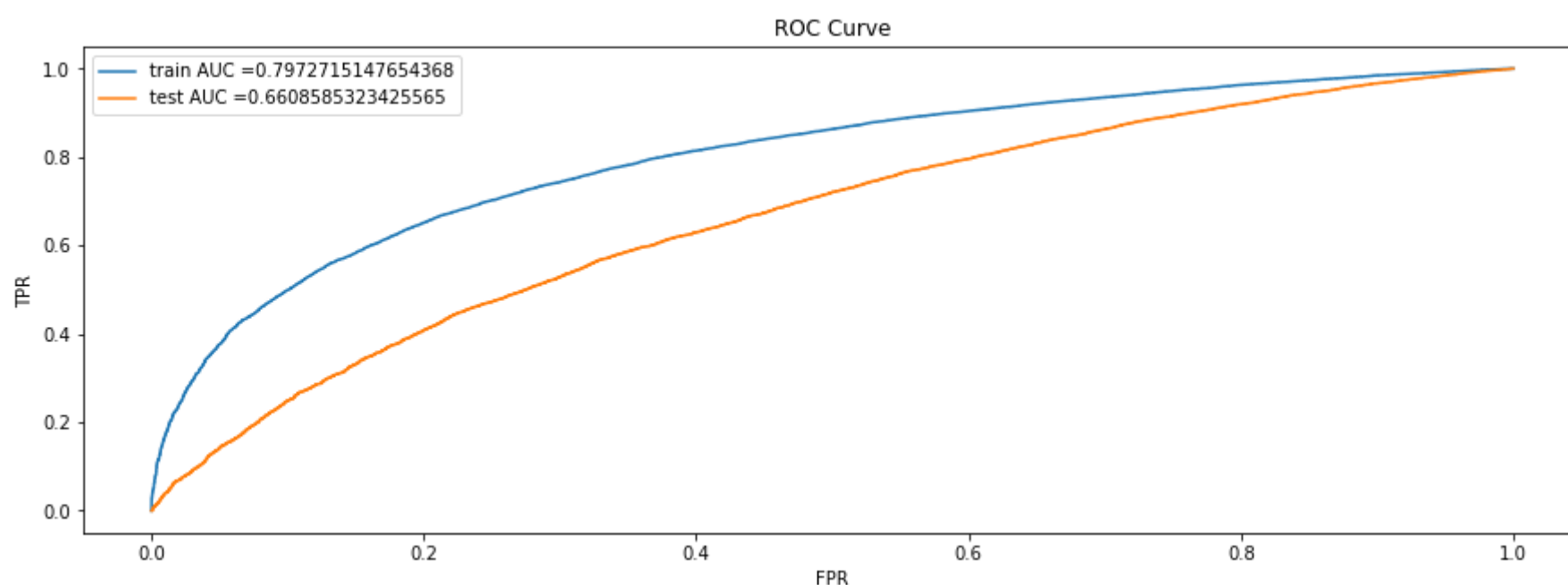
SGD = SGDClassifier(loss="log",alpha= best_alpha_bow,class_weight = 'balanced')
SGD.fit(X_train_bow, y_train)

# 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 = batch_predict(SGD, X_train_bow)
y_test_pred_bow = batch_predict(SGD, 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 [108]: [#https://stackoverflow.com/questions/28719067/roc-curve-and-cut-off-point-python](https://stackoverflow.com/questions/28719067/roc-curve-and-cut-off-point-python)

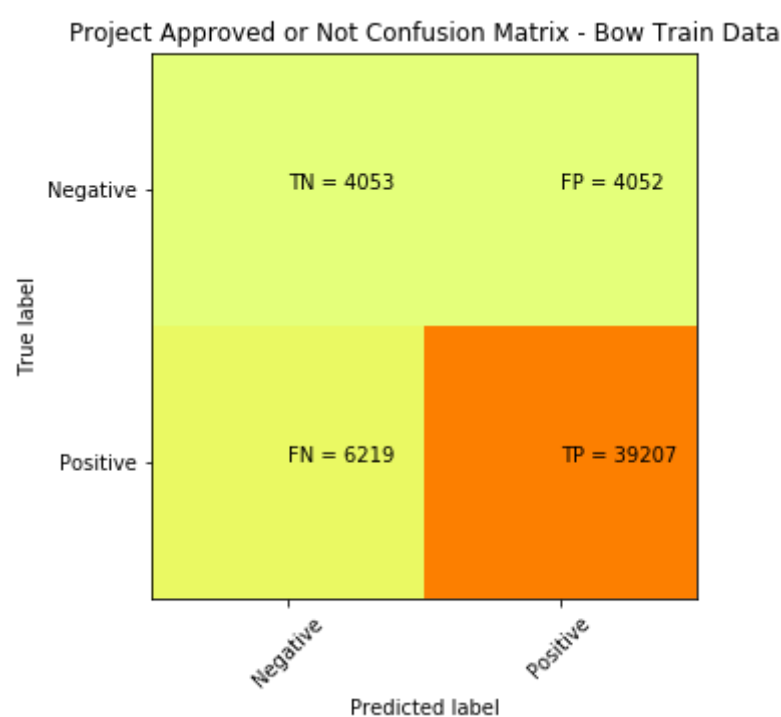
```
def predict(proba, threshold, fpr, tpr):
    t = threshold[np.argmax(fpr*(1-tpr))]
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
           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

In [109]: `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) 0.2499999961943051 for threshold 0.385
[[ 4053  4052]
 [ 6219 39207]]
```

In [110]: [#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/](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()
```



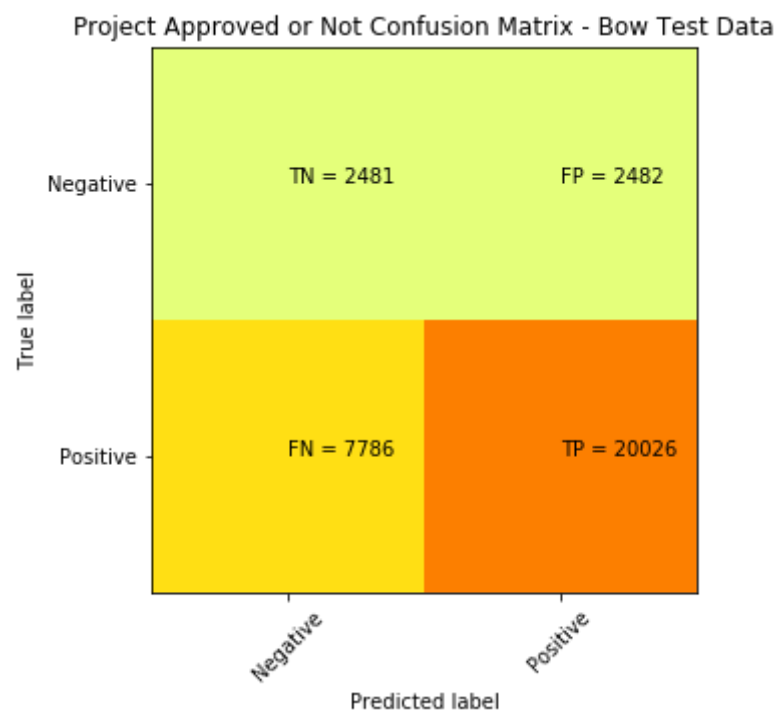
### Test confusion matrix

In [111]: `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 0.465
[[ 2481  2482]
 [ 7786 20026]]
```

In [112]: <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 = 'log') on TFIDF, SET 2

TFIDF with bi-grams with min\_df=10 and max\_features=5000

In [113]: `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[113]: (53531, 6965)

In [114]: `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[114]: (22942, 6965)

In [115]: `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[115]: (32775, 6965)

