

## 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>	<ul style="list-style-type: none"><li>•</li><li>•</li></ul>	Title of the project. <b>Examples:</b> Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li></ul>	Grade level of students for which the project is targeted. One of the following enumerated values:  Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	<ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li></ul>	One or more (comma-separated) subject categories for the project from the following enumerated list of values:  Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth
	<ul style="list-style-type: none"><li>•</li><li>•</li></ul>	<b>Examples:</b>  Music & The Arts 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 (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)</a> ). <b>Example:</b> WY	
<code>project_subject_subcategories</code>	<ul style="list-style-type: none"><li>•</li><li>•</li></ul>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b>  Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	<ul style="list-style-type: none"><li>•</li></ul>	An explanation of the resources needed for the project. <b>Example:</b>  My students need hands on literacy materials to manage sensory needs!</code

Feature	Description
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <div><div><div></div><div></div><div></div><div></div><div></div><div></div></div><div>nan</div><div>Dr.</div><div>Mr.</div><div>Mrs.</div><div>Ms.</div><div>Teacher.</div></div>
teacher_number_of_previously_posted_projects	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
id	A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502
description	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. <b>Example:</b> 3
price	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
project_is_approved	A binary flag indicating whether DonorsChoose approved the project. A value of <code>0</code> indicates the project was not approved, and a value of <code>1</code> indicates the project was approved.

Notes on the Essay Data

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

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

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

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

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

```
In [1]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 import pandas as pd
6 import numpy as np
7 import nltk
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 from sklearn.feature_extraction.text import CountVectorizer
12 from sklearn.metrics import confusion_matrix
13 from sklearn import metrics
14 from sklearn.metrics import roc_curve, auc
15
16 from sklearn.preprocessing import Normalizer
17 import re
18 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
19
20 import pickle
21 from tqdm import tqdm
22 import os
23
24 # from plotly import plotly
25 # import plotly.offline as offline
26 # import plotly.graph_objs as go
27 # offline.init_notebook_mode()
28 from collections import Counter
```

## 1. READING FILES

```
In [2]: 1 project_data = pd.read_csv('train_data.csv')
2 resource_data = pd.read_csv('resources.csv')
```

```
In [3]: 1 print("Number of data points in train data (project data) ", project_data.shape)
2 print('-'*50)
3 print("The attributes of data :", project_data.columns.values)
```

Number of data points in train data (project data) (109248, 17)

-----

The attributes of data : ['Unnamed: 0' 'id' 'teacher\_id' 'teacher\_prefix' 'school\_state'  
 'project\_submitted\_datetime' 'project\_grade\_category'  
 'project\_subject\_categories' 'project\_subject\_subcategories'  
 'project\_title' 'project\_essay\_1' 'project\_essay\_2' 'project\_essay\_3'  
 'project\_essay\_4' 'project\_resource\_summary'  
 'teacher\_number\_of\_previously\_posted\_projects' 'project\_is\_approved']

```
In [4]: 1 print("Number of data points in train data (resource data) ", resource_data.shape)
2 print('-'*50)
3 print("The attributes of data :", resource_data.columns.values)
4 resource_data.head(2)
```

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

-----  
The attributes of data : ['id' 'description' 'quantity' 'price']

Out[4]:

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

## 1.1 Preprocessing Categorical Features: project\_grade\_category

```
In [5]: 1 print("Project grade", project_data['project_grade_category'].value_counts())
2 ## visulaize how project grade looks like
3 print('-'*50)
4 print(project_data['project_grade_category'].values[1000])
5 print(project_data['project_grade_category'].values[1500])
```

Project grade Grades PreK-2 44225

Grades 3-5 37137

Grades 6-8 16923

Grades 9-12 10963

Name: project\_grade\_category, dtype: int64

-----  
Grades 3-5  
Grades PreK-2

Remove spaces , convert "-" to "\_" and convert letters to lower case and remove grades as it is repeated

```
In [6]: 1 # https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-other-column-value
2 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('Grades ', '')
3 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ', '_')
4 project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-', '_')
5 project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
6 project_data['project_grade_category'].value_counts()
```

Out[6]: prek\_2 44225  
3\_5 37137  
6\_8 16923  
9\_12 10963  
Name: project\_grade\_category, dtype: int64

```

In [7]: 1 categories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3 # reference from course material : reference_EDA.ipynb
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 cat_list = []
8 for i in categories:
9     temp = ""
10    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
11    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
12        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
13            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
14            j = j.replace(' ','') # we are replacing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
15            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
16            temp = temp.replace('&','_') # we are replacing the & value into
17    cat_list.append(temp.strip())
18
19 project_data['clean_categories'] = cat_list
20 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
21 project_data.head(2)
22

```

Out[7]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_subject_subcategories	project_title	project_essay_1	project_essay_2
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	prek_2	ESL, Literacy	Educational Support for English Learners at Home	My students are English learners that are work...	\\"The limits of your language are the limits o...
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	6_8	Civics & Government, Team Sports	Wanted: Projector for Hungry Learners	Our students arrive to our school eager to lea...	The projector we need for our school is very c...

```

In [8]: 1 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
2 from collections import Counter
3 my_counter = Counter()
4 for word in project_data['clean_categories'].values:
5     my_counter.update(word.split())
6 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
7 cat_dict = dict(my_counter)
8 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))

```

```
In [9]: 1 #examples after preprocessing subject
2
3 print(project_data['clean_categories'].values[100])
4 print('-'*50)
5 print(project_data['clean_categories'].values[200])
6 print('-'*50)
```

Literacy\_Language Math\_Science

-----  
AppliedLearning  
-----

### 1.3 Preprocessing Categorical Features: project\_subject\_subcategories

```
In [10]: 1 sub_catogories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3 # reference from course material : reference_EDA.ipynb
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7
8 sub_cat_list = []
9 for i in sub_catogories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the catogory based on space "Math & Science"=> "Math", "&", "Science"
14             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15             j = j.replace(' ', '') # we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&', '_')
18         sub_cat_list.append(temp.strip())
19
20 project_data['clean_subcategories'] = sub_cat_list
21 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
22 project_data.head(2)
23
24 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
25 from collections import Counter
26 my_counter = Counter()
27 for word in project_data['clean_subcategories'].values:
28     my_counter.update(word.split())
29
30 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
31 sub_cat_dict = dict(my_counter)
32 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
```

### 1.3 Preprocessing Categorical Features: school\_state

```
In [11]: 1 project_data['school_state'].value_counts()
2 ## Convert it to Lower
3 project_data['school_state'] = project_data['school_state'].str.lower()
4 project_data['school_state'].value_counts()
```

```
Out[11]: 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
```

## 1.4 Preprocessing Categorical Features: Teacher\_prefix

```
In [12]: 1 print(project_data['teacher_prefix'].value_counts())
2 # try to remove the dots from the teacher prefix
3 project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
4 project_data['teacher_prefix']=project_data['teacher_prefix'].str.replace('.', '')
5 project_data['teacher_prefix']=project_data['teacher_prefix'].str.lower()
6 project_data['teacher_prefix']=project_data['teacher_prefix'].str.strip()
```

```
Mrs.      57269
Ms.       38955
Mr.       10648
Teacher   2360
Dr.        13
Name: teacher_prefix, dtype: int64
```

## 1.5 Combining all the essays

```
In [13]: 1 # merge two column text dataframe:
2 project_data["essay"] = project_data["project_essay_1"].map(str) + \
3     project_data["project_essay_2"].map(str) + \
4     project_data["project_essay_3"].map(str) + \
5     project_data["project_essay_4"].map(str)
```

## 1.6 Number of Words in the Essay and Title

```
In [14]: 1 source:'''https://www.geeksforgeeks.org/python-program-to-count-words-in-a-sentence/'''
2 words_counter=[]
3 for string in project_data['essay']:
4     res = len(re.findall(r'\w+', string))
5     words_counter.append(res)
6
7 project_data["words_in_essay"] = words_counter
8
9 words_counter=[]
10
11 for string in project_data['project_title']:
12     res = len(re.findall(r'\w+', string))
13     words_counter.append(res)
14 project_data["words_in_title"] = words_counter
```

## 1.7. Preprocessing Numerical Values: price

```
In [15]: 1 ## calculate the overall count of resources and the total price for each project id
2 price_data=resource_data.groupby('id').agg({'price':'sum','quantity':'sum' })
```

```
In [16]: 1 project_data=pd.merge(project_data, price_data, on='id', how='left')
```

## 1.8 Preprocessing Text Features: project\_title



```
In [17]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [18]: 1 # https://gist.github.com/sebleier/554280
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'
3 stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
11            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
12            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
13            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
14            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
15            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
16            "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
17            'won', "won't", 'wouldn', "wouldn't"]
```

```
In [19]: 1 project_data['project_title'].head(5)
```

```
Out[19]: 0    Educational Support for English Learners at Home
1          Wanted: Projector for Hungry Learners
2    Soccer Equipment for AWESOME Middle School Stu...
3          Techie Kindergarteners
4          Interactive Math Tools
Name: project_title, dtype: object
```

```
In [20]: 1 print("printing some random reviews")
2 print(9, project_data['project_title'].values[9])
3 print(34, project_data['project_title'].values[34])
4 print(147, project_data['project_title'].values[147])
```

```
printing some random reviews
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

