

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

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

Feature	Description
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

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

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

Feature	Description
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description 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 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

```

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

In [5]:

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

1.2 preprocessing of project_subject_categories

In [6]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [7]:

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

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

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=>
            "Math", "&", "Science"
            j=j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science"=>
            "Math&Science"
```

```

temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
temp = temp.replace('&','_')
sub_cat_list.append(temp.strip())

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

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

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

```

1.3 Text preprocessing

In [8]:

```

# 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 [9]:

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

Out[9]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_c
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Gra

In [10]:

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

In [11]:

```

# 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 s

chool. \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 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\n\nannan

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\n\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\n\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 a lot 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 learners 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 them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\n\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 d

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. The cart will allow 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 [12]:

```
# 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 [13]:

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

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. 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 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 them because they develop their core, 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. 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

In [14]:

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

My 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. 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 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 them because they develop their core, 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. 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

In [15]:

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

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 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 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 They want to be able to move as they learn or so they say Wobble chairs are the answer and I love them because they develop their core 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 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

In [16]:

```
# 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 [17]:

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

100% | 109248/109248 [02:06<00:00, 865.14it/s]


```
# after preprocessing
preprocessed_essays[20000]
```

In [19]:

Out[19]:

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

In [20]:

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248 [00:05<00:00]
0, 20200.00it/s]
```

```
# after preprocessing
preprocessed_title[20000]
```

```
'need move input'
```

```
# Updating dataframe for clean project title and remove old project title
project_data['clean_project_title'] = preprocessed_title
project_data.drop(['project_title'], axis=1, inplace=True)
project_data.head(2)
```

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades 1-5
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grades 1-5

```
# similarly you can preprocess the titles also
# Combining all the above students
from tqdm import tqdm
preprocessed_grade = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_grade_category'].values):
    sent = decontracted(sentence)
    sent = sent.replace(' ', '_')
    sent = sent.replace('-', '_')
    sent = sent.lower()
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_grade.append(sent.strip())
```

```
preprocessed grade[:10]
```

```
['grades_prek_2',
'grades_6_8',
'grades_6_8',
'grades_prek_2',
'grades_prek_2',
'grades_3_5',
'grades_6_8',
'grades_3_5',
'grades_prek_2',
'grades_prek_2']
```

In [25]:

```
# Updating dataframe for clean project title and remove old project title
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data['project_grade_category'] = preprocessed_grade
project_data.head(2)
```

Out[25]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_essay_1
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	My students are English learners that are work...
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Our students arrive to our school eager to lea...

In [26]:

```
# remove unnecessary column: https://cmdlinetips.com/2018/04/how-to-drop-one-or-more-columns-in-pandas-dataframe/
project_data = project_data.drop(['Unnamed: 0', 'id', 'teacher_id', 'project_submitted_datetime', \
                                  'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4', \
                                  'project_resource_summary'], axis=1)
```

In [27]:

```
project_data
```

Out[27]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_cat
0	Mrs.	IN	0	0	154.60	23	Literacy_Lai
1	Mr.	FL	7	1	299.00	1	History_Health
2	Ms.	AZ	1	0	516.85	22	Health
3	Mrs.	KY	4	1	232.90	4	Literacy_Lai Math_S
4	Mrs.	TX	1	1	67.98	4	Math_S
5	Mrs.	FL	1	1	113.22	11	Literacy_Lai Specia
6	Mrs.	CT	1	1	159.99	3	Literacy_Lai Specia

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_cat
7	Ms.	GA	7	1	229.00	4	Math_S
8	Mrs.	SC	28	1	241.98	6	Health_
9	Ms.	NC	36	1	125.36	14	Literacy_Lai
10	Mrs.	CA	37	1	100.21	10	Literacy_Lai
11	Ms.	CA	32	1	431.77	8	Literacy_Lai AppliedLo
12	Mrs.	NY	5	0	219.46	22	Math_S
13	Mrs.	OK	30	1	399.99	1	Specia
14	Ms.	MA	15	0	91.94	10	Literacy_Lai
15	Ms.	TX	3	1	435.84	24	Health_
16	Mrs.	FL	1	1	298.43	7	Literacy_Lai Specia
17	Ms.	NV	0	1	158.63	12	Math_S Literacy_Lai
18	Mrs.	GA	0	1	59.98	4	AppliedLo
19	Ms.	OH	9	1	749.42	7	Health_
20	Mrs.	PA	23	1	213.85	1	Literacy_Lai
21	Mrs.	NC	0	1	250.91	4	Math_S Specia

22	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	27009	price	quantity	clean_cat
23	Mr.	AL	2	1	299.98	2	Mus	
24	Mrs.	FL	0	1	250.00	6	Math_5	
25	Mrs.	AL	11	0	268.99	2	Math_5	
26	Ms.	TX	2	1	280.83	4	Literacy_Lai Math_5	
27	Teacher	LA	2	1	660.84	7	Literacy_Lai Math_5	
28	Mrs.	GA	5	0	129.98	3	Literacy_Lai Specia	
29	Mrs.	VA	0	1	86.74	53	Literacy_Lai AppliedL	
...	
109218	Mrs.	IL	4	0	747.00	3	Literacy_Lai	
109219	Teacher	FL	0	0	300.18	14	Literacy_Lai History	
109220	Mrs.	WI	3	1	121.59	14	Literacy_Lai	
109221	Teacher	NY	1	1	289.52	32	Literacy_Lai	
109222	Ms.	NC	34	1	241.08	40	Literacy_Lai Math_5	
109223	Ms.	GA	12	1	692.17	46	Care_	
109224	Ms.	NY	7	1	915.27	7	Math_5 Literacy_Lai	

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_cat
109225	Mrs.	NC	1	0	737.95	2	AppliedL Literacy_Lai
109226	Ms.	CA	47	1	379.96	9	Health_
109227	Mrs.	NY	0	1	428.24	5	Literacy_Lai
109228	Mrs.	LA	8	1	159.43	4	Literacy_Lai Math_§
109229	Mrs.	CO	0	1	688.00	2	Literacy_Lai Math_§
109230	Ms.	NY	0	1	309.60	4	Specia
109231	Mrs.	AZ	7	1	155.70	5	Health_
109232	Mrs.	MD	0	1	43.20	20	Literacy_Lai
109233	Ms.	AZ	1	1	490.05	1	Math_§ History
109234	Ms.	NY	9	1	273.72	6	AppliedL
109235	Mrs.	TX	1	1	11.86	24	AppliedL Literacy_Lai
109236	Mrs.	OH	6	1	269.00	1	Specia
109237	Mrs.	IN	4	1	30.76	30	Literacy_Lai
109238	Mrs.	WI	41	1	267.50	12	Health_
109239	Mrs.	MN	6	1	178.98	2	Literacy_Lai Math_§

109240	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_data
109241	Mrs.	MD	0	1	659.00	1	AppliedL Literacy_Lai
109242	Mrs.	SC	26	1	592.16	2	Math_5 Literacy_Lai
109243	Mr.	MO	0	1	59.98	8	Literacy_Lai Math_5
109244	Ms.	NJ	0	1	846.32	4	Literacy_Lai Math_5
109245	Mrs.	NJ	3	1	239.96	4	Literacy_Lai Math_5
109246	Mrs.	NY	0	1	73.05	16	Health_ Specia
109247	Ms.	VA	0	1	109.90	5	AppliedL Math_5

109248 rows × 11 columns

◀		▶
---	--	---

Check whether each column contain NaN or Not

In [28]:

```
project_data['teacher_prefix'].isnull().values.any()
```

Out[28]:

True

In [29]:

```
project_data['school_state'].isnull().values.any()
```

Out[29]:

False

In [30]:

```
project_data['teacher_number_of_previously_posted_projects'].isnull().values.any()
```

Out[30]:

False

In [31]:

```
project_data['project_is_approved'].isnull().values.any()
```

Out[31]:

False

In [32]:

```
project_data['price'].isnull().values.any()
```

Out[32]:

False

In [33]:

```
project_data['quantity'].isnull().values.any()
```

Out[33]:

False

In [34]:

```
project_data['clean_categories'].isnull().values.any()
```

Out[34]:

False

In [35]:

```
project_data['clean_subcategories'].isnull().values.any()
```

Out[35]:

False

In [36]:

```
project_data['clean_essay'].isnull().values.any()
```

Out[36]:

False

In [37]:

```
project_data['clean_project_title'].isnull().values.any()
```

Out[37]:

False

In [38]:

```
project_data['project_grade_category'].isnull().values.any()
```

Out[38]:

False

Since we got 'teacher_prefix' attributes which contain NaN. Let check how many NaN are contain in this attributes

In [39]:

```
project_data['teacher_prefix'].isnull().sum().sum()
```

Out[39]:

3

1.5 Preparing data for models

In [40]:

```
project_data.columns
```

Out[40]:

```
Index(['teacher_prefix', 'school_state',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'price', 'quantity', 'clean_categories', 'clean_subcategories',  
      'clean_essay', 'clean_project_title', 'project_grade_category'],  
      dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)
- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

1.5.1 Vectorizing Categorical data

- <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

In [0]:

```
# we use count vectorizer to convert the values into one  
from sklearn.feature_extraction.text import CountVectorizer  
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)  
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)  
print(vectorizer.get_feature_names())  
print("Shape of matrix after one hot encodig ",categories_one_hot.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sp  
orts', 'Math_Science', 'Literacy_Language']  
Shape of matrix after one hot encodig (109248, 9)
```

In [0]:

```
# we use count vectorizer to convert the values into one  
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
```

```
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", sub_categories_one_hot.shape)
```

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

Shape of matrix after one hot encoding (109248, 30)

In [0]:

```
# you can do the similar thing with state, teacher_prefix and project_grade_category also
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [0]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encoding ", text_bow.shape)
```

Shape of matrix after one hot encoding (109248, 16623)

In [0]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
```

1.5.2.2 TFIDF vectorizer

In [0]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
text_tfidf = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encoding ", text_tfidf.shape)
```

Shape of matrix after one hot encoding (109248, 16623)

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

```
'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile, 'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.", len(model), " words loaded!")
    return model
model = loadGloveModel('glove_42B_300d.txt')
```

```

model = loadGloveModel( glove.42B.300d.txt )

# =====
Output:

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

# =====

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

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

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

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

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

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

'''

```

Out[0]:

```

'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\n
def loadGloveModel(
gloveFile):\n
    print ("Loading Glove Model")\n
    f = open(gloveFile,\r', encoding="utf8")\n
    model = {}\n
    for line in tqdm(f):\n
        splitLine = line.split()\n
        word = splitLine[0]\n
        embedding = np.array([float(val) for val in splitLine[1:]])\n
        model[word] = embedding\n
    print ("Done.",len(model)," words loaded!")\n
    return model\n
model = loadGloveModel(\'glove.42B.300d.txt\')\n
\n\n# =====\n
Output:\n
\nLoading Glove Model\n
1917495it [06:32, 4879.69it/s]\n
Done. 1917495 words loaded!\n
\n# =====\n
\n\nwords = []\n
for i in preprocod_texts:\n
    words.extend(i.split(\' \'))\n
\nfor i in preprocod_titles:\n
    words.extend(i.split(\' \'))\n
\nprint("all the words in the coupus", len(words))\n
\nwords = set(words)\n
\nprint("the unique words in the coupus", len(words))\n
\n\ninter_words = set(model.keys()).intersection(words)\n
\nprint("The number of words that are present in both glove vectors and our coupus", len(inter_words), "(" , np.round(len(inter_words)/len(words)*100,3) , "%) ")
\n\nwords_courpus = {}\n
words_glove = set(model.keys())\n
for i in words:\n
    if i in words_glove:\n
        words_courpus[i] = model[i]\n
\nprint("word 2 vec length", len(words_courpus))\n
\n\n# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n
\nimport pickle\n
with open(\'glove_vectors\', \'wb\') as f:\n
    pickle.dump(words_courpus, f)\n
\n\n\n'

```

In [0]:

```

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

```

In [0]:

```

# average Word2Vec
# compute average word2vec for each word

```

```
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors.append(vector)

print(len(avg_w2v_vectors))
print(len(avg_w2v_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248 [01:00<00:00, 1806.88it/s]
```

```
109248
300
```

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

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

In [0]:

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

print(len(tfidf_w2v_vectors))
print(len(tfidf_w2v_vectors[0]))
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248 [08:03<00:00, 225.81it/s]
```

```
109248
300
```

In [0]:

```
# Similarly you can vectorize for title also
```

1.5.3 Vectorizing Numerical features

In [0]:

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

In [0]:

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

# price_standardized = standardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

# Now standardize the data with above mean and variance.
price_standardized = price_scalar.transform(project_data['price'].values.reshape(-1, 1))
```

In [0]:

```
price_standardized
```

Out[0]:

```
array([[0.00098843, 0.00191166, 0.00330448, ..., 0.00153418, 0.00046704,
        0.00070265]])
```

1.5.4 Merging all the above features

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

In [0]:

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 1)
```

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

Out[0]:

```
(109248, 16663)
```

In [0]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Assignment 4: Naive Bayes

1. Apply Multinomial NaiveBayes on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)

2. The hyper paramter tuning(find best Alpha)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- 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. Feature importance

- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using values of `feature_log_prob_` parameter of [MultinomialNB](#) and print their corresponding feature names

4. 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. Here on X-axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.
- 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](#).

5. 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](#)

2. Naive Bayes

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

In [41]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [42]:

```
# Combine the train.csv and resource.csv
```

```

# Combine the trainset and testset
# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
from sklearn.model_selection import train_test_split

# https://www.geeksforgeeks.org/python-pandas-dataframe-sample/
# Take 50k dataset
project_data = project_data.sample(n=50000)
# Remove that row which contain NaN. We observed that only 3 rows that contain NaN
project_data = project_data[pd.notnull(project_data['teacher_prefix'])]
project_data.shape

```

Out[42]:

(49998, 11)

In [43]:

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

Out[43]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_category
85924	Mrs.	TX	0	1	86.90	8	Literacy_Lang
30192	Mr.	NY	3	1	79.99	3	Literacy_Lang

In [114]:

```

# Split train and test
tr_X, ts_X, tr_y, ts_y, = train_test_split(project_data, project_data['project_is_approved'].values, te
st_size=0.3, random_state=1, stratify=project_data['project_is_approved'].values)
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)

# After train data, We are going to perform KFold Cross validation at the time of training model

# Reset index of df
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)
tr_X.drop(['project_is_approved'], axis=1, inplace=True)
ts_X.drop(['project_is_approved'], axis=1, inplace=True)

print('Shape of train data:', tr_X.shape)
print('Shape of test data:', ts_X.shape)

```

Shape of train data: (34998, 10)

Shape of test data: (15000, 10)

In [115]:

```
tr_X.head(2)
```

Out[115]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	price	quantity	clean_categories	clean_subcategorie
0	Mrs.	AZ	1	34.61	7	Literacy_Language Math_Science	LiteracyMathematic

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	price	quantity	clean_categories	clean_subcategory
1	Mrs.	TX	3	310.56	11	Math_Science	Health_LifeScienc Mathematic

In [116]:

```
ts_X.head(2)
```

Out[116]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	price	quantity	clean_categories	clean_subcategory
0	Ms.	CT	7	334.44	12	Math_Science Literacy_Language	EnvironmentalScienc Literature_Writin
1	Ms.	TX	0	12.40	30	Literacy_Language	Litera

In []:

In [117]:

```
print('Shape of Train Data',[tr_X.shape, tr_y.shape])
print('Shape of Test Data',[ts_X.shape, ts_y.shape])
```

```
Shape of Train Data [(34998, 10), (34998,)]
Shape of Test Data [(15000, 10), (15000,)]
```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [118]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label

# For Numerical with train data
### 1) quantity

from sklearn.preprocessing import Normalizer
# normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.Normal
izer.html

quantity_scalar = Normalizer()
quantity_scalar.fit(tr_X['quantity'].values.reshape(1,-1)) # finding the mean and standard deviation of
this data
quantity_normalized = quantity_scalar.transform(tr_X['quantity'].values.reshape(1, -1))

### 2) price
```



```
# the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1
```

```
price_scalar = Normalizer()
price_scalar.fit(tr_X['price'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
price_normalized = price_scalar.transform(tr_X['price'].values.reshape(1, -1))

### 3) For teacher_number_of_previously_projects

# We are going to represent the teacher_number_of_previously_posted_projects, as numerical values within the range 0-1

teacher_number_of_previously_posted_projects_scalar = Normalizer()
teacher_number_of_previously_posted_projects_scalar.fit(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
teacher_number_of_previously_posted_projects_normalized = teacher_number_of_previously_posted_projects_scalar.transform(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
```

In [119]:

```
print('Shape of quantity:', quantity_normalized.T.shape)
print('Shape of price:', price_normalized.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', teacher_number_of_previously_posted_projects_normalized.T.shape)
```

Shape of quantity: (34998, 1)

Shape of price: (34998, 1)

Shape of teacher_number_of_previously_posted_projects: (34998, 1)

In [120]:

```
quantity_normalized.T
```

Out[120]:

```
array([[0.00119224],
       [0.00187352],
       [0.00153288],
       ...,
       [0.0017032 ],
       [0.0017032 ],
       [0.00306577]])
```

In [121]:

```
price_normalized.T
```

Out[121]:

```
array([[0.00040346],
       [0.00362029],
       [0.00314666],
       ...,
       [0.00200517],
       [0.0021566 ],
       [0.00196764]])
```

In [122]:

```
teacher_number_of_previously_posted_projects_normalized.T
```

Out[122]:

```
array([[0.00018164],
       [0.00054493],
       [0.00181642],
       ...,
       [0.01852746],
       [0.00000000]])
```

```
[0.00090821],
[0.         ]])
```

In [123]:

```
# Transform numerical attributes for test data
ts_price = price_scalar.transform(ts_X['price'].values.reshape(1,-1))
ts_quantity = quantity_scalar.transform(ts_X['quantity'].values.reshape(1,-1))
ts_teacher_number_of_previously_posted_projects = \
teacher_number_of_previously_posted_projects_scalar.transform(ts_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
```

In [124]:

```
print('-----Test data-----')
print('Shape of quantity:', ts_quantity.T.shape)
print('Shape of price:', ts_price.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', ts_teacher_number_of_previously_posted_projects.T.shape)
```

```
-----Test data-----
Shape of quantity: (15000, 1)
Shape of price: (15000, 1)
Shape of teacher_number_of_previously_posted_projects: (15000, 1)
```

In [125]:

```
# For categorical with train data
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# One hot encoding for school state

### 1) school_state
print('=====\\n')
# Count Vectorize with vocabulary contains unique code of school state and we are doing boolean BoW
vectorizer_school_state = CountVectorizer(vocabulary=tr_X['school_state'].unique(), lowercase=False, binary=True)
vectorizer_school_state.fit(tr_X['school_state'].values)
print('List of feature in school_state',vectorizer_school_state.get_feature_names())

school_state_one_hot = vectorizer_school_state.transform(tr_X['school_state'].values)
print("\\nShape of school_state matrix after one hot encoding ",school_state_one_hot.shape)

### 2) project_subject_categories
print('=====\\n')
vectorizer_categories = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_categories.fit(tr_X['clean_categories'].values)
print('List of features in project_subject_categories',vectorizer_categories.get_feature_names())

categories_one_hot = vectorizer_categories.transform(tr_X['clean_categories'].values)
print("\\nShape of project_subject_categories matrix after one hot encoding ",categories_one_hot.shape)

### 3) project_subject_subcategories
print('=====\\n')
vectorizer_subcategories = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_subcategories.fit(tr_X['clean_subcategories'].values)
print('List of features in project_subject_subcategories',vectorizer_subcategories.get_feature_names())

subcategories_one_hot = vectorizer_subcategories.transform(tr_X['clean_subcategories'].values)
print("\\nShape of project_subject_subcategories matrix after one hot encoding ",subcategories_one_hot.shape)

### 4) project_grade_category
print('=====\\n')
# One hot encoding for project_grade_category

# Count Vectorize with vocabulary contains unique code of project grade category and we are doing bool
```

```

# Some of the data is filled with nan. So we update the nan to 'None' as a string
vectorizer_grade_category = CountVectorizer(vocabulary=tr_X['project_grade_category'].unique(), lowercase=False, binary=True)
vectorizer_grade_category.fit(tr_X['project_grade_category'].values)
print('List of features in project_grade_category',vectorizer_grade_category.get_feature_names())

project_grade_category_one_hot = vectorizer_grade_category.transform(tr_X['project_grade_category'].values)
print("\nShape of project_grade_category matrix after one hot encoding ",project_grade_category_one_hot.shape)

### 5) teacher_prefix
print('=====\\n')
# One hot encoding for teacher_prefix

# Count Vectorize with vocabulary contains unique code of teacher_prefix and we are doing boolean BoW
# Since some of the data is filled with nan. So we update the nan to 'None' as a string
# tr_X['teacher_prefix'] = tr_X['teacher_prefix'].fillna('None')
vectorizer_teacher_prefix = CountVectorizer(vocabulary=tr_X['teacher_prefix'].unique(), lowercase=False, binary=True)
vectorizer_teacher_prefix.fit(tr_X['teacher_prefix'].values)
print('List of features in teacher_prefix',vectorizer_teacher_prefix.get_feature_names())

teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(tr_X['teacher_prefix'].values)
print("\nShape of teacher_prefix matrix after one hot encoding ",teacher_prefix_one_hot.shape)

```

```

=====

List of feature in school_state ['AZ', 'TX', 'UT', 'SC', 'PA', 'FL', 'GA', 'IA', 'VA', 'NC', 'OH', 'IN', 'IL', 'KY', 'OK', 'MA', 'MN', 'KS', 'AL', 'TN', 'MD', 'MO', 'CA', 'MI', 'ID', 'CT', 'NY', 'NJ', 'OR', 'HI', 'NE', 'WV', 'DE', 'WI', 'ME', 'MS', 'LA', 'WA', 'AR', 'NV', 'DC', 'NM', 'SD', 'RI', 'CO', 'WY', 'MT', 'ND', 'NH', 'AK', 'VT']

```

```

Shape of school_state matrix after one hot encoding (34998, 51)
=====

```

```

List of features in project_subject_categories ['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']

```

```

Shape of project_subject_categories matrix after one hot encoding (34998, 9)
=====

```

```

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

```

```

Shape of project_subject_subcategories matrix after one hot encoding (34998, 30)
=====

```

```

List of features in project_grade_category ['grades_3_5', 'grades_6_8', 'grades_prek_2', 'grades_9_12']

```

```

Shape of project_grade_category matrix after one hot encoding (34998, 4)
=====

```

```

List of features in teacher_prefix ['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.']

```

```

Shape of teacher_prefix matrix after one hot encoding (34998, 5)

```

In [126]:

```

list_features = ['quantity','price','teacher_number_of_previously_posted_projects'] # storing all features names for further print feature importance
list_features

```

Out[126]:

```

['quantity', 'price', 'teacher_number_of_previously_posted_projects']

```

In [127]:

```
vectorizer_school_state.get_feature_names()[0], len(vectorizer_school_state.get_feature_names())
```

Out[127]:

```
('AZ', 51)
```

In [128]:

```
for i in range(len(vectorizer_school_state.get_feature_names())):
    list_features.append(vectorizer_school_state.get_feature_names()[i])
for i in range(len(vectorizer_categories.get_feature_names())):
    list_features.append(vectorizer_categories.get_feature_names()[i])
for i in range(len(vectorizer_subcategories.get_feature_names())):
    list_features.append(vectorizer_subcategories.get_feature_names()[i])
for i in range(len(vectorizer_grade_category.get_feature_names())):
    list_features.append(vectorizer_grade_category.get_feature_names()[i])
for i in range(len(vectorizer_teacher_prefix.get_feature_names())):
    list_features.append(vectorizer_teacher_prefix.get_feature_names()[i])

list_features
```

Out[128]:

```
['quantity',
 'price',
 'teacher_number_of_previously_posted_projects',
 'AZ',
 'TX',
 'UT',
 'SC',
 'PA',
 'FL',
 'GA',
 'IA',
 'VA',
 'NC',
 'OH',
 'IN',
 'IL',
 'KY',
 'OK',
 'MA',
 'MN',
 'KS',
 'AL',
 'TN',
 'MD',
 'MO',
 'CA',
 'MI',
 'ID',
 'CT',
 'NY',
 'NJ',
 'OR',
 'HI',
 'NE',
 'WV',
 'DE',
 'WI',
 'ME',
 'MS',
 'LA',
 'WA',
 'AR',
 'NV',
 'DC',
 'NM',
 'SD',
 'RI',
 'CO',
 'WY']
```

```

'WI',
'MT',
'ND',
'NH',
'AK',
'VT',
'Warmth',
'Care_Hunger',
'History_Civics',
'Music_Arts',
'AppliedLearning',
'SpecialNeeds',
'Health_Sports',
'Math_Science',
'Literacy_Language',
'Economics',
'CommunityService',
'FinancialLiteracy',
'ParentInvolvement',
'Extracurricular',
'Civics_Government',
'ForeignLanguages',
'NutritionEducation',
'Warmth',
'Care_Hunger',
'SocialSciences',
'PerformingArts',
'CharacterEducation',
'TeamSports',
'Other',
'College_CareerPrep',
'Music',
'History_Geography',
'Health_LifeScience',
'EarlyDevelopment',
'ESL',
'Gym_Fitness',
'EnvironmentalScience',
'VisualArts',
'Health_Wellness',
'AppliedSciences',
'SpecialNeeds',
'Literature_Writing',
'Mathematics',
'Literacy',
'grades_3_5',
'grades_6_8',
'grades_prek_2',
'grades_9_12',
'Mrs.',
'Ms.',
'Mr.',
'Teacher',
'Dr.']

```

In [129]:

```

# Transform categorical for test data
ts_school_state = vectorizer_school_state.transform(ts_X['school_state'].values)
ts_project_subject_category = vectorizer_categories.transform(ts_X['clean_categories'].values)
ts_project_subject_subcategory = vectorizer_subcategories.transform(ts_X['clean_subcategories'].values)
ts_project_grade_category = vectorizer_grade_category.transform(ts_X['project_grade_category'].values)
ts_teacher_prefix = vectorizer_teacher_prefix.transform(ts_X['teacher_prefix'].values)

```

In [130]:

```

print('-----Test data-----')
print('Shape of school_state:', ts_school_state.shape)
print('Shape of project_subject_categories:', ts_project_subject_category.shape)
print('Shape of project_subject_subcategories:', ts_project_subject_subcategory.shape)
print('Shape of project_grade_category:', ts_project_grade_category.shape)
print('Shape of teacher_prefix:', ts_teacher_prefix.shape)

```

-----Test data-----

```
Shape of school_state: (15000, 51)
Shape of project_subject_categories: (15000, 9)
Shape of project_subject_subcategories: (15000, 30)
Shape of project_grade_category: (15000, 4)
Shape of teacher_prefix: (15000, 5)
```

2.3 Make Data Model Ready: encoding eassay, and project_title

In [131]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Note:

We already have preprocessed both essay and project_title in Text processing section (1.3 and 1.4) above

2.4 Appling NB() on different kind of featurization as mentioned in the instructions

Apply Naive Bayes 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 instructions

BoW

In [62]:

```
### BoW in Essay and Title on Train

# # We are considering only the words which appeared in at least 10 documents with max feature = 5000(r
ows or projects).
vectorizer_bow = CountVectorizer(min_df=10, ngram_range=(1,1), max_features=5000)
text_bow = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on train",text_bow.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = CountVectorizer(min_df=10, ngram_range=(1,1), max_features=5000)
title_bow = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig ",title_bow.shape)

### BoW in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on test",ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig on test",ts_title.shape)
```

```
Shape of essay matrix after one hot encodig on train (34998, 5000)
Shape of title matrix after one hot encodig (34998, 1586)
```

```
=====
Shape of essay matrix after one hot encodig on test (15000, 5000)
Shape of title matrix after one hot encodig on test (15000, 1586)
```

In [63]:

```
for i in tqdm(range(len(vectorizer_bow.get_feature_names()))):  
    list_features.append(vectorizer_bow.get_feature_names()[i])  
for i in tqdm(range(len(vectorizer_bowt.get_feature_names()))):  
    list_features.append(vectorizer_bowt.get_feature_names()[i])
```

list_features

```
100% |████████████████████████████████████████████████████████████████████████████████| 5000/5000 [00:27<00  
:00, 185.13it/s]  
100% |████████████████████████████████████████████████████████████████████████████████| 1586/1586 [00:02<00  
:00, 701.45it/s]
```

Out[63]:

```
['quantity',  
'price',  
'teacher_number_of_previously_posted_projects',  
'AZ',  
'TX',  
'UT',  
'SC',  
'PA',  
'FL',  
'GA',  
'IA',  
'VA',  
'NC',  
'OH',  
'IN',  
'IL',  
'KY',  
'OK',  
'MA',  
'MN',  
'KS',  
'AL',  
'TN',  
'MD',  
'MO',  
'CA',  
'MI',  
'ID',  
'CT',  
'NY',  
'NJ',  
'OR',  
'HI',  
'NE',  
'WV',  
'DE',  
'WI',  
'ME',  
'MS',  
'LA',  
'WA',  
'AR',  
'NV',  
'DC',  
'NM',  
'SD',  
'RI',  
'CO',  
'WY',  
'MT',  
'ND',  
'NH',  
'AK',  
'VT',  
'Warmth',  
'Care_Hunger',  
'History_Civics',  
'Music_Arts',  
'AppliedLearning',
```

'SpecialNeeds',
'Health_Sports',
'Math_Science',
'Literacy_Language',
'Economics',
'CommunityService',
'FinancialLiteracy',
'ParentInvolvement',
'Extracurricular',
'Civics_Government',
'ForeignLanguages',
'NutritionEducation',
'Warmth',
'Care_Hunger',
'SocialSciences',
'PerformingArts',
'CharacterEducation',
'TeamSports',
'Other',
'College_CareerPrep',
'Music',
'History_Geography',
'Health_LifeScience',
'EarlyDevelopment',
'ESL',
'Gym_Fitness',
'EnvironmentalScience',
'VisualArts',
'Health_Wellness',
'AppliedSciences',
'SpecialNeeds',
'Literature_Writing',
'Mathematics',
'Literacy',
'grades_3_5',
'grades_6_8',
'grades_prek_2',
'grades_9_12',
'Mrs.',
'Ms.',
'Mr.',
'Teacher',
'Dr.',
'00',
'000',
'10',
'100',
'1000',
'10th',
'11',
'110',
'11th',
'12',
'120',
'12th',
'13',
'14',
'15',
'150',
'16',
'17',
'18',
'180',
'19',
'1st',
'20',
'200',
'2015',
'2016',
'2017',
'21',
'21st',
'22',
'23',
'24',
'25',
'250'.

'26',
'27',
'28',
'29',
'2nd',
'30',
'300',
'31',
'32',
'33',
'34',
'35',
'36',
'3d',
'3doodler',
'3rd',
'40',
'400',
'42',
'44',
'45',
'450',
'48',
'4th',
'50',
'500',
'55',
'5th',
'60',
'600',
'65',
'6th',
'70',
'700',
'74',
'75',
'7th',
'80',
'800',
'84',
'85',
'8th',
'90',
'900',
'92',
'94',
'95',
'96',
'97',
'98',
'99',
'9th',
'abc',
'abilities',
'ability',
'able',
'absent',
'absolute',
'absolutely',
'absorb',
'abstract',
'abundance',
'abuse',
'academic',
'academically',
'academics',
'academy',
'accelerated',
'accept',
'acceptance',
'accepted',
'accepting',
'access',
'accessed',
'accessibility',
'accessible',
'accessing'.

'accessing',
'accessories',
'accommodate',
'accommodations',
'accompany',
'accomplish',
'accomplished',
'accomplishing',
'accomplishment',
'accomplishments',
'according',
'account',
'accountability',
'accountable',
'accounts',
'accuracy',
'accurate',
'accurately',
'accustomed',
'achieve',
'achieved',
'achievement',
'achievements',
'achievers',
'achieving',
'acquire',
'acquired',
'acquiring',
'acquisition',
'across',
'act',
'acting',
'action',
'actions',
'activate',
'active',
'actively',
'activities',
'activity',
'actual',
'actually',
'adapt',
'adaptations',
'adapted',
'adaptive',
'add',
'added',
'adding',
'addition',
'additional',
'additionally',
'additions',
'address',
'addressed',
'addressing',
'adds',
'adequate',
'adhd',
'adjust',
'administration',
'administrators',
'adolescents',
'adopted',
'adorable',
'adore',
'adult',
'adulthood',
'adults',
'advance',
'advanced',
'advancement',
'advances',
'advancing',
'advantage',
'advantages',
'adventure',
'adventures',
'adventurous'

'adventurous',
'adversity',
'advocate',
'affect',
'affected',
'affects',
'affluent',
'afford',
'affordable',
'afforded',
'afraid',
'africa',
'african',
'afternoon',
'age',
'aged',
'ages',
'agility',
'ago',
'agree',
'agreed',
'agricultural',
'agriculture',
'ahead',
'aid',
'aide',
'aides',
'aids',
'aim',
'aims',
'air',
'alabama',
'alaska',
'albert',
'alert',
'algebra',
'align',
'aligned',
'alike',
'alive',
'alleviate',
'allow',
'allowed',
'allowing',
'allows',
'almost',
'alone',
'along',
'alongside',
'aloud',
'alouds',
'alphabet',
'already',
'also',
'alternate',
'alternative',
'alternatives',
'although',
'always',
'amaze',
'amazed',
'amazes',
'amazing',
'amazingly',
'amazon',
'ambassadors',
'ambition',
'ambitious',
'america',
'american',
'americans',
'among',
'amongst',
'amount',
'amounts',
'ample',
'analysis',
'analytical'

analytical ,
'analyze',
'analyzing',
'anatomy',
'anchor',
'ancient',
'angeles',
'anger',
'angles',
'animal',
'animals',
'animation',
'annotate',
'announcements',
'annual',
'another',
'answer',
'answering',
'answers',
'anxiety',
'anxious',
'anymore',
'anyone',
'anything',
'anytime',
'anywhere',
'ap',
'apart',
'apartment',
'apartments',
'app',
'apparent',
'appeal',
'appealing',
'appear',
'apple',
'applicable',
'application',
'applications',
'applied',
'apply',
'applying',
'appreciate',
'appreciated',
'appreciation',
'appreciative',
'approach',
'approaches',
'approaching',
'appropriate',
'appropriately',
'approved',
'approximately',
'apps',
'april',
'aquarium',
'ar',
'arabic',
'archery',
'architects',
'architecture',
'area',
'areas',
'arise',
'arizona',
'arkansas',
'arms',
'around',
'arrangement',
'arrangements',
'array',
'arrive',
'arrived',
'art',
'article',
'articles',
'articulation',
'artifacts'

'alliacs',
'artist',
'artistic',
'artists',
'arts',
'artwork',
'artworks',
'asd',
'asia',
'asian',
'aside',
'ask',
'asked',
'asking',
'aspect',
'aspects',
'aspirations',
'aspire',
'assemblies',
'assess',
'assessed',
'assessment',
'assessments',
'asset',
'assign',
'assigned',
'assignment',
'assignments',
'assist',
'assistance',
'assistant',
'assisting',
'associated',
'association',
'assortment',
'assure',
'athlete',
'athletes',
'athletic',
'athletics',
'atlanta',
'atmosphere',
'attach',
'attached',
'attain',
'attempt',
'attempting',
'attend',
'attendance',
'attended',
'attending',
'attention',
'attentive',
'attitude',
'attitudes',
'attractive',
'attributes',
'audience',
'audio',
'audiobooks',
'auditory',
'august',
'authentic',
'author',
'authors',
'autism',
'autistic',
'availability',
'available',
'avenue',
'avenues',
'average',
'avid',
'avoid',
'await',
'award',
'awarded',
'awards'

'awards',
'aware',
'awareness',
'away',
'awe',
'awesome',
'babies',
'baby',
'baccalaureate',
'back',
'background',
'backgrounds',
'backpack',
'backpacks',
'backs',
'bad',
'bag',
'bags',
'balance',
'balanced',
'balancing',
'ball',
'balls',
'baltimore',
'band',
'bands',
'bank',
'bar',
'bare',
'barely',
'barrier',
'barriers',
'bars',
'base',
'baseball',
'based',
'bases',
'basic',
'basically',
'basics',
'basis',
'basketball',
'basketballs',
'baskets',
'bass',
'bathroom',
'batteries',
'battery',
'battle',
'bay',
'beach',
'beads',
'bean',
'beanbag',
'beanbags',
'bear',
'bears',
'beat',
'beautiful',
'beautifully',
'beauty',
'became',
'become',
'becomes',
'becoming',
'bed',
'beds',
'bee',
'bees',
'beg',
'began',
'begging',
'begin',
'beginner',
'beginning',
'begins',
'begun',
'behind',

'benaved',
'behavior',
'behavioral',
'behaviorally',
'behaviors',
'behind',
'beings',
'belief',
'beliefs',
'believe',
'believer',
'believes',
'believing',
'bell',
'bellies',
'bells',
'belong',
'belonging',
'belongings',
'belongs',
'beloved',
'bench',
'benches',
'beneficial',
'benefit',
'benefits',
'benjamin',
'besides',
'best',
'better',
'beyond',
'bi',
'big',
'bigger',
'biggest',
'bike',
'bikes',
'bilingual',
'bin',
'binder',
'binders',
'bingo',
'bins',
'biographies',
'biology',
'bird',
'birds',
'birthday',
'bit',
'bits',
'black',
'blank',
'blast',
'blend',
'blended',
'blending',
'blessed',
'blessing',
'block',
'blocks',
'blog',
'blood',
'blossom',
'blow',
'blue',
'bluetooth',
'board',
'boards',
'bodies',
'body',
'bond',
'bonds',
'bonus',
'boogie',
'book',
'books',
'bookshelf',

'bookshelves',
'boost',
'bored',
'boring',
'born',
'borrow',
'borrowed',
'boston',
'bot',
'bots',
'bottle',
'bottles',
'bottom',
'bought',
'bounce',
'bouncing',
'bouncy',
'bound',
'boundaries',
'bounds',
'bowling',
'box',
'boxes',
'boy',
'boys',
'brain',
'brainpop',
'brains',
'brainstorm',
'brainstormed',
'brainstorming',
'brand',
'brave',
'break',
'breakfast',
'breaking',
'breakout',
'breaks',
'breath',
'breathing',
'bricks',
'bridge',
'bridges',
'bright',
'brighten',
'brighter',
'brightest',
'brightly',
'brilliant',
'bring',
'bringing',
'brings',
'broad',
'broaden',
'broader',
'broke',
'broken',
'bronx',
'brooklyn',
'brothers',
'brought',
'brushes',
'bubble',
'bubbles',
'bucket',
'buckets',
'buddies',
'budding',
'buddy',
'budget',
'budgets',
'build',
'builders',
'building',
'buildings',
'builds',
'built',
'built',

'bulletin',
'bullying',
'bunch',
'bundle',
'burden',
'burn',
'burning',
'bursting',
'bus',
'business',
'businesses',
'busy',
'butterflies',
'butterfly',
'button',
'buttons',
'buy',
'buying',
'buzz',
'ca',
'cabinet',
'caddies',
'cafeteria',
'calculate',
'calculator',
'calculators',
'calculus',
'calendar',
'california',
'call',
'called',
'calm',
'calming',
'calories',
'came',
'camera',
'cameras',
'camp',
'campus',
'candy',
'cannot',
'canvas',
'capabilities',
'capability',
'capable',
'capacity',
'capture',
'car',
'card',
'cardboard',
'cardiovascular',
'cards',
'care',
'cared',
'career',
'careers',
'carefully',
'cares',
'caring',
'carolina',
'carpet',
'carpets',
'carry',
'carrying',
'cars',
'cart',
'cartridges',
'carts',
'case',
'cases',
'cat',
'catch',
'catching',
'cater',
'caucasian',
'caught',
'cause',

'caused',
'causes',
'causing',
'cd',
'cds',
'cease',
'celebrate',
'celebrated',
'celebrating',
'celebration',
'cell',
'cells',
'center',
'centered',
'centers',
'central',
'century',
'ceramics',
'cerebral',
'certain',
'certainly',
'certificates',
'certified',
'chain',
'chair',
'chairs',
'chalk',
'challenge',
'challenged',
'challenges',
'challenging',
'chance',
'chances',
'change',
'changed',
'changer',
'changes',
'changing',
'channel',
'chaos',
'chaotic',
'chapter',
'character',
'characteristics',
'characters',
'charge',
'charged',
'charging',
'charismatic',
'chart',
'charter',
'charts',
'check',
'checked',
'checking',
'cheer',
'cheerful',
'chemical',
'chemicals',
'chemistry',
'chess',
'chicago',
'chickens',
'chicks',
'child',
'childhood',
'children',
'china',
'chinese',
'chips',
'choice',
'choices',
'choir',
'choose',
'choosing',
'chore',
'chorus',

'chose',
'chosen',
'christmas',
'chrome',
'chromebook',
'chromebooks',
'circle',
'circles',
'circuit',
'circuits',
'circumstances',
'cities',
'citizen',
'citizens',
'citizenship',
'city',
'civil',
'class',
'classes',
'classic',
'classical',
'classics',
'classified',
'classmate',
'classmates',
'classroom',
'classrooms',
'classwork',
'clay',
'clean',
'cleaner',
'cleaning',
'clear',
'clearly',
'clever',
'click',
'climate',
'climb',
'climbing',
'clip',
'clipboard',
'clipboards',
'clips',
'clock',
'close',
'closed',
'closely',
'closer',
'closet',
'closing',
'clothes',
'clothing',
'cloud',
'club',
'clubs',
'clues',
'clutter',
'co',
'coach',
'coaches',
'coaching',
'coast',
'coaster',
'coats',
'code',
'coded',
'codes',
'coding',
'coffee',
'cognitive',
'coins',
'cold',
'collaborate',
'collaborating',
'collaboration',
'collaborative',
'collaboratively',

```
'collaborators',
'collar',
'colleagues',
'collect',
'collected',
'collecting',
'collection',
'college',
'colleges',
'color',
'colored',
'colorful',
'coloring',
'colors',
'com',
'combat',
'combination',
'combine',
...]
```

In [64]:

```
print('Shape of normalized essay in train data', text_bow.shape)
print('Shape of normalized title in train data', title_bow.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

```
Shape of normalized essay in train data (34998, 5000)
Shape of normalized title in train data (34998, 1586)
=====
```

```
Shape of normalized essay in test data (15000, 5000)
Shape of normalized title in test data (15000, 1586)
```

TFIDF

In [132]:

```
### TFIDF in Essay and Title on Train

# # We are considering only the words which appeared in at least 10 documents with max feature 5000 (rows or projects).
vectorizer_bow = TfidfVectorizer(min_df=10, ngram_range=(1,1), max_features=5000)
text_bow = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on train",text_bow.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = TfidfVectorizer(min_df=10, ngram_range=(1,1), max_features=5000)
title_bow = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding ",title_bow.shape)

### TFIDF in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on test",ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding on test",ts_title.shape)
```

```
Shape of essay matrix after one hot encoding on train (34998, 5000)
Shape of title matrix after one hot encoding (34998, 1586)
=====
```

```
Shape of essay matrix after one hot encoding on test (15000, 5000)
Shape of title matrix after one hot encoding on test (15000, 1586)
```

In [133]:

```

for i in tqdm(range(len(vectorizer_bow.get_feature_names()))):
    list_features.append(vectorizer_bow.get_feature_names()[i])
for i in tqdm(range(len(vectorizer_bowt.get_feature_names()))):
    list_features.append(vectorizer_bowt.get_feature_names()[i])

list_features

```

```

100% |████████████████████████████████████████████████████████████████████████████████| 5000/5000 [00:27<00
:00, 184.88it/s]
100% |████████████████████████████████████████████████████████████████████████████████| 1586/1586 [00:02<00
:00, 710.11it/s]

```

Out[133]:

```

['quantity',
 'price',
 'teacher_number_of_previously_posted_projects',
 'AZ',
 'TX',
 'UT',
 'SC',
 'PA',
 'FL',
 'GA',
 'IA',
 'VA',
 'NC',
 'OH',
 'IN',
 'IL',
 'KY',
 'OK',
 'MA',
 'MN',
 'KS',
 'AL',
 'TN',
 'MD',
 'MO',
 'CA',
 'MI',
 'ID',
 'CT',
 'NY',
 'NJ',
 'OR',
 'HI',
 'NE',
 'WV',
 'DE',
 'WI',
 'ME',
 'MS',
 'LA',
 'WA',
 'AR',
 'NV',
 'DC',
 'NM',
 'SD',
 'RI',
 'CO',
 'WY',
 'MT',
 'ND',
 'NH',
 'AK',
 'VT',
 'Warmth',
 'Care_Hunger',
 'History_Civics',
 'Music_Arts',
 'AppliedLearning',
 'SpecialNeeds',
 'Health_Sports',
 'Math_Science'.

```

'Math_Science',
'Literacy_Language',
'Economics',
'CommunityService',
'FinancialLiteracy',
'ParentInvolvement',
'Extracurricular',
'Civics_Government',
'ForeignLanguages',
'NutritionEducation',
'Warmth',
'Care_Hunger',
'SocialSciences',
'PerformingArts',
'CharacterEducation',
'TeamSports',
'Other',
'College_CareerPrep',
'Music',
'History_Geography',
'Health_LifeScience',
'EarlyDevelopment',
'ESL',
'Gym_Fitness',
'EnvironmentalScience',
'VisualArts',
'Health_Wellness',
'AppliedSciences',
'SpecialNeeds',
'Literature_Writing',
'Mathematics',
'Literacy',
'grades_3_5',
'grades_6_8',
'grades_prek_2',
'grades_9_12',
'Mrs.',
'Ms.',
'Mr.',
'Teacher',
'Dr.',
'00',
'000',
'10',
'100',
'1000',
'10th',
'11',
'110',
'11th',
'12',
'120',
'12th',
'13',
'14',
'15',
'150',
'16',
'17',
'18',
'180',
'19',
'1st',
'20',
'200',
'2015',
'2016',
'2017',
'21',
'21st',
'22',
'23',
'24',
'25',
'250',
'26',
'27',
'28'

'28',
'29',
'2nd',
'30',
'300',
'31',
'32',
'33',
'34',
'35',
'36',
'3d',
'3doodler',
'3rd',
'40',
'400',
'42',
'44',
'45',
'450',
'48',
'4th',
'50',
'500',
'55',
'5th',
'60',
'600',
'65',
'6th',
'70',
'700',
'74',
'75',
'7th',
'80',
'800',
'84',
'85',
'8th',
'90',
'900',
'92',
'94',
'95',
'96',
'97',
'98',
'99',
'9th',
'abc',
'abilities',
'ability',
'able',
'absent',
'absolute',
'absolutely',
'absorb',
'abstract',
'abundance',
'abuse',
'academic',
'academically',
'academics',
'academy',
'accelerated',
'accept',
'acceptance',
'accepted',
'accepting',
'access',
'accessed',
'accessibility',
'accessible',
'accessing',
'accessories',
'accommodate',
'accommodations'

accommodations',
'accompany',
'accomplish',
'accomplished',
'accomplishing',
'accomplishment',
'accomplishments',
'according',
'account',
'accountability',
'accountable',
'accounts',
'accuracy',
'accurate',
'accurately',
'accustomed',
'achieve',
'achieved',
'achievement',
'achievements',
'achievers',
'achieving',
'acquire',
'acquired',
'acquiring',
'acquisition',
'across',
'act',
'acting',
'action',
'actions',
'activate',
'active',
'actively',
'activities',
'activity',
'actual',
'actually',
'adapt',
'adaptations',
'adapted',
'adaptive',
'add',
'added',
'adding',
'addition',
'additional',
'additionally',
'additions',
'address',
'addressed',
'addressing',
'adds',
'adequate',
'adhd',
'adjust',
'administration',
'administrators',
'adolescents',
'adopted',
'adorable',
'adore',
'adult',
'adulthood',
'adults',
'advance',
'advanced',
'advancement',
'advances',
'advancing',
'advantage',
'advantages',
'adventure',
'adventures',
'adventurous',
'adversity',
'advocate',
'affect'

'affect',
'affected',
'affects',
'affluent',
'afford',
'affordable',
'afforded',
'afraid',
'africa',
'african',
'afternoon',
'age',
'aged',
'ages',
'agility',
'ago',
'agree',
'agreed',
'agricultural',
'agriculture',
'ahead',
'aid',
'aide',
'aides',
'aids',
'aim',
'aims',
'air',
'alabama',
'alaska',
'albert',
'alert',
'algebra',
'align',
'aligned',
'alike',
'alive',
'alleviate',
'allow',
'allowed',
'allowing',
'allows',
'almost',
'alone',
'along',
'alongside',
'aloud',
'alouds',
'alphabet',
'already',
'also',
'alternate',
'alternative',
'alternatives',
'although',
'always',
'amaze',
'amazed',
'amazes',
'amazing',
'amazingly',
'amazon',
'ambassadors',
'ambition',
'ambitious',
'america',
'american',
'americans',
'among',
'amongst',
'amount',
'amounts',
'ample',
'analysis',
'analytical',
'analyze',
'analyzing',
'anatomy'

'anatomy',
'anchor',
'ancient',
'angeles',
'anger',
'angles',
'animal',
'animals',
'animation',
'annotate',
'announcements',
'annual',
'another',
'answer',
'answering',
'answers',
'anxiety',
'anxious',
'anymore',
'anyone',
'anything',
'anytime',
'anywhere',
'ap',
'apart',
'apartment',
'apartments',
'app',
'apparent',
'appeal',
'appealing',
'appear',
'apple',
'applicable',
'application',
'applications',
'applied',
'apply',
'applying',
'appreciate',
'appreciated',
'appreciation',
'appreciative',
'approach',
'approaches',
'approaching',
'appropriate',
'appropriately',
'approved',
'approximately',
'apps',
'april',
'aquarium',
'ar',
'arabic',
'archery',
'architects',
'architecture',
'area',
'areas',
'arise',
'arizona',
'arkansas',
'arms',
'around',
'arrangement',
'arrangements',
'array',
'arrive',
'arrived',
'art',
'article',
'articles',
'articulation',
'artifacts',
'artist',
'artistic',
'artistic'

'artists',
'arts',
'artwork',
'artworks',
'asd',
'asia',
'asian',
'aside',
'ask',
'asked',
'asking',
'aspect',
'aspects',
'aspirations',
'aspire',
'assemblies',
'assess',
'assessed',
'assessment',
'assessments',
'asset',
'assign',
'assigned',
'assignment',
'assignments',
'assist',
'assistance',
'assistant',
'assisting',
'associated',
'association',
'assortment',
'assure',
'athlete',
'athletes',
'athletic',
'athletics',
'atlanta',
'atmosphere',
'attach',
'attached',
'attain',
'attempt',
'attempting',
'attend',
'attendance',
'attended',
'attending',
'attention',
'attentive',
'attitude',
'attitudes',
'attractive',
'attributes',
'audience',
'audio',
'audiobooks',
'auditory',
'august',
'authentic',
'author',
'authors',
'autism',
'autistic',
'availability',
'available',
'avenue',
'avenues',
'average',
'avid',
'avoid',
'await',
'award',
'awarded',
'awards',
'aware',
'awareness',
'

away',
'awe',
'awesome',
'babies',
'baby',
'baccalaureate',
'back',
'background',
'backgrounds',
'backpack',
'backpacks',
'backs',
'bad',
'bag',
'bags',
'balance',
'balanced',
'balancing',
'ball',
'balls',
'baltimore',
'band',
'bands',
'bank',
'bar',
'bare',
'barely',
'barrier',
'barriers',
'bars',
'base',
'baseball',
'based',
'bases',
'basic',
'basically',
'basics',
'basis',
'basketball',
'basketballs',
'baskets',
'bass',
'bathroom',
'batteries',
'battery',
'battle',
'bay',
'beach',
'beads',
'bean',
'beanbag',
'beanbags',
'bear',
'bears',
'beat',
'beautiful',
'beautifully',
'beauty',
'became',
'become',
'becomes',
'becoming',
'bed',
'beds',
'bee',
'bees',
'beg',
'began',
'begging',
'begin',
'beginner',
'beginning',
'begins',
'begun',
'behaved',
'behavior',
'behavioral',

'behaviorally',
'behaviors',
'behind',
'beings',
'belief',
'beliefs',
'believe',
'believer',
'believes',
'believing',
'bell',
'bellies',
'bells',
'belong',
'belonging',
'belongings',
'belongs',
'beloved',
'bench',
'benches',
'beneficial',
'benefit',
'benefits',
'benjamin',
'besides',
'best',
'better',
'beyond',
'bi',
'big',
'bigger',
'biggest',
'bike',
'bikes',
'bilingual',
'bin',
'binder',
'binders',
'bingo',
'bins',
'biographies',
'biology',
'bird',
'birds',
'birthday',
'bit',
'bits',
'black',
'blank',
'blast',
'blend',
'blended',
'blending',
'blessed',
'blessing',
'block',
'blocks',
'blog',
'blood',
'blossom',
'blow',
'blue',
'bluetooth',
'board',
'boards',
'bodies',
'body',
'bond',
'bonds',
'bonus',
'boogie',
'book',
'books',
'bookshelf',
'bookshelves',
'boost',
'bored',

'boring',
'born',
'borrow',
'borrowed',
'boston',
'bot',
'bots',
'bottle',
'bottles',
'bottom',
'bought',
'bounce',
'bouncing',
'bouncy',
'bound',
'boundaries',
'bounds',
'bowling',
'box',
'boxes',
'boy',
'boys',
'brain',
'brainpop',
'brains',
'brainstorm',
'brainstormed',
'brainstorming',
'brand',
'brave',
'break',
'breakfast',
'breaking',
'breakout',
'breaks',
'breath',
'breathing',
'bricks',
'bridge',
'bridges',
'bright',
'brighten',
'brighter',
'brightest',
'brightly',
'brilliant',
'bring',
'bringing',
'brings',
'broad',
'broaden',
'broader',
'broke',
'broken',
'bronx',
'brooklyn',
'brothers',
'brought',
'brushes',
'bubble',
'bubbles',
'bucket',
'buckets',
'buddies',
'budding',
'buddy',
'budget',
'budgets',
'build',
'builders',
'building',
'buildings',
'builds',
'built',
'bulletin',
'bullying',
'bunch',

'bundle',
'burden',
'burn',
'burning',
'bursting',
'bus',
'business',
'businesses',
'busy',
'butterflies',
'butterfly',
'button',
'buttons',
'buy',
'buying',
'buzz',
'ca',
'cabinet',
'caddies',
'cafeteria',
'calculate',
'calculator',
'calculators',
'calculus',
'calendar',
'california',
'call',
'called',
'calm',
'calming',
'calories',
'came',
'camera',
'cameras',
'camp',
'campus',
'candy',
'cannot',
'canvas',
'capabilities',
'capability',
'capable',
'capacity',
'capture',
'car',
'card',
'cardboard',
'cardiovascular',
'cards',
'care',
'cared',
'career',
'careers',
'carefully',
'cares',
'caring',
'carolina',
'carpet',
'carpets',
'carry',
'carrying',
'cars',
'cart',
'cartridges',
'carts',
'case',
'cases',
'cat',
'catch',
'catching',
'cater',
'caucasian',
'caught',
'cause',
'caused',
'causes',
'causing',

'cd',
'cds',
'cease',
'celebrate',
'celebrated',
'celebrating',
'celebration',
'cell',
'cells',
'center',
'centered',
'centers',
'central',
'century',
'ceramics',
'cerebral',
'certain',
'certainly',
'certificates',
'certified',
'chain',
'chair',
'chairs',
'chalk',
'challenge',
'challenged',
'challenges',
'challenging',
'chance',
'chances',
'change',
'changed',
'changer',
'changes',
'changing',
'channel',
'chaos',
'chaotic',
'chapter',
'character',
'characteristics',
'characters',
'charge',
'charged',
'charging',
'charismatic',
'chart',
'charter',
'charts',
'check',
'checked',
'checking',
'cheer',
'cheerful',
'chemical',
'chemicals',
'chemistry',
'chess',
'chicago',
'chickens',
'chicks',
'child',
'childhood',
'children',
'china',
'chinese',
'chips',
'choice',
'choices',
'choir',
'choose',
'choosing',
'chore',
'chorus',
'chose',
'chosen',
'christmas',

'chrome',
'chromebook',
'chromebooks',
'circle',
'circles',
'circuit',
'circuits',
'circumstances',
'cities',
'citizen',
'citizens',
'citizenship',
'city',
'civil',
'class',
'classes',
'classic',
'classical',
'classics',
'classified',
'classmate',
'classmates',
'classroom',
'classrooms',
'classwork',
'clay',
'clean',
'cleaner',
'cleaning',
'clear',
'clearly',
'clever',
'click',
'climate',
'climb',
'climbing',
'clip',
'clipboard',
'clipboards',
'clips',
'clock',
'close',
'closed',
'closely',
'closer',
'closet',
'closing',
'clothes',
'clothing',
'cloud',
'club',
'clubs',
'clues',
'clutter',
'co',
'coach',
'coaches',
'coaching',
'coast',
'coaster',
'coats',
'code',
'coded',
'codes',
'coding',
'coffee',
'cognitive',
'coins',
'cold',
'collaborate',
'collaborating',
'collaboration',
'collaborative',
'collaboratively',
'collaborators',
'collar',
'colleagues',

```
'collect',
'collected',
'collecting',
'collection',
'college',
'colleges',
'color',
'colored',
'colorful',
'coloring',
'colors',
'com',
'combat',
'combination',
'combine',
...]
```

In [134]:

```
print('Shape of normalized essay in train data', text_bow.shape)
print('Shape of normalized title in train data', title_bow.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

```
Shape of normalized essay in train data (34998, 5000)
Shape of normalized title in train data (34998, 1586)
```

```
Shape of normalized essay in test data (15000, 5000)
Shape of normalized title in test data (15000, 1586)
```

Merge them

In [135]:

```
# for train data
from scipy.sparse import hstack
tr_X = hstack((quantity_normalized.T, price_normalized.T, teacher_number_of_previously_posted_projects_
normalized.T, \
               school_state_one_hot, categories_one_hot, subcategories_one_hot, project_grade_category_o
ne_hot, \
               teacher_prefix_one_hot, text_bow, title_bow))
tr_X.shape
```

Out[135]:

(34998, 6688)

In [136]:

```
tr X = tr X.toarray()
```

In [137]:

 $\text{tr } X$

Out[137]:

```
array([[0.00119224, 0.00040346, 0.00018164, ..., 0.          , 0.          ,
        0.          ],
       [0.00187352, 0.00362029, 0.00054493, ..., 0.          , 0.          ,
        0.          ],
       [0.00153288, 0.00314666, 0.00181642, ..., 0.          , 0.          ,
        0.          ],
       ...,
       [0.0017032 , 0.00200517, 0.01852746, ..., 0.          , 0.          ,
        0.          ],
```

```
[0.0017032 , 0.0021566 , 0.00090821, ..., 0.      , 0.      ,
 0.      ],
[0.00306577, 0.00196764, 0.      , ..., 0.      , 0.      ,
 0.      ]])
```

In [138]:

```
tr_X.shape, tr_y.shape
```

Out[138]:

```
((34998, 6688), (34998,))
```

In [139]:

```
# for test data
ts_X = hstack((ts_quantity.T, ts_price.T, ts_teacher_number_of_previously_posted_projects.T, ts_school_
state, \
               ts_project_subject_category, ts_project_subject_subcategory,ts_project_grade_category, \
               ts_teacher_prefix, ts_essay, ts_title))
ts_X.shape
```

Out[139]:

```
(15000, 6688)
```

In [140]:

```
ts_X = ts_X.toarray()
```

In [141]:

```
ts_X.shape, ts_y.shape
```

Out[141]:

```
((15000, 6688), (15000,))
```

In [142]:

```
len(list_features)
```

Out[142]:

```
6688
```

In [143]:

```
# check whether data still contain NaN or infinity or not
```

In [144]:

```
np.any(np.isnan(tr_X)), np.any(np.isnan(ts_X))
```

Out[144]:

```
(False, False)
```

In [145]:

```
np.all(np.isfinite(tr_X)), np.all(np.isfinite(ts_X))
```

Out[145]:

(True, True)

In []:

In [146]:

```
from sklearn.naive_bayes import MultinomialNB
from sklearn.model_selection import KFold
```

In [151]:

```
def kcross(tr_X, tr_y):
    kf = KFold(n_splits=10, random_state=1)
    tr_score = []
    cv_score = []
    kfoldno = 0
    for train_index, cv_index in kf.split(tr_X):
        kfoldno += 1
        print('*****')
        print('KFold', kfoldno)
        print('TRAIN:', train_index, 'CV:', cv_index)
        print('*****')
        X_train, X_cv = tr_X[train_index], tr_X[cv_index]
        y_train, y_cv = tr_y[train_index], tr_y[cv_index]
        for i in [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]:
            mnb = MultinomialNB(class_prior=[0.5, 0.5], alpha=i)
            mnb.fit(X_train, y_train)
            tr_pred = np.argmax(mnb.predict_proba(X_train), axis=1)
            cv_pred = np.argmax(mnb.predict_proba(X_cv), axis=1)
            print('====Alpha={0}===='.format(i))
            print('Train score', metrics.roc_auc_score(y_train, tr_pred))
            print('CV score', metrics.roc_auc_score(y_cv, cv_pred))
    return kf
```

In [152]:

```
def bestK_tr_cv_data(kf, k, tr_X, tr_y):
    kfoldno = 0
    for train_index, cv_index in kf.split(tr_X):
        kfoldno += 1
        print('*****')
        print('KFold', kfoldno)
        print('TRAIN:', train_index, 'CV:', cv_index)
        print('*****')
        if kfoldno == k:
            print('Storing....')
            X_train, X_cv = tr_X[train_index], tr_X[cv_index]
            y_train, y_cv = tr_y[train_index], tr_y[cv_index]
            break
        else:
            continue
    return X_train, y_train, X_cv, y_cv
```

In [153]:

```
def multinb(X, y, cv_X, cv_y):
    """
    Parameters:
    k - alpha
    X - train feature data
    y - train class data
    cv_X - valid feature data
    cv_y - valid class data

    Return:
    Print the AUC score of Train and CV data
    and Store as well
```

```

"""
tr_score = []
cv_score = []
index = 0
# Create knn model instance
for i in [0.0001,0.001,0.01,0.1,1,10,100,1000,10000]:
    nb = MultinomialNB(alpha=i, class_prior=[0.5,0.5])
    # Fit the model with train data
    nb.fit(X,y)
    # Predict the cv data
    tr_pred = np.argmax(nb.predict_proba(X),axis=1)
    cv_pred = np.argmax(nb.predict_proba(cv_X),axis=1)
    # Evaluate accuracy to see how much it corrected class label
    tr_score.append(metrics.roc_auc_score(y, tr_pred))
    cv_score.append(metrics.roc_auc_score(cv_y, cv_pred))
    print('\nTrain AUC and CV AUC score for alpha:{0} is {1} , {2}'.format(i,tr_score[index],cv_score[index]))
    index += 1
return tr_score, cv_score

```

In [154]:

```

def plotauc_tr_cv(feature_name, n_list, X_score, cv_X_score):
    """
    Parameters:
    feature_name - name the feature you want to print in graph
    n_list = range of list on x axis
    X_score - Train AUC score
    cv_X_score - CV AUC score
    Return:
    Plot the graph of Train and CV AUC score
    """
    n_list = np.log(n_list)
    plt.figure(figsize=(10,10))
    plt.plot(n_list, X_score, label='Train AUC')
    plt.plot(n_list, cv_X_score, label='CV AUC')
    plt.scatter(n_list, X_score)
    plt.scatter(n_list, cv_X_score)
    plt.legend()
    plt.xlabel('Log-Hyperparameter(alpha) ')
    plt.ylabel('AUC Score')
    plt.title('Train AUC vs CV AUC plot with {0} features'.format(feature_name))
    plt.show()

```

In [155]:

```

def plotauc_tr_ts(k, feature_name, X, y, ts_X, ts_y):
    """
    Parameters:
    k = alpha
    feature_name - (string) Write feature to print the plot title
    X - train feature data
    y - train class data
    ts_X - test feature data
    ts_y - test class data

    Return:
    Save the prediction of test data, train threshold value, FPR and TPR for train data
    and plot the graph for Train and Test data
    """
    nb = MultinomialNB(alpha=k, class_prior=[0.5,0.5])
    # Fit the model with train data
    nb.fit(X,y)

    tr_predict = np.argmax(nb.predict_proba(tr_X),axis=1)
    ts_predict = np.argmax(nb.predict_proba(ts_X),axis=1)

    # Compute ROC curve and ROC area for each class
    fpr = dict()
    tpr = dict()
    roc_auc = dict()
    fpr, tpr, tr_thre = roc_curve(y, tr_predict)
    roc_auc = auc(fpr, tpr)

```

```

fpr_t = dict()
tpr_t = dict()
roc_auc_t = dict()
fpr_t, tpr_t, _ = roc_curve(ts_y, ts_predict)
roc_auc_t = auc(fpr_t, tpr_t)

plt.figure()
lw = 2
plt.plot(fpr, tpr, color='darkorange',
         lw=lw, label='ROC curve for train data (area = %0.2f)' % roc_auc)
plt.plot(fpr_t, tpr_t, color='blue',
         lw=lw, label='ROC curve for test data (area = %0.2f)' % roc_auc_t)
plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC With Maximum AUC on MultinomialNB Classifier for alpha={0} on {1} features'.format(k
, feature_name))
plt.legend(loc="lower right")
plt.show()

return tr_thre, fpr, tpr, tr_predict, ts_predict