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

```

In [117]: from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import GridSearchCV

sgd = SGDClassifier(loss="log",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}]

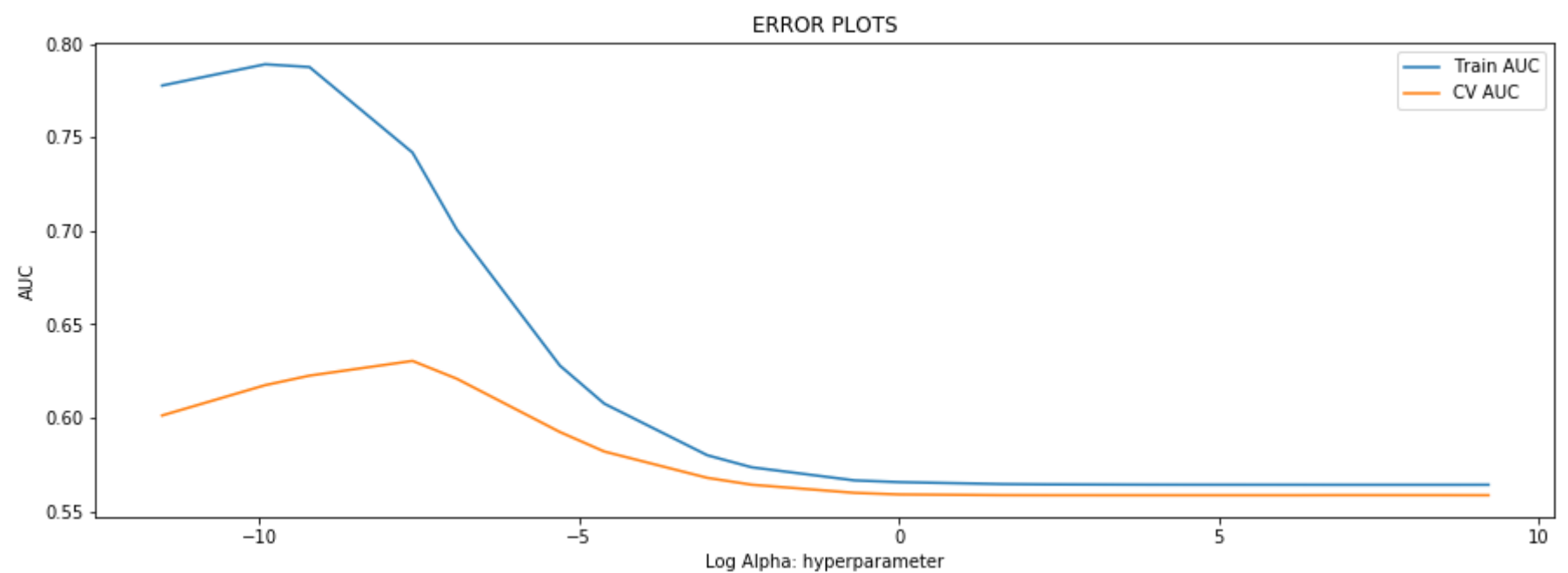
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf, y_train)

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 [118]: #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 | 0.7775765887622192 | 0.6011666559996026 |
| 5e-05  | -9.903487552536127  | 0.7889507438991665 | 0.6173669569936115 |
| 0.0001 | -9.210340371976182  | 0.7874695552461676 | 0.6224663966815855 |
| 0.0005 | -7.600902459542082  | 0.7417677141837135 | 0.6303519167216632 |
| 0.001  | -6.907755278982137  | 0.7005614770969005 | 0.6208049318982388 |
| 0.005  | -5.298317366548036  | 0.6278600817706149 | 0.592337751185932  |
| 0.01   | -4.605170185988091  | 0.6075098748198854 | 0.5819061342942915 |
| 0.05   | -2.995732273553991  | 0.5799204790709885 | 0.5678368625243985 |
| 0.1    | -2.3025850929940455 | 0.5734695857516744 | 0.5641483346035544 |
| 0.5    | -0.6931471805599453 | 0.5664506956765468 | 0.5597968198181246 |
| 1      | 0.0                 | 0.5655352865964174 | 0.5589511213165848 |
| 5      | 1.6094379124341003  | 0.5644765013148305 | 0.5585614285965628 |
| 10     | 2.302585092994046   | 0.5643122784341624 | 0.5585225317925335 |
| 50     | 3.912023005428146   | 0.5641744601918767 | 0.5585153693191578 |
| 100    | 4.605170185988092   | 0.5641554499711553 | 0.558513121902649  |
| 500    | 6.214608098422191   | 0.5641361891268772 | 0.5585169770750567 |
| 1000   | 6.907755278982137   | 0.5641160750129518 | 0.5585654861995436 |
| 2500   | 7.824046010856292   | 0.564116023596264  | 0.558554785336744  |
| 5000   | 8.517193191416238   | 0.5641171212244821 | 0.5585506531122707 |
| 10000  | 9.210340371976184   | 0.5641256919093339 | 0.5585224785466505 |

### Using Best K Value – Training the Model

```
In [119]: best_alpha_tfidf = 0.001
```

```
In [120]: #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="log",alpha= best_alpha_tfidf,class_weight ='balanced')
SGD.fit(X_train_tfidf, y_train)

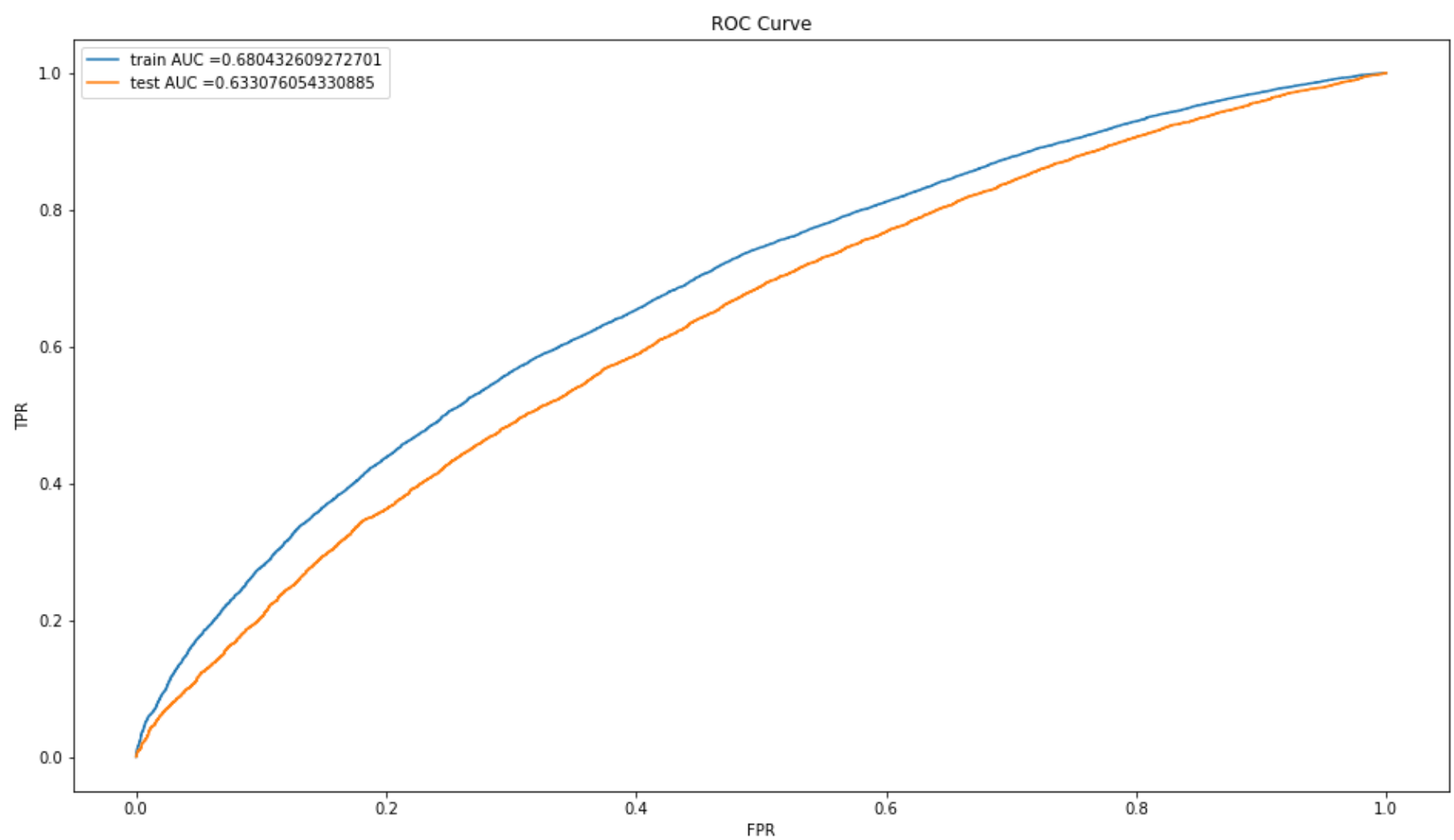
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_tfidf = batch_predict(SGD, X_train_tfidf)
y_test_pred_tfidf = batch_predict(SGD, 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

