

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

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

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

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

| Feature                  | Description   |
|--------------------------|---|
| <code>id</code>          | A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502 |
| <code>description</code> | Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25                 |
| <code>quantity</code>    | Quantity of the resource required. <b>Example:</b> 3  |
| <code>price</code>       | Price of the resource required. <b>Example:</b> 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   |
|----------------------------------|---|
| <code>project_is_approved</code> | A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project was not approved, and a value of 1 indicates the project was approved. |

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- `__project_essay_1:` "Introduce us to your classroom"
- `__project_essay_2:` "Tell us more about your students"
- `__project_essay_3:` "Describe how your students will use the materials you're requesting"
- `__project_essay_3:` "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- `__project_essay_1:` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- `__project_essay_2:` "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

```

In [2]: %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

```

## 1.1 Reading Data

```

In [3]: project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')

```

```

In [4]: 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 [5]: 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[5]:

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

## 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"
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e remove)
        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"
            j=j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e remove)
        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]))
```

## 1.5 Preprocessing project grade category

```
In [8]: #Replacing spaces & hyphens in the text of project grade category with underscore
#converting Capital Letters in the string to smaller letters
#Performing a value count of project grade category
# https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
project_data['project_grade_category'].value_counts()
```

```
Out[8]: grades_prek_2    44225
grades_3_5             37137
grades_6_8             16923
grades_9_12            10963
Name: project_grade_category, dtype: int64
```

## 1.2 preprocessing of Teacher prefix

```
In [9]: # check if we have any nan values are there in the column
print(project_data['teacher_prefix'].isnull().values.any())
print("number of nan values", project_data['teacher_prefix'].isnull().values.sum())
```

```
True
number of nan values 3
```

```
In [10]: #Replacing the Nan values with most frequent value in the column
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
```

```
In [11]: #Converting teacher prefix text into smaller case
project_data['teacher_prefix'] = project_data['teacher_prefix'].str.lower()
project_data['teacher_prefix'].value_counts()
```

```
Out[11]: mrs.      57272
ms.       38955
mr.       10648
teacher   2360
dr.        13
Name: teacher_prefix, dtype: int64
```

## Preprocessing State feature

```
In [12]: #Converting states text into smaller case
project_data['school_state'] = project_data['school_state'].str.lower()
project_data['school_state'].value_counts()
```

```
Out[12]: ca      15388
tx       7396
ny       7318
fl       6185
nc       5091
il       4350
ga       3963
sc       3936
mi       3161
pa       3109
in       2620
mo       2576
oh       2467
la       2394
ma       2389
wa       2334
ok       2276
nj       2237
az       2147
va       2045
wi       1827
al       1762
ut       1731
tn       1688
ct       1663
md       1514
nv       1367
ms       1323
ky       1304
or       1242
mn       1208
co       1111
ar       1049
id        693
ia        666
ks        634
nm        557
dc        516
hi        507
me        505
wv        503
nh        348
ak        345
de        343
ne        309
sd        300
ri        285
mt        245
nd        143
wy         98
vt         80
Name: school_state, dtype: int64
```

## 1.3 Text preprocessing

```
In [13]: # 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 [14]: project_data.head(2)
```

Out[14]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | project_submitted_datetime | project_grade |
|---|------------|---------|----------------------------------|----------------|--------------|----------------------------|---------------|
| 0 | 160221     | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | mrs.           | in           | 2016-12-05 13:43:57        | grade         |
| 1 | 140945     | p258326 | 897464ce9ddc600bcd1151f324dd63a  | mr.            | fl           | 2016-10-25 09:22:10        | gr            |

```
In [15]: # printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

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

=====

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. nnnnn

=====

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\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\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one. \r\n\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! nnnnn

=====

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. 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\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\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time. \r\n\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 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 [16]: # 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 [17]: sent = decontracted(project_data['essay'].values[2000])
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\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. 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\n\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 [22]: # after preprocessing
preprocessed_essays[20000]
```

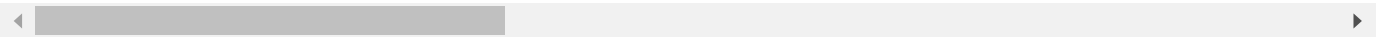
Out[22]: 'my kindergarten students varied disabilities ranging speech language delays cognitive delays gross fine motor delays autism they eager beavers always strive work hardest working past limitations the materials ones i seek students i teach title i school students receive free reduced price lunch despite disabili es limitations students love coming school come eager learn explore have ever felt like ants pants neede d groove move meeting this kids feel time the want able move learn say wobble chairs answer i love devel op core enhances gross motor turn fine motor skills they also want learn games kids not want sit workshe ets they want learn count jumping playing physical engagement key success the number toss color shape ma ts make happen my students forget work fun 6 year old deserves nannan'

## 1.4 Preprocessing of `project\_title`

```
In [23]: # similarly you can preprocess the titles also
project_data.head(2)
```

Out[23]:

|   | Unnamed: 0 | id      | teacher_id                       | teacher_prefix | school_state | project_submitted_datetime | project_grade |
|---|------------|---------|----------------------------------|----------------|--------------|----------------------------|---------------|
| 0 | 160221     | p253737 | c90749f5d961ff158d4b4d1e7dc665fc | mrs.           | in           | 2016-12-05 13:43:57        | grade         |
| 1 | 140945     | p258326 | 897464ce9ddc600bcd1151f324dd63a  | mr.            | fl           | 2016-10-25 09:22:10        | gr            |



```
In [24]: # printing some random project titles.
print(project_data['project_title'].values[54])
print("="*50)
print(project_data['project_title'].values[89])
print("="*50)
print(project_data['project_title'].values[999])
print("="*50)
print(project_data['project_title'].values[11156])
print("="*50)
print(project_data['project_title'].values[89436])
print("="*50)
```

Swim For Life At YMCA!  
=====  
Education Through Technology  
=====  
Focus Pocus  
=====  
Making Math Interactive!  
=====  
Classroom Supplies: Help a New Teacher Organize the Classroom!  
=====

In [25]: *#Removing phrases from the title features*  
`import re`

```
def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)
    phrase = re.sub(r"Gotta", "Got to", 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 [26]: *#Checkingt titles after removing phrases*  
`sent = decontracted(project_data['project_title'].values[89436])`  
`print(sent)`  
`print("="*50)`

Classroom Supplies: Help a New Teacher Organize the Classroom!  
=====

In [27]: *# Remove \\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)`

Classroom Supplies: Help a New Teacher Organize the Classroom!

In [28]: *#Removing numbers & symbols form the titles*  
`sent = re.sub('[^A-Za-z0-9]+', ' ', sent)`  
`print(sent)`

Classroom Supplies Help a New Teacher Organize the Classroom

In [29]: *# <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', 'the`  
`'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 't`  
`'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do`  
`'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while`  
`'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before`  
`'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'aga`  
`'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each',`  
`'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`  
`'ve', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't`  
`"hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'r`  
`"mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'we`  
`'won', "won't", 'wouldn', "wouldn't"]`

```
In [30]: #Combining all the above preprocessed statements
from tqdm import tqdm
preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].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_titles.append(sent.lower().strip())
```

```
In [31]: #checking cleaned text after preprocessing
print(preprocessed_titles[54])
print("="*50)
print(preprocessed_titles[89])
print("="*50)
print(preprocessed_titles[999])
print("="*50)
print(preprocessed_titles[11156])
print("="*50)
print(preprocessed_titles[89436])
```

