# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

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

| Feature | Description |
|---|---|
| project_essay_4 | Fourth application essay |
| project_submitted_datetime | Datetime when project application was submitted. **Example:** `2016-04-28 12:43:56.245` |
| teacher_id | A unique identifier for the teacher of the proposed project. **Example:** `bdf8baa8fedef6bfeec7ae4ff1c15c56` |
| teacher_prefix | Teacher's title. One of the following enumerated values: <br>• `nan`<br>• `Dr.`<br>• `Mr.`<br>• `Mrs.`<br>• `Ms.`<br>• `Teacher.` |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. **Example:** `2` |

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

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

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

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

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

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

## Notes on the Essay Data

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

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

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

In [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
from sklearn.model_selection import train_test_split
import sklearn.model_selection as model_selection
```

## 1.1 Reading Data

In [2]:

```
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

In [3]:

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

```
Number of data points in train data (109248, 17)
--------------------------------------------------
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

```
Number of data points in train data (1541272, 4)
['id' 'description' 'quantity' 'price']
```

Out[4]:

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

## 1.2 preprocessing of `project_subject_categories`

```python
catogories = 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 catogories:
    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 catogory 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 placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
        temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

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

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

## 1.3 preprocessing of `project_subject_subcategories`

```python
sub_catogories = 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_catogories:
    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 catogory 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 placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
        temp +=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

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]))
```

## 1.3 Text preprocessing

In [7]:

```
# 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 [8]:

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

Out[8]:

| | Unnamed: 0 | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_cate |
|---|---|---|---|---|---|---|---|
| 0 | 160221 | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | Mrs. | IN | 2016-12-05 13:43:57 | Grades P |
| 1 | 140945 | p258326 | 897464ce9ddc600bced1151f324dd63a | Mr. | FL | 2016-10-25 09:22:10 | Grade |

In [9]:

```
#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

In [10]:

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery.  We also have over 40 countries represented with the families within our school.  Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein  Our English learner's have a strong support system at home that begs for more resources.  Many times our parents are learning to read and speak English along side of their children.  Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist.  All families with students within the Level 1 proficiency status, will be a offered to be a part of this program.  These educational videos will be specially chosen by the En

offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills.\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\nnannan

==================================================

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity.My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.nannan

==================================================

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more.With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

==================================================

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say.Wobble chairs are the answer and I love then because they develop their core, which enhances gross motor and in Turn fine motor skills. \r\nThey also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

==================================================

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character.In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to

ow me to have more room for storage of things that are needed for the day and has an extra part to
it I can use.  The table top chart has all of the letter, words and pictures for students to learn
about different letters and it is more accessible.nannan
==================================================


In [11]:

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

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

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

In [12]:

```python
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cogniti
ve delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work th
eir hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out
for my students. I teach in a Title I school where most of the students receive free or reduced pr
ice lunch.  Despite their disabilities and limitations, my students love coming to school and come
eager to learn and explore.Have you ever felt like you had ants in your pants and you needed to gr
oove and move as you were in a meeting? This is how my kids feel all the time. The want to be able
to move as they learn or so they say.Wobble chairs are the answer and I love then because they dev
elop their core, which enhances gross motor and in Turn fine motor skills. \r\nThey also want to l
earn through games, my kids do not want to sit and do worksheets. They want to learn to count by j
umping and playing. Physical engagement is the key to our success. The number toss and color and s
hape mats can make that happen. My students will forget they are doing work and just have the fun
a 6 year old deserves.nannan
==================================================


In [13]:

```python
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\\r', ' ')
sent = sent.replace('\\"', ' ')
sent = sent.replace('\\n', ' ')
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cogniti
ve delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work th
eir hardest working past their limitations.     The materials we have are the ones I seek out for
my students. I teach in a Title I school where most of the students receive free or reduced price
lunch.  Despite their disabilities and limitations, my students love coming to school and come eag
er to learn and explore.Have you ever felt like you had ants in your pants and you needed to groov
e and move as you were in a meeting? This is how my kids feel all the time. The want to be able to
move as they learn or so they say.Wobble chairs are the answer and I love then because they develo
p their core, which enhances gross motor and in Turn fine motor skills.   They also want to learn t
hrough games, my kids do not want to sit and do worksheets. They want to learn to count by jumping
and playing. Physical engagement is the key to our success. The number toss and color and shape ma
ts can make that happen. My students will forget they are doing work and just have the fun a 6 yea
r old deserves.nannan

In [14]:

```python
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
```

```python
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays cognitiv
e delays gross fine motor delays to autism They are eager beavers and always strive to work their
hardest working past their limitations The materials we have are the ones I seek out for my studen
ts I teach in a Title I school where most of the students receive free or reduced price lunch
Despite their disabilities and limitations my students love coming to school and come eager to lea
rn and explore Have you ever felt like you had ants in your pants and you needed to groove and mov
e as you were in a meeting This is how my kids feel all the time The want to be able to move as th
ey learn or so they say Wobble chairs are the answer and I love then because they develop their co
re which enhances gross motor and in Turn fine motor skills They also want to learn through games
my kids do not want to sit and do worksheets They want to learn to count by jumping and playing Ph
ysical engagement is the key to our success The number toss and color and shape mats can make that
happen My students will forget they are doing work and just have the fun a 6 year old deserves nan
nan

In [15]:

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

In [16]:

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

```
100%|████████████████████████████████████████████████████████████| 109248/109248
[01:14<00:00, 1464.62it/s]
```

In [17]:

```python
# after preprocesing
project_data['processed_essay'] = preprocessed_essays;
```

```
project_data.drop(['essay'], axis=1, inplace=True)
preprocessed_essays[20000]
```

Out[17]:

'my kindergarten students varied disabilities ranging speech language delays cognitive delays gros
s fine motor delays autism they eager beavers always strive work hardest working past limitations
the materials ones i seek students i teach title i school students receive free reduced price lunc
h despite disabilities limitations students love coming school come eager learn explore have ever
felt like ants pants needed groove move meeting this kids feel time the want able move learn say w
obble chairs answer i love develop core enhances gross motor turn fine motor skills they also want
learn games kids not want sit worksheets they want learn count jumping playing physical engagement
key success the number toss color shape mats make happen my students forget work fun 6 year old de
serves nannan'

## 1.4 Preprocessing of `project_title`

In [18]:

```python
# similarly you can preprocess the titles also

processed_titles = [];
for title in tqdm(project_data['project_title'].values):
    sent = decontracted(title)
    sent = re.sub('\S*\d\S*', '', sent);
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    processed_titles.append(sent.strip())
```

```
100%|████████████████████████████████████████████████████████| 109248/109248
[00:02<00:00, 37603.21it/s]
```

In [19]:

```python
project_data.drop(['project_title'], axis=1, inplace=True)
project_data['processed_titles'] = processed_titles

#testing after preprocessing project_title column
print(processed_titles[3])

print(processed_titles[40]);

print(processed_titles[500]);

print(processed_titles[4000]);

project_data.columns
```

```
Techie Kindergarteners
Leveling Books in a Multi Age Class
Classroom Chromebooks for College Bound Seniors
Inspire Summer Reading
```

Out[19]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
       'project_submitted_datetime', 'project_grade_category',
       'project_essay_1', 'project_essay_2', 'project_essay_3',
       'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approved',
       'clean_categories', 'clean_subcategories', 'processed_essay',
       'processed_titles'],
      dtype='object')
```

## Preprocessing of project_grade_category

In [20]:

```python
print(project_data['project_grade_category'][1])
print(project_data['project_grade_category'][223])
```

```
print(project_data['project_grade_category'][134])
```

```
Grades 6-8
Grades PreK-2
Grades PreK-2
```

```
processed_grades = [];
for grades in project_data['project_grade_category']:
    grades = grades.replace('-', '');
    processed_grades.append(grades)
```

```
print(processed_grades[1])
print(processed_grades[223])
print(processed_grades[134])

project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data['processed_grades'] = processed_grades
```

```
Grades 68
Grades PreK2
Grades PreK2
```

## Preprocessing of teacher_prefix

```
print(project_data['teacher_prefix'][2]);
print(project_data['teacher_prefix'][234]);
print(project_data['teacher_prefix'][425]);
```

```
Ms.
Ms.
Ms.
```

```
preprocessed_teacher_prefix = [];
for prefix in project_data['teacher_prefix']:
    prefix = str(prefix).replace('.', '');
    preprocessed_teacher_prefix.append(prefix);
```

```
project_data.drop(['teacher_prefix'], axis=1, inplace=True)
project_data['processed_teacher_prefix'] = preprocessed_teacher_prefix

print(preprocessed_teacher_prefix[321])
print(preprocessed_teacher_prefix[310])
```

```
Mrs
Ms
```

## 1.5 Preparing data for models

```
project_data.columns
```

```
Index(['Unnamed: 0', 'id', 'teacher id', 'school state',
```

```
       'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
       'project_essay_3', 'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approved',
       'clean_categories', 'clean_subcategories', 'processed_essay',
       'processed_titles', 'processed_grades', 'processed_teacher_prefix'],
      dtype='object')
```

we are going to consider

```
      - school_state : categorical data
      - clean_categories : categorical data
      - clean_subcategories : categorical data
      - project_grade_category : categorical data
      - teacher_prefix : categorical data

      - project_title : text data
      - text : text data
      - project_resource_summary: text data (optinal)

      - quantity : numerical (optinal)
      - teacher_number_of_previously_posted_projects : numerical
      - price : numerical
```

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

In [27]:

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

Out[27]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
       'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
       'project_essay_3', 'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approved',
       'clean_categories', 'clean_subcategories', 'processed_essay',
       'processed_titles', 'processed_grades', 'processed_teacher_prefix',
       'price', 'quantity'],
      dtype='object')
```

In [28]:

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

**Computing Sentiment Scores**

In [29]:

```python
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# import nltk
# nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students w
ith the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelli
gences i use a wide range\
```

```
of techniques to help all my students succeed students in my class come from a variety of differen
t backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a carin
g community of successful \
learners which can be seen through collaborative student project based learning in and out of the
classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice
a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the ki
ndergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role pla
y in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food
i will take their idea \
and create common core cooking lessons where we learn important math and writing concepts while co
oking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into maki
ng the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project woul
d expand our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make homemade apple
sauce make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create our own cook
books to be printed and \
shared with families students will gain math and literature skills as well as a life long enjoymen
t for healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

for k in ss:
    print('{0}: {1}, '.format(k, ss[k]), end='')

# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
```

neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,

# Assignment 5: Logistic Regression

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

   - Set 1: categorical, numerical features + project_title(BOW) + preprocessed_eassay (`BOW with bi-grams` with `min_df=10` and `max_features=5000`)
   - Set 2: categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (`TFIDF with bi-grams` with `min_df=10` and `max_features=5000`)
   - Set 3: categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
   - Set 4: categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

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

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

3. **Representation of results**

   - You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
   - Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
   - Along with plotting ROC curve, you need to print the confusion matrix with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.

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

   - **school_state** : categorical data
   - **clean_categories** : categorical data

- **clean_subcategories** : categorical data
- **project_grade_category** :categorical data
- **teacher_prefix** : categorical data
- **quantity** : numerical data
- **teacher_number_of_previously_posted_projects** : numerical data
- **price** : numerical data
- **sentiment score's of each of the essay** : numerical data
- **number of words in the title** : numerical data
- **number of words in the combine essays** : numerical data

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

6. **Conclusion**

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link

---

**Note: Data Leakage**

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

# 2. Logistic Regression

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

In [30]:

```
#splitting project_data into x and y, y=project_is_approved.

#fetching all the columns except project_is_approved.
cols_to_select = [col for col in project_data.columns if col != 'project_is_approved'];
X = project_data[cols_to_select]
print(X.columns)
y = project_data['project_is_approved'];
print(y.shape)
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
       'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
       'project_essay_3', 'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'clean_categories',
       'clean_subcategories', 'processed_essay', 'processed_titles',
       'processed_grades', 'processed_teacher_prefix', 'price', 'quantity'],
      dtype='object')
(109248,)
```

In [31]:

```
#splitting project_data into train and test and CV data.
X_1, X_test, y_1, y_test = model_selection.train_test_split(X, y, test_size=0.3, random_state=1)
X_train, X_cv, y_train, y_cv = model_selection.train_test_split(X_1, y_1, test_size=0.3, random_sta
te=1);

print('shape of train data ', X_train.shape);
print('shape of test data ', X_test.shape);
print('shape of cross validation data ', X_cv.shape)
```

```
shape of train data  (53531, 19)
shape of test data   (32775, 19)
shape of cross validation data  (22942, 19)
```

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

## Vectorizing Categorical features

In [32]:

```python
#vectorizing school_state
from sklearn.feature_extraction.text import CountVectorizer

#creating dictionary for school_state as state as keys along with no. of projects from that state
as values.
school_state_dict = dict(X_train['school_state'].value_counts());

#configuring CountVectorizer for school state, in which vocabulary will be name of states.
vectorizer = CountVectorizer(vocabulary=list(school_state_dict.keys()), lowercase=False, binary=True);

#applying vectorizer on school_state column to obtain numerical value for each state.
vectorizer.fit(X_train['school_state'].values);

school_state_vector = vectorizer.transform(X_train['school_state'].values);
test_school_state_vector = vectorizer.transform(X_test['school_state'].values);
cv_school_state_vector = vectorizer.transform(X_cv['school_state'].values);

print('shape of matrix after one hot encoding of school_state for train data ',
school_state_vector.shape);
print('shape of matrix after one hot encoding of school_state for test data ',
test_school_state_vector.shape);
print('shape of matrix after one hot encoding of school_state for cv data ',
cv_school_state_vector.shape);

features_name_list = vectorizer.get_feature_names();
```

```
shape of matrix after one hot encoding of school_state for train data  (53531, 51)
shape of matrix after one hot encoding of school_state for test data  (32775, 51)
shape of matrix after one hot encoding of school_state for cv data  (22942, 51)
```

In [33]:

```python
#vectorizing categories



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

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

categories_vector = vectorizer.transform(X_train['clean_categories'].values);
test_categories_vector = vectorizer.transform(X_test['clean_categories'].values);
cv_categories_vector = vectorizer.transform(X_cv['clean_categories'].values);

print('shape of matrix after one hot encoding of clean_categories for train data',
categories_vector.shape)
print('shape of matrix after one hot encoding of clean_categories for test data',
test_categories_vector.shape)
print('shape of matrix after one hot encoding of clean_categories for cv data',
cv_categories_vector.shape)

features_name_list.extend( vectorizer.get_feature_names());
```

```
shape of matrix after one hot encoding of clean_categories for train data (53531, 9)
shape of matrix after one hot encoding of clean_categories for test data (32775, 9)
shape of matrix after one hot encoding of clean_categories for cv data (22942, 9)
```

In [34]:

```python
#vectorizing subcategories

subcategories_dict = dict(X_train['clean_subcategories'].value_counts());
```

```python
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=
True);

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

subcategories_vector = vectorizer.transform(X_train['clean_subcategories'].values);
test_subcategories_vector = vectorizer.transform(X_test['clean_subcategories'].values);
cv_subcategories_vector = vectorizer.transform(X_cv['clean_subcategories'].values);

print('shape of matrix after one hot encoding of clean_subcategories for train data',
subcategories_vector.shape)
print('shape of matrix after one hot encoding of clean_subcategories for test data',
test_subcategories_vector.shape)
print('shape of matrix after one hot encoding of clean_subcategories for cv data',
cv_subcategories_vector.shape)

features_name_list.extend( vectorizer.get_feature_names());
```

```
shape of matrix after one hot encoding of clean_subcategories for train data (53531, 30)
shape of matrix after one hot encoding of clean_subcategories for test data (32775, 30)
shape of matrix after one hot encoding of clean_subcategories for cv data (22942, 30)
```

In [35]:

```python
#vectorizing project_grade_category

grade_dict = dict(X_train['processed_grades'].value_counts());

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

vectorizer.fit(X_train['processed_grades'].values);

grade_vector = vectorizer.transform(X_train['processed_grades'].values);
test_grade_vector = vectorizer.transform(X_test['processed_grades'].values);
cv_grade_vector = vectorizer.transform(X_cv['processed_grades'].values);

print('shape of matrix after one hot encoding of grade_category for train data', grade_vector.shap
e)
print('shape of matrix after one hot encoding of grade_category for test data', test_grade_vector.
shape)
print('shape of matrix after one hot encoding of grade_category for cv data', cv_grade_vector.shap
e)

features_name_list.extend( vectorizer.get_feature_names());
```

```
shape of matrix after one hot encoding of grade_category for train data (53531, 4)
shape of matrix after one hot encoding of grade_category for test data (32775, 4)
shape of matrix after one hot encoding of grade_category for cv data (22942, 4)
```

In [36]:

```python
#vectorizing teacher_prefix

teacher_prefix_dict = dict(X_train['processed_teacher_prefix'].value_counts());

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

vectorizer.fit(X_train['processed_teacher_prefix'].values.astype('U'));

teacher_prefix_vector = vectorizer.transform(X_train['processed_teacher_prefix'].values.astype('U')
);
test_teacher_prefix_vector = vectorizer.transform(X_test['processed_teacher_prefix'].values.astype(
'U'));
cv_teacher_prefix_vector = vectorizer.transform(X_cv['processed_teacher_prefix'].values.astype('U')
);

print('shape of matrix after one hot encoding of teacher_prefix for train data',
teacher_prefix_vector.shape)
print('shape of matrix after one hot encoding of teacher_prefix for test data',
test_teacher_prefix_vector.shape)
print('shape of matrix after one hot encoding of teacher_prefix for cv data',
cv_teacher_prefix_vector.shape)
```

```
features_name_list.extend( vectorizer.get_feature_names());
```

```
shape of matrix after one hot encoding of teacher_prefix for train data (53531, 6)
shape of matrix after one hot encoding of teacher_prefix for test data (32775, 6)
shape of matrix after one hot encoding of teacher_prefix for cv data (22942, 6)
```

## Encoding Numerical data

In [37]:

```python
#vectorizing price

from sklearn.preprocessing import StandardScaler
price_normalizer = StandardScaler()
#configuring StandarScaler to obtain the mean and variance.
price_normalizer.fit(X_train['price'].values.reshape(-1, 1));

# Now standardize the data with maen and variance obtained above.
price_standardized = price_normalizer.transform(X_train['price'].values.reshape(-1, 1))
test_price_standardized = price_normalizer.transform(X_test['price'].values.reshape(-1, 1))
cv_price_standardized = price_normalizer.transform(X_cv['price'].values.reshape(-1, 1))

features_name_list.append('price');
```

In [38]:

```python
#vectorizing teacher_number_of_previously_posted_projects

teacher_normalizer = StandardScaler();

teacher_normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1
));

teacher_number_standardized =
teacher_normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape
(-1,1));

test_teacher_number_standardized =
teacher_normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape
(-1,1));

cv_teacher_number_standardized =
teacher_normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-
1,1));

features_name_list.append('teacher number of previously posted projects');
```

In [39]:

```python
#vectorizing quantity:

quantity_normalizer = StandardScaler();

quantity_normalizer.fit(X_train['quantity'].values.reshape(-1, 1));

quantity_standardized = quantity_normalizer.transform(X_train['quantity'].values.reshape(-1, 1))

test_quantity_standardized = quantity_normalizer.transform(X_test['quantity'].values.reshape(-1, 1)
)

cv_quantity_standardized = quantity_normalizer.transform(X_cv['quantity'].values.reshape(-1, 1))

features_name_list.append('quantity');
```

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

## Vectorizing using BOW on train data

In [40]:

```python
#vectorizing essay

#configure CountVectorizer with word to occur in at least 10 documents.
vectorizer = CountVectorizer(min_df=10, ngram_range=(1,2), max_features=5000);

vectorizer.fit(X_train['processed_essay']);

#transforming essay into vector
essay_bow = vectorizer.transform(X_train['processed_essay']);
cv_essay_bow = vectorizer.transform(X_cv['processed_essay']);
test_essay_bow = vectorizer.transform(X_test['processed_essay']);

print('Shape of matrix after one hot encoding for train data: ', essay_bow.shape);
print('Shape of matrix after one hot encoding for test data: ', test_essay_bow.shape);
print('Shape of matrix after one hot encoding for cv data: ', cv_essay_bow.shape);
```

```
Shape of matrix after one hot encoding for train data:  (53531, 5000)
Shape of matrix after one hot encoding for test data:   (32775, 5000)
Shape of matrix after one hot encoding for cv data:   (22942, 5000)
```

In [41]:

```python
bow_features_name = vectorizer.get_feature_names()
len(bow_features_name)
```

Out[41]:

```
5000
```

In [42]:

```python
#vectorizing project_title

#configure CountVectorizer with word to occur in at least 10 documents.
vectorizer = CountVectorizer();

vectorizer.fit(X_train['processed_titles']);

#transforming title into vector
title_bow = vectorizer.transform(X_train['processed_titles']);
cv_title_bow = vectorizer.transform(X_cv['processed_titles']);
test_title_bow = vectorizer.transform(X_test['processed_titles']);

print('Shape of matrix after one hot encoding for train data: ', title_bow.shape);
print('Shape of matrix after one hot encoding for test data: ', test_title_bow.shape);
print('Shape of matrix after one hot encoding for cv data: ', cv_title_bow.shape);
```

```
Shape of matrix after one hot encoding for train data:  (53531, 12188)
Shape of matrix after one hot encoding for test data:   (32775, 12188)
Shape of matrix after one hot encoding for cv data:   (22942, 12188)
```

In [43]:

```python
bow_features_name.extend(vectorizer.get_feature_names())
print(len(bow_features_name))
```

```
17188
```

In [44]:

```python
len(features_name_list)
```

Out[44]:

```
103
```

In [45]:

```
final_bow_featues_name = [];
final_bow_featues_name.extend(features_name_list);
final_bow_featues_name.extend(bow_features_name);
print(len(final_bow_featues_name))
```

```
17291
```

## Vectorizing using tf-idf

In [46]:

```
#vectorizing essay

#importing TfidfVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer

#configuring TfidfVectorizer with a word to occur atleast in 10 documnets.
vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)

vectorizer.fit(X_train['processed_essay']);

#vectorizing essay using tfidf
essay_tfidf = vectorizer.transform(X_train['processed_essay']);
test_essay_tfidf = vectorizer.transform(X_test['processed_essay']);
cv_essay_tfidf = vectorizer.transform(X_cv['processed_essay']);

print("Shape of matrix after one hot encoding for train data: ",essay_tfidf.shape)
print("Shape of matrix after one hot encoding for test data: ",test_essay_tfidf.shape)
print("Shape of matrix after one hot encoding for cv data: ",cv_essay_tfidf.shape)
```

```
Shape of matrix after one hot encoding for train data:  (53531, 5000)
Shape of matrix after one hot encoding for test data:  (32775, 5000)
Shape of matrix after one hot encoding for cv data:  (22942, 5000)
```

In [47]:

```
#vectorizing project_title

vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,2), max_features=5000);

vectorizer.fit(X_train['processed_titles']);

title_tfidf = vectorizer.transform(X_train['processed_titles']);
test_title_tfidf = vectorizer.transform(X_test['processed_titles']);
cv_title_tfidf = vectorizer.transform(X_cv['processed_titles']);

print('Shape of title_tfidf after one hot encoding for train data ', title_tfidf.shape)
print('Shape of title_tfidf after one hot encoding for test data ', test_title_tfidf.shape)
print('Shape of title_tfidf after one hot encoding for cv data ', cv_title_tfidf.shape)
```

```
Shape of title_tfidf after one hot encoding for train data   (53531, 5000)
Shape of title_tfidf after one hot encoding for test data   (32775, 5000)
Shape of title_tfidf after one hot encoding for cv data   (22942, 5000)
```

## Vectorizing using avg w2v on train

In [48]:

```
#vectorizing essay

essay_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['processed_essay']): # 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
    essay_avg_w2v.append(vector)

#printing number of documents
print(len(essay_avg_w2v))

#printing dimension of each essay avg w2v
print(len(essay_avg_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████████| 53531/53531
[00:15<00:00, 3371.67it/s]
```

```
53531
300
```

In [49]:

```python
#vectorizing project_title

title_avg_w2v = [];
for sentence in tqdm(X_train['processed_titles']):
    vector = np.zeros(300);
    cnt_words = 0;
    for word in sentence.split():
        if word in glove_words:
            vector += model[word];
            cnt_words += 1;
    if cnt_words != 0:
        vector /= cnt_words;
    title_avg_w2v.append(vector);

print(len(title_avg_w2v));
print(len(title_avg_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████████| 53531/53531
[00:00<00:00, 128350.70it/s]
```

```
53531
300
```

## Vectorizing using avg w2v on CV

In [50]:

```python
#vectorizing essay

cv_essay_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['processed_essay']): # 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
    cv_essay_avg_w2v.append(vector)

#printing number of documents
print(len(cv_essay_avg_w2v))

#printing dimension of each essay avg w2v
print(len(cv_essay_avg_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████████| 22942/22942
[00:06<00:00, 3371.35it/s]
```

```
22942
```

```
300
```

```
#vectorizing project_title

cv_title_avg_w2v = [];
for sentence in tqdm(X_cv['processed_titles']):
    vector = np.zeros(300);
    cnt_words = 0;
    for word in sentance.split():
        if word in glove_words:
            vector += model[word];
            cnt_words += 1;
    if cnt_words != 0:
        vector /= cnt_words;
    cv_title_avg_w2v.append(vector);

print(len(cv_title_avg_w2v));
print(len(cv_title_avg_w2v[0]))
```

```
100%|████████████████████████████████████████| 22942/22942
[00:00<00:00, 116520.86it/s]
```

```
22942
300
```

## Vectorizing using avg w2v on test data

```
#vectorizing essay

test_essay_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['processed_essay']): # 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
    test_essay_avg_w2v.append(vector)

#printing number of documents
print(len(test_essay_avg_w2v))

#printing dimension of each essay avg w2v
print(len(test_essay_avg_w2v[0]))
```

```
100%|████████████████████████████████████████| 32775/32775
[00:10<00:00, 3072.54it/s]
```

```
32775
300
```

```
#vectorizing project_title

test_title_avg_w2v = [];
for sentence in tqdm(X_test['processed_titles']):
    vector = np.zeros(300);
    cnt_words = 0;
    for word in sentance.split():
        if word in glove_words:
            vector += model[word];
            cnt_words += 1;
    if cnt_words != 0:
```

```
        vector /= cnt_words;
    test_title_avg_w2v.append(vector);

print(len(test_title_avg_w2v));
print(len(test_title_avg_w2v[0]))
```

```
32775
300
```

## Vectorizing using tfidf weighted w2v

In [54]:

```python
#finding out tfidf words and corresponding idf value for essay

tfidf_model = TfidfVectorizer()

tfidf_model.fit(X_train['processed_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 = set(tfidf_model.get_feature_names())
```

In [55]:

```python
#vectorizing essay

essay_tfidf_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['processed_essay']): # 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 tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    essay_tfidf_w2v.append(vector)

print(len(essay_tfidf_w2v))
print(len(essay_tfidf_w2v[0]))
```

```
53531
300
```

In [56]:

```python
#vectorizing essay

cv_essay_tfidf_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['processed_essay']): # 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 tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    cv_essay_tfidf_w2v.append(vector)

print(len(cv_essay_tfidf_w2v))
print(len(cv_essay_tfidf_w2v[0]))
```

100%|██████████████████████████████████████████████████████████| 22942/22942 [00:
55<00:00, 413.16it/s]

22942
300

In [57]:

```
#vectorizing essay

test_essay_tfidf_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['processed_essay']): # 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 tf
idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    test_essay_tfidf_w2v.append(vector)

print(len(test_essay_tfidf_w2v))
print(len(test_essay_tfidf_w2v[0]))
```

100%|██████████████████████████████████████████████████████████| 32775/32775 [01:
18<00:00, 418.97it/s]

32775
300

In [58]:

```
#finding out tfidf words and corresponding idf value for project_title

tfidf_model = TfidfVectorizer()

tfidf_model.fit(X_train['processed_titles'])

# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))

tfidf_words = set(tfidf_model.get_feature_names())
```

In [59]:

```
#vectorizing project_tile

title_tfidf_w2v = [];

for sentence in tqdm(X_train['processed_titles']):
    vector = np.zeros(300);
    tfidf_weight = 0;
    for word in sentence.split():
        if (word in glove_words) and (word in tfidf_words):
```

```
                tfidf = dictionary[word] * (sentance.count(word) / len(sentance.split()));
                vector = tfidf * model[word];
                tfidf_weight += tfidf;
        if tfidf_weight != 0:
            vector /= tfidf_weight;
        title_tfidf_w2v.append(vector);

print(len(title_tfidf_w2v))
print(len(title_tfidf_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████| 53531/53531
[00:00<00:00, 77626.86it/s]
```

```
53531
300
```

In [60]:

```
#vectorizing project_tile

cv_title_tfidf_w2v = [];

for sentance in tqdm(X_cv['processed_titles']):
    vector = np.zeros(300);
    tfidf_weight = 0;
    for word in sentance.split():
        if (word in glove_words) and (word in tfidf_words):
            tfidf = dictionary[word] * (sentance.count(word) / len(sentance.split()));
            vector = tfidf * model[word];
            tfidf_weight += tfidf;
    if tfidf_weight != 0:
        vector /= tfidf_weight;
    cv_title_tfidf_w2v.append(vector);

print(len(cv_title_tfidf_w2v))
print(len(cv_title_tfidf_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████| 22942/22942
[00:00<00:00, 65964.51it/s]
```

```
22942
300
```

In [61]:

```
#vectorizing project_tile

test_title_tfidf_w2v = [];

for sentance in tqdm(X_test['processed_titles']):
    vector = np.zeros(300);
    tfidf_weight = 0;
    for word in sentance.split():
        if (word in glove_words) and (word in tfidf_words):
            tfidf = dictionary[word] * (sentance.count(word) / len(sentance.split()));
            vector = tfidf * model[word];
            tfidf_weight += tfidf;
    if tfidf_weight != 0:
        vector /= tfidf_weight;
    test_title_tfidf_w2v.append(vector);

print(len(test_title_tfidf_w2v))
print(len(test_title_tfidf_w2v[0]))
```

```
100%|████████████████████████████████████████████████████████| 32775/32775
[00:00<00:00, 70373.00it/s]
```

```
32775
300
```

## Finding count of words in essay and project_title

In [62]:

```python
train_essay_words_counts = []
for i in X_train['processed_essay']:
    train_essay_words_counts.append(len(i.split()))
train_essay_words_counts = np.array(train_essay_words_counts).reshape(-1, 1);
print(train_essay_words_counts.shape)

test_essay_words_counts = []
for i in X_test['processed_essay']:
    test_essay_words_counts.append(len(i.split()))
test_essay_words_counts = np.array(test_essay_words_counts).reshape(-1, 1);
print(test_essay_words_counts.shape)

cv_essay_words_counts = []
for i in X_cv['processed_essay']:
    cv_essay_words_counts.append(len(i.split()))
cv_essay_words_counts = np.array(cv_essay_words_counts).reshape(-1,1)
print(cv_essay_words_counts.shape)
```

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

In [63]:

```python
train_project_title_words_counts = []
for i in X_train['processed_titles']:
    train_project_title_words_counts.append(len(i.split()))
train_project_title_words_counts = np.array(train_project_title_words_counts).reshape(-1,1);
print(train_project_title_words_counts.shape)

test_project_title_words_counts = []
for i in X_test['processed_titles']:
    test_project_title_words_counts.append(len(i.split()))
test_project_title_words_counts = np.array(test_project_title_words_counts).reshape(-1, 1);
print(test_project_title_words_counts.shape)

cv_project_title_words_counts = []
for i in X_cv['processed_titles']:
    cv_project_title_words_counts.append(len(i.split()))
cv_project_title_words_counts = np.array(cv_project_title_words_counts).reshape(-1, 1);

print(cv_project_title_words_counts.shape)
```

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

In [64]:

```python
project_data.columns
```

Out[64]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
       'project_submitted_datetime', 'project_essay_1', 'project_essay_2',
       'project_essay_3', 'project_essay_4', 'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'project_is_approved',
       'clean_categories', 'clean_subcategories', 'processed_essay',
       'processed_titles', 'processed_grades', 'processed_teacher_prefix',
       'price', 'quantity'],
      dtype='object')
```

## Finding Sentiments score for each essay

In [65]:

```python
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

train_neg_sentiments = [];
train_pos_sentiments = [];
train_neu_sentiments = [];
train_comp_sentiments = [];

sid = SentimentIntensityAnalyzer()

for essay in X_train['processed_essay']:
    ss = sid.polarity_scores(essay);
    train_neg_sentiments.append(ss['neg']);
    train_pos_sentiments.append(ss['pos']);
    train_neu_sentiments.append(ss['neu']);
    train_comp_sentiments.append(ss['compound']);

train_neg_sentiments = np.array(train_neg_sentiments).reshape(-1, 1);
print(train_neg_sentiments.shape);

train_pos_sentiments = np.array(train_pos_sentiments).reshape(-1,1);
print(train_pos_sentiments.shape);

train_neu_sentiments = np.array(train_neu_sentiments).reshape(-1,1);
print(train_neu_sentiments.shape);

train_comp_sentiments = np.array(train_comp_sentiments).reshape(-1, 1);
print(train_comp_sentiments.shape);
```

```
(53531, 1)
(53531, 1)
(53531, 1)
(53531, 1)
```

In [66]:

```python
cv_neg_sentiments = [];
cv_pos_sentiments = [];
cv_neu_sentiments = [];
cv_comp_sentiments = [];

sid = SentimentIntensityAnalyzer()

for essay in X_cv['processed_essay']:
    ss = sid.polarity_scores(essay);
    cv_neg_sentiments.append(ss['neg']);
    cv_pos_sentiments.append(ss['pos']);
    cv_neu_sentiments.append(ss['neu']);
    cv_comp_sentiments.append(ss['compound']);

cv_neg_sentiments = np.array(cv_neg_sentiments).reshape(-1, 1);
print(cv_neg_sentiments.shape);

cv_pos_sentiments = np.array(cv_pos_sentiments).reshape(-1, 1);
print(cv_pos_sentiments.shape);

cv_neu_sentiments = np.array(cv_neu_sentiments).reshape(-1, 1);
print(cv_neu_sentiments.shape);

cv_comp_sentiments = np.array(cv_comp_sentiments).reshape(-1, 1);
print(cv_comp_sentiments.shape);
```

```
(22942, 1)
(22942, 1)
(22942, 1)
(22942, 1)
```

In [67]:

```python
test_neg_sentiments = [];
test_pos_sentiments = [];
test_neu_sentiments = [];
test_comp_sentiments = [];
```

```python
sid = SentimentIntensityAnalyzer()

for essay in X_test['processed_essay']:
    ss = sid.polarity_scores(essay);
    test_neg_sentiments.append(ss['neg']);
    test_pos_sentiments.append(ss['pos']);
    test_neu_sentiments.append(ss['neu']);
    test_comp_sentiments.append(ss['compound']);

test_neg_sentiments = np.array(test_neg_sentiments).reshape(-1, 1);
print(test_neg_sentiments.shape);

test_pos_sentiments = np.array(test_pos_sentiments).reshape(-1, 1);
print(test_pos_sentiments.shape);

test_neu_sentiments = np.array(test_neu_sentiments).reshape(-1, 1);
print(test_neu_sentiments.shape);

test_comp_sentiments = np.array(test_comp_sentiments).reshape(-1, 1);
print(test_comp_sentiments.shape);
```

```
(32775, 1)
(32775, 1)
(32775, 1)
(32775, 1)
```

## Merging data

In [68]:

```python
from scipy.sparse import hstack

#concatinating train data
#with bow
train_set_1 = hstack((school_state_vector, categories_vector, subcategories_vector, grade_vector, t
eacher_prefix_vector, price_standardized, teacher_number_standardized, quantity_standardized,
essay_bow, title_bow)).tocsr()

#with tfidf
train_set_2 = hstack((school_state_vector, categories_vector, subcategories_vector, grade_vector, t
eacher_prefix_vector, price_standardized, teacher_number_standardized, quantity_standardized,
essay_tfidf, title_tfidf)).tocsr()

#with avg w2v
train_set_3 = hstack((school_state_vector, categories_vector, subcategories_vector, grade_vector, t
eacher_prefix_vector, price_standardized, teacher_number_standardized, quantity_standardized,
essay_avg_w2v, title_avg_w2v)).tocsr()

#with tfidf wt w2v
train_set_4 = hstack((school_state_vector, categories_vector, subcategories_vector, grade_vector, t
eacher_prefix_vector, price_standardized, teacher_number_standardized, quantity_standardized,
essay_tfidf_w2v, title_tfidf_w2v)).tocsr()

train_set_5 = hstack((school_state_vector, categories_vector, subcategories_vector, grade_vector, t
eacher_prefix_vector, price_standardized, teacher_number_standardized, quantity_standardized,
train_essay_words_counts, train_project_title_words_counts, train_comp_sentiments)).tocsr()

#concatinating cv data
#with bow
cv_set_1 = hstack((cv_school_state_vector, cv_categories_vector, cv_subcategories_vector,
cv_grade_vector, cv_teacher_prefix_vector, cv_price_standardized, cv_teacher_number_standardized,
cv_quantity_standardized, cv_essay_bow, cv_title_bow)).tocsr()

#with tfidf
cv_set_2 = hstack((cv_school_state_vector, cv_categories_vector, cv_subcategories_vector,
cv_grade_vector, cv_teacher_prefix_vector, cv_price_standardized, cv_teacher_number_standardized,
cv_quantity_standardized, cv_essay_tfidf, cv_title_tfidf)).tocsr()

#with avg w2v
cv_set_3 = hstack((cv_school_state_vector, cv_categories_vector, cv_subcategories_vector,
cv_grade_vector, cv_teacher_prefix_vector, cv_price_standardized, cv_teacher_number_standardized,
cv_quantity_standardized, cv_essay_avg_w2v, cv_title_avg_w2v)).tocsr()
```

```
#with tfidf wt w2v
cv_set_4 = hstack((cv_school_state_vector, cv_categories_vector, cv_subcategories_vector,
cv_grade_vector, cv_teacher_prefix_vector, cv_price_standardized, cv_teacher_number_standardized,
cv_quantity_standardized, cv_essay_tfidf_w2v, cv_title_tfidf_w2v)).tocsr()

cv_set_5 = hstack((cv_school_state_vector, cv_categories_vector, cv_subcategories_vector,
cv_grade_vector, cv_teacher_prefix_vector, cv_price_standardized, cv_teacher_number_standardized,
cv_quantity_standardized, cv_essay_words_counts, cv_project_title_words_counts, cv_comp_sentiments
)).tocsr()

#concatinating test data
#with bow
test_set_1 = hstack((test_school_state_vector, test_categories_vector, test_subcategories_vector,
test_grade_vector, test_teacher_prefix_vector, test_price_standardized,
test_teacher_number_standardized, test_quantity_standardized, test_essay_bow,
test_title_bow)).tocsr()

#with tfidf
test_set_2 = hstack((test_school_state_vector, test_categories_vector, test_subcategories_vector,
test_grade_vector, test_teacher_prefix_vector, test_price_standardized,
test_teacher_number_standardized, test_quantity_standardized, test_essay_tfidf, test_title_tfidf))
.tocsr()

#with avg w2v
test_set_3 = hstack((test_school_state_vector, test_categories_vector, test_subcategories_vector,
test_grade_vector, test_teacher_prefix_vector, test_price_standardized,
test_teacher_number_standardized, test_quantity_standardized, test_essay_avg_w2v,
test_title_avg_w2v)).tocsr()

#with tfidf wt w2v
test_set_4 = hstack((test_school_state_vector, test_categories_vector, test_subcategories_vector,
test_grade_vector, test_teacher_prefix_vector, test_price_standardized,
test_teacher_number_standardized, test_quantity_standardized, test_essay_tfidf_w2v,
test_title_tfidf_w2v)).tocsr()

test_set_5 = hstack((test_school_state_vector, test_categories_vector, test_subcategories_vector,
test_grade_vector, test_teacher_prefix_vector, test_price_standardized,
test_teacher_number_standardized, test_quantity_standardized, test_essay_words_counts,
test_project_title_words_counts, test_comp_sentiments)).tocsr()
```

In [69]:

```
print(train_set_1.shape, cv_set_1.shape, test_set_1.shape)
print(train_set_2.shape, cv_set_2.shape, test_set_2.shape)
print(train_set_3.shape, cv_set_3.shape, test_set_3.shape)
print(train_set_4.shape, cv_set_4.shape, test_set_4.shape)
```

```
(53531, 17291) (22942, 17291) (32775, 17291)
(53531, 10103) (22942, 10103) (32775, 10103)
(53531, 703) (22942, 703) (32775, 703)
(53531, 703) (22942, 703) (32775, 703)
```

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

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

In [70]:

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

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

    # (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
```

## Applying LogisticRegression on set 1

```python
from sklearn.linear_model import SGDClassifier;
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score;

#creating list for holding auc value for train, cv
train_auc = [];
cv_auc = [];

#defining list of lambda's
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000];

for a in alpha:
    #using SGDClassifier, and passing log in loss, which makes it LogisticRegression
    LR = SGDClassifier(loss='log', penalty='l2', alpha=a); #usig L2 Regularization
    LR.fit(train_set_1, y_train); #training model using training data.

    y_train_pred = LR.predict_proba(train_set_1)[:, 1]; #predicting probability for training data
    y_cv_pred = LR.predict_proba(cv_set_1)[:, 1]; #predicting probability for cv data

    train_auc.append(roc_auc_score(y_train, y_train_pred));
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred));

#plotting error plot
plt.plot(np.log(alpha), train_auc, label='Train AUC');
plt.plot(np.log(alpha), cv_auc, label='CV AUC');

plt.scatter(np.log(alpha), train_auc, label='Train AUC points');
plt.scatter(np.log(alpha), cv_auc, label='CV AUC points');

plt.xlabel('lambda: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');

plt.xticks(np.log(alpha))

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

```python
optimal_alpha = 0.01
set1_alpha = optimal_alpha;
print(optimal_alpha)
#training model using optimal_alpha
LR = SGDClassifier(loss='log', penalty='l2', alpha=optimal_alpha);
```

```
LR.fit(train_set_1, y_train);

y_train_pred = LR.predict_proba(train_set_1)[:, 1];
y_test_pred = LR.predict_proba(test_set_1)[:, 1];

train_fpr, train_tpr, train_thresholds = metrics.roc_curve(y_train, y_train_pred);
test_fpr, test_tpr, test_thresholds = metrics.roc_curve(y_test, y_test_pred);

train_auc = roc_auc_score(y_train, y_train_pred);
test_auc = roc_auc_score(y_test, y_test_pred);

set1_auc = test_auc;

#plotting ROC curve
plt.plot(train_fpr, train_tpr, label="train AUC: "+str(train_auc))
plt.plot(test_fpr, test_tpr, label="test AUC: "+str(test_auc))


plt.grid();
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC curve')
plt.legend();
plt.show()
```

0.01



In [73]:

```
#Testing whether choose optimal alpha is correct or not
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.multiclass import OneVsRestClassifier

tuned_parameters=[{'estimator__C': [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]}]
log_reg_clf = OneVsRestClassifier(LogisticRegression())
model = GridSearchCV(log_reg_clf, tuned_parameters, scoring = 'roc_auc', cv=5)
model.fit(train_set_1, y_train);
```

In [74]:

```
print(model.best_estimator_)
print(model.score(test_set_1, y_test))
```

```
OneVsRestClassifier(estimator=LogisticRegression(C=0.01, class_weight=None, dual=False,
fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=None, solver='warn',
          tol=0.0001, verbose=0, warm_start=False),
          n_jobs=None)
0.7284055363911082
```

In [75]:

```
import seaborn as sns
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_thresholds, train_fpr,
train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.5011600384969563 for threshold 0.821

```
print("Test confusion matrix")
test_cm = confusion_matrix(y_test, predict(y_test_pred, train_thresholds, test_fpr, test_tpr))
sns.heatmap(test_cm, annot=True, fmt="d");
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4524596127295528 for threshold 0.857

```
#Finding top 20 positive and negative features
top_features_weights= LR.coef_
top_features_indices = np.argsort(top_features_weights[0, :])

top_negative_features_indices = top_features_indices[:20];
top_negative_features = [final_bow_featues_name[i] for i in top_negative_features_indices]

print('Negative Features');
print(top_negative_features)

top_positive_features_indices = top_features_indices[-20:];
top_positive_features = [final_bow_featues_name[i] for i in top_positive_features_indices]

print('*'*100);
print('Positive Features');
print(top_positive_features)
```

Negative Features
['supplies', 'price', 'these materials', 'materials', 'items', 'the materials', 'the students', 'q
```

uantity', 'manipulatives', 'materials help', 'going', 'options', 'equipment', 'materials allow', '
taught', 'breakout', 'ways', 'resources', 'supplies', 'today']
*********************************************************************************************

Positive Features
['requesting', 'pencils', 'rug', 'paper', 'carpet', 'markers', 'set', 'Mrs', 'kits', 'used', '3d',
'wobble', 'headphones', 'chairs', 'books', 'chromebooks', 'balls', 'stools', 'nannan', 'teacher nu
mber of previously posted projects']

## Performing perturbation test

In [78]:

```python
#how to find index of non zero value in sparse matrix https://docs.scipy.org/doc/scipy-
0.19.0/reference/generated/scipy.sparse.find.html

from scipy.sparse import csr_matrix, find
print(type(train_set_1))

#scipy find return row indices, column indices, value for non-zero element
r, c, v = find(train_set_1);

print(train_set_1[r[1], c[1]])
```

```
<class 'scipy.sparse.csr.csr_matrix'>
1.0
```

In [79]:

```python
#adding noise
t = train_set_1;
t[t.nonzero()] = t[t.nonzero()]+np.random.normal(0, 1);
print(type(t))
print(t[r[1], c[1]])
```

```
<class 'scipy.sparse.csr.csr_matrix'>
1.8817895013686061
```

In [80]:

```python
t.shape
```

Out[80]:

```
(53531, 17291)
```

In [81]:

```python
#training model on new data with noise

LR = SGDClassifier(loss='log', penalty='l2', alpha=set1_alpha);
LR.fit(t, y_train)
```

Out[81]:

```
SGDClassifier(alpha=0.01, average=False, class_weight=None,
       early_stopping=False, epsilon=0.1, eta0=0.0, fit_intercept=True,
       l1_ratio=0.15, learning_rate='optimal', loss='log', max_iter=None,
       n_iter=None, n_iter_no_change=5, n_jobs=None, penalty='l2',
       power_t=0.5, random_state=None, shuffle=True, tol=None,
       validation_fraction=0.1, verbose=0, warm_start=False)
```

In [82]:

```python
#finding weights for each feature
new_feature_weights = LR.coef_

print(new_feature_weights.shape);
```

```
#adding small values to feature weights to avoid division by zero
new_feature_weights = new_feature_weights + 0.00001
top_features_weights = top_features_weights + 0.00001
```

(1, 17291)

```
#finding % change
w = (abs(top_features_weights - new_feature_weights)/top_features_weights)*100
```

```
w[0, 7180]
```

65.10528703270327

## Applying Logistic Regression on set 2

```
#creating list for holding auc value for train, cv
train_auc = [];
cv_auc = [];

#defining list of lambda's
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000];

for a in alpha:
    #using SGDClassifier, and passing log in loss, which makes it LogisticRegression
    LR = SGDClassifier(loss='log', penalty='l2', alpha=a); #usig L1 Regularization
    LR.fit(train_set_2, y_train); #training model using training data.

    y_train_pred = LR.predict_proba(train_set_2)[:, 1]; #predicting probability for training data
    y_cv_pred = LR.predict_proba(cv_set_2)[:, 1]; #predicting probability for cv data

    train_auc.append(roc_auc_score(y_train, y_train_pred));
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred));
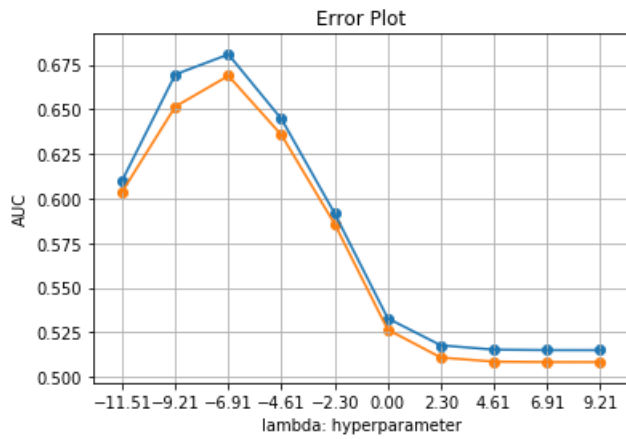
#plotting error plot
plt.plot(np.log(alpha), train_auc, label='Train AUC');
plt.plot(np.log(alpha), cv_auc, label='CV AUC');

plt.scatter(np.log(alpha), train_auc, label='Train AUC points');
plt.scatter(np.log(alpha), cv_auc, label='CV AUC points');

plt.xlabel('lambda: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');

plt.xticks(np.log(alpha))