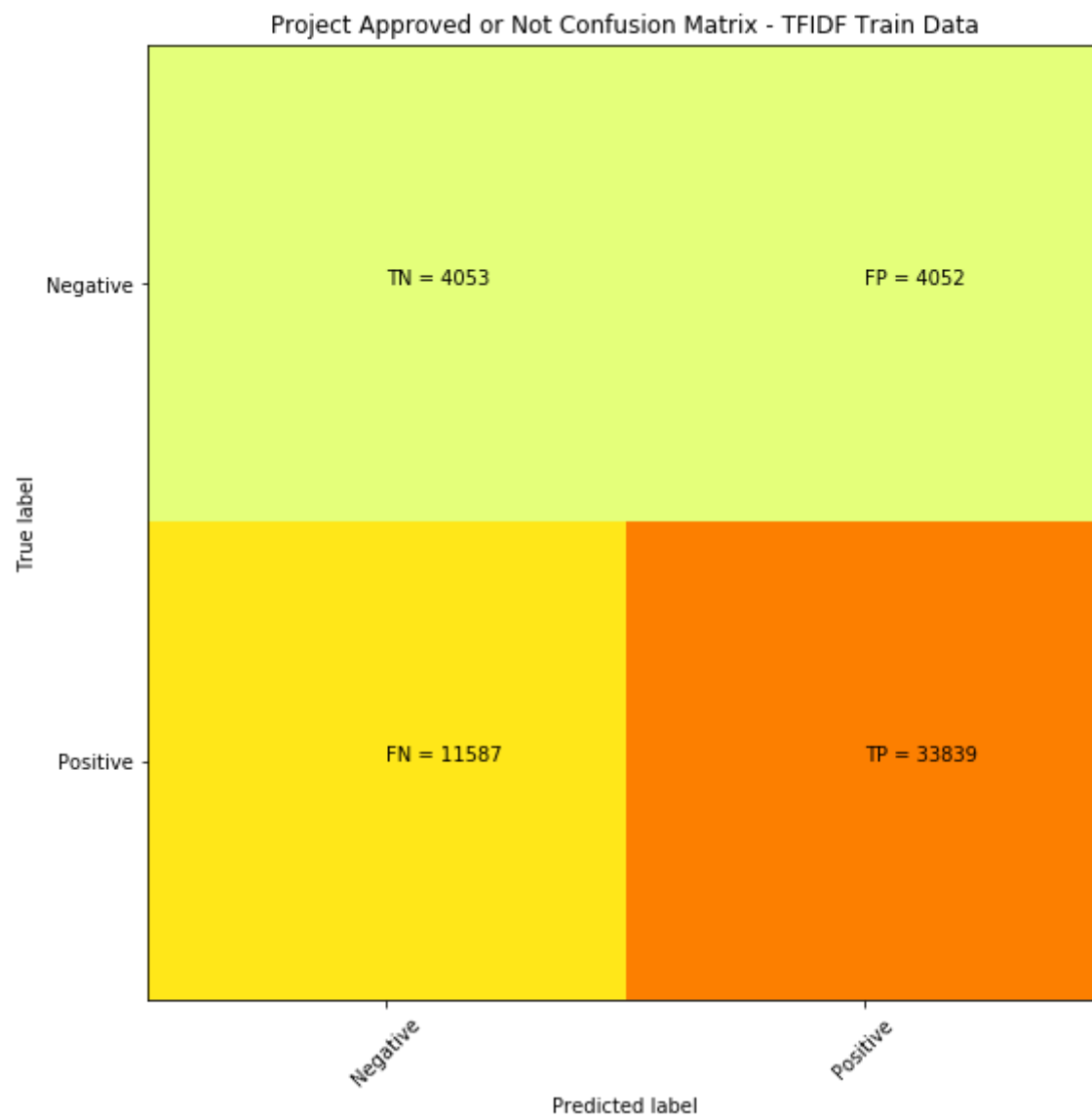
### Train confusion matrix

```
In [121]: 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*(1-fpr) 0.2499999961943051 for threshold 0.496
[[ 4053  4052]
 [11587 33839]]
```

In [122]: <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()
```



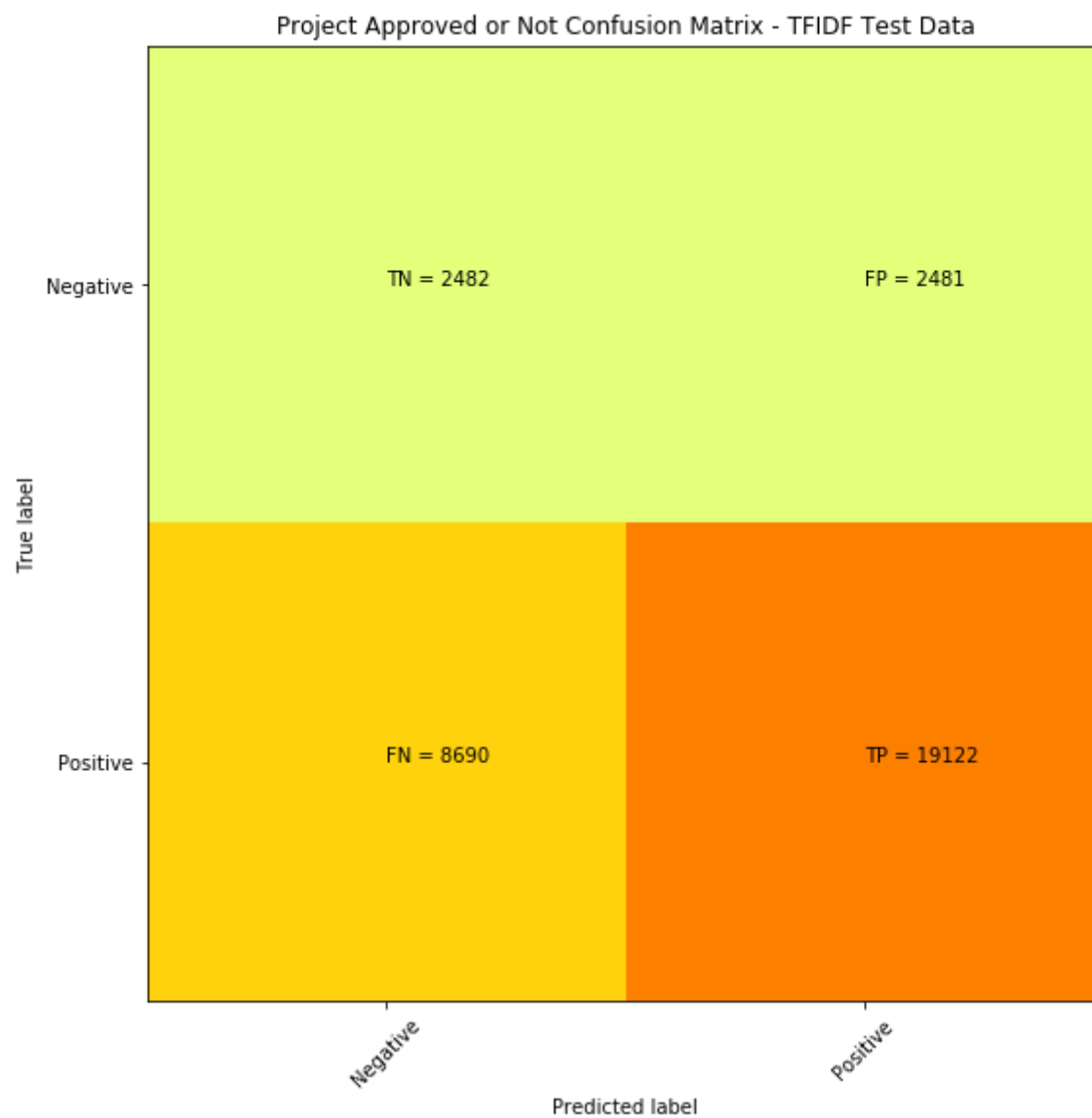
### Test confusion matrix