```
In [32]: project.data.columns
```

```
Out[32]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
               'project_submitted_datetime', 'project_grade_category', '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',
               'clean_categories', 'clean_subcategories', 'essay'],
              dtype='object')
```

```
In [33]: # 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 encoding ", categories_one_hot.shape)
```

```
In [34]: # 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', 'LiteratureWriting', 'Mathematics', 'Literacy']
Shape of matrix after one hot encoding (109248, 30)
```

```
In [35]: # Applying count vectorizer on school state feature & one hot encoding School_state feature
vectorizer = CountVectorizer(binary=True)
school_state_count = vectorizer.fit_transform(project_data['school_state'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", school_state_count.shape)

['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
Shape of matrix after one hot encoding (109248, 51)
```

```
In [36]: #Replacing spaces & hyphens in the text of project grade category with underscore
#converting Capital Letters in the string to smaller letters
#Performing a value count of project grade category
# https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
project_data['project_grade_category'].value_counts()
```

```
Out[36]: grades_prek_2    44225
grades_3_5    37137
grades_6_8    16923
grades_9_12    10963
Name: project_grade_category, dtype: int64
```

```
In [37]: #One hot encoding project grade category feature
vectorizer = CountVectorizer(binary=True)
project_grade_one = vectorizer.fit_transform(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", project_grade_one.shape)

['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
Shape of matrix after one hot encoding (109248, 4)
```

```
In [38]: #One hot encoding the teacher prefix column
vectorizer = CountVectorizer(binary=True)
teacher_prefix_one = vectorizer.fit_transform(project_data['teacher_prefix'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encoding ", teacher_prefix_one.shape)

['dr', 'mr', 'mrs', 'ms', 'teacher']
Shape of matrix after one hot encoding (109248, 5)
```

## 1.5.2 Vectorizing Text data

### 1.5.2.1 Bag of words

```
In [39]: # 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 [40]: # you can vectorize the title also
vectorizer = CountVectorizer(min_df=10)
title_bow = vectorizer.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ",title_bow.shape)
# before you vectorize the title make sure you preprocess it
```

Shape of matrix after one hot encoding (109248, 3328)

### 1.5.2.2 TFIDF vectorizer

```
In [41]: 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)

```
In [42]: # you can vectorize the title also
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
title_tfidf = vectorizer.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ",title_tfidf.shape)
```

Shape of matrix after one hot encoding (109248, 3328)

### 1.5.2.3 Using Pretrained Models: Avg W2V

In [43]:

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

# =====
Output:

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

# =====

words = []
for i in 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[43]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef (https://stackov
erflow.com/a/38230349/4084039\ndef) loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f
= open(gloveFile,\r', encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine =
line.split()\n        word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLin
e[1:]])\n        model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return mo
del\nmodel = loadGloveModel('glove.42B.300d.txt')\n\n# =====\nOutput:\n    \nLo
ading Glove Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n# =====
=====
\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split(' '))\n\nfor i in preproc
od_titles:\n    words.extend(i.split(' '))\nprint("all the words in the coupus", len(words))\nwords =
set(words)\nprint("the unique words in the coupus", len(words))\n\ninter_words = set(model.keys()).inte
rsection(words)\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 = {}\nwords_glov
e = set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\np
rint("word 2 vec length", len(words_courpus))\n\n# stronging variables into pickle files python: htt
p://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport (http://www.je
ssicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport) pickle\nwith open('glo
ve_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```







```
In [53]: #Normalizing price
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(project_data['price'].values.reshape(1, -1))
price_norm = normalizer.transform(project_data['quantity'].values.reshape(1, -1))
print("After vectorizations")
print(price_norm.shape)
```

After vectorizations  
(1, 109248)

```
In [54]: price_norm = price_norm.T
```

```
In [55]: price_norm.shape
```

```
Out[55]: (109248, 1)
```

```
In [56]: #Normalizing quantity
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(project_data['quantity'].values.reshape(1, -1))
quantity_norm = normalizer.transform(project_data['quantity'].values.reshape(1, -1))
print("After vectorizations")
print(quantity_norm.shape)
```

After vectorizations  
(1, 109248)

```
In [57]: quan_norm = quantity_norm.T
quan_norm.shape
```

```
Out[57]: (109248, 1)
```

```
In [58]: # Normalizing teacher previously posted projects
#Normalizing quantity
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(project_data['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
tpp_norm = normalizer.transform(project_data['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
print("After vectorizations")
print(tpp_norm.shape)
```

After vectorizations  
(1, 109248)

```
In [59]: tpp_norm = quantity_norm.T
tpp_norm.shape
```

```
Out[59]: (109248, 1)
```

## 1.5.4 Merging all the following features

```
In [60]: 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 [61]: # 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[61]: (109248, 16663)
```

## Assignment 4: Naive Bayes

### 1. Apply Multinomial Naive Bayes 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 parameter tuning(find best Alpha)

- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- Find the best hyper parameter 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 \(https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html) 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 \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](https://seaborn.pydata.org/generated/seaborn.heatmap.html).



(<https://seaborn.pydata.org/generated/seaborn.heatmap.html>)

### 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 \(http://zetcode.com/python/prettytable/\)](http://zetcode.com/python/prettytable/)



## 1.7 Test Train split

```
In [62]: data = project_data
data.head(5)
```

Out[62]:

|   | 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        | grades          |
| 2 | 21895      | p182444 | 3465aaf82da834c0582ebd0ef8040ca0 | ms.            | az           | 2016-08-31 12:03:56        | grades          |
| 3 | 45         | p246581 | f3cb9bffbba169bef1a77b243e620b60 | mrs.           | ky           | 2016-10-06 21:16:17        | grades          |
| 4 | 172407     | p104768 | be1f7507a41f8479dc06f047086a39ec | mrs.           | tx           | 2016-07-11 01:10:09        | grades          |

```
In [63]: y = data['project_is_approved'].values
data.drop(['project_is_approved'], axis=1, inplace=True)
data.head(1)
```

Out[63]:

|   | 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          |

```
In [64]: X = data
```

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

```
In [65]: # 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
#Splitting data into test & train set
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size = 0.33,stratify=y)
```

```
In [66]: #Splitting training data into training & cross validation sets  
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train,  
                                                stratify=y_train,  
                                                test_size = 0.33)
```

## 1.3 Text preprocessing

```
In [67]: # printing some random reviews
print(X_train['essay'].values[0])
print("="*50)
print(X_train['essay'].values[100])
print("="*50)
print(X_train['essay'].values[300])
print("="*50)
print(X_train['essay'].values[5000])
print("="*50)
print(X_train['essay'].values[20000])
print("="*50)
```

I am a Kindergarten ESL (English as a Second Language) teacher at a small East Texas Title I school with over 80% of our population labeled as economically disadvantaged. This year I have a group of 19 students, all at various economic and academic levels. Our school is designed to love and nurture our students while providing them with a positive educational experience. \r\n\r\nWe are an active group of learners and I strive to keep them focused and having fun. I am constantly exploring new meaningful activities for my students to keep their little hands and minds busy. My room is a student lead learning environment and these little ones are always ready to learn.\r\n\r\n\r\n\r\nWobble chairs are fun and safe stools that are perfect for my overly active kids. They provide a safe learning, sitting environment for students with restlessness and extra energy. A wobble chair means a student doesn't have to sit still. Since research shows that students learn better when they can pay attention and focus, we really need these chairs! \r\n\r\nWobble Chairs transform the boring conventional seat into a chair that really ROCKS! Some students need a little more help than others when it comes to paying attention, focusing and learning. Wobble chairs allow children to move and flex without leaving their seat.nannan