```

In [156]:

```

# we are writing our own function for predict, with defined threshould
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[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))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

def plot_cm(feature_names, tr_thresholds, train_fpr, train_tpr, y_train, y_train_pred, y_test, y_test_p
red):
    """
    Parameters:
    feature_name - (string) Write feature to print the plot title
    tr_thresholds - train threshold value
    train_fpr = FPR for train data
    train_tpr - TPR for train data
    y_true - test class data
    y_pred - test prediction value

    Return:
    Plot the confusion matrix for Train and Test Data
    """
    best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
    print("Train confusion matrix")
    cm = metrics.confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
    plt.figure(figsize = (10,7))
    sns.heatmap(cm, annot=True, fmt="d")
    plt.xlabel('Predicted Class')
    plt.ylabel('True Class')
    plt.title('Confusion matrix for Train Data when MultinomialNB with {0} features'.format(feature_nam
es))

    print("Test confusion matrix")
    cm = metrics.confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
    plt.figure(figsize = (10,7))
    sns.heatmap(cm, annot=True, fmt="d")
    plt.xlabel('Predicted Class')
    plt.ylabel('True Class')
    plt.title('Confusion matrix for Test Data when MultinomialNB with {0} features'.format(feature_name

```

```
print('CONFUSION MATRIX FOR TEST DATA WHEN MULTICOLLINEAR WITH 10 FEATURES', format(feature_names))
```

2.4.1 Applying Naive Bayes on BOW, SET 1

In [83]:

```
kf = kcross(tr_X, tr_y)
```

```
*****
KFold 1
TRAIN: [ 3500  3501  3502 ... 34995 34996 34997] CV: [   0   1   2 ... 3497 3498 3499]
*****
=====Alpha=0.0001=====
Train score 0.7205799453475253
CV score 0.6447112635791881
=====Alpha=0.001=====
Train score 0.7202243931187743
CV score 0.6443745632424878
=====Alpha=0.01=====
Train score 0.7189859141860966
CV score 0.6448129089638523
=====Alpha=0.1=====
Train score 0.7168823721112101
CV score 0.6466329966329967
=====Alpha=1=====
Train score 0.7074695868554447
CV score 0.647716155263325
=====Alpha=10=====
Train score 0.6506380553336707
CV score 0.6246045359252906
=====Alpha=100=====
Train score 0.5007544868337132
CV score 0.49966329966329964
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 2
TRAIN: [   0   1   2 ... 34995 34996 34997] CV: [3500 3501 3502 ... 6997 6998 6999]
*****
=====Alpha=0.0001=====
Train score 0.7197320171770937
CV score 0.6432451285099287
=====Alpha=0.001=====
Train score 0.7195078856866192
CV score 0.6422173792807406
=====Alpha=0.01=====
Train score 0.718480616355278
CV score 0.6437713856453564
=====Alpha=0.1=====
Train score 0.7165941763104519
CV score 0.643612476627992
=====Alpha=1=====
Train score 0.7089131850753457
CV score 0.6443306628363403
=====Alpha=10=====
Train score 0.6493960585472411
CV score 0.5940266719498756
=====Alpha=100=====
Train score 0.5007650004076841
CV score 0.5006892936597366
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 3
TRAIN: [   0   1   2 ... 34995 34996 34997] CV: [ 7000  7001  7002 ... 10497 10498 10499]
```

```

*****
=====Alpha=0.0001=====
Train score 0.7191230883345257
CV score 0.6614111676726921
=====Alpha=0.001=====
Train score 0.7187863124027789
CV score 0.6602307292241257
=====Alpha=0.01=====
Train score 0.7177572748335528
CV score 0.6618540290933161
=====Alpha=0.1=====
Train score 0.7153961828737674
CV score 0.6634142881908874
=====Alpha=1=====
Train score 0.7064242387873195
CV score 0.6619249499613875
=====Alpha=10=====
Train score 0.6479989837935577
CV score 0.623323509479756
=====Alpha=100=====
Train score 0.5009836857668456
CV score 0.49983136593591904
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 4
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [10500 10501 10502 ... 13997 13998 13999]
*****
=====Alpha=0.0001=====
Train score 0.7179394436857456
CV score 0.6363795518207283
=====Alpha=0.001=====
Train score 0.7176769272784701
CV score 0.6381497043261749
=====Alpha=0.01=====
Train score 0.7167018663371613
CV score 0.6393168378462497
=====Alpha=0.1=====
Train score 0.7145532373240738
CV score 0.6415343915343916
=====Alpha=1=====
Train score 0.7063046460577591
CV score 0.6396086212262684
=====Alpha=10=====
Train score 0.6506046409092374
CV score 0.5978349673202614
=====Alpha=100=====
Train score 0.5006213803943131
CV score 0.5019354964207906
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 5
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [14000 14001 14002 ... 17497 17498 17499]
*****
=====Alpha=0.0001=====
Train score 0.7180777352091073
CV score 0.6504581977875841
=====Alpha=0.001=====
Train score 0.7176659044341166
CV score 0.6530015933672624
=====Alpha=0.01=====
Train score 0.7168343531360537
CV score 0.6539139298183897
=====Alpha=0.1=====
Train score 0.7142717801617483
CV score 0.6533156922513899
=====Alpha=1=====
Train score 0.7043290398716122

```



```
CV score 0.6486161286394395
=====Alpha=10=====
Train score 0.6488643219805411
CV score 0.6040105434136203
=====Alpha=100=====
Train score 0.5006669456619209
CV score 0.5015837026814057
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 6
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [17500 17501 17502 ... 20997 20998 20999]
*****
=====Alpha=0.0001=====
Train score 0.7193362061752424
CV score 0.6496338001701185
=====Alpha=0.001=====
Train score 0.7190182946668457
CV score 0.6489562933950507
=====Alpha=0.01=====
Train score 0.7183120014031674
CV score 0.6499342274444642
=====Alpha=0.1=====
Train score 0.7162175255831428
CV score 0.648983492572152
=====Alpha=1=====
Train score 0.7081846794485971
CV score 0.6482948588609974
=====Alpha=10=====
Train score 0.6509338786562289
CV score 0.6019771329100152
=====Alpha=100=====
Train score 0.5007207714605469
CV score 0.49983062330623307
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 7
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [21000 21001 21002 ... 24497 24498 24499]
*****
=====Alpha=0.0001=====
Train score 0.7201338111650821
CV score 0.6522324674681588
=====Alpha=0.001=====
Train score 0.7198529438468032
CV score 0.6520650193770671
=====Alpha=0.01=====
Train score 0.718758851908904
CV score 0.6546484111326267
=====Alpha=0.1=====
Train score 0.7167472600646271
CV score 0.6550549777039935
=====Alpha=1=====
Train score 0.7084481335785832
CV score 0.6534202412816229
=====Alpha=10=====
Train score 0.651230351149166
CV score 0.6276273713125584
=====Alpha=100=====
Train score 0.5006469574152037
CV score 0.5008053145548226
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 8
```

```

TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [24500 24501 24502 ... 27997 27998 27999]
*****
=====Alpha=0.0001=====
Train score 0.7169035451057898
CV score 0.6490323198055712
=====Alpha=0.001=====
Train score 0.7166602087775479
CV score 0.6499883427501219
=====Alpha=0.01=====
Train score 0.7157242998227711
CV score 0.6499366397961169
=====Alpha=0.1=====
Train score 0.713793284596295
CV score 0.6480891423154316
=====Alpha=1=====
Train score 0.7057701943261437
CV score 0.6481536907238478
=====Alpha=10=====
Train score 0.6506131959705275
CV score 0.6065161136591497
=====Alpha=100=====
Train score 0.5007343424035386
CV score 0.4996640913671481
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 9
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [28000 28001 28002 ... 31496 31497 31498]
*****
=====Alpha=0.0001=====
Train score 0.7178194993696366
CV score 0.648137462635741
=====Alpha=0.001=====
Train score 0.717330634482148
CV score 0.6515383984764211
=====Alpha=0.01=====
Train score 0.71634148253147
CV score 0.651292141208521
=====Alpha=0.1=====
Train score 0.7147884452069853
CV score 0.6544603781199945
=====Alpha=1=====
Train score 0.7041593124278179
CV score 0.6474664825255086
=====Alpha=10=====
Train score 0.6474753644918707
CV score 0.5951051244214058
=====Alpha=100=====
Train score 0.5006507810666053
CV score 0.5026350473595923
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 10
TRAIN: [    0    1    2 ... 31496 31497 31498] CV: [31499 31500 31501 ... 34995 34996 34997]
*****
=====Alpha=0.0001=====
Train score 0.7179945829694359
CV score 0.6479956435933408
=====Alpha=0.001=====
Train score 0.7176392568174311
CV score 0.6476570049208045
=====Alpha=0.01=====
Train score 0.716734358346234
CV score 0.6506427312387354
=====Alpha=0.1=====
Train score 0.7141450952591496
CV score 0.6519352641939842
=====Alpha=1=====

```

```

Train score 0.7062656120959325
CV score 0.6506259853703131
=====Alpha=10=====
Train score 0.6523680881942935
CV score 0.6104690827853713
=====Alpha=100=====
Train score 0.5007390675593446
CV score 0.4998306806637318
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5

```

In [84]:

```

# From the above observation, we got best KFold Number and Alpha value
kfoldno = 3
alpha = 0.1
tr_X, tr_y, cv_X, cv_y = bestK_tr_cv_data(kf, kfoldno, tr_X, tr_y)

*****
KFold 1
TRAIN: [ 3500  3501  3502 ... 34995 34996 34997] CV: [   0    1    2 ... 3497 3498 3499]
*****
KFold 2
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [3500 3501 3502 ... 6997 6998 6999]
*****
KFold 3
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [ 7000  7001  7002 ... 10497 10498 10499]
*****
Storing.....

```

In [85]:

```

# We got to know that KFold 1 have got high cv score.
# We need to compile again for all alpha [0.0001,0.001,0.01,0.1,1,10,100,1000,10000] to store the roc_auc_score value
# So that we can plot graph of ROC_AUC Score with Log-Hyperparameter
tr_score, cv_score = multirun(tr_X, tr_y, cv_X, cv_y)

```

```

Train AUC and CV AUC score for alpha:0.0001 is 0.7191230883345257 , 0.6614111676726921
Train AUC and CV AUC score for alpha:0.001 is 0.7187863124027789 , 0.6602307292241257
Train AUC and CV AUC score for alpha:0.01 is 0.7177572748335528 , 0.6618540290933161
Train AUC and CV AUC score for alpha:0.1 is 0.7153961828737674 , 0.6634142881908874
Train AUC and CV AUC score for alpha:1 is 0.7064242387873195 , 0.6619249499613875
Train AUC and CV AUC score for alpha:10 is 0.6479989837935577 , 0.623323509479756
Train AUC and CV AUC score for alpha:100 is 0.5009836857668456 , 0.49983136593591904
Train AUC and CV AUC score for alpha:1000 is 0.5 , 0.5
Train AUC and CV AUC score for alpha:10000 is 0.5 , 0.5

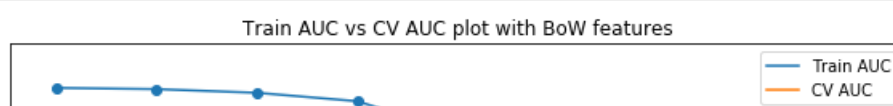
```

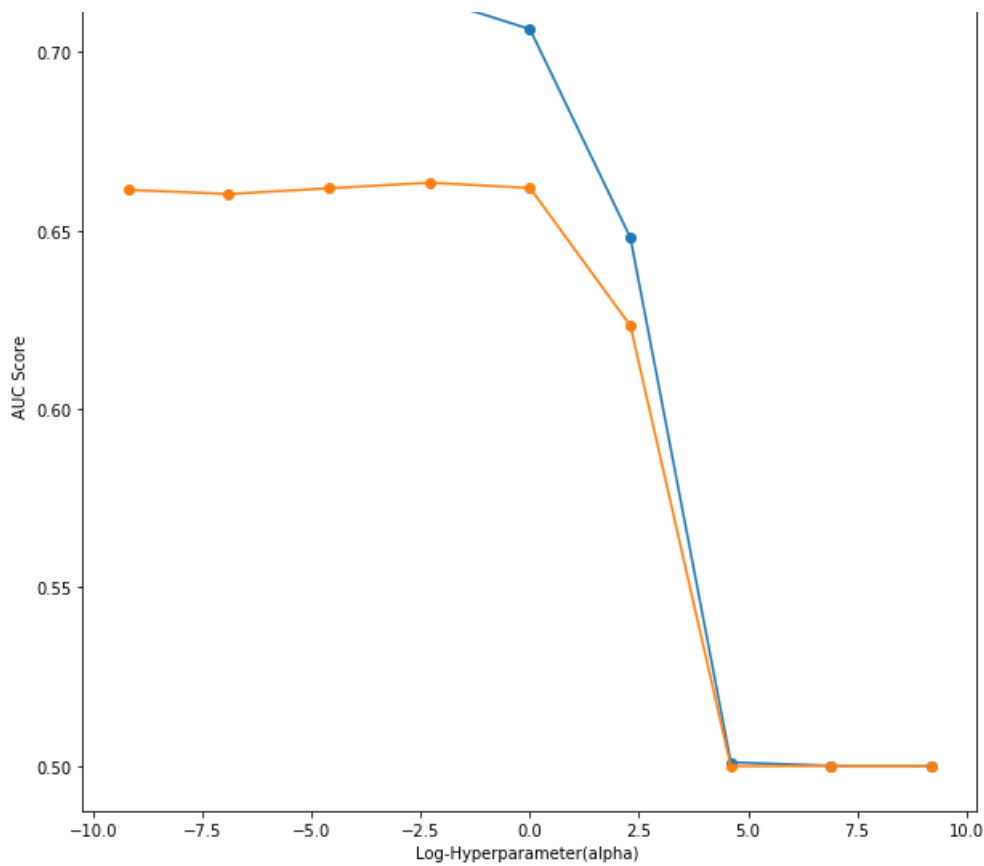
In [86]:

```

plotauc_tr_cv('BoW', [0.0001,0.001,0.01,0.1,1,10,100,1000,10000], tr_score, cv_score)

```

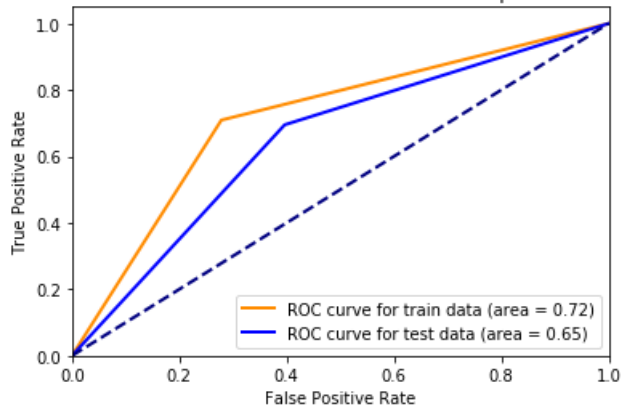




In [87]:

```
tr_thre, fpr, tpr, tr_predict, ts_predict = plotauc_tr_ts(alpha, 'BoW', tr_X, tr_y, ts_X, ts_y)
```

ROC With Maximum AUC on MultinomialNB Classifier for alpha=0.1 on BoW features



In [88]:

```
plot_cm('BoW', tr_thre, fpr, tpr, tr_y, tr_predict, ts_y, ts_predict)
```

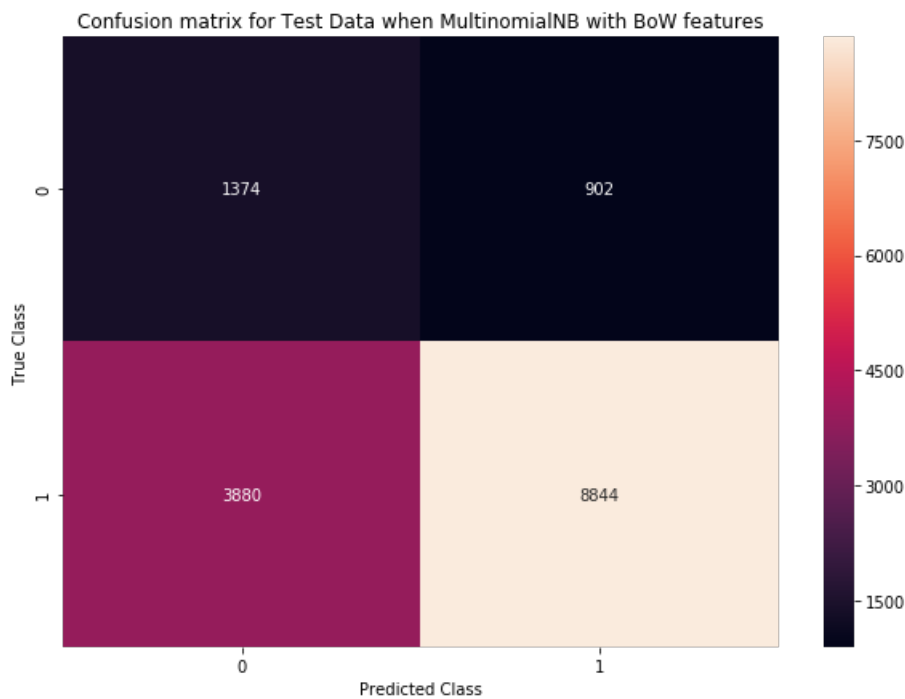
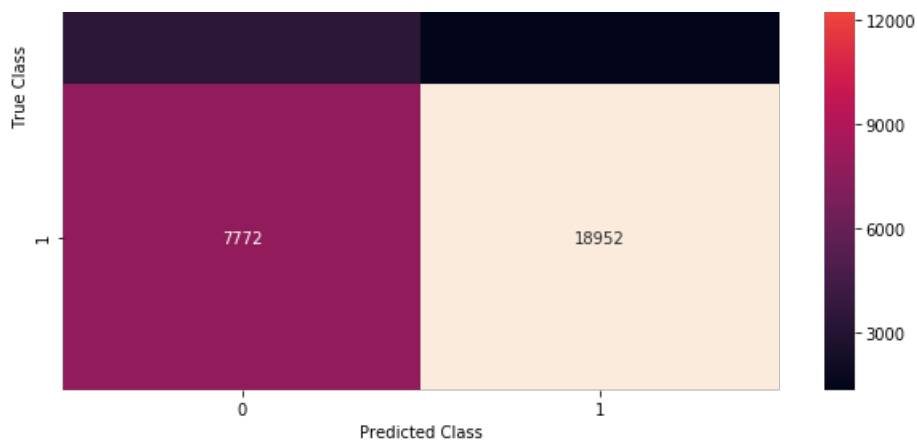
the maximum value of $tpr \cdot (1 - fpr)$ 0.5117529987527236 for threshold 1

Train confusion matrix

Test confusion matrix

Confusion matrix for Train Data when MultinomialNB with BoW features





2.4.1.1 Top 10 important features of positive class from SET 1

In [89]:

```
# Please write all the code with proper documentation
multinb = MultinomialNB(alpha=alpha, class_prior=[0.5,0.5])
multinb.fit(tr_X, tr_y)
```

Out[89]:

```
MultinomialNB(alpha=0.1, class_prior=[0.5, 0.5], fit_prior=True)
```

In [107]:

```
top10_pos = np.argsort(multinb.feature_log_prob_[1])
```

In [108]:

```
# reverse the order list
top10_pos = np.flip(top10_pos)[:10]
```

In [109]:

```
print('-----Top 10 feature names for positive class-----')
for i in top10_pos:
```

```
print(list_features[i])
```

```
-----Top 10 feature names for positive class-----  
students  
school  
learning  
classroom  
not  
learn  
help  
many  
nannan  
need
```

2.4.1.2 Top 10 important features of negative class from SET 1

In [110]:

```
# Please write all the code with proper documentation  
top10_neg = np.argsort(multinb.feature_log_prob_[0])
```

In [111]:

```
# reverse the order list  
top10_neg = np.flip(top10_neg)[:10]
```

In [112]:

```
print('-----Top 10 feature names for negative class-----')  
for i in top10_neg:  
    print(list_features[i])
```

```
-----Top 10 feature names for negative class-----  
students  
school  
learning  
classroom  
learn  
not  
help  
nannan  
many  
work
```

2.4.2 Applying Naive Bayes on TFIDF, SET 2

In [147]:

```
# Please write all the code with proper documentation  
kf = Kcross(tr_X, tr_y)
```

```
*****  
KFold 1  
TRAIN: [ 3500  3501  3502 ... 34995 34996 34997] CV: [   0   1   2 ... 3497 3498 3499]  
*****  
=====Alpha=0.0001=====  
Train score 0.7045851894188727  
CV score 0.6175655930372912  
=====Alpha=0.001=====  
Train score 0.704510336318083  
CV score 0.6175655930372912  
=====Alpha=0.01=====  
Train score 0.7037660577933132  
CV score 0.6172288927005908  
=====Alpha=0.1=====  
Train score 0.7016744234543286  
CV score 0.6104250047646274
```

```

CV score 0.6104230047040274
=====Alpha=1=====
Train score 0.6704742327676598
CV score 0.5967124070897656
=====Alpha=10=====
Train score 0.5075669882977067
CV score 0.5032050060351947
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 2
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [3500 3501 3502 ... 6997 6998 6999]
*****
=====Alpha=0.0001=====
Train score 0.7051915894436976
CV score 0.6211072450129397
=====Alpha=0.001=====
Train score 0.7050982013226667
CV score 0.6211072450129397
=====Alpha=0.01=====
Train score 0.7045565502206868
CV score 0.6217965386726763
=====Alpha=0.1=====
Train score 0.7018051481347342
CV score 0.6199123317524982
=====Alpha=1=====
Train score 0.6697835465616422
CV score 0.6124931380197046
=====Alpha=10=====
Train score 0.5076726479308922
CV score 0.5007099311944592
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 3
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [ 7000 7001 7002 ... 10497 10498 10499]
*****
=====Alpha=0.0001=====
Train score 0.7039098081980883
CV score 0.6250130021591463
=====Alpha=0.001=====
Train score 0.7037601300062011
CV score 0.6256103134702369
=====Alpha=0.01=====
Train score 0.703187651680851
CV score 0.6260389907172464
=====Alpha=0.1=====
Train score 0.7001869707610833
CV score 0.6242754251312036
=====Alpha=1=====
Train score 0.6688021675057612
CV score 0.6077618950055949
=====Alpha=10=====
Train score 0.50698544735438
CV score 0.5082851334100329
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5

```

```

CV SCORE 0.5
*****
KFold 4
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [10500 10501 10502 ... 13997 13998 13999]
*****
=====Alpha=0.0001=====
Train score 0.7014752279869758
CV score 0.6112959072517896
=====Alpha=0.001=====
Train score 0.701343969783338
CV score 0.6112959072517896
=====Alpha=0.01=====
Train score 0.7014105707801648
CV score 0.610634531590414
=====Alpha=0.1=====
Train score 0.6991507305953986
CV score 0.6096424680983505
=====Alpha=1=====
Train score 0.6677177216499681
CV score 0.6078042328042328
=====Alpha=10=====
Train score 0.5083204339004893
CV score 0.5042697634609399
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 5
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [14000 14001 14002 ... 17497 17498 17499]
*****
=====Alpha=0.0001=====
Train score 0.699092804080587
CV score 0.6290392311117039
=====Alpha=0.001=====
Train score 0.6990661947518241
CV score 0.6290392311117039
=====Alpha=0.01=====
Train score 0.6982957518593689
CV score 0.6290392311117039
=====Alpha=0.1=====
Train score 0.6950159546501347
CV score 0.6270251324208681
=====Alpha=1=====
Train score 0.6652947262703119
CV score 0.6004894143843759
=====Alpha=10=====
Train score 0.5087567836976692
CV score 0.5048026524764844
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 6
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [17500 17501 17502 ... 20997 20998 20999]
*****
=====Alpha=0.0001=====
Train score 0.706183707013395
CV score 0.627948638063023
=====Alpha=0.001=====
Train score 0.7061089043055369
CV score 0.627948638063023
=====Alpha=0.01=====
Train score 0.7058586113970025
CV score 0.627779261369256
=====Alpha=0.1=====
Train score 0.7002041000050750

```



```

Train score 0.7002841088858759
CV score 0.6274677071588235
=====Alpha=1=====
Train score 0.6699949276033987
CV score 0.6128802939489248
=====Alpha=10=====
Train score 0.508412071981824
CV score 0.5014749371946274
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 7
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [21000 21001 21002 ... 24497 24498 24499]
*****
=====Alpha=0.0001=====
Train score 0.7038777613884815
CV score 0.6225628809932734
=====Alpha=0.001=====
Train score 0.703661139174413
CV score 0.6223954329021816
=====Alpha=0.01=====
Train score 0.7034470981671206
CV score 0.6220605367199981
=====Alpha=0.1=====
Train score 0.7005961416152349
CV score 0.6237428362188265
=====Alpha=1=====
Train score 0.6675789807294756
CV score 0.6211196999747198
=====Alpha=10=====
Train score 0.5074470978547204
CV score 0.5043373616435715
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 8
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [24500 24501 24502 ... 27997 27998 27999]
*****
=====Alpha=0.0001=====
Train score 0.7017145288863507
CV score 0.6295627856909345
=====Alpha=0.001=====
Train score 0.7016770925281597
CV score 0.6297307400073604
=====Alpha=0.01=====
Train score 0.7013588834835356
CV score 0.6293948313745086
=====Alpha=0.1=====
Train score 0.6988632801045577
CV score 0.6285679052467901
=====Alpha=1=====
Train score 0.6703748685738214
CV score 0.6096025552177915
=====Alpha=10=====
Train score 0.5077743097911988
CV score 0.5065498972042511
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5

```

```

Train score 0.5
CV score 0.5
*****
KFold 9
TRAIN: [    0    1    2 ... 34995 34996 34997] CV: [28000 28001 28002 ... 31496 31497 31498]
*****
=====Alpha=0.0001=====
Train score 0.7006242184141337
CV score 0.6150160177582705
=====Alpha=0.001=====
Train score 0.7005979888725294
CV score 0.6145099448837766
=====Alpha=0.01=====
Train score 0.7002238073945126
CV score 0.615444524323029
=====Alpha=0.1=====
Train score 0.6968752515476919
CV score 0.6161464048330748
=====Alpha=1=====
Train score 0.6669525215121037
CV score 0.6005823779434207
=====Alpha=10=====
Train score 0.5080665265510916
CV score 0.4997944177481807
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5
*****
KFold 10
TRAIN: [    0    1    2 ... 31496 31497 31498] CV: [31499 31500 31501 ... 34995 34996 34997]
*****
=====Alpha=0.0001=====
Train score 0.7025330835712095
CV score 0.6093933157935867
=====Alpha=0.001=====
Train score 0.7026193580503415
CV score 0.6093933157935867
=====Alpha=0.01=====
Train score 0.7026537649706814
CV score 0.6096317893642649
=====Alpha=0.1=====
Train score 0.7003594873358661
CV score 0.6083463889085292
=====Alpha=1=====
Train score 0.669555473740142
CV score 0.5906394937041737
=====Alpha=10=====
Train score 0.5081461750781163
CV score 0.4998068022957965
=====Alpha=100=====
Train score 0.5
CV score 0.5
=====Alpha=1000=====
Train score 0.5
CV score 0.5
=====Alpha=10000=====
Train score 0.5
CV score 0.5

```

In [148]:

```

# From the above observation, we got best KFold Number and Alpha value
kfoldno = 8
alpha = 0.001
# Printing KFold with train and cv index value so that I can observe everything is flowing well.
tr_X, tr_y, cv_X, cv_y = bestK_tr_cv_data(kf, kfoldno, tr_X, tr_y)

```

```

*****
KFold 1

```

```

TRAIN: [ 3500 3501 3502 ... 34995 34996 34997] CV: [ 0 1 2 ... 3497 3498 3499]
*****
*****
KFold 2
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [3500 3501 3502 ... 6997 6998 6999]
*****
*****
KFold 3
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [ 7000 7001 7002 ... 10497 10498 10499]
*****
*****
KFold 4
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [10500 10501 10502 ... 13997 13998 13999]
*****
*****
KFold 5
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [14000 14001 14002 ... 17497 17498 17499]
*****
*****
KFold 6
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [17500 17501 17502 ... 20997 20998 20999]
*****
*****
KFold 7
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [21000 21001 21002 ... 24497 24498 24499]
*****
*****
KFold 8
TRAIN: [ 0 1 2 ... 34995 34996 34997] CV: [24500 24501 24502 ... 27997 27998 27999]
*****
Storing.....

```

In [157]:

```

# We got to know that KFold 1 have got high cv score.
# We need to compile again for all alpha [0.0001,0.001,0.01,0.1,1,10,100,1000,10000] to store the roc_auc_score value
# So that we can plot graph of ROC_AUC Score with Log-Hyperparameter
tr_score, cv_score = multintb(tr_X, tr_y, cv_X, cv_y)

```

Train AUC and CV AUC score for alpha:0.0001 is 0.7017145288863507 , 0.6295627856909345

Train AUC and CV AUC score for alpha:0.001 is 0.7016770925281597 , 0.6297307400073604

Train AUC and CV AUC score for alpha:0.01 is 0.7013588834835356 , 0.6293948313745086

Train AUC and CV AUC score for alpha:0.1 is 0.6988632801045577 , 0.6285679052467901

Train AUC and CV AUC score for alpha:1 is 0.6703748685738214 , 0.6096025552177915

Train AUC and CV AUC score for alpha:10 is 0.5077743097911988 , 0.5065498972042511

Train AUC and CV AUC score for alpha:100 is 0.5 , 0.5

Train AUC and CV AUC score for alpha:1000 is 0.5 , 0.5

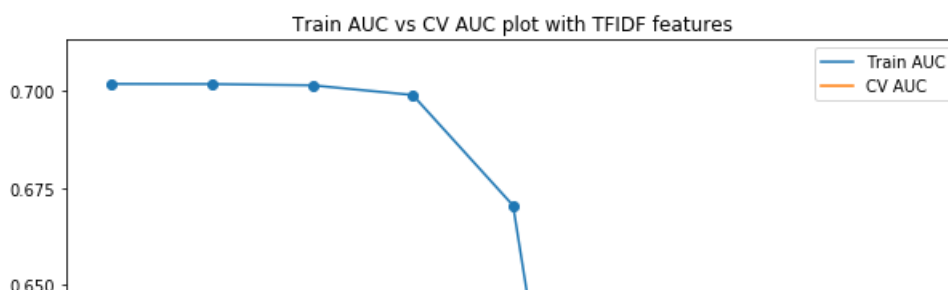
Train AUC and CV AUC score for alpha:10000 is 0.5 , 0.5

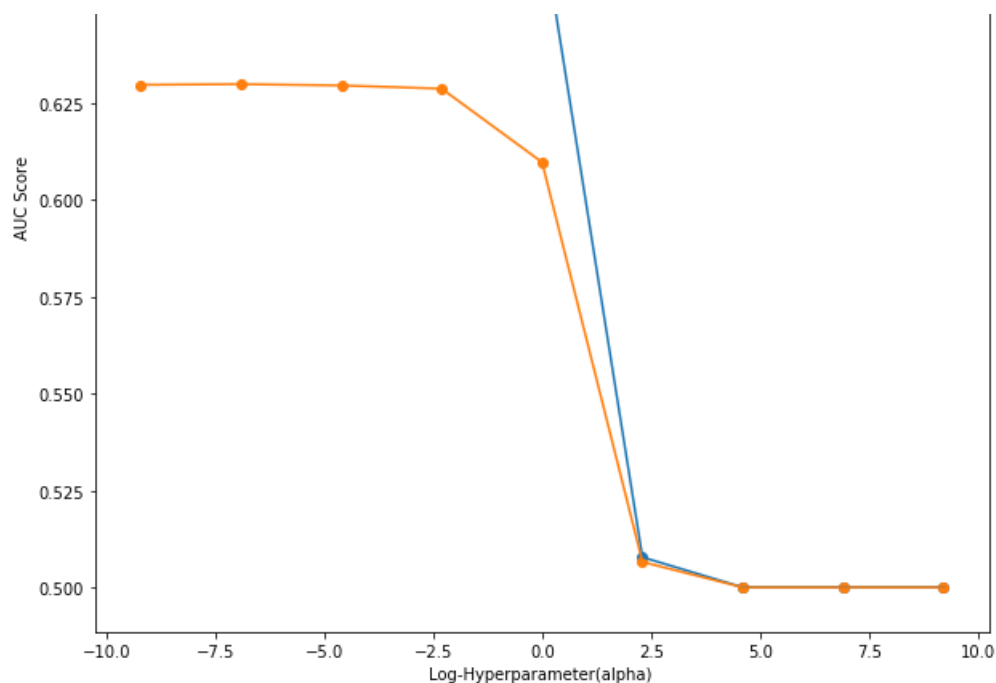
In [158]:

```

plotauc_tr_cv('TFIDF', [0.0001,0.001,0.01,0.1,1,10,100,1000,10000], tr_score, cv_score)

```

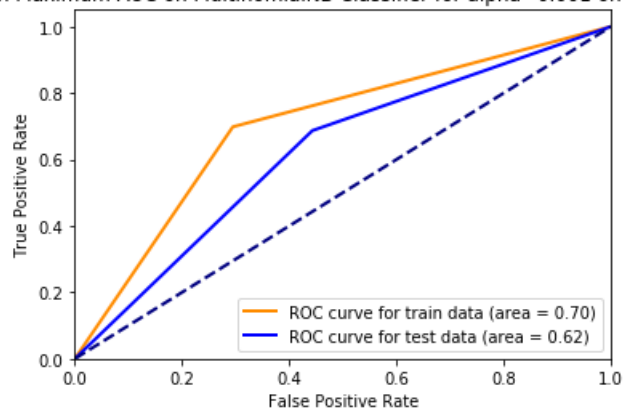




In [159]:

```
tr_thre, fpr, tpr, tr_predict, ts_predict = plotauc_tr_ts(alpha, 'TFIDF', tr_X, tr_y, ts_X, ts_y)
```

ROC With Maximum AUC on MultinomialNB Classifier for alpha=0.001 on TFIDF features



In [160]:

```
plot_cm('TFIDF', tr_thre, fpr, tpr, tr_y, tr_predict, ts_y, ts_predict)
```

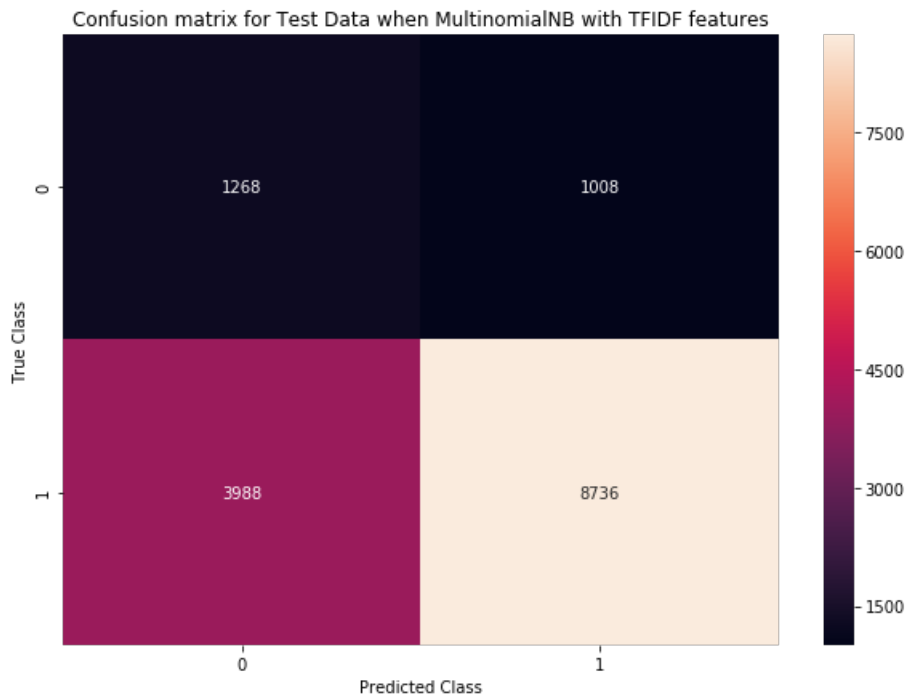
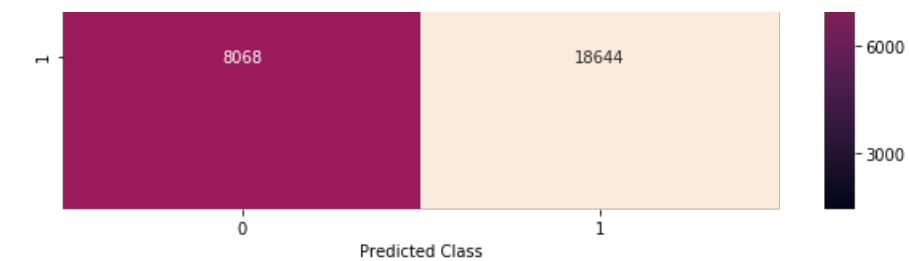
the maximum value of $tpr \cdot (1 - fpr)$ 0.4923369511279216 for threshold 1

Train confusion matrix

Test confusion matrix

Confusion matrix for Train Data when MultinomialNB with TFIDF features





2.4.2.1 Top 10 important features of positive class from SET 2

In [161]:

```
# Please write all the code with proper documentation
```

In [162]:

```
# Please write all the code with proper documentation
multinb = MultinomialNB(alpha=alpha, class_prior=[0.5,0.5])
multinb.fit(tr_X, tr_y)
```

Out[162]:

```
MultinomialNB(alpha=0.001, class_prior=[0.5, 0.5], fit_prior=True)
```

In [163]:

```
top10_pos = np.argsort(multinb.feature_log_prob_[1])
```

In [164]:

```
# Reverse the order of list
top10_pos = np.flip(top10_pos)[0:10]
```

In [165]:

```
print('-----Top 10 feature names for positive class-----')
for i in top10_pos:
    print(list_features[i])
```

```
-----Top 10 feature names for positive class-----
Literacy_Language
grades_prek_2
Math_Science
grades_3_5
Literacy
Mathematics
Literature_Writing
grades_6_8
CA
students
```

2.4.2.2 Top 10 important features of negative class from SET 2

In [166]:

```
top10_neg = np.argsort(multinb.feature_log_prob_[0])
```

In [167]:

```
# Reverse the order of list
top10_neg = np.flip(top10_neg)[0:10]
```

In [168]:

```
print('-----Top 10 feature names for negative class-----')
for i in top10_neg:
    print(list_features[i])
```

```
-----Top 10 feature names for negative class-----
Literacy_Language
grades_prek_2
Math_Science
grades_3_5
Mathematics
Literacy
Literature_Writing
grades_6_8
students
AppliedLearning
```

3. Conclusions

In [169]:

```
# Please compare all your models using Prettytable library
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Features", "Model", "Alpha", "Maximum AUC score",]

x.add_row(["BoW", "MultinomialNB", 0.1, 0.6634142881908874])
x.add_row(["TFIDF", "MultinomialNB", 0.001, 0.6297307400073604])

print(x)
```

Features	Model	Alpha	Maximum AUC score
BoW	MultinomialNB	0.1	0.6634142881908874
TFIDF	MultinomialNB	0.001	0.6297307400073604