```
In [123]: 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 0.508
[[ 2482  2481]
 [ 8690 19122]]
```

In [124]: <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 = 'log') on AVG W2V, SET 3

In [125]: `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[125]: (53531, 701)

In [126]: `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[126]: (22942, 701)

In [127]: `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_essays_test, avg_w2v_title_test, price_test, prev_post_test)).tocsr()`  
`X_test_avg_w2v.shape`

Out[127]: (32775, 701)

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

```

In [128]: sgd = SGDClassifier(loss="log",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}]

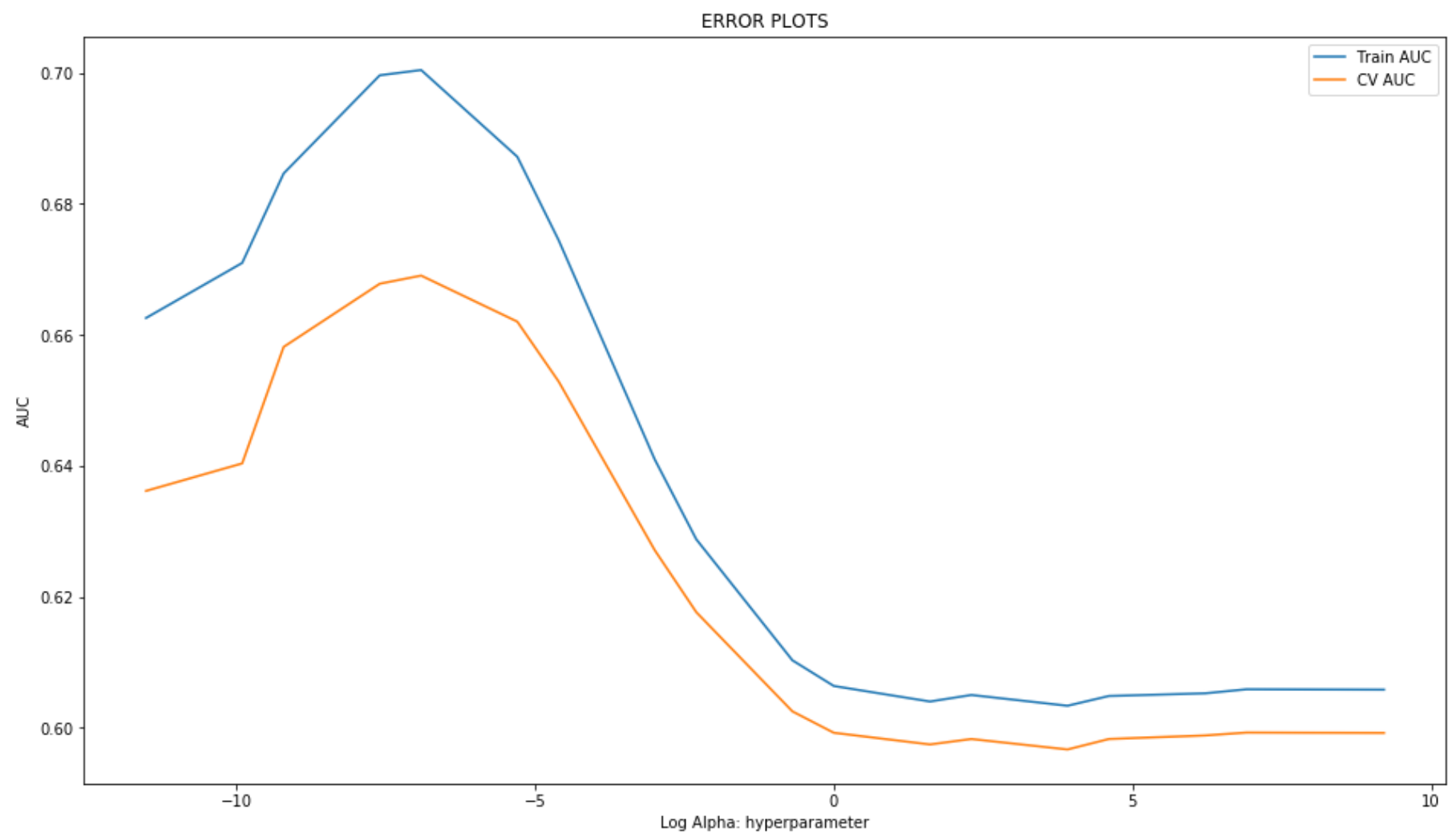
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_avg_w2v, y_train)

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 [129]: #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 | 0.6625528419060726 | 0.6361612337178001 |
| 5e-05  | -9.903487552536127  | 0.6709577606750036 | 0.6403565210472756 |
| 0.0001 | -9.210340371976182  | 0.6846039162661511 | 0.6581332795534267 |
| 0.0005 | -7.600902459542082  | 0.6995952563314108 | 0.6677806017498972 |
| 0.001  | -6.907755278982137  | 0.7003972470465012 | 0.6690210012465461 |
| 0.005  | -5.298317366548036  | 0.6871629947072511 | 0.6619975439889921 |
| 0.01   | -4.605170185988091  | 0.6744060373543411 | 0.6528451269776412 |
| 0.05   | -2.995732273553991  | 0.6409266676048916 | 0.6270871769648324 |
| 0.1    | -2.3025850929940455 | 0.6287797887647979 | 0.617656927178138  |
| 0.5    | -0.6931471805599453 | 0.6103062296490743 | 0.6025095692637262 |
| 1      | 0.0                 | 0.6063849534878107 | 0.5992425629114563 |
| 5      | 1.6094379124341003  | 0.6040147084542631 | 0.5974636025497657 |
| 10     | 2.302585092994046   | 0.6050136991545316 | 0.5982951035894813 |
| 50     | 3.912023005428146   | 0.603367344214437  | 0.5966947556757757 |
| 100    | 4.605170185988092   | 0.6048589385447124 | 0.5982975960090093 |
| 500    | 6.214608098422191   | 0.6052660968915241 | 0.5988324757716188 |
| 1000   | 6.907755278982137   | 0.6058855708412606 | 0.5992754465608966 |
| 2500   | 7.824046010856292   | 0.6058607843624836 | 0.5992513113221826 |
| 5000   | 8.517193191416238   | 0.6058460545930959 | 0.5992396353993357 |
| 10000  | 9.210340371976184   | 0.6058292753198907 | 0.5992229318740522 |