=====

My students are all new to personal finance with many of them never having any type of finance discussion at home. This class is an elective so I try to make it as fun as possible by using interactive lessons. The majority of my students get free or reduced-price lunch and I hope to break the cycle of poverty for some of them by making sure they know all the components of what the world holds for them after they graduate whether they go to college or not. My students are learning to create spending plans during the budgeting portion of the personal finance class. One way people create spending plans is to use Financial Technology to make it easier for every day life. My students need to explore different apps available to consumers such as YNAB and Mint. The best way to see these apps is through a tablet. My students will share this tablet during group time to compare and contrast different financial technology avenues. The iPad will get students excited about creating budgets and hopefully will allow them to continue using these techniques once they have jobs and money to spend.nannan

=====

My students come to school eager to learn and excited about what the day has in store. \r\n\r\nI teach at a Title 1 school in Oklahoma City where 100% of our students receive free breakfast and lunch, for some this will be the only meals they eat. \r\n\r\nMany of my students are ELL students, where Spanish is their first language. My students may face many obstacles but they don't let that stop them. They want to be in school because they want to learn so that they can be successful one day.\r\n\r\nHands on activities are the best for all students, especially kindergarteners. Having these stem items in my classroom will allow my students to explore science through hands on learning and exploration.\r\n\r\nThe art of teaching is the art of assisting discovery. - Mark Van Doren. \r\n\r\nWater, magnets, and how things move are very intriguing to young children. The stem science station will provide my students the opportunity to explore magnets, explore objects that sink and float, and explore motion. \r\n\r\nThe fairy tale stem kit will give my students the chance to see fairy tales in a whole new way. They will be able to retell the fairytales through hands on activities.nannan

=====

Majority, if not all, of our students come from low-socioeconomic backgrounds. All students are eligible for free breakfast and lunch. Yet, once they come in our school they are the richest students around. \r\n\r\n\r\nOur boys and girls become the future of our community regardless of their backgrounds.\r\n\r\n\r\nOur students work their hardest through cold winters at home, through hungry nights, and through what many adults couldn't face on a daily basis, the what if's. It's only fair we try to work our hardest too.\r\n\r\n\r\nMy students will be able to create concrete objects by using a 3Doodler pen in various aspects in the classroom. My student's will be able to engage in reading by using the pen to create objects found in the story, they will be able to use it to create 3D figures in math, and even depict various concepts in science.\r\n\r\n\r\nMy students will enjoy drawing not with paper and pencil, but with a very evolving aspect in our generation, technology. Who would have thought that we would be able to create these vocabulary words that are necessary for learning into concrete objects they can actually touch?!nannan

=====

My fourth grade students enter the classroom knowing their school year will be full of excitement. They know that Mr. B. makes it a priority to make learning fun and relevant. Throughout the year I try and apply outdoor sports within the curriculum. We take what we learned in class and apply it to team sports in a way that is engaging. The outdoor events help keep the students motivated while retaining the basics of what was taught in class through constant repetition. My students always try and get me to take them outside. I usually take them to read and write outside along with doing various science experiments too. However, most kids just want to play! So we came up with a great compromise. \r\n\r\n\r\nMy students said we need to "Sharpen the Saw" based on our Leader in Me training from Franklin Covey's 7 Habits of Highly Effective People. I couldn't argue with their reasoning. We began playing Multiplication Baseball. Basically it is whiffleball using multiplication facts as my pitches. The kids have to get the problem right before I pitch the ball. They love it! So we decided to take it a step further. This is

where Academic Athletes was born. It is my goal to form an after-school health and sports program while integrating academics as well. Kids will learn about nutrition, team sports, and physical exercise through various units. They will know that it is OK to PLAY!!nannan  
=====

```
In [68]: # 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 [69]: 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 [70]: # \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 [71]: #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 many 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 fun a 6 year old deserves fun

```
In [72]: # 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', 'the', \
'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', \
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', \
'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', \
'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', \
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", \
"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", 'won', \
"won't", 'wouldn', "wouldn't"]
```

## Preprocessing for Train Data

```
In [73]: # Combining all the above students
from tqdm import tqdm
preprocessed_essays_xtr = []
# tqdm is for printing the status bar
for sentence in tqdm(X_train['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
    preprocessed_essays_xtr.append(sent.lower().strip())
```

100%|███| 49041/49041 [00:32<00:00  
0, 1498.16it/s]

```
In [74]: # after preprocessing
preprocessed_essays_xtr[300]
```

```
Out[74]: 'students come school eager learn excited day store teach title 1 school oklahoma city 100 students rece
ive free breakfast lunch meals eat many students ell students spanish first language students may face m
any obstacles not let stop want school want learn successful one day hands activities best students espe
cially kindergarteners stem items class room allow students explore science hands leaning exploration ar
t teaching art assisting discovery mark van doren water magents things move intriguing young children st
em science station provide students opportunity explore magnets explore objects sink float explore motio
n fairy tale stem kit give students chance see fairy tales whole new way able retell fairytales hands ac
tivities nannan'
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 24155/24155 [00:16<00:00  
0, 1485.88it/s]
```

Out[76]: 'students come walks life many different experiences come range prior knowledge skills beliefs backgrou  
ds languages half class speaks spoken home another language would like offer students ability learn life  
despite differences school distinctive fact despite differences students teachers parents family science  
important classroom answer questions young minds children exposed science able learn must communicate ot  
hers listen others patience set minds knowing help solve problems within world students also learn impor  
tance life cycles water animals plants life would different not donation project would teach students ev  
en though may small able help world big way nannan'

```
100%|██████████████████████████████████████████████████████████████████████████████| 36052/36052 [00:22<00:00  
0, 1578.74it/s]
```

```
Out[78]: 'participate no excuses university class college bound every morning start day eagles chant school also
heavily focused arts academy creative expression students busily preparing middle school want everything
need successful enjoy positive role models younger students campus school serves diverse population majo
rity students low income families always arrive eager learn full curiosity desire showcase knowledge new
creative ways students love athletics morning full learning ready play teachers school rotate running di
fferent sporting events tournaments lunch students materials project go towards lunchtime soccer games b
asketball tournaments kick ball tournaments football games lunchtime yoga classes need equipment expand
number students participate types tournaments hold students must work completed exemplary behavior parti
cipate team tournaments contributing project promoting staying active learning teamwork promoting good c
itizenship nannan'
```

## 1.4 Preprocessing of `project\_title`

```
In [79]: # similarly you can preprocess the titles also
#printing random titles
print(data['project_title'].values[49])
print("="*50)
print(data['project_title'].values[89])
print("="*50)
print(data['project_title'].values[999])
print("="*50)
print(data['project_title'].values[11156])
print("="*50)
print(data['project_title'].values[20000])
print("="*50)
```

Rainy Day Run Around!  
=====

Education Through Technology  
=====

Focus Pocus  
=====

Making Math Interactive!  
=====

We Need To Move It While We Input It!  
=====

```
In [80]: #Removing phrases from the title features
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)
    phrase = re.sub(r"Gotta", "Got to", 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 [81]: #Checkingt titles after removing phrases
sent = decontracted(project_data['project_title'].values[89436])
print(sent)
print("="*50)
```

Classroom Supplies: Help a New Teacher Organize the Classroom!  
=====

```
In [82]: # Remove \\r \\n \\t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\\r', ' ')
sent = sent.replace('\\\\', ' ')
sent = sent.replace('\\n', ' ')
print(sent)
```

Classroom Supplies: Help a New Teacher Organize the Classroom!





```
In [91]: # 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
```

```
#We use fit only for train data
vectorizer_state = CountVectorizer(binary=True)
vectorizer_state.fit(X_train['school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_ohe = vectorizer_state.transform(X_train['school_state'].values)
X_cv_state_ohe = vectorizer_state.transform(X_cv['school_state'].values)
X_test_state_ohe = vectorizer_state.transform(X_test['school_state'].values)

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer_state.get_feature_names())
print("=*75")
```

```
After vectorizations
(49041, 51) (49041,)
(24155, 51) (24155,)
(36052, 51) (36052,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks',
'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm', 'nv', 'n
y', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
=====
```

## 2.2.2 One hot encoding the categorical features : teacher\_prefix

```
In [92]: # we use count vectorizer to convert the values into one
#We use fit only for train data
vectorizer_tp = CountVectorizer(binary=True)
vectorizer_tp.fit(X_train['teacher_prefix'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer_tp.transform(X_train['teacher_prefix'].values)
X_cv_teacher_ohe = vectorizer_tp.transform(X_cv['teacher_prefix'].values)
X_test_teacher_ohe = vectorizer_tp.transform(X_test['teacher_prefix'].values)

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer_tp.get_feature_names())
print("=*50")
```

```
After vectorizations
(49041, 5) (49041,)
(24155, 5) (24155,)
(36052, 5) (36052,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
=====
```

## 2.2.3 One hot encoding the categorical features : grades

```
In [93]: #Replacing spaces & hyphens in the text of project grade category with underscore
#converting Capital Letters in the string to smaller letters
#Performing a value count of project grade category
# https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
project_data['project_grade_category'].value_counts()
```

```
Out[93]: grades_prek_2      44225
grades_3_5      37137
grades_6_8      16923
grades_9_12     10963
Name: project_grade_category, dtype: int64
```

```
In [94]: #We use fit only for train data
vectorizer_grade = CountVectorizer()
vectorizer_grade.fit(X_train['project_grade_category'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer_grade.transform(X_train['project_grade_category'].values)
X_cv_grade_ohe = vectorizer_grade.transform(X_cv['project_grade_category'].values)
X_test_grade_ohe = vectorizer_grade.transform(X_test['project_grade_category'].values)

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer_grade.get_feature_names())
print("=*70")
```

```
After vectorizations
(49041, 4) (49041,)
(24155, 4) (24155,)
(36052, 4) (36052,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
=====
```

## 2.2.4 One hot encoding the categorical features : project subject category

```
In [95]: #We use fit only for train data
vectorizer_category = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), binary=True)
vectorizer_category.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_cat_ohe = vectorizer_category.transform(X_train['clean_categories'].values)
X_cv_cat_ohe = vectorizer_category.transform(X_cv['clean_categories'].values)
X_test_cat_ohe = vectorizer_category.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_cat_ohe.shape, y_train.shape)
print(X_cv_cat_ohe.shape, y_cv.shape)
print(X_test_cat_ohe.shape, y_test.shape)
print(vectorizer_category.get_feature_names())
print("=*70")
```

```
After vectorizations
(49041, 9) (49041,)
(24155, 9) (24155,)
(36052, 9) (36052,)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
=====
```

## 2.2.5 One hot encoding the categorical features : project subject sub-category

```
In [96]: #We use fit only for train data
vectorizer_subcat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), binary=True)
vectorizer_subcat.fit(X_train['clean_subcategories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_subcat_ohe = vectorizer_subcat.transform(X_train['clean_subcategories'].values)
X_cv_subcat_ohe = vectorizer_subcat.transform(X_cv['clean_subcategories'].values)
X_test_subcat_ohe = vectorizer_subcat.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_subcat_ohe.shape, y_train.shape)
print(X_cv_subcat_ohe.shape, y_cv.shape)
print(X_test_subcat_ohe.shape, y_test.shape)
print(vectorizer_subcat.get_feature_names())
print("=*70")
```

```
After vectorizations
(49041, 30) (49041,)
(24155, 30) (24155,)
(36052, 30) (36052,)
['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']
=====
```

## 1.5.2 Vectorizing Text data

### i) BoW encoding

#### 1.5.2.1 Bag of words on Essay Feature

```
In [97]: # We are considering only the words which appeared in at least 10 documents(rows or projects).
#Applying BoW on essays feature
#Considering only the words which appear atleast in 10 documents or reviews
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("=*100")

from sklearn.feature_extraction.text import CountVectorizer
vectorizer_essay_bow = CountVectorizer(min_df=10)
vectorizer_essay_bow.fit(preprocessed_essays_xtr) # fitting only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer_essay_bow.transform(preprocessed_essays_xtr)
X_cv_essay_bow = vectorizer_essay_bow.transform(preprocessed_essays_xcv)
X_test_essay_bow = vectorizer_essay_bow.transform(preprocessed_essays_xte)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("=*100")
```

```
(49041, 19) (49041,)
(24155, 19) (24155,)
(36052, 19) (36052,)
=====
After vectorizations
(49041, 12056) (49041,)
(24155, 12056) (24155,)
(36052, 12056) (36052,)
=====
```



### 1.5.2.2 Bag of words on Project Title feature

```
In [98]: # you can vectorize the title also
# before you vectorize the title make sure you preprocess it
#Applying BoW on project titles feature
#Considering only the words which appear atleast in 10 documents or reviews
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("=*100)

from sklearn.feature_extraction.text import CountVectorizer
vectorizer_titles_bow = CountVectorizer(min_df=10)
vectorizer_titles_bow.fit(preprocessed_titles_xtr) # fitting only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_titles_bow = vectorizer_titles_bow.transform(preprocessed_titles_xtr)
X_cv_titles_bow = vectorizer_titles_bow.transform(preprocessed_titles_xcv)
X_test_titles_bow = vectorizer_titles_bow.transform(preprocessed_titles_xte)

print("After vectorizations")
print(X_train_titles_bow.shape, y_train.shape)
print(X_cv_titles_bow.shape, y_cv.shape)
print(X_test_titles_bow.shape, y_test.shape)
print("=*100)

(49041, 19) (49041,)
(24155, 19) (24155,)
(36052, 19) (36052,)
=====
After vectorizations
(49041, 2103) (49041,)
(24155, 2103) (24155,)
(36052, 2103) (36052,)
=====
```

### ii) TFIDF Vectorization

TFIDF vectorizer on essay feature

```
In [99]: #Applying TF-IDF on essays feature
#Considering only the words which appear atleast in 10 documents or reviews
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_essay_tfidf = TfidfVectorizer(min_df=10)
vectorizer_essay_tfidf.fit(preprocessed_essays_xtr) # fitting only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer_essay_tfidf.transform(preprocessed_essays_xtr)
X_cv_essay_tfidf = vectorizer_essay_tfidf.transform(preprocessed_essays_xcv)
X_test_essay_tfidf = vectorizer_essay_tfidf.transform(preprocessed_essays_xte)

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
(49041, 19) (49041,)
(24155, 19) (24155,)
(36052, 19) (36052,)
```

```
=====
After vectorizations
```

```
(49041, 12056) (49041,)
(24155, 12056) (24155,)
(36052, 12056) (36052,)
```

## TFIDF on Project Title feature

```
In [100]: #Applying Tfidf on project titles feature
#Considering only the words which appear atleast in 10 documents or reviews
print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_title = TfidfVectorizer(min_df=10)
vectorizer_tfidf_title.fit(preprocessed_titles_xtr) # fitting only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_titles_tfidf = vectorizer_tfidf_title.transform(preprocessed_titles_xtr)
X_cv_titles_tfidf = vectorizer_tfidf_title.transform(preprocessed_titles_xcv)
X_test_titles_tfidf = vectorizer_tfidf_title.transform(preprocessed_titles_xte)

print("After vectorizations")
print(X_train_titles_tfidf.shape, y_train.shape)
print(X_cv_titles_tfidf.shape, y_cv.shape)
print(X_test_titles_tfidf.shape, y_test.shape)
print("="*100)
```

```
(49041, 19) (49041,)
(24155, 19) (24155,)
(36052, 19) (36052,)
```

```
=====
After vectorizations
```

```
(49041, 2103) (49041,)
(24155, 2103) (24155,)
(36052, 2103) (36052,)
```

## 1.5.3 Vectorizing Numerical features

### For Price feature

```
In [101]: from sklearn.preprocessing import Normalizer
price_normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
price_normalizer.fit(X_train['price'].values.reshape(1,-1))

X_train_price_norm = price_normalizer.transform(X_train['price'].values.reshape(1,-1))
X_cv_price_norm = price_normalizer.transform(X_cv['price'].values.reshape(1,-1))
X_test_price_norm = price_normalizer.transform(X_test['price'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("=*100)
```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

```
In [102]: X_train_price_norm = X_train_price_norm.T
X_cv_price_norm = X_cv_price_norm.T
X_test_price_norm = X_test_price_norm.T

print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("=*100)
```

(49041, 1) (49041,)

(24155, 1) (24155,)

(36052, 1) (36052,)

=====

### For Quantity

```
In [103]: #Normalizing quantity
from sklearn.preprocessing import Normalizer
quan_normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
quan_normalizer.fit(X_train['quantity'].values.reshape(1,-1))

X_train_quantity_norm = quan_normalizer.transform(X_train['quantity'].values.reshape(1,-1))
X_cv_quantity_norm = quan_normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
X_test_quantity_norm = quan_normalizer.transform(X_test['quantity'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("=*100")
```

```
After vectorizations
(1, 49041) (49041,)
(1, 24155) (24155,)
(1, 36052) (36052,)
=====
```

```
In [104]: X_train_quantity_norm = X_train_quantity_norm.T
X_cv_quantity_norm = X_cv_quantity_norm.T
X_test_quantity_norm = X_test_quantity_norm.T

print("Final Matrix")
print(X_train_quantity_norm.shape, y_train.shape)
print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("=*100")
```

```
Final Matrix
(49041, 1) (49041,)
(24155, 1) (24155,)
(36052, 1) (36052,)
=====
```

## For teacher previously posted projects

```
In [105]: # Normalizing teacher previously posted projects
from sklearn.preprocessing import Normalizer
tpp_normalizer = Normalizer()
# normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
tpp_normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))

X_train_tpp_norm = tpp_normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_cv_tpp_norm = tpp_normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_test_tpp_norm = tpp_normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_tpp_norm.shape, y_train.shape)
print(X_cv_tpp_norm.shape, y_cv.shape)
print(X_test_tpp_norm.shape, y_test.shape)
print("=*100")
```

```
After vectorizations
(1, 49041) (49041,)
(1, 24155) (24155,)
(1, 36052) (36052,)
=====
```

```
In [106]: X_train_tpp_norm = X_train_tpp_norm.T
X_cv_tpp_norm = X_cv_tpp_norm.T
X_test_tpp_norm = X_test_tpp_norm.T

print(X_train_tpp_norm.shape, y_train.shape)
print(X_cv_tpp_norm.shape, y_cv.shape)
print(X_test_tpp_norm.shape, y_test.shape)
print("=*100)
```

```
(49041, 1) (49041,)
(24155, 1) (24155,)
(36052, 1) (36052,)
=====
```

## 1.5.4 Merging Numerical & Categorical features

- we need to merge all the numerical vectors & catogorical features

```
In [107]: # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_numcat = hstack((X_train_state_ohe, X_train_teacher_ohe, X_train_grade_ohe, X_train_price_norm, X_train_cat_ohe, X_train_tpp_norm))
X_cv_numcat = hstack((X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_price_norm, X_cv_cat_ohe, X_cv_tpp_norm))
X_te_numcat = hstack((X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe, X_test_price_norm, X_test_cat_ohe, X_test_tpp_norm))

print("Final Data matrix")
print(X_tr_numcat.shape, y_train.shape)
print(X_cv_numcat.shape, y_cv.shape)
print(X_te_numcat.shape, y_test.shape)
print("=*100)
```

```
Final Data matrix
(49041, 102) (49041,)
(24155, 102) (24155,)
(36052, 102) (36052,)
=====
```

```
In [108]: """
# 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[108]: '\n# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039\nfrom (https://stackoverflow.com/a/19710648/4084039\nfrom) scipy.sparse import hstack\n# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)\nX = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))\nX.shape'
```

```
In [109]: # 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
```

## 2.4.1 Applying Naive Bayes on BOW, SET 1

Consider Set 1 :- categorical, numerical features + project\_title(BOW) + preprocessed\_essay (BOW)

```
In [110]: # Please write all the code with proper documentation
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_set1 = hstack((X_train_essay_bow, X_train_titles_bow, X_tr_numcat)).tocsr()
X_cv_set1 = hstack((X_cv_essay_bow, X_cv_titles_bow, X_cv_numcat)).tocsr()
X_te_set1 = hstack((X_test_essay_bow, X_test_titles_bow, X_te_numcat)).tocsr()

print("Final Data matrix")
print(X_tr_set1.shape, y_train.shape)
print(X_cv_set1.shape, y_cv.shape)
print(X_te_set1.shape, y_test.shape)
print("=*100)
```

```
Final Data matrix
(49041, 14261) (49041,)
(24155, 14261) (24155,)
(36052, 14261) (36052,)
=====
```

```
In [111]: def batch_predict(clf, data):
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
# consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
# in this for loop we will iterate until the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
# we will be predicting for the last data points
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

```
In [112]: import matplotlib.pyplot as plt
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
import math

alpha = [0.00001, 0.00007, 0.00012, 0.00019, 0.001, 0.008, 0.04, 0.09, 0.1, 0.5, 3, 9, 19, 58, 99, 825, 2000]

train_auc = []
cv_auc = []
a = []
b = []

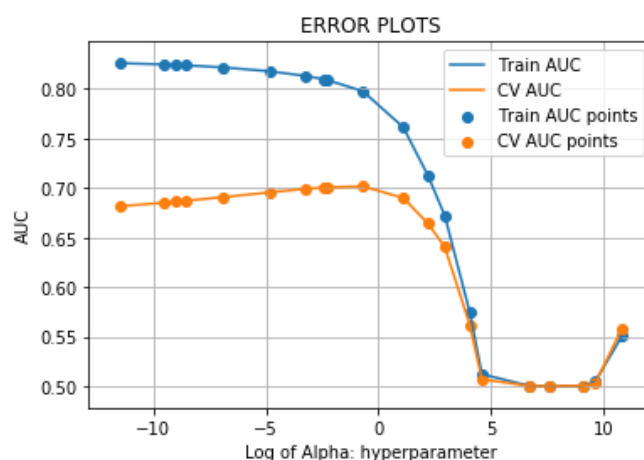
for i in tqdm(alpha):
    naive = MultinomialNB(alpha=i, class_prior=[0.5,0.5])
    naive.fit(X_tr_set1, y_train)
    y_train_pred = batch_predict(naive, X_tr_set1)
    y_cv_pred = batch_predict(naive, X_cv_set1)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

log_alpha = []
for z in tqdm(alpha):
    y= math.log(z)
    log_alpha.append(y)

log_alpha = np.array(log_alpha)
alpha = np.array(alpha)

plt.plot(log_alpha, train_auc, label='Train AUC')
plt.plot(log_alpha, cv_auc, label='CV AUC')
plt.scatter(log_alpha, train_auc, label='Train AUC points')
plt.scatter(log_alpha, cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log of Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 20/20 [00:05<0
0:00, 5.11it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 20/20 [00:00<0:0
0, 20006.22it/s]
```



```
In [113]: # from the error plot we choose alpha such that, we will have maximum AUC on cv data and gap between the t
best_a = 0.5
```

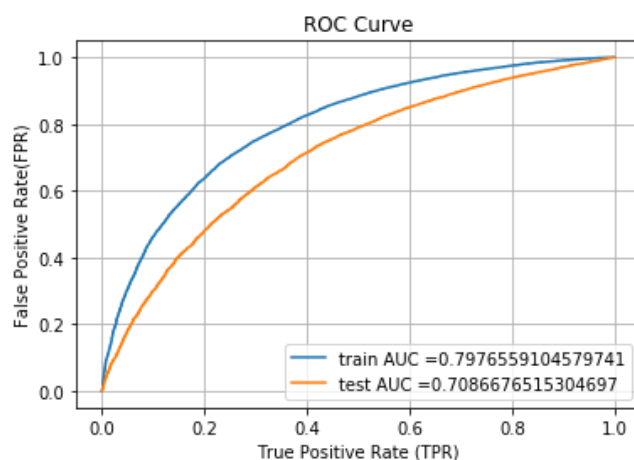
```
In [114]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

naive = MultinomialNB(alpha=best_a, class_prior=[0.5,0.5])
naive.fit(X_tr_set1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(naive, X_tr_set1)
y_test_pred = batch_predict(naive, X_te_set1)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate (TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



```
In [115]: # we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the Least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(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
```

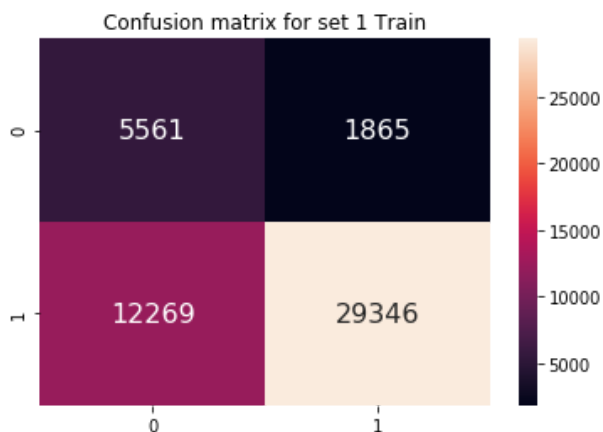


```
In [117]: print("=*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, train_fpr, train_tpr)))
```

```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.52807664968067 for threshold 0.575
[[ 5561 1865]
 [12269 29346]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.52807664968067 for threshold 0.575
[[ 3465 1994]
 [ 9774 20819]]
```

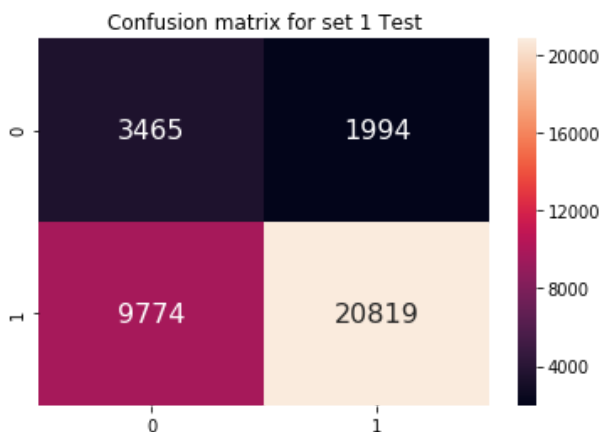
```
In [118]: #How to plot confusion matrix using heat map - https://stackoverflow.com/questions/35572000/how-can-i-plot
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
array = [[5561,1865],
         [12269,29346]]
df_cm = pd.DataFrame(array, index = [i for i in "01"],
                     columns = [i for i in "01"])
plt.figure
plt.title('Confusion matrix for set 1 Train')
sn.heatmap(df_cm, annot=True, annot_kws={"size":16}, fmt='g')
```

Out[118]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1e436530358>



```
In [119]: #How to plot confusion matrix using heat map - https://stackoverflow.com/questions/35572000/how-can-i-plot
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
array = [[3465,1994],
         [9774,20819]]
df_cm = pd.DataFrame(array, index = [i for i in "01"],
                     columns = [i for i in "01"])
plt.figure
plt.title('Confusion matrix for set 1 Test')
sn.heatmap(df_cm, annot=True, annot_kws={"size":16}, fmt='g')
```

Out[119]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1e4921ffeb8>



#### 2.4.1.1 Top 10 important features of positive class from SET 1

```
In [120]: #how to calculate feature probability score in naive bayes python - https://www.datacamp.com/community/tutorials/naive-bayes-classification
#how to calculate feature probability score in naive bayes python- https://towardsdatascience.com/naive-bayes-classification
nb_bow = MultinomialNB(alpha = 0.5)
nb_bow.fit(X_tr_set1, y_train)
```

Out[120]: MultinomialNB(alpha=0.5, class\_prior=None, fit\_prior=True)

```
In [121]: bow_features_probs = []

for x in range(14261) :
    y = naive.feature_log_prob_[1,x]
    bow_features_probs.append(y)

len(bow_features_probs)
```

Out[121]: 14261

```
In [122]: bow_features_names = []
```

```
In [123]: for _ in vectorizer_state.get_feature_names() :
    bow_features_names.append(_)
```

```
In [124]: for _ in vectorizer_tp.get_feature_names() :
    bow_features_names.append(_)
```

```
In [125]: for _ in vectorizer_grade.get_feature_names() :
    bow_features_names.append(_)
```

```
In [126]: for _ in vectorizer_category.get_feature_names() :
    bow_features_names.append(_)
```

```
In [127]: for _ in vectorizer_subcat.get_feature_names() :
    bow_features_names.append(_)
```

```
In [128]: for _ in vectorizer_essay_bow.get_feature_names() :  
         bow_features_names.append(_)
```

```
In [129]: for _ in vectorizer_titles_bow.get_feature_names() :  
         bow_features_names.append(_)
```

```
In [130]: bow_features_names.append('price')  
         bow_features_names.append('quantity')  
         bow_features_names.append('teacher_number_of_previously_posted')
```

```
In [131]: len(bow_features_names)
```

Out[131]: 14261

```
In [132]: final_bow_features_set1_pos = pd.DataFrame({'feature_prob_estimates' : bow_features_probs, 'feature_names'  
         < [REDACTED]
```

```
In [133]: final_bow_features_set1_pos.sort_values(by = ['feature_prob_estimates'], ascending = True, inplace=True)
```

```
In [134]: final_bow_features_set1_pos.head(10)
```

Out[134]:

|       | feature_prob_estimates | feature_names |
|-------|------------------------|---------------|
| 14246 | -16.297731             | year          |
| 14230 | -16.297731             | working       |
| 14247 | -16.297731             | yearbook      |
| 14220 | -16.297731             | wobbles       |
| 14236 | -16.297731             | worms         |
| 14245 | -16.297731             | ye            |
| 14221 | -16.297731             | wobbling      |
| 14222 | -16.297731             | wobbly        |
| 14244 | -16.297731             | xylophone     |
| 14223 | -16.297731             | wonder        |

#### 2.4.1.2 Top 10 important features of negative class from SET 1

```
In [135]: bow_features_neg_probs1 = []  
         for z in range(14261) :  
             z1 = naive.feature_log_prob_[0,z]  
             bow_features_neg_probs1.append(z1)
```

```
In [136]: final_features_bow_set1_neg = pd.DataFrame({'feature_prob_estimates' : bow_features_neg_probs1, 'feature_names'  
         < [REDACTED]
```

```
In [137]: final_features_bow_set1_neg.sort_values(by = ['feature_prob_estimates'], ascending=True, inplace=True)
```

```
In [138]: final_features_bow_set1_neg.head(10)
```

```
Out[138]:
```

|       | feature_prob_estimates | feature_names |
|-------|------------------------|---------------|
| 11122 | -14.524942             | touchpad      |
| 8206  | -14.524942             | portland      |
| 9927  | -14.524942             | skin          |
| 2703  | -14.524942             | cozy          |
| 13619 | -14.524942             | pretend       |
| 6860  | -14.524942             | menu          |
| 13618 | -14.524942             | presentations |
| 2692  | -14.524942             | cousins       |
| 796   | -14.524942             | appalachian   |
| 5145  | -14.524942             | hardware      |

## 2.4.2 Applying Naive Bayes on TFIDF, SET 2

Consider Set 2 :- categorical, numerical features + project\_title(TFIDF) + preprocessed\_essay (TFIDF)

```
In [139]: # Please write all the code with proper documentation
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_set2 = hstack((X_train_essay_tfidf, X_train_titles_tfidf, X_tr_numcat)).tocsr()
X_cv_set2 = hstack((X_cv_essay_tfidf, X_cv_titles_tfidf, X_cv_numcat)).tocsr()
X_te_set2 = hstack((X_test_essay_tfidf, X_test_titles_tfidf, X_te_numcat)).tocsr()

print("Final Data matrix")
print(X_tr_set2.shape, y_train.shape)
print(X_cv_set2.shape, y_cv.shape)
print(X_te_set2.shape, y_test.shape)
print("=*100")
```

Final Data matrix

```
(49041, 14261) (49041,)
(24155, 14261) (24155,)
(36052, 14261) (36052,)
```

=====

```
In [140]: def batch_predict(clf, data):
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

    y_data_pred = []
    tr_loop = data.shape[0] - data.shape[0]%1000
    # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
    # in this for loop we will iterate until the last 1000 multiplier
    for i in range(0, tr_loop, 1000):
        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
    # we will be predicting for the last data points
    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred
```

```
In [141]: import matplotlib.pyplot as plt
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
import math

alpha = [0.00001, 0.00007, 0.00012, 0.00019, 0.001, 0.008, 0.04, 0.09, 0.1, 0.5, 3, 9, 19, 58, 99, 825, 2000]

train_auc = []
cv_auc = []
a = []
b = []

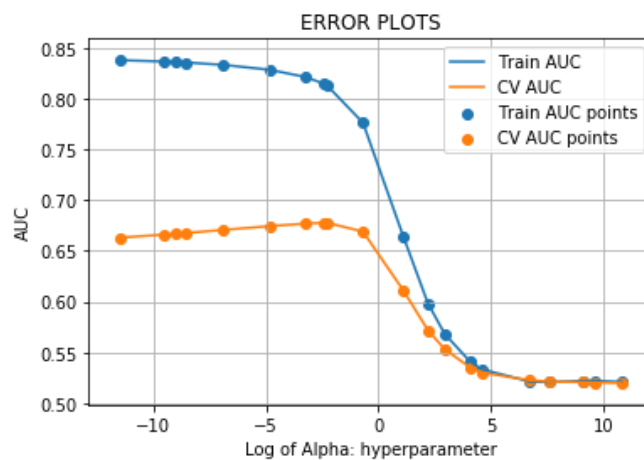
for i in tqdm(alpha):
    naive = MultinomialNB(alpha=i, class_prior=[0.5,0.5])
    naive.fit(X_tr_set2, y_train)
    y_train_pred = batch_predict(naive, X_tr_set2)
    y_cv_pred = batch_predict(naive, X_cv_set2)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
    a.append(y_train_pred)
    b.append(y_cv_pred)

log_alpha = []
for z in tqdm(alpha):
    y= math.log(z)
    log_alpha.append(y)

log_alpha = np.array(log_alpha)
alpha = np.array(alpha)

plt.plot(log_alpha, train_auc, label='Train AUC')
plt.plot(log_alpha, cv_auc, label='CV AUC')
plt.scatter(log_alpha, train_auc, label='Train AUC points')
plt.scatter(log_alpha, cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log of Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 20/20 [00:04<0
0:00, 5.09it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 20/20 [00:00<00:0
0, 19987.15it/s]
```



```
In [142]: # from the error plot we choose alpha such that, we will have maximum AUC on cv data and gap between the t
best_a = 0.5
```

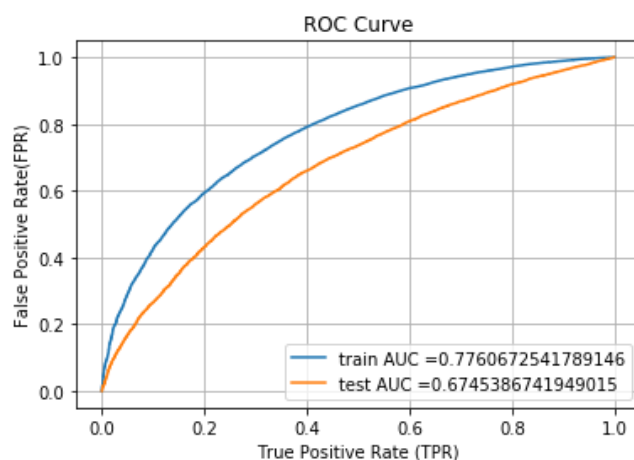
```
In [143]: # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
from sklearn.metrics import roc_curve, auc

naive = MultinomialNB(alpha=best_a, class_prior=[0.5,0.5])
naive.fit(X_tr_set2, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(naive, X_tr_set2)
y_test_pred = batch_predict(naive, X_te_set2)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate (TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



```
In [144]: # we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the Least fpr
def predict(proba, threshold, fpr, tpr):

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

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

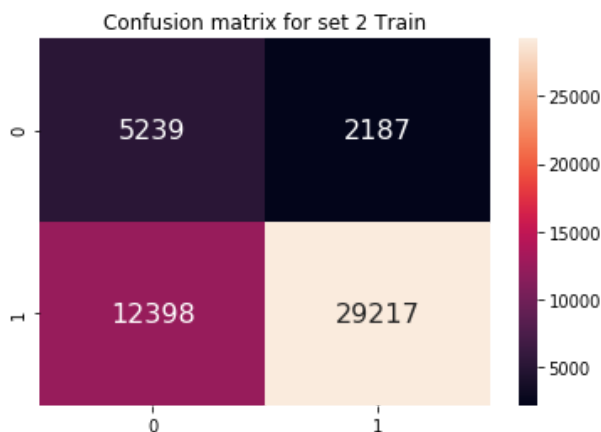
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions
```

```
In [146]: print("=*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, train_fpr, train_tpr)))
```

```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.49531237101902936 for threshold 0.52
[[ 5239  2187]
 [12398 29217]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.49531237101902936 for threshold 0.52
[[ 3151  2308]
 [ 9779 20814]]
```

```
In [147]: #How to plot confusion matrix using heat map - https://stackoverflow.com/questions/35572000/how-can-i-plot
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
array = [[5239,2187],
         [12398,29217]]
df_cm = pd.DataFrame(array, index = [i for i in "01"],
                     columns = [i for i in "01"])
plt.figure
plt.title('Confusion matrix for set 2 Train')
sn.heatmap(df_cm, annot=True, annot_kws={"size":16}, fmt='g')
```

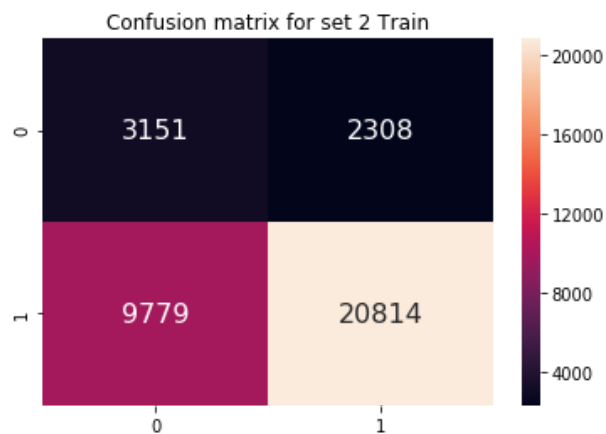
Out[147]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1e49227a940>



```
In [148]: #How to plot confusion matrix using heat map - https://stackoverflow.com/questions/35572000/how-can-i-plot
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
array = [[3151,2308],
         [9779,20814]]
df_cm = pd.DataFrame(array, index = [i for i in "01"],
                     columns = [i for i in "01"])

plt.figure
plt.title('Confusion matrix for set 2 Train')
sn.heatmap(df_cm, annot=True, annot_kws={"size":16}, fmt='g')
```

```
Out[148]: <matplotlib.axes. subplots.AxesSubplot at 0x1e402087390>
```



#### 2.4.2.1 Top 10 important features of positive class from SET 2

```
In [149]: #how to calculate feature probability score in naive bayes python - https://www.datacamp.com/community/tutorials/naive-bayes-classification-python#comment=67807
          #how to calculate feature probability score in naive bayes python- https://towardsdatascience.com/naive-bayes-classification-in-python-part-2-features-probability-score-4e0d3c0f0000
          nb_tfidf = MultinomialNB(alpha = 0.5)
          nb_tfidf.fit(X_train, y_train)
```

```
Out[149]: MultinomialNB(alpha=0.5, class_prior=None, fit_prior=True)
```

```
In [150]: tfidf_features_probs = []
```

```
In [151]: for c in range(14261) :
           d = naive.feature_log_prob_[1,c]
           tfidf.features_probs.append(d)
```

```
In [152]: len(tfidf_features_probs)
```

Out[152]: 14261

```
In [153]: tfidf_features_names = []
```

```
In [154]: for _ in vectorizer_state.get_feature_names():
          tfidf_features_names.append(_)
```

```
In [155]: for _ in vectorizer_tp.get_feature_names():
          tfidf_features.names.append( )
```

```
In [156]: for _ in vectorizer_grade.get_feature_names():
          tfidf_features.names.append(_)
```

```
In [157]: for _ in vectorizer_category.get_feature_names():
          tfidf.features.names.append( )
```



```
In [158]: for _ in vectorizer_subcat.get_feature_names() :  
         tfidf_features_names.append(_)
```

```
In [159]: for _ in vectorizer_essay_tfidf.get_feature_names() :  
         tfidf_features_names.append(_)
```

```
In [160]: for _ in vectorizer_tfidf_title.get_feature_names() :  
         tfidf_features_names.append(_)
```

```
In [161]: tfidf_features_names.append('price')  
         tfidf_features_names.append('quantity')  
         tfidf_features_names.append('teacher_number_of_previously_posted')
```

```
In [162]: len(tfidf_features_names)
```

Out[162]: 14261

```
In [163]: final_tfidf_features_set2_pos = pd.DataFrame({'feature_prob_estimates' : tfidf_features_probs, 'feature_names' : tfidf_features_names})
```

```
In [164]: final_tfidf_features_set2_pos.sort_values(by = ['feature_prob_estimates'], ascending = True, inplace=True)
```

```
In [165]: final_tfidf_features_set2_pos.head(10)
```

Out[165]:

|       | feature_prob_estimates | feature_names |
|-------|------------------------|---------------|
| 14245 | -13.94023              | ye            |
| 14233 | -13.94023              | workshop      |
| 14244 | -13.94023              | xylophone     |
| 14243 | -13.94023              | writing       |
| 14242 | -13.94023              | writers       |
| 14241 | -13.94023              | writer        |
| 14220 | -13.94023              | wobbles       |
| 14240 | -13.94023              | write         |
| 14239 | -13.94023              | wow           |
| 14238 | -13.94023              | would         |

#### 2.4.2.2 Top 10 important features of negative class from SET 2

```
In [168]: tfidf_features_neg_probs = []  
         for a in range(14261) :  
             bn = naive.feature_log_prob_[0,a]  
             tfidf_features_neg_probs.append(bn)
```

```
In [169]: final_features_tfidf_neg = pd.DataFrame({'feature_prob_estimates' : tfidf_features_neg_probs, 'feature_names' : tfidf_features_names})
```

```
In [170]: final_features_tfidf_neg.sort_values(by = ['feature_prob_estimates'], ascending = True, inplace=True)
```

```
In [171]: final_features_tfidf_neg.head(10)
```

```
Out[171]:
```

|       | feature_prob_estimates | feature_names |
|-------|------------------------|---------------|
| 5067  | -12.261929             | gun           |
| 6838  | -12.261929             | membership    |
| 12728 | -12.261929             | easier        |
| 12725 | -12.261929             | ease          |
| 6771  | -12.261929             | may           |
| 3329  | -12.261929             | dismissal     |
| 6767  | -12.261929             | maximized     |
| 11264 | -12.261929             | tricycles     |
| 3331  | -12.261929             | disorder      |
| 708   | -12.261929             | amplified     |

### 3. Conclusions

```
In [172]: # Please compare all your models using Prettytable Library
from prettytable import PrettyTable
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter-Alpha value", "AUC"]
x.add_row(["BOW", "Naive Bayes", 0.5, 0.70])
x.add_row(["TFIDF", "Naive Bayes", 0.5, 0.67])
print(x)
```

| Vectorizer | Model       | Hyper Parameter-Alpha value | AUC  |
|------------|-------------|-----------------------------|------|
| BOW        | Naive Bayes | 0.5                         | 0.7  |
| TFIDF      | Naive Bayes | 0.5                         | 0.67 |