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

$$-11.51\ -9.21\ -6.91\ -4.61\ -2.30\ 0.00\ 2.30\ 4.61\ 6.91\ 9.21$$
lambda: hyperparameter

In [86]:

```python
optimal_alpha = 0.001
set2_alpha = optimal_alpha;
print(optimal_alpha)
#training model using optimal_alpha
LR = SGDClassifier(loss='log', penalty='l2', alpha=optimal_alpha);

LR.fit(train_set_2, y_train);

y_train_pred = LR.predict_proba(train_set_2)[:, 1];
y_test_pred = LR.predict_proba(test_set_2)[:, 1];

train_fpr, train_tpr, train_thresholds = metrics.roc_curve(y_train, y_train_pred);
test_fpr, test_tpr, test_thresholds = metrics.roc_curve(y_test, y_test_pred);

train_auc = roc_auc_score(y_train, y_train_pred);
test_auc = roc_auc_score(y_test, y_test_pred);

set2_auc = test_auc;

#plotting ROC curve
plt.plot(train_fpr, train_tpr, label="train AUC: "+str(train_auc))
plt.plot(test_fpr, test_tpr, label="test AUC: "+str(test_auc))


plt.grid();
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC curve')
plt.legend();
plt.show()
```

0.001



In [87]:

```python
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_thresholds, train_fpr,
train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.4287335595889016 for threshold 0.833

```
print("Test confusion matrix")
test_cm = confusion_matrix(y_test, predict(y_test_pred, train_thresholds, test_fpr, test_tpr))
sns.heatmap(test_cm, annot=True, fmt="d");
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4060542353223939 for threshold 0.851



## Applying Logistic Regression on set 3

```
#creating list for holding auc value for train, cv
train_auc = [];
cv_auc = [];

#defining list of lambda's
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000];

for a in alpha:
    #using SGDClassifier, and passing log in loss, which makes it LogisticRegression
    LR = SGDClassifier(loss='log', penalty='l2', alpha=a); #usig L1 Regularization
    LR.fit(train_set_3, y_train); #training model using training data.

    y_train_pred = LR.predict_proba(train_set_3)[:, 1]; #predicting probability for training data
    y_cv_pred = LR.predict_proba(cv_set_3)[:, 1]; #predicting probability for cv data

    train_auc.append(roc_auc_score(y_train, y_train_pred));
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred));

#plotting error plot
plt.plot(np.log(alpha), train_auc, label='Train AUC');
plt.plot(np.log(alpha), cv_auc, label='CV AUC');

plt.scatter(np.log(alpha), train_auc, label='Train AUC points');
plt.scatter(np.log(alpha), cv_auc, label='CV AUC points');

plt.xlabel('lambda: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');

plt.xticks(np.log(alpha))

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

Error Plot

```python
optimal_alpha = 0.001
set3_alpha = optimal_alpha;
print(optimal_alpha)
#training model using optimal_alpha
LR = SGDClassifier(loss='log', penalty='l2', alpha=optimal_alpha);

LR.fit(train_set_3, y_train);

y_train_pred = LR.predict_proba(train_set_3)[:, 1];
y_test_pred = LR.predict_proba(test_set_3)[:, 1];

train_fpr, train_tpr, train_thresholds = metrics.roc_curve(y_train, y_train_pred);
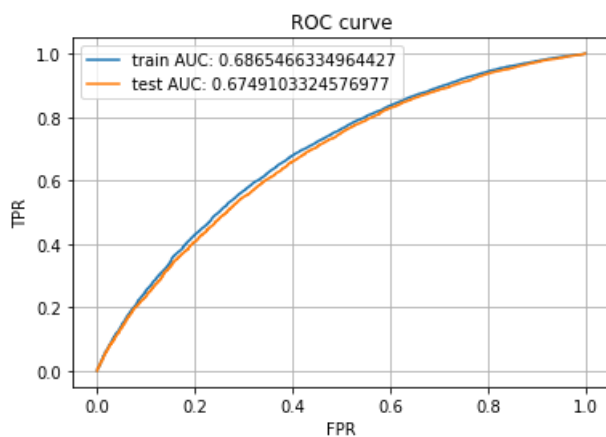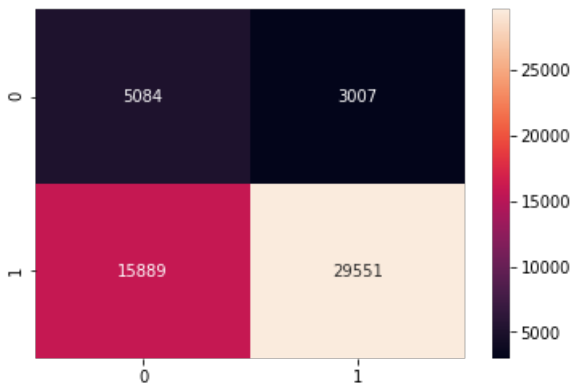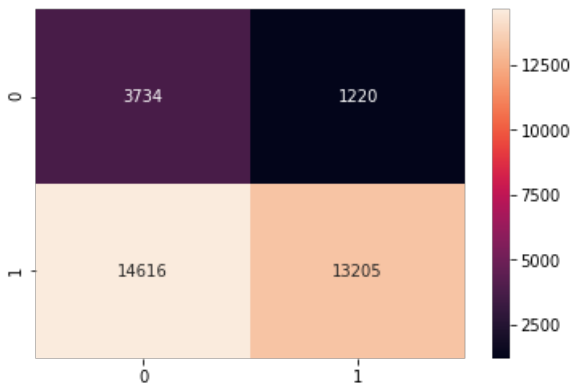test_fpr, test_tpr, test_thresholds = metrics.roc_curve(y_test, y_test_pred);

train_auc = roc_auc_score(y_train, y_train_pred);
test_auc = roc_auc_score(y_test, y_test_pred);

set3_auc = test_auc;

#plotting ROC curve
plt.plot(train_fpr, train_tpr, label="train AUC: "+str(train_auc))
plt.plot(test_fpr, test_tpr, label="test AUC: "+str(test_auc))


plt.grid();
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC curve')
plt.legend();
plt.show()
```

```
0.001
```



ROC curve
train AUC: 0.6865466334964427
test AUC: 0.6749103324576977

```python
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_thresholds, train_fpr,
```

```
train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.4086365414710485 for threshold 0.839

```
print("Test confusion matrix")
test_cm = confusion_matrix(y_test, predict(y_test_pred, train_thresholds, test_fpr, test_tpr))
sns.heatmap(test_cm, annot=True, fmt="d");
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.3969897994151057 for threshold 0.867



## Applying Logistic Regression on set 4

In [93]:

```
#creating list for holding auc value for train, cv
train_auc = [];
cv_auc = [];

#defining list of lambda's
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000];

for a in alpha:
    #using SGDClassifier, and passing log in loss, which makes it LogisticRegression
    LR = SGDClassifier(loss='log', penalty='l2', alpha=a); #usig L2 Regularization
    LR.fit(train_set_4, y_train); #training model using training data.

    y_train_pred = LR.predict_proba(train_set_4)[:, 1]; #predicting probability for training data
    y_cv_pred = LR.predict_proba(cv_set_4)[:, 1]; #predicting probability for cv data

    train_auc.append(roc_auc_score(y_train, y_train_pred));
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred));

#plotting error plot
plt.plot(np.log(alpha), train_auc, label='Train AUC');
```

```
plt.plot(np.log(alpha), cv_auc, label='CV AUC');

plt.scatter(np.log(alpha), train_auc, label='Train AUC points');
plt.scatter(np.log(alpha), cv_auc, label='CV AUC points');

plt.xlabel('lambda: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');

plt.xticks(np.arange(0, 10, 1))