### Using Best K Value – Training the Model

```
In [130]: best_alpha_avg_w2v = 0.001
```

```
In [131]: #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="log",alpha= best_alpha_avg_w2v,class_weight = 'balanced')
SGD.fit(X_train_avg_w2v, y_train)

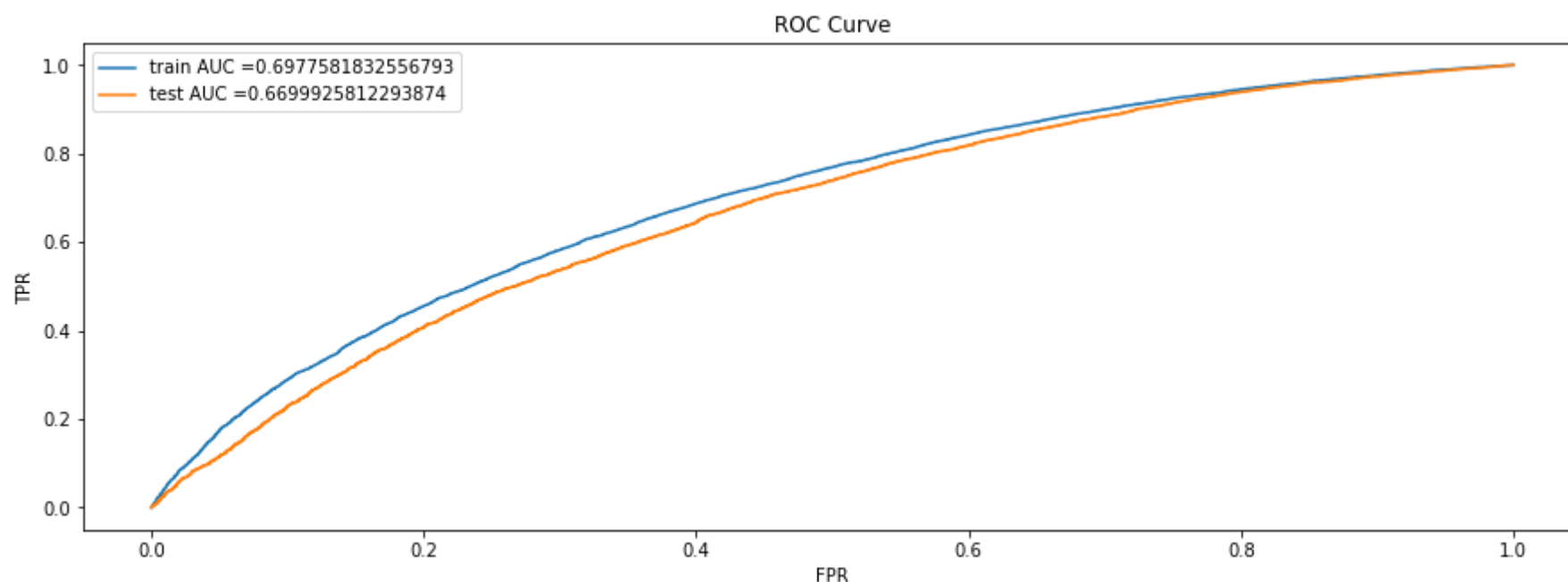
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_avg_w2v = batch_predict(SGD, X_train_avg_w2v)
y_test_pred_avg_w2v = batch_predict(SGD, 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

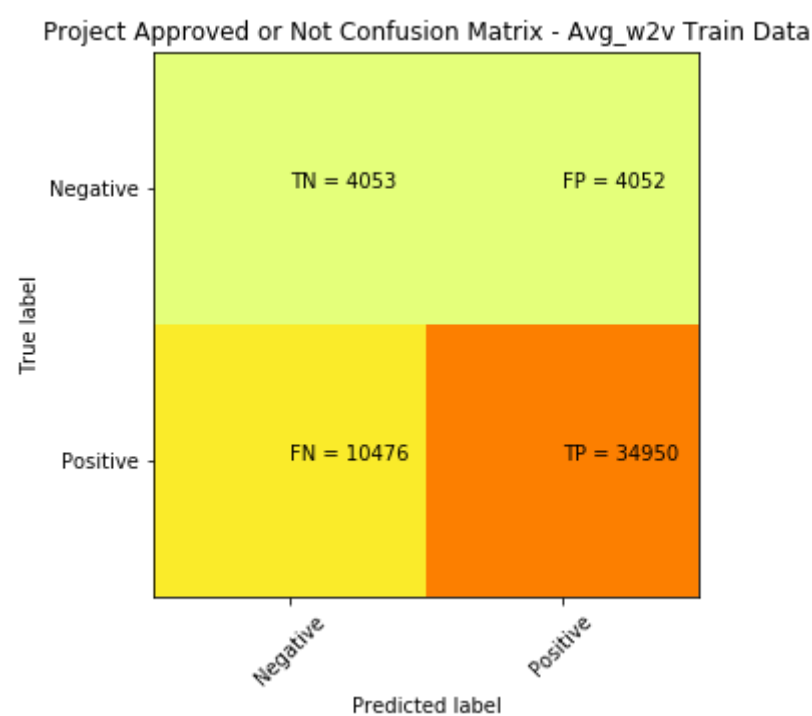
### Train confusion matrix

```
In [132]: 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 0.413  
[[ 4053 4052]  
[10476 34950]]



### Test confusion matrix

```
In [133]: 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/

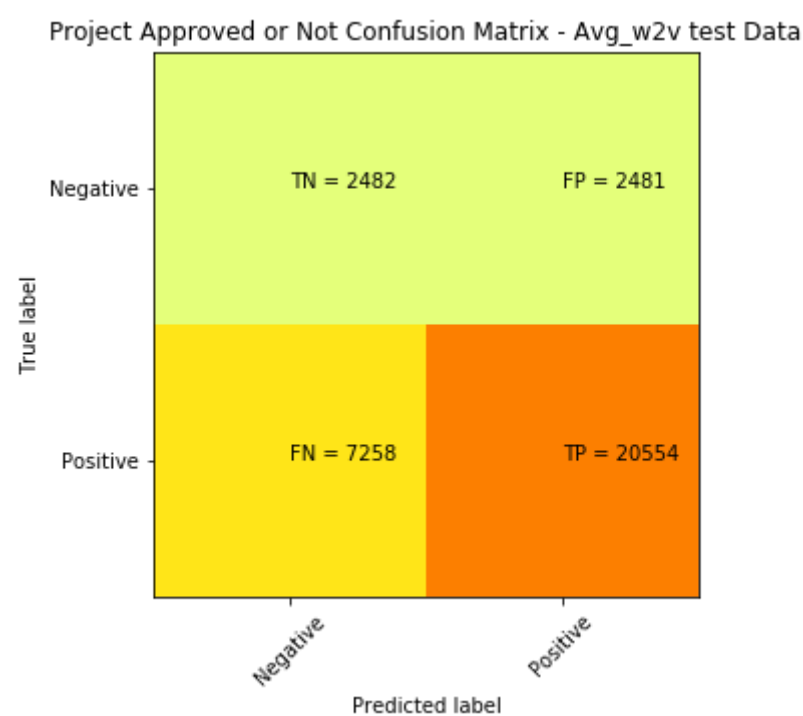
plt.clf()
plt.imshow(avg_w2v_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Avg_w2v test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
tick_marks = np.arange(len(classNames))
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 0.423

```
[[ 2482  2481]
```

```
 [ 7258 20554]]
```



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

```
In [134]: 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[134]: (53531, 701)

```
In [135]: 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[135]: (22942, 701)

```
In [136]: 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_w2v_essays_test, tfidf_w2v_title_test, price_test, prev_post_test)).tocsr()
X_test_tfidf_w2v.shape
```

Out[136]: (32775, 701)

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



```

In [137]: sgd = SGDClassifier(loss="log",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}]

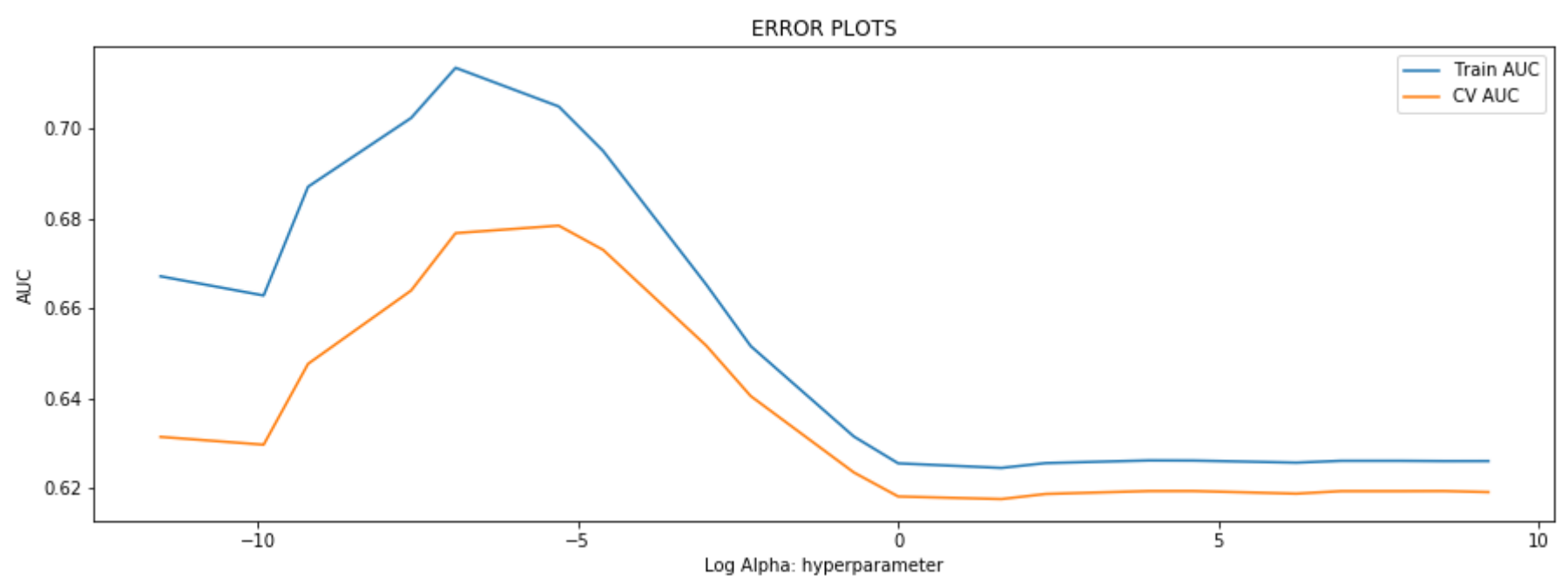
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_tfidf_w2v, y_train)

