

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care & Hunger</code> <code>Health & Sports</code> <code>History & Civics</code> <code>Literacy & Language</code> <code>Math & Science</code> <code>Music & The Arts</code> <code>Special Needs</code> <code>Warmth</code> Examples: <ul style="list-style-type: none">• <code>Music & The Arts</code>• <code>Literacy & Language, Math & Science</code>
<code>school_state</code>		State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: <code>Literacy</code> <code>Literature & Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

Feature	Description
project_essay_4	Fourth application essay
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <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 [237]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
from sklearn.model_selection import train_test_split
import sklearn.model_selection as model_selection

```

1.1 Reading Data

In [322]:

```

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

```

In [323]:

```

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 [324]:

```

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[324]:

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

id	description	quantity	price
----	-------------	----------	-------

1.2 preprocessing of project_subject_categories

In [325]:

```
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 [326]:

```
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 [327]:

```
# 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 [328]:

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

Out[328]:

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

In [329]:

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

In [330]:

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

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein Our English learner's have a strong support system at home that at begs for more resources. Many times our parents are learning to read and speak English alongside 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 En

ordered 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\nnnannan

=====

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

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs 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 beavers and always strive to work through 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 are 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. \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 a makeup of 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 on smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to

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

In [331]:

```
# 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 [332]:

```
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. \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\nThey 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 [333]:

```
# \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 [334]:

```
#remove special 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 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 nan nan

In [335]:

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

In [336]:

```
# 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)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_essays.append(sent.lower().strip())
```

100%|

109248/109248 [01:11<00:00, 1556.30it/s]

In [337]:

```
# after preprocessing
```



```
preprocessed_essays[20000]

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

Out[337]:

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

1.4 Preprocessing of `project_title`

In [338]:

```
# similarly you can preprocess the titles also

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

100%|

109248/109248 [00:02<00:00, 44063.42it/s]

In [339]:

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

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

print(processed_titles[40]);

print(processed_titles[500]);

print(processed_titles[4000]);

project_data.columns
```

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

Out[339]:

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

Preprocessing of project_grade_category

In [340]:

```
print(project_data['project_grade_category'][1])
print(project_data['project_grade_category'][223])
print(project_data['project_grade_category'][134])
```

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

In [341]:

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

In [342]:

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

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

Grades 68
Grades PreK2
Grades PreK2

Preprocessing of teacher_prefix

In [343]:

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

Ms.
Ms.
Ms.

In [344]:

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

In [345]:

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

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

Mrs
Ms

1.5 Preparing data for models

In [346]:

```
project_data.columns
```

Out[346]:

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

we are going to consider

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

Merging Price of each project.

In [347]:

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

Out[347]:

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

1.5.2.3 Using Pretrained Models: Avg W2V

In [348]:

```
# 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())
```

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 [349]:

```
project_data.columns
```

Out[349]:

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

In [350]:

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

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

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

In [351]:

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

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

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

Vectorizing categorical values

We have following columns with categorical values:

- school_state
- clean_categories
- clean_subcategories
- project_grade_category
- teacher_prefix

In [352]:

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

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

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

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

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

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

cat_features = vectorizer.get_feature_names();
print(len(cat_features));
```

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

In [353]:

```
#vectorizing categories

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

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

```

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

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

cat_features.extend(vectorizer.get_feature_names());
print(len(cat_features), cat_features[23])

```

```

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

```

In [354]:

```

#vectorizing subcategories

vectorizer = CountVectorizer(min_df=10, vocabulary=list(sorted_sub_cat_dict.keys()),
lowercase=False, binary=True);

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

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

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

cat_features.extend(vectorizer.get_feature_names());
print(len(cat_features))

```

```

shape of matrix after one hot encoding of clean_subcategories for train data (53531, 30)
shape of matrix after one hot encoding of clean_subcategories for test data (32775, 30)
shape of matrix after one hot encoding of clean_subcategories for cv data (22942, 30)
90

```

In [355]:

```

#vectorizing project_grade_category

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

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

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

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

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

cat_features.extend(vectorizer.get_feature_names());
print(len(cat_features))

```

```
shape of matrix after one hot encoding of grade_category for train data (53531, 4)
shape of matrix after one hot encoding of grade_category for test data (32775, 4)
shape of matrix after one hot encoding of grade_category for cv data (22942, 4)
94
```

In [356]:

```
#vectorizing teacher_prefix

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

vectorizer = CountVectorizer(min_df=10, vocabulary=list(teacher_prefix_dict.keys()),
lowercase=False, binary=True);

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

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

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

cat_features.extend(vectorizer.get_feature_names());
print(len(cat_features))
```

```
shape of matrix after one hot encoding of teacher_prefix for train data (53531, 6)
shape of matrix after one hot encoding of teacher_prefix for test data (32775, 6)
shape of matrix after one hot encoding of teacher_prefix for cv data (22942, 6)
100
```

2.2 Make Data Model Ready: encoding numerical, categorical features

In [357]:

```
#vectorizing price

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

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

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

In [358]:

```
#vectorizing teacher_number_of_previously_posted_projects

teacher_normalizer = Normalizer();

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

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

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

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

In [359]:

```
#vectorizing quantity:  
  
quantity_normalizer = Normalizer();  
  
quantity_normalizer.fit(X_train['quantity'].values.reshape(-1, 1));  
  
quantity_standardized = quantity_normalizer.transform(X_train['quantity'].values.reshape(-1, 1))  
  
test_quantity_standardized = quantity_normalizer.transform(X_test['quantity'].values.reshape(-1, 1)  
)  
  
cv_quantity_standardized = quantity_normalizer.transform(X_cv['quantity'].values.reshape(-1, 1))  
  
cat_features.append('quantity');  
  
print(len(cat_features))
```

103

2.3 Make Data Model Ready: encoding eassay, and project_title

Vectorizing using BOW on train data.

In [360]:

```
#vectorizing essay  
  
#configure CountVectorizer with word to occur in at least 10 documents.  
vectorizer = CountVectorizer(min_df=10);  
  
vectorizer.fit(X_train['processed_essay']);  
  
#transforming essay into vector  
essay_bow = vectorizer.transform(X_train['processed_essay']);  
cv_essay_bow = vectorizer.transform(X_cv['processed_essay']);  
test_essay_bow = vectorizer.transform(X_test['processed_essay']);  
  
print('Shape of matrix after one hot encoding for train data: ', essay_bow.shape);  
print('Shape of matrix after one hot encoding for test data: ', test_essay_bow.shape);  
print('Shape of matrix after one hot encoding for cv data: ', cv_essay_bow.shape);
```

```
Shape of matrix after one hot encoding for train data:  (53531, 12591)  
Shape of matrix after one hot encoding for test data:  (32775, 12591)  
Shape of matrix after one hot encoding for cv data:  (22942, 12591)
```

In [361]:

```
bow_features = vectorizer.get_feature_names();  
print(len(bow_features), bow_features[546])
```

12591 alcohol

In [362]:

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



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

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

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

Shape of matrix after one hot encoding for train data: (53531, 2191)
 Shape of matrix after one hot encoding for test data: (32775, 2191)
 Shape of matrix after one hot encoding for cv data: (22942, 2191)

In [363]:

```
bow_features.extend(vectorizer.get_feature_names())
print(len(bow_features))
```

14782

Vectorizing using tf-idf

In [364]:

```
#vectorizing essay

#importing TfidfVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer

#configuring TfidfVectorizer with a word to occur atleast in 10 documnets.
vectorizer = TfidfVectorizer(min_df=10)

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

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

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

Shape of matrix after one hot encoding for train data: (53531, 12591)
 Shape of matrix after one hot encoding for test data: (32775, 12591)
 Shape of matrix after one hot encoding for cv data: (22942, 12591)

In [365]:

```
tfidf_features = vectorizer.get_feature_names();
print(len(tfidf_features))
```

12591

In [366]:

```
#vectorizing project_title

vectorizer = TfidfVectorizer(min_df = 10);

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

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

print('Shape of title tfidf after one hot encoding for train data: ', title_tfidf.shape)
```

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

Shape of title_tfidf after one hot encoding for train data (53531, 2191)
 Shape of title_tfidf after one hot encoding for test data (32775, 2191)
 Shape of title_tfidf after one hot encoding for cv data (22942, 2191)

In [367]:

```
tfidf_features.extend(vectorizer.get_feature_names())
print(len(tfidf_features))
```

14782

In [368]:

```
tfidf_features[10855]
```

Out[368]:

'students'

Merging Data

In [369]:

```
from scipy.sparse import hstack

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

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

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

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

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

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

In [370]:

```
#merging all the feature names corresponding to the index value.
final_features_list = []
final_features_list.extend(cat_features);
```

```
final_features_list.extend(tridf_features);

#checking no. of features and no of names present in feature list.
print(len(final_features_list))
print(train_set_2.shape);
```

```
14885
(53531, 14885)
```

2.4 Applying 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

In [371]:

```
# we are writing our own function for predict, with defined thresould
# we will pick a threshold that will give the least fpr
def predict(proba, 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))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

2.4.1 Applying Naive Bayes on BOW, SET 1

In [372]:

```
from sklearn.naive_bayes import MultinomialNB;
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score;

#creating list of alpha on which we train naive bayes.
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]

train_auc = []; #list for storing auc of train data
cv_auc = []; #list for storing auc of cv data

for i in alpha:
    NB = MultinomialNB(alpha=i, class_prior = [0.5, 0.5]);
    NB.fit(train_set_1, y_train);

    y_train_pred = NB.predict_proba(train_set_1)[:,-1];
    y_cv_pred = NB.predict_proba(cv_set_1)[:,-1];

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

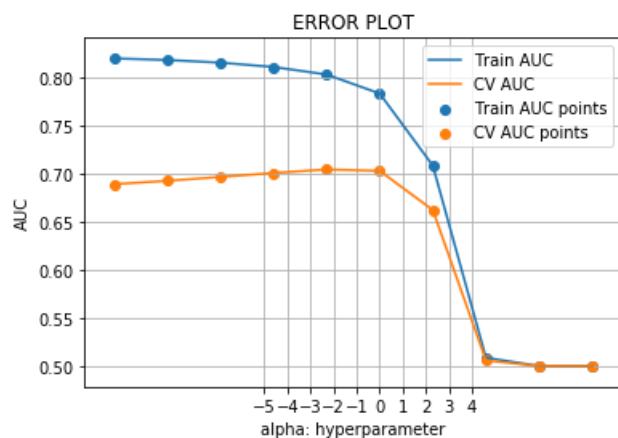
plt.plot(np.log(alpha), train_auc, label='Train AUC')
plt.plot(np.log(alpha), cv_auc, label='CV AUC')

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

plt.xticks(np.log10(alpha))

plt.legend()
```

```
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOT")
plt.grid()
plt.show()
```



In [373]:

```
#obtaining optimal alpha
optimal_alpha = 10;

set1_alpha = optimal_alpha

#using Naive Bayes using optimal alpha
NB = MultinomialNB(alpha=optimal_alpha, class_prior = [0.5, 0.5]);

NB.fit(train_set_1, y_train);

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

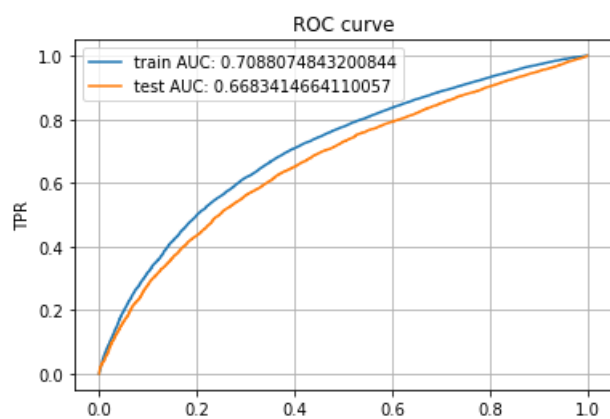
#obtaining fpr, tp, threshold and auc for train and test.
train_fpr, train_tpr, train_threshold = metrics.roc_curve(y_train, y_train_pred)
train_auc = metrics.roc_auc_score(y_train, y_train_pred)

test_fpr, test_tpr, test_threshold = metrics.roc_curve(y_test, y_test_pred)
test_auc = metrics.roc_auc_score(y_test, y_test_pred);

set1_auc = test_auc;

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

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



In [374]:

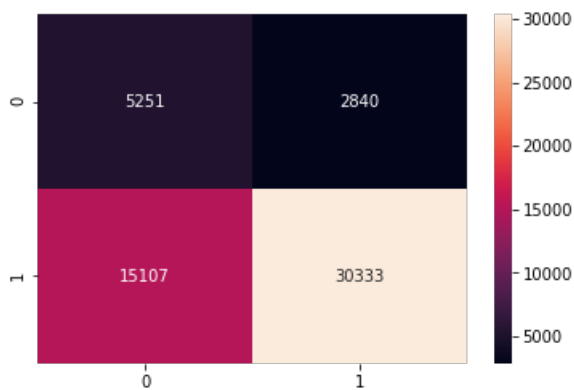
```
import seaborn as sns
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_threshold, train_fpr, train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d")
```

Train confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.4332283408925932 for threshold 0.992

Out[374]:

<matplotlib.axes._subplots.AxesSubplot at 0x2dd05df0e48>

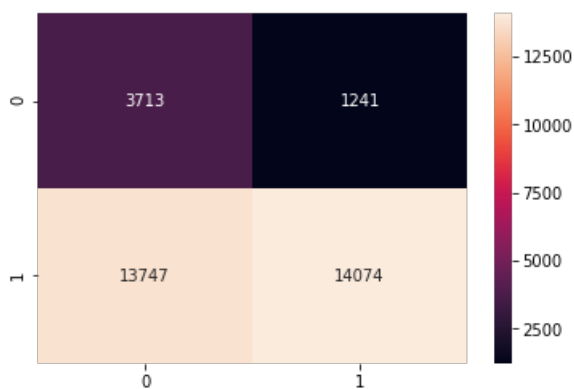


In [375]:

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

Test confusion matrix

the maximum value of $tpr \cdot (1 - fpr)$ 0.39670235568038287 for threshold 1.0



2.4.1.1 Top 10 important features of positive class from SET 1

In [376]:

```
#obtaining log probability of all features of positive class.
pos_features_prob = NB.feature_log_prob_[1, :];

#sortin features of positive class and getting top 10 from it.
features_positive_indices = np.argsort(NB.feature_log_prob_[1, :])[:, :-1][:10]
print(features_positive_indices)

features = [final_features_list[i] for i in features_positive_indices];
print(features)
```

```
[10958 9950 7470 6584 2199 11381 7667 11416 6580 5432]
['students', 'school', 'my', 'learning', 'classroom', 'the', 'not', 'they', 'learn', 'help']
```

2.4.1.2 Top 10 important features of negative class from SET 1

In [377]:

```
#how to find best features in Naive Bayes https://stackoverflow.com/questions/50526898/how-to-get-feature-importance-in-naive-bayes

#obtaining log probability of all features
neg_features_prob = NB.feature_log_prob_[0, :];

#sorting features of negative class and getting top 10 from it
features_negative_indices = np.argsort(NB.feature_log_prob_[0, :])[-10:]
print(features_negative_indices)

features = [final_features_list[i] for i in features_negative_indices];
print(features)

[ 5432 11381 11416 7667 6580 2199 7470 6584 9950 10958]
['help', 'the', 'they', 'not', 'learn', 'classroom', 'my', 'learning', 'school', 'students']
```

2.4.2 Applying Naive Bayes on TFIDF, SET 2

In [378]:

```
#list for containing auc values for train and test for each alpha.
train_auc = [];
cv_auc = [];

#creating list of alpha.
alpha = [0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]

for a in alpha:
    NB = MultinomialNB(alpha=a, class_prior = [0.5, 0.5]);

    #training model using train data.
    NB.fit(train_set_2, y_train);

    y_train_pred = NB.predict_proba(train_set_2)[: , 1];
    y_cv_pred = NB.predict_proba(cv_set_2)[: , 1];

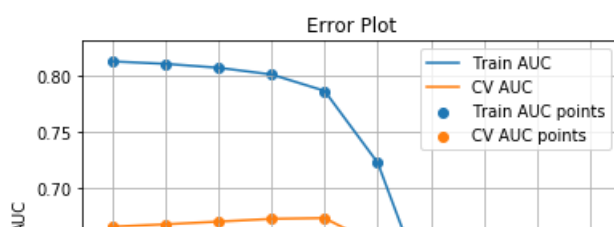
    train_auc.append( metrics.roc_auc_score(y_train, y_train_pred) );
    cv_auc.append( metrics.roc_auc_score(y_cv, y_cv_pred) );

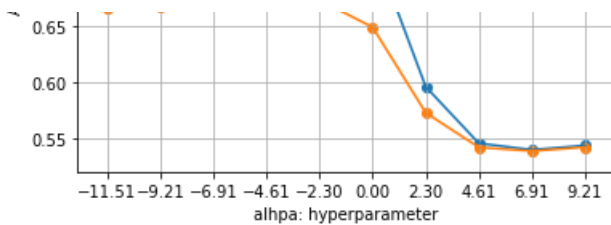
plt.plot(np.log(alpha), train_auc, label="Train AUC");
plt.plot(np.log(alpha), cv_auc, label="CV AUC");

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

plt.xlabel('alpha: hyperparameter');
plt.ylabel('AUC');
plt.title('Error Plot');
plt.xticks(np.log(alpha))
plt.grid();
plt.legend();

plt.show();
```





In [379]:

```
#getting optimal alpha
optimal_alpha = 10;

set2_alpha = optimal_alpha;

#training Naive Bayes using Optimal alpha
NB = MultinomialNB(alpha=optimal_alpha,class_prior = [0.5, 0.5]);

NB.fit(train_set_2, y_train);

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

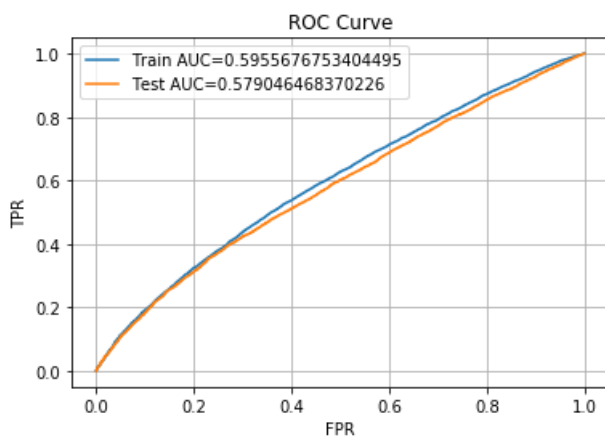
train_fpr, train_tpr, train_treshholds = metrics.roc_curve(y_train, y_train_pred);
test_fpr, test_tpr, test_treshholds = metrics.roc_curve(y_test, y_test_pred);

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

set2_auc = test_auc;

plt.plot(train_fpr, train_tpr, label='Train AUC='+str(train_auc));
plt.plot(test_fpr, test_tpr, label='Test AUC='+str(test_auc));

plt.xlabel('FPR');
plt.ylabel('TPR');
plt.title('ROC Curve')
plt.grid();
plt.legend();
plt.plot();
```

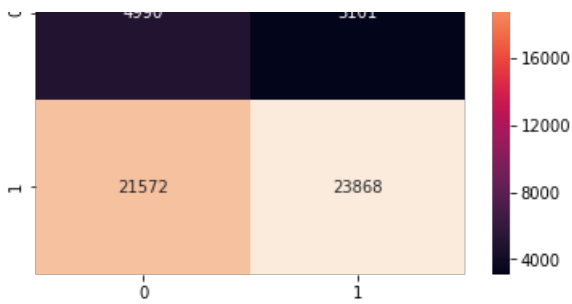


In [380]:

```
#printing confusion matrix for train and test
print("Train confusion matrix")
train_cm = confusion_matrix(y_train, predict(y_train_pred, train_treshholds, train_fpr, train_tpr))
sns.heatmap(train_cm, annot=True, fmt="d");
```

Train confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.32394855786554694 for threshold 0.996

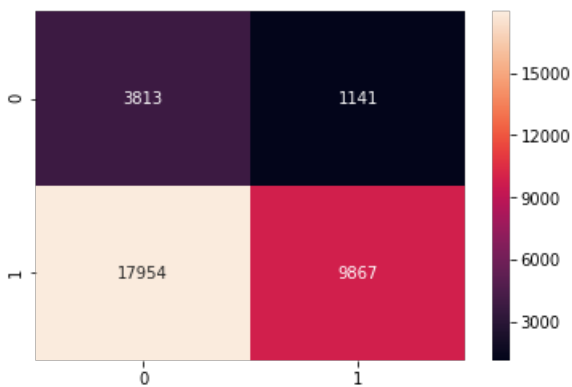




In [381]:

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

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.3075501181445482 for threshold 0.998



2.4.2.1 Top 10 important features of positive class from SET 2

In [382]:

```
#obtaining log probability of all features of positive class.
pos_features_prob = NB.feature_log_prob_[1, :];

#sortin features of positive class and getting top 10 from it.
features_positive_indices = np.argsort(NB.feature_log_prob_[1, :])[-10:]
print(features_positive_indices)

features = [final_features_list[i] for i in features_positive_indices];
print(features)
```

```
[ 87  88  89  95  58  59  94 101 100 102]
['Literature_Writing', 'Mathematics', 'Literacy', 'Ms', 'Math_Science', 'Literacy_Language',
'Mrs', 'teacher_number_of_previously_posted_projects', 'price', 'quantity']
```

2.4.2.2 Top 10 important features of negative class from SET 2

In [383]:

```
#obtaining log probability of all features
neg_features_prob = NB.feature_log_prob_[0, :];

#sorting features of negative class and getting top 10 from it
features_negative_indices = np.argsort(NB.feature_log_prob_[0, :])[-10:]
print(features_negative_indices)

features = [final_features_list[i] for i in features_negative_indices];
print(features)
```

```
[ 87  88  89  95  58  59  94 101 100 102]
```



```
[ 8/ 89 88 95 58 59 94 101 100 102]
['Literature_Writing', 'Literacy', 'Mathematics', 'Ms', 'Math_Science', 'Literacy_Language',
'Mrs', 'teacher_number_of_previously_posted_projects', 'price', 'quantity']
```

3. Conclusions

In [384]:

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

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

table.add_row(['BOW', 'Brute', set1_alpha, set1_auc]);
table.add_row(['TFIDF', 'Brute', set2_alpha, set2_auc]);

print(table)
```

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	10	0.6683414664110057
TFIDF	Brute	10	0.579046468370226

In []: