

Table of Contents

- [1 Author : SAURABH SINGHAI](#)
- [2 About the DonorsChoose Data Set](#)
 - [2.1 Notes on the Essay Data](#)
- [3 Import Libraries](#)
- [4 Reading Data](#)
- [5 Approved And Non Approved Projects](#)
- [6 Project Dataframe shape and Column Values](#)
- [7 Resource data shape and Column Values](#)
- [8 Preprocessing of project_subject_categories](#)
- [9 Preprocessing of project_subject_subcategories](#)
- [10 Preprocessing of project_grade_category](#)
- [11 Create new column 'Essay' by merging all project Essays](#)
- [12 Use Decontraction function to decontract project essay](#)
- [13 Remove line breaks](#)
- [14 Remove Special Chars](#)
- [15 Remove Stopwards and Join the essays](#)
- [16 Drop essay columns 1, 2, 3, 4](#)
- [17 Preprocessing of project_title](#)
- [18 Drop column project_title and use Cleaned_Title](#)
- [19 Add up the price based on project id](#)
- [20 Adding the word count of essay and title as new columns](#)
- [21 Separate out the Dependant and independant variables](#)
- [22 Splitting data into Test,Train,CV](#)
- [23 Vectorize the features](#)
 - [23.1 Vectorize the Categorical Features - categories](#)
 - [23.2 Vectorize the Categorical Features - subcategories](#)
 - [23.3 Vectorize the Categorical Features - school state](#)
 - [23.4 Vectorize the Categorical Features - teacher prefix](#)
 - [23.5 Vectorize the Categorical Features - project_grade_category](#)
 - [23.6 Vectorize the Numerical Features - price](#)
 - [23.7 Vectorize the Numerical Features - quantity](#)
 - [23.8 Vectorize the Numerical Features - essay count](#)
 - [23.9 Vectorize the Numerical Features - title count](#)
 - [23.10 Vectorizing Text data](#)
 - [23.10.1 Bag of words - essay](#)
 - [23.10.2 Bag of words - cleaned title](#)
 - [23.10.3 TFIDF vectorizer - essay](#)
 - [23.10.4 TFIDF vectorizer - cleaned tittle](#)
 - [23.10.5 Using Pretrained Models: Avg W2V](#)
 - [23.10.6 Using Pretrained Models: TFIDF weighted W2V - Essay](#)
 - [23.10.7 Using Pretrained Models: TFIDF weighted W2V - Cleaned Title](#)
- [24 Merging all the above features](#)
 - [24.1 BOW](#)
 - [24.2 TFIDF](#)
 - [24.3 Word2Vec](#)
 - [24.4 TFIDF- WORD2VEC](#)
- [25 Apply SVM on BOW](#)
 - [25.1 Find best Hyper-Parameter value to train model](#)
 - [25.2 Use best Hyper-Parameter value to train model](#)
 - [25.3 Plot Confusion Matrix for Test Data](#)
 - [25.4 Plot Confusion Matrix for Train Data](#)
- [26 Apply SVM on TFIDF](#)
 - [26.1 Find best Hyper-Parameter value to train model](#)
 - [26.2 Use best Hyper-Parameter value to train model](#)
 - [26.3 Plot Confusion Matrix for Test Data](#)
 - [26.4 Plot Confusion Matrix for Train Data](#)
- [27 Apply SVM on W2V](#)
 - [27.1 Find best Hyper-Parameter value to train model](#)
 - [27.2 Use best Hyper-Parameter value to train model](#)

- [27.3 Plot Confusion Matrix for Test Data](#)
- [27.4 Plot Confusion Matrix for Train Data](#)
- [28 Apply SVM on TFIDFW2V](#)
 - [28.1 Find best Hyper-Parameter value to train model](#)
 - [28.2 Use best Hyper-Parameter value to train model](#)
 - [28.3 Plot Confusion Matrix for Test Data](#)
 - [28.4 Plot Confusion Matrix for Train Data](#)
- [29 Apply TruncatedSVD on TfidfVectorizer of essay text. choose the number of components \(n_components \) using elbow method :numerical data](#)
 - [29.1 Find the best dimension](#)
 - [29.2 Use best dimension value to train TruncatedSVD](#)
 - [29.3 Merge data](#)
- [30 Apply SVM on Set 5](#)
 - [30.1 Find best Hper-parameter to train the model](#)
 - [30.2 Use best Hper-parameter to train the model](#)
 - [30.3 Plot Confusion Matrix for Test Data](#)
 - [30.4 Plot Confusion Matrix for Train Data](#)
- [31 Conclusion](#)

Author : SAURABH SINGHAI

DonorsChoose

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

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

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

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school.

DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

About the DonorsChoose Data Set

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

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

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

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

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

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

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

remove 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.

Import Libraries

In [1]:

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

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

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

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

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

from tqdm import tqdm
import os

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

IOPub data rate exceeded.
The notebook server will temporarily stop sending output
to the client in order to avoid crashing it.
To change this limit, set the config variable
`--NotebookApp.iopub_data_rate_limit`.

Reading Data

In [2]:

```
# *****PLEASE NOTE--Considering 50K points as system becomes Unresponsive with higher no of points

project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
project_data= project_data.sample(n=50000, random_state=0)
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna(' ')
```

Approved And Non Approved Projects

In [3]:

```
y_value_counts = project_data['project_is_approved'].value_counts()

print("Project Not Approved & Approved Count=\n",y_value_counts)

print("Number of projects approved for funding ", y_value_counts[1], "=", (y_value_counts[1]/(y_value_counts[1]+y_value_counts[0]))*100,"%")

print("Number of projects not approved for funding ", y_value_counts[0], "=", (y_value_counts[0]/(y_value_counts[1]+y_value_counts[0]))*100,"%")
```

```
Project Not Approved & Approved Count=
1      42460
0       7540
Name: project_is_approved, dtype: int64
Number of projects approved for funding  42460 = 84.92 %
Number of projects not approved for funding  7540 = 15.08 %
```

Project Dataframe shape and Column Values

In [4]:

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

```
Number of data points in project data (50000, 17)
*****
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

Resource data shape and Column Values

In [5]:

```
print("Number of data points in resource data", resource_data.shape)
print('*'*50)
print(resource_data.columns.values)
resource_data.head(5)
```

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

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

Out[5]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95
2	p069063	Cory Stories: A Kid's Book About Living With Adhd	1	8.45
3	p069063	Dixon Ticonderoga Wood-Cased #2 HB Pencils, Bo...	2	13.59
4	p069063	EDUCATIONAL INSIGHTS FLUORESCENT LIGHT FILTERS...	3	24.95

Preprocessing of project_subject_categories

In [6]:

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

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

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

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

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

Preprocessing of project_subject_subcategories

In [7]:

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

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

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
```

```

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

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

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

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

```

Preprocessing of project_grade_category

In [8]:

```
print(project_data["project_grade_category"].values[0:10])
```

```
['Grades 3-5' 'Grades 6-8' 'Grades 3-5' 'Grades 3-5' 'Grades PreK-2'
 'Grades PreK-2' 'Grades 3-5' 'Grades 3-5' 'Grades PreK-2' 'Grades 9-12']
```

In [9]:

```

project_data["project_grade_category"] =
project_data["project_grade_category"].str.replace("Grades ", "")
project_data["project_grade_category"] = project_data["project_grade_category"].str.replace("-", "_")

print(project_data["project_grade_category"].values[0:10])

```

```
['3_5' '6_8' '3_5' '3_5' 'PreK_2' 'PreK_2' '3_5' '3_5' 'PreK_2' '9_12']
```

Create new column 'Essay' by merging all project Essays

In [10]:

```

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

```

In [11]:

```
project_data['essay'].head(5)
```

Out[11]:

```

75155    Starting the new year off right sets the tone ...
77488    Have you ever worked so hard on a project only...
7803     My students come to class every day ready to l...
56268    \"We love science in your class!\" CJ exclaime...
46902    My students are caring, outgoing, and creative...
Name: essay, dtype: object

```

In [12]:

In [12]:

```
# printing some random reviews
print(project_data['essay'].values[0])
```

Starting the new year off right sets the tone for months to come. My class will be thrilled to receive basic supplies to help them be successful.\r\n\r\nMy students are curious, inquisitive, and enthusiastic learners who enjoy school.\r\n\r\nOur school is a public community school in New York City that receives Title I funding, which means that many students are eligible for free or reduced price lunch. Most of my students are English language learners. Our self-contained class is comprised of students with disabilities in second and third grade. We need printer ink so we can showcase our wonderful work, and other supplies such as pocket charts for subject-specific word walls.\r\n\r\nThe poetry book will align with our specialized phonics and reading program, and the Reciprocal Teaching Strategies book will help us get where we need to be.\r\n\r\nChart paper is a staple for any literacy or math lesson, and folders will help keep us organized. Ziplock pouches will attach to students' homework folders, making it simple and easy to transport school books home and back. \r\n\r\nPlease help us meet our needs with your support and generous donations. Thank you!nannan

Use Decontraction function to decontract project essay

In [13]:

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

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

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

In [14]:

```
sent = decontracted(project_data['essay'].values[1])
print(sent)
print("="*50)
```

Have you ever worked so hard on a project only to get it back from a teacher with a dismal grade? Or have a loved one that tries so very hard in school but just does not seem to grasp the concepts? That is how my students with special needs feel everyday! I create a classroom where everyone succeeds.\r\nMy students all have mild to moderate disabilities. The disabilities range from various levels of autism, moderate learning disabilities, challenges with attention, to being classified as intellectually impaired. \r\nThese students are just wonderful people but face daily challenges that you and I could never fathom. Most of my students come from low socioeconomic homes. Suffering from disabilities makes it difficult for them to read, comprehend, write, and solve math equations using typical learning styles. Technology is a way to bridge that learning gap these students struggle with each and every day.\r\nThese Chromebooks will be used in my classroom to help students complete their Common Core assignments in all subject areas. Students will be able to use this technology to help with the 21st century skills needed to be successful with the new Common Core State Standards and daily life.\r\nThis technology will make a huge impact on their lives. We currently have a teacher computer, document camera, projector, printer, and one chromebook per two students. The school itself has a few computer labs that it shares with the entire student body. These Chromebooks will be a huge benefit to my students' lives.nannan

=====

Remove line breaks

In [15]:

```
# \r\n\nt remove from string python: http://texthandler.com/info/remove-line-breaks-python/
```



```
# If you're removing from string python: http://texendriest.com/info/remove-the-breaks-python/
sent = sent.replace('\\r', ' ')
sent = sent.replace('\\n', ' ')
sent = sent.replace('\\t', ' ')
print(sent)
```

Have you ever worked so hard on a project only to get it back from a teacher with a dismal grade? Or have a loved one that tries so very hard in school but just does not seem to grasp the concepts? That is how my students with special needs feel everyday! I create a classroom where everyone succeeds. My students all have mild to moderate disabilities. The disabilities range from various levels of autism, moderate learning disabilities, challenges with attention, to being classified as intellectually impaired. These students are just wonderful people but face daily challenges that you and I could never fathom. Most of my students come from low socioeconomic homes. Suffering from disabilities makes it difficult for them to read, comprehend, write, and solve math equations using typical learning styles. Technology is a way to bridge that learning gap these students struggle with each and every day. These Chromebooks will be used in my classroom to help students complete their Common Core assignments in all subject areas. Students will be able to use this technology to help with the 21st century skills needed to be successful with the new Common Core State Standards and daily life. This technology will make a huge impact on their lives. We currently have a teacher computer, document camera, projector, printer, and one chromebook per two students. The school itself has a few computer labs that it shares with the entire student body. These Chromebooks will be a huge benefit to my students' lives.

Remove Special Chars

In [16]:

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

Have you ever worked so hard on a project only to get it back from a teacher with a dismal grade Or have a loved one that tries so very hard in school but just does not seem to grasp the concepts That is how my students with special needs feel everyday I create a classroom where everyone succeeds My students all have mild to moderate disabilities The disabilities range from various levels of autism moderate learning disabilities challenges with attention to being classified as intellectually impaired These students are just wonderful people but face daily challenges that you and I could never fathom Most of my students come from low socioeconomic homes Suffering from disabilities makes it difficult for them to read comprehend write and solve math equations using typical learning styles Technology is a way to bridge that learning gap these students struggle with each and every day These Chromebooks will be used in my classroom to help students complete their Common Core assignments in all subject areas Students will be able to use this technology to help with the 21st century skills needed to be successful with the new Common Core State Standards and daily life This technology will make a huge impact on their lives We currently have a teacher computer document camera projector printer and one chromebook per two students The school itself has a few computer labs that it shares with the entire student body These Chromebooks will be a huge benefit to my students lives

Remove Stopwords and Join the essays

In [17]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', \
            'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', \
            'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", \
            'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', \
            'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', \
            'while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', \
            'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', \
            'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', '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 [18]:

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

In [19]:

```
# after preprocessing
preprocessed_essays=Text_cleaner(project_data['essay'])
```

```
100% |████████████████████████████████████████████████████████████████████████████████| 50000/50000 [00:
51<00:00, 973.13it/s]
```

In [20]:

```
preprocessed_essays[1]
```

Out[20]:

'ever worked hard project get back teacher dismal grade loved one tries hard school not seem grasp concepts students special needs feel everyday create classroom everyone succeeds students mild moderate disabilities disabilities range various levels autism moderate learning disabilities challenges attention classified intellectually impaired students wonderful people face daily challenges could never fathom students come low socioeconomic homes suffering disabilities makes difficult read comprehend write solve math equations using typical learning styles technology way bridge learning gap students struggle every day chromebooks used classroom help students complete common core assignments subject areas students able use technology help 21st century skills needed successful new common core state standards daily life technology make huge impact lives currently teacher computer document camera projector printer one chromebook per two students school computer labs shares entire student body chromebooks huge benefit students lives nannan'

Drop essay columns 1, 2, 3, 4

In [21]:

```
project_data['essay'] = preprocessed_essays
project_data.drop(['project_essay_1'], axis=1, inplace=True)
project_data.drop(['project_essay_2'], axis=1, inplace=True)
project_data.drop(['project_essay_3'], axis=1, inplace=True)
project_data.drop(['project_essay_4'], axis=1, inplace=True)
project_data.head()
```

Out[21]:

Unnamed: 0				
------------	--	--	--	--

Add up the price based on project id

In [25]:

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

In [26]:

```
project_data.columns
```

Out[26]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay', 'Cleaned_title',
      'price', 'quantity'],
      dtype='object')
```

In [27]:

```
project_data.drop(['project_resource_summary'], axis=1, inplace=True)
project_data.drop(['Unnamed: 0'], axis=1, inplace=True)
project_data.drop(['id'], axis=1, inplace=True)
project_data.drop(['teacher_id'], axis=1, inplace=True)
```

Adding the word count of essay and title as new columns

In [28]:

```
project_data['essay_count']=project_data['essay'].str.len()
project_data['title_count']=project_data['Cleaned_title'].str.len()
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)
- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical
- Length of words in essay
- Length of words in title

Separate out the Dependant and independant variables

In [29]:

```
#https://stackoverflow.com/questions/29763620/how-to-select-all-columns-except-one-column-in-panda
s
X=project_data.loc[:, project_data.columns != 'project_is_approved']
y=project_data['project_is_approved']
X.shape
```

```
Out[29]:  
(50000, 13)
```

Splitting data into Test,Train,CV

```
In [30]:
```

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html  
  
from sklearn.model_selection import train_test_split  
  
X_train,X_test,y_train, y_test=train_test_split(X, y, test_size=0.3, random_state=0,stratify=y)  
X_train,X_cv,y_train,y_cv=train_test_split(X_train, y_train, test_size=0.01,  
random_state=0,stratify=y_train)  
  
print(X_train.shape)  
print(X_test.shape)  
print(X_cv.shape)  
print(y_train.shape)  
print(y_test.shape)  
print(y_cv.shape)
```

```
(34650, 13)  
(15000, 13)  
(350, 13)  
(34650,)   
(15000,)   
(350,)
```

```
In [31]:
```

```
X.head(2)
```

```
Out[31]:
```

	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	teacher_number_of_previously_posts
0	Ms.	NY	2017-02-05 17:49:01	3_5	7
1	Mrs.	CA	2016-05-27 14:44:25	6_8	3

Vectorize the features

Vectorize the Categorical Features - categories

```
In [32]:
```

```
# we use count vectorizer to convert the values into one  
from sklearn.feature_extraction.text import CountVectorizer
```

```

from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['clean_categories'].values)

# we use the fitted CountVectorizer to transform the text to vector
X_train_clean_categories=vectorizer.transform(X_train['clean_categories'].values)
X_test_clean_categories=vectorizer.transform(X_test['clean_categories'].values)
X_cv_clean_categories=vectorizer.transform(X_cv['clean_categories'].values)

print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",X_train_clean_categories.shape)

['AppliedLearning', 'Care_Hunger', 'Health_Sports', 'History_Civics', 'Literacy_Language',
'Math_Science', 'Music_Arts', 'SpecialNeeds', 'Warmth']
Shape of matrix after one hot encodig  (34650, 9)

```

Vectorize the Categorical Features - subcategories

In [33]:

```

vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['clean_subcategories'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_clean_sub_categories=vectorizer.transform(X_train['clean_subcategories'].values)
X_test_clean_sub_categories=vectorizer.transform(X_test['clean_subcategories'].values)
X_cv_clean_sub_categories=vectorizer.transform(X_cv['clean_subcategories'].values)

print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",X_train_clean_sub_categories.shape)

['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government',
'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics',
'EnvironmentalScience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness',
'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writing', 'M
athematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'Socia
lSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix after one hot encodig  (34650, 30)

```

Vectorize the Categorical Features - school state

In [34]:

```

vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['school_state'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_skl_state=vectorizer.transform(X_train['school_state'].values)
X_test_skl_state=vectorizer.transform(X_test['school_state'].values)
X_cv_skl_state=vectorizer.transform(X_cv['school_state'].values)

print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",X_train_skl_state.shape)

['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'K
S', 'KY', 'LA', 'MA', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH', 'NJ', 'NM',
'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC', 'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV
', 'WY']
Shape of matrix after one hot encodig  (34650, 51)

```

Vectorize the Categorical Features - teacher prefix

In [35]:

```

vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['teacher_prefix'].values)

```

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_prefix=vectorizer.transform(X_train['teacher_prefix'].values)
X_test_teacher_prefix=vectorizer.transform(X_test['teacher_prefix'].values)
X_cv_teacher_prefix=vectorizer.transform(X_cv['teacher_prefix'].values)

print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",X_train_teacher_prefix.shape)
```

```
['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
Shape of matrix after one hot encodig (34650, 5)
```

Vectorize the Categorical Features - project_grade_category

In [36]:

```
vectorizer = CountVectorizer(lowercase=False, binary=True)
vectorizer.fit(X_train['project_grade_category'].values)

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_grade_category=vectorizer.transform(X_train['project_grade_category'].values)
X_test_project_grade_category=vectorizer.transform(X_test['project_grade_category'].values)
X_cv_project_grade_category=vectorizer.transform(X_cv['project_grade_category'].values)

print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",X_train_project_grade_category.shape)
```

```
['3_5', '6_8', '9_12', 'PreK_2']
Shape of matrix after one hot encodig (34650, 4)
```

Vectorize the Numerical Features - price

In [37]:

```
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler

# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.
73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

price_scalar = StandardScaler()
price_scalar.fit(X_train['price'].values.reshape(-1,1)) # finding the mean and standard deviation
of this data
print(f"Mean : {price_scalar.mean_[0]}, Standard deviation : {np.sqrt(price_scalar.var_[0])}")

# Now standardize the data with above mean and variance.
X_train_price_standardized = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
X_test_price_standardized = price_scalar.transform(X_test['price'].values.reshape(-1, 1))
X_cv_price_standardized = price_scalar.transform(X_cv['price'].values.reshape(-1, 1))
```

```
Mean : 301.45576450216447, Standard deviation : 379.23697871265836
```

Vectorize the Numerical Features - quantity

In [38]:

```
quantity_scalar = StandardScaler()
quantity_scalar.fit(X_train['quantity'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation :
{np.sqrt(quantity_scalar.var_[0])}")
```

```
# Now standardize the data with above mean and variance.
X_train_quantity_standardized = quantity_scaler.transform(X_train['quantity'].values.reshape(-1, 1)
)
X_test_quantity_standardized = quantity_scaler.transform(X_test['quantity'].values.reshape(-1, 1))
X_cv_quantity_standardized = quantity_scaler.transform(X_cv['quantity'].values.reshape(-1, 1))
```

Mean : 16.944848484848485, Standard deviation : 26.452525257704806

Vectorize the Numerical Features - essay count

In [39]:

```
count_scaler = StandardScaler()
count_scaler.fit(X_train['essay_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {count_scaler.mean_[0]}, Standard deviation : {np.sqrt(count_scaler.var_[0])}")

# Now standardize the data with above mean and variance.
X_train_essay_count_standardized = count_scaler.transform(X_train['essay_count'].values.reshape(-1
, 1))
X_test_essay_count_standardized = count_scaler.transform(X_test['essay_count'].values.reshape(-1, 1
))
X_cv_essay_count_standardized = count_scaler.transform(X_cv['essay_count'].values.reshape(-1, 1))
```

Mean : 1015.1746897546898, Standard deviation : 277.8146315663551

Vectorize the Numerical Features - title count

In [40]:

```
count_scaler = StandardScaler()
count_scaler.fit(X_train['title_count'].values.reshape(-1,1)) # finding the mean and standard
deviation of this data
print(f"Mean : {count_scaler.mean_[0]}, Standard deviation : {np.sqrt(count_scaler.var_[0])}")

# Now standardize the data with above mean and variance.
X_train_title_count_standardized = count_scaler.transform(X_train['title_count'].values.reshape(-1
, 1))
X_test_title_count_standardized = count_scaler.transform(X_test['title_count'].values.reshape(-1, 1
))
X_cv_title_count_standardized = count_scaler.transform(X_cv['title_count'].values.reshape(-1, 1))
```

Mean : 25.69772005772006, Standard deviation : 11.646602313602617

Vectorizing Text data

Bag of words - essay

In [41]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['essay'])

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow=vectorizer.transform(X_train['essay'].values)
X_test_essay_bow=vectorizer.transform(X_test['essay'].values)
X_cv_essay_bow=vectorizer.transform(X_cv['essay'].values)

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

Shape of matrix after one hot encodig (34650, 5000)

Bag of words - cleaned title

In [42]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer = CountVectorizer(min_df=5)
vectorizer.fit(X_train['Cleaned_title'])

# we use the fitted CountVectorizer to convert the text to vector
X_train_cleaned_title_bow=vectorizer.transform(X_train['Cleaned_title'].values)
X_test_cleaned_title_bow=vectorizer.transform(X_test['Cleaned_title'].values)
X_cv_cleaned_title_bow=vectorizer.transform(X_cv['Cleaned_title'].values)

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

Shape of matrix after one hot encodig (34650, 2587)

TFIDF vectorizer - essay

In [43]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['essay'])

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf=vectorizer.transform(X_train['essay'].values)
X_test_essay_tfidf=vectorizer.transform(X_test['essay'].values)
X_cv_essay_tfidf=vectorizer.transform(X_cv['essay'].values)

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

Shape of matrix after one hot encodig (34650, 5000)

TFIDF vectorizer - cleaned title

In [44]:

```
# Similarly you can vectorize for title alsovectorizer = TfidfVectorizer(min_df=10)
vectorizer = TfidfVectorizer(min_df=5)

# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer.fit(X_train['Cleaned_title'])

# we use the fitted CountVectorizer to convert the text to vector
X_train_cleaned_title_tfidf=vectorizer.transform(X_train['Cleaned_title'].values)
X_test_cleaned_title_tfidf=vectorizer.transform(X_test['Cleaned_title'].values)
X_cv_cleaned_title_tfidf=vectorizer.transform(X_cv['Cleaned_title'].values)

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

Shape of matrix after one hot encodig (34650, 2587)

Using Pretrained Models: Avg W2V

In [45]:

```
'''
# 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):
```

```

for line in tqdm(1):
    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)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(", np.round(len(inter_words)/len(words)*100,3), "%) ")

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

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

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

'''

```

Out[45]:

```

'\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
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel('\glove.42B.300d.txt')\n\n# =====\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('\ '))\n\nprint("all the words in the
coupus", len(words))\n\nwords = set(words)\n\nprint("the unique words in the coupus",
len(words))\n\n\ninter_words = set(model.keys()).intersection(words)\n\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 = {}
\n\nwords_glove =
set(model.keys())\n\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\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\n\nimport pic
kle\n\nwith open('\glove_vectors', '\wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'

```

In [46]:

```

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

```

```
model = pickle.load(f)
glove_words = set(model.keys())
```

In [47]:

```
# average Word2Vec
# compute average word2vec for each review.
def avg_w2v_vectors(preprocessed_essays):
    avg_w2v_vectors_text = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(preprocessed_essays): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        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_text.append(vector)
    return avg_w2v_vectors_text

X_train_essay_w2v=avg_w2v_vectors(X_train['essay'])
X_test_essay_w2v=avg_w2v_vectors(X_test['essay'])
X_cv_essay_w2v=avg_w2v_vectors(X_cv['essay'])

X_train_cleaned_title_w2v=avg_w2v_vectors(X_train['Cleaned_title'])
X_test_cleaned_title_w2v=avg_w2v_vectors(X_test['Cleaned_title'])
X_cv_cleaned_title_w2v=avg_w2v_vectors(X_cv['Cleaned_title'])
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 34650/34650
[00:19<00:00, 1781.15it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 15000/15000
[00:08<00:00, 1812.74it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 350/350
[00:00<00:00, 1471.45it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 34650/34650
[00:00<00:00, 37562.38it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 15000/15000
[00:00<00:00, 42159.35it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 350/350
[00:00<00:00, 50042.83it/s]
```

Using Pretrained Models: TFIDF weighted W2V - Essay

In [48]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['essay'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words_essay = set(tfidf_model.get_feature_names())
print(len(tfidf_words_essay))
```

36002

In [49]:

```
# average Word2Vec
# compute average word2vec for each review.
def tfidf_w2v_vectors(tfidf_words,preprocessed_essays):
    tfidf_w2v_vectors_text = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(preprocessed_essays): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if (word in glove_words) and (word in tfidf_words):
                vec = model[word] # getting the vector for each word
                # here we are multiplying idf value(dictionary[word]) and the tf
                value((sentence.count(word)/len(sentence.split())))
                tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
                the tfidf value for each word
            vector += tf_idf*vec
```

```
X_train_essay_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_essay,X_train['essay'])
X_test_essay_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_essay,X_test['essay'])
X_cv_essay_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_essay,X_cv['essay'])
```

Using Pretrained Models: TFIDF weighted W2V - Cleaned Title

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['Cleaned_title'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words_Cleaned_title = set(tfidf_model.get_feature_names())
print(len(tfidf words Cleaned title))
```

```
X_train_cleaned_title_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_Cleaned_title,X_train['Cleaned_title'])
X_test_cleaned_title_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_Cleaned_title,X_test['Cleaned_title'])
X_cv_cleaned_title_tfidf_w2v=tfidf_w2v_vectors(tfidf_words_Cleaned_title,X_cv['Cleaned title'])
```

```
X_train_prev_proj=X_train['teacher_number_of_previously_posted_projects'][:,np.newaxis]
X_test_prev_proj=X_test['teacher_number_of_previously_posted_projects'][:,np.newaxis]
X_cv_prev_proj=X_cv['teacher number of previously posted projects'][:,np.newaxis]
```

BOW

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train_bow = hstack((X_train_clean_categories,
X_train_clean_sub_categories,X_train_skl_state,X_train_teacher_prefix,
```

```

        X_train_project_grade_category,X_train_price_standardized,X_train_quantity_standardized,
,X_train_prev_proj,
        X_train_essay_bow,X_train_cleaned_title_bow,X_train_essay_count_standardized,X_train_title_count_standardized
    )).toarray()

X_test_bow = hstack((X_test_clean_categories,
X_test_clean_sub_categories,X_test_skl_state,X_test_teacher_prefix,
        X_test_project_grade_category,X_test_price_standardized,X_test_quantity_standardized,X_test_prev_proj,
        X_test_essay_bow,X_test_cleaned_title_bow,X_test_essay_count_standardized,X_test_title_count_standardized
    )).toarray()

X_cv_bow = hstack((X_cv_clean_categories,
X_cv_clean_sub_categories,X_cv_skl_state,X_cv_teacher_prefix,
        X_cv_project_grade_category,X_cv_price_standardized,X_cv_quantity_standardized,X_cv_prev_proj,
X_cv_essay_bow,X_cv_cleaned_title_bow,X_cv_essay_count_standardized,X_cv_title_count_standardized
    )).toarray()

print(X_train_bow.shape)
print(X_test_bow.shape)
#print(X_cv_bow.shape)

```

```

(34650, 7691)
(15000, 7691)

```

TFIDF

In [54]:

```

# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train_tfidf = hstack((X_train_clean_categories, X_train_clean_sub_categories,X_train_skl_state,X_train_teacher_prefix,
        X_train_project_grade_category,X_train_price_standardized,X_train_quantity_standardized,
,X_train_prev_proj,
        X_train_essay_tfidf,X_train_cleaned_title_tfidf,X_train_essay_count_standardized,X_train_title_count_standardized
    )).toarray()

```

In [55]:

```

X_test_tfidf = hstack((X_test_clean_categories,
X_test_clean_sub_categories,X_test_skl_state,X_test_teacher_prefix,
        X_test_project_grade_category,X_test_price_standardized,X_test_quantity_standardized,X_test_prev_proj,
        X_test_essay_tfidf,X_test_cleaned_title_tfidf,X_test_essay_count_standardized,X_test_title_count_standardized
    )).toarray()

```

In [56]:

```

X_cv_tfidf = hstack((X_cv_clean_categories,
X_cv_clean_sub_categories,X_cv_skl_state,X_cv_teacher_prefix,
        X_cv_project_grade_category,X_cv_price_standardized,X_cv_quantity_standardized,X_cv_prev_proj,
X_cv_essay_tfidf,X_cv_cleaned_title_tfidf,X_cv_essay_count_standardized,X_cv_title_count_standardized
    )).toarray()

print(X_train_tfidf.shape)
print(X_test_tfidf.shape)
#print(X_cv_tfidf.shape)

```

```

(34650, 7691)

```

```
(15000, 7691)
```

Word2Vec

In [57]:

```
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train_w2v = hstack((X_train_clean_categories,
X_train_clean_sub_categories,X_train_skl_state,X_train_teacher_prefix,
                      X_train_project_grade_category,X_train_price_standardized,X_train_quantity_standardized,
,X_train_prev_proj,
                      X_train_essay_w2v,X_train_cleaned_title_w2v,X_train_essay_count_standardized,X_train_title_count_standardized
                      )).toarray()

X_test_w2v = hstack((X_test_clean_categories,
X_test_clean_sub_categories,X_test_skl_state,X_test_teacher_prefix,
                      X_test_project_grade_category,X_test_price_standardized,X_test_quantity_standardized,X_test_prev_proj,
                      X_test_essay_w2v,X_test_cleaned_title_w2v,X_test_essay_count_standardized,X_test_title_count_standardized
                      )).toarray()

X_cv_w2v = hstack((X_cv_clean_categories,
X_cv_clean_sub_categories,X_cv_skl_state,X_cv_teacher_prefix,
                   X_cv_project_grade_category,X_cv_price_standardized,X_cv_quantity_standardized,X_cv_prev_proj,
X_cv_essay_w2v,X_cv_cleaned_title_w2v,X_cv_essay_count_standardized,X_cv_title_count_standardized
                   )).toarray()

print(X_train_w2v.shape)
print(X_test_w2v.shape)
#print(X_cv_w2v.shape)
```

```
(34650, 704)
(15000, 704)
```

TFIDF- WORD2VEC

In [58]:

```
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train_tfidf_w2v = hstack((X_train_clean_categories,
X_train_clean_sub_categories,X_train_skl_state,X_train_teacher_prefix,
                      X_train_project_grade_category,X_train_price_standardized,X_train_quantity_standardized,
,X_train_prev_proj,
                      X_train_essay_tfidf_w2v,X_train_cleaned_title_tfidf_w2v,X_train_essay_count_standardized,X_train_title_count_standardized
                      )).toarray()

X_test_tfidf_w2v = hstack((X_test_clean_categories, X_test_clean_sub_categories,X_test_skl_state,X_test_teacher_prefix,
                      X_test_project_grade_category,X_test_price_standardized,X_test_quantity_standardized,X_test_prev_proj,
                      X_test_essay_tfidf_w2v,X_test_cleaned_title_tfidf_w2v,X_test_essay_count_standardized,X_test_title_count_standardized
                      )).toarray()

X_cv_tfidf_w2v = hstack((X_cv_clean_categories,
X_cv_clean_sub_categories,X_cv_skl_state,X_cv_teacher_prefix,
                   X_cv_project_grade_category,X_cv_price_standardized,X_cv_quantity_standardized,X_cv_prev_proj,
                   X_cv_essay_tfidf_w2v,X_cv_cleaned_title_tfidf_w2v,X_cv_essay_count_standardized,X_cv_title_count_standardized
                   )).toarray()
```

```
)).toarray())

print(X_train_tfidf_w2v.shape)
print(X_test_tfidf_w2v.shape)
#print(X_cv_tfidf_w2v.shape)
```

```
(34650, 704)
(15000, 704)
```

Apply SVM on BOW

Find best Hyper-Parameter value to train model

In [59]:

```
## By using "l2" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc',return_train_score=True)

    grid.fit(X_train_bow, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)

print("Best CV accuracy", grid.best_score_)
print("%"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc= grid.cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("***50)

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
plt.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()

plt.xscale('log')
plt.xlabel("log")
```

[illegible]

This plot shows the ROC_AUC score versus the log scale for both training and cross-validation data. The x-axis represents the log scale, ranging from 10^{-4} to 10^4 . The y-axis represents the ROC_AUC score, ranging from 0.60 to 0.80. The training AUC (blue line) starts at approximately 0.77, peaks at 10^{-2} (around 0.79), and then decreases to about 0.58 at 10^3 . The cross-validation AUC (orange line) starts at approximately 0.65, peaks at 10^{-1} (around 0.70), and then decreases to about 0.58 at 10^3 . Both curves converge to a similar value of approximately 0.58 for 10^3 and 10^4 .

log	Train AUC	CV AUC
10^{-4}	0.77	0.65
10^{-3}	0.76	0.64
10^{-2}	0.79	0.68
10^{-1}	0.75	0.70
10^0	0.69	0.67
10^1	0.66	0.65
10^2	0.60	0.59
10^3	0.58	0.58
10^4	0.58	0.58

```
## By using "l1" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc',return_train_score=True)

    grid.fit(X_train_bow, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)
```



```

train_auc = grid.cv_results_['mean_train_score']
cv_auc= grid.cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("'"*50)

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
plt.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()

plt.xscale('log')
plt.xlabel("log")

plt.ylabel("ROC_AUC score")

plt.title("ROC_AUC vs log plot")

plt.grid()

plt.show()

```

```

0%|
[00:00<?, ?it/s]
100%| 1/1 [08:
04<00:00, 484.26s/it]

```

```
{'alpha': 0.0001}
```

```

SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              epsilon=0.1, eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)

```

```

!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
[ 0.78302377 0.69882848 0.59193237 0.5624955 0.54573417 0.49793022
 0.5         0.5         0.5         ]

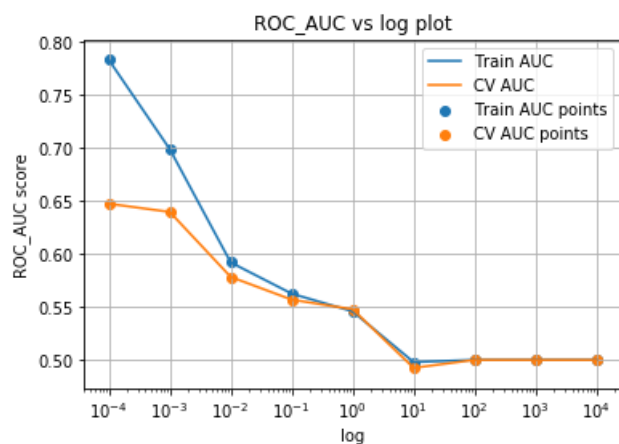
```

```

=====
[ 0.64729572 0.63954861 0.57827044 0.55680975 0.5477497 0.49244837
 0.5         0.5         0.5         ]

```

```
*****
```



Observation: l2 regularization gives better result than l1; best alpha is 0.1

Use best Hyper-Parameter value to train model

In [72]:

```
#
https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

Classifier = SGDClassifier(loss = 'hinge', penalty = 'l2', alpha = 0.1)

Classifier.fit(X_train_bow ,y_train)

#https://scikitlearn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function

y_train_pred = Classifier.decision_function(X_train_bow)

y_test_pred = Classifier.decision_function(X_test_bow)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)

test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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

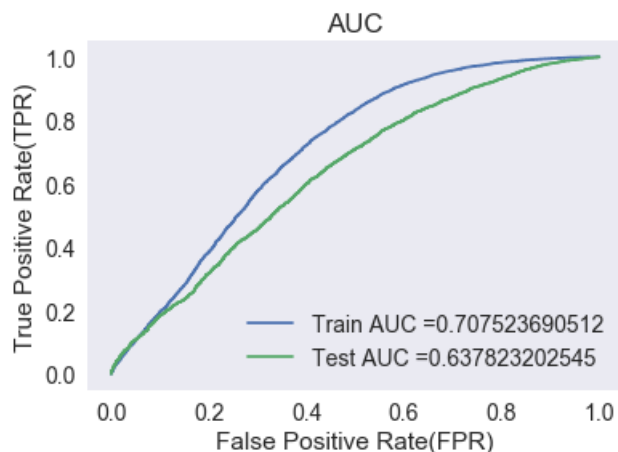
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))

plt.legend()

plt.ylabel("True Positive Rate(TPR) ")
plt.xlabel("False Positive Rate(FPR) ")

plt.title("AUC")

plt.grid()
plt.show()
```



In [64]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(fpr*(1-tpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

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

Plot Confusion Matrix for Test Data

In [73]:

```
import seaborn as sea

pred = predict(y_test_pred, te_thresholds, test_fpr, test_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_test, pred)
print(cm)
print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

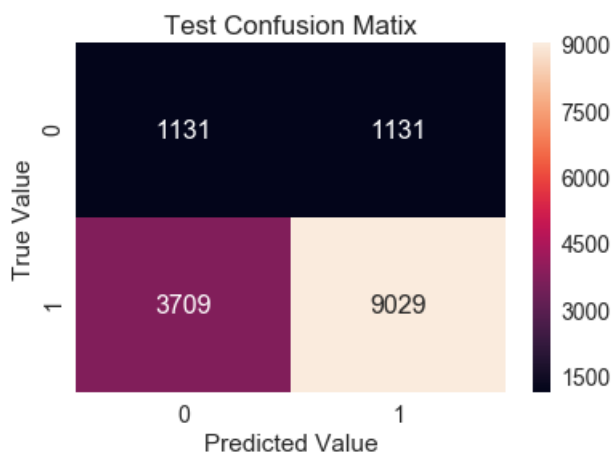
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 1.029

```
=====
[[1131 1131]
 [3709 9029]]
=====
```

Out[73]:

Text(0.5,1,'Test Confusion Matix')



Plot Confusion Matrix for Train Data

In [74]:

```
import seaborn as sea

pred = predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_train, pred)
print(cm)
# print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test confusion matrix, annot = True, annot_kws={"size":16}, fmt = 'd')
```

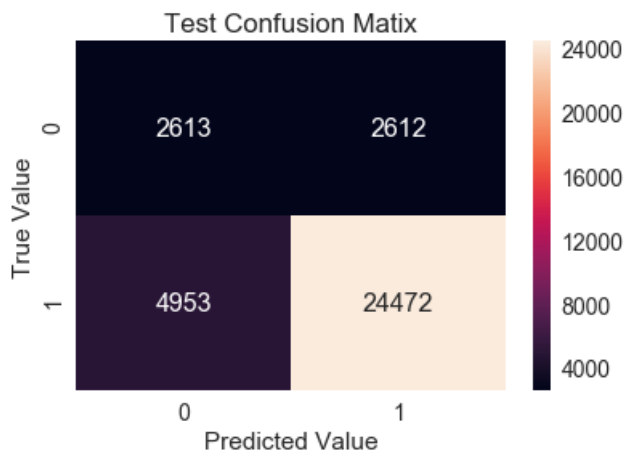
```
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999990843 for threshold 1.013

```
[[ 2613  2612]
 [ 4953 24472]]
```

Out[74]:

Text(0.5,1,'Test Confusion Matix')



Apply SVM on TFIDF

Find best Hyper-Parameter value to train model

In [67]:

```
## By using "l2" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm
import qtconsole

# hyperparameter tuning with l2 reg
parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_train_tfidf, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)

print("Best CV accuracy", grid.best_score_)
print("%"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc = grid.cv_results_['mean_test_score']
```

```

cv_auc= grid_cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("***50)

plt1.close()

plt1.plot(parameters['alpha'], train_auc, label='Train AUC')
plt1.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt1.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt1.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt1.legend()

plt1.xscale('log')
plt1.xlabel("log")

plt1.ylabel("ROC_AUC score")

plt1.title("ROC_AUC vs log plot")

plt1.grid()

plt1.show()

```

```

0%|
[00:00<?, ?it/s]
100%| 1/1 [04:
42<00:00, 282.39s/it]

```

```
{'alpha': 0.0001}
```

```

SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
              epsilon=0.1, eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)

```

```
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
```

```
Best CV accuracy 0.632997231629
```

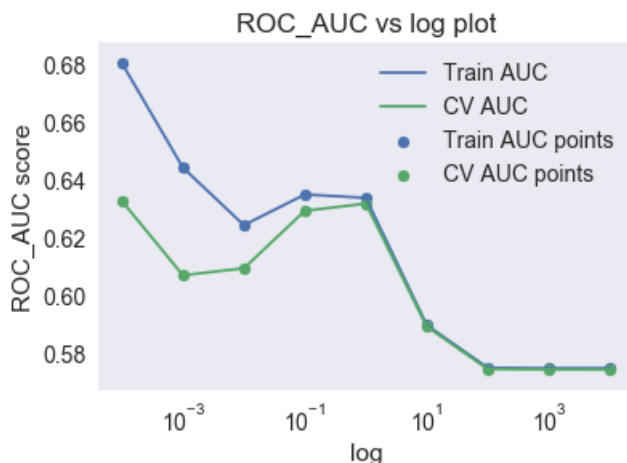
```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```
[ 0.68074857  0.64459919  0.62464874  0.63528817  0.63411627  0.5899937
 0.57524266  0.57512428  0.57512228]
```

```
=====
```

```
[ 0.63299723  0.60729722  0.60970164  0.62962507  0.63212504  0.58951378
 0.57478728  0.57466816  0.57466507]
```

```
*****
```



In [68]:

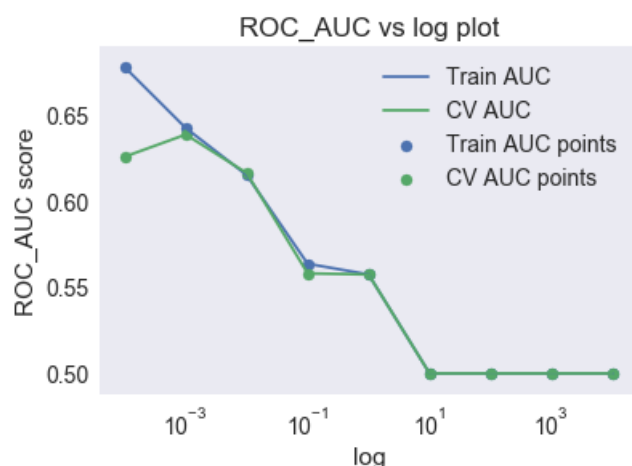
```
## Bringing "l1" Regularization
```



```

0.5      0.5      ]
=====
[ 0.62629821  0.63920401  0.61680067  0.55829979  0.55765554  0.5      0.5
0.5      0.5      ]
*****

```



Observation: l2 regularization gives better result than l1; best alpha is 0.0001

Use best Hyper-Parameter value to train model

In [75]:

```

#
https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

Classifier = SGDClassifier(loss = 'hinge', penalty = 'l2', alpha = 0.0001)

Classifier.fit(X_train_tfidf, y_train)

#https://scikitlearn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function

y_train_pred = Classifier.decision_function(X_train_tfidf)
y_test_pred = Classifier.decision_function(X_test_tfidf)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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

plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))

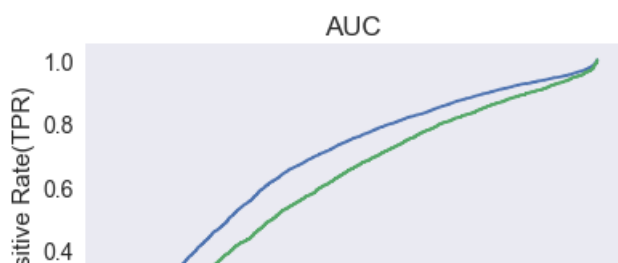
plt.legend()

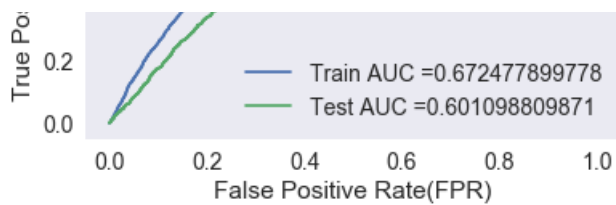
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")

plt.title("AUC")

plt.grid()
plt.show()

```





Plot Confusion Matrix for Test Data

In [76]:

```
import seaborn as sea

pred = predict(y_test_pred, te_thresholds, test_fpr, test_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_test, pred)
print(cm)
print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

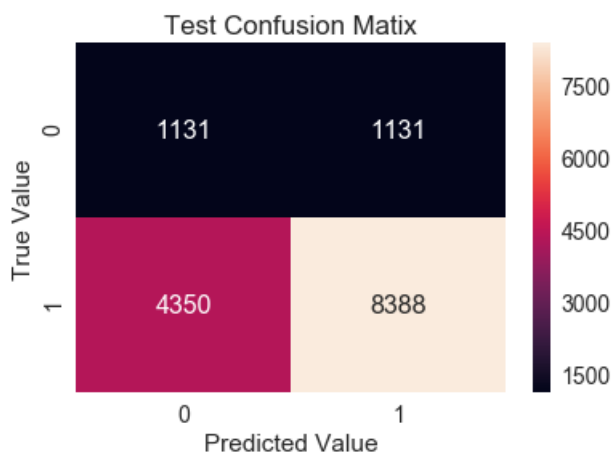
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 0.586

```
=====
[[1131 1131]
 [4350 8388]]
=====
```

Out[76]:

Text(0.5,1,'Test Confusion Matix')



Plot Confusion Matrix for Train Data

In [77]:

```
import seaborn as sea

pred = predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_train, pred)
```



```

print(cm)
# print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")

```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999990843 for threshold 0.201

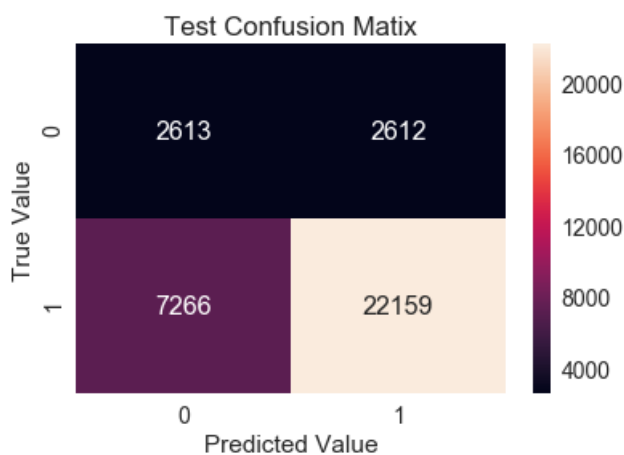
```

[[ 2613  2612]
 [ 7266 22159]]

```

Out[77]:

Text(0.5,1,'Test Confusion Matix')



Apply SVM on W2V

Find best Hyper-Parameter value to train model

In [78]:

```

## By using "l2" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_train_w2v, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)

```

```

print(grid.best_estimator_)
print("!"*50)

print("Best CV accuracy", grid.best_score_)
print("%"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc= grid.cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("***50)

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
plt.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()

plt.xscale('log')
plt.xlabel("log")

plt.ylabel("ROC_AUC score")

plt.title("ROC_AUC vs log plot")

plt.grid()

plt.show()

```

```

0%|
[00:00<?, ?it/s]
100%| 1/1 [00:30<00:00, 30.44s/it]

```

```
{'alpha': 0.001}
```

```

SGDClassifier(alpha=0.001, average=False, class_weight='balanced',
              epsilon=0.1, eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l2', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!

```

```
Best CV accuracy 0.658165199823
```

```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```

[ 0.65055143  0.68240046  0.67268392  0.66178741  0.62408732  0.59688571
  0.57465823  0.57791571  0.57792005]

```

```
=====
```

```

[ 0.63098966  0.6581652  0.65537588  0.65398121  0.62120382  0.59341604
  0.57437182  0.57716868  0.57717519]

```

```
*****
```

ROC_AUC vs log plot



log

```
## By using "l1" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv = 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_train_w2v, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc = grid.cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("*"*50)

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
plt.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()

plt.xscale('log')
plt.xlabel("log")

plt.ylabel("ROC_AUC score")

plt.title("ROC_AUC vs log plot")

plt.grid()

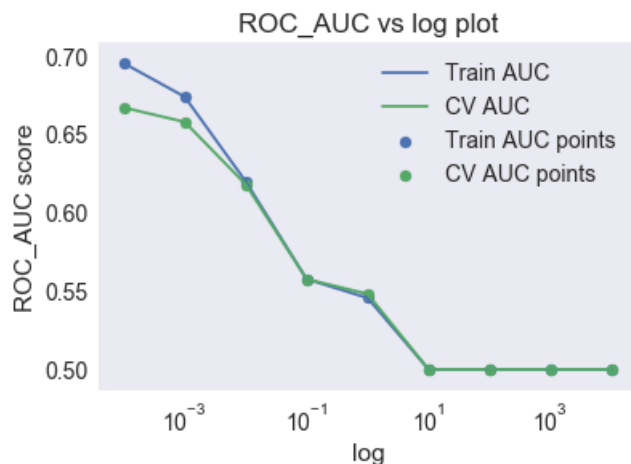
plt.show()
```

```
{'alpha': 0.0001}
-----
SGDClassifier(alpha=0.0001, average=False, class_weight='balanced',
```

```

epsilon=0.1, eta=0.0, fit_intercept=True, l1_ratio=0.15,
learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
shuffle=True, tol=None, verbose=0, warm_start=False)
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
[ 0.69547684  0.67406645  0.62004185  0.55767555  0.54561415  0.5          0.5
  0.5          0.5          ]
=====
[ 0.66716126  0.65800399  0.61815827  0.55765554  0.54820379  0.5          0.5
  0.5          0.5          ]
*****

```



Observation: l2 regularization gives better result than l1; best alpha is 0.001

Use best Hyper-Parameter value to train model

In [80]:

```

#
https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

Classifier = SGDClassifier(loss = 'hinge', penalty = 'l2', alpha = 0.001)

Classifier.fit(X_train_w2v ,y_train)

#https://scikitlearn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function

y_train_pred = Classifier.decision_function(X_train_w2v)

y_test_pred = Classifier.decision_function(X_test_w2v)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)

test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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

plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))

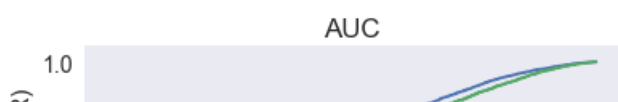
plt.legend()

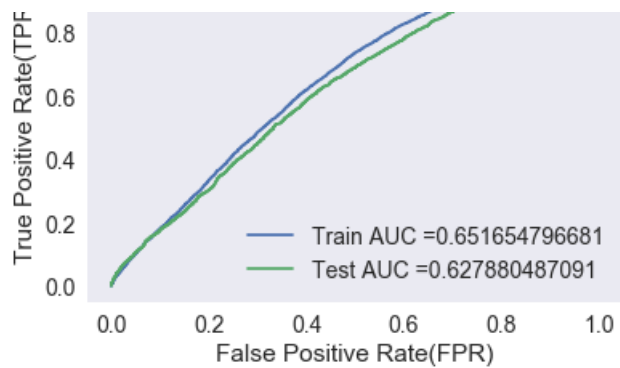
plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")

plt.title("AUC")

plt.grid()
plt.show()

```





Plot Confusion Matrix for Test Data

In [81]:

```
import seaborn as sea

pred = predict(y_test_pred, te_thresholds, test_fpr, test_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_test, pred)
print(cm)
print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

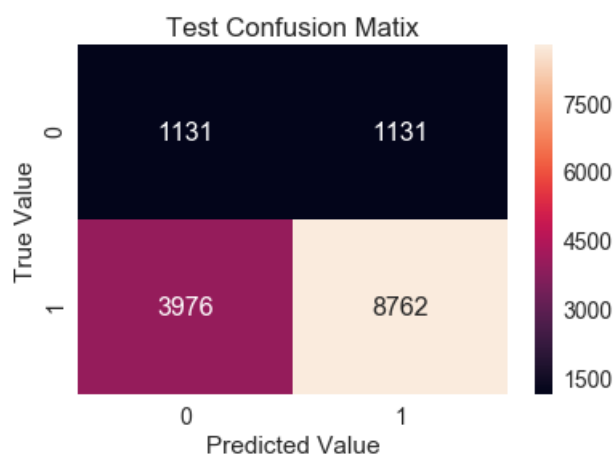
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 1.18

```
=====
[[1131 1131]
 [3976 8762]]
=====
```

Out[81]:

Text(0.5,1,'Test Confusion Matix')



Plot Confusion Matrix for Train Data

In [82]:

```
import seaborn as sea
```

```

pred = predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_train, pred)
print(cm)
# print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")

```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999990843 for threshold 1.086

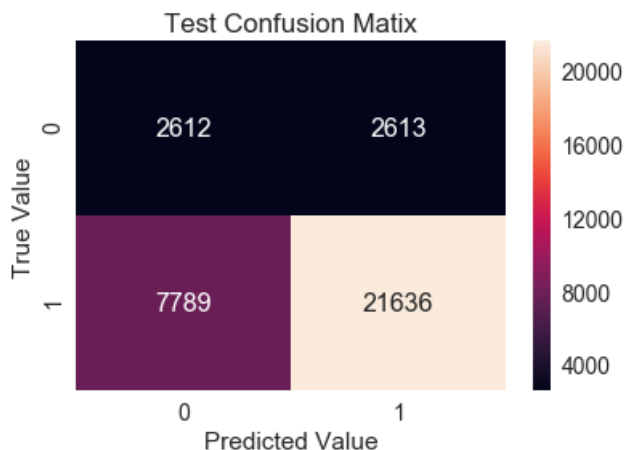
```

=====
[[ 2612  2613]
 [ 7789 21636]]

```

Out[82]:

Text(0.5,1,'Test Confusion Matix')



Apply SVM on TFIDFW2V

Find best Hyper-Parameter value to train model

In [83]:

```

## By using "l2" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_train_tfidf_w2v, y_train)

# https://www.geeksforgeeks.org

```

```
0%|
[00:00<?, ?it/s]
100%| 1/1 [00
:40<00:00, 40.37s/it]
```

ROC_AUC vs log plot

AUC score

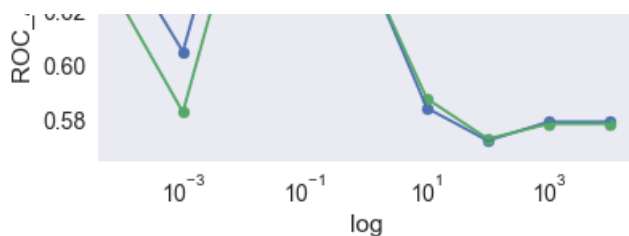
Train AUC

CV AUC

Train AUC points

CV AUC points

log scale	Train AUC	CV AUC
0	0.642	0.622
1	0.618	0.618
2	0.685	0.665
3	0.675	0.665
4	0.638	0.638
5	0.618	0.618



In [84]:

```
## By using "l1" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha': [10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv= 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_train_tfidf_w2v, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc= grid.cv_results_['mean_test_score']

# print train_auc
print(train_auc)

print("="*50)

# print cv_auc
print(cv_auc)

print("*"*50)

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
plt.plot(parameters['alpha'], cv_auc, label='CV AUC')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()

plt.xscale('log')
plt.xlabel("log")

plt.ylabel("ROC_AUC score")

plt.title("ROC_AUC vs log plot")

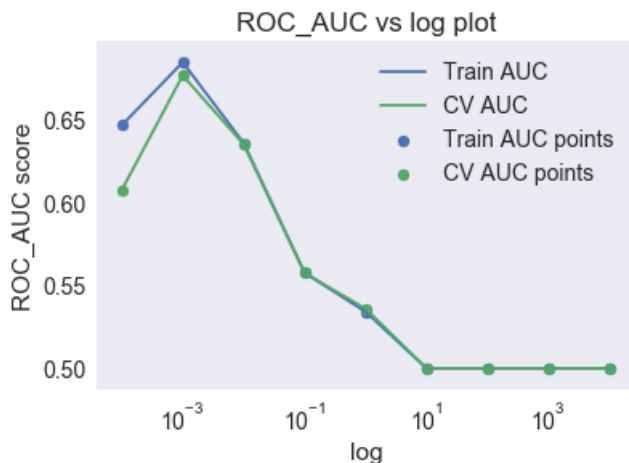
plt.grid()

plt.show()
```


:09<00:00, 69.68s/it]

```
{'alpha': 0.001}
```

```
-----
SGDClassifier(alpha=0.001, average=False, class_weight='balanced',
              epsilon=0.1, eta0=0.0, fit_intercept=True, l1_ratio=0.15,
              learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
              n_jobs=1, penalty='l1', power_t=0.5, random_state=None,
              shuffle=True, tol=None, verbose=0, warm_start=False)
!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!
[ 0.6468363  0.68493974 0.63609844 0.55767555 0.53425517 0.5          0.5
  0.5          0.5          ]
=====
[ 0.60768643 0.6769831 0.63536777 0.55765554 0.53600556 0.5          0.5
  0.5          0.5          ]
*****
```



Observation: l2 regularization gives better result than l1; best alpha is 0.001

Use best Hyper-Parameter value to train model

In [100]:

```
#
https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

Classifier = SGDClassifier(loss = 'hinge', penalty = 'l2', alpha = 0.001)

Classifier.fit(X_train_tfidf_w2v ,y_train)

#https://scikitlearn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function

y_train_pred = Classifier.decision_function(X_train_tfidf_w2v)

y_test_pred = Classifier.decision_function(X_test_tfidf_w2v)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)

test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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

plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))

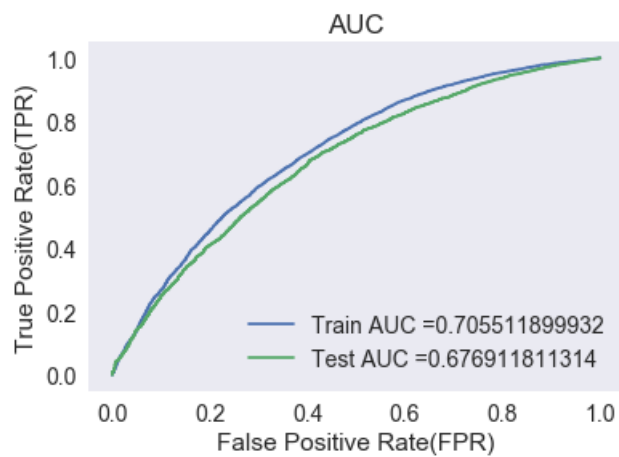
plt.legend()

plt.ylabel("True Positive Rate(TPR)")
plt.xlabel("False Positive Rate(FPR)")

plt.title("AUC")

plt.grid()
```

```
plt.show()
```



Plot Confusion Matrix for Test Data

In [86]:

```
import seaborn as sea

pred = predict(y_test_pred, te_thresholds, test_fpr, test_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_test, pred)
print(cm)
print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

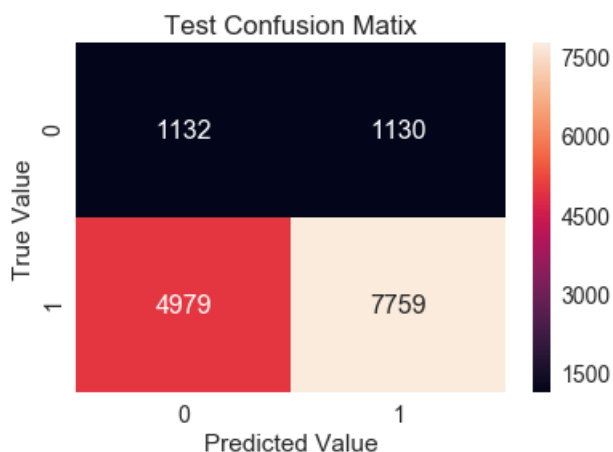
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999804559 for threshold 0.993

```
=====
[[1132 1130]
 [4979 7759]]
=====
```

Out[86]:

Text(0.5,1,'Test Confusion Matix')



Plot Confusion Matrix for Train Data

In [87]:

```
import seaborn as sea

pred = predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_train, pred)
print(cm)
# print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

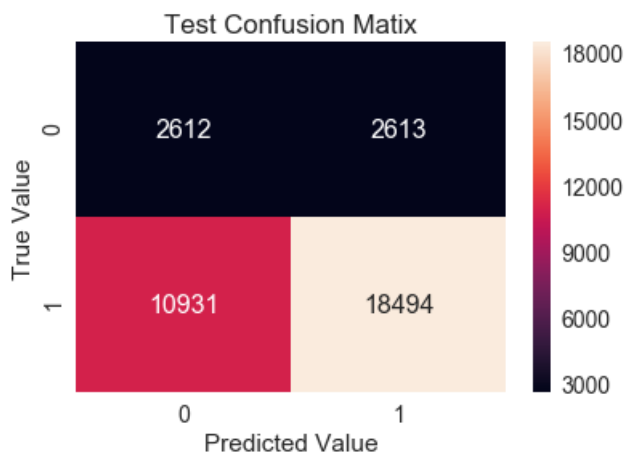
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999990843 for threshold 0.992

```
=====
[[ 2612  2613]
 [10931 18494]]
```

Out[87]:

Text(0.5,1,'Test Confusion Matix')



Apply TruncatedSVD on TfidfVectorizer of essay text, choose the number of components (n_components) using elbow method :numerical data

Find the best dimension

In [88]:

```
from sklearn.decomposition import TruncatedSVD
#https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html
#declaring index as Dimensions in train_text_tfidf
Dim = [25,50,100,200,500,1500,2000,2500,3000,3500]
Variance_sum = []
for i in tqdm(Dim):
    svd = TruncatedSVD(n_components = i, random_state = 42)
    svd.fit(X_train_essay_tfidf)
    Variance_sum.append(svd.explained_variance_ratio_.sum())
```

```
0%|
[00:00<?, ?it/s]
10%|
```

1/10 100%

```

00:45, 5.09s/it] | 2/10
20% | 3/10 [00:
[00:10<00:42, 5.27s/it]
30% | 4/10
<00:48, 6.91s/it]
40% | 5/10 [01:
[00:40<01:03, 10.58s/it]
50% | 6/10 [04:
8<01:49, 21.83s/it]
60% | 7/10
6<04:34, 68.69s/it]
70% | 8/10 [14:0
[08:38<06:10, 123.57s/it]
80% | 9/10 [21:2
6<06:10, 185.02s/it]
90% | 10/10
0<04:19, 259.54s/it]
100% |
[30:36<00:00, 183.65s/it]

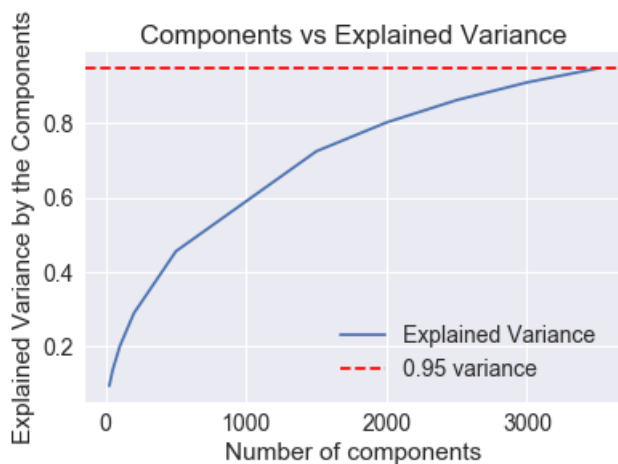
```

In [89]:

```

plt.plot(Dim, Variance_sum, label = 'Explained Variance ')
plt.axhline(0.95, linestyle = '--', color = 'r', label = '0.95 variance')
plt.xlabel("Number of components")
plt.ylabel("Explained Variance by the Components")
plt.legend()
plt.title("Components vs Explained Variance")
sns.despine()
plt.show()

```



OBSERVATION: At 3500 dimensions we have Accuracy of 95%

Use best dimension value to train TruncatedSVD

In [90]:

```

svd = TruncatedSVD(n_components= 3500)

svd.fit(X_train_essay_tfidf)
#Transforms:
#Train SVD
X_tr_essay_tfidf_trunSVD= svd.transform(X_train_essay_tfidf )
#Test SVD
X_te_essay_tfidf_trunSVD = svd.transform(X_test_essay_tfidf )
#CV SVD
#X_cval_essay_tfidf_trunSVD = svd.transform(X_cv_essay_tfidf )

```

Merge data

In [91]:

```
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set5_train = hstack((X_train_clean_categories, X_train_clean_sub_categories, X_train_skl_state, X_train_teacher_prefix,
                        X_train_project_grade_category, X_train_price_standardized, X_train_quantity_standardized,
                        X_train_prev_proj,
                        X_train_essay_tfidf_trunSVD, X_train_cleaned_title_tfidf, X_train_essay_count_standardized, X_train_title_count_standardized
                        )).toarray()
```

In [92]:

```
X_set5_test = hstack((X_test_clean_categories,
X_test_clean_sub_categories, X_test_skl_state, X_test_teacher_prefix,
                        X_test_project_grade_category, X_test_price_standardized, X_test_quantity_standardized, X_test_prev_proj,
                        X_test_essay_tfidf_trunSVD, X_test_cleaned_title_tfidf, X_test_essay_count_standardized, X_test_title_count_standardized
                        )).toarray()
```

In [93]:

```
print(X_set5_train.shape)
print(X_set5_test.shape)
```

```
(34650, 6191)
(15000, 6191)
```

Apply SVM on Set 5

Find best Hyper-parameter to train the model

In [95]:

```
## By using "l2" Regularizer
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
from sklearn import linear_model
from sklearn.linear_model import SGDClassifier
from sklearn import svm

# hyperparameter tuning with l2 reg
parameters = {'alpha': [10**-5, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4, 10**5]}

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

for alpha in tqdm(parameters):
    grid = GridSearchCV(sd, parameters, cv = 5, scoring='roc_auc', return_train_score=True)

    grid.fit(X_set5_train, y_train)

# https://www.geeksforgeeks.org

# print best hyper parameter after tuning
print(grid.best_params_)
print("-"*50)

# print how our model looks after hyper-parameter tuning
print(grid.best_estimator_)
print("!"*50)

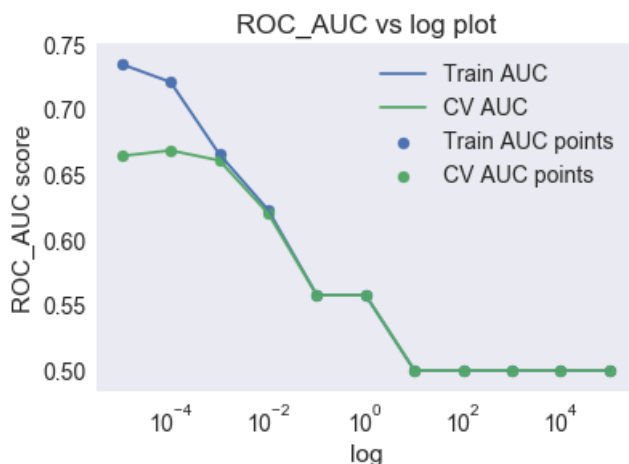
print("Best CV accuracy", grid.best_score_)
print("%"*50)

train_auc = grid.cv_results_['mean_train_score']
cv_auc = grid.cv_results_['mean_test_score']

# print train auc
print(train_auc)
```



```
[ 0.73483322  0.72119238  0.66581167  0.62316093  0.55765554  0.55765554
  0.5          0.5          0.5          0.5          0.5          ]
=====
[ 0.66451013  0.66871438  0.66105744  0.62053252  0.55765554  0.55765554
  0.5          0.5          0.5          0.5          0.5          ]
*****
```



Observation: l2 regularizatiOn works better than l1 with best alpha as 0.00001

Use best Hper-parameter to train the model

In [97]:

```
#
https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

Classifier = SGDClassifier(loss = 'hinge', penalty = 'l2', alpha =0.00001)

Classifier.fit(X_set5_train ,y_train)

#https://scikitlearn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function

y_train_pred = Classifier.decision_function(X_set5_train)

y_test_pred = Classifier.decision_function(X_set5_test)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)

test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

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

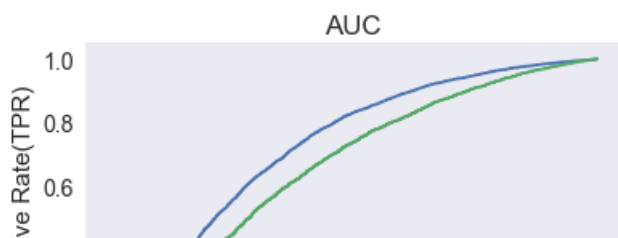
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))

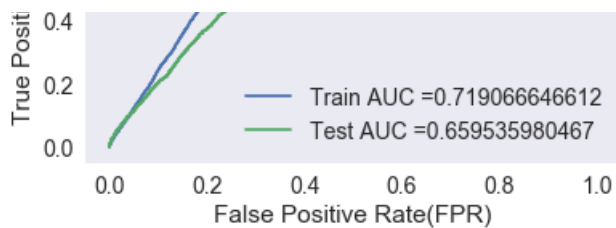
plt.legend()

plt.ylabel("True Positive Rate(TPR) ")
plt.xlabel("False Positive Rate(FPR) ")

plt.title("AUC")

plt.grid()
plt.show()
```





Plot Confusion Matrix for Test Data

In [98]:

```
import seaborn as sea

pred = predict(y_test_pred, te_thresholds, test_fpr, test_fpr)

#print(pred)
print("="*50)

cm = confusion_matrix(y_test, pred)
print(cm)
print("="*50)

test_confusion_matrix = pd.DataFrame(cm, range(2), range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

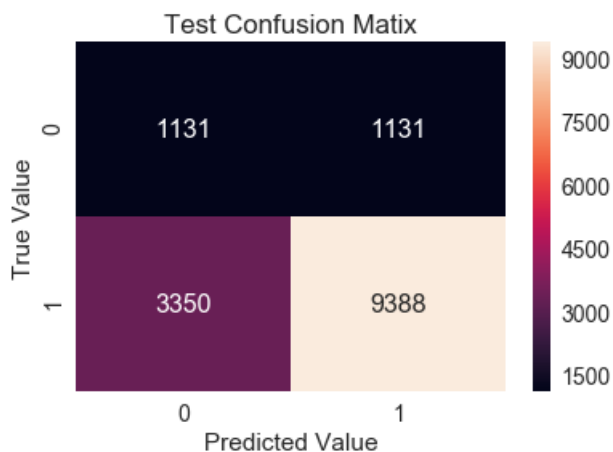
plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 5.137

```
=====
[[1131 1131]
 [3350 9388]]
=====
```

Out[98]:

Text(0.5,1,'Test Confusion Matix')



Plot Confusion Matrix for Train Data

In [99]:

```
import seaborn as sea

pred = predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)

#print(pred)
print("="*50)
```

```

cm = confusion_matrix(y_train,pred)
print(cm)
# print("="*50)

test_confusion_matrix = pd.DataFrame(cm,range(2),range(2))

sea.set(font_scale=1.4)

sea.heatmap(test_confusion_matrix, annot = True, annot_kws={"size":16}, fmt = 'd')

plt.xlabel("Predicted Value")
plt.ylabel("True Value")
plt.title("Test Confusion Matix")

```

the maximum value of $tpr \cdot (1 - fpr)$ 0.249999990843 for threshold 1.398

=====

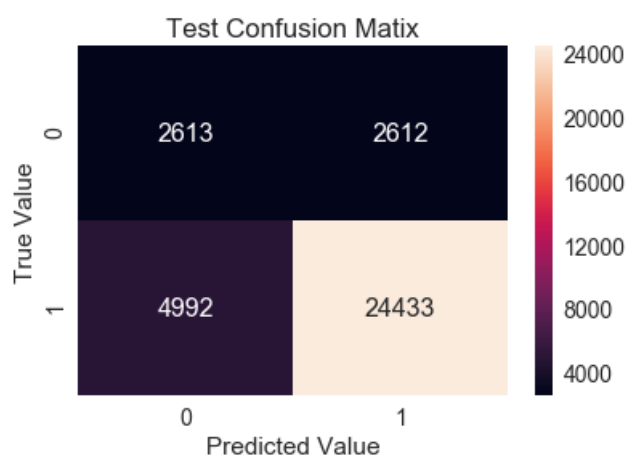
```

[[ 2613  2612]
 [ 4992 24433]]

```

Out[99]:

Text(0.5,1,'Test Confusion Matix')



Conclusion

In []:

```

# Please compare all your models using Prettytable library
#how to use pretty table http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
tb = PrettyTable()
tb.field_names= ("Vectorizer", " Alpha ", " AUC ")
tb.add_row(["BOW ", 0.1, 63])
tb.add_row(["Tf - Idf ", 0.0001, 60])
tb.add_row(["AVG - W2V", 0.001, 62])
tb.add_row(["TfIdf - W2V", 0.001, 67])
tb.add_row(["SVD-Top 3500 Features", 0.00001, 65])
print(tb.get_string(titles = "SVM- Observations"))

```

	Vectorizer	Alpha	AUC
	BOW	0.1	63
	Tf - Idf	0.0001	60
	AVG - W2V	0.001	62
	TfIdf - W2V	0.001	67
	SVD-Top 3500 Features	1e-05	65