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 [138]: #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 | 0.6671149153207473 | 0.6313269639939523 |
| 5e-05  | -9.903487552536127  | 0.6628393977616932 | 0.6295873290910916 |
| 0.0001 | -9.210340371976182  | 0.6870726947428034 | 0.6476056069247839 |
| 0.0005 | -7.600902459542082  | 0.7024106970105071 | 0.6639562319046739 |
| 0.001  | -6.907755278982137  | 0.7135757946153304 | 0.6767348136210685 |
| 0.005  | -5.298317366548036  | 0.7049586829831367 | 0.6784073061486604 |
| 0.01   | -4.605170185988091  | 0.695036654174232  | 0.6730005737472593 |
| 0.05   | -2.995732273553991  | 0.6652580074455495 | 0.6516520805335198 |
| 0.1    | -2.3025850929940455 | 0.6515647874904074 | 0.6404261242384088 |
| 0.5    | -0.6931471805599453 | 0.6314158870694279 | 0.6234055311333121 |
| 1      | 0.0                 | 0.625425148467147  | 0.6180379594227331 |
| 5      | 1.6094379124341003  | 0.6244102945270252 | 0.6174663816197169 |
| 10     | 2.302585092994046   | 0.6254830301763005 | 0.6186082149657617 |
| 50     | 3.912023005428146   | 0.6260795343415834 | 0.6192526778519115 |
| 100    | 4.605170185988092   | 0.6260495634093509 | 0.6192556027479295 |
| 500    | 6.214608098422191   | 0.6255848058515042 | 0.6186719833473787 |
| 1000   | 6.907755278982137   | 0.6259919052876963 | 0.6192304321999001 |
| 2500   | 7.824046010856292   | 0.625993129749927  | 0.619231498987801  |
| 5000   | 8.517193191416238   | 0.6259291656771503 | 0.6192563490079852 |
| 10000  | 9.210340371976184   | 0.6259365389026553 | 0.6190136065459421 |

## Using Best K Value – Training the Model

```
In [139]: best_alpha_tfidf_w2v = 0.001
```

```
In [140]: #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="log",alpha= best_alpha_tfidf_w2v,class_weight = 'balanced')
SGD.fit(X_train_tfidf_w2v, y_train)

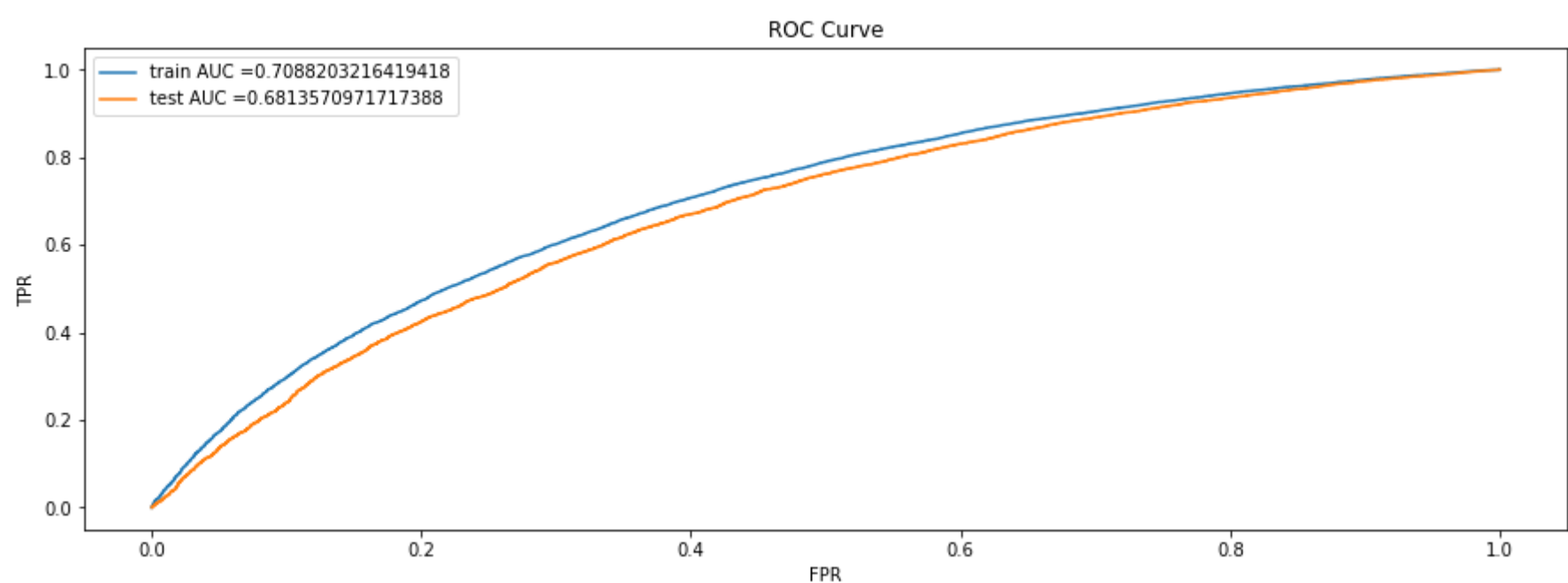
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_tfidf_w2v = batch_predict(SGD, X_train_tfidf_w2v)
y_test_pred_tfidf_w2v = batch_predict(SGD, 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 [141]: 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)

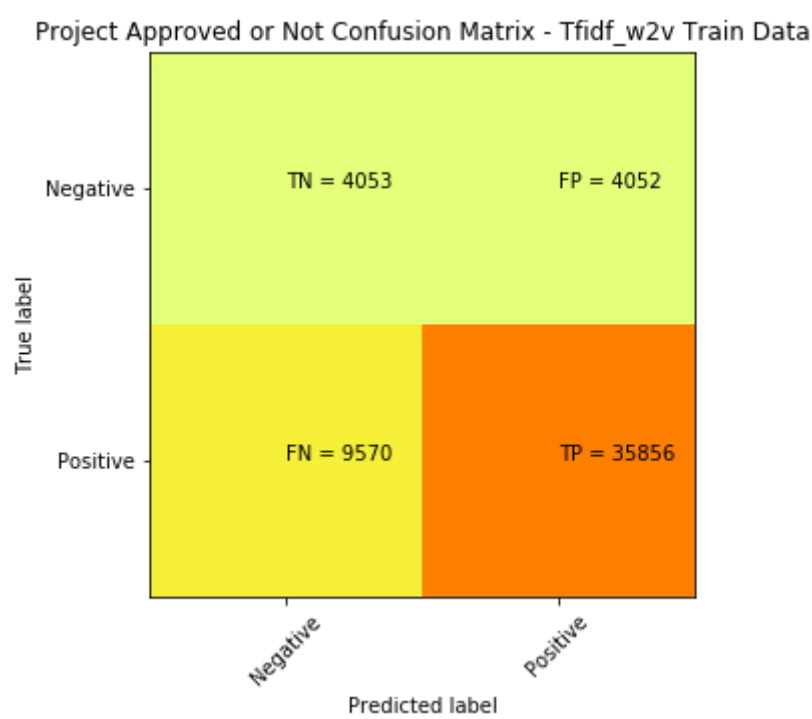
#http://www.tarekatwan.com/index.php/2017/12/how-to-plot-a-confusion-matrix-in-python/

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 0.422

```
[[ 4053  4052]
 [ 9570 35856]]
```



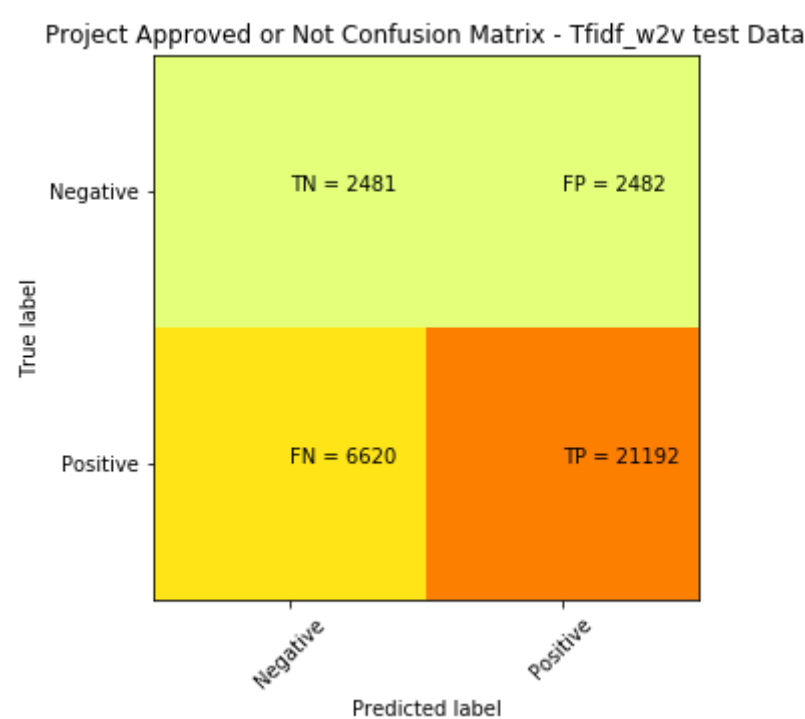
**Test confusion matrix**

```
In [142]: 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 0.435  
[[ 2481 2482]  
[ 6620 21192]]



## 2.5 Categorical and Numerical (ONLY)

```
In [143]: 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)).tocsr()
X_train_cat_num_freatures.shape
```

Out[143]: (53531, 108)

```
In [144]: 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)).tocsr()
X_cv_cat_num_freatures.shape
```

Out[144]: (22942, 108)

```
In [145]: 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)).tocsr()
X_test_cat_num_freatures.shape
```

Out[145]: (32775, 108)

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

```

In [146]: sgd = SGDClassifier(loss="log",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}]

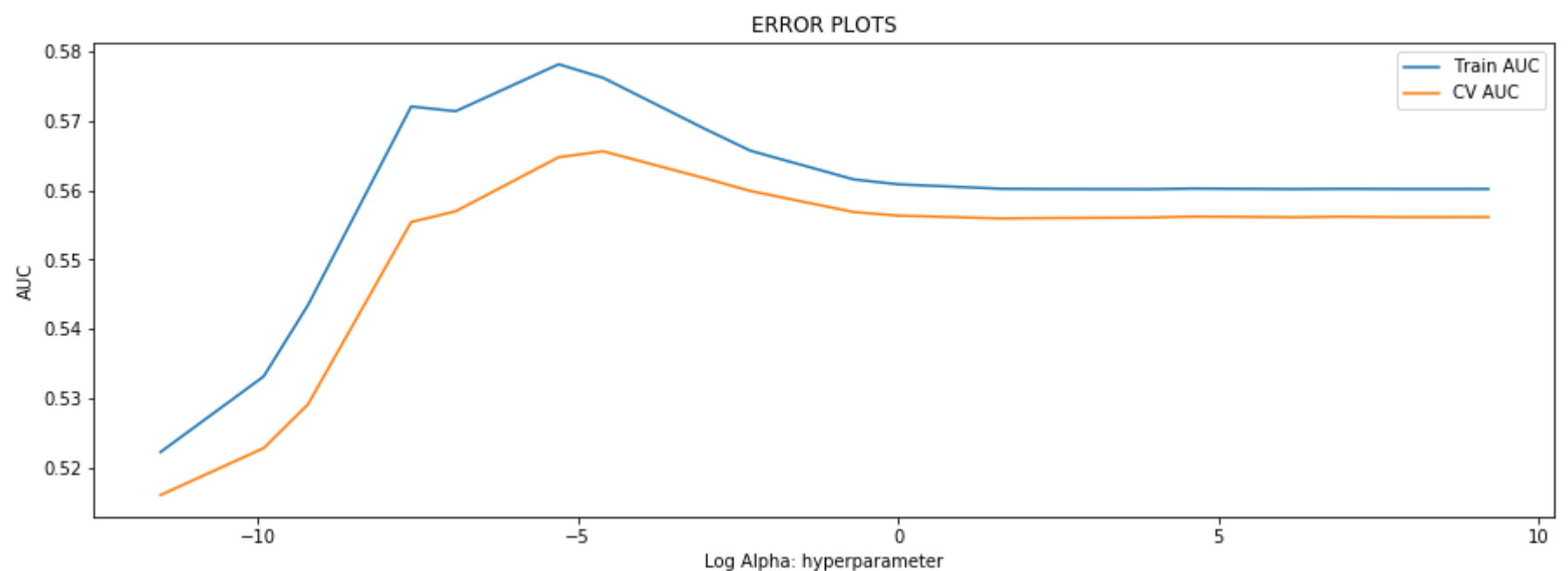
clf = GridSearchCV(sgd,tuned_parameters, cv=3, scoring='roc_auc')
clf.fit(X_train_cat_num_freatures, y_train)

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 [147]: #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 | 0.522299936179675  | 0.5161133407421102 |
| 5e-05  | -9.903487552536127  | 0.5331843048982967 | 0.5228646627409833 |
| 0.0001 | -9.210340371976182  | 0.5434748403011258 | 0.5291909845526561 |
| 0.0005 | -7.600902459542082  | 0.5720476975221271 | 0.5554140184463552 |
| 0.001  | -6.907755278982137  | 0.5713644244739279 | 0.5569613245495302 |
| 0.005  | -5.298317366548036  | 0.5781270260097442 | 0.5647525234654343 |
| 0.01   | -4.605170185988091  | 0.5761717016203762 | 0.5655998342475677 |
| 0.05   | -2.995732273553991  | 0.5687240193729055 | 0.5616567862919305 |
| 0.1    | -2.3025850929940455 | 0.5656651619222745 | 0.5598742615428395 |
| 0.5    | -0.6931471805599453 | 0.5615471027880935 | 0.5568497865502098 |
| 1      | 0.0                 | 0.5608449045569582 | 0.5563417340893506 |
| 5      | 1.6094379124341003  | 0.5602039725941016 | 0.5559317847519096 |
| 10     | 2.302585092994046   | 0.5601646983061783 | 0.5559894714442356 |
| 50     | 3.912023005428146   | 0.5601415575719434 | 0.5560545684742675 |
| 100    | 4.605170185988092   | 0.5602350066290861 | 0.5562067593749908 |
| 500    | 6.214608098422191   | 0.5601540254809223 | 0.556102242764248  |
| 1000   | 6.907755278982137   | 0.5601960951248908 | 0.5561830824296998 |
| 2500   | 7.824046010856292   | 0.5601651395573543 | 0.556119437467091  |
| 5000   | 8.517193191416238   | 0.5601648766040012 | 0.5561213603576012 |
| 10000  | 9.210340371976184   | 0.560167868926162  | 0.5561284248141191 |

## Using Best K Value – Training the Model

```
In [148]: best_alpha_cat_num_freatures = 0.01
```

```
In [149]: SGD = SGDClassifier(loss="log",alpha= best_alpha_cat_num_freatures,class_weight = 'balanced')
SGD.fit(X_train_cat_num_freatures, y_train)

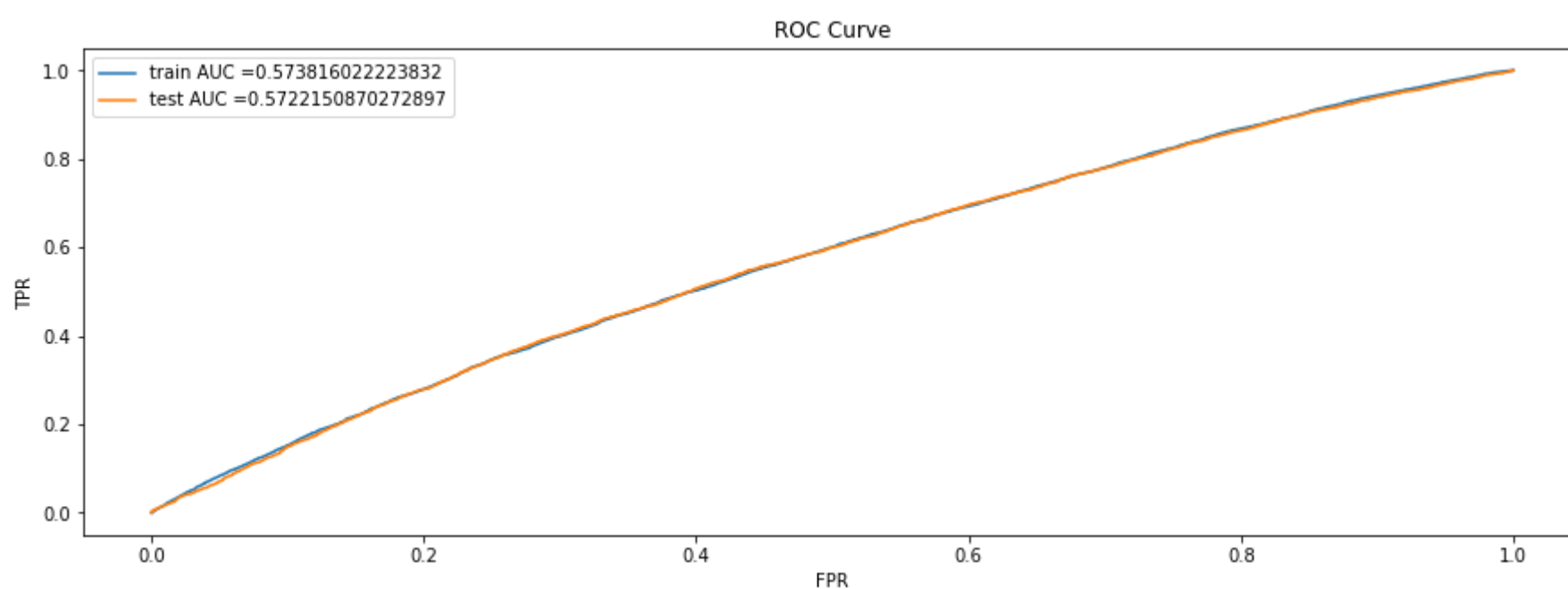
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred_cat_num_freatures = batch_predict(SGD, X_train_cat_num_freatures)
y_test_pred_cat_num_freatures = batch_predict(SGD, 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 [150]: 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_thresh
olds, 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/

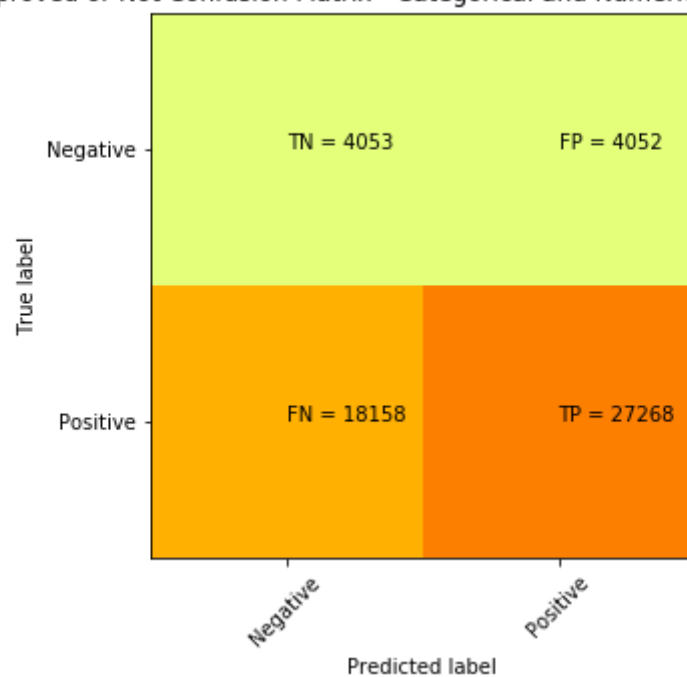
plt.clf()
plt.imshow(cat_num_freatures_train_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Train Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
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  $tpr \cdot (1 - fpr)$  0.2499999961943051 for threshold 0.488

```
[[ 4053  4052]
 [18158 27268]]
```

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



## Test confusion matrix

```
In [151]: print("Train confusion matrix")
cat_num_freatures_test_confusion_matrix = confusion_matrix(y_test, predict(y_test_pred_cat_num_freatures, te_threshold
s, 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/

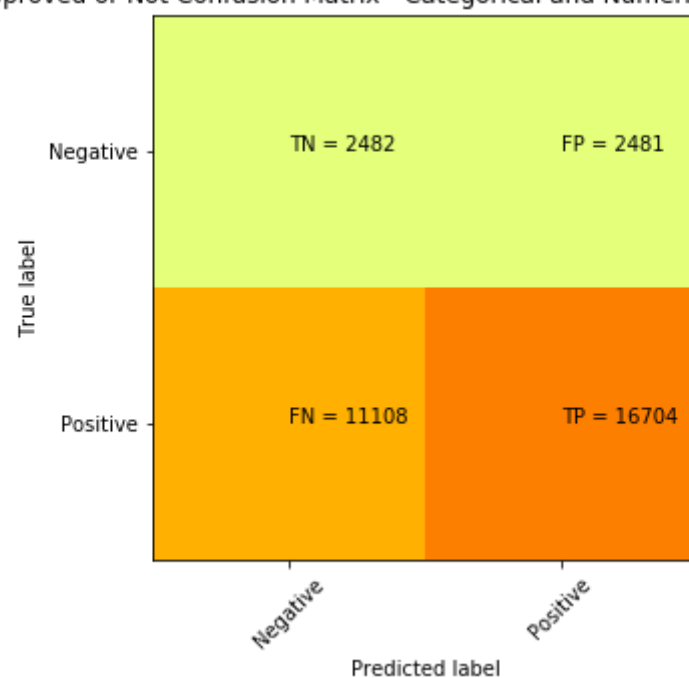
plt.clf()
plt.imshow(cat_num_freatures_test_confusion_matrix, interpolation='nearest', cmap=plt.cm.Wistia)
classNames = ['Negative', 'Positive']
plt.title('Project Approved or Not Confusion Matrix - Categorical and Numerical freatures Test Data')
plt.ylabel('True label')
plt.xlabel('Predicted label')
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\*(1-fpr) 0.24999998985034083 for threshold 0.488

```
[[ 2482  2481]
 [11108 16704]]
```

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



### 3. Conclusions

```
In [1]: 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(Bi-gram)", "SGD Classifier (with Loss = 'log')", "0.001", 0.66])
x.add_row(["TFIDF(Bi-gram)", "SGD Classifier (with Loss = 'log')", "0.001", 0.63])
x.add_row(["AVG W2V", "SGD Classifier (with Loss = 'log')", "0.001", 0.67])
x.add_row(["TFIDF W2V", "SGD Classifier (with Loss = 'log')", "0.001", 0.68])
x.add_row(["Categorical and Numerical freatures", "SGD Classifier (with Loss = 'log')", "0.01", 0.57])
print(x)
```

| Vectorizer                          | Model                              | Hyper Parameter | AUC  |
|-------------------------------------|------------------------------------|-----------------|------|
| BOW(Bi-gram)                        | SGD Classifier (with Loss = 'log') | 0.001           | 0.66 |
| TFIDF(Bi-gram)                      | SGD Classifier (with Loss = 'log') | 0.001           | 0.63 |
| AVG W2V                             | SGD Classifier (with Loss = 'log') | 0.001           | 0.67 |
| TFIDF W2V                           | SGD Classifier (with Loss = 'log') | 0.001           | 0.68 |
| Categorical and Numerical freatures | SGD Classifier (with Loss = 'log') | 0.01            | 0.57 |