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

```
optimal_alpha = 0.001
set4_alpha = optimal_alpha;
print(optimal_alpha)
#training model using optimal_alpha
LR = SGDClassifier(loss='log', penalty='l2', alpha=optimal_alpha);

LR.fit(train_set_4, y_train);

y_train_pred = LR.predict_proba(train_set_4)[:, 1];
y_test_pred = LR.predict_proba(test_set_4)[:, 1];

train_fpr, train_tpr, train_thresholds = metrics.roc_curve(y_train, y_train_pred);
test_fpr, test_tpr, test_thresholds = metrics.roc_curve(y_test, y_test_pred);

train_auc = roc_auc_score(y_train, y_train_pred);
test_auc = roc_auc_score(y_test, y_test_pred);

set4_auc = test_auc;

#plotting ROC curve
plt.plot(train_fpr, train_tpr, label="train AUC: "+str(train_auc))
plt.plot(test_fpr, test_tpr, label="test AUC: "+str(test_auc))

plt.grid();
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC curve')
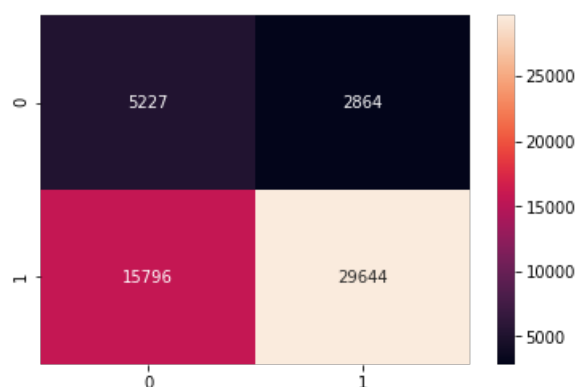plt.legend();
plt.show()
```

0.001

```
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_thresholds, train_fpr,
train_tpr))
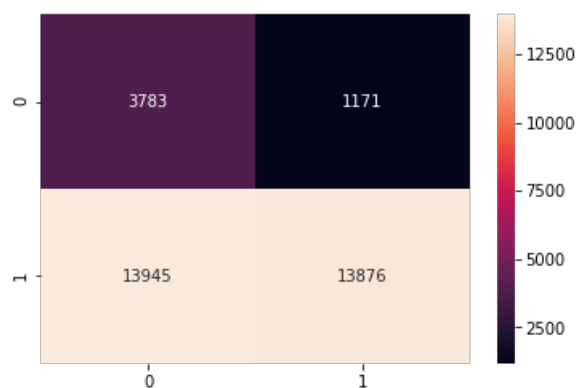sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.42145264212888245 for threshold 0.866

```
print("Test confusion matrix")
test_cm = confusion_matrix(y_test, predict(y_test_pred, train_thresholds, test_fpr, test_tpr))
sns.heatmap(test_cm, annot=True, fmt="d");
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4142935828427471 for threshold 0.892



## 2.5 Logistic Regression with added Features `Set 5`

```
#creating list for holding auc value for train, cv
train_auc = [];
cv_auc = [];

#defining list of lambda's
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000];
```

```
for a in alpha:
    #using SGDClassifier, and passing log in loss, which makes it LogisticRegression
    LR = SGDClassifier(loss='log', penalty='l2', alpha=a); #usig L2 Regularization
    LR.fit(train_set_5, y_train); #training model using training data.

    y_train_pred = LR.predict_proba(train_set_5)[:, 1]; #predicting probability for training data
    y_cv_pred = LR.predict_proba(cv_set_5)[:, 1]; #predicting probability for cv data

    train_auc.append(roc_auc_score(y_train, y_train_pred));
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred));

#plotting error plot
plt.plot(np.log(alpha), train_auc, label='Train AUC');
plt.plot(np.log(alpha), cv_auc, label='CV AUC');

plt.scatter(np.log(alpha), train_auc, label='Train AUC points');
plt.scatter(np.log(alpha), cv_auc, label='CV AUC points');

plt.xlabel('lambda: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');

plt.xticks(np.arange(0, 10, 1))

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

```
optimal_alpha = 0.01
set5_alpha = optimal_alpha;
print(optimal_alpha)
#training model using optimal_alpha
LR = SGDClassifier(loss='log', penalty='l2', alpha=optimal_alpha);

LR.fit(train_set_5, y_train);

y_train_pred = LR.predict_proba(train_set_5)[:, 1];
y_test_pred = LR.predict_proba(test_set_5)[:, 1];

train_fpr, train_tpr, train_thresholds = metrics.roc_curve(y_train, y_train_pred);
test_fpr, test_tpr, test_thresholds = metrics.roc_curve(y_test, y_test_pred);

train_auc = roc_auc_score(y_train, y_train_pred);
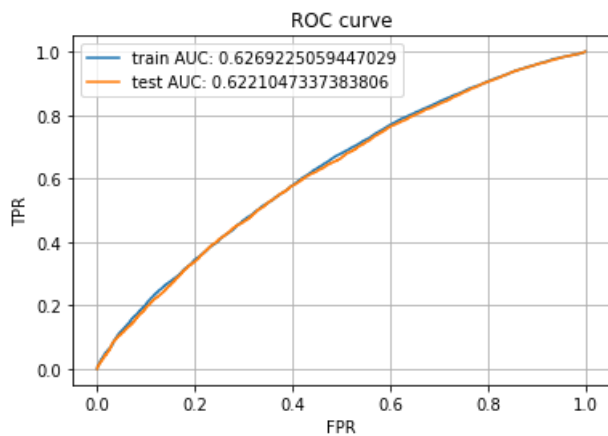test_auc = roc_auc_score(y_test, y_test_pred);

set5_auc = test_auc;

#plotting ROC curve
plt.plot(train_fpr, train_tpr, label="train AUC: "+str(train_auc))
plt.plot(test_fpr, test_tpr, label="test AUC: "+str(test_auc))


plt.grid();
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC curve')
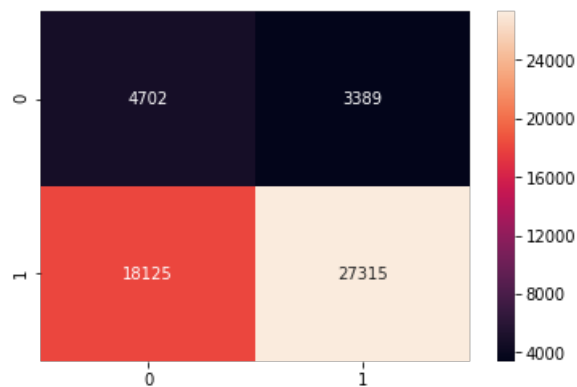```

```
plt.legend();
plt.show()
```

0.01



ROC curve

```
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_thresholds, train_fpr,
train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.34933596993529586 for threshold 0.92

```
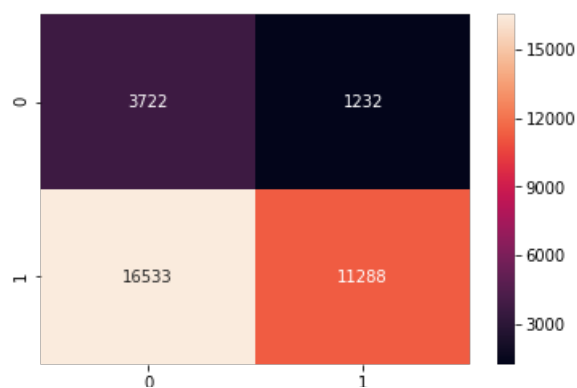print("Test confusion matrix")
test_cm = confusion_matrix(y_test, predict(y_test_pred, train_thresholds, test_fpr, test_tpr))
sns.heatmap(test_cm, annot=True, fmt="d");
```

Test confusion matrix
the maximum value of tpr*(1-fpr) 0.34686040148497044 for threshold 0.955

# 3. Conclusion

In [101]:

```python
# Please compare all your models using Prettytable library
from prettytable import PrettyTable

table = PrettyTable();
table.field_names = ['Vectorizer', 'Model', 'Hyper parameter', 'AUC'];

table.add_row(['BOW', 'Brute', set1_alpha, set1_auc]);
table.add_row(['TFIDF', 'Brute', set2_alpha, set2_auc]);
table.add_row(['W2V', 'Brute', set3_alpha, set3_auc]);
table.add_row(['TFIDFW2V', 'Brute', set4_alpha, set4_auc]);
table.add_row(['Data containing counts', 'Brute', set5_alpha, set5_auc]);

print(table)
```

```
+------------------------+-------+-----------------+--------------------+
|       Vectorizer       | Model | Hyper parameter |        AUC         |
+------------------------+-------+-----------------+--------------------+
|          BOW           | Brute |       0.01      | 0.7281684317691781 |
|         TFIDF          | Brute |      0.001      | 0.6811503291189769 |
|          W2V           | Brute |      0.001      | 0.6749103324576977 |
|        TFIDFW2V        | Brute |      0.001      | 0.6900659171019438 |
| Data containing counts | Brute |       0.01      | 0.6221047337383806 |
+------------------------+-------+-----------------+--------------------+
```

In [ ]:

In [ ]: