# 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 |
|---|---|
| `project_id` | A unique identifier for the proposed project. **Example:** `p036502` |
| `project_title` | Title of the project. **Examples:** <br> - `Art Will Make You Happy!` <br> - `First Grade Fun` |
| `project_grade_category` | Grade level of students for which the project is targeted. One of the following enumerated values: <br> - `Grades PreK-2` <br> - `Grades 3-5` <br> - `Grades 6-8` <br> - `Grades 9-12` |
| `project_subject_categories` | One or more (comma-separated) subject categories for the project from the following enumerated list of values: <br> - `Applied Learning` <br> - `Care & Hunger` <br> - `Health & Sports` <br> - `History & Civics` <br> - `Literacy & Language` <br> - `Math & Science` <br> - `Music & The Arts` <br> - `Special Needs` <br> - `Warmth` <br><br> **Examples:** <br> - `Music & The Arts` <br> - `Literacy & Language, Math & Science` |
| `school_state` | State where school is located ([Two-letter U.S. postal code](#)). **Example:** `WY` |
| `project_subject_subcategories` | One or more (comma-separated) subject subcategories for the project. **Examples:** <br> - `Literacy` <br> - `Literature & Writing, Social Sciences` |
| `project_resource_summary` | An explanation of the resources needed for the project. **Example:** <br> - `My students need hands on literacy materials to manage sensory needs!` |
| `project_essay_1` | First application essay[*] |
| `project_essay_2` | Second application essay[*] |
| `project_essay_3` | Third application essay[*] |

| Feature | Description |
|---|---|
| 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:<br><br>• `nan`<br>• `Dr.`<br>• `Mr.`<br>• `Mrs.`<br>• `Ms.`<br>• `Teacher.` |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. **Example:** `2` |

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

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

| Feature | Description |
|---|---|
| id | A `project_id` value from the `train.csv` file. **Example:** `p036502` |
| description | Desciption of the resource. **Example:** `Tenor Saxophone Reeds, Box of 25` |
| quantity | Quantity of the resource required. **Example:** `3` |
| price | Price of the resource required. **Example:** `9.95` |

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

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

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

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:
- __project_essay_1:__ "Introduce us to your classroom"
- __project_essay_2:__ "Tell us more about your students"
- __project_essay_3:__ "Describe how your students will use the materials you're requesting"
- __project_essay_3:__ "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:
- __project_essay_1:__ "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2:__ "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [2]:

```python
from google.colab import drive
drive.mount('/content/gdrive')
```

```
Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6
qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%
b&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2
www.googleapis.com%2fauth%2fdrive%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly
ttps%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly

Enter your authorization code:
```

```
..........
Mounted at /content/gdrive
```

In [3]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer
import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle
from tqdm import tqdm
import os
import chart_studio.plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
from scipy.sparse import hstack,vstack
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn import model_selection
from sklearn.metrics import roc_auc_score
from sklearn.model_selection import GridSearchCV
from prettytable import PrettyTable
from sklearn.preprocessing import Normalizer
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer
nltk.download('vader_lexicon')
import pdb
from sklearn.linear_model import SGDClassifier
from sklearn.decomposition import TruncatedSVD
```

Output hidden; open in https://colab.research.google.com to view.

## 1.1 Reading Data

In [4]:

```
Project_data = pd.read_csv('/content/gdrive/My Drive/Colab Notebooks/train_data.csv')
Resource_data = pd.read_csv('/content/gdrive/My Drive/Colab Notebooks/resources.csv')
print(Project_data.shape)
print(Resource_data.shape)
```

```
(109248, 17)
(1541272, 4)
```

In [5]:

```
# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(Project_data.columns)]
#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
```

```
Project_data['Date'] = pd.to_datetime(Project_data['project_submitted_datetime'])
Project_data.drop('project_submitted_datetime', axis=1, inplace=True)
Project_data.sort_values(by=['Date'], inplace=True)
# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
Project_data = Project_data[cols]
Project_data.head(2)
```

Out[5]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | Date | project_grade_category | project_s |
|---|---|---|---|---|---|---|---|---|
| **55660** | 8393 | p205479 | 2bf07ba08945e5d8b2a3f269b2b3cfe5 | Mrs. | CA | 2016-04-27 00:27:36 | Grades PreK-2 | |
| **76127** | 37728 | p043609 | 3f60494c61921b3b43ab61bdde2904df | Ms. | UT | 2016-04-27 00:31:25 | Grades 3-5 | |

# 2. Support Vector Machines

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

In [6]:

```
y = Project_data['project_is_approved'].values
Project_data.drop(['project_is_approved'], axis=1, inplace=True)
n_z = len(Project_data)
y_z = np.zeros(n_z, dtype=np.int32)

X = Project_data
# train test split
X_train, X_Test, y_train, y_Test = train_test_split(X, y, test_size=0.33, random_state=0, stratify=
y_z)
print('Shape of X_train: ',X_train.shape)
print('Shape of y_train: ',y_train.shape)
print('Shape of X_Test: ',X_Test.shape)
print('Shape of y_Test: ',y_Test.shape)
```

```
Shape of X_train:  (73196, 16)
Shape of y_train:  (73196,)
Shape of X_Test:  (36052, 16)
Shape of y_Test:  (36052,)
```

## 1.2 preprocessing of `project_subject_categories`

In [7]:

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
catogories_train = list(X_train['project_subject_categories'].values)
cat_list = []
for i in catogories_train:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
```

```
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

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

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

print(len(sorted_cat_dict_train))


catogories_Test = list(X_Test['project_subject_categories'].values)
cat_list = []
for i in catogories_Test:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

9

## 1.3 preprocessing of `project_subject_subcategories`

In [8]:

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
sub_catogories_train = list(X_train['project_subject_subcategories'].values)
sub_cat_list = []
for i in sub_catogories_train:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
        temp +=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

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

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

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

print(len(sorted_sub_cat_dict_train))

sub_catogories_Test = list(X_Test['project_subject_subcategories'].values)
sub_cat_list = []
for i in sub_catogories_Test:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & H
unger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Scienc
e"=> "Math","&", "Science"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i
.e removing 'The')
        j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
Science"=>"Math&Science"
        temp +=j.strip()+" "#" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

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

30

## 1.3 Text preprocessing

In [0]:

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


X_Test["essay"] = X_Test["project_essay_1"].map(str) +\
                        X_Test["project_essay_2"].map(str) + \
                        X_Test["project_essay_3"].map(str) + \
                        X_Test["project_essay_4"].map(str)
```

In [0]:

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

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

In [0]:

```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've"
,\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
```

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

In [12]:

```python
# Combining all the above stundents
# tqdm is for printing the status bar

preprocessed_essays_train = []

preprocessed_essays_Test = []

for sentance in tqdm(X_train['essay'].values):
    sent = decontracted(sentance)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays_train.append(sent.lower().strip())


for sentance in tqdm(X_Test['essay'].values):
    sent = decontracted(sentance)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays_Test.append(sent.lower().strip())

print("Shape of preprocessed_essays_train after preprocessing",len(preprocessed_essays_train))

print("Shape of preprocessed_essays_Test after preprocessing",len(preprocessed_essays_Test))
# pdb.set_trace()
```

```
100%|██████████| 73196/73196 [00:40<00:00, 1824.41it/s]
100%|██████████| 36052/36052 [00:19<00:00, 1814.90it/s]
```

```
Shape of preprocessed_essays_train after preprocessing 73196
Shape of preprocessed_essays_Test after preprocessing 36052
```

In [13]:

```python
word_count_essay_train = []
for a in tqdm(X_train["essay"]) :
    b = len(a.split())
    word_count_essay_train.append(b)
```

```
X_train["word_count_essay_train"] = word_count_essay_train


word_count_essay_Test = []
for a in tqdm(X_Test["essay"]) :
    b = len(a.split())
    word_count_essay_Test.append(b)

X_Test["word_count_essay_Test"] = word_count_essay_Test
```

```
100%|████████| 73196/73196 [00:01<00:00, 64405.29it/s]
100%|████████| 36052/36052 [00:00<00:00, 65019.05it/s]
```

## 1.4 Preprocessing of `project_title`

In [14]:

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



preprocessed_titles_Test = []
for sentence in tqdm(X_Test['project_title'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_titles_Test.append(sent.lower().strip())

print("Shape of preprocessed_titles_train after preprocessing",len(preprocessed_titles_train))

print("Shape of preprocessed_titles_Test after preprocessing",len(preprocessed_titles_Test))
```

```
100%|████████| 73196/73196 [00:01<00:00, 42666.74it/s]
100%|████████| 36052/36052 [00:00<00:00, 42310.39it/s]
```

```
Shape of preprocessed_titles_train after preprocessing 73196
Shape of preprocessed_titles_Test after preprocessing 36052
```

In [15]:

```
word_count_title_train = []
for a in tqdm(X_train["project_title"]) :
    b = len(a.split())
    word_count_title_train.append(b)

X_train["word_count_title_train"] = word_count_title_train


word_count_title_Test = []
for a in tqdm(X_Test["project_title"]) :
    b = len(a.split())
    word_count_title_Test.append(b)

X_Test["word_count_title_Test"] = word_count_title_Test
```

# Make Data Model Ready: encoding numerical, categorical features

## 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/

In [16]:

```python
# we use count vectorizer to convert the values into one
vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict_train.keys()), lowercase=False, binary=True)
vectorizer_cat.fit(X_train['clean_categories'].values)
categories_one_hot_train = vectorizer_cat.transform(X_train['clean_categories'].values)

categories_one_hot_Test = vectorizer_cat.transform(X_Test['clean_categories'].values)
print(vectorizer_cat.get_feature_names())
print("Shape of categories_one_hot_train matrix after one hot encodig ",categories_one_hot_train.shape)

print("Shape of categories_one_hot_Test matrix after one hot encodig ",categories_one_hot_Test.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of categories_one_hot_train matrix after one hot encodig  (73196, 9)
Shape of categories_one_hot_Test matrix after one hot encodig  (36052, 9)
```

In [17]:

```python
# we use count vectorizer to convert the values into one
vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict_train.keys()), lowercase=False, binary=True)
sub_categories_one_hot_train = vectorizer_sub_cat.fit_transform(X_train['clean_subcategories'].values)

sub_categories_one_hot_Test = vectorizer_sub_cat.transform(X_Test['clean_subcategories'].values)
print(vectorizer_sub_cat.get_feature_names())
print("Shape of sub_categories_one_hot_train matrix after one hot encodig ",sub_categories_one_hot_train.shape)

print("Shape of sub_categories_one_hot_Test matrix after one hot encodig ",sub_categories_one_hot_Test.shape)
```

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Civics_Government', '
Extracurricular', 'ForeignLanguages', 'Warmth', 'Care_Hunger', 'NutritionEducation',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL
', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of sub_categories_one_hot_train matrix after one hot encodig  (73196, 30)
Shape of sub_categories_one_hot_Test matrix after one hot encodig  (36052, 30)
```

**School State**

In [18]:

```python
sch1_catogories = list(X_train['school_state'].values)
school_list = []
for sent in sch1_catogories:
    school_list.append(sent.lower().strip())
X_train['school_categories'] = school_list
```

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

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
sch_dict = dict(my_counter_sch)
sorted_sch_dict = dict(sorted(sch_dict.items(), key=lambda kv: kv[1]))

vectorizer_sch = CountVectorizer(vocabulary=list(sorted_sch_dict.keys()), lowercase=False, binary=
True)
vectorizer_sch.fit(X_train['school_categories'].values)
#print(vectorizer.get_feature_names())

sch_one_hot_train = vectorizer_sch.transform(X_train['school_categories'].values)
print("Shape of sch_one_hot_train matrix after one hot encodig ",sch_one_hot_train.shape)
#----------------------------------------------------------------------------------------


sch1_catogories_Test = list(X_Test['school_state'].values)
school_list_Test = []
for sent in sch1_catogories_Test:
    school_list_Test.append(sent.lower().strip())
X_Test['school_categories'] = school_list_Test
X_Test.drop(['school_state'], axis=1, inplace=True)

sch_one_hot_Test = vectorizer_sch.transform(X_Test['school_categories'].values)

print("Shape of sch_one_hot_Test matrix after one hot encodig ",sch_one_hot_Test.shape)
```

```
Shape of sch_one_hot_train matrix after one hot encodig  (73196, 51)
Shape of sch_one_hot_Test matrix after one hot encodig  (36052, 51)
```

**Prefix**

In [19]:

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
prefix_catogories_train = list(X_train['teacher_prefix'].values)
prefix_list_train = []
for sent in prefix_catogories_train:
    sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split())
    prefix_list_train.append(sent.lower().strip())
X_train['prefix_catogories'] = prefix_list_train
X_train.drop(['teacher_prefix'], axis=1, inplace=True)

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
prefix_dict_train = dict(my_counter_prefix_train)
sorted_prefix_dict_train = dict(sorted(prefix_dict_train.items(), key=lambda kv: kv[1]))

vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict_train.keys()), lowercase=Fa
lse, binary=True)
vectorizer_prefix.fit(X_train['prefix_catogories'].values)
#print(vectorizer.get_feature_names())

prefix_one_hot_train = vectorizer_prefix.transform(X_train['prefix_catogories'].values)
print("Shape of prefix_one_hot_train matrix after one hot encodig ",prefix_one_hot_train.shape)


#----------------------------------------------------------------------------------------

prefix_catogories_Test = list(X_Test['teacher_prefix'].values)
```

```
prefix_list_Test = []
for sent in prefix_catogories_Test:
    sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split())
    prefix_list_Test.append(sent.lower().strip())
X_Test['prefix_catogories'] = prefix_list_Test
X_Test.drop(['teacher_prefix'], axis=1, inplace=True)

prefix_one_hot_Test = vectorizer_prefix.transform(X_Test['prefix_catogories'])

print("Shape of prefix_one_hot_Test matrix after one hot encodig ",prefix_one_hot_Test.shape)
```

```
Shape of prefix_one_hot_train matrix after one hot encodig  (73196, 6)
Shape of prefix_one_hot_Test matrix after one hot encodig   (36052, 6)
```

**project_grade_category**

In [20]:

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
grade_catogories_train = list(X_train['project_grade_category'].values)
grade_list_train = []
for sent in grade_catogories_train:
    sent = sent.replace('-','_')
    sent = sent.replace(' ','_')
    # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split())
    grade_list_train.append(sent.lower().strip())

# temp = temp.replace('-','_')
X_train['new_grade_category'] = grade_list_train
X_train.drop(['project_grade_category'], axis=1, inplace=True)

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

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
grade_dict_train = dict(my_counter_grade_train)
sorted_grade_dict_train = dict(sorted(grade_dict_train.items(), key=lambda kv: kv[1]))


vectorizer_grade = CountVectorizer(vocabulary=list(sorted_grade_dict_train.keys()), lowercase=Fals
e, binary=True)
vectorizer_grade.fit(X_train['new_grade_category'].values)
#print(vectorizer.get_feature_names())

grade_one_hot_train = vectorizer_grade.transform(X_train['new_grade_category'].values)
print("Shape of grade_one_hot_train matrix after one hot encodig ",grade_one_hot_train.shape)


#-------------------------------------------------------------------------------------

grade_catogories_Test = list(X_Test['project_grade_category'].values)
grade_list_Test = []
for sent in grade_catogories_Test:
    sent = sent.replace('-','_')
    sent = sent.replace(' ','_')
    # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split())
    grade_list_Test.append(sent.lower().strip())

# temp = temp.replace('-','_')
X_Test['new_grade_category'] = grade_list_Test
X_Test.drop(['project_grade_category'], axis=1, inplace=True)

grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)
```

```
grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)

print("Shape of grade_one_hot_Test matrix after one hot encodig ",grade_one_hot_Test.shape)
```

```
Shape of grade_one_hot_train matrix after one hot encodig  (73196, 4)
Shape of grade_one_hot_Test matrix after one hot encodig  (36052, 4)
```

## 2.2 Make Data Model Ready: encoding numerical, categorical features

### 1.5.2 Vectorizing Numerical features

**Price and Quantity data**

In [0]:

```
price_data = Resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
X_train = pd.merge(X_train, price_data, on='id', how='left')

X_Test = pd.merge(X_Test, price_data, on='id', how='left')
```

In [22]:

```
price_norm = Normalizer(norm='l2', copy=False)
price_norm.fit(X_train['price'].values.reshape(1,-1))

price_norm.transform(X_train['price'].values.reshape(1,-1))

price_norm.transform(X_Test['price'].values.reshape(1,-1))

price_norm_train = (X_train['price'].values.reshape(-1,1))

price_norm_Test = (X_Test['price'].values.reshape(-1,1))

print("Shape of price_norm_train matrix after one hot encodig ",price_norm_train.shape)

print("Shape of price_norm_Test matrix after one hot encodig ",price_norm_Test.shape)
```

```
Shape of price_norm_train matrix after one hot encodig  (73196, 1)
Shape of price_norm_Test matrix after one hot encodig  (36052, 1)
```

In [23]:

```
quantity_norm = Normalizer(norm='l2', copy=False)
quantity_norm.fit(X_train['quantity'].values.reshape(1,-1))

quantity_norm_train = quantity_norm.transform(X_train['quantity'].values.reshape(1,-1))

quantity_norm_Test = quantity_norm.transform(X_Test['quantity'].values.reshape(1,-1))

quantity_norm_train = (X_train['quantity'].values.reshape(-1,1))

quantity_norm_Test = (X_Test['quantity'].values.reshape(-1,1))

print("Shape of quantity_norm_train matrix after one hot encodig ",quantity_norm_train.shape)

print("Shape of quantity_norm_Test matrix after one hot encodig ",quantity_norm_Test.shape)
```

```
Shape of quantity_norm_train matrix after one hot encodig  (73196, 1)
Shape of quantity_norm_Test matrix after one hot encodig  (36052, 1)
```

**teacher_number_of_previously_posted_projects**

In [24]:

```
teacher_prev_post_norm = Normalizer(norm='l2', copy=False)
teacher_prev_post_norm.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(
1,-1))
```

```
teacher_prev_post_norm_train =
teacher_prev_post_norm.transform(X_train['teacher_number_of_previously_posted_projects'].values.re
shape(1,-1))

teacher_prev_post_norm_Test =
teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.res
hape(1,-1))

teacher_prev_post_norm_train =
(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

teacher_prev_post_norm_Test =
(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
print("Shape of teacher_prev_post_norm_train matrix after one hot encodig
",teacher_prev_post_norm_train.shape)

print("Shape of teacher_prev_post_norm_Test matrix after one hot encodig
",teacher_prev_post_norm_Test.shape)
```

```
Shape of teacher_prev_post_norm_train matrix after one hot encodig  (73196, 1)
Shape of teacher_prev_post_norm_Test matrix after one hot encodig  (36052, 1)
```

**Title word count**

In [25]:

```
title_norm = Normalizer(norm='l2', copy=False)
title_norm.fit(X_train['word_count_title_train'].values.reshape(1,-1))
word_count_title_train = title_norm.transform(X_train['word_count_title_train'].values.reshape(1,-1
))

word_count_title_Test = title_norm.transform(X_Test['word_count_title_Test'].values.reshape(1,-1))

word_count_title_train = (X_train['word_count_title_train'].values.reshape(-1,1))

word_count_title_Test = (X_Test['word_count_title_Test'].values.reshape(-1,1))

print(word_count_title_train.shape)

print(word_count_title_Test.shape)
```

```
(73196, 1)
(36052, 1)
```

**Essay word count**

In [26]:

```
essay_norm = Normalizer(norm='l2', copy=False)
essay_norm.fit(X_train['word_count_essay_train'].values.reshape(1,-1))
word_count_essay_train = essay_norm.transform(X_train['word_count_essay_train'].values.reshape(1,-1
))

word_count_essay_Test = essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(1,-1))

word_count_essay_train = (X_train['word_count_essay_train'].values.reshape(-1,1))

word_count_essay_Test = (X_Test['word_count_essay_Test'].values.reshape(-1,1))

print(word_count_essay_train.shape)

print(word_count_essay_Test.shape)
```

```
(73196, 1)
(36052, 1)
```

**Sentiment Scores**

In [27]:

```python
# https://www.geeksforgeeks.org/python-sentiment-analysis-using-vader/
sid = SentimentIntensityAnalyzer()
essays = X_train['essay']

sentiment_pos_essay_Train = []
sentiment_neg_essay_Train = []
sentiment_neut_essay_Train = []
sentiment_com_essay_Train = []

for essay in tqdm(essays):
    res = sid.polarity_scores(essay)
    sentiment_pos_essay_Train.append(res['pos'])
    sentiment_neg_essay_Train.append(res['neg'])
    sentiment_neut_essay_Train.append(res['neu'])
    sentiment_com_essay_Train.append(res['compound'])
X_train['sentiment_pos_essay_Train'] = sentiment_pos_essay_Train
X_train['sentiment_neg_essay_Train'] = sentiment_neg_essay_Train
X_train['sentiment_neut_essay_Train'] = sentiment_neut_essay_Train
X_train['sentiment_com_essay_Train'] = sentiment_com_essay_Train

essays = X_Test['essay']

sentiment_pos_essay_Test = []
sentiment_neg_essay_Test = []
sentiment_neut_essay_Test = []
sentiment_com_essay_Test = []

for essay in tqdm(essays):
    res = sid.polarity_scores(essay)
    sentiment_pos_essay_Test.append(res['pos'])
    sentiment_neg_essay_Test.append(res['neg'])
    sentiment_neut_essay_Test.append(res['neu'])
    sentiment_com_essay_Test.append(res['compound'])
X_Test['sentiment_pos_essay_Test'] = sentiment_pos_essay_Test
X_Test['sentiment_neg_essay_Test'] = sentiment_neg_essay_Test
X_Test['sentiment_neut_essay_Test'] = sentiment_neut_essay_Test
X_Test['sentiment_com_essay_Test'] = sentiment_com_essay_Test

sentiment_norm_pos = Normalizer(norm='l2', copy=False)
sentiment_norm_neg = Normalizer(norm='l2', copy=False)
sentiment_norm_neut = Normalizer(norm='l2', copy=False)
sentiment_norm_com = Normalizer(norm='l2', copy=False)

sentiment_norm_pos.fit(X_train['sentiment_pos_essay_Train'].values.reshape(1,-1))
sentiment_norm_neg.fit(X_train['sentiment_neg_essay_Train'].values.reshape(1,-1))
sentiment_norm_neut.fit(X_train['sentiment_neut_essay_Train'].values.reshape(1,-1))
sentiment_norm_com.fit(X_train['sentiment_com_essay_Train'].values.reshape(1,-1))

senti_pos_ess_Tr_norm = sentiment_norm_pos.transform(X_train['sentiment_pos_essay_Train'].values.reshape(1,-1))
senti_pos_ess_Tr_norm = (X_train['sentiment_pos_essay_Train'].values.reshape(-1,1))

senti_neg_ess_Tr_norm = sentiment_norm_neg.transform(X_train['sentiment_neg_essay_Train'].values.reshape(1,-1))
senti_neg_ess_Tr_norm = (X_train['sentiment_neg_essay_Train'].values.reshape(-1,1))

senti_neut_ess_Tr_norm =
sentiment_norm_neut.transform(X_train['sentiment_neut_essay_Train'].values.reshape(1,-1))
senti_neut_ess_Tr_norm = (X_train['sentiment_neut_essay_Train'].values.reshape(-1,1))

senti_com_ess_Tr_norm = sentiment_norm_com.transform(X_train['sentiment_com_essay_Train'].values.reshape(1,-1))
senti_com_ess_Tr_norm = (X_train['sentiment_com_essay_Train'].values.reshape(-1,1))


senti_pos_ess_Ts_norm =
sentiment_norm_pos.transform(X_Test['sentiment_pos_essay_Test'].values.reshape(1,-1))
senti_pos_ess_Ts_norm = (X_Test['sentiment_pos_essay_Test'].values.reshape(-1,1))

senti_neg_ess_Ts_norm =
sentiment_norm_neg.transform(X_Test['sentiment_neg_essay_Test'].values.reshape(1,-1))
senti_neg_ess_Ts_norm = (X_Test['sentiment_neg_essay_Test'].values.reshape(-1,1))

senti_neut_ess_Ts_norm = sentiment_norm_neut.transform(X_Test['sentiment_neut_essay_Test'].values.reshape(1,-1))
senti_neut_ess_Ts_norm = (X_Test['sentiment_neut_essay_Test'].values.reshape(-1,1))
```

```
senti_com_ess_Ts_norm =
sentiment_norm_com.transform(X_Test['sentiment_com_essay_Test'].values.reshape(1,-1))
senti_com_ess_Ts_norm = (X_Test['sentiment_com_essay_Test'].values.reshape(-1,1))


print("Shape of senti_pos_ess_Tr_norm matrix after one hot encodig ",senti_pos_ess_Tr_norm.shape)
print("Shape of senti_neg_ess_Tr_norm matrix after one hot encodig ",senti_neg_ess_Tr_norm.shape)
print("Shape of senti_neut_ess_Tr_norm matrix after one hot encodig ",senti_neut_ess_Tr_norm.shape
)
print("Shape of senti_com_ess_Tr_norm matrix after one hot encodig ",senti_com_ess_Tr_norm.shape)
print("Shape of senti_pos_ess_Ts_norm matrix after one hot encodig ",senti_pos_ess_Ts_norm.shape)
print("Shape of senti_neg_ess_Ts_norm matrix after one hot encodig ",senti_neg_ess_Ts_norm.shape)
print("Shape of senti_neut_ess_Ts_norm matrix after one hot encodig ",senti_neut_ess_Ts_norm.shape
)
print("Shape of senti_com_ess_Ts_norm matrix after one hot encodig ",senti_com_ess_Ts_norm.shape)
```

```
100%|███████████| 73196/73196 [03:20<00:00, 364.58it/s]
100%|███████████| 36052/36052 [01:39<00:00, 361.19it/s]
```

```
Shape of senti_pos_ess_Tr_norm matrix after one hot encodig  (73196, 1)
Shape of senti_neg_ess_Tr_norm matrix after one hot encodig  (73196, 1)
Shape of senti_neut_ess_Tr_norm matrix after one hot encodig  (73196, 1)
Shape of senti_com_ess_Tr_norm matrix after one hot encodig  (73196, 1)
Shape of senti_pos_ess_Ts_norm matrix after one hot encodig  (36052, 1)
Shape of senti_neg_ess_Ts_norm matrix after one hot encodig  (36052, 1)
Shape of senti_neut_ess_Ts_norm matrix after one hot encodig  (36052, 1)
Shape of senti_com_ess_Ts_norm matrix after one hot encodig  (36052, 1)
```

## 2.3 Make Data Model Ready: encoding essay, and project_title

### 1.5.3 Vectorizing Text data

#### 1.5.3.1 Bag of words

In [28]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer_essays_bow = CountVectorizer(min_df=10)
text_bow_train = vectorizer_essays_bow.fit_transform(preprocessed_essays_train)

text_bow_Test = vectorizer_essays_bow.transform(preprocessed_essays_Test)
print("Shape of matrix after one hot encodig ",text_bow_train.shape)

print("Shape of text_bow_Test ",text_bow_Test.shape)
```

```
Shape of matrix after one hot encodig  (73196, 14144)
Shape of text_bow_Test  (36052, 14144)
```

#### Bag of Words for Project Title

In [29]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer_titles_bow = CountVectorizer(min_df=10)
title_bow_train = vectorizer_titles_bow.fit_transform(preprocessed_titles_train)

title_bow_Test = vectorizer_titles_bow.transform(preprocessed_titles_Test)
print("Shape of matrix (title) after one hot encoding ",title_bow_train.shape)

print("Shape of title_bow_test ",title_bow_Test.shape)
```

```
Shape of matrix (title) after one hot encoding  (73196, 2631)
Shape of title_bow_test  (36052, 2631)
```

**1.5.2.2 TFIDF vectorizer**

```python
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_essays_tfidf = TfidfVectorizer(min_df=10)
text_tfidf_train = vectorizer_essays_tfidf.fit_transform(preprocessed_essays_train)

text_tfidf_Test = vectorizer_essays_tfidf.transform(preprocessed_essays_Test)
print("Shape of matrix after one hot encodig ",text_tfidf_train.shape)

print("Shape of text_tfidf_test ",text_tfidf_Test.shape)
```

```
Shape of matrix after one hot encodig  (73196, 14144)
Shape of text_tfidf_test  (36052, 14144)
```

**TFIDF vectorizer for Project Title**

```python
vectorizer_titles_tfidf = TfidfVectorizer(min_df=10)
title_tfidf_train = vectorizer_titles_tfidf.fit_transform(preprocessed_titles_train)

title_tfidf_Test = vectorizer_titles_tfidf.transform(preprocessed_titles_Test)
print("Shape of matrix(title) after one hot encoding ",title_tfidf_train.shape)

print("Shape of title_tfidf_test ",title_tfidf_Test.shape)
```

```
Shape of matrix(title) after one hot encoding  (73196, 2631)
Shape of title_tfidf_test  (36052, 2631)
```

**1.5.2.3 Using Pretrained Models: Avg W2V**

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

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

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

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

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

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

inter_words = set(model.keys()).intersection(words)
```

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

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


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

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


'''
```

Out[32]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\'r\',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n
word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        m
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\'glove.42B.300d.txt\')\n# ============================\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495  words loaded!\n\n#
============================\n\nwords = []\nfor i in preproced_texts:\n    words.extend(i.split(\'
\'))\n\nfor i in preproced_titles:\n    words.extend(i.split(\' \'))\nprint("all the words in the
coupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus",
len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha
t are present in both glove vectors and our coupus",        len(inter_words),"
(",np.round(len(inter_words)/len(words)*100,3),"%)")\n\nwords_courpus = {}\nwords_glove =
set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\nwith open(\'glove_vectors\', \'wb\') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```

In [0]:

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

**Average Word2Vec for Project_Essays**

In [34]:

```python
avg_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays_train): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train.append(vector)
#-------------------------------------------------------------------------------------


avg_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove words:
```

```
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_Test.append(vector)

print(len(avg_w2v_vectors_Test))
print(len(avg_w2v_vectors_Test[1]))
```

```
100%|██████████| 73196/73196 [00:24<00:00, 2931.05it/s]
100%|██████████| 36052/36052 [00:12<00:00, 2966.94it/s]
```

```
36052
300
```

**AVG W2V on project_title**

In [35]:

```python
# Similarly you can vectorize for title also
# compute average word2vec for each title.
avg_w2v_vectors_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles_train): # for each review/sentence
    vector_title = np.zeros(300) # as word vectors are of zero length
    cnt_title_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector_title += model[word]
            cnt_title_words += 1
    if cnt_title_words != 0:
        vector_title /= cnt_title_words
    avg_w2v_vectors_title_train.append(vector_title)


#-------------------------------------------------------------------------------------------
avg_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
    vector_title = np.zeros(300) # as word vectors are of zero length
    cnt_title_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector_title += model[word]
            cnt_title_words += 1
    if cnt_title_words != 0:
        vector_title /= cnt_title_words
    avg_w2v_vectors_title_Test.append(vector_title)

print(len(avg_w2v_vectors_title_Test))
print(len(avg_w2v_vectors_title_Test[0]))
```

```
100%|██████████| 73196/73196 [00:01<00:00, 50420.89it/s]
100%|██████████| 36052/36052 [00:00<00:00, 50524.73it/s]
```

```
36052
300
```

**1.5.2.3 Using Pretrained Models: TFIDF weighted W2V**

In [0]:

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

**TFIDF weighted W2V for Project_Essays**

```python
tfidf_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays_train): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words_essays):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_train.append(vector)

#-----------------------------------------------------------------------------------------
tfidf_w2v_vectors_Test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words_essays):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_Test.append(vector)

print(len(tfidf_w2v_vectors_Test))
print(len(tfidf_w2v_vectors_Test[0]))
```

```
100%|██████████| 73196/73196 [02:14<00:00, 543.55it/s]
100%|██████████| 36052/36052 [01:06<00:00, 543.00it/s]
```

```
36052
300
```

**TFIDF weighted W2V on project_title**

```python
# Similarly you can vectorize for title also
tfidf_model_title = TfidfVectorizer()
tfidf_model_title.fit(preprocessed_titles_train)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model_title.get_feature_names(), list(tfidf_model_title.idf_)))
tfidf_words_title = set(tfidf_model_title.get_feature_names())

# compute tfidf word2vec for each title.
tfidf_w2v_vectors_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles_train): # for each review/sentence
    vector_title = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words_title):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
```

```
idf value for each word
            vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector_title /= tf_idf_weight
    tfidf_w2v_vectors_title_train.append(vector_title)

#----------------------------------------------------------------------------------------
--------------------------------

tfidf_w2v_vectors_title_Test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
    vector_title = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words_title):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
idf value for each word
            vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector_title /= tf_idf_weight
    tfidf_w2v_vectors_title_Test.append(vector_title)

print(len(tfidf_w2v_vectors_title_Test))
print(len(tfidf_w2v_vectors_title_Test[0]))
```

```
100%|████████| 73196/73196 [00:02<00:00, 29654.23it/s]
100%|████████| 36052/36052 [00:01<00:00, 29392.09it/s]
```

```
36052
300
```

### 1.5.4 Merging all the above features

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

Merging vectorised Train data

In [39]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X1_train =
hstack((categories_one_hot_train,sub_categories_one_hot_train,sch_one_hot_train,grade_one_hot_train
,prefix_one_hot_train, text_bow_train,title_bow_train, price_norm_train, quantity_norm_train, teach
er_prev_post_norm_train, word_count_essay_train, word_count_title_train, senti_pos_ess_Tr_norm,sen
ti_neg_ess_Tr_norm,senti_neut_ess_Tr_norm,senti_com_ess_Tr_norm))
print(X1_train.shape)
X2_train =
hstack((categories_one_hot_train,sub_categories_one_hot_train,sch_one_hot_train,grade_one_hot_train
,prefix_one_hot_train, text_tfidf_train,title_tfidf_train, price_norm_train, quantity_norm_train, t
eacher_prev_post_norm_train, word_count_essay_train, word_count_title_train, senti_pos_ess_Tr_norm
,senti_neg_ess_Tr_norm,senti_neut_ess_Tr_norm,senti_com_ess_Tr_norm))
print(X2_train.shape)
X3_train =
hstack((categories_one_hot_train,sub_categories_one_hot_train,sch_one_hot_train,grade_one_hot_train
,prefix_one_hot_train, avg_w2v_vectors_train,avg_w2v_vectors_title_train, price_norm_train,
quantity_norm_train, teacher_prev_post_norm_train, word_count_essay_train, word_count_title_train,
senti_pos_ess_Tr_norm,senti_neg_ess_Tr_norm,senti_neut_ess_Tr_norm,senti_com_ess_Tr_norm))
print(X3_train.shape)
X4_train =
hstack((categories_one_hot_train,sub_categories_one_hot_train,sch_one_hot_train,grade_one_hot_train
,prefix_one_hot_train, tfidf_w2v_vectors_train,tfidf_w2v_vectors_title_train, price_norm_train,
quantity_norm_train, teacher_prev_post_norm_train, word_count_essay_train, word_count_title_train,
senti_pos_ess_Tr_norm,senti_neg_ess_Tr_norm,senti_neut_ess_Tr_norm,senti_com_ess_Tr_norm))
print(X4_train.shape)
```

```
print(y_train.shape)
```

```
(73196, 16884)
(73196, 16884)
(73196, 709)
(73196, 709)
(73196,)
```

Merging vectorised Test data

In [40]:

```
X1_Test =
hstack((categories_one_hot_Test,sub_categories_one_hot_Test,sch_one_hot_Test,grade_one_hot_Test,pr
efix_one_hot_Test,text_bow_Test,title_bow_Test,price_norm_Test, quantity_norm_Test,
teacher_prev_post_norm_Test, word_count_essay_Test, word_count_title_Test,senti_pos_ess_Ts_norm,se
nti_neg_ess_Ts_norm,senti_neut_ess_Ts_norm,senti_com_ess_Ts_norm))
print(X1_Test.shape)
X2_Test =
hstack((categories_one_hot_Test,sub_categories_one_hot_Test,sch_one_hot_Test,grade_one_hot_Test,pr
efix_one_hot_Test,text_tfidf_Test,title_tfidf_Test,price_norm_Test, quantity_norm_Test,
teacher_prev_post_norm_Test, word_count_essay_Test, word_count_title_Test,senti_pos_ess_Ts_norm,se
nti_neg_ess_Ts_norm,senti_neut_ess_Ts_norm,senti_com_ess_Ts_norm))
print(X2_Test.shape)
X3_Test =
hstack((categories_one_hot_Test,sub_categories_one_hot_Test,sch_one_hot_Test,grade_one_hot_Test,pr
efix_one_hot_Test,avg_w2v_vectors_Test,avg_w2v_vectors_title_Test,price_norm_Test,
quantity_norm_Test, teacher_prev_post_norm_Test, word_count_essay_Test, word_count_title_Test,
senti_pos_ess_Ts_norm,senti_neg_ess_Ts_norm,senti_neut_ess_Ts_norm,senti_com_ess_Ts_norm))
print(X3_Test.shape)
X4_Test =
hstack((categories_one_hot_Test,sub_categories_one_hot_Test,sch_one_hot_Test,grade_one_hot_Test,pr
efix_one_hot_Test,tfidf_w2v_vectors_Test,tfidf_w2v_vectors_title_Test,price_norm_Test,
quantity_norm_Test, teacher_prev_post_norm_Test, word_count_essay_Test, word_count_title_Test,
senti_pos_ess_Ts_norm,senti_neg_ess_Ts_norm,senti_neut_ess_Ts_norm,senti_com_ess_Ts_norm))
print(X4_Test.shape)


print(y_Test.shape)
```

```
(36052, 16884)
(36052, 16884)
(36052, 709)
(36052, 709)
(36052,)
```

# 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 value
   - Find the best hyper paramter using k-fold cross validation or simple cross validation data
   - Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. **Representation of results**

   - You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
   - Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.

- Along with plotting ROC curve, you need to print the confusion matrix with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.

4. **[Task-2] Apply 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** on **TfidfVectorizer** of essay text, choose the number of components (`n_components`) using **elbow method** : 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

**Note: Data Leakage**

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

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

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

## 2.4.1 Applying Support Vector Machines on `Set 1 - BOW`

In [41]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X1_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)


for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
```
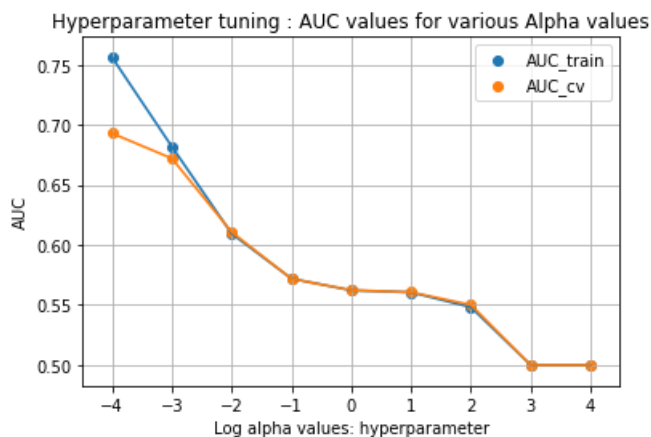
```
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:  3.3min finished
100%|██████████| 9/9 [00:00<00:00, 21998.10it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l1',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6928279877473916
```



```
CPU times: user 2min 36s, sys: 7.93 s, total: 2min 44s
Wall time: 5min 50s
```

In [42]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X1_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter
```
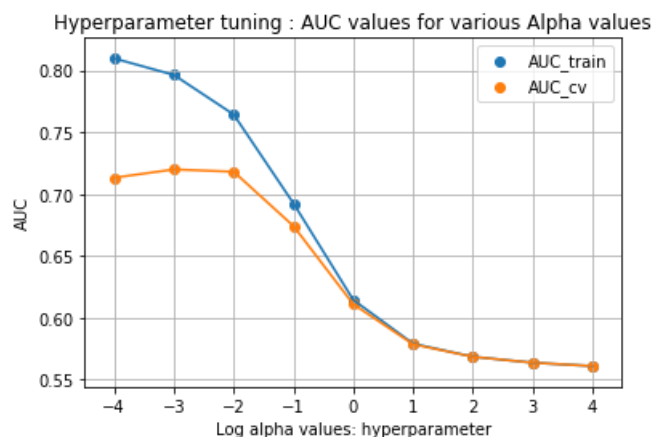
```
# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:   48.7s finished
100%|████████████| 9/9 [00:00<00:00, 26621.11it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.001}
GridSearch_SGDclf.best_score_:  0.7198648775926259
```



CPU times: user 20.5 s, sys: 4.15 s, total: 24.6 s
Wall time: 1min 7s

**Observations:**

- Alpha values are chosen from 0.00001 to 10000 and for the sake of graphical representation the Alpha values are scaled down by applying a log function on them without losing the relationship with their corresponding AUC values.
- Alpha = 0.001 gives the maximum cv score, hence it is considered as the optimal Alpha value.

In [43]:

```
%%time
a1=GridSearch_SGDclf.best_params_['alpha']

SGDclf_opt = SGDClassifier(alpha=a1,loss = 'hinge',penalty='l2', n_jobs=-1)
SGDclf_opt.fit(X1_train, y_train)
pred = SGDclf_opt.predict(X1_Test)
acc = accuracy_score(y_Test, pred, normalize=True) * float(100)
print('\nThe optimal alpha value = {0}'.format(a1))
print('\nTest accuracy for (alpha value = {0}) is {1}%'.format(a1,acc))

fpr_train, tpr_train, thresholds = roc_curve(y_train, SGDclf_opt.decision_function(X1_train.tocsr(
)))
fpr_Test, tpr_Test, thresholds = roc_curve(y_Test, SGDclf_opt.decision_function(X1_Test.tocsr()))
```
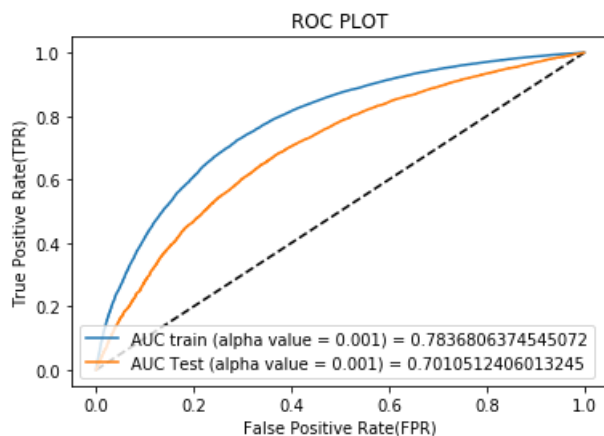
```
#ROC plot
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr_train, tpr_train, label="AUC train (alpha value = {0}) = ".format(a1)+str(auc(fpr_trai
n, tpr_train)))
plt.plot(fpr_Test, tpr_Test, label="AUC Test (alpha value = {0}) = ".format(a1)+str(auc(fpr_Test, t
pr_Test)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("ROC PLOT")
plt.show()
print("AUC value for train data (alpha value = {0}) = ".format(a1), auc(fpr_train, tpr_train))
print("AUC value for Test data (alpha value = {0}) = ".format(a1), auc(fpr_Test, tpr_Test))
print("="*115)

AUC1 = auc(fpr_Test, tpr_Test)
pred0 = SGDclf_opt.predict(X1_train)
pred2 = SGDclf_opt.predict(X1_Test)
```

The optimal alpha value = 0.001

Test accuracy for (alpha value = 0.001) is 31.751359147897478%



ROC PLOT

```
AUC value for train data (alpha value = 0.001) =  0.7836806374545072
AUC value for Test data (alpha value = 0.001) =  0.7010512406013245
===================================================================================================
============
CPU times: user 25.1 s, sys: 3.83 s, total: 28.9 s
Wall time: 24.3 s
```

In [44]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
%matplotlib inline
Train = confusion_matrix(y_train, pred0)
sns.heatmap(Train,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Out[44]:

Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



Project is APPROVED or NOT Confusion Matrix - Train Data

**Observations for train data:** Here we got 10695 - true positives, 13387 - true negatives, 393 - false negatives, 48721 - false positives.

```python
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
Test = confusion_matrix(y_Test, pred2)
sns.heatmap(Test,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

Out[45]:

```
Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```



**Observations for Test data:** Here we got 5070 - true positives, 6377 - true negatives, 384 - false negatives, 24221 - false positives.

## 2.4.2 Applying Support Vector Machines on `Set 2 - TFIDF`

In [46]:

```python
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X2_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
```
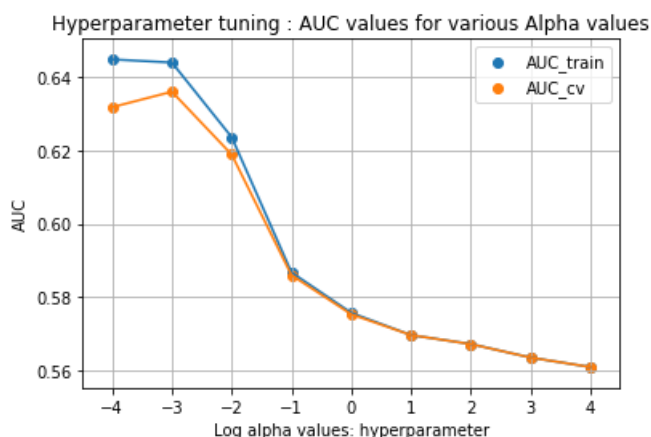
```
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:  2.7min finished
100%|████████| 9/9 [00:00<00:00, 18078.90it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
               early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
               l1_ratio=0.15, learning_rate='optimal', loss='hinge',
               max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l1',
               power_t=0.5, random_state=None, shuffle=True, tol=0.001,
               validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6730503313221478
```



```
CPU times: user 1min 15s, sys: 4.16 s, total: 1min 19s
Wall time: 3min 55s
```

In [47]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X2_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
```

```
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done   90 out of   90 | elapsed:    52.4s finished
100%|████████████| 9/9 [00:00<00:00, 14654.01it/s]
```

```
GridSearch_SGDclf.best_estimator_ :  SGDClassifier(alpha=0.001, average=False, class_weight=None,
               early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
               l1_ratio=0.15, learning_rate='optimal', loss='hinge',
               max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l2',
               power_t=0.5, random_state=None, shuffle=True, tol=0.001,
               validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_ :  {'alpha': 0.001}
GridSearch_SGDclf.best_score_ :  0.6359272282669068
```



CPU times: user 27.7 s, sys: 4.19 s, total: 31.9 s
Wall time: 1min 18s

**Observations:**

- Alpha values are chosen from 0.00001 to 10000 and for the sake of graphical representation the Alpha values are scaled down by applying a log function on them without losing the relationship with their corresponding AUC values.
- Alpha = 0.001 gives the maximum cv score, hence it is considered as the optimal Alpha value.

In [48]:

```
%%time
a2=GridSearch_SGDclf.best_params_['alpha']

SGDclf_opt = SGDClassifier(alpha=a2,loss = 'hinge',penalty='l2', n_jobs=-1)
SGDclf_opt.fit(X2_train, y_train)
pred = SGDclf_opt.predict(X2_Test)
acc = accuracy_score(y_Test, pred, normalize=True) * float(100)
print('\nThe optimal alpha value = {0}'.format(a2))
print('\nTest accuracy for (alpha value = {0}) is {1}%'.format(a2,acc))

fpr_train, tpr_train, thresholds = roc_curve(y_train, SGDclf_opt.decision_function(X2_train.tocsr(
)))
fpr_Test, tpr_Test, thresholds = roc_curve(y_Test, SGDclf_opt.decision_function(X2_Test.tocsr()))
```

```
#ROC plot
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr_train, tpr_train, label="AUC train (alpha value = {0}) = ".format(a2)+str(auc(fpr_trai
n, tpr_train)))
plt.plot(fpr_Test, tpr_Test, label="AUC Test (alpha value = {0}) = ".format(a2)+str(auc(fpr_Test, t
pr_Test)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("ROC PLOT")
plt.show()
print("AUC value for train data (alpha value = {0}) = ".format(a2), auc(fpr_train, tpr_train))
print("AUC value for Test data (alpha value = {0}) = ".format(a2), auc(fpr_Test, tpr_Test))
print("="*115)

AUC2 = auc(fpr_Test, tpr_Test)
pred0 = SGDclf_opt.predict(X2_train)
pred2 = SGDclf_opt.predict(X2_Test)
```
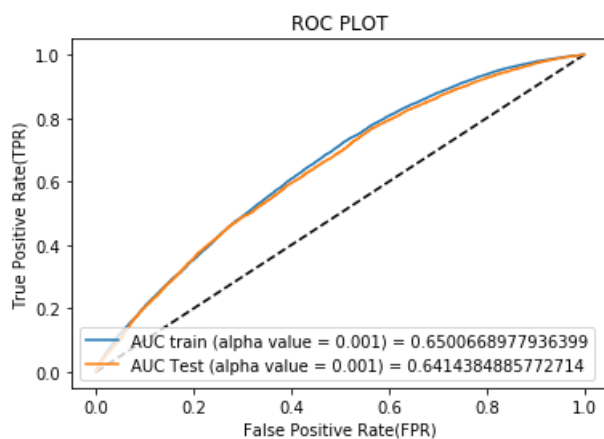
The optimal alpha value = 0.001

Test accuracy for (alpha value = 0.001) is 84.86907799844668%



ROC PLOT

```
AUC value for train data (alpha value = 0.001) =  0.6500668977936399
AUC value for Test data (alpha value = 0.001) =  0.6414384885772714
===========================================================================================================
============
CPU times: user 30.2 s, sys: 3.74 s, total: 33.9 s
Wall time: 29.4 s
```

In [49]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
%matplotlib inline
Train = confusion_matrix(y_train, pred0)
sns.heatmap(Train,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Out[49]:

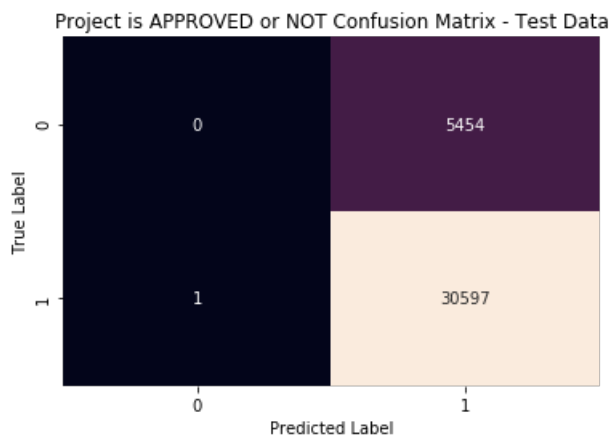Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



Project is APPROVED or NOT Confusion Matrix - Train Data

**Observations for train data:** Here we got 0 - true positives, 62100 - true negatives, 11088 - false negatives, 8 - false positives.

```python
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
Test = confusion_matrix(y_Test, pred2)
sns.heatmap(Test,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

Out[50]:

```
Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```



**Observations for Test data:** Here we got 0 - true positives, 30597 - true negatives, 5454 - false negatives, 1 - false positives.

## 2.4.3 Applying Support Vector Machines on `Set 3 - AVGW2V`

In [51]:

```python
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X3_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
```

```
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

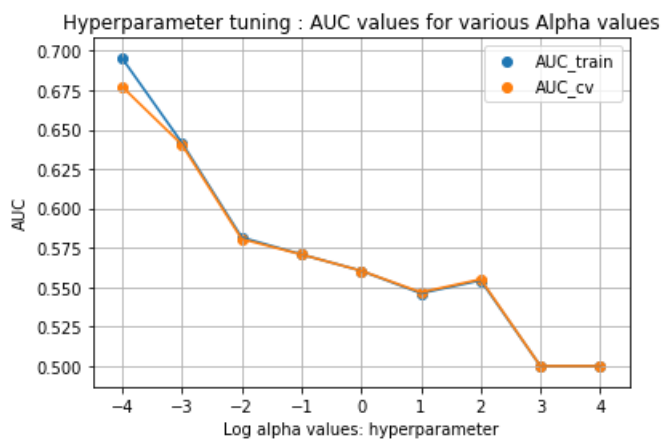Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:  6.4min finished
100%|██████████| 9/9 [00:00<00:00, 24244.53it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l1',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6771048085020802
```



```
CPU times: user 2min 36s, sys: 4.88 s, total: 2min 40s
Wall time: 8min 55s
```

In [52]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X3_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
```
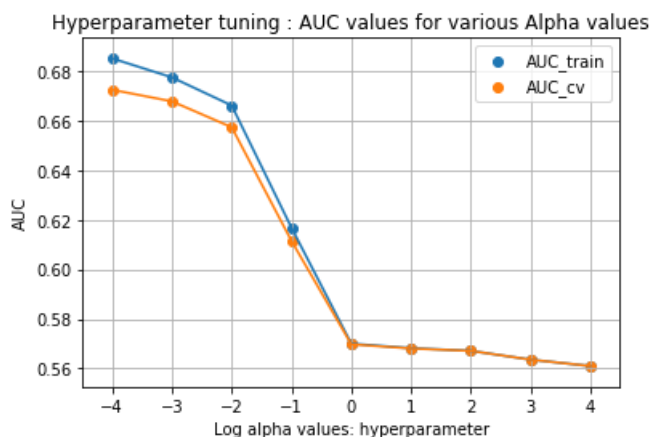
```
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:  1.5min finished
100%|██████████| 9/9 [00:00<00:00, 28926.23it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
            early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
            l1_ratio=0.15, learning_rate='optimal', loss='hinge',
            max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l2',
            power_t=0.5, random_state=None, shuffle=True, tol=0.001,
            validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6724263467862011
```



```
CPU times: user 58 s, sys: 4.77 s, total: 1min 2s
Wall time: 2min 25s
```

**Observations:**

- Alpha values are chosen from 0.00001 to 10000 and for the sake of graphical representation the Alpha values are scaled down by applying a log function on them without losing the relationship with their corresponding AUC values.
- Alpha = 0.0001 gives the maximum cv score, hence it is considered as the optimal Alpha value.

In [53]:

```
%%time
a3=GridSearch_SGDclf.best_params_['alpha']

SGDclf_opt = SGDClassifier(alpha=a1,loss = 'hinge',penalty='l2', n_jobs=-1)
SGDclf_opt.fit(X3_train, y_train)
pred = SGDclf_opt.predict(X3_Test)
acc = accuracy_score(y_Test, pred, normalize=True) * float(100)
print('\nThe optimal alpha value = {0}'.format(a3))
print('\nTest accuracy for (alpha value = {0}) is {1}%'.format(a3,acc))

fpr_train, tpr_train, thresholds = roc_curve(y_train, SGDclf_opt.decision_function(X3_train.tocsr(
)))
fpr_Test, tpr_Test, thresholds = roc_curve(y_Test, SGDclf_opt.decision_function(X3_Test.tocsr()))

#ROC plot
```
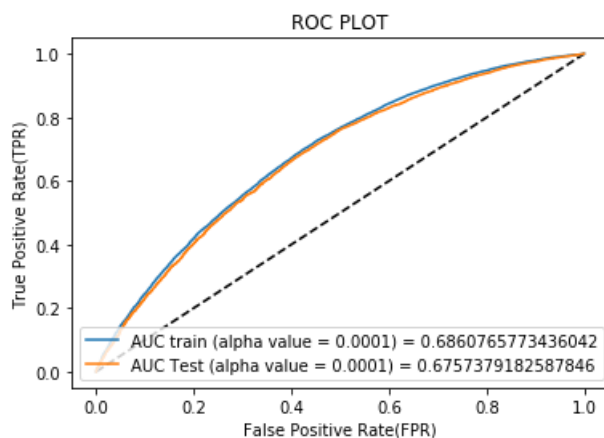
```
#ROC plot
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr_train, tpr_train, label="AUC train (alpha value = {0}) = ".format(a3)+str(auc(fpr_trai
n, tpr_train)))
plt.plot(fpr_Test, tpr_Test, label="AUC Test (alpha value = {0}) = ".format(a3)+str(auc(fpr_Test, t
pr_Test)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("ROC PLOT")
plt.show()
print("AUC value for train data (alpha value = {0}) = ".format(a3), auc(fpr_train, tpr_train))
print("AUC value for Test data (alpha value = {0}) = ".format(a3), auc(fpr_Test, tpr_Test))
print("="*115)

AUC3 = auc(fpr_Test, tpr_Test)
pred0 = SGDclf_opt.predict(X3_train)
pred2 = SGDclf_opt.predict(X3_Test)
```

The optimal alpha value = 0.0001

Test accuracy for (alpha value = 0.0001) is 84.80250748918229%



ROC PLOT

AUC value for train data (alpha value = 0.0001) =  0.6860765773436042
AUC value for Test data (alpha value = 0.0001) =  0.6757379182587846
=================================================================================================
============
CPU times: user 37.5 s, sys: 8.03 ms, total: 37.5 s
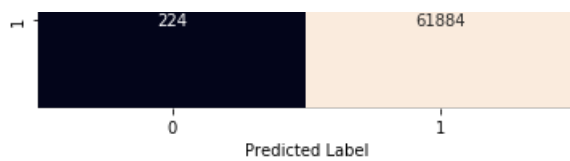Wall time: 37.5 s

In [54]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
%matplotlib inline
Train = confusion_matrix(y_train, pred0)
sns.heatmap(Train,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Out[54]:

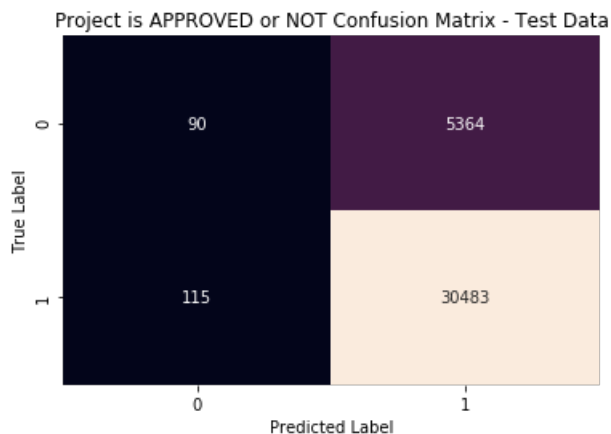Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



Project is APPROVED or NOT Confusion Matrix - Train Data

**Observations for train data:** Here we got 226 - true positives, 61884 - true negatives, 10862 - false negatives, 224 - false positives.

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
Test = confusion_matrix(y_Test, pred2)
sns.heatmap(Test,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

```
Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```



**Observations for Test data:** Here we got 90 - true positives, 30483 - true negatives, 5364 - false negatives, 115 - false positives.

## 2.4.4 Applying Support Vector Machines on `Set 4 - TFIDF weighted W2V`

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X4_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
```
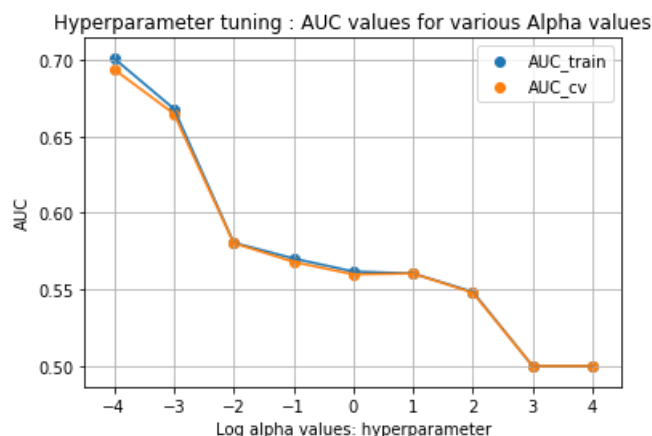
```
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:  3.7min finished
100%|██████████| 9/9 [00:00<00:00, 13914.02it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
                early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
                l1_ratio=0.15, learning_rate='optimal', loss='hinge',
                max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l1',
                power_t=0.5, random_state=None, shuffle=True, tol=0.001,
                validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6936736999012151
```



```
CPU times: user 1min 6s, sys: 415 ms, total: 1min 7s
Wall time: 4min 50s
```

In [57]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-1)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-1, scoring='roc_auc',return_t
rain_score=True,verbose=1)
GridSearch_SGDclf.fit(X4_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
```
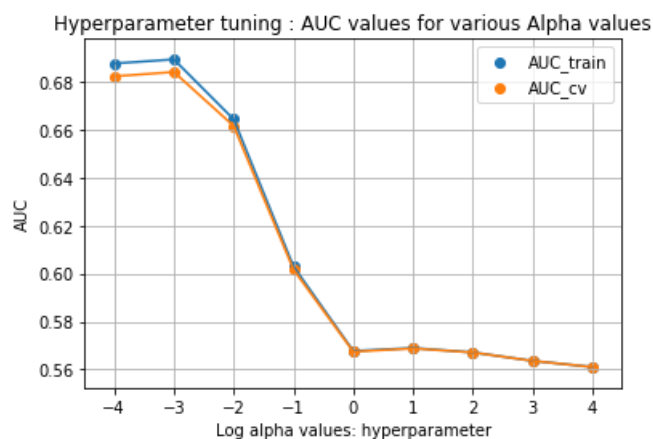
```
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 40 concurrent workers.
[Parallel(n_jobs=-1)]: Done  90 out of  90 | elapsed:   55.4s finished
100%|██████████| 9/9 [00:00<00:00, 28255.04it/s]
```

```
GridSearch_SGDclf.best_estimator_ :  SGDClassifier(alpha=0.001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=-1, penalty='l2',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_ :  {'alpha': 0.001}
GridSearch_SGDclf.best_score_ :  0.6840629457497863
```



```
CPU times: user 27.9 s, sys: 392 ms, total: 28.3 s
Wall time: 1min 22s
```

**Observations:**

- Alpha values are chosen from 0.00001 to 10000 and for the sake of graphical representation the Alpha values are scaled down by applying a log function on them without losing the relationship with their corresponding AUC values.
- Alpha = 0.001 gives the maximum cv score, hence it is considered as the optimal Alpha value.

In [58]:

```
%%time
a4=GridSearch_SGDclf.best_params_['alpha']

SGDclf_opt = SGDClassifier(alpha=a4,loss = 'hinge',penalty='l2', n_jobs=-1)
SGDclf_opt.fit(X4_train, y_train)
pred = SGDclf_opt.predict(X4_Test)
acc = accuracy_score(y_Test, pred, normalize=True) * float(100)
print('\nThe optimal alpha value = {0}'.format(a4))
print('\nTest accuracy for (alpha value = {0}) is {1}%'.format(a4,acc))

fpr_train, tpr_train, thresholds = roc_curve(y_train, SGDclf_opt.decision_function(X4_train.tocsr(
)))
fpr_Test, tpr_Test, thresholds = roc_curve(y_Test, SGDclf_opt.decision_function(X4_Test.tocsr()))

#ROC plot
plt.plot([0,1],[0,1],'k--')
```
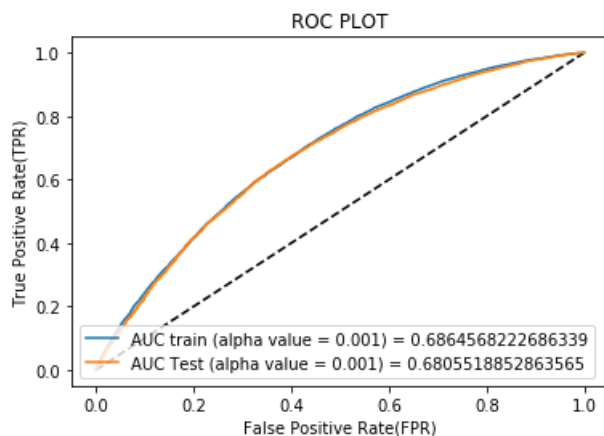
```python
plt.plot(fpr_train, tpr_train, label="AUC train (alpha value = {0}) = ".format(a4)+str(auc(fpr_trai
n, tpr_train)))
plt.plot(fpr_Test, tpr_Test, label="AUC Test (alpha value = {0}) = ".format(a4)+str(auc(fpr_Test, t
pr_Test)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("ROC PLOT")
plt.show()
print("AUC value for train data (alpha value = {0}) = ".format(a4), auc(fpr_train, tpr_train))
print("AUC value for Test data (alpha value = {0}) = ".format(a4), auc(fpr_Test, tpr_Test))
print("="*115)

AUC4 = auc(fpr_Test, tpr_Test)
pred0 = SGDclf_opt.predict(X4_train)
pred2 = SGDclf_opt.predict(X4_Test)
```

The optimal alpha value = 0.001

Test accuracy for (alpha value = 0.001) is 84.96893376234328%



```
AUC value for train data (alpha value = 0.001) =  0.6864568222686339
AUC value for Test data (alpha value = 0.001) =  0.6805518852863565
=================================================================================================
============
CPU times: user 24.1 s, sys: 7.7 ms, total: 24.1 s
Wall time: 24.1 s
```
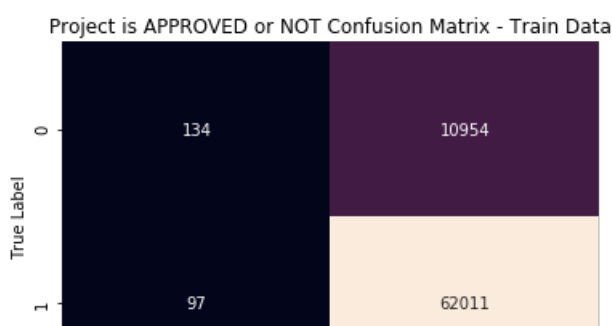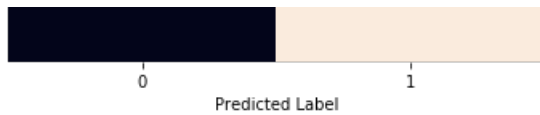
In [59]:

```python
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
%matplotlib inline
Train = confusion_matrix(y_train, pred0)
sns.heatmap(Train,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Out[59]:

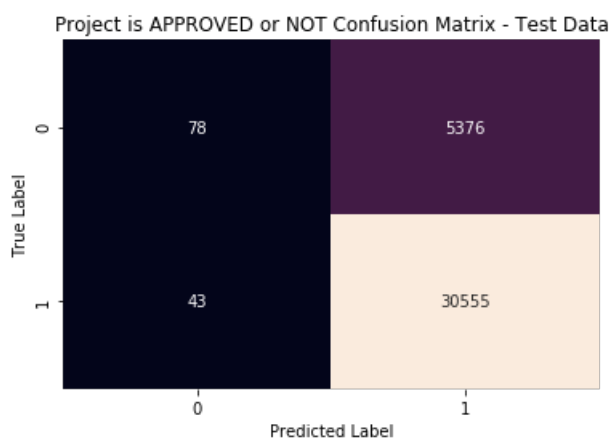Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')

**Observations for train data:** Here we got 134 - true positives, 62011 - true negatives, 10954 - false negatives, 97 - false positives.

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
Test = confusion_matrix(y_Test, pred2)
sns.heatmap(Test,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

```
Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```



**Observations for Test data:** Here we got 78 - true positives, 30555 - true negatives, 5376 - false negatives, 43 - false positives.
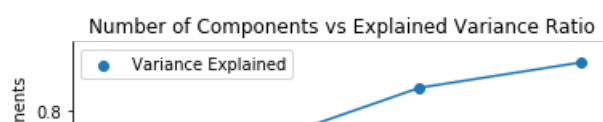
## 2.4.5 Applying Support Vector Machines on 'Set 5' (TFIDF-Truncated SVD)
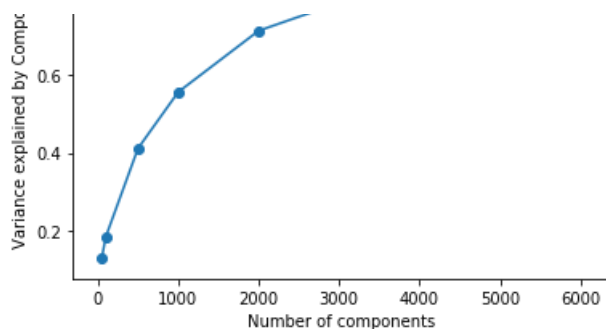
```
# https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/
# https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html
no_of_comp = [50,100,500,1000,2000,4000,6000]
agg_variance = []
for i in tqdm(no_of_comp):
    tsvd1 = TruncatedSVD(n_components =i,random_state = 0,n_iter = 3 )
    tsvd1.fit(text_tfidf_train)
    agg_variance.append(tsvd1.explained_variance_ratio_.sum())

plt.plot(no_of_comp,agg_variance)
plt.scatter(no_of_comp,agg_variance,label = 'Variance Explained')
plt.xlabel("Number of components")
plt.ylabel("Variance explained by Components")
plt.legend()
plt.title("Number of Components vs Explained Variance Ratio")
sns.despine()
plt.show()
```

```
100%|██████████| 7/7 [18:40<00:00, 260.88s/it]
```

**Observations:**

1. Using Elbow method the number of components is chosen as 2000 since a sharp break can be observed at 2000 with 70% of the Variance explained.
2. The Range for number of components and the iteration are kept to a minimum due to memory (RAM) constraint.

In [62]:

```
tsvd2 = TruncatedSVD(n_components = 2000,n_iter = 3)
tsvd2.fit(text_tfidf_train)
tsvd_train = tsvd2.transform(text_tfidf_train)
tsvd_Test = tsvd2.transform(text_tfidf_Test)
print(len(tsvd_train))
print(len(tsvd_Test))
```

```
73196
36052
```

In [63]:

```
X5_train =
hstack((categories_one_hot_train,sub_categories_one_hot_train,sch_one_hot_train,grade_one_hot_train
,prefix_one_hot_train, price_norm_train, quantity_norm_train, teacher_prev_post_norm_train,
word_count_essay_train, word_count_title_train,
senti_pos_ess_Tr_norm,senti_neg_ess_Tr_norm,senti_neut_ess_Tr_norm,senti_com_ess_Tr_norm,tsvd_train
))
print(X5_train.shape)
X5_Test =
hstack((categories_one_hot_Test,sub_categories_one_hot_Test,sch_one_hot_Test,grade_one_hot_Test,pr
efix_one_hot_Test,price_norm_Test, quantity_norm_Test, teacher_prev_post_norm_Test,
word_count_essay_Test, word_count_title_Test,
senti_pos_ess_Ts_norm,senti_neg_ess_Ts_norm,senti_neut_ess_Ts_norm,senti_com_ess_Ts_norm,tsvd_Test
))
print(X5_Test.shape)
```

```
(73196, 2109)
(36052, 2109)
```

In [64]:

```
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l1', n_jobs=-2)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, n_jobs=-2, cv = 10, scoring='roc_auc',return_t
rain_score=True,verbose=1,pre_dispatch='4')
GridSearch_SGDclf.fit(X5_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
    b = np.log10(av)
```

```python
    log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

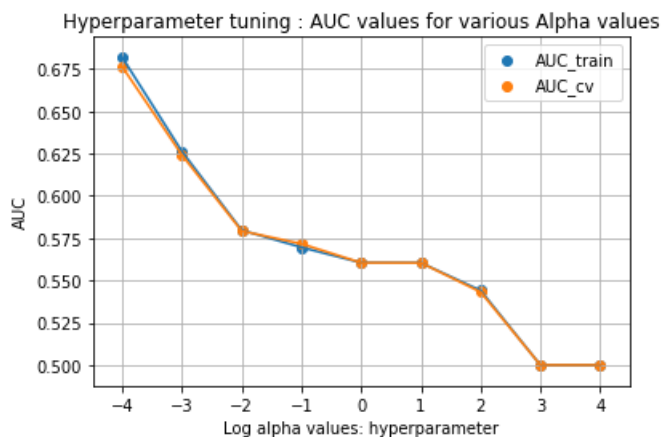Fitting 10 folds for each of 9 candidates, totalling 90 fits

```
[Parallel(n_jobs=-2)]: Using backend LokyBackend with 39 concurrent workers.
[Parallel(n_jobs=-2)]: Done   46 tasks      | elapsed: 32.3min
[Parallel(n_jobs=-2)]: Done   90 out of  90 | elapsed: 94.7min finished
100%|██████████| 9/9 [00:00<00:00, 17166.32it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
              early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
              l1_ratio=0.15, learning_rate='optimal', loss='hinge',
              max_iter=1000, n_iter_no_change=5, n_jobs=-2, penalty='l1',
              power_t=0.5, random_state=None, shuffle=True, tol=0.001,
              validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6763627276686598
```



```
CPU times: user 6min 31s, sys: 36.9 s, total: 7min 8s
Wall time: 1h 41min 14s
```

In [65]:

```python
%%time
alphas=[1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]
log_alphas = []
auc_scores_train = []
auc_scores_cv = []
SGDclf=SGDClassifier(loss = 'hinge',penalty='l2', n_jobs=-2)
parameters={'alpha' : [1e-4,1e-3,1e-2,1e-1,1,1e+1,1e+2,1e+3,1e+4]}
GridSearch_SGDclf = GridSearchCV(SGDclf, parameters, cv = 10, n_jobs=-2, scoring='roc_auc',return_t
rain_score=True,verbose=1,pre_dispatch='4')
GridSearch_SGDclf.fit(X5_train, y_train)
auc_scores_train= GridSearch_SGDclf.cv_results_['mean_train_score']
auc_scores_cv = GridSearch_SGDclf.cv_results_['mean_test_score']
print('GridSearch_SGDclf.best_estimator_: ', GridSearch_SGDclf.best_estimator_)
print('GridSearch_SGDclf.best_params_: ', GridSearch_SGDclf.best_params_)
print('GridSearch_SGDclf.best_score_: ', GridSearch_SGDclf.best_score_)

for av in tqdm(alphas):
```

```
        b = np.log10(av)
        log_alphas.append(b)

# Performance of model on Train data and Test data for each hyper parameter.
plt.plot(log_alphas, auc_scores_train)
plt.gca()
plt.plot(log_alphas, auc_scores_cv)
plt.gca()
plt.scatter(log_alphas, auc_scores_train, label='AUC_train')
plt.scatter(log_alphas, auc_scores_cv, label='AUC_cv')
plt.legend()
plt.xlabel("Log alpha values: hyperparameter")
plt.ylabel("AUC")
plt.title("Hyperparameter tuning : AUC values for various Alpha values")
plt.grid()
plt.show()
```

Fitting 10 folds for each of 9 candidates, totalling 90 fits
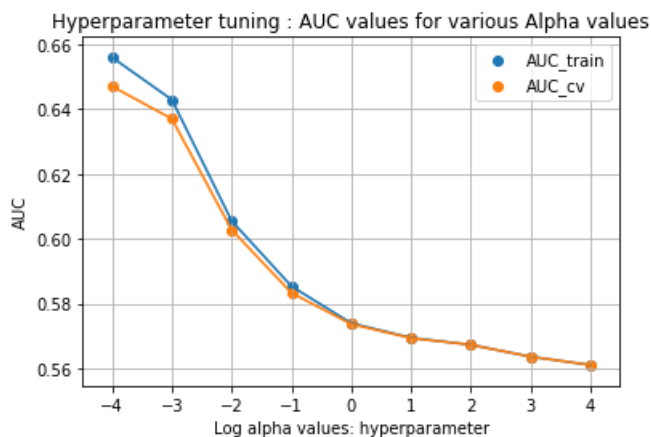
```
[Parallel(n_jobs=-2)]: Using backend LokyBackend with 39 concurrent workers.
[Parallel(n_jobs=-2)]: Done   46 tasks      | elapsed: 13.3min
[Parallel(n_jobs=-2)]: Done   90 out of  90 | elapsed: 14.1min finished
100%|████████████| 9/9 [00:00<00:00, 19113.28it/s]
```

```
GridSearch_SGDclf.best_estimator_:  SGDClassifier(alpha=0.0001, average=False, class_weight=None,
               early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
               l1_ratio=0.15, learning_rate='optimal', loss='hinge',
               max_iter=1000, n_iter_no_change=5, n_jobs=-2, penalty='l2',
               power_t=0.5, random_state=None, shuffle=True, tol=0.001,
               validation_fraction=0.1, verbose=0, warm_start=False)
GridSearch_SGDclf.best_params_:  {'alpha': 0.0001}
GridSearch_SGDclf.best_score_:  0.6470792168058934
```



```
CPU times: user 2min 54s, sys: 17 s, total: 3min 11s
Wall time: 17min
```

**Observations:**

- Alpha values are chosen from 0.00001 to 10000 and for the sake of graphical representation the Alpha values are scaled down by applying a log function on them without losing the relationship with their corresponding AUC values.
- Alpha = 0.0001 gives the maximum cv score, hence it is considered as the optimal Alpha value.

In [66]:

```
%%time
a5=GridSearch_SGDclf.best_params_['alpha']

SGDclf_opt = SGDClassifier(alpha=a5,loss = 'hinge',penalty='l2', n_jobs=-2)
SGDclf_opt.fit(X5_train, y_train)
pred = SGDclf_opt.predict(X5_Test)
acc = accuracy_score(y_Test, pred, normalize=True) * float(100)
print('\nThe optimal alpha value = {0}'.format(a5))
```

```
print('\nTest accuracy for (alpha value = {0}) is {1}%'.format(a5,acc))

fpr_train, tpr_train, thresholds = roc_curve(y_train, SGDclf_opt.decision_function(X5_train.tocsr(
)))
fpr_Test, tpr_Test, thresholds = roc_curve(y_Test, SGDclf_opt.decision_function(X5_Test.tocsr()))

#ROC plot
plt.plot([0,1],[0,1],'k--')
plt.plot(fpr_train, tpr_train, label="AUC train (alpha value = {0}) = ".format(a5)+str(auc(fpr_trai
n, tpr_train)))
plt.plot(fpr_Test, tpr_Test, label="AUC Test (alpha value = {0}) = ".format(a5)+str(auc(fpr_Test, t
pr_Test)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("ROC PLOT")
plt.show()
print("AUC value for train data (alpha value = {0}) = ".format(a5), auc(fpr_train, tpr_train))
print("AUC value for Test data (alpha value = {0}) = ".format(a5), auc(fpr_Test, tpr_Test))
print("="*115)

AUC5 = auc(fpr_Test, tpr_Test)
pred0 = SGDclf_opt.predict(X5_train)
pred2 = SGDclf_opt.predict(X5_Test)
```
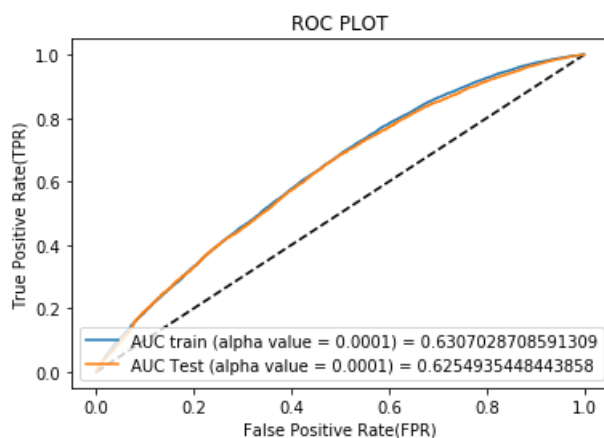
The optimal alpha value = 0.0001

Test accuracy for (alpha value = 0.0001) is 84.87185176966604%



```
AUC value for train data (alpha value = 0.0001) =  0.6307028708591309
AUC value for Test data (alpha value = 0.0001) =  0.6254935448443858
================================================================================================
============
CPU times: user 2min 51s, sys: 442 ms, total: 2min 51s
Wall time: 2min 51s
```

In [67]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
%matplotlib inline
Train = confusion_matrix(y_train, pred0)
sns.heatmap(Train,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```
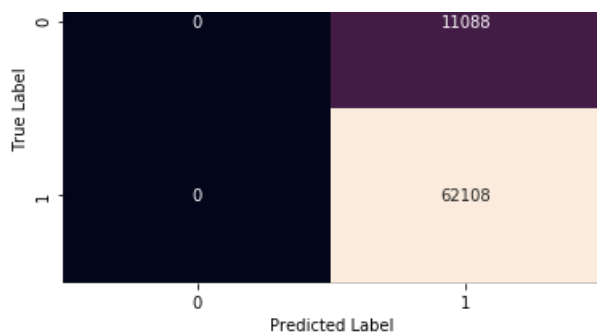
Out[67]:

Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



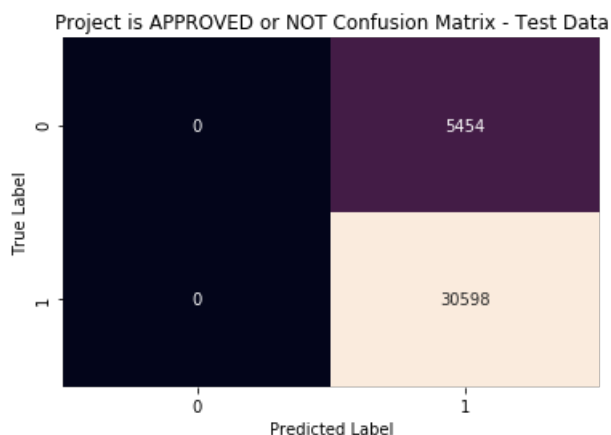Project is APPROVED or NOT Confusion Matrix - Train Data

**Observations for train data:** Here we got 0 - true positives, 62108 - true negatives, 11088 - false negatives, 0 - false positives.

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
#https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
#https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
Test = confusion_matrix(y_Test, pred2)
sns.heatmap(Test,annot=True,cbar=False,fmt='g')
plt.ylabel('True Label')
plt.xlabel('Predicted Label')
plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

Out[68]:

```
Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```

Project is APPROVED or NOT Confusion Matrix - Test Data



**Observations for Test data:** Here we got 0 - true positives, 30598 - true negatives, 5454 - false negatives, 0 - false positives.

# 3. Conclusions

```
from prettytable import PrettyTable
pt = PrettyTable()
pt.field_names = ["Vectorizer", "Model", "Alpha : Hyper Parameter", " Test AUC"]
pt.add_row(["BOW", "SVM", a1, AUC1])
pt.add_row(["TFIDF", "SVM", a2, AUC2])
pt.add_row(["AVGW2V", "SVM", a3, AUC3])
pt.add_row(["TFIDFW2V", "SVM", a4, AUC4])
pt.add_row(["Set 5 (TFIDF-Truncated SVD)", "SVM", a5, AUC5])
print(pt)
```

| Vectorizer | Model | Alpha : Hyper Parameter | Test AUC |
|---|---|---|---|
| BOW | SVM | 0.001 | 0.7010512406013245 |
| TFIDF | SVM | 0.001 | 0.6414384885772714 |
| AVGW2V | SVM | 0.0001 | 0.6757379182587846 |
| TFIDFW2V | SVM | 0.001 | 0.6805518852863565 |

```
| Set 5 (TFIDF-Truncated SVD) |  SVM  |         0.0001         | 0.6254935448443858 |
+-----------------------------+-------+------------------------+--------------------+
```

**SUMMARY:**

1. 'SVM' model's space and time consumption is higher in comparison with Logistic regression model.
2. The Bag of Words model gives the highest AUC score out of all the five models.
3. Since the intervals between the range of number of components for Truncated SVD method can vary the results (the knee point of the Elbow curve), it has to be chosen carefully. Thus, this way does not give a reliable value for the number of components.