

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12

Feature	Description	
	One or more (comma-separated) subject categories for the project from the following enumerated list of values:	
project_subject_categories	<ul style="list-style-type: none"> • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth 	
	Examples:	
	<ul style="list-style-type: none"> • Music & The Arts • Literacy & Language, Math & Science 	
school_state	State where school is located (Two-letter U.S. postal code (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). Example: WY	
project_subject_subcategories	One or more (comma-separated) subject subcategories for the project. Examples:	
	<ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Sciences 	
project_resource_summary	An explanation of the resources needed for the project. Example:	
	<ul style="list-style-type: none"> • My students need hands on literacy materials to manage sensory needs!</code 	
project_essay_1	First application essay*	
project_essay_2	Second application essay*	
project_essay_3	Third application essay*	
project_essay_4	Fourth application essay*	
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245	
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56	
teacher_prefix	Teacher's title. One of the following enumerated values:	
	<ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher. 	
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2	

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

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

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

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

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

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved, and a value of <code>1</code> 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 [0]:

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4 import sqlite3
5 import pandas as pd
6 import numpy as np
7 import nltk
8 import string
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 from sklearn.feature_extraction.text import TfidfTransformer
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn.metrics import confusion_matrix
15 from sklearn import metrics
16 from sklearn.metrics import roc_curve, auc
17 from nltk.stem.porter import PorterStemmer
18 import re
19 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
20 import string
21 from nltk.corpus import stopwords
22 from nltk.stem import PorterStemmer
23 from nltk.stem.wordnet import WordNetLemmatizer
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27 from tqdm import tqdm
28 import os
29 import chart_studio.plotly
30 # from plotly import plotly
31 import plotly.offline as offline
32 import plotly.graph_objs as go
33 offline.init_notebook_mode()
34 from collections import Counter
35 from scipy.sparse import hstack, vstack
36 from sklearn.model_selection import train_test_split
37 from sklearn.neighbors import KNeighborsClassifier
38 from sklearn.metrics import accuracy_score
39 from sklearn.model_selection import cross_val_score
40 from sklearn import model_selection
41 from sklearn.preprocessing import StandardScaler
42 from sklearn.model_selection import RandomizedSearchCV
43 #from sklearn.impute import SimpleImputer
44 from sklearn.datasets import load_digits
45 #from sklearn.feature_selection import SelectKBest, chi2
```

```

46 #from sklearn.model_selection import GridSearchCV
47 from sklearn.feature_selection import SelectKBest,f_classif
48 from prettytable import PrettyTable
49 from sklearn.naive_bayes import MultinomialNB
50 from sklearn.preprocessing import Normalizer
51 from sklearn.metrics import confusion_matrix
52 #import math
53 #from sklearn.linear_model import LogisticRegression
54 import nltk
55 from nltk.sentiment.vader import SentimentIntensityAnalyzer
56 from sklearn.linear_model import SGDClassifier
57 nltk.download('vader_lexicon')
58 import pdb
59 from sklearn.decomposition import TruncatedSVD

```

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...

1.1 Reading Data

In [0]:

```

1 Project_data = pd.read_csv('train_data.csv')
2 Resource_data = pd.read_csv('resources.csv')
3 print(Project_data.shape)
4 print(Resource_data.shape)

```

(109248, 17)

(1541272, 4)

```
In [0]: 1 # how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
2 cols = ['Date' if x=='project_submitted_datetime' else x for x in list(Project_data.columns)]
3 #sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
4 Project_data['Date'] = pd.to_datetime(Project_data['project_submitted_datetime'])
5 Project_data.drop('project_submitted_datetime', axis=1, inplace=True)
6 Project_data.sort_values(by=['Date'], inplace=True)
7 # how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
8 Project_data = Project_data[cols]
9 Project_data.head(2)
```

```
Out[5]:
```

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject_categories	project
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	Math & Science	Applie
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	Special Needs	

1.2 preprocessing of project_subject_categories

```
In [0]: 1 y = Project_data['project_is_approved'].values
2 Project_data.drop(['project_is_approved'], axis=1, inplace=True)
3 lpd = len(Project_data)
4 ys = np.zeros(lpd, dtype=np.int32)
5 X = Project_data
```

```
In [0]: 1 #Splitting the Dataset into three Train,CV and Test
2 X1, X_Test, Y1, Y_Test = train_test_split(X, y, test_size=0.33, random_state=0, stratify=ys)
3 nx1 = len(X1)
4 ys1 = np.zeros(nx1, dtype=np.int32)
5 X_Train, X_CV, Y_Train, Y_CV = train_test_split(X1, Y1, test_size=0.33, random_state=0, stratify=ys1)
6 print('Shape of the X_Train data is {0} and Y_Train data is: {1}'.format(X_Train.shape,Y_Train.shape[0]))
7 print('Shape of the X_CV data is {0} and Y_CV data is : {1}'.format(X_CV.shape,Y_CV.shape[0]))
8 print('Shape of the X_Test data is {0} and Y_Test data is : {1}'.format(X_Test.shape,Y_Test.shape[0]))
```

Shape of the X_Train data is (49041, 16) and Y_Train data is: 49041

Shape of the X_CV data is (24155, 16) and Y_CV data is : 24155

Shape of the X_Test data is (36052, 16) and Y_Test data is : 36052

In [0]:

```
1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6
7 #*****Train Data*****
8 categories = list(X_Train['project_subject_categories'].values)
9 cat_list = []
10 for i in categories:
11     temp = ""
12     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
13     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
14         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
15             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
16             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
17             temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
18             temp = temp.replace('&', '_') # we are replacing the & value into
19     cat_list.append(temp.strip())
20
21 X_Train['clean_categories'] = cat_list
22 X_Train.drop(['project_subject_categories'], axis=1, inplace=True)
23
24 from collections import Counter
25 my_counter = Counter()
26 for word in X_Train['clean_categories'].values:
27     my_counter.update(word.split())
28
29 cat_dict = dict(my_counter)
30 sorted_cat_dict_Train = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
31 print(len(sorted_cat_dict_Train))
32 #*****CV Data*****
33 categories = list(X_CV['project_subject_categories'].values)
34 cat_list = []
35 for i in categories:
36     temp = ""
37     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
38     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
39         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
40             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
41             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
42             temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
43             temp = temp.replace('&', '_') # we are replacing the & value into
44     cat_list.append(temp.strip())
45
```



```

46 X_CV['clean_categories'] = cat_list
47 X_CV.drop(['project_subject_categories'], axis=1, inplace=True)
48
49 *****Test Data*****
50 categories = list(X_Test['project_subject_categories'].values)
51 cat_list = []
52 for i in categories:
53     temp = ""
54     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
55     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
56         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
57             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
58             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
59             temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
60             temp = temp.replace('&','_') # we are replacing the & value into
61         cat_list.append(temp.strip())
62
63 X_Test['clean_categories'] = cat_list
64 X_Test.drop(['project_subject_categories'], axis=1, inplace=True)
65

```

9

1.3 preprocessing of project_subject_subcategories

In [0]:

```
1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6 #*****Train Data*****
7 sub_categories = list(X_Train['project_subject_subcategories'].values)
8 sub_cat_list = []
9 for i in sub_categories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
14             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&','_')
18     sub_cat_list.append(temp.strip())
19
20 X_Train['clean_subcategories'] = sub_cat_list
21 X_Train.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in X_Train['clean_subcategories'].values:
26     my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict_Train = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
30 print(len(sorted_sub_cat_dict_Train))
31 #*****CV Data*****
32 sub_categories = list(X_CV['project_subject_subcategories'].values)
33 sub_cat_list = []
34 for i in sub_categories:
35     temp = ""
36     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
37     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
38         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
39             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
40             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
41             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
42             temp = temp.replace('&','_')
43     sub_cat_list.append(temp.strip())
44
45 X_CV['clean_subcategories'] = sub_cat_list
```

```

46 X_CV.drop(['project_subject_subcategories'], axis=1, inplace=True)
47
48 *****Test Data*****
49 sub_categories = list(X_Test['project_subject_subcategories'].values)
50 sub_cat_list = []
51 for i in sub_categories:
52     temp = ""
53     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
54     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
55         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
56             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
57             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
58             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
59             temp = temp.replace('&','_')
60     sub_cat_list.append(temp.strip())
61
62 X_Test['clean_subcategories'] = sub_cat_list
63 X_Test.drop(['project_subject_subcategories'], axis=1, inplace=True)

```

30

1.3 Text preprocessing

In [0]:

```

1 # merge two column text dataframe:
2 X_Train["essay"] = X_Train["project_essay_1"].map(str) + \
3                     X_Train["project_essay_2"].map(str) + \
4                     X_Train["project_essay_3"].map(str) + \
5                     X_Train["project_essay_4"].map(str)
6
7 X_CV["essay"] = X_CV["project_essay_1"].map(str) + \
8                 X_CV["project_essay_2"].map(str) + \
9                 X_CV["project_essay_3"].map(str) + \
10                X_CV["project_essay_4"].map(str)
11
12 X_Test["essay"] = X_Test["project_essay_1"].map(str) + \
13                  X_Test["project_essay_2"].map(str) + \
14                  X_Test["project_essay_3"].map(str) + \
15                  X_Test["project_essay_4"].map(str)

```

```
In [0]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8     # general
9     phrase = re.sub(r"n't", " not", phrase)
10    phrase = re.sub(r"\'re", " are", phrase)
11    phrase = re.sub(r"\s", " is", phrase)
12    phrase = re.sub(r"\d", " would", phrase)
13    phrase = re.sub(r"\ll", " will", phrase)
14    phrase = re.sub(r"\t", " not", phrase)
15    phrase = re.sub(r"\ve", " have", phrase)
16    phrase = re.sub(r"\m", " am", phrase)
17    return phrase
```

```
In [0]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
11            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
12            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
13            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
14            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
15            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
16            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
17            'won', "won't", 'wouldn', "wouldn't"]
```

In [0]:

```
1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Essays in Train data set-----
5 preprocessed_essays_Train = []
6 for sentence in tqdm(X_Train['essay'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\\r', ' ')
9     sent = sent.replace('\\\"', ' ')
10    sent = sent.replace('\\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14    preprocessed_essays_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Essays in CV data set-----
18 preprocessed_essays_CV = []
19 for sentence in tqdm(X_CV['essay'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\\"', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
27     preprocessed_essays_CV.append(sent.lower().strip())
28 # pdb.set_trace()
29
30 #-----PreProcessing of Essays in Test data set-----
31 preprocessed_essays_Test = []
32 for sentence in tqdm(X_Test['essay'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\\r', ' ')
35     sent = sent.replace('\\\"', ' ')
36     sent = sent.replace('\\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
40     preprocessed_essays_Test.append(sent.lower().strip())
41 # pdb.set_trace()
```

100%|██████████| 49041/49041 [00:26<00:00, 1836.95it/s]

100%|██████████| 24155/24155 [00:13<00:00, 1820.77it/s]

100%|██████████| 36052/36052 [00:19<00:00, 1823.46it/s]

In [0]:

```
1 word_count_essay_Train = []
2 for a in tqdm(X_Train["essay"]):
3     b = len(a.split())
4     word_count_essay_Train.append(b)
5
6 X_Train["word_count_essay_Train"] = word_count_essay_Train
7
8 word_count_essay_CV = []
9 for a in tqdm(X_CV["essay"]):
10     b = len(a.split())
11     word_count_essay_CV.append(b)
12
13 X_CV["word_count_essay_CV"] = word_count_essay_CV
14
15 word_count_essay_Test = []
16 for a in tqdm(X_Test["essay"]):
17     b = len(a.split())
18     word_count_essay_Test.append(b)
19
20 X_Test["word_count_essay_Test"] = word_count_essay_Test
```

100%|██████████| 49041/49041 [00:00<00:00, 66616.32it/s]
100%|██████████| 24155/24155 [00:00<00:00, 67353.42it/s]
100%|██████████| 36052/36052 [00:00<00:00, 65733.43it/s]

1.4 Preprocessing of project_title

In [0]:

```
1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Project Title in Train data set-----
5 preprocessed_titles_Train = []
6 for sentence in tqdm(X_Train['project_title'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\\r', ' ')
9     sent = sent.replace('\\\"', ' ')
10    sent = sent.replace('\\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e not in stopwords)
14    preprocessed_titles_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Project Title in CV data set-----
18 preprocessed_titles_CV = []
19 for sentence in tqdm(X_CV['project_title'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\\"', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e not in stopwords)
27     preprocessed_titles_CV.append(sent.lower().strip())
28 # pdb.set_trace()
29
30 #-----PreProcessing of Project Title in Test data set-----
31 preprocessed_titles_Test = []
32 for sentence in tqdm(X_Test['project_title'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\\r', ' ')
35     sent = sent.replace('\\\"', ' ')
36     sent = sent.replace('\\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e not in stopwords)
40     preprocessed_titles_Test.append(sent.lower().strip())
41 # pdb.set_trace()
```

```
100%|██████████| 49041/49041 [00:01<00:00, 42155.30it/s]
100%|██████████| 24155/24155 [00:00<00:00, 42563.89it/s]
100%|██████████| 36052/36052 [00:00<00:00, 42638.68it/s]
```

```
In [0]: 1 word_count_title_Train = []
2 for a in tqdm(X_Train["project_title"]) :
3     b = len(a.split())
4     word_count_title_Train.append(b)
5
6 X_Train["word_count_title_Train"] = word_count_title_Train
7
8 word_count_title_CV = []
9 for a in tqdm(X_CV["project_title"]) :
10     b = len(a.split())
11     word_count_title_CV.append(b)
12
13 X_CV["word_count_title_CV"] = word_count_title_CV
14
15 word_count_title_Test = []
16 for a in tqdm(X_Test["project_title"]) :
17     b = len(a.split())
18     word_count_title_Test.append(b)
19
20 X_Test["word_count_title_Test"] = word_count_title_Test
```

```
100%|██████████| 49041/49041 [00:00<00:00, 862093.25it/s]
100%|██████████| 24155/24155 [00:00<00:00, 862351.90it/s]
100%|██████████| 36052/36052 [00:00<00:00, 901682.45it/s]
```

1.5 Preparing data for models

1.5.1 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/>)

In [0]:

```
1 #-----Vectorizing categorical data for Train,CV and Test-----
2
3 # we use count vectorizer to convert the values into one hot encoding
4 vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict_Train.keys()), lowercase=False, binary=True)
5 vectorizer_cat.fit(X_Train['clean_categories'].values)
6 categories_one_hot_Train = vectorizer_cat.transform(X_Train['clean_categories'].values)
7 categories_one_hot_CV = vectorizer_cat.transform(X_CV['clean_categories'].values)
8 categories_one_hot_Test = vectorizer_cat.transform(X_Test['clean_categories'].values)
9 print(vectorizer_cat.get_feature_names())
10 print("-"*120)
11 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Train.shape))
12 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(categories_one_hot_CV.shape))
13 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Test.shape))
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
```

```
Shape of Train dataset matrix after one hot encoding is: (49041, 9)
```

```
Shape of CV dataset matrix after one hot encoding is: (24155, 9)
```

```
Shape of Test dataset matrix after one hot encoding is: (36052, 9)
```

```
In [0]: 1 #-----Vectorizing Sub categorical data for Train,CV and Test-----
2
3 # we use count vectorizer to convert the values into one hot encoding
4 vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict_Train.keys()), lowercase=False, binary=True)
5 vectorizer_sub_cat.fit(X_Train['clean_categories'].values)
6 sub_categories_one_hot_Train = vectorizer_sub_cat.transform(X_Train['clean_subcategories'].values)
7 sub_categories_one_hot_CV = vectorizer_sub_cat.transform(X_CV['clean_subcategories'].values)
8 sub_categories_one_hot_Test = vectorizer_sub_cat.transform(X_Test['clean_subcategories'].values)
9 print(vectorizer_sub_cat.get_feature_names())
10 print("-"*120)
11 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Train.shape))
12 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_CV.shape))
13 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Test.shape))
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'Warmth', 'Care_Hunger', 'NutritionEducation', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'EarlyDevelopment', 'Health_LifeScience', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
```

```
-----
Shape of Train dataset matrix after one hot encoding is: (49041, 30)
Shape of CV dataset matrix after one hot encoding is: (24155, 30)
Shape of Test dataset matrix after one hot encoding is: (36052, 30)
```

School State

```

In [0]: 1 #-----Vectorizing categorical data of School state for Train dataset-----
2
3 school_catogories_Train = list(X_Train['school_state'].values)
4 school_list_Train = []
5 for sent in school_catogories_Train:
6     school_list_Train.append(sent.lower().strip())
7 X_Train['school_categories'] = school_list_Train
8 X_Train.drop(['school_state'], axis=1, inplace=True)
9
10 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
11 my_counter_school_Train = Counter()
12 for word in X_Train['school_categories'].values:
13     my_counter_school_Train.update(word.split())
14
15 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
16 school_dict_Train = dict(my_counter_school_Train)
17 sorted_school_dict_Train = dict(sorted(school_dict_Train.items(), key=lambda kv: kv[1]))
18
19 vectorizer_school = CountVectorizer(vocabulary=list(sorted_school_dict_Train.keys()), lowercase=False, binary=True)
20 vectorizer_school.fit(X_Train['school_categories'].values)
21 #print(vectorizer.get_feature_names())
22
23 school_one_hot_Train = vectorizer_school.transform(X_Train['school_categories'].values)
24
25 #-----Vectorizing categorical data of School state for CV dataset-----
26
27 school_catogories_CV = list(X_CV['school_state'].values)
28 school_list_CV = []
29 for sent in school_catogories_CV:
30     school_list_CV.append(sent.lower().strip())
31 X_CV['school_categories'] = school_list_CV
32 X_CV.drop(['school_state'], axis=1, inplace=True)
33
34 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
35 my_counter_school_CV = Counter()
36 for word in X_CV['school_categories'].values:
37     my_counter_school_CV.update(word.split())
38
39 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
40 school_dict_CV = dict(my_counter_school_CV)
41 sorted_school_dict_CV = dict(sorted(school_dict_CV.items(), key=lambda kv: kv[1]))
42 school_one_hot_CV = vectorizer_school.transform(X_CV['school_categories'].values)
43
44 #-----Vectorizing categorical data of School state for Test dataset-----
45

```

```

46 school_catogories_Test = list(X_Test['school_state'].values)
47 school_list_Test = []
48 for sent in school_catogories_Test:
49     school_list_Test.append(sent.lower().strip())
50 X_Test['school_categories'] = school_list_Test
51 X_Test.drop(['school_state'], axis=1, inplace=True)
52
53 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
54 my_counter_school_Test = Counter()
55 for word in X_Test['school_categories'].values:
56     my_counter_school_Test.update(word.split())
57
58 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
59 school_dict_Test = dict(my_counter_school_Test)
60 sorted_school_dict_Test = dict(sorted(school_dict_Test.items(), key=lambda kv: kv[1]))
61 school_one_hot_Test = vectorizer_school.transform(X_Test['school_categories'].values)
62 print("-"*120)
63 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(school_one_hot_Train.shape))
64 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(school_one_hot_CV.shape))
65 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(school_one_hot_Test.shape))

```

```

Shape of Train dataset matrix after one hot encoding is: (49041, 51)
Shape of CV dataset matrix after one hot encoding is: (24155, 51)
Shape of Test dataset matrix after one hot encoding is: (36052, 51)

```

Prefix

In [0]:

```
1 #-----Vectorizing categorical data of Teacher Prefix for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 prefix_catogories_Train = list(X_Train['teacher_prefix'].values)
8 prefix_list_Train = []
9 for sent in prefix_catogories_Train:
10     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
11     # https://gist.github.com/sebleier/554280
12     sent = ' '.join(e for e in sent.split())
13     prefix_list_Train.append(sent.lower().strip())
14 X_Train['prefix_catogories'] = prefix_list_Train
15 X_Train.drop(['teacher_prefix'], axis=1, inplace=True)
16
17 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
18 my_counter_prefix_Train = Counter()
19 for word in X_Train['prefix_catogories'].values:
20     my_counter_prefix_Train.update(word.split())
21
22 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
23 prefix_dict_Train = dict(my_counter_prefix_Train)
24 sorted_prefix_dict_Train = dict(sorted(prefix_dict_Train.items(), key=lambda kv: kv[1]))
25
26
27 vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict_Train.keys()), lowercase=False, binary=True)
28 vectorizer_prefix.fit(X_Train['prefix_catogories'].values)
29 #print(vectorizer.get_feature_names())
30
31 prefix_one_hot_Train = vectorizer_prefix.transform(X_Train['prefix_catogories'].values)
32 #print("Shape of matrix after one hot encodig ",prefix_one_hot.shape)
33
34 #-----Vectorizing categorical data of Teacher Prefix for CV dataset-----
35
36 prefix_catogories_CV = list(X_CV['teacher_prefix'].values)
37 prefix_list_CV = []
38 for sent in prefix_catogories_CV:
39     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
40     # https://gist.github.com/sebleier/554280
41     sent = ' '.join(e for e in sent.split())
42     prefix_list_CV.append(sent.lower().strip())
43 X_CV['prefix_catogories'] = prefix_list_CV
44 X_CV.drop(['teacher_prefix'], axis=1, inplace=True)
45
```

```

46 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
47 my_counter_prefix_CV = Counter()
48 for word in X_CV['prefix_catogories'].values:
49     my_counter_prefix_CV.update(word.split())
50
51 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
52 prefix_dict_CV = dict(my_counter_prefix_CV)
53 sorted_prefix_dict_CV = dict(sorted(prefix_dict_CV.items(), key=lambda kv: kv[1]))
54 prefix_one_hot_CV = vectorizer_prefix.transform(X_CV['prefix_catogories'].values)
55
56 #-----Vectorizing categorical data of Teacher Prefix for Test dataset-----
57
58 prefix_catogories_Test = list(X_Test['teacher_prefix'].values)
59 prefix_list_Test = []
60 for sent in prefix_catogories_Test:
61     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
62     # https://gist.github.com/sebleier/554280
63     sent = ' '.join(e for e in sent.split())
64     prefix_list_Test.append(sent.lower().strip())
65 X_Test['prefix_catogories'] = prefix_list_Test
66 X_Test.drop(['teacher_prefix'], axis=1, inplace=True)
67
68 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
69 my_counter_prefix_Test = Counter()
70 for word in X_Test['prefix_catogories'].values:
71     my_counter_prefix_Test.update(word.split())
72
73 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
74 prefix_dict_Test = dict(my_counter_prefix_Test)
75 sorted_prefix_dict_Test = dict(sorted(prefix_dict_Test.items(), key=lambda kv: kv[1]))
76 prefix_one_hot_Test = vectorizer_prefix.transform(X_Test['prefix_catogories'].values)
77 print("-"*120)
78 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Train.shape))
79 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_CV.shape))
80 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Test.shape))

```

```

-----
Shape of Train dataset matrix after one hot encoding is: (49041, 6)
Shape of CV dataset matrix after one hot encoding is: (24155, 6)
Shape of Test dataset matrix after one hot encoding is: (36052, 6)

```

project_grade_category

In [0]:

```
1 #-----Vectorizing categorical data of Project Grade for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 grade_catogories_Train = list(X_Train['project_grade_category'].values)
8 grade_list_Train = []
9 for sent in grade_catogories_Train:
10     sent = sent.replace('-', '_')
11     sent = sent.replace(' ', '_')
12     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
13     # https://gist.github.com/sebleier/554280
14     sent = ' '.join(e for e in sent.split())
15     grade_list_Train.append(sent.lower().strip())
16
17 # temp = temp.replace('-', '_')
18 X_Train['new_grade_category'] = grade_list_Train
19 X_Train.drop(['project_grade_category'], axis=1, inplace=True)
20
21 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
22 my_counter_grade_Train = Counter()
23 for word in X_Train['new_grade_category'].values:
24     my_counter_grade_Train.update(word.split())
25
26 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
27 grade_dict_Train = dict(my_counter_grade_Train)
28 sorted_grade_dict_Train = dict(sorted(grade_dict_Train.items(), key=lambda kv: kv[1]))
29
30 vectorizer_grade = CountVectorizer(vocabulary=list(sorted_grade_dict_Train.keys()), lowercase=False, binary=True)
31 vectorizer_grade.fit(X_Train['new_grade_category'].values)
32 #print(vectorizer.get_feature_names())
33
34 grade_one_hot_Train = vectorizer_grade.transform(X_Train['new_grade_category'].values)
35
36 #-----Vectorizing categorical data of Project Grade for CV dataset-----
37
38 grade_catogories_CV = list(X_CV['project_grade_category'].values)
39 grade_list_CV = []
40 for sent in grade_catogories_CV:
41     sent = sent.replace('-', '_')
42     sent = sent.replace(' ', '_')
43     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
44     # https://gist.github.com/sebleier/554280
45     sent = ' '.join(e for e in sent.split())
```

```

46     grade_list_CV.append(sent.lower().strip())
47
48     # temp = temp.replace('-', '_')
49     X_CV['new_grade_category'] = grade_list_CV
50     X_CV.drop(['project_grade_category'], axis=1, inplace=True)
51
52     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
53     my_counter_grade_CV = Counter()
54     for word in X_CV['new_grade_category'].values:
55         my_counter_grade_CV.update(word.split())
56
57     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
58     grade_dict_CV = dict(my_counter_grade_CV)
59     sorted_grade_dict_CV = dict(sorted(grade_dict_CV.items(), key=lambda kv: kv[1]))
60
61     grade_one_hot_CV = vectorizer_grade.transform(X_CV['new_grade_category'].values)
62
63     #-----Vectorizing categorical data of Project Grade for Train dataset-----
64
65     grade_categories_Test = list(X_Test['project_grade_category'].values)
66     grade_list_Test = []
67     for sent in grade_categories_Test:
68         sent = sent.replace('-', '_')
69         sent = sent.replace(' ', '_')
70         # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
71         # https://gist.github.com/sebleier/554280
72         sent = ' '.join(e for e in sent.split())
73         grade_list_Test.append(sent.lower().strip())
74
75     # temp = temp.replace('-', '_')
76     X_Test['new_grade_category'] = grade_list_Test
77     X_Test.drop(['project_grade_category'], axis=1, inplace=True)
78
79     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
80     my_counter_grade_Test = Counter()
81     for word in X_Test['new_grade_category'].values:
82         my_counter_grade_Test.update(word.split())
83
84     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
85     grade_dict_Test = dict(my_counter_grade_Test)
86     sorted_grade_dict_Test = dict(sorted(grade_dict_Test.items(), key=lambda kv: kv[1]))
87
88     grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)
89     print("-"*120)
90     print('Shape of Train dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Train.shape))
91     print('Shape of CV dataset matrix after one hot encoding is: {}'.format(grade_one_hot_CV.shape))

```



```
92 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Test.shape))
```

```
-----  
Shape of Train dataset matrix after one hot encoding is: (49041, 4)  
Shape of CV dataset matrix after one hot encoding is: (24155, 4)  
Shape of Test dataset matrix after one hot encoding is: (36052, 4)
```

1.5.2 Vectorizing Numerical features

```
In [0]: 1 price_data = Resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()  
2 X_Train = pd.merge(X_Train, price_data, on='id', how='left')  
3 X_CV = pd.merge(X_CV, price_data, on='id', how='left')  
4 X_Test = pd.merge(X_Test, price_data, on='id', how='left')
```

```
In [0]: 1 price_norm = Normalizer(norm='l2', copy=False)  
2 price_norm.fit(X_Train['price'].values.reshape(1,-1))  
3  
4 p=price_norm.transform(X_Train['price'].values.reshape(1,-1))  
5 price_norm.transform(X_CV['price'].values.reshape(1,-1))  
6 price_norm.transform(X_Test['price'].values.reshape(1,-1))  
7 price_norm_Train = (X_Train['price'].values.reshape(-1,1))  
8 price_norm_CV = (X_CV['price'].values.reshape(-1,1))  
9 price_norm_Test = (X_Test['price'].values.reshape(-1,1))  
10 print("-"*120)  
11 print('Shape of Train normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Train.shape))  
12 print('Shape of CV normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_CV.shape))  
13 print('Shape of Test normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Test.shape))
```

```
-----  
Shape of Train normalized price dataset matrix after one hot encoding is: (49041, 1)  
Shape of CV normalized price dataset matrix after one hot encoding is: (24155, 1)  
Shape of Test normalized price dataset matrix after one hot encoding is: (36052, 1)
```

```
In [0]: 1 quantity_norm = Normalizer(norm='l2', copy=False)
2 quantity_norm.fit(X_Train['quantity'].values.reshape(1,-1))
3
4 quantity_norm.transform(X_Train['quantity'].values.reshape(1,-1))
5 quantity_norm.transform(X_CV['quantity'].values.reshape(1,-1))
6 quantity_norm.transform(X_Test['quantity'].values.reshape(1,-1))
7 quantity_norm_Train = quantity_norm.transform(X_Train['quantity'].values.reshape(-1,1))
8 quantity_norm_CV = quantity_norm.transform(X_CV['quantity'].values.reshape(-1,1))
9 quantity_norm_Test = quantity_norm.transform(X_Test['quantity'].values.reshape(-1,1))
10 print("-"*120)
11 print('Shape of Train normalized quantity dataset matrix after one hot encoding is: {}'.format(quantity_norm_Train.shape))
12 print('Shape of CV normalized quantity dataset matrix after one hot encoding is: {}'.format(quantity_norm_CV.shape))
13 print('Shape of Test normalized quantity dataset matrix after one hot encoding is: {}'.format(quantity_norm_Test.shape))
```

```
-----
Shape of Train normalized quantity dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized quantity dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized quantity dataset matrix after one hot encoding is: (36052, 1)
```

```
In [0]: 1 teacher_prev_post_norm = Normalizer(norm='l2', copy=False)
2 teacher_prev_post_norm.fit(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
3
4 teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
5 teacher_prev_post_norm.transform(X_CV['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
6 teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
7 teacher_prev_post_norm_Train = teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
8 teacher_prev_post_norm_CV = teacher_prev_post_norm.transform(X_CV['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
9 teacher_prev_post_norm_Test = teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
10 print("-"*120)
11 print('Shape of Train normalized previously posted project dataset matrix after one hot encoding is: {}'.format(teacher_prev_post_norm_Train.shape))
12 print('Shape of CV normalized previously posted project dataset matrix after one hot encoding is: {}'.format(teacher_prev_post_norm_CV.shape))
13 print('Shape of Test normalized previously posted project dataset matrix after one hot encoding is: {}'.format(teacher_prev_post_norm_Test.shape))
```

```
-----
Shape of Train normalized previously posted project dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized previously posted project dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized previously posted project dataset matrix after one hot encoding is: (36052, 1)
```

```
In [0]: 1 title_norm = Normalizer(norm='l2', copy=False)
2 title_norm.fit(X_Train['word_count_title_Train'].values.reshape(1,-1))
3 title_norm.transform(X_Train['word_count_title_Train'].values.reshape(1,-1))
4 title_norm.transform(X_CV['word_count_title_CV'].values.reshape(1,-1))
5 title_norm.transform(X_Test['word_count_title_Test'].values.reshape(1,-1))
6 word_count_title_Train = title_norm.transform(X_Train['word_count_title_Train'].values.reshape(-1,1))
7 word_count_title_CV = title_norm.transform(X_CV['word_count_title_CV'].values.reshape(-1,1))
8 word_count_title_Test = title_norm.transform(X_Test['word_count_title_Test'].values.reshape(-1,1))
9 print("-"*120)
10 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Train.shape))
11 print('Shape of CV normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_CV.shape))
12 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Test.shape))
```

```
-----
Shape of Train normalized title dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized title dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)
```

```
In [0]: 1 essay_norm = Normalizer(norm='l2', copy=False)
2 essay_norm.fit(X_Train['word_count_essay_Train'].values.reshape(1,-1))
3 essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(1,-1))
4 essay_norm.transform(X_CV['word_count_essay_CV'].values.reshape(1,-1))
5 essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(1,-1))
6 word_count_essay_Train = essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(-1,1))
7 word_count_essay_CV = essay_norm.transform(X_CV['word_count_essay_CV'].values.reshape(-1,1))
8 word_count_essay_Test = essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(-1,1))
9 print("-"*120)
10 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Train.shape))
11 print('Shape of CV normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_CV.shape))
12 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Test.shape))
```

```
-----
Shape of Train normalized title dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized title dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)
```

1.5.3 Vectorizing Text data

1.5.3.1 Bag of words

```
In [0]: 1 # We are considering only the words which appeared in at least 10 documents(rows or projects).
2 vectorizer_essays_bow = CountVectorizer(min_df=10)
3 text_bow_Train = vectorizer_essays_bow.fit_transform(preprocessed_essays_Train)
4 text_bow_CV = vectorizer_essays_bow.transform(preprocessed_essays_CV)
5 text_bow_Test = vectorizer_essays_bow.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying Bag Of Words for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(text_bow_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(text_bow_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(text_bow_Test.shape))
```

Applying Bag Of Words for Text Data

Shape of Train dataset matrix after one hot encoding is: (49041, 12097)
Shape of CV dataset matrix after one hot encoding is: (24155, 12097)
Shape of Test dataset matrix after one hot encoding is: (36052, 12097)

Bag of Words for Project Title

```
In [0]: 1 # you can vectorize the title also
2 # before you vectorize the title make sure you preprocess it
3 vectorizer_titles_bow = CountVectorizer(min_df=10)
4 title_bow_Train = vectorizer_titles_bow.fit_transform(preprocessed_titles_Train)
5 title_bow_CV = vectorizer_titles_bow.transform(preprocessed_titles_CV)
6 title_bow_Test = vectorizer_titles_bow.transform(preprocessed_titles_Test)
7 print("-"*120)
8 print("Applying Bag Of Words for Project Title Data")
9 print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(title_bow_Train.shape))
11 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(title_bow_CV.shape))
12 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(title_bow_Test.shape))
13
```

Applying Bag Of Words for Project Title Data

Shape of Train dataset matrix after one hot encoding is: (49041, 2083)
Shape of CV dataset matrix after one hot encoding is: (24155, 2083)
Shape of Test dataset matrix after one hot encoding is: (36052, 2083)

1.5.2.2 TFIDF vectorizer

```
In [0]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer_essays_tfidf = TfidfVectorizer(min_df=10)
3 text_tfidf_Train = vectorizer_essays_tfidf.fit_transform(preprocessed_essays_Train)
4 text_tfidf_CV = vectorizer_essays_tfidf.transform(preprocessed_essays_CV)
5 text_tfidf_Test = vectorizer_essays_tfidf.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying TFIDF for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(text_tfidf_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Test.shape))
```

Applying TFIDF for Text Data

Shape of Train dataset matrix after one hot encoding is: (49041, 12097)
Shape of CV dataset matrix after one hot encoding is: (24155, 12097)
Shape of Test dataset matrix after one hot encoding is: (36052, 12097)

TFIDF vectorizer for Project Title

```
In [0]: 1 vectorizer_titles_tfidf = TfidfVectorizer(min_df=10)
2 title_tfidf_Train = vectorizer_titles_tfidf.fit_transform(preprocessed_titles_Train)
3 title_tfidf_CV = vectorizer_titles_tfidf.transform(preprocessed_titles_CV)
4 title_tfidf_Test = vectorizer_titles_tfidf.transform(preprocessed_titles_Test)
5 print("-"*120)
6 print("Applying TFIDF for Project Title")
7 print("-"*120)
8 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Train.shape))
9 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(title_tfidf_CV.shape))
10 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Test.shape))
```

Applying TFIDF for Project Title

Shape of Train dataset matrix after one hot encoding is: (49041, 2083)
Shape of CV dataset matrix after one hot encoding is: (24155, 2083)
Shape of Test dataset matrix after one hot encoding is: (36052, 2083)

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

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



In [0]:

```
1 # average Word2Vec
2 # compute average word2vec for each review.
3 avg_w2v_vectors_Train = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays_Train): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     cnt_words = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if word in glove_words:
9             vector += model[word]
10            cnt_words += 1
11    if cnt_words != 0:
12        vector /= cnt_words
13    avg_w2v_vectors_Train.append(vector)
14 #-----
15
16 avg_w2v_vectors_CV = []; # the avg-w2v for each sentence/review is stored in this list
17 for sentence in tqdm(preprocessed_essays_CV): # for each review/sentence
18     vector = np.zeros(300) # as word vectors are of zero length
19     cnt_words = 0; # num of words with a valid vector in the sentence/review
20     for word in sentence.split(): # for each word in a review/sentence
21         if word in glove_words:
22             vector += model[word]
23             cnt_words += 1
24     if cnt_words != 0:
25         vector /= cnt_words
26     avg_w2v_vectors_CV.append(vector)
27 #-----
28
29 avg_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
30 for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
31     vector = np.zeros(300) # as word vectors are of zero length
32     cnt_words = 0; # num of words with a valid vector in the sentence/review
33     for word in sentence.split(): # for each word in a review/sentence
34         if word in glove_words:
35             vector += model[word]
36             cnt_words += 1
37     if cnt_words != 0:
38         vector /= cnt_words
39     avg_w2v_vectors_Test.append(vector)
40
41 print(len(avg_w2v_vectors_Test))
42 print(len(avg_w2v_vectors_Test[1]))
```

100%|██████████| 49041/49041 [00:12<00:00, 3941.77it/s]

100%|██████████| 24155/24155 [00:06<00:00, 3925.89it/s]

100%|██████████| 36052/36052 [00:09<00:00, 3947.34it/s]

36052

300

AVG W2V on project_title

In [0]:

```
1 # Similarly you can vectorize for title also
2 # compute average word2vec for each title.
3 avg_w2v_vectors_title_Train = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_titles_Train): # for each review/sentence
5     vector_title = np.zeros(300) # as word vectors are of zero length
6     cnt_title_words = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if word in glove_words:
9             vector_title += model[word]
10            cnt_title_words += 1
11    if cnt_title_words != 0:
12        vector_title /= cnt_title_words
13    avg_w2v_vectors_title_Train.append(vector_title)
14
15
16 #-----
17 avg_w2v_vectors_title_CV = []; # the avg-w2v for each sentence/review is stored in this list
18 for sentence in tqdm(preprocessed_titles_CV): # for each review/sentence
19     vector_title = np.zeros(300) # as word vectors are of zero length
20     cnt_title_words = 0; # num of words with a valid vector in the sentence/review
21     for word in sentence.split(): # for each word in a review/sentence
22         if word in glove_words:
23             vector_title += model[word]
24             cnt_title_words += 1
25     if cnt_title_words != 0:
26         vector_title /= cnt_title_words
27     avg_w2v_vectors_title_CV.append(vector_title)
28
29 #-----
30 avg_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
31 for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
32     vector_title = np.zeros(300) # as word vectors are of zero length
33     cnt_title_words = 0; # num of words with a valid vector in the sentence/review
34     for word in sentence.split(): # for each word in a review/sentence
35         if word in glove_words:
36             vector_title += model[word]
37             cnt_title_words += 1
38     if cnt_title_words != 0:
39         vector_title /= cnt_title_words
40     avg_w2v_vectors_title_Test.append(vector_title)
41
42 print(len(avg_w2v_vectors_title_Test))
43 print(len(avg_w2v_vectors_title_Test[0]))
```

100%|██████████| 24155/24155 [00:00<00:00, 66048.26it/s]
100%|██████████| 36052/36052 [00:00<00:00, 66043.95it/s]

36052
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

```
1 tfidf_model_essays = TfidfVectorizer()  
2 tfidf_model_essays.fit(preprocessed_essays_Train)  
3 # we are converting a dictionary with word as a key, and the idf as a value  
4 dictionary = dict(zip(tfidf_model_essays.get_feature_names(), list(tfidf_model_essays.idf_)))  
5 tfidf_words_essays = set(tfidf_model_essays.get_feature_names())
```

In [0]:

```
1 # average Word2Vec
2 # compute average word2vec for each review.
3 tfidf_w2v_vectors_Train = []; # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays_Train): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if (word in glove_words) and (word in tfidf_words_essays):
9             vec = model[word] # getting the vector for each word
10             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
11             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
12             vector += (vec * tf_idf) # calculating tfidf weighted w2v
13             tf_idf_weight += tf_idf
14     if tf_idf_weight != 0:
15         vector /= tf_idf_weight
16     tfidf_w2v_vectors_Train.append(vector)
17
18 #-----
19 tfidf_w2v_vectors_CV = []; # the avg-w2v for each sentence/review is stored in this list
20 for sentence in tqdm(preprocessed_essays_CV): # for each review/sentence
21     vector = np.zeros(300) # as word vectors are of zero length
22     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
23     for word in sentence.split(): # for each word in a review/sentence
24         if (word in glove_words) and (word in tfidf_words_essays):
25             vec = model[word] # getting the vector for each word
26             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
27             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
28             vector += (vec * tf_idf) # calculating tfidf weighted w2v
29             tf_idf_weight += tf_idf
30     if tf_idf_weight != 0:
31         vector /= tf_idf_weight
32     tfidf_w2v_vectors_CV.append(vector)
33 #-----
34 tfidf_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
35 for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
36     vector = np.zeros(300) # as word vectors are of zero length
37     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
38     for word in sentence.split(): # for each word in a review/sentence
39         if (word in glove_words) and (word in tfidf_words_essays):
40             vec = model[word] # getting the vector for each word
41             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))
42             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
43             vector += (vec * tf_idf) # calculating tfidf weighted w2v
44             tf_idf_weight += tf_idf
45     if tf_idf_weight != 0:
```

```
46         vector /= tf_idf_weight
47         tfidf_w2v_vectors_Test.append(vector)
48
49     print(len(tfidf_w2v_vectors_Test))
50     print(len(tfidf_w2v_vectors_Test[0]))
51
```

```
100%|██████████| 49041/49041 [01:21<00:00, 599.75it/s]
100%|██████████| 24155/24155 [00:40<00:00, 597.22it/s]
100%|██████████| 36052/36052 [01:00<00:00, 596.10it/s]
```

36052

300

Using Pretrained Models: TFIDF weighted W2V on project_title

In [0]:

```
1 # Similarly you can vectorize for title also
2 tfidf_model_title = TfidfVectorizer()
3 tfidf_model_title.fit(preprocessed_titles_Train)
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(tfidf_model_title.get_feature_names(), list(tfidf_model_title.idf_)))
6 tfidf_words_title = set(tfidf_model_title.get_feature_names())
7
8 # compute tfidf word2vec for each title.
9 tfidf_w2v_vectors_title_Train = []; # the avg-w2v for each sentence/review is stored in this list
10 for sentence in tqdm(preprocessed_titles_Train): # for each review/sentence
11     vector_title = np.zeros(300) # as word vectors are of zero length
12     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
13     for word in sentence.split(): # for each word in a review/sentence
14         if (word in glove_words) and (word in tfidf_words_title):
15             vec = model[word] # getting the vector for each word
16             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
17             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
18             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
19             tf_idf_weight += tf_idf
20     if tf_idf_weight != 0:
21         vector_title /= tf_idf_weight
22     tfidf_w2v_vectors_title_Train.append(vector_title)
23 #-----
24
25 tfidf_w2v_vectors_title_CV = []; # the avg-w2v for each sentence/review is stored in this list
26 for sentence in tqdm(preprocessed_titles_CV): # for each review/sentence
27     vector_title = np.zeros(300) # as word vectors are of zero length
28     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
29     for word in sentence.split(): # for each word in a review/sentence
30         if (word in glove_words) and (word in tfidf_words_title):
31             vec = model[word] # getting the vector for each word
32             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
33             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
34             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
35             tf_idf_weight += tf_idf
36     if tf_idf_weight != 0:
37         vector_title /= tf_idf_weight
38     tfidf_w2v_vectors_title_CV.append(vector_title)
39 #-----
40
41
42 tfidf_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
43 for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
44     vector_title = np.zeros(300) # as word vectors are of zero length
45     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
```

```

46     for word in sentence.split(): # for each word in a review/sentence
47         if (word in glove_words) and (word in tfidf_words_title):
48             vec = model[word] # getting the vector for each word
49             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
50             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
51             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
52             tf_idf_weight += tf_idf
53         if tf_idf_weight != 0:
54             vector_title /= tf_idf_weight
55             tfidf_w2v_vectors_title_Test.append(vector_title)
56
57 print(len(tfidf_w2v_vectors_title_Test))
58 print(len(tfidf_w2v_vectors_title_Test[0]))
59
60

```

```

100%|██████████| 49041/49041 [00:01<00:00, 51119.08it/s]
100%|██████████| 24155/24155 [00:00<00:00, 29048.77it/s]
100%|██████████| 36052/36052 [00:01<00:00, 32121.15it/s]

```

36052

300

Calculating the sentiment score's of each of the essay

In [0]:

```
1
2 sid = SentimentIntensityAnalyzer()
3
4 essays = X_Train['essay']
5 essays_sentiment_TR_P = []
6 essays_sentiment_TR_N = []
7 essays_sentiment_TR_NE = []
8 essays_sentiment_TR_C = []
9 for essay in tqdm(essays):
10     res = sid.polarity_scores(essay)
11     essays_sentiment_TR_P.append(res['pos'])
12     essays_sentiment_TR_N.append(res['neg'])
13     essays_sentiment_TR_NE.append(res['neu'])
14     essays_sentiment_TR_C.append(res['compound'])
15 X_Train['sentiment_essay_TR_P'] = essays_sentiment_TR_P
16 X_Train['sentiment_essay_TR_N'] = essays_sentiment_TR_N
17 X_Train['sentiment_essay_TR_NE'] = essays_sentiment_TR_NE
18 X_Train['sentiment_essay_TR_C'] = essays_sentiment_TR_C
19
20 essays = X_CV['essay']
21 essays_sentiment_CV_P = []
22 essays_sentiment_CV_N = []
23 essays_sentiment_CV_NE = []
24 essays_sentiment_CV_C = []
25 for essay in tqdm(essays):
26     res = sid.polarity_scores(essay)
27     essays_sentiment_CV_P.append(res['pos'])
28     essays_sentiment_CV_N.append(res['neg'])
29     essays_sentiment_CV_NE.append(res['neu'])
30     essays_sentiment_CV_C.append(res['compound'])
31 X_CV['sentiment_essay_CV_P'] = essays_sentiment_CV_P
32 X_CV['sentiment_essay_CV_N'] = essays_sentiment_CV_N
33 X_CV['sentiment_essay_CV_NE'] = essays_sentiment_CV_NE
34 X_CV['sentiment_essay_CV_C'] = essays_sentiment_CV_C
35
36 essays = X_Test['essay']
37 essays_sentiment_TS_P = []
38 essays_sentiment_TS_N = []
39 essays_sentiment_TS_NE = []
40 essays_sentiment_TS_C = []
41 for essay in tqdm(essays):
42     res = sid.polarity_scores(essay)
43     essays_sentiment_TS_P.append(res['pos'])
44     essays_sentiment_TS_N.append(res['neg'])
45     essays_sentiment_TS_NE.append(res['neu'])
```

```

46 essays_sentiment_TS_C.append(res['compound'])
47 X_Test['sentiment_essay_TS_P'] = essays_sentiment_TS_P
48 X_Test['sentiment_essay_TS_N'] = essays_sentiment_TS_N
49 X_Test['sentiment_essay_TS_NE'] = essays_sentiment_TS_NE
50 X_Test['sentiment_essay_TS_C'] = essays_sentiment_TS_C
51
52 sentiment_norm_P = Normalizer(norm='l2', copy=False)
53 sentiment_norm_N = Normalizer(norm='l2', copy=False)
54 sentiment_norm_NE = Normalizer(norm='l2', copy=False)
55 sentiment_norm_C = Normalizer(norm='l2', copy=False)
56
57
58 sentiment_norm_P.fit(X_Train['sentiment_essay_TR_P'].values.reshape(1,-1))
59 sentiment_norm_N.fit(X_Train['sentiment_essay_TR_N'].values.reshape(1,-1))
60 sentiment_norm_NE.fit(X_Train['sentiment_essay_TR_NE'].values.reshape(1,-1))
61 sentiment_norm_C.fit(X_Train['sentiment_essay_TR_C'].values.reshape(1,-1))
62
63 sentiment_Train_P = sentiment_norm_P.transform(X_Train['sentiment_essay_TR_P'].values.reshape(1,-1))
64 sentiment_CV_P = sentiment_norm_P.transform(X_CV['sentiment_essay_CV_P'].values.reshape(1,-1))
65 sentiment_Test_P = sentiment_norm_P.transform(X_Test['sentiment_essay_TS_P'].values.reshape(1,-1))
66 sentiment_Train_N = sentiment_norm_N.transform(X_Train['sentiment_essay_TR_N'].values.reshape(1,-1))
67 sentiment_CV_N = sentiment_norm_N.transform(X_CV['sentiment_essay_CV_N'].values.reshape(1,-1))
68 sentiment_Test_N = sentiment_norm_N.transform(X_Test['sentiment_essay_TS_N'].values.reshape(1,-1))
69 sentiment_Train_NE = sentiment_norm_NE.transform(X_Train['sentiment_essay_TR_NE'].values.reshape(1,-1))
70 sentiment_CV_NE = sentiment_norm_NE.transform(X_CV['sentiment_essay_CV_NE'].values.reshape(1,-1))
71 sentiment_Test_NE = sentiment_norm_NE.transform(X_Test['sentiment_essay_TS_NE'].values.reshape(1,-1))
72 sentiment_Train_C = sentiment_norm_C.transform(X_Train['sentiment_essay_TR_C'].values.reshape(1,-1))
73 sentiment_CV_C = sentiment_norm_C.transform(X_CV['sentiment_essay_CV_C'].values.reshape(1,-1))
74 sentiment_Test_C = sentiment_norm_C.transform(X_Test['sentiment_essay_TS_C'].values.reshape(1,-1))
75
76 sentiment_Train_P = (X_Train['sentiment_essay_TR_P'].values.reshape(-1,1))
77 sentiment_CV_P = (X_CV['sentiment_essay_CV_P'].values.reshape(-1,1))
78 sentiment_Test_P = (X_Test['sentiment_essay_TS_P'].values.reshape(-1,1))
79 sentiment_Train_N = (X_Train['sentiment_essay_TR_N'].values.reshape(-1,1))
80 sentiment_CV_N = (X_CV['sentiment_essay_CV_N'].values.reshape(-1,1))
81 sentiment_Test_N = (X_Test['sentiment_essay_TS_N'].values.reshape(-1,1))
82 sentiment_Train_NE = (X_Train['sentiment_essay_TR_NE'].values.reshape(-1,1))
83 sentiment_CV_NE = (X_CV['sentiment_essay_CV_NE'].values.reshape(-1,1))
84 sentiment_Test_NE = (X_Test['sentiment_essay_TS_NE'].values.reshape(-1,1))
85 sentiment_Train_C = (X_Train['sentiment_essay_TR_C'].values.reshape(-1,1))
86 sentiment_CV_C = (X_CV['sentiment_essay_CV_C'].values.reshape(-1,1))
87 sentiment_Test_C = (X_Test['sentiment_essay_TS_C'].values.reshape(-1,1))
88
89
90 print("Shape of sentiment Train matrix after one hot encodig ",sentiment_Train_P.shape)
91 print("Shape of sentiment CV matrix after one hot encodig ",sentiment_CV_P.shape)

```



```

92 print("Shape of sentiment Test matrix after one hot encodig ",sentiment_Test_P.shape)
93 print("Shape of sentiment Train matrix after one hot encodig ",sentiment_Train_N.shape)
94 print("Shape of sentiment CV matrix after one hot encodig ",sentiment_CV_N.shape)
95 print("Shape of sentiment Test matrix after one hot encodig ",sentiment_Test_N.shape)
96 print("Shape of sentiment Train matrix after one hot encodig ",sentiment_Train_NE.shape)
97 print("Shape of sentiment CV matrix after one hot encodig ",sentiment_CV_NE.shape)
98 print("Shape of sentiment Test matrix after one hot encodig ",sentiment_Test_NE.shape)
99 print("Shape of sentiment Train matrix after one hot encodig ",sentiment_Train_C.shape)
100 print("Shape of sentiment CV matrix after one hot encodig ",sentiment_CV_C.shape)
101 print("Shape of sentiment Test matrix after one hot encodig ",sentiment_Test_C.shape)

100%|██████████| 49041/49041 [02:03<00:00, 397.69it/s]
100%|██████████| 24155/24155 [01:00<00:00, 396.71it/s]
100%|██████████| 36052/36052 [01:31<00:00, 395.80it/s]

```

```

Shape of sentiment Train matrix after one hot encodig (49041, 1)
Shape of sentiment CV matrix after one hot encodig (24155, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (49041, 1)
Shape of sentiment CV matrix after one hot encodig (24155, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (49041, 1)
Shape of sentiment CV matrix after one hot encodig (24155, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (49041, 1)
Shape of sentiment CV matrix after one hot encodig (24155, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)

```

Assignment 7: SVM

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

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_essay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_essay (TFIDF)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_essay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

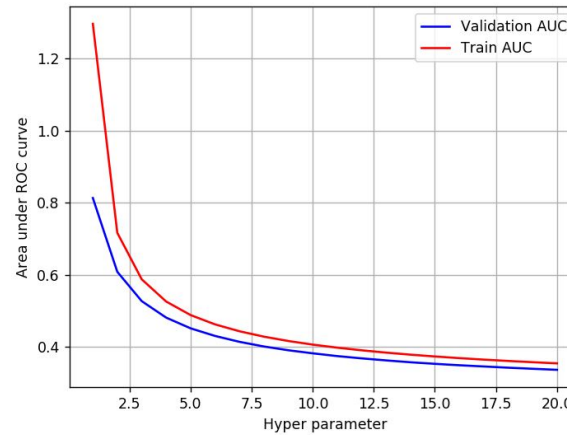
2. The hyper paramter tuning (best alpha in range 10^{-4} to 10^4 , and the best penalty among 'l1', 'l2')

- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](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

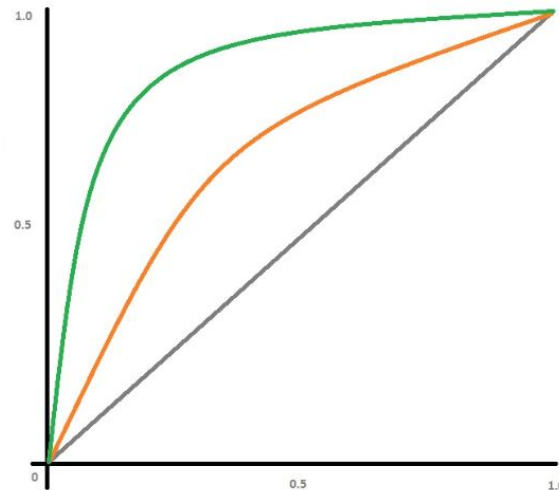
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

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/) (<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](#).

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

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

4. [Task-2] Apply the Support Vector Machines on these features by finding the best hyper paramter as suggested in step 2 and step 3

- Consider these set of features **Set 5** :
 - school_state** : categorical data
 - clean_categories** : categorical data
 - clean_subcategories** : categorical data
 - project_grade_category** :categorical data
 - teacher_prefix** : categorical data
 - quantity** : numerical data
 - teacher_number_of_previously_posted_projects** : numerical data
 - price** : numerical data
 - sentiment score's of each of the essay** : numerical data
 - number of words in the title** : numerical data
 - number of words in the combine essays** : numerical data
 - Apply **TruncatedSVD** (<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>) on **TfidfVectorizer** (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html) of essay text, choose the number of components (**n_components**) using **elbow method** (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/>) : numerical data
- Conclusion**
 - You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link](http://zetcode.com/python/prettytable/) (<http://zetcode.com/python/prettytable/>)

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	7	0.78
TFIDF	Brute	12	0.79
W2V	Brute	10	0.78
TFIDFW2V	Brute	6	0.78

Note: Data Leakage

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

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

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

1.5.4 Merging all the above features

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

In [0]:

```
1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 BOW_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one_
3 print(BOW_Train.shape)
4 TFIDF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one
5 print(TFIDF_Train.shape)
6 AVG_W2V_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_c
7 print(AVG_W2V_Train.shape)
8 TFIDF_W2V_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefi
9 print(TFIDF_W2V_Train.shape)
10 #FiveF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_or
11 #print(FiveF_Train.shape)
```

(49041, 14289)

(49041, 14289)

(49041, 709)

(49041, 709)

In [0]:

```
1 BOW_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,text_bow
2 print(BOW_CV.shape)
3 TFIDF_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,text_t
4 print(TFIDF_CV.shape)
5 AVG_W2V_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,avg_v
6 print(AVG_W2V_CV.shape)
7 TFIDF_W2V_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,tf
8 print(TFIDF_W2V_CV.shape)
9 #FiveF_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV, price
10 #print(FiveF_CV.shape)
```

(24155, 14289)

(24155, 14289)

(24155, 709)

(24155, 709)

```
In [0]: 1 BOW_Test = hstack((categories_one_hot_Test, sub_categories_one_hot_Test, school_one_hot_Test, grade_one_hot_Test, prefix_one_hot_Test))
2 print(BOW_Test.shape)
3 TFIDF_Test = hstack((categories_one_hot_Test, sub_categories_one_hot_Test, school_one_hot_Test, grade_one_hot_Test, prefix_one_hot_Test))
4 print(TFIDF_Test.shape)
5 AVG_W2V_Test = hstack((categories_one_hot_Test, sub_categories_one_hot_Test, school_one_hot_Test, grade_one_hot_Test, prefix_one_hot_Test))
6 print(AVG_W2V_Test.shape)
7 TFIDF_W2V_Test = hstack((categories_one_hot_Test, sub_categories_one_hot_Test, school_one_hot_Test, grade_one_hot_Test, prefix_one_hot_Test))
8 print(TFIDF_W2V_Test.shape)
9 #FiveF_Test = hstack((categories_one_hot_Test, sub_categories_one_hot_Test, school_one_hot_Test, grade_one_hot_Test, prefix_one_hot_Test))
10 #print(FiveF_Test.shape)
```

(36052, 14289)

(36052, 14289)

(36052, 709)

(36052, 709)

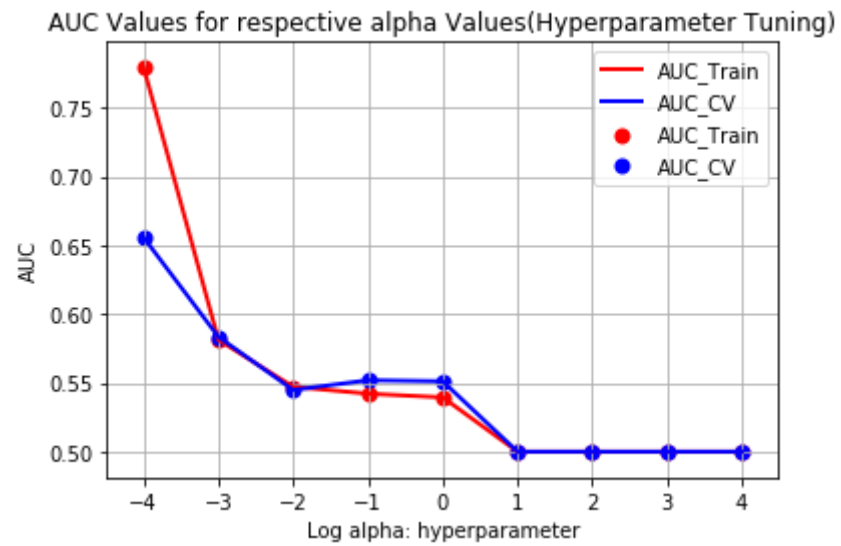
```
In [0]: 1 import pdb
2
3 def batch_predict(clf, data):
4     y_data_pred = []
5     tr_loop = data.shape[0] - data.shape[0]%1000
6     for i in range(0, tr_loop, 1000):
7         y_data_pred.extend(clf.decision_function(data[i:i+1000]))
8     y_data_pred.extend(clf.decision_function(data[tr_loop:]))
9     #pdb.set_trace()
10    return y_data_pred
```

Applying Support Vector Machines(SVM) on BOW with L1 regularization, SET 1

In [0]:

```
1 %%time
2
3 BOW_TR_CSR = BOW_Train.tocsr()
4 BOW_CV_CSR = BOW_CV.tocsr()
5 BOW_TS_CSR = BOW_Test.tocsr()
6
7 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
8 L_C_Para= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(C_Para):
14     SGC = SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1,alpha=i)
15     SGC.fit(BOW_Train, Y_Train)
16     pred = SGC.predict(BOW_CV)
17     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18     ACCV.append(acc)
19
20     Train_pred = batch_predict(SGC, BOW_TR_CSR)
21     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
22     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
23
24     CV_pred = batch_predict(SGC, BOW_CV_CSR)
25     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
26
27     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
28
29
30 for av in tqdm(C_Para):
31     b = np.log10(av)
32     L_C_Para.append(b)
33
34 # Performance of model on Train data and Test data for each hyper parameter.
35 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
36 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
37 plt.gca()
38 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
39 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
40 plt.gca()
41 plt.legend()
42 plt.xlabel("Log alpha: hyperparameter")
43 plt.ylabel("AUC")
44 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
45 plt.grid()
```

```
46 plt.show()
100%|██████████| 9/9 [00:20<00:00, 1.53s/it]
100%|██████████| 9/9 [00:00<00:00, 32044.77it/s]
```



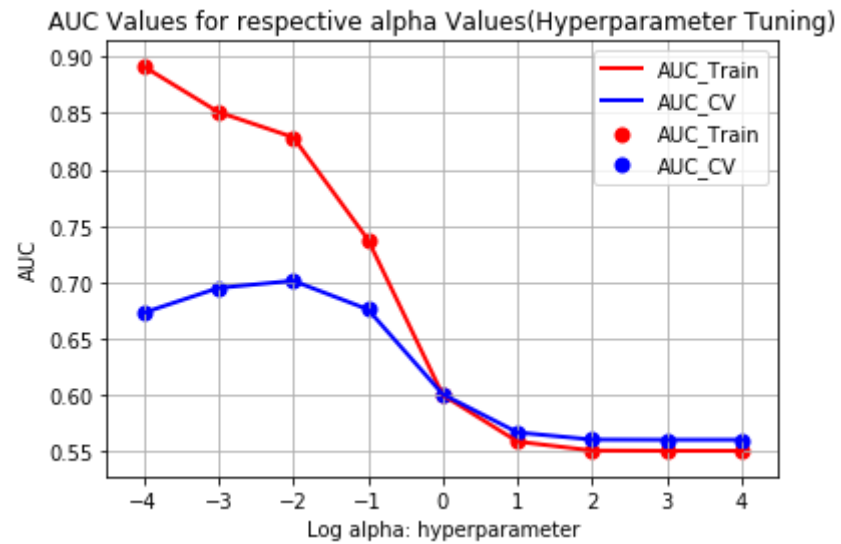
CPU times: user 26.3 s, sys: 30.5 s, total: 56.8 s
Wall time: 21.3 s

Applying Support Vector Machines(SVM) on BOW with L2 regularization, SET 1

In [0]:

```
1 %%time
2
3 BOW_TR_CSR = BOW_Train.tocsr()
4 BOW_CV_CSR = BOW_CV.tocsr()
5 BOW_TS_CSR = BOW_Test.tocsr()
6
7 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
8 L_C_Para= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(C_Para):
14
15     SGC = SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1,alpha=i)
16     SGC.fit(BOW_Train, Y_Train)
17     pred = SGC.predict(BOW_CV)
18     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
19     ACCV.append(acc)
20
21     Train_pred = batch_predict(SGC, BOW_TR_CSR)
22     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
23     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
24
25     CV_pred = batch_predict(SGC, BOW_CV_CSR)
26     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
27
28     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
29
30
31 for av in tqdm(C_Para):
32     b = np.log10(av)
33     L_C_Para.append(b)
34
35 # Performance of model on Train data and Test data for each hyper parameter.
36 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
37 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
38 plt.gca()
39 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
40 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
41 plt.gca()
42 plt.legend()
43 plt.xlabel("Log alpha: hyperparameter")
44 plt.ylabel("AUC")
45 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
```

```
46 plt.grid()
47 plt.show()
100%|██████████| 9/9 [00:08<00:00, 1.25it/s]
100%|██████████| 9/9 [00:00<00:00, 43389.35it/s]
```



CPU times: user 13.8 s, sys: 28.3 s, total: 42.1 s
Wall time: 9.25 s

OBSERVATION:

1. α value has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha = 0.01$ is having the highest AUC CV score.
3. Hence $\alpha = 0.01$ is choosen as a best α value.

In [0]:

```
1 SGC_OPT= 0.01
2 BOW_opt = SGDClassifier(penalty='l2', alpha=SGC_OPT, class_weight='balanced', n_jobs=-1)
3 BOW_opt.fit(BOW_Train, Y_Train)
4 pred = BOW_opt.predict(BOW_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha Value = {0} is {1}%'.format(SGC_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, BOW_opt.decision_function(BOW_TR_CSR))
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, BOW_opt.decision_function(BOW_TS_CSR))
```

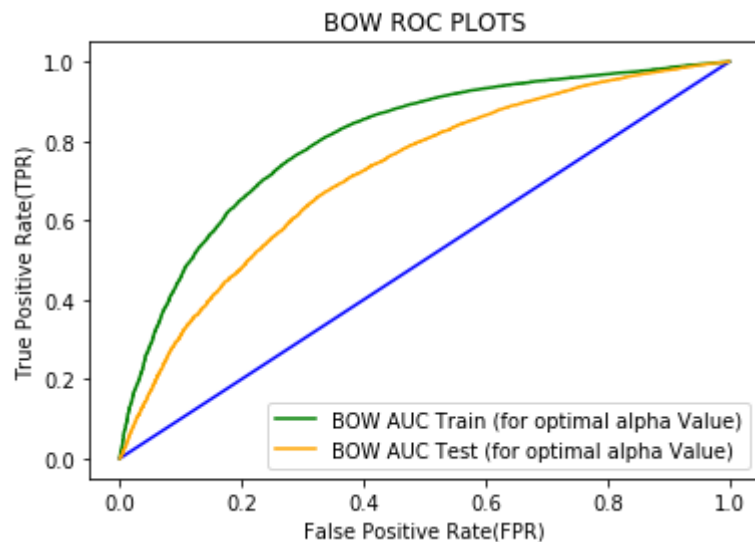
Test accuracy for alpha Value = 0.01 is 66.8118273604793%

BOW ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="BOW AUC Train (for optimal alpha Value)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="BOW AUC Test (for optimal alpha Value)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("BOW ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal alpha Value) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal alpha Value) =", auc(a_fpr_Test, a_tpr_Test))
16 BOW_AOPT=SGC_OPT
17 BOW_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred1 = BOW_opt.predict(BOW_Train)
19 pred2 = BOW_opt.predict(BOW_Test)
20

```



```

AUC Train (for optimal alpha Value) = 0.8043252832934277
AUC Test (for optimal alpha Value) = 0.7185840057086738
CPU times: user 651 ms, sys: 7.82 ms, total: 659 ms
Wall time: 656 ms

```

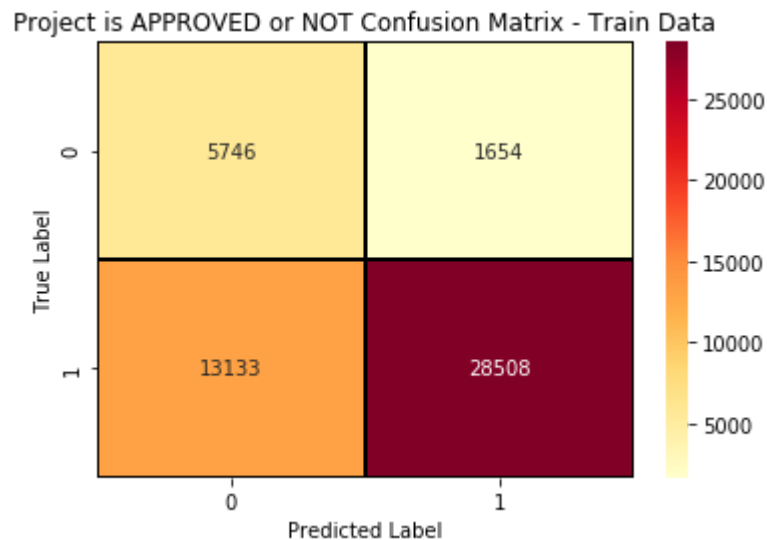
BOW CONFUSION MATRIX

In [0]:

```
1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred1)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

CPU times: user 104 ms, sys: 34.1 ms, total: 139 ms

Wall time: 94.2 ms



OBSERVATION:

True Negative =5746; False Negative =13133; True Positive = 28508; False Positive = 1654

Accuracy (Overall, how often is the classifier correct) = 0.70

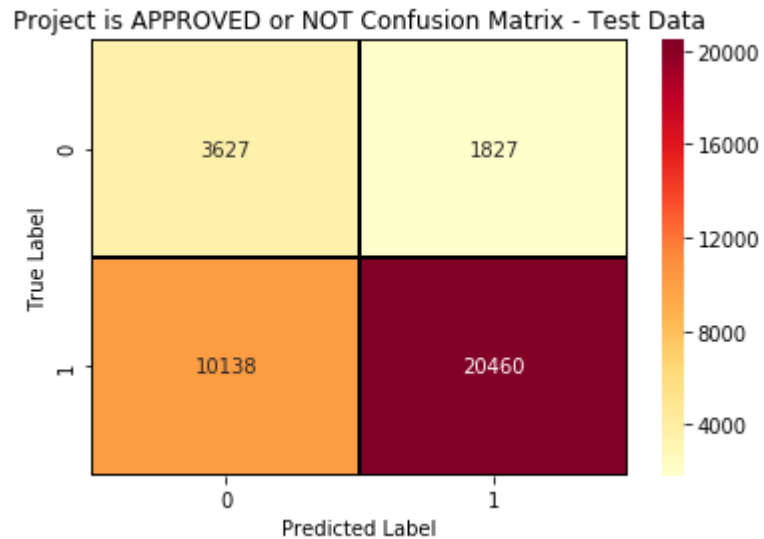
Precision(When it predicts yes, how often is it correct) =0.95

Misclassification (Overall, how often is it wrong) =0.31

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred2)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

CPU times: user 93.9 ms, sys: 37.3 ms, total: 131 ms

Wall time: 85.3 ms



OBSERVATION:

True Negative =3627; False Negative =10138; True Positive = 20460; False Positive = 1827

Accuracy (Overall, how often is the classifier correct) = 0.67

Precision(When it predicts yes, how often is it correct) =0.92

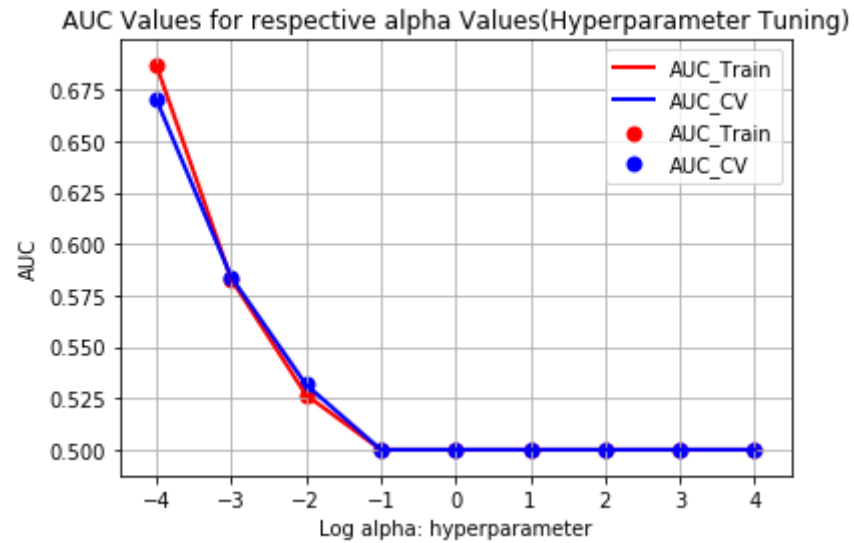
Misclassification (Overall, how often is it wrong) =0.34

Applying Support Vector Machines(SVM) on TFIDF with L1 regularization, SET 2

In [0]:

```
1 %%time
2
3 TFIDF_TR_CSR = TFIDF_Train.tocsr()
4 TFIDF_CV_CSR = TFIDF_CV.tocsr()
5 TFIDF_TS_CSR = TFIDF_Test.tocsr()
6
7 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
8 L_C_Para= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(C_Para):
14     SGC = SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1,alpha=i)
15     SGC.fit(TFIDF_Train, Y_Train)
16     pred = SGC.predict(TFIDF_CV)
17     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18     ACCV.append(acc)
19     #a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, SGC.predict_proba(TFIDF_TR_CSR)[:,1])
20     Train_pred = batch_predict(SGC, TFIDF_TR_CSR)
21     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
22     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
23
24     #a_fpr_cv, a_tpr_cv, thresholds = roc_curve(Y_CV, SGC.predict_proba(TFIDF_CV_CSR)[:,1])
25     CV_pred = batch_predict(SGC, TFIDF_CV_CSR)
26     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
27     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
28
29
30 for av in tqdm(C_Para):
31     b = np.log10(av)
32     L_C_Para.append(b)
33
34 # Performance of model on Train data and Test data for each hyper parameter.
35 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
36 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
37 plt.gca()
38 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
39 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
40 plt.gca()
41 plt.legend()
42 plt.xlabel("Log alpha: hyperparameter")
43 plt.ylabel("AUC")
44 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
45 plt.grid()
```

```
46 plt.show()
100%|██████████| 9/9 [00:08<00:00, 1.13it/s]
100%|██████████| 9/9 [00:00<00:00, 24120.60it/s]
```



CPU times: user 13.4 s, sys: 28.4 s, total: 41.8 s
Wall time: 8.99 s

Applying Support Vector Machines(SVM) on TFIDF with L2 regularization, SET 2

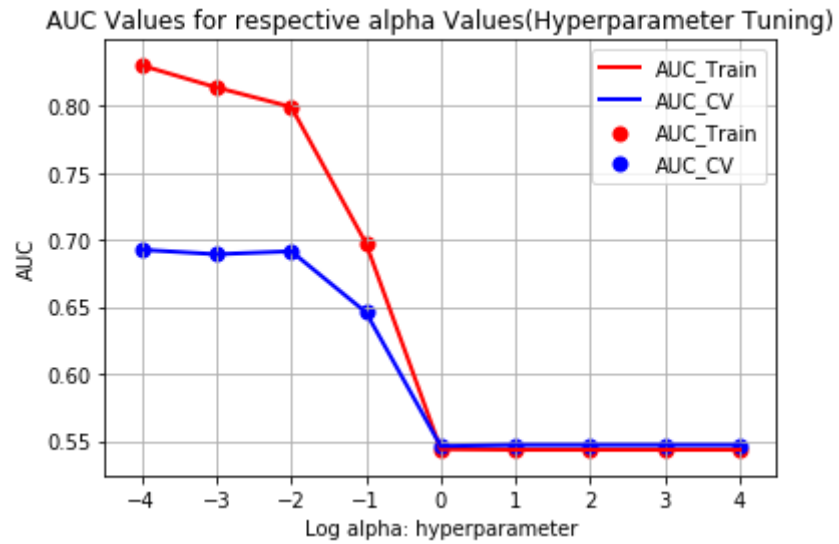
In [0]:

```
1 %%time
2
3 TFIDF_TR_CSR = TFIDF_Train.tocsr()
4 TFIDF_CV_CSR = TFIDF_CV.tocsr()
5 TFIDF_TS_CSR = TFIDF_Test.tocsr()
6
7 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
8 L_C_Para= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(C_Para):
14     SGC = SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1,alpha=i)
15     SGC.fit(TFIDF_Train, Y_Train)
16     pred = SGC.predict(TFIDF_CV)
17     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
18     ACCV.append(acc)
19     #a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, SGC.predict_proba(TFIDF_TR_CSR)[:,1])
20     Train_pred = batch_predict(SGC, TFIDF_TR_CSR)
21     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
22     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
23
24     #a_fpr_cv, a_tpr_cv, thresholds = roc_curve(Y_CV, SGC.predict_proba(TFIDF_CV_CSR)[:,1])
25     CV_pred = batch_predict(SGC, TFIDF_CV_CSR)
26     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
27     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
28
29
30 for av in tqdm(C_Para):
31     b = np.log10(av)
32     L_C_Para.append(b)
33
34 # Performance of model on Train data and Test data for each hyper parameter.
35 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
36 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
37 plt.gca()
38 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
39 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
40 plt.gca()
41 plt.legend()
42 plt.xlabel("Log alpha: hyperparameter")
43 plt.ylabel("AUC")
44 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
45 plt.grid()
```

```
46 plt.show()
```

```
100%|██████████| 9/9 [00:06<00:00, 1.37it/s]
```

```
100%|██████████| 9/9 [00:00<00:00, 20482.22it/s]
```



CPU times: user 12 s, sys: 29.8 s, total: 41.9 s

Wall time: 7.42 s

OBSERVATION:

1. α value has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha = 0.0001$ is having the highest AUC CV score.
3. Hence $\alpha = 0.0001$ is choosen as a best α value.

```
In [0]: 1 SGC_OPT=0.0001
2 TFIDF_opt =SGDClassifier(penalty='l2', alpha=SGC_OPT, class_weight='balanced', n_jobs=-1)
3
4 TFIDF_opt.fit(TFIDF_Train, Y_Train)
5 pred = TFIDF_opt.predict(TFIDF_Test)
6 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
7 print('\nTest accuracy for alpha Value = {0} is {1}%'.format(SGC_OPT,acc))
8
9 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, TFIDF_opt.decision_function(TFIDF_TR_CSR))
10 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, TFIDF_opt.decision_function(TFIDF_TS_CSR))
```

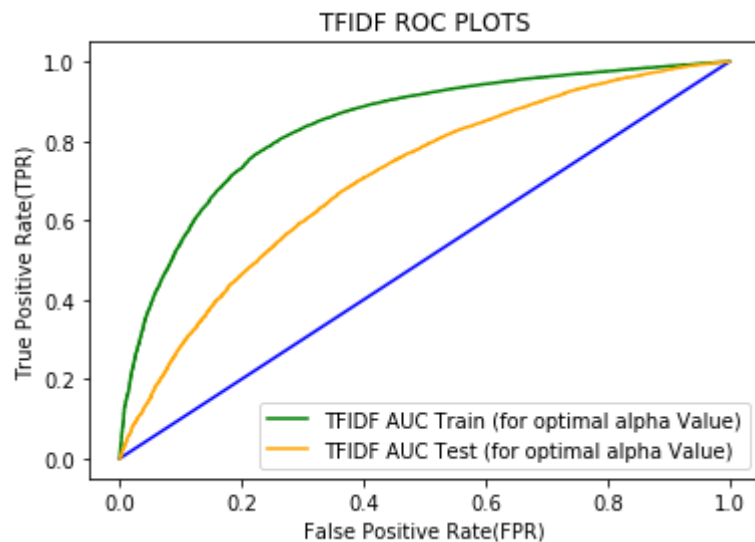
Test accuracy for alpha Value = 0.0001 is 74.9112393209808%

TFIDF ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="TFIDF AUC Train (for optimal alpha Value)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="TFIDF AUC Test (for optimal alpha Value)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal alpha Value) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal alpha Value) =", auc(a_fpr_Test, a_tpr_Test))
16 TFIDF_AOPT=SGC_OPT
17 TFIDF_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred3 = TFIDF_opt.predict(TFIDF_Train)
19 pred4 = TFIDF_opt.predict(TFIDF_Test)
20

```



AUC Train (for optimal alpha Value) = 0.8392740895959478
 AUC Test (for optimal alpha Value) = 0.7051983991130664
 CPU times: user 578 ms, sys: 7.82 ms, total: 586 ms
 Wall time: 581 ms

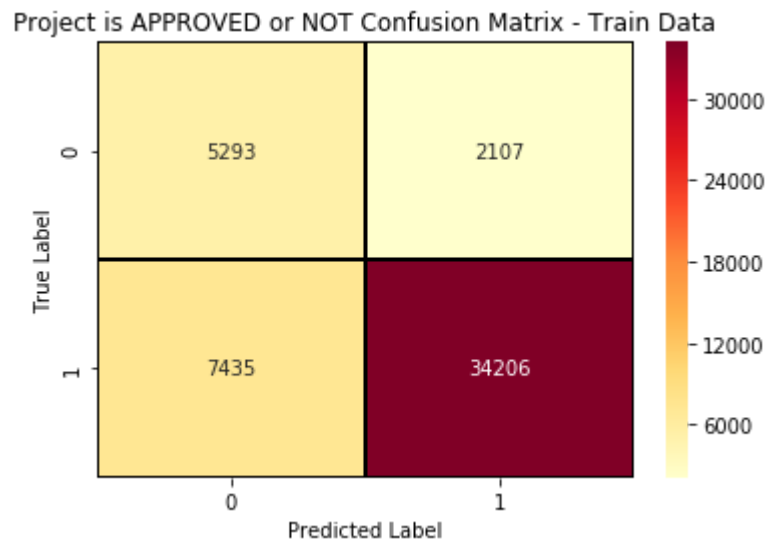
TFIDF CONFUSION MATRIX

In [0]:

```
1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred3)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

CPU times: user 100 ms, sys: 35.3 ms, total: 136 ms

Wall time: 92.2 ms



OBSERVATION:

True Negative =5293; False Negative =7435; True Positive = 34206; False Positive = 2107

Accuracy (Overall, how often is the classifier correct) = 0.81

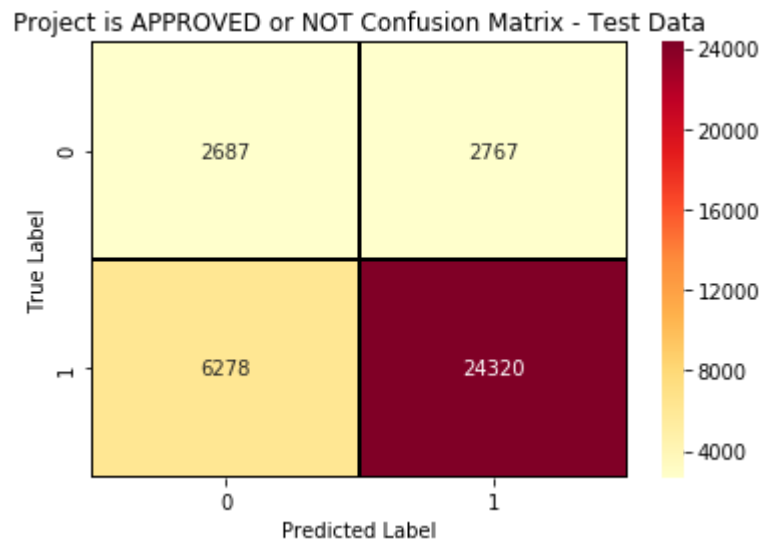
Precision(When it predicts yes, how often is it correct) =0.95

Misclassification (Overall, how often is it wrong) =0.20

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred4)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

CPU times: user 91.9 ms, sys: 38 ms, total: 130 ms

Wall time: 83.1 ms



OBSERVATION:

True Negative =2687; False Negative =6278; True Positive = 24320; False Positive = 2767

Accuracy (Overall, how often is the classifier correct) = 0.75

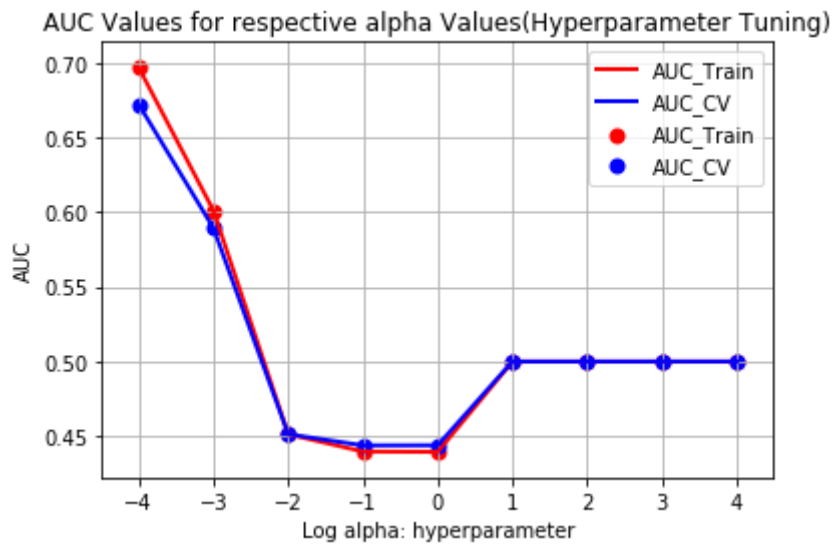
Precision(When it predicts yes, how often is it correct) =0.90

Misclassification (Overall, how often is it wrong) =0.26

Applying Support Vector Machines(SVM) on AVG W2V with L1 regularization, SET 3

In [0]:

```
1 %%time
2 AVG_W2V_TR_CSR = AVG_W2V_Train.tocsr()
3 AVG_W2V_CV_CSR = AVG_W2V_CV.tocsr()
4 AVG_W2V_TS_CSR = AVG_W2V_Test.tocsr()
5
6 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
7 L_C_Para= []
8 ACCV = []
9 AUC_TR = []
10 AUC_CV = []
11
12 for i in tqdm(C_Para):
13     SGC = SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1,alpha=i)
14     SGC.fit(AVG_W2V_Train, Y_Train)
15     pred = SGC.predict(AVG_W2V_CV)
16     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
17     ACCV.append(acc)
18     Train_pred = batch_predict(SGC, AVG_W2V_TR_CSR)
19     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
20     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
21
22     CV_pred = batch_predict(SGC, AVG_W2V_CV_CSR)
23     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
24     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
25
26
27 for av in tqdm(C_Para):
28     b = np.log10(av)
29     L_C_Para.append(b)
30
31 # Performance of model on Train data and Test data for each hyper parameter.
32 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
34 plt.gca()
35 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
37 plt.gca()
38 plt.legend()
39 plt.xlabel("Log alpha: hyperparameter")
40 plt.ylabel("AUC")
41 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
42 plt.grid()
43 plt.show()
```

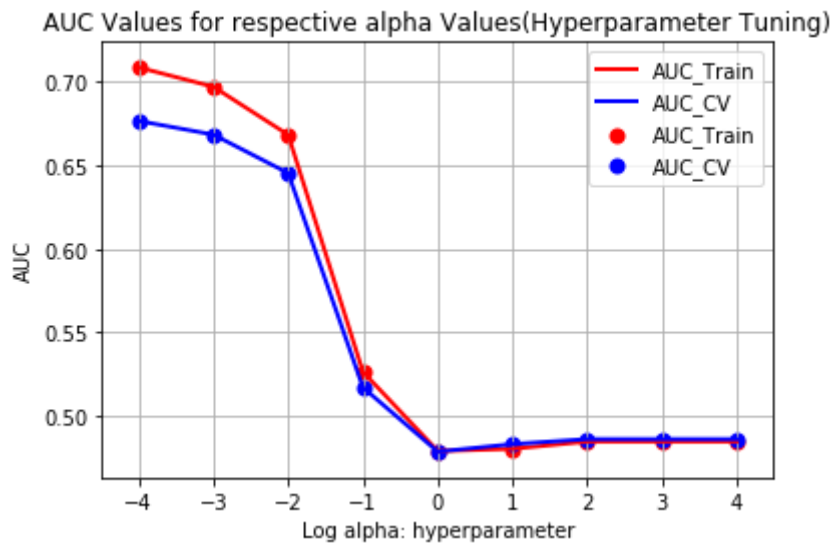


CPU times: user 34.9 s, sys: 684 ms, total: 35.6 s
Wall time: 35.4 s

Applying Support Vector Machines(SVM) on AVG W2V with L2 regularization, SET 3

In [0]:

```
1 %%time
2 AVG_W2V_TR_CSR = AVG_W2V_Train.tocsr()
3 AVG_W2V_CV_CSR = AVG_W2V_CV.tocsr()
4 AVG_W2V_TS_CSR = AVG_W2V_Test.tocsr()
5
6 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
7 L_C_Para= []
8 ACCV = []
9 AUC_TR = []
10 AUC_CV = []
11
12 for i in tqdm(C_Para):
13     SGC = SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1,alpha=i)
14     SGC.fit(AVG_W2V_Train, Y_Train)
15     pred = SGC.predict(AVG_W2V_CV)
16     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
17     ACCV.append(acc)
18     Train_pred = batch_predict(SGC, AVG_W2V_TR_CSR)
19     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
20     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
21
22     CV_pred = batch_predict(SGC, AVG_W2V_CV_CSR)
23     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
24     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
25
26
27 for av in tqdm(C_Para):
28     b = np.log10(av)
29     L_C_Para.append(b)
30
31 # Performance of model on Train data and Test data for each hyper parameter.
32 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
34 plt.gca()
35 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
37 plt.gca()
38 plt.legend()
39 plt.xlabel("Log alpha: hyperparameter")
40 plt.ylabel("AUC")
41 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
42 plt.grid()
43 plt.show()
```



CPU times: user 18.5 s, sys: 265 ms, total: 18.8 s
Wall time: 18.5 s

OBSERVATION:

1. α value has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha = 0.0001$ is having the highest AUC CV score.
3. Hence $\alpha = 0.0001$ is choosen as a best α value.

```
In [0]: 1 SGC_OPT=0.0001
2 AVG_W2V_opt =SGDClassifier(penalty='l2', alpha=SGC_OPT, class_weight='balanced', n_jobs=-1)
3 AVG_W2V_opt.fit(AVG_W2V_Train, Y_Train)
4 pred = AVG_W2V_opt.predict(AVG_W2V_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha Value = {0} is {1}%'.format(SGC_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, AVG_W2V_opt.decision_function(AVG_W2V_TR_CSR))
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, AVG_W2V_opt.decision_function(AVG_W2V_TS_CSR))
```

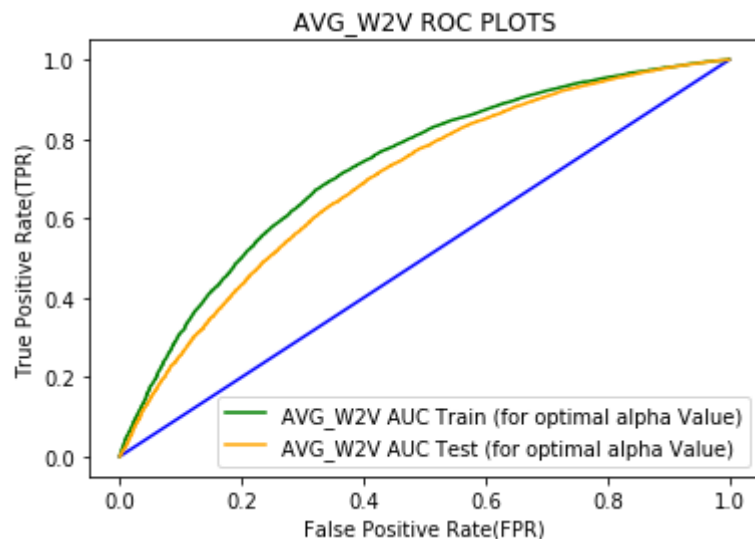
Test accuracy for alpha Value = 0.0001 is 73.63530456007989%

AVG_W2V ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="AVG_W2V AUC Train (for optimal alpha Value)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="AVG_W2V AUC Test (for optimal alpha Value)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("AVG_W2V ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal alpha Value) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal alpha Value) =", auc(a_fpr_Test, a_tpr_Test))
16 AVG_W2V_AOPT=SGC_OPT
17 AVG_W2V_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred5 = AVG_W2V_opt.predict(AVG_W2V_Train)
19 pred6 = AVG_W2V_opt.predict(AVG_W2V_Test)

```



AUC Train (for optimal alpha Value) = 0.7281203881050187
 AUC Test (for optimal alpha Value) = 0.6950187921378363
 CPU times: user 731 ms, sys: 2.77 ms, total: 734 ms
 Wall time: 731 ms

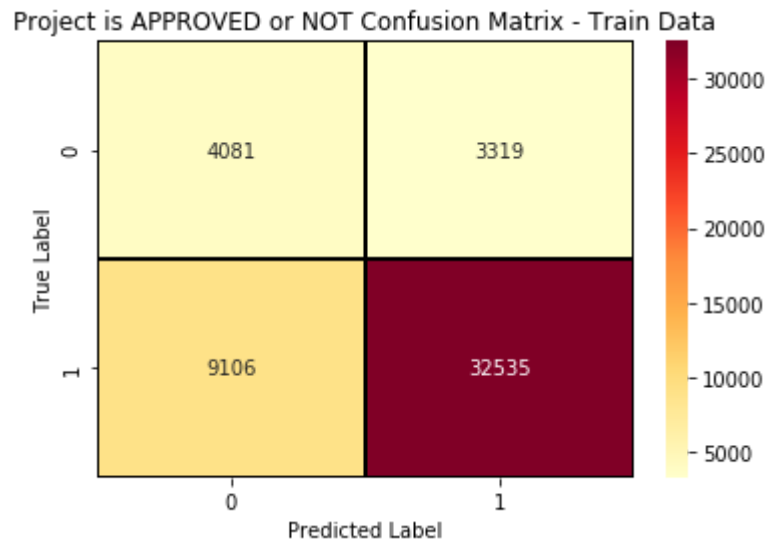
AVG_W2V CONFUSION MATRIX

In [0]:

```
1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred5)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
11
```

CPU times: user 119 ms, sys: 25.9 ms, total: 144 ms

Wall time: 98.2 ms



OBSERVATION:

True Negative =4081; False Negative =9106; True Positive = 32535; False Positive = 3319

Accuracy (Overall, how often is the classifier correct) = 0.75

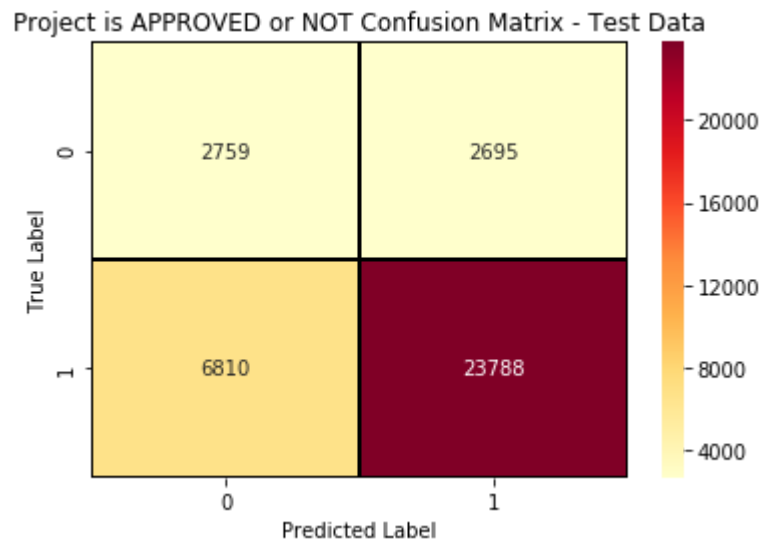
Precision(When it predicts yes, how often is it correct) =0.91

Misclassification (Overall, how often is it wrong) =0.26

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred6)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

CPU times: user 90 ms, sys: 40 ms, total: 130 ms

Wall time: 84.4 ms



OBSERVATION:

True Negative =2759; False Negative =6810; True Positive = 23788; False Positive = 2695

Accuracy (Overall, how often is the classifier correct) = 0.74

Precision(When it predicts yes, how often is it correct) =0.90

Misclassification (Overall, how often is it wrong) =0.27

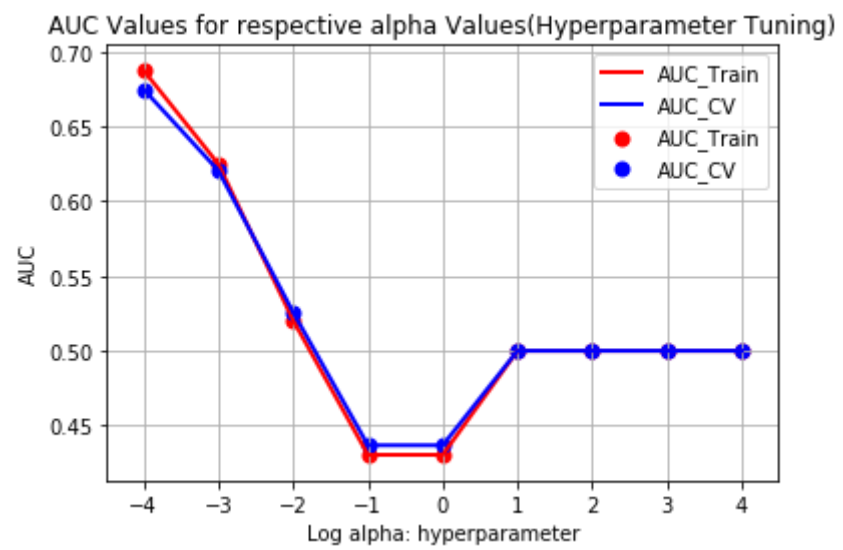
Applying Support Vector Machines(SVM) on TFIDF W2V with L1 regularization, SET 4

In [0]:

```
1 TFIDF_W2V_TR_CSR = TFIDF_W2V_Train.tocsr()
2 TFIDF_W2V_CV_CSR = TFIDF_W2V_CV.tocsr()
3 TFIDF_W2V_TS_CSR = TFIDF_W2V_Test.tocsr()
4
5 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
6 L_C_Para= []
7 ACCV = []
8 AUC_TR = []
9 AUC_CV = []
10
11 for i in tqdm(C_Para):
12     SGC = SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1,alpha=i)
13     SGC.fit(TFIDF_W2V_Train, Y_Train)
14     pred = SGC.predict(TFIDF_W2V_CV)
15     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
16     ACCV.append(acc)
17     Train_pred = batch_predict(SGC, TFIDF_W2V_TR_CSR)
18     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
19     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
20
21     CV_pred = batch_predict(SGC, TFIDF_W2V_CV_CSR)
22     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
23     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
24
25
26 for av in tqdm(C_Para):
27     b = np.log10(av)
28     L_C_Para.append(b)
29
30 # Performance of model on Train data and Test data for each hyper parameter.
31 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
32 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.gca()
34 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
35 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.gca()
37 plt.legend()
38 plt.xlabel("Log alpha: hyperparameter")
39 plt.ylabel("AUC")
40 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
41 plt.grid()
42 plt.show()
```

100%|██████████| 9/9 [00:17<00:00, 1.85s/it]

100%|██████████| 9/9 [00:00<00:00, 46374.37it/s]



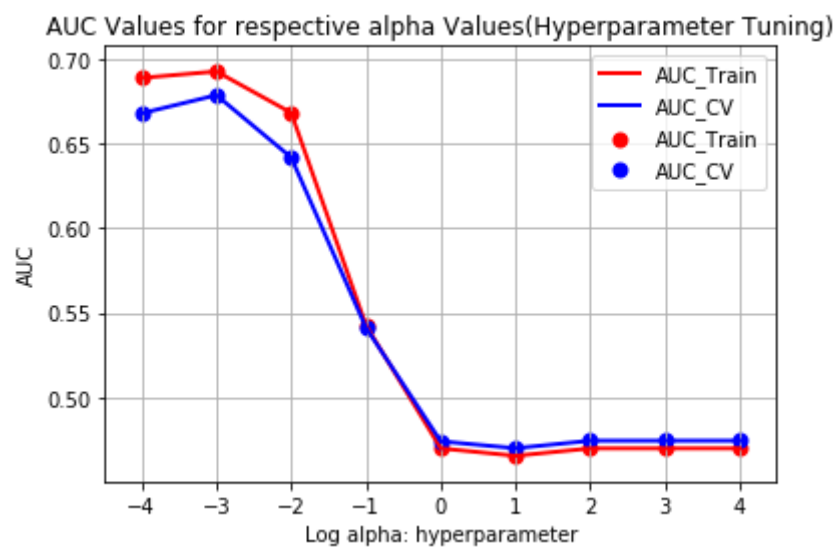
Applying Support Vector Machines(SVM) on TFIDF W2V with L2 regularization, **SET 4**

In [0]:

```
1 TFIDF_W2V_TR_CSR = TFIDF_W2V_Train.tocsr()
2 TFIDF_W2V_CV_CSR = TFIDF_W2V_CV.tocsr()
3 TFIDF_W2V_TS_CSR = TFIDF_W2V_Test.tocsr()
4
5 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
6 L_C_Para= []
7 ACCV = []
8 AUC_TR = []
9 AUC_CV = []
10
11 for i in tqdm(C_Para):
12     SGC = SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1,alpha=i)
13     SGC.fit(TFIDF_W2V_Train, Y_Train)
14     pred = SGC.predict(TFIDF_W2V_CV)
15     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
16     ACCV.append(acc)
17     Train_pred = batch_predict(SGC, TFIDF_W2V_TR_CSR)
18     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
19     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
20
21     CV_pred = batch_predict(SGC, TFIDF_W2V_CV_CSR)
22     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
23     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
24
25
26 for av in tqdm(C_Para):
27     b = np.log10(av)
28     L_C_Para.append(b)
29
30 # Performance of model on Train data and Test data for each hyper parameter.
31 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
32 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.gca()
34 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
35 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.gca()
37 plt.legend()
38 plt.xlabel("Log alpha: hyperparameter")
39 plt.ylabel("AUC")
40 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
41 plt.grid()
42 plt.show()
```

100%|██████████| 9/9 [00:09<00:00, 1.02it/s]

100%|██████████| 9/9 [00:00<00:00, 36972.32it/s]



OBSERVATION:

1. α value has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha = 0.001$ is having the highest AUC CV score.
3. Hence $\alpha = 0.001$ is choosen as a best α value.

```
In [0]: 1 SGC_OPT=0.001
2 TFIDF_W2V_opt =SGDClassifier(penalty='l2', alpha=SGC_OPT, class_weight='balanced', n_jobs=-1)
3 TFIDF_W2V_opt.fit(TFIDF_W2V_Train, Y_Train)
4 pred = TFIDF_W2V_opt.predict(TFIDF_W2V_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha Value = {0} is {1}%'.format(SGC_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, TFIDF_W2V_opt.decision_function(TFIDF_W2V_TR_CSR))
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, TFIDF_W2V_opt.decision_function(TFIDF_W2V_TS_CSR))
```

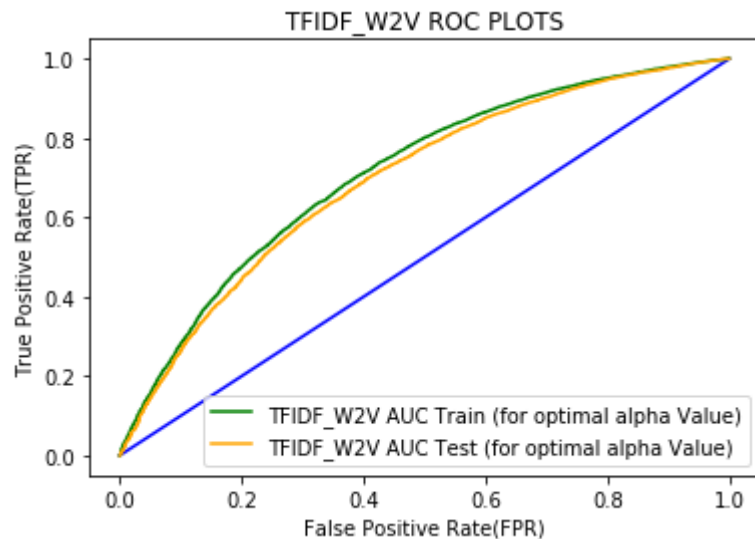
Test accuracy for alpha Value = 0.001 is 69.52734938422279%

TFIDF_W2V ROC PLOT

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="TFIDF_W2V AUC Train (for optimal alpha Value)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="TFIDF_W2V AUC Test (for optimal alpha Value)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF_W2V ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal alpha Value) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal alpha Value) =", auc(a_fpr_Test, a_tpr_Test))
16 TFIDF_W2V_AOPT=SGC_OPT
17 TFIDF_W2V_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred7 = TFIDF_W2V_opt.predict(TFIDF_W2V_Train)
19 pred8 = TFIDF_W2V_opt.predict(TFIDF_W2V_Test)

```



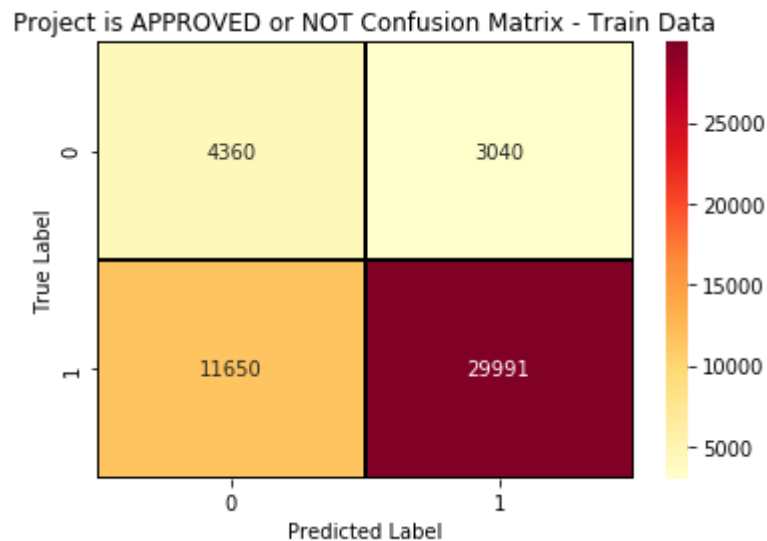
AUC Train (for optimal alpha Value) = 0.7113998596108175
AUC Test (for optimal alpha Value) = 0.6958268266201744

CPU times: user 519 ms, sys: 3.2 ms, total: 522 ms
Wall time: 521 ms

TFIDF_W2V CONFUSION MATRIX

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred7)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

CPU times: user 102 ms, sys: 41 ms, total: 143 ms
Wall time: 96.8 ms



OBSERVATION:

True Negative =4360; False Negative =11650; True Positive = 29991; False Positive = 3040

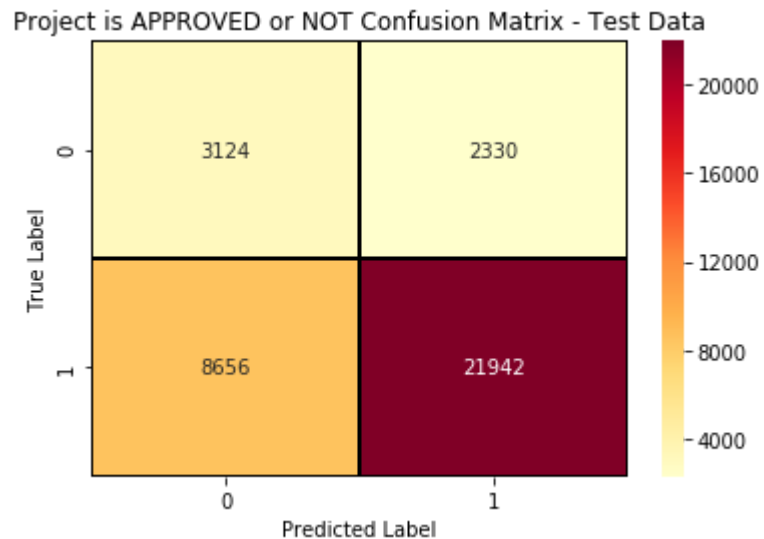
Accuracy (Overall, how often is the classifier correct) = 0.71

Precision(When it predicts yes, how often is it correct) =0.91

Misclassification (Overall, how often is it wrong) =0.30

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred8)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

CPU times: user 1 s, sys: 21.3 ms, total: 1.02 s
Wall time: 980 ms



OBSERVATION:

True Negative =3124; False Negative =8656; True Positive = 21942; False Positive = 2330

Accuracy (Overall, how often is the classifier correct) = 0.70

Precision(When it predicts yes, how often is it correct) =0.91

Misclassification (Overall, how often is it wrong) =0.31

- **[Task-2] Apply the Support Vector Machines on these features by finding the best hyper paramter as suggested in step 2 and step 3**
- Consider these set of features **Set 5** :
 - **school_state** : categorical data
 - **clean_categories** : categorical data
 - **clean_subcategories** : categorical data

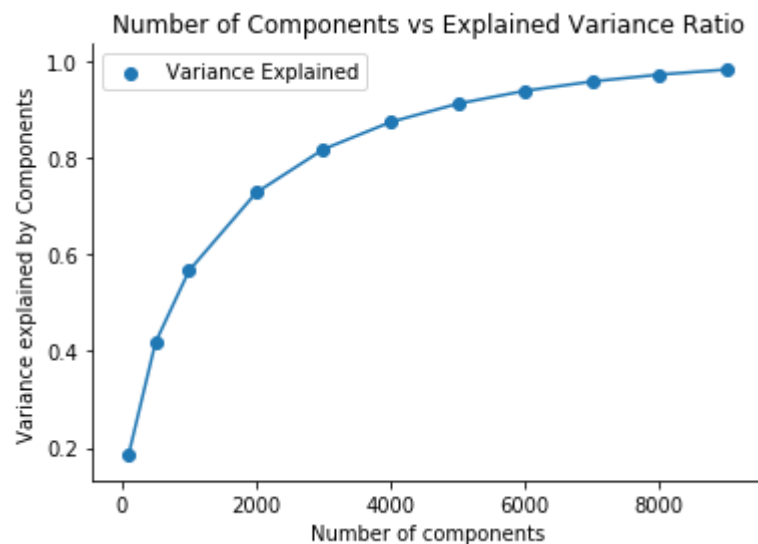
- **project_grade_category** :categorical data
- **teacher_prefix** : categorical data
- **quantity** : numerical data
- **teacher_number_of_previously_posted_projects** : numerical data
- **price** : numerical data
- **sentiment score's of each of the essay** : numerical data
- **number of words in the title** : numerical data
- **number of words in the combine essays** : numerical data
- Apply **TruncatedSVD** (<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>) on **TfidfVectorizer** (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html) of essay text, choose the number of components (**n_components**) using **elbow method** (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/>) : numerical data

Choosing the no of components using the Elbow Method

In [0]:

```
1 %%time
2 #https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html
3 #https://www.analyticsvidhya.com/blog/2016/03/practical-guide-principal-component-analysis-python/
4 #https://medium.com/@bioturing/how-to-read-pca-biplots-and-scree-plots-186246aae063
5 no_of_comp = [100,500,1000,2000,3000,4000,5000,6000,7000,8000,9000]
6 agg_variance = []
7 for i in tqdm(no_of_comp):
8     tsvd1 = TruncatedSVD(n_components = i,random_state = 0,n_iter = 3 )
9     tsvd1.fit(text_tfidf_Train)
10    agg_variance.append(tsvd1.explained_variance_ratio_.sum())
11
12 plt.plot(no_of_comp,agg_variance)
13 plt.scatter(no_of_comp,agg_variance,label = 'Variance Explained')
14 plt.xlabel("Number of components")
15 plt.ylabel("Variance explained by Components")
16 plt.legend()
17 plt.title("Number of Components vs Explained Variance Ratio")
18 sns.despine()
19 plt.show()
```

100%|██████████| 11/11 [45:44<00:00, 436.29s/it]

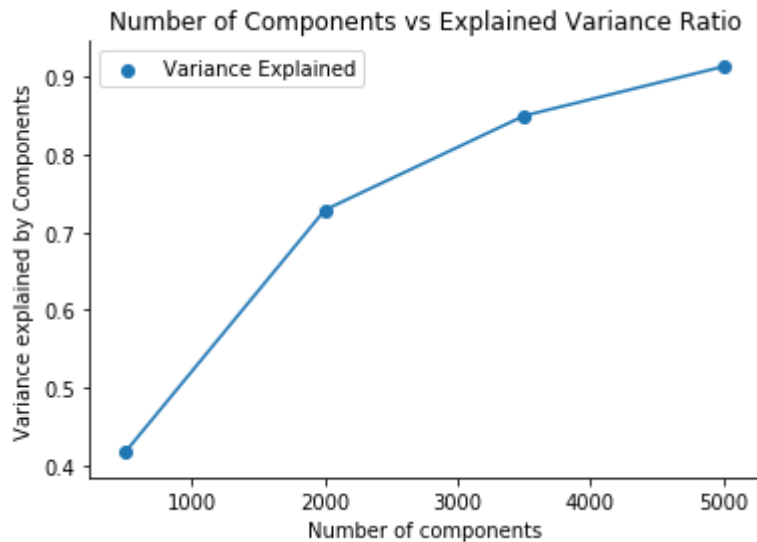


CPU times: user 7h 26min 9s, sys: 2h 6min 3s, total: 9h 32min 13s
Wall time: 45min 45s

Zoom In the Elbow Curve

```
In [0]: 1 %%time
2 no_of_comp = [500,2000,3500,5000]
3 agg_variance = []
4 for i in tqdm(no_of_comp):
5     tsvd1 = TruncatedSVD(n_components = i, random_state = 0, n_iter = 3 )
6     tsvd1.fit(text_tfidf_Train)
7     agg_variance.append(tsvd1.explained_variance_ratio_.sum())
8
9 plt.plot(no_of_comp, agg_variance)
10 plt.scatter(no_of_comp, agg_variance, label = 'Variance Explained')
11 plt.xlabel("Number of components")
12 plt.ylabel("Variance explained by Components")
13 plt.legend()
14 plt.title("Number of Components vs Explained Variance Ratio")
15 sns.despine()
16 plt.show()
```

100%|██████████| 4/4 [09:39<00:00, 143.96s/it]



CPU times: user 1h 2min 21s, sys: 34min 21s, total: 1h 36min 42s
Wall time: 9min 39s

```
In [0]: 1 from sklearn.decomposition import TruncatedSVD
2 tsvd2 = TruncatedSVD(n_components = 2000,n_iter = 5)
3 tsvd2.fit(text_tfidf_Train)
4 tsvd_Train = tsvd2.transform(text_tfidf_Train)
5 tsvd_CV = tsvd2.transform(text_tfidf_CV)
6 tsvd_Test = tsvd2.transform(text_tfidf_Test)
7 print(len(tsvd_Train))
8 print(len(tsvd_Test))
```

49041

36052

```
In [0]: 1 FiveF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one
2 print(FiveF_Train.shape)
3 FiveF_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,text_t
4 print(FiveF_CV.shape)
5 FiveF_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_
6 print(FiveF_Test.shape)
```

(49041, 16289)

(24155, 16289)

(36052, 16289)

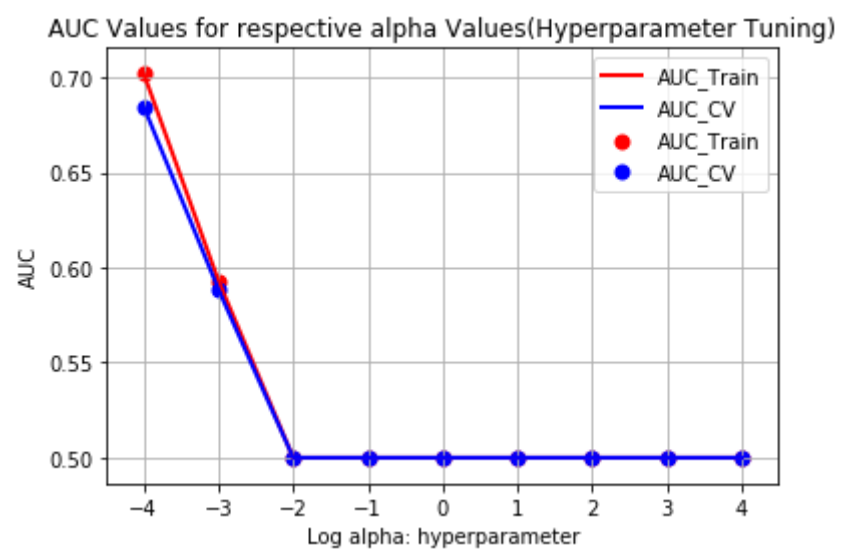
Applying Support Vector Machines(SVM) on TFIDF for Selected features with L1 regularization, SET 5

In [0]:

```
1 FiveF_TR_CSR = FiveF_Train.tocsr()
2 FiveF_CV_CSR = FiveF_CV.tocsr()
3 FiveF_TS_CSR = FiveF_Test.tocsr()
4
5 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
6 L_C_Para= []
7 ACCV = []
8 AUC_TR = []
9 AUC_CV = []
10
11 for i in tqdm(C_Para):
12     SGC = SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1,alpha=i)
13     SGC.fit(FiveF_Train, Y_Train)
14     pred = SGC.predict(FiveF_CV)
15     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
16     ACCV.append(acc)
17     Train_pred = batch_predict(SGC, FiveF_TR_CSR)
18     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
19     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
20
21     CV_pred = batch_predict(SGC, FiveF_CV_CSR)
22     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
23     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
24
25
26 for av in tqdm(C_Para):
27     b = np.log10(av)
28     L_C_Para.append(b)
29
30 # Performance of model on Train data and Test data for each hyper parameter.
31 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
32 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.gca()
34 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
35 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.gca()
37 plt.legend()
38 plt.xlabel("Log alpha: hyperparameter")
39 plt.ylabel("AUC")
40 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
41 plt.grid()
42 plt.show()
```

100%|██████████| 9/9 [02:02<00:00, 13.72s/it]

100%|██████████| 9/9 [00:00<00:00, 26123.69it/s]



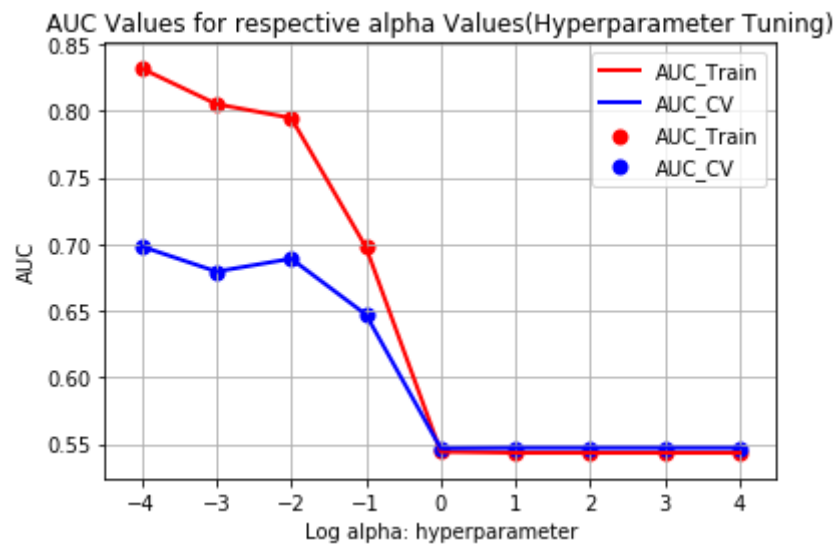
Applying Support Vector Machines(SVM) on TFIDF for Selected features with L2 regularization, **SET 5**

In [0]:

```
1 FiveF_TR_CSR = FiveF_Train.tocsr()
2 FiveF_CV_CSR = FiveF_CV.tocsr()
3 FiveF_TS_CSR = FiveF_Test.tocsr()
4
5 C_Para = [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
6 L_C_Para= []
7 ACCV = []
8 AUC_TR = []
9 AUC_CV = []
10
11 for i in tqdm(C_Para):
12     SGC = SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1,alpha=i)
13     SGC.fit(FiveF_Train, Y_Train)
14     pred = SGC.predict(FiveF_CV)
15     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
16     ACCV.append(acc)
17     Train_pred = batch_predict(SGC, FiveF_TR_CSR)
18     a_fpr_train,a_tpr_train,c = roc_curve(Y_Train, Train_pred)
19     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
20
21     CV_pred = batch_predict(SGC, FiveF_CV_CSR)
22     a_fpr_cv,a_tpr_cv,c = roc_curve(Y_CV, CV_pred)
23     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
24
25
26 for av in tqdm(C_Para):
27     b = np.log10(av)
28     L_C_Para.append(b)
29
30 # Performance of model on Train data and Test data for each hyper parameter.
31 plt.plot(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
32 plt.scatter(L_C_Para, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.gca()
34 plt.plot(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
35 plt.scatter(L_C_Para, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.gca()
37 plt.legend()
38 plt.xlabel("Log alpha: hyperparameter")
39 plt.ylabel("AUC")
40 plt.title("AUC Values for respective alpha Values(Hyperparameter Tuning)")
41 plt.grid()
42 plt.show()
```

100%|██████████| 9/9 [01:08<00:00, 7.46s/it]

100%|██████████| 9/9 [00:00<00:00, 17839.67it/s]



OBSERVATION:

1. α value has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha = 0.0001$ is having the highest AUC CV score.
3. Hence $\alpha = 0.0001$ is choosen as a best α value.

```
In [0]: 1 SGC_OPT=0.0001
2 FiveF_opt =SGDClassifier(penalty='l2', alpha=SGC_OPT, class_weight='balanced', n_jobs=-1)
3 FiveF_opt.fit(FiveF_Train, Y_Train)
4 pred = FiveF_opt.predict(FiveF_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha Value = {0} is {1}%'.format(SGC_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, FiveF_opt.decision_function(FiveF_TR_CSR))
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, FiveF_opt.decision_function(FiveF_TS_CSR))
```

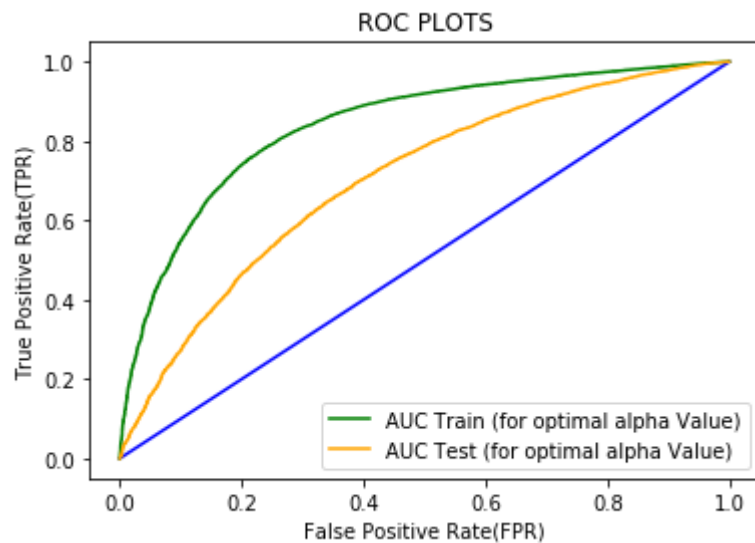
Test accuracy for alpha Value = 0.0001 is 73.23310773327417%

ROC PLOT FOR SELECTED FEATURES

```

In [0]: 1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="AUC Train (for optimal alpha Value)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="AUC Test (for optimal alpha Value)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal alpha Value) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal alpha Value) =", auc(a_fpr_Test, a_tpr_Test))
16 FiveF_AOPT=SGC_OPT
17 FiveF_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred9 = FiveF_opt.predict(FiveF_Train)
19 pred10 = FiveF_opt.predict(FiveF_Test)

```



AUC Train (for optimal alpha Value) = 0.839568955577046

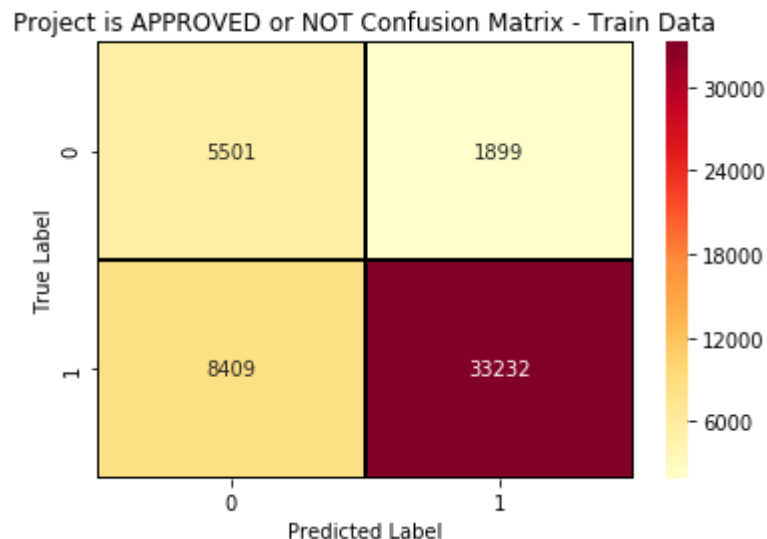
AUC Test (for optimal alpha Value) = 0.7041598897018491

CPU times: user 4.48 s, sys: 7.29 ms, total: 4.49 s
Wall time: 4.48 s

CONFUSION MATRIX FOR SELECTED FEATURES

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred9)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

CPU times: user 104 ms, sys: 45.9 ms, total: 150 ms
Wall time: 104 ms



OBSERVATION:

True Negative =5501; False Negative =8409; True Positive = 33232; False Positive = 1899

Accuracy (Overall, how often is the classifier correct) = 0.80

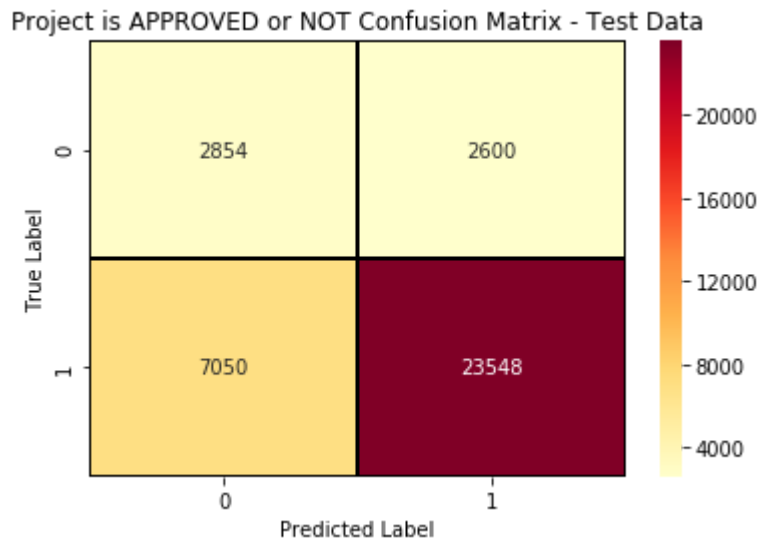
Precision(When it predicts yes, how often is it correct) =0.95

Misclassification (Overall, how often is it wrong) =0.22

```
In [0]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred10)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

CPU times: user 92.8 ms, sys: 34.2 ms, total: 127 ms

Wall time: 82.9 ms



OBSERVATION:

True Negative =2854; False Negative =7050; True Positive = 23548; False Positive = 2600

Accuracy (Overall, how often is the classifier correct) = 0.74

Precision(When it predicts yes, how often is it correct) =0.91

Misclassification (Overall, how often is it wrong) =0.27

3. Conclusions

```
In [0]: 1 # Please compare all your models using Prettytable Library
2
3 pt = PrettyTable()
4 pt.field_names= ("S.No", "Vectorizer", "Model", "HyperParameter", "AUC")
5 pt.add_row(["1", "BOW", "SVM", BOW_AOPT, BOW_AUC])
6 pt.add_row(["2", "TFIDF", "SVM", TFIDF_AOPT, TFIDF_AUC])
7 pt.add_row(["3", "AVG_W2V", "SVM", AVG_W2V_AOPT, AVG_W2V_AUC])
8 pt.add_row(["4", "TFIDF_W2V", "SVM", TFIDF_W2V_AOPT, TFIDF_W2V_AUC])
9 pt.add_row(["5", "SET 5", "SVM", FiveF_AOPT, FiveF_AUC])
10 print(pt)
```

```
+-----+-----+-----+-----+-----+
| S.No | Vectorizer | Model | HyperParameter | AUC |
+-----+-----+-----+-----+-----+
| 1    | BOW       | SVM   | 0.01           | 72.0 |
| 2    | TFIDF     | SVM   | 0.0001         | 71.0 |
| 3    | AVG_W2V   | SVM   | 0.0001         | 70.0 |
| 4    | TFIDF_W2V | SVM   | 0.001          | 70.0 |
| 5    | SET 5     | SVM   | 0.0001         | 70.0 |
+-----+-----+-----+-----+-----+
```

SUMMARY:

1. **"BOW and TFIDF"** performing very well in accuracy and precision scores. When compared to other vectorizers such as AVG_W2V and TFIDF_W2V.
2. Hyperparameter such as regularization and the number of iterations needs to be chosen very carefully.
3. Linear SVM is fast when compared to the Kernel SVM.
4. For graphical representation applied Log function on alpha values.