```
In [21]: 1 # Combining all the above stundents
2 from tqdm import tqdm
3 def preprocess_text(text_data):
4     preprocessed_text = []
5     # tqdm is for printing the status bar
6     for sentence in tqdm(text_data):
7         sent = decontracted(sentence)
8         sent = sent.replace('\\r', ' ')
9         sent = sent.replace('\\n', ' ')
10        sent = sent.replace('\\\"', ' ')
11        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12        # https://gist.github.com/sebleier/554280
13        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14        preprocessed_text.append(sent.lower().strip())
15    return preprocessed_text
```

```
In [22]: 1 preprocessed_titles = preprocess_text(project_data['project_title'].values)
```

100%|██| 109248/109248 [00:03<00:00, 36099.69it/s]

```
In [23]: 1 print("printing some random reviews")
2 print(9, preprocessed_titles[9])
3 print(34, preprocessed_titles[34])
4 print(147, preprocessed_titles[147])
5 #merge the column in the project_data
6 project_data['processed_title']=preprocessed_titles
```

printing some random reviews  
9 love reading pure pleasure  
34 ball  
147 needs chromebook

## 1.9 Preprocessing Categorical Features: essay

In [24]:

```

1 print("printing some random essay")
2 print(9, project_data['essay'].values[9])
3 print('-'*50)
4 print(34, project_data['essay'].values[34])
5 print('-'*50)
6 print(147, project_data['essay'].values[147])

```

printing some random essay

9 Over 95% of my students are on free or reduced lunch. I have a few who are homeless, but despite that, they come to school with an eagerness to learn. My students are inquisitive eager learners who embrace the challenge of not having great books and other resources every day. Many of them are not afforded the opportunity to engage with these big colorful pages of a book on a regular basis at home and they don't travel to the public library. \r\nIt is my duty as a teacher to do all I can to provide each student an opportunity to succeed in every aspect of life. \r\nReading is Fundamental! My students will read these books over and over again while boosting their comprehension skills. These books will be used for read alouds, partner reading and for Independent reading. \r\nThey will engage in reading to build their "Love for Reading" by reading for pure enjoyment. They will be introduced to some new authors as well as some old favorites. I want my students to be ready for the 21st Century and know the pleasure of holding a good hard back book in hand. There's nothing like a good book to read! \r\nMy students will soar in Reading, and more because of your consideration and generous funding contribution. This will help build stamina and prepare for 3rd grade. Thank you so much for reading our proposal!nannan

34 My students mainly come from extremely low-income families, and the majority of them come from homes where both parents work full time. Most of my students are at school from 7:30 am to 6:00 pm (2:30 to 6:00 pm in the after-school program), and they all receive free and reduced meals for breakfast and lunch. \r\n\r\n\r\nI want my students to feel as comfortable in my classroom as they do at home. Many of my students take on multiple roles both at home as well as in school. They are sometimes the caretakers of younger siblings, cooks, babysitters, academics, friends, and most of all, they are developing who they are going to become as adults. I consider it an essential part of my job to model helping others gain knowledge in a positive manner. As a result, I have a community of students who love helping each other in and outside of the classroom. They consistently look for opportunities to support each other's learning in a kind and helpful way. I am excited to be experimenting with alternative seating in my classroom this school year. Studies have shown that giving students the option of where they sit in a classroom increases focus as well as motivation. \r\n\r\nBy allowing students choice in the classroom, they are able to explore and create in a welcoming environment. Alternative classroom seating has been experimented with more frequently in recent years. I believe (along with many others), that every child learns differently. This does not only apply to how multiplication is memorized, or a paper is written, but applies to the space in which they are asked to work. I have had students in the past ask "Can I work in the library? Can I work on the carpet?" My answer was always, "As long as you're learning, you can work wherever you want!" \r\n\r\nWith the yoga balls and the lap-desks, I will be able to increase the options for seating in my classroom and expand its imaginable space.nannan

147 My students are eager to learn and make their mark on the world.\r\n\r\nThey come from a Title 1 school and need extra love.\r\n\r\nMy fourth grade students are in a high poverty area and still come to school every day to get their education. I am trying to make it fun and educational for them so they can get the most out of their schooling. I created a caring environment for the students to bloom! They deserve the best.\r\nThank you!\r\nI am requesting 1 Chromebook to access online interventions, differentiate instruction, and get extra practice. The Chromebook will be used to supplement ELA and math instruction. Students will play ELA and math games that are engaging and fun, as well as participate in assignments online. This in turn will help my students improve their skills. Having a Chromebook in the classroom would not only allow students to use the programs at their own pace, but would ensure more students are getting adequate time to use the programs. The online programs have been especially beneficial to my students with special needs. They are able to work at their level as well as be challenged with some different materials. This is making these students more confident in their abilities.\r\n\r\nThe Chromebook would allow my students to have daily access to computers and increase their computing skills.\r\nThis will change their lives for the better as they become more successful in school. Having access to technology in the classroom would help bridge the achievement gap.nannan

In [25]:

```

1 preprocessed_essays = preprocess_text(project_data['essay'].values)

```

100%|██| 109248/109248 [01:13<00:00, 1488.95it/s]

```
In [26]: 1 print("printing some random essay")
2 print(9, preprocessed_essays[9])
3 print('-'*50)
4 print(34, preprocessed_essays[34])
5 print('-'*50)
6 print(147, preprocessed_essays[147])
7
8 #merge the column in the project_data
9 project_data['processed_essay']=preprocessed_essays
```

printing some random essay

9 95 students free reduced lunch homeless despite come school eagerness learn students inquisitive eager learners embrace challenge not great books resources every day many not afforded opportunity engage big colorful pages book regular basis home not travel public library duty teacher provide student opportunity succeed every aspect life reading fund amental students read books boosting comprehension skills books used read alouds partner reading independent reading engage reading build love reading reading pure enjoyment in troduced new authors well old favorites want students ready 21st century know pleasure holding good hard back book hand nothing like good book read students soar reading consid eration generous funding contribution help build stamina prepare 3rd grade thank much reading proposal nannan

34 students mainly come extremely low income families majority come homes parents work full time students school 7 30 6 00 pm 2 30 6 00 pm school program receive free reduced m eals breakfast lunch want students feel comfortable classroom home many students take multiple roles home well school sometimes caretakers younger siblings cooks babysitters ac ademics friends developing going become adults consider essential part job model helping others gain knowledge positive manner result community students love helping outside cl assroom consistently look opportunities support learning kind helpful way excited experimenting alternative seating classroom school year studies shown giving students option s it classroom increases focus well motivation allowing students choice classroom able explore create welcoming environment alternative classroom seating experimented frequently recent years believe along many others every child learns differently not apply multiplication memorized paper written applies space asked work students past ask work library w ork carpet answer always long learning work wherever want yoga balls lap desks able increase options seating classroom expand imaginable space nannan

147 students eager learn make mark world come title 1 school need extra love fourth grade students high poverty area still come school every day get education trying make fun e ducational get schooling created caring environment students bloom deserve best thank requesting 1 chromebook access online interventions differentiate instruction get extra pr actice chromebook used supplement ela math instruction students play ela math games engaging fun well participate assignments online turn help students improve skills chromeboo k classroom would not allow students use programs pace would ensure students getting adequate time use programs online programs especially beneficial students special needs abl e work level well challenged different materials making students confident abilities chromebook would allow students daily access computers increase computing skills change liv es better become successful school access technology classroom would help bridge achievement gap nannan

## Train,Test,CV Split

```
In [27]: 1 # train test split using sklearn.model selection
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'], test_size=0.33, stratify = project_data['project_is_approved'],random_
4 X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train,random_state=0)
```

```
In [28]: 1 ## drop the y labels from splits
2 X_train.drop(['project_is_approved'], axis=1, inplace=True)
3 X_test.drop(['project_is_approved'], axis=1, inplace=True)
4 X_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

In [29]:

▶

1 X\_train.head(2)

Out[29]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_title	project_essay_1	project_essay_2	...	teacher_number_of_previ
75742	118221	p186156	f50f55a2b44b65b54f38f03c5df21922	mrs	tx	2017-03-01 16:21:46	9_12	I-Waste: a Multi-Media Art Installation on El...	My students are creative human beings. They a...	It's no secret that the arts are underfunded i...	...	
61001	57644	p180433	9e0fb5827f551d7e6966f8b3985e387b	ms	ny	2017-03-09 10:19:06	6_8	Authentic Listening and Speaking Activities fo...	The Hyde Park Central School District is locat...	One area my students struggle the most with is...	...	

2 rows × 23 columns

VECTORIZING DATA

One hot encoding on Categorical

In [30]:

▶

```
1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean categories
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 cat_vectorize = CountVectorizer(lowercase=False, binary=True)
6 cat_vectorize.fit(X_train['clean_categories'].values)
7
8 train_categories = cat_vectorize.transform(X_train['clean_categories'].values)
9 test_categories = cat_vectorize.transform(X_test['clean_categories'].values)
10 cv_categories = cat_vectorize.transform(X_cv['clean_categories'].values)
11
12 print(cat_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",train_categories.shape)
14 print("Shape of matrix of Test data after one hot encoding ",test_categories.shape)
15 print("Shape of matrix of CV data after one hot encoding ",cv_categories.shape)
```

```
['AppliedLearning', 'Care_Hunger', 'Health_Sports', 'History_Civics', 'Literacy_Language', 'Math_Science', 'Music_Arts', 'SpecialNeeds', 'Warmth']
Shape of matrix of Train data after one hot encoding (49041, 9)
Shape of matrix of Test data after one hot encoding (36052, 9)
Shape of matrix of CV data after one hot encoding (24155, 9)
```

```
In [31]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean subcategories
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 subcat_vectorize = CountVectorizer(lowercase=False, binary=True)
6 subcat_vectorize.fit(X_train['clean_subcategories'].values)
7
8 train_subcategories = subcat_vectorize.transform(X_train['clean_subcategories'].values)
9 test_subcategories = subcat_vectorize.transform(X_test['clean_subcategories'].values)
10 cv_subcategories = subcat_vectorize.transform(X_cv['clean_subcategories'].values)
11
12 print(subcat_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",train_subcategories.shape)
14 print("Shape of matrix of Test data after one hot encoding ",test_subcategories.shape)
15 print("Shape of matrix of CV data after one hot encoding ",cv_subcategories.shape)
16
```

```
['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government', 'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics', 'EnvironmentalSc
ience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness', 'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writin
g', 'Mathematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'SocialSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix of Train data after one hot encoding (49041, 30)
Shape of matrix of Test data after one hot encoding (36052, 30)
Shape of matrix of CV data after one hot encoding (24155, 30)
```

```
In [32]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## school state
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 sklstate_vectorize = CountVectorizer(lowercase=False, binary=True)
6 sklstate_vectorize.fit(X_train['school_state'].values)
7
8 sklstate_train = sklstate_vectorize.transform(X_train['school_state'].values)
9 sklstate_test = sklstate_vectorize.transform(X_test['school_state'].values)
10 sklstate_cv = sklstate_vectorize.transform(X_cv['school_state'].values)
11
12 print(sklstate_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",sklstate_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",sklstate_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",sklstate_cv.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 of Train data after one hot encoding (49041, 51)
Shape of matrix of Test data after one hot encoding (36052, 51)
Shape of matrix of CV data after one hot encoding (24155, 51)
```

```
In [33]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## teacher_prefix
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 teacher_prefix_vectorize = CountVectorizer(lowercase=False, binary=True)
6 teacher_prefix_vectorize.fit(X_train['teacher_prefix'].values)
7
8 teacher_prefix_train = teacher_prefix_vectorize.transform(X_train['teacher_prefix'].values)
9 teacher_prefix_test = teacher_prefix_vectorize.transform(X_test['teacher_prefix'].values)
10 teacher_prefix_cv = teacher_prefix_vectorize.transform(X_cv['teacher_prefix'].values)
11
12 print(teacher_prefix_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",teacher_prefix_cv.shape)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
Shape of matrix of Train data after one hot encoding (49041, 5)
Shape of matrix of Test data after one hot encoding (36052, 5)
Shape of matrix of CV data after one hot encoding (24155, 5)
```

```
In [34]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## project_grade
3 from sklearn.feature_extraction.text import CountVectorizer
4
5 proj_grade_vectorize = CountVectorizer(lowercase=False, binary=True)
6 proj_grade_vectorize.fit(X_train['teacher_prefix'].values)
7
8 proj_grade_train = proj_grade_vectorize.transform(X_train['teacher_prefix'].values)
9 proj_grade_test = proj_grade_vectorize.transform(X_test['teacher_prefix'].values)
10 proj_grade_cv = proj_grade_vectorize.transform(X_cv['teacher_prefix'].values)
11
12 print(proj_grade_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",proj_grade_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",proj_grade_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",proj_grade_cv.shape)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
Shape of matrix of Train data after one hot encoding (49041, 5)
Shape of matrix of Test data after one hot encoding (36052, 5)
Shape of matrix of CV data after one hot encoding (24155, 5)
```

## Vectorizing Text data

### A. BOW on Essay

#### *Train data*



```
In [35]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  bow_essay = CountVectorizer(min_df=10,max_features=5000)
4  bow_essay.fit(X_train['processed_essay'])
5
6  bow_essay_train = bow_essay.transform(X_train['processed_essay'])
7
8  print("Shape of matrix after one hot encoding ",bow_essay_train.shape)
```

Shape of matrix after one hot encoding (49041, 5000)

### Test data

```
In [36]: 1
2  bow_essay_test = bow_essay.transform(X_test['processed_essay'])
3
4  print("Shape of matrix after one hot encoding ",bow_essay_test.shape)
```

Shape of matrix after one hot encoding (36052, 5000)

### CV data

```
In [37]: 1  bow_essay_cv = bow_essay.transform(X_cv['processed_essay'])
2  print("Shape of matrix after one hot encoding ",bow_essay_cv.shape)
```

Shape of matrix after one hot encoding (24155, 5000)

## A. BOW on Titles

### Train data

```
In [38]: 1  ##Considering the words that appeared in atleast 10 documents
2
3  bowtitle = CountVectorizer(min_df=10)
4  bowtitle.fit(X_train['processed_title'])
5
6  bow_title_train = bowtitle.transform(X_train['processed_title'])
7
8  print("Shape of matrix after one hot encoding ",bow_title_train.shape)
```

Shape of matrix after one hot encoding (49041, 2008)

### Test data



```
In [39]: 1  ##Considering the words that appeared in atleast 10 documents
        2
        3
        4 bow_title_test = bowtitle.transform(X_test['processed_title'])
        5
        6 print("Shape of matrix after one hot encoding ",bow_title_test.shape)
```

Shape of matrix after one hot encoding (36052, 2008)

### CV data

```
In [40]: 1  ##Considering the words that appeared in atleast 10 documents
        2
        3 bow_title_cv = bowtitle.transform(X_cv['processed_title'])
        4
        5 print("Shape of matrix after one hot encoding ",bow_title_cv.shape)
```

Shape of matrix after one hot encoding (24155, 2008)

### A. TFIDF on Essay

#### Train data

```
In [41]: 1  ##Considering the words that appeared in atleast 10 documents
        2
        3 tfidf_essay = TfidfVectorizer(min_df=10,max_features=5000) #selecting top 5000 features
        4 tfidf_essay.fit(X_train['processed_essay'])
        5
        6 tfidf_essay_train = tfidf_essay.transform(X_train['processed_essay'])
        7
        8 print("Shape of matrix after one hot encoding ",tfidf_essay_train.shape)
```

Shape of matrix after one hot encoding (49041, 5000)

#### Test data

```
In [42]: 1
        2 tfidf_essay_test = tfidf_essay.transform(X_test['processed_essay'])
        3
        4 print("Shape of matrix after one hot encoding ",tfidf_essay_test.shape)
```

Shape of matrix after one hot encoding (36052, 5000)

### CV data

```
In [43]: 1
          2
          3 tfidf_essay_cv = tfidf_essay.transform(X_cv['processed_essay'])
          4
          5 print("Shape of matrix after one hot encoding ",tfidf_essay_cv.shape)
```

Shape of matrix after one hot encoding (24155, 5000)

## A. TFIDF on Titles

### Train data

```
In [44]: 1 ##Considering the words that appeared in atleast 10 documents
          2
          3 tfidftitle = TfidfVectorizer(min_df=10)
          4 tfidftitle.fit(X_train['processed_title'])
          5
          6 tfidf_title_train = tfidftitle.transform(X_train['processed_title'])
          7
          8 print("Shape of matrix after one hot encoding ",tfidf_title_train.shape)
```

Shape of matrix after one hot encoding (49041, 2008)

### Test data

```
In [45]: 1 ##Considering the words that appeared in atleast 10 documents
          2
          3 tfidf_title_test = tfidftitle.transform(X_test['processed_title'])
          4
          5 print("Shape of matrix after one hot encoding ",tfidf_title_test.shape)
```

Shape of matrix after one hot encoding (36052, 2008)

### CV data

```
In [46]: 1 ##Considering the words that appeared in atleast 10 documents
          2
          3 tfidf_title_cv = tfidftitle.transform(X_cv['processed_title'])
          4
          5 print("Shape of matrix after one hot encoding ",tfidf_title_cv.shape)
```

Shape of matrix after one hot encoding (24155, 2008)

Type *Markdown* and LaTeX:  $\alpha^2$

## Vectorizing Numerical Features

### Price

```

In [47]: ▶ 1 normalizer = Normalizer()
2 # normalizer.fit(X_train['price'].values)
3 # this will rise an error Expected 2D array, got 1D array instead:
4 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
5 # Reshape your data either using
6 # array.reshape(-1, 1) if your data has a single feature
7 # array.reshape(1, -1) if it contains a single sample.
8 normalizer.fit(X_train['price'].values.reshape(1,-1))
9
10 X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1))
11 X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
12 X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))
13
14 print("After vectorizations")
15 print(X_train_price_norm.shape, y_train.shape)
16 print(X_cv_price_norm.shape, y_cv.shape)
17 print(X_test_price_norm.shape, y_test.shape)
18 print("="*100)
19
20 ## reshaping
21 X_train_price_norm=X_train_price_norm.reshape(-1,1)
22 X_cv_price_norm=X_cv_price_norm.reshape(-1,1)
23 X_test_price_norm=X_test_price_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

**Quantity**

In [48]:

```

1
2 normalizer = Normalizer()
3
4 # normalizer.fit(X_train['price'].values)
5 # this will rise an error Expected 2D array, got 1D array instead:
6 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
7 # Reshape your data either using
8 # array.reshape(-1, 1) if your data has a single feature
9 # array.reshape(1, -1) if it contains a single sample.
10
11 normalizer.fit(X_train['quantity'].values.reshape(1,-1))
12
13 quantity_train_norm = normalizer.transform(X_train['quantity'].values.reshape(1,-1))
14 quantity_cv_norm = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
15 quantity_test_norm = normalizer.transform(X_test['quantity'].values.reshape(1,-1))
16
17 print("After vectorizations")
18 print(quantity_train_norm.shape, y_train.shape)
19 print(quantity_cv_norm.shape, y_cv.shape)
20 print(quantity_test_norm.shape, y_test.shape)
21 print("=*100)
22
23 ## reshaping
24 quantity_train_norm=quantity_train_norm.reshape(-1,1)
25 quantity_cv_norm=quantity_cv_norm.reshape(-1,1)
26 quantity_test_norm=quantity_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

## Number of Previously posted projects

```

In [49]: ▶ 1 normalizer = Normalizer()
2
3 # normalizer.fit(X_train['price'].values)
4 # this will rise an error Expected 2D array, got 1D array instead:
5 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
6 # Reshape your data either using
7 # array.reshape(-1, 1) if your data has a single feature
8 # array.reshape(1, -1) if it contains a single sample.
9
10 normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
11
12 prev_projects_train_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
13 prev_projects_cv_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
14 prev_projects_test_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1, -1))
15
16 print("After vectorizations")
17 print(prev_projects_train_norm.shape, y_train.shape)
18 print(prev_projects_cv_norm.shape, y_cv.shape)
19 print(prev_projects_test_norm.shape, y_test.shape)
20 print("=="*100)
21
22 ## reshaping
23 prev_projects_train_norm=prev_projects_train_norm.reshape(-1,1)
24 prev_projects_cv_norm=prev_projects_cv_norm.reshape(-1,1)
25 prev_projects_test_norm=prev_projects_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

## Title Word counts

```

In [50]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['words_in_title'].values.reshape(1,-1))
4
5 title_word_count_train_norm = normalizer.transform(X_train['words_in_title'].values.reshape(1,-1))
6 title_word_count_cv_norm = normalizer.transform(X_cv['words_in_title'].values.reshape(1,-1))
7 title_word_count_test_norm = normalizer.transform(X_test['words_in_title'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(title_word_count_train_norm.shape, y_train.shape)
11 print(title_word_count_cv_norm.shape, y_cv.shape)
12 print(title_word_count_test_norm.shape, y_test.shape)
13 print("="*100)
14
15 ## reshaping
16 title_word_count_train_norm=title_word_count_train_norm.reshape(-1,1)
17 title_word_count_cv_norm=title_word_count_cv_norm.reshape(-1,1)
18 title_word_count_test_norm=title_word_count_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

=====

### Essay Words Counts

```

In [51]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['words_in_essay'].values.reshape(1,-1))
4
5 essay_word_count_train_norm = normalizer.transform(X_train['words_in_essay'].values.reshape(1,-1))
6 essay_word_count_cv_norm = normalizer.transform(X_cv['words_in_essay'].values.reshape(1,-1))
7 essay_word_count_test_norm = normalizer.transform(X_test['words_in_essay'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(essay_word_count_train_norm.shape, y_train.shape)
11 print(essay_word_count_cv_norm.shape, y_cv.shape)
12 print(essay_word_count_test_norm.shape, y_test.shape)
13
14 ## reshaping
15 essay_word_count_train_norm=essay_word_count_train_norm.reshape(-1,1)
16 essay_word_count_cv_norm=essay_word_count_cv_norm.reshape(-1,1)
17 essay_word_count_test_norm=essay_word_count_test_norm.reshape(-1,1)

```

After vectorizations

(1, 49041) (49041,)

(1, 24155) (24155,)

(1, 36052) (36052,)

## Assignment 6: Apply NB

### 1. Apply Multinomial NB on these feature sets

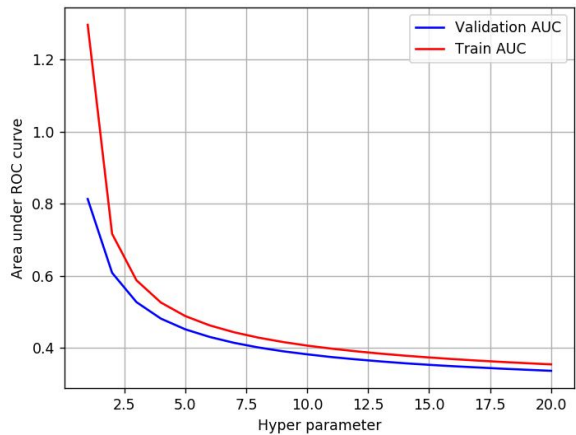
- **Set 1**: categorical, numerical features + preprocessed\_eassay (BOW)
- **Set 2**: categorical, numerical features + preprocessed\_eassay (TFIDF)

2. The hyper paramter tuning(find best alpha:smoothing parameter)

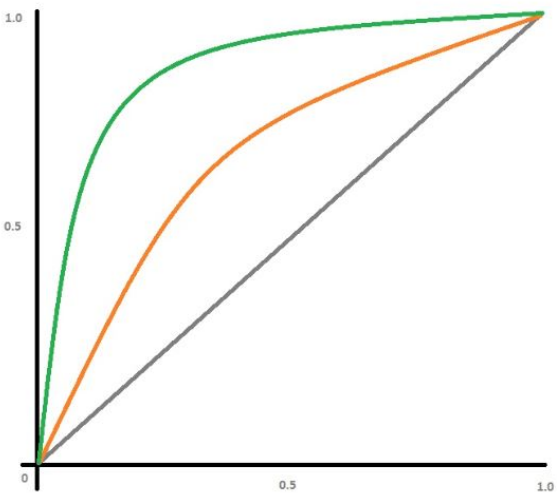
- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- find the best hyper paramter using k-fold cross validation(use GridsearchCV or RandomsearchCV)/simple cross validation data (write for loop to iterate over hyper parameter values)
- 

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure



- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.



- Along with plotting ROC curve, you need to print the [confusion matrix \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/\)](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

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

- 4. fine the top 20 features from either from feature **Set 1** or feature **Set 2** using absolute 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
- 5. You need to summarize the results at the end of the notebook, summarize it in the table format

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	7	0.78
TFIDF	Brute	12	0.79
W2V	Brute	10	0.78
TFIDFW2V	Brute	6	0.78

SET1 :Apply Multinomial NB on these feature sets (categorical, numerical features + preprocessed\_eassay (BOW))

Step 1 : Find best hyperparameter with maximum AUC

```
In [52]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 from scipy.sparse import hstack
3
4
5 X_train_bow = hstack((train_categories, train_subcategories,sklstate_train,teacher_prefix_train,
6                       proj_grade_train,bow_essay_train,bow_title_train,
7                       X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
8                       essay_word_count_train_norm)).tocsr()
9
10 X_test_bow = hstack((test_categories, test_subcategories,sklstate_test,teacher_prefix_test,
11                     proj_grade_test,bow_essay_test,bow_title_test,
12                     X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
13                     essay_word_count_test_norm)).tocsr()
14
15 X_cv_bow = hstack((cv_categories, cv_subcategories,sklstate_cv,teacher_prefix_cv,
16                   proj_grade_cv,bow_essay_cv,bow_title_cv,
17                   X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
18                   essay_word_count_cv_norm)).tocsr()
19
20
21 print(X_train_bow.shape)
22 print(X_test_bow.shape)
23 print(X_cv_bow.shape)
```

(49041, 7113)
(36052, 7113)
(24155, 7113)

```
In [53]: 1 X_train_bow
```

Out[53]: <49041x7113 sparse matrix of type '<class 'numpy.float64''>'
with 5398319 stored elements in Compressed Sparse Row format>



In [54]:

```

1 print("Final Data matrix")
2 print(X_train_bow.shape, y_train.shape)
3 print(X_cv_bow.shape, y_cv.shape)
4 print(X_test_bow.shape, y_test.shape)
5 print("="*100)

```

```

Final Data matrix
(49041, 7113) (49041,)
(24155, 7113) (24155,)
(36052, 7113) (36052,)
=====

```

In [55]:

```

1 def batch_predict(clf, data):
2     # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
3     # not the predicted outputs
4
5     y_data_pred = []
6     tr_loop = data.shape[0] - data.shape[0]%1000
7     # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 49041%1000 = 49000
8     # in this for loop we will iterate until the last 1000 multiplier
9     for i in range(0, tr_loop, 1000):
10         y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
11         # we will be predicting for the last data points
12         y_data_pred.extend(clf.predict_proba(data[tr_loop:][:,1])
13
14     return y_data_pred

```

In [56]:

```

1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.metrics import roc_auc_score
3 import seaborn as sns
4 import math
5
6 ### Consider alphas values
7 train_auc=[]
8 cv_auc=[]
9 alphas=[0.00001, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000]
10
11 for i in tqdm(alphas):
12     MB=MultinomialNB(class_prior=[0.5,0.5],alpha=i)
13     MB.fit(X_train_bow, y_train)
14     y_train_pred = batch_predict(MB, X_train_bow)
15     y_cv_pred = batch_predict(MB, X_cv_bow)
16     train_auc.append(roc_auc_score(y_train,y_train_pred))
17     cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
18
19 ## Scaling the alphas.
20 # min_alpha=min(alphas)
21 # max_alpha=max(alphas)
22 # for i in tqdm(alphas):
23 #     scaled_alpha.append((max_alpha-i)/(max_alpha -min_alpha))
24
25
26

```

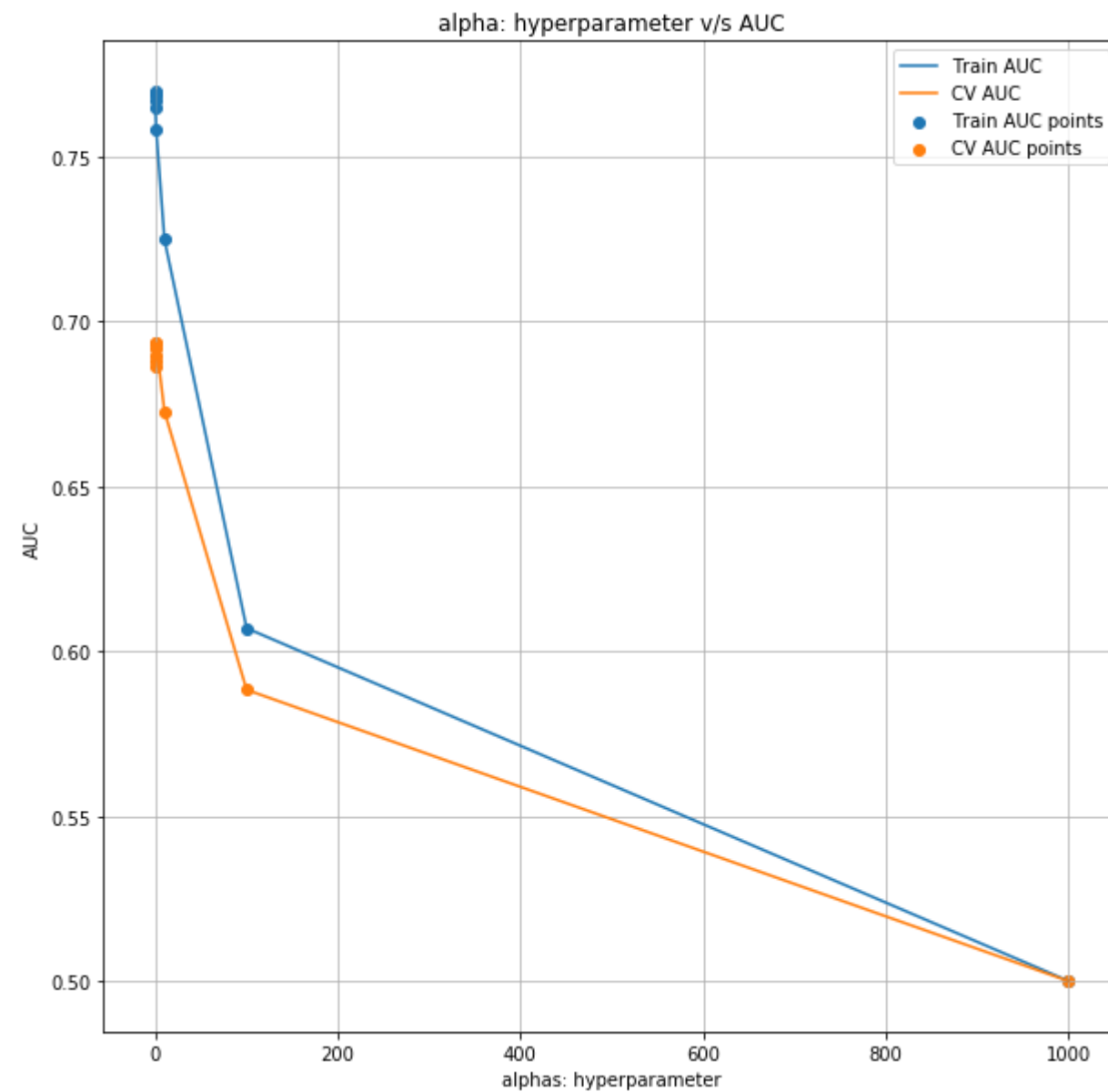
100% | 9/9 [00:01<00:00, 5.75it/s]

In [57]:

```

1 plt.figure(figsize=(10,10))
2 plt.plot(alphas, train_auc, label='Train AUC')
3 plt.plot(alphas, cv_auc, label='CV AUC')
4
5
6 plt.scatter(alphas, train_auc, label='Train AUC points')
7 plt.scatter(alphas, cv_auc, label='CV AUC points')
8
9
10 plt.legend()
11 plt.xlabel("alphas: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("alpha: hyperparameter v/s AUC")
14 plt.grid(which='major', alpha=0.9)
15 plt.show()

```



### Observations :

1. There is a steep fall when alpha is greater than 1 which shows that as alpha increases beyond 1 performance reduces steeply.
2. Best AUC on cross validation is achieved at alpha 0.0001

## RandomSearch CV using K-fold Crossvalidation with k=10

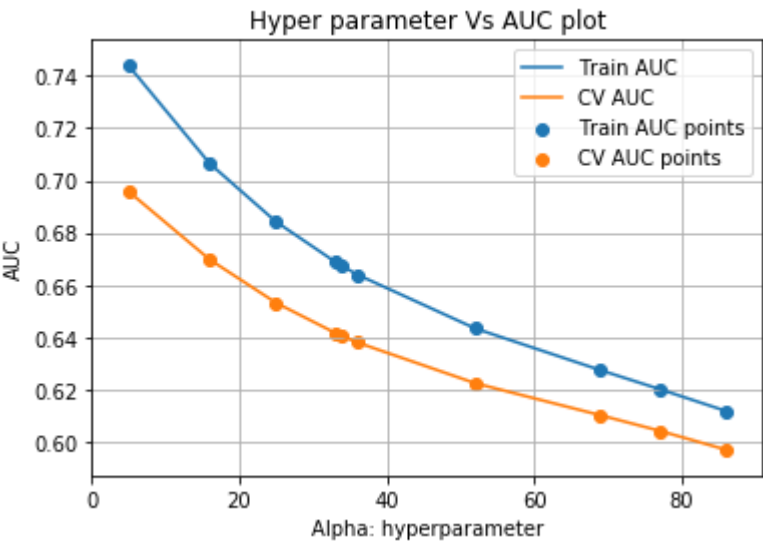
```
In [59]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4
5 parameters={"alpha" : sp_randint(0.0001,100) }
6
7 #RS_Log_alphas =[]
8
9 #RS_alphas=[0.00001, 0.0001, 0.001, 0.1,0.7,0.8, 1,100,1000]
10
11 # for a in tqdm(RS_alphas):
12 #     b = math.log(a)
13 #     RS_Log_alphas.append(b)
14
15 clf = RandomizedSearchCV(MB, parameters,return_train_score=True ,cv=10, scoring='roc_auc',verbose=1,n_jobs=10)
16 clf.fit(X_train_bow, y_train)
17 results = pd.DataFrame.from_dict(clf.cv_results_)
18 results = results.sort_values(['param_alpha'])
19
20 RS_alphas=results['param_alpha']
21 train_auc= results['mean_train_score']
22 train_auc_std= results['std_train_score']
23 cv_auc = results['mean_test_score']
24 cv_auc_std= results['std_test_score']
25
```

Fitting 10 folds for each of 10 candidates, totalling 100 fits

```
[Parallel(n_jobs=10)]: Using backend LokyBackend with 10 concurrent workers.
[Parallel(n_jobs=10)]: Done 30 tasks      | elapsed:    2.2s
[Parallel(n_jobs=10)]: Done 100 out of 100 | elapsed:    6.2s finished
```

In [60]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12
13 plt.legend()
14 plt.xlabel("Alpha: hyperparameter")
15 plt.ylabel("AUC")
16 plt.title("Hyper parameter Vs AUC plot")
17 plt.grid()
18 plt.show()
19
20 results.head()
```



Out[60]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	split2_train_score	split3_train_score	split4_train_score
9	0.380720	0.114470	0.017478	0.008547	5	{'alpha': 5}	0.690013	0.711586	0.710972	0.687821	...	0.743296	0.744372	0.743006
4	0.404020	0.119674	0.020120	0.013498	16	{'alpha': 16}	0.660111	0.686049	0.689977	0.665957	...	0.705205	0.706741	0.705881
6	0.410510	0.109627	0.015958	0.010345	25	{'alpha': 25}	0.643747	0.668834	0.675012	0.652397	...	0.682529	0.684328	0.684084
7	0.351259	0.126504	0.018051	0.010802	33	{'alpha': 33}	0.633130	0.656826	0.664340	0.643136	...	0.666977	0.668928	0.669185
3	0.366248	0.121903	0.024679	0.022986	34	{'alpha': 34}	0.632005	0.655523	0.663208	0.642104	...	0.665275	0.667236	0.667556

5 rows × 31 columns

Observations :

1. We can see a steady fall in the AUC value as alpha increases.

### Train the Model using the best hyper parameter value

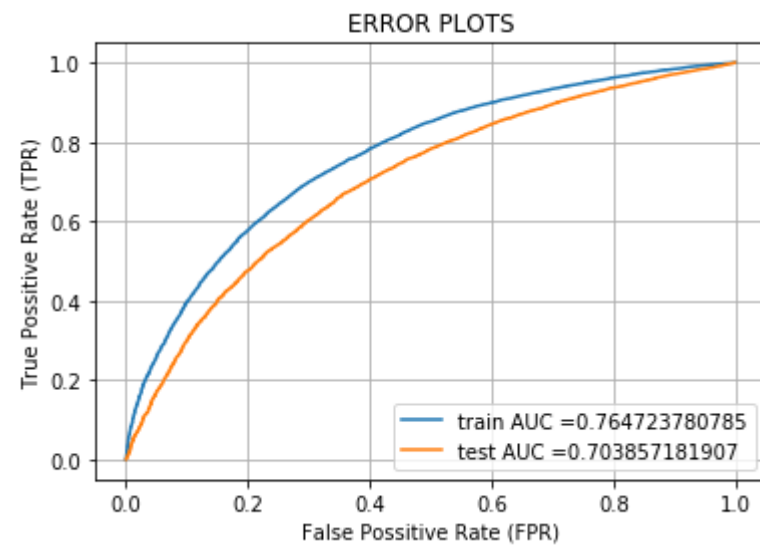
In [61]:

```
1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_alpha = clf.best_params_
3 print('Best Alpha as a result of Random Search',best_alpha)
```

Best Alpha as a result of Random Search {'alpha': 5}

In [62]:

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 MB_bow = MultinomialNB(alpha = 0.1,class_prior=[0.5,0.5])
6 MB_bow.fit(X_train_bow, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(MB_bow, X_train_bow)
11 y_test_pred = batch_predict(MB_bow, X_test_bow)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("False Positive Rate (FPR)")
20 plt.ylabel("True Positive Rate (TPR)")
21 plt.title("ERROR PLOTS")
22 plt.grid()
23 plt.show()
```



### Observations:

1. Train AUC observed to be 0.76 and Test AUC observed to be 0.704 when considering the best alpha.

## CONFUSION MATRIX

```
In [63]: ► 1 # we are writing our own function for predict, with defined threshould
2 # we will pick a threshold that will give the least fpr
3 def find_best_threshold(threshould, fpr, tpr):
4     t = threshould[np.argmax(tpr*(1-fpr))]
5     # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
6     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
7     return t
8
9 def predict_with_best_t(proba, threshould):
10    predictions = []
11    for i in proba:
12        if i>=threshould:
13            predictions.append(1)
14        else:
15            predictions.append(0)
16    return predictions
```

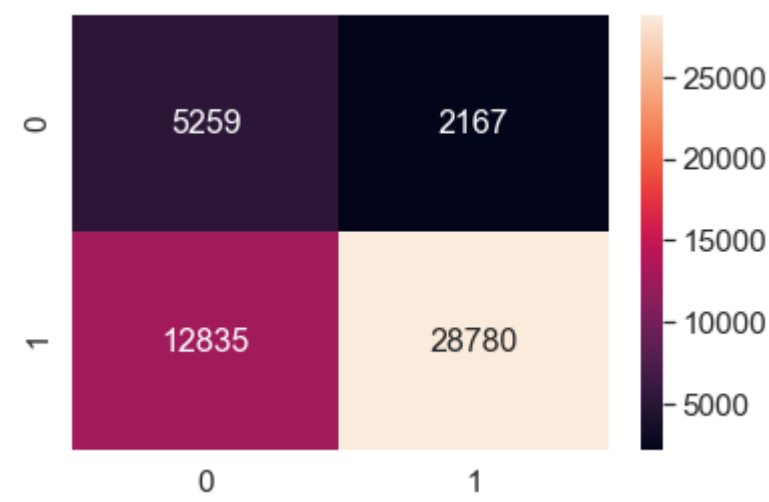
```
In [64]: ► 1 print("="*100)
2 from sklearn.metrics import confusion_matrix
3 best_t = find_best_threshold(tr_threshholds, train_fpr, train_tpr)
4 print("Train confusion matrix")
5 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
6 print("Test confusion matrix")
7 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.489766545636 for threshold 0.506
Train confusion matrix
[[ 5259  2167]
 [12835 28780]]
Test confusion matrix
[[ 3426  2033]
 [ 9774 20819]]
```

### Ploting Confusion Matrix on Train data

In [65]:

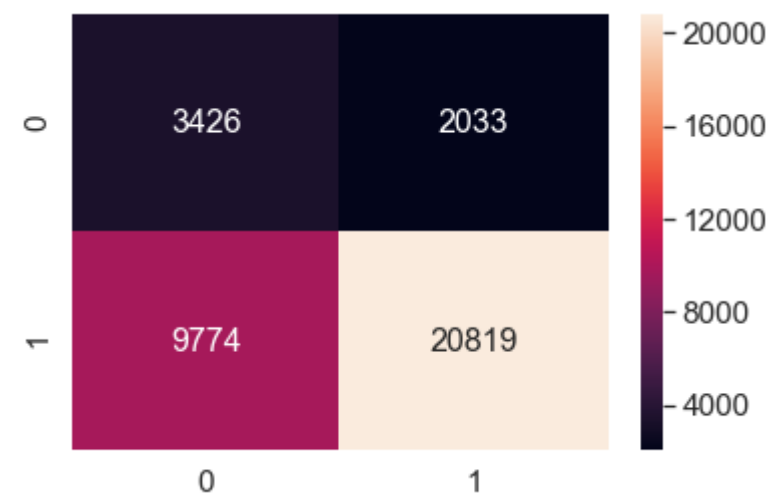
```
1  ### PLOT the matrix for Train
2  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3  df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4  # plt.figure(figsize=(10,7))
5  sns.set(font_scale=1.4) # for label size
6  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7  plt.show()
```



Plotting Cnfusion Matrix on Test Data

In [66]:

```
1  ### PLOT the matrix for Train
2  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3  df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4  # plt.figure(figsize=(10,7))
5  sns.set(font_scale=1.4) # for label size
6  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7  plt.show()
```



Observations :

- 1.We can observe from train and test we are getting majority True positives
- 2.Least number of data falls in False negative,which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
- 3.For a model to perform well we need High True Positive Rate and Low False Positive Rate.From the above our train data has True Positive Rate as 92% and False Positive Rate as 70%
- 4.In our test data : True Positive Rate as 91% and False Positive Rate as 73%.

**SET2 :Apply Multinomial NB on these feature sets (categorical, numerical features + preprocessed\_eassay (TFIDF))**

**Step 1 : Find best hyperparameter with maximum AUC**

```
In [67]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 from scipy.sparse import hstack
3
4
5 X_train_tfidf = hstack((train_categories, train_subcategories,sklstate_train,teacher_prefix_train,
6                        proj_grade_train,tfidf_essay_train,tfidf_title_train,
7                        X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
8                        essay_word_count_train_norm)).tocsr()
9
10 X_test_tfidf = hstack((test_categories, test_subcategories,sklstate_test,teacher_prefix_test,
11                       proj_grade_test,tfidf_essay_test,tfidf_title_test,
12                       X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
13                       essay_word_count_test_norm)).tocsr()
14
15 X_cv_tfidf = hstack((cv_categories, cv_subcategories,sklstate_cv,teacher_prefix_cv,
16                     proj_grade_cv,tfidf_essay_cv,tfidf_title_cv,
17                     X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
18                     essay_word_count_cv_norm)).tocsr()
19
20
21 print(X_train_tfidf.shape)
22 print(X_test_tfidf.shape)
23 print(X_cv_tfidf.shape)
```

(49041, 7113)

(36052, 7113)

(24155, 7113)

```
In [68]: 1 print("Final Data matrix")
2 print(X_train_tfidf.shape, y_train.shape)
3 print(X_cv_tfidf.shape, y_cv.shape)
4 print(X_test_tfidf.shape, y_test.shape)
5 print("=="*100)
```

Final Data matrix

(49041, 7113) (49041,)

(24155, 7113) (24155,)

(36052, 7113) (36052,)

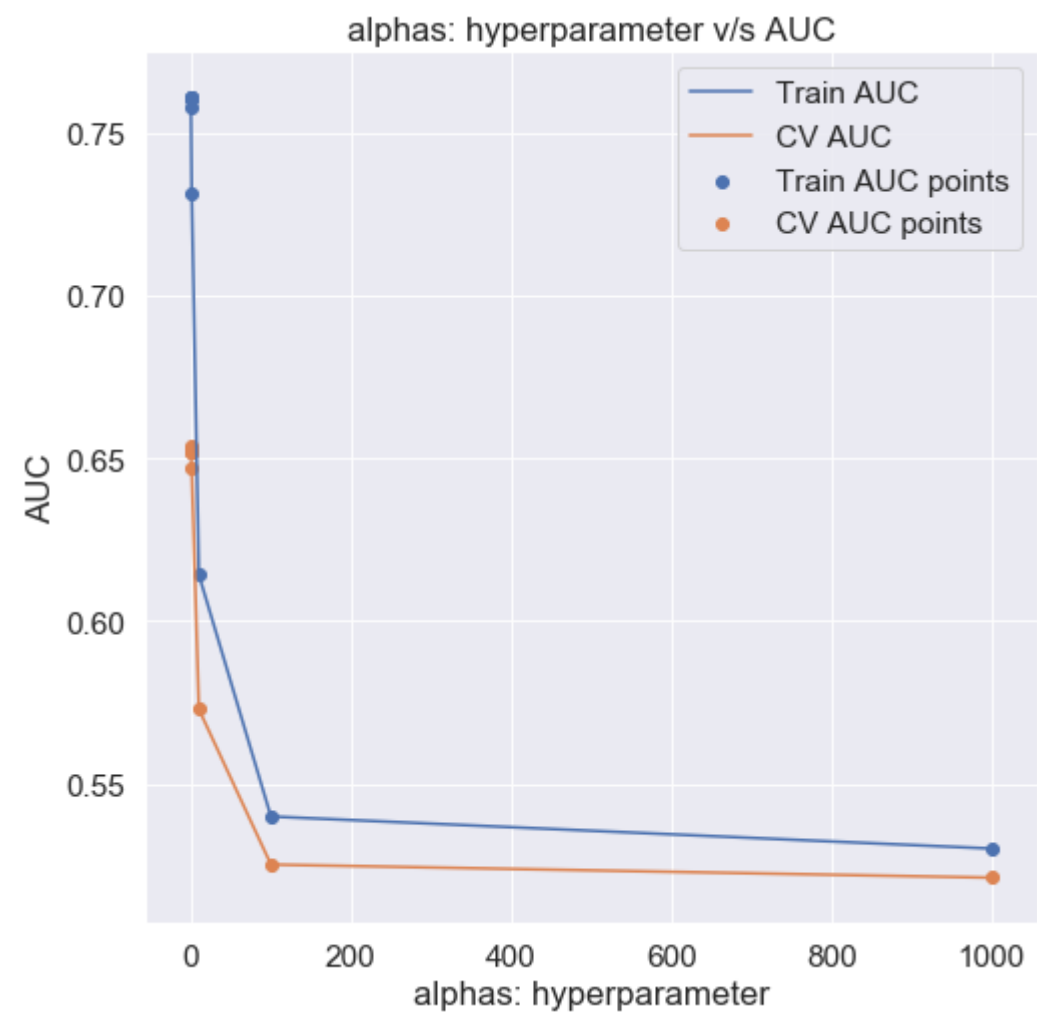
=====



100% |

In [73]:

```
1 plt.figure(figsize=(8,8))
2 plt.plot(alphas, train_auc, label='Train AUC')
3 plt.plot(alphas, cv_auc, label='CV AUC')
4
5
6 plt.scatter(alphas, train_auc, label='Train AUC points')
7 plt.scatter(alphas, cv_auc, label='CV AUC points')
8
9
10 plt.legend()
11 plt.xlabel("alphas: hyperparameter")
12 plt.ylabel("AUC")
13 plt.title("alphas: hyperparameter v/s AUC")
14 plt.grid(which='major', alpha=0.9)
15 plt.show()
```



**Observations :**

1. Above plot shows that Train\_data has higher AUC than Cross validation data.
2. There is a steep fall when alpha is greater than 0.001 which shows that as alpha increases beyond 0.001 performance reduces steeply.

**RandomSearch CV using K-fold Crossvalidation with k=10**

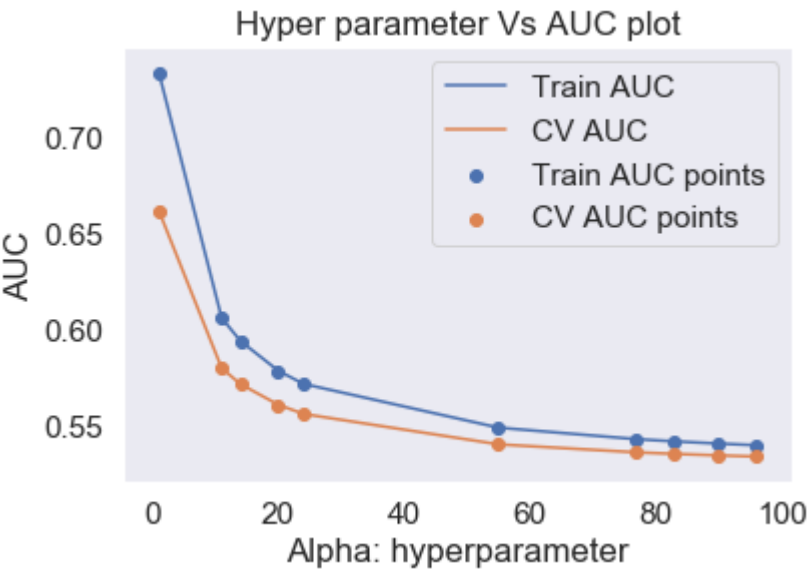
```
In [75]: 1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import RandomizedSearchCV
4
5 parameters={"alpha" : sp_randint(0.001,100) }
6
7 #RS_Log_alphas =[]
8
9 #RS_alphas=[0.00001, 0.0001, 0.001, 0.1,0.7,0.8, 1,100,1000]
10
11 # for a in tqdm(RS_alphas):
12 #     b = math.log(a)
13 #     RS_Log_alphas.append(b)
14
15 clf = RandomizedSearchCV(MB, parameters, cv=10,return_train_score=True,scoring='roc_auc',verbose=1,n_jobs=10)
16 clf.fit(X_train_tfidf, y_train)
17 results = pd.DataFrame.from_dict(clf.cv_results_)
18 results = results.sort_values(['param_alpha'])
19
20 RS_alphas=results['param_alpha']
21 train_auc= results['mean_train_score']
22 train_auc_std= results['std_train_score']
23 cv_auc = results['mean_test_score']
24 cv_auc_std= results['std_test_score']
25
26
27
```

Fitting 10 folds for each of 10 candidates, totalling 100 fits

```
[Parallel(n_jobs=10)]: Using backend LokyBackend with 10 concurrent workers.
[Parallel(n_jobs=10)]: Done 30 tasks      | elapsed:    1.8s
[Parallel(n_jobs=10)]: Done 100 out of 100 | elapsed:    5.2s finished
```

In [76]: ▶

```
1 plt.plot(RS_alphas, train_auc, label='Train AUC')
2 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
3 # plt.gca().fill_between(K, train_auc - train_auc_std,train_auc + train_auc_std,alpha=0.2,color='darkblue')
4
5 plt.plot(RS_alphas, cv_auc, label='CV AUC')
6 # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
7 # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')
8
9 plt.scatter(RS_alphas, train_auc, label='Train AUC points')
10 plt.scatter(RS_alphas, cv_auc, label='CV AUC points')
11
12
13 plt.legend()
14 plt.xlabel("Alpha: hyperparameter")
15 plt.ylabel("AUC")
16 plt.title("Hyper parameter Vs AUC plot")
17 plt.grid()
18 plt.show()
19
20 results.head()
```



Out[76]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	split2_train_score	split3_train_score	split4_train_score
1	0.353951	0.095472	0.024635	0.019147	1	{'alpha': 1}	0.656654	0.678948	0.671653	0.658960	...	0.734155	0.733256	0.7334
7	0.377292	0.057856	0.014861	0.010077	11	{'alpha': 11}	0.569192	0.589660	0.594393	0.585483	...	0.604553	0.605507	0.6077

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_alpha	params	split0_test_score	split1_test_score	split2_test_score	split3_test_score	...	split2_train_score	split3_train_score	split4_train_sc
3	0.320244	0.096062	0.014062	0.009307	14	{'alpha': 14}	0.560949	0.580587	0.586078	0.577900	...	0.592443	0.593452	0.5955
0	0.294511	0.080867	0.011968	0.003567	20	{'alpha': 20}	0.550755	0.568740	0.575343	0.568131	...	0.576866	0.577944	0.5807
2	0.333807	0.116230	0.012368	0.004305	24	{'alpha': 24}	0.546300	0.563554	0.570532	0.563775	...	0.569999	0.571092	0.5740

5 rows × 31 columns

Observations :

- 1. Above plot shows that Train\_data has higher AUC than Cross validation data.
- 2. There is a steep fall when alpha is greater than 0.001 which shows that as alpha increases beyond 0.001 performance reduces steeply.
- 3. Alpha at which cross validation AUC is maximum is 1.

Train the Model using the best hyper parameter value

In [77]: ▶

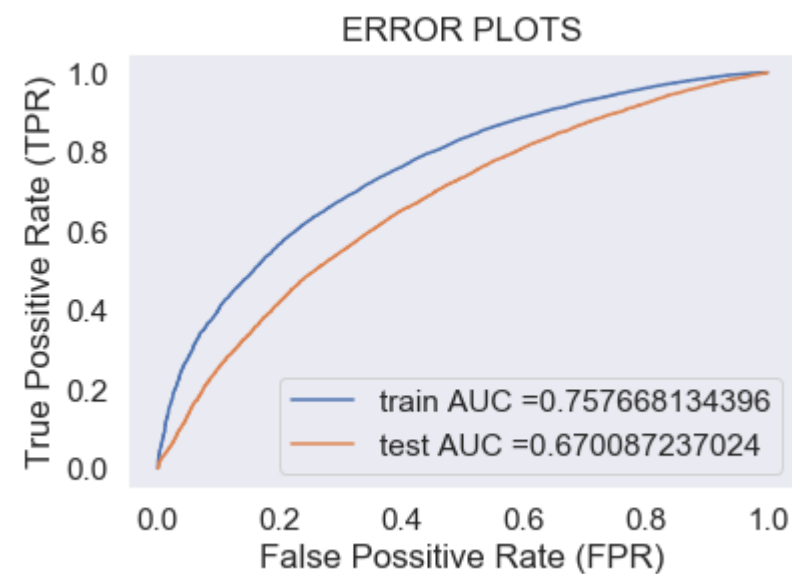
```
1 ### https://forums.fast.ai/t/hyperparameter-random-search-interpretation/8591 ---to get the best hyper parameter as a result of Random search
2 best_alpha = clf.best_params_
3 print('Best Alpha as a result of Random Search',best_alpha)
```

Best Alpha as a result of Random Search {'alpha': 1}

```

In [78]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3
4
5 MB_tfidf = MultinomialNB(alpha = 0.1,class_prior=[0.5,0.5])
6 MB_tfidf.fit(X_train_tfidf, y_train)
7 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
8 # not the predicted outputs
9
10 y_train_pred = batch_predict(MB_tfidf, X_train_tfidf)
11 y_test_pred = batch_predict(MB_tfidf, X_test_tfidf)
12
13 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
14 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
15
16 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
17 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
18 plt.legend()
19 plt.xlabel("False Possitive Rate (FPR)")
20 plt.ylabel("True Possitive Rate (TPR)")
21 plt.title("ERROR PLOTS")
22 plt.grid()
23 plt.show()

```



### Observations:

1. Train AUC observed to be 0.75 and Test AUC observed to be 0.67 when considering the best alpha.

### CONFUSION MATRIX

```

In [79]: ▶ 1 # we are writing our own function for predict, with defined threshould
           2 # we will pick a threshold that will give the least fpr
           3 def find_best_threshold(threshould, fpr, tpr):
           4     t = threshould[np.argmax(tpr*(1-fpr))]
           5     # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
           6     print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
           7     return t
           8
           9 def predict_with_best_t(proba, threshould):
          10     predictions = []
          11     for i in proba:
          12         if i>=threshould:
          13             predictions.append(1)
          14         else:
          15             predictions.append(0)
          16     return predictions

```

```

In [80]: ▶ 1 print("="*100)
           2 from sklearn.metrics import confusion_matrix
           3 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
           4 print("Train confusion matrix")
           5 print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
           6 print("Test confusion matrix")
           7 print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))

```

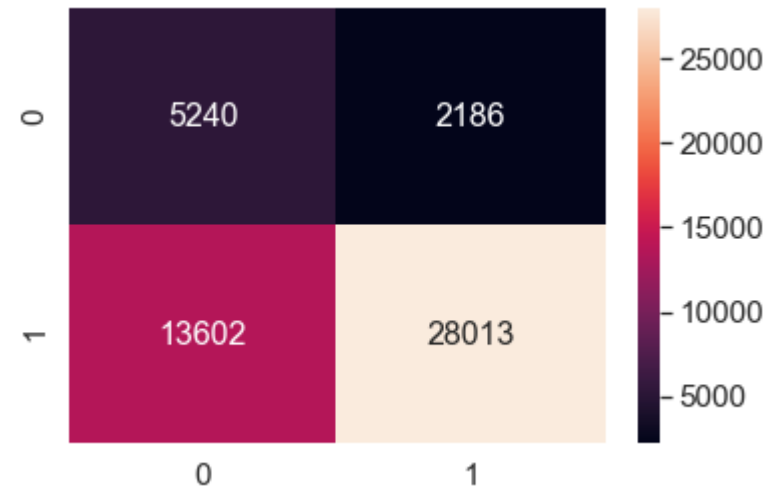
```

=====
the maximum value of tpr*(1-fpr) 0.474991747645 for threshold 0.505
Train confusion matrix
[[ 5240  2186]
 [13602 28013]]
Test confusion matrix
[[ 3240  2219]
 [10493 20100]]

```

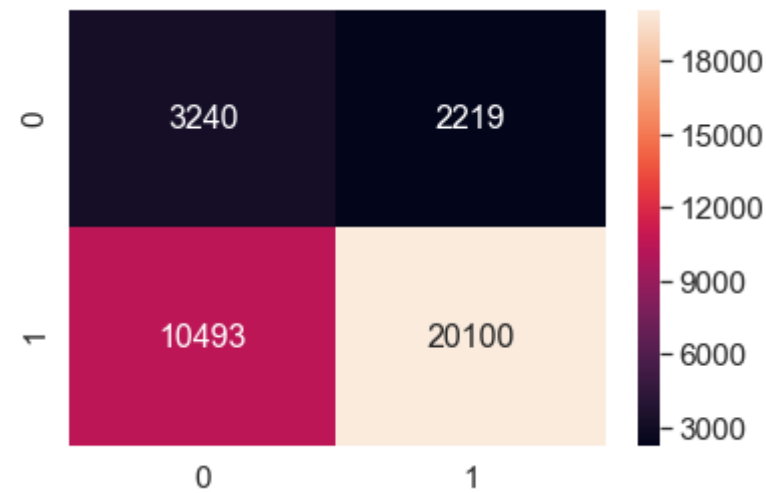
### Ploting Confusion Matrix on Train data

```
In [81]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



#### Plotting Confusion Matrix on Test Data

```
In [82]: 1 ### PLOT the matrix for Train
2 # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
3 df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)), range(2), range(2))
4 # plt.figure(figsize=(10,7))
5 sns.set(font_scale=1.4) # for label size
6 sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g') # font size
7 plt.show()
```



#### Observations :



1. We can observe from train and test we are getting majority True positives
2. Least number of data falls in False negative, which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
3. For a model to perform well we need High True Positive Rate and Low False Positive Rate. From the above our train data has True Positive Rate as 93% and False Positive Rate as 72%
4. In our test data : True Positive Rate as 90% and False Positive Rate as 76%.

#### Step 4: finding the top 20 features from either from feature Set 1 absolute values of feature\_log\_prob\_ parameter of MultinomialNB

```
In [83]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 ## Reference : https://datascience.stackexchange.com/questions/65219/find-the-top-n-features-from-feature-set-using-absolute-values-of-feature-log-p
3
4 from scipy.sparse import hstack
5
6
7 X_train_bow = hstack((train_categories, train_subcategories, sklstate_train, teacher_prefix_train,
8                       proj_grade_train, bow_essay_train, bow_title_train,
9                       X_train_price_norm, quantity_train_norm, prev_projects_train_norm, title_word_count_train_norm,
10                      essay_word_count_train_norm)).tocsr()
11
12 X_test_bow = hstack((test_categories, test_subcategories, sklstate_test, teacher_prefix_test,
13                     proj_grade_test, bow_essay_test, bow_title_test,
14                     X_test_price_norm, quantity_test_norm, prev_projects_test_norm, title_word_count_test_norm,
15                     essay_word_count_test_norm)).tocsr()
16
17 X_cv_bow = hstack((cv_categories, cv_subcategories, sklstate_cv, teacher_prefix_cv,
18                   proj_grade_cv, bow_essay_cv, bow_title_cv,
19                   X_cv_price_norm, quantity_cv_norm, prev_projects_cv_norm, title_word_count_cv_norm,
20                   essay_word_count_cv_norm)).tocsr()
21
22
23 print(X_train_bow.shape)
24 print(X_test_bow.shape)
25 print(X_cv_bow.shape)
```

(49041, 7113)  
(36052, 7113)  
(24155, 7113)

```
In [84]: 1 ## Train the model with the chosen hyper parameter
2 MB=MultinomialNB(class_prior=[0.5,0.5],alpha=0.1)
3 MB.fit(X_train_bow,y_train)
```

Out[84]: MultinomialNB(alpha=0.1, class\_prior=[0.5, 0.5])

```
In [85]: 1 ### sort all the features based on the log probabilities using argsort
2 # Positive class
3 class_1_sorted_prob=MB.feature_log_prob_[1,:].argsort()
4 class_0_sorted_prob=MB.feature_log_prob_[0,:].argsort()
```

```
In [86]: 1 features_lst=list(cat_vectorize.get_feature_names()+ subcat_vectorize.get_feature_names()+
2             sklstate_vectorize.get_feature_names()+teacher_prefix_vectorize.get_feature_names()+
3             proj_grade_vectorize.get_feature_names()+bow_essay.get_feature_names()+bowtitle.get_feature_names()+
4             ["price"]+['Quantity']+['teacher_number_of_previously_posted_projects']+['words_in_title']+
5             ['words_in_essay'])
```

```
In [87]: 1 Most_imp_words_1 = []
2 Most_imp_words_0 = []
3
4 for index in class_1_sorted_prob[-20:-1]:
5     Most_imp_words_1.append(features_lst[index])
6
7 for index in class_0_sorted_prob[-20:-1]:
8     Most_imp_words_0.append(features_lst[index])
9
10 print("20 most imp features for positive class:\n")
11 print(Most_imp_words_1)
12
13 print("\n" + "-"*100)
14
15 print("\n20 most imp features for negative class:\n")
16 print(Most_imp_words_0)
```

20 most imp features for positive class:

```
['technology', 'would', 'class', 'come', 'able', 'day', 'love', 'use', 'reading', 'work', 'need', 'nannan', 'many', 'help', 'learn', 'not', 'classroom', 'learning', 'school']
```

-----

20 most imp features for negative class:

```
['class', 'use', 'day', 'skills', 'able', 'reading', 'materials', 'love', 'come', 'work', 'need', 'many', 'nannan', 'help', 'learn', 'not', 'classroom', 'learning', 'school']
```

## Observations:

1. Both the negative and positive class tends to have similar words/features as most important with difference in ordering.

## Final Representation

In [88]:

```
1  ##http://zetcode.com/python/prettytable/
2
3  from prettytable import PrettyTable
4
5  x = PrettyTable()
6  x.field_names = ["Vectorizer", "Model", "Hyper Parameter", " Train AUC" ,"Test AUC "]
7
8  x.add_row(["BOW", "Multinomial Naive Bayes", 0.1, 0.76,0.70])
9  x.add_row(["TFIDF", "Multinomial Naive Bayes",0.1, 0.75, 0.67])
10
11 print(x)
```

Vectorizer	Model	Hyper Parameter	Train AUC	Test AUC
BOW	Multinomial Naive Bayes	0.1	0.76	0.7
TFIDF	Multinomial Naive Bayes	0.1	0.75	0.67