

## ▼ DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects. A large number of volunteers is needed to manually screen each submission before it's approved to be posted.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are several challenges to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be screened as quickly as possible
- How to increase the consistency of project vetting across different volunteers to improve the quality of the review process
- 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 will be approved based on the text of project descriptions as well as additional metadata about the project, teacher, and school. The goal is to 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> <ul style="list-style-type: none"><li>• Art Will Make You Happy!</li><li>• First Grade Fun</li></ul>
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following categories: <ul style="list-style-type: none"><li>• Grades PreK-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>

Feature	Description
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project for <ul style="list-style-type: none"> <li>• Applied Learning</li> <li>• Care &amp; Hunger</li> <li>• Health &amp; Sports</li> <li>• History &amp; Civics</li> <li>• Literacy &amp; Language</li> <li>• Math &amp; Science</li> <li>• Music &amp; The Arts</li> <li>• Special Needs</li> <li>• Warmth</li> </ul> <p><b>Examples:</b></p> <ul style="list-style-type: none"> <li>• Music &amp; The Arts</li> <li>• Literacy &amp; Language, Math &amp; Science</li> </ul>
<code>school_state</code>	State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> <code>CA</code>
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project <ul style="list-style-type: none"> <li>• Literacy</li> <li>• Literature &amp; Writing, Social Sciences</li> </ul>
<code>project_resource_summary</code>	An explanation of the resources needed for the project. <b>Example:</b> <ul style="list-style-type: none"> <li>• My students need hands on literacy materi</li> </ul>
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*
<code>project_essay_4</code>	Fourth application essay*
<code>project_submitted_datetime</code>	Datetime when project application was submitted. <b>Example:</b> <code>2016-04-15 14:30:00</code>
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. <b>Example:</b> <code>000001</code>
<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

\* 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 row in the data set 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> <code>p036502</code>

## ▼ 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 following:

- \_\_project\_essay\_1:\_\_ "Describe your students: What makes your students special? Specific details about your neighborhood, and your school are all helpful."
- \_\_project\_essay\_2:\_\_ "About your project: How will these materials make a difference in your students' lives?"

For all projects with project\_submitted\_datetime of 2016-05-17 and later, the values of project\_essay\_1 and project\_essay\_2 are:

```
%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

```

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

```

```

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']

```

```

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']
```

	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

```

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

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in catogories:
    temp = ""
    # consider we have text like this "Math & Science. Warmth. Care & Hunger"

```

```

for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Ca
    if 'The' in j.split(): # this will split each of the category based on space "Math &
        j=j.replace('The','') # if we have the words "The" we are going to replace it wit
    j = j.replace(' ', '') # we are replacing all the ' '(space) with ''(empty) ex:"Math &
    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

```

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

# 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", "Ca
        if 'The' in j.split(): # this will split each of the category based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
        j = j.replace(' ', '') # we are replacing all the ' '(space) with ''(empty) ex:"Math &
        temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
        temp = temp.replace('&','_')
    sub_cat_list.append(temp.strip())

```

```

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

```

```

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

```

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

## ▼ 1.3 Text preprocessing

```
# 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)
```

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



Unnamed:  
0

id

teacher\_id

teacher\_prefix

school\_state

pr

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
---	--------	---------	----------------------------------	------	----

1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL
---	--------	---------	----------------------------------	-----	----

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

```
# 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 la
=====
The 51 fifth grade students that will cycle through my classroom this year all love lear
=====
How do you remember your days of school? Was it in a sterile environment with plain wall
=====
My kindergarten students have varied disabilities ranging from speech and language delay
=====
The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates
=====

```

# <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
```

```
sent = decontracted(project_data['essay'].values[20000])
```

```
print(sent)
```

```
print("="*50)
```



My kindergarten students have varied disabilities ranging from speech and language delay  
=====

# \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 delay

#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 delay

```
# 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", "yo  
"you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',  
'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they',  
'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll"  
'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'h  
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'unt  
'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'dur  
'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', '  
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'bo  
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'ver  
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd  
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'does  
'hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "  
'mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',  
'won', "won't", 'wouldn', "wouldn't"]
```

```
# Combining all the above stundents
```

```
from tqdm import tqdm
```

```
preprocessed_essays = []
```

```
# tqdm is for printing the status bar
```

```
for sentence in tqdm(project_data['essay'].values):
```

```
    sent = decontracted(sentence)
```

```
    sent = sent.replace('\r', ' ')
```

```
    sent = sent.replace('\\"', ' ')
```

```
    sent = sent.replace('\n', ' ')
```

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

```
    # https://gist.github.com/sebleier/554280
```


```
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
```

```
    preprocessed_essays.append(sent.lower().strip())
```

 100% |  | 109248/10

```
# after preprocesing
```

```
preprocessed_essays[20000]
```

 'my kindergarten students varied disabilities ranging speech language delays cognitive d

## 1.4 Preprocessing of `project\_title`

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



## ▼ 1.5 Preparing data for models

project\_data.columns

```
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')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
  
- project\_title : text data
- text : text data
- project\_resource\_summary: text data (optinal)
  
- quantity : numerical (optinal)
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

### ▼ 1.5.1 Vectorizing Categorical data


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

```
# 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
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNee
Shape of matrix after one hot encodig (109248, 9)
```

```
# we use count vectorizer to convert the values into one
```

```
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, bi
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.shape)
```


 ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurric  
Shape of matrix after one hot encodig (109248, 30)

# you can do the similar thing with state, teacher\_prefix and project\_grade\_category also

## ▼ 1.5.2 Vectorizing Text data

### ▼ 1.5.2.1 Bag of words


```
# 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 encodig ",text_bow.shape)
```

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

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

### ▼ 1.5.2.2 TFIDF vectorizer

```
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 encodig ",text_tfidf.shape)
```

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

### ▼ 1.5.2.3 Using Pretrained Models: Avg W2V

```
...
# 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-

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

...

 '\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\n\ndef

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-
# make sure you have the glove_vectors file

```

```

with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors.append(vector)

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

```



100% | 109248/10  
109248  
300

### ▼ 1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

```

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

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

```



```
print(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)
```

```
# 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
```

```
(109248, 16663)
```

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

## ▼ Assignment 4: Naive Bayes

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

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




### 2. The hyper parameter tuning(find best Alpha)

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


### 3. Feature importance

- Find the top 10 features of positive class and top 10 features of negative class for both fe  
`feature\_log\_prob\_` parameter of [MultinomialNB](#) and print their corresponding feature name

### 4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for figure. Here on X-axis you will have alpha values, since they have a wide range, just to represent 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 curve on both train and test. 
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and actual and visualize your confusion matrices using [seaborn heatmaps](#). 

## 5. [Conclusion](#)

- [You need to summarize the results at the end of the notebook, summarize it in the table for this prettytable library link](#) 

# 2. Naive Bayes

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

```
%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_notebook as tqdm1
from tqdm import tqdm
import time
import os
```


```
from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
```

```
from sklearn.model_selection import train_test_split
```

 C:\Users\LENOVO\Anaconda3\lib\site-packages\smart\_open\ssh.py:34: UserWarning: paramiko warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip in C:\Users\LENOVO\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected warnings.warn("detected Windows; aliasing chunkize to chunkize\_serial")

```
project_data = pd.read_csv('train_data.csv', nrows=50000)
resource_data = pd.read_csv('resources.csv')
```

```
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 (50000, 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']
```

```
# not_accepted = project_data[project_data.project_is_approved==0]
# accepted = project_data[project_data.project_is_approved==1]
```

```
# print(accepted.shape)
# print(not_accepted.shape)
```

```
# # https://towardsdatascience.com/methods-for-dealing-with-imbalanced-data-5b761be45a18
# # Upsampling minority class
# from sklearn.utils import resample
# not_accepted_upsampled = resample(not_accepted,
```



```

# replace=True, # sample with replacement
# n_samples=len(accepted), # match number in majority class
# random_state=27) # reproducible results

# # combine majority and upsampled minority
# project_data = pd.concat([accepted, not_accepted_upsampled])

# # check new class counts
# print(project_data.project_is_approved.value_counts())
# print(project_data.shape)

```

## ▼ Text preprocessing(1)

```

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

# 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", "Ca
        if 'The' in j.split(): # this will split each of the category based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math &
            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)
project_data.head(5)

```



	Unnamed: 0		id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a		Mr.	FL	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0		Ms.	AZ	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60		Mrs.	KY	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec		Mrs.	TX	

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

```
Counter({'Literacy_Language': 23998,
        'History_Civics': 2689,
        'Health_Sports': 6538,
        'Math_Science': 18874,
        'SpecialNeeds': 6233,
        'AppliedLearning': 5569,
        'Music_Arts': 4699,
        'Warmth': 643,
        'Care_Hunger': 643})
```

```
# dict sort by value python: https://stackoverflow.com/a/613218/4084039
cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

```
# ind = np.arange(len(sorted_cat_dict))
# plt.figure(figsize=(20,5))
# p1 = plt.bar(ind, list(sorted_cat_dict.values()))

# plt.ylabel('Projects')
# plt.title('% of projects aproved category wise')
# plt.xticks(ind, list(sorted_cat_dict.keys()))
# plt.show()
# print(sorted_cat_dict)

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

# 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", "Ca
        if 'The' in j.split(): # this will split each of the catogory based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
            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)
project_data.head(2)
```



Unnamed:  
0

id

teacher\_id teacher\_prefix school\_state pr

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL

```
# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
from collections import Counter
my_counter = Counter()
```

```

for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

# dict sort by value python: https://stackoverflow.com/a/613218/4084039
sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

# ind = np.arange(len(sorted_sub_cat_dict))
# plt.figure(figsize=(20,5))
# p1 = plt.bar(ind, list(sorted_sub_cat_dict.values()))

# plt.ylabel('Projects')
# plt.title('% of projects aproved state wise')
# plt.xticks(ind, list(sorted_sub_cat_dict.keys()))
# plt.show()

# 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)

# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-grou
price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
price_data.head(2)

```

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

```

# join two dataframes in python:
project_data = pd.merge(project_data, price_data, on='id', how='left')

#presence of the numerical digits in a strings with numeric : https://stackoverflow.com/a/198
def hasNumbers(inputString):
    return any(i.isdigit() for i in inputString)
p1 = project_data[['id','project_resource_summary']]
p1 = pd.DataFrame(data=p1)
p1.columns = ['id','digits_in_summary']
p1['digits_in_summary'] = p1['digits_in_summary'].map(hasNumbers)
# https://stackoverflow.com/a/17383325/8089731
p1['digits_in_summary'] = p1['digits_in_summary'].astype(int)
project_data = pd.merge(project_data, p1, on='id', how='left')
project_data.head(5)

```

	Unnamed: 0		id	teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a		Mr.	FL	
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0		Ms.	AZ	
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60		Mrs.	KY	
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec		Mrs.	TX	

5 rows × 21 columns

## ▼ Text preprocessing(2)

```
# 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", " ", 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

```

# <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", "yo
    "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they',
    'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll"
    'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'h
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'unt
    'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'dur
    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', '
    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'bo
    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'ver
    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd
    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'does
    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "
    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',
    'won', "won't", 'wouldn', "wouldn't"]

```

# Combining all the above statemennts

```
from tqdm import tqdm
```

```
preprocessed_essays = []
```

# tqdm is for printing the status bar

```
for sentence in tqdm(project_data['essay'].values):
```

```
    sent = decontracted(sentence)
```

```
    sent = sent.replace('\r', ' ')
```

```
    sent = sent.replace('\n', ' ')
```

```
    sent = sent.replace('\n', ' ')
```

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

```
    sent = re.sub('nannan', '', sent)
```

# <https://gist.github.com/sebleier/554280>

```
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
```

```
    preprocessed_essays.append(sent.lower().strip())
```




HBox(children=(IntProgress(value=0, max=50000), HTML(value='')))

# preprocessed\_essays

```
from tqdm import tqdm
```

```
preprocessed_titles = []
```

```
# tqdm is for printing the status bar
for title in tqdm1(project_data['project_title'].values):
    _title = decontracted(title)
    _title = _title.replace('\\r', ' ')
    _title = _title.replace('\\\"', ' ')
    _title = _title.replace('\\n', ' ')
    _title = re.sub('[^A-Za-z0-9]+', ' ', _title)
    # https://gist.github.com/sebleier/554280
    _title = ' '.join(e for e in _title.split() if e not in stopwords)
    preprocessed_titles.append(_title.lower().strip())
```

 HBox(children=(IntProgress(value=0, max=50000), HTML(value='')))


```
preprocessed_titles[1000]
```

 'sailing into super 4th grade year'

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

# 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
```

```
project_grade_cat_list = []
for i in tqdm1(project_grade_catogories):
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Ca
        if 'The' in j.split(): # this will split each of the catogory based on space "Math &
            j=j.replace('The','') # if we have the words "The" we are going to replace it wit
            j = j.replace(' ','') # we are placeing all the ' '(space) with ''(empty) ex:"Math &
            temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&','_')
    project_grade_cat_list.append(temp.strip())
```

 HBox(children=(IntProgress(value=0, max=50000), HTML(value='')))

```
project_data['clean_project_grade_category'] = project_grade_cat_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data.head(2)
```



	Unnamed: 0	id		teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a		Mr.	FL	

2 rows × 21 columns

```
project_data.drop(['project_essay_1','project_essay_2','project_essay_3','project_essay_4'],
project_data.head(2)
```



	Unnamed: 0	id		teacher_id	teacher_prefix	school_state	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	
1	140945	p258326	897464ce9ddc600bced1151f324dd63a		Mr.	FL	

```
#Replacing Nan's with maximum occured value: https://stackoverflow.com/a/51053916/8089731
project_data['teacher_prefix'].value_counts().argmax()
project_data.fillna(value=project_data['teacher_prefix'].value_counts().argmax(),axis=1,inpla

project_data['preprocessed_essays'] = preprocessed_essays
project_data['preprocessed_titles'] = preprocessed_titles
```



```
project_data.columns
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_title',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay', 'price', 'quantity',
      'digits_in_summary', 'clean_project_grade_category',
      'preprocessed_essays', 'preprocessed_titles'],
      dtype='object')
```

```
# 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
```

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

```
X_train, X_test, y_train, y_test = train_test_split(project_data,project_data['project_is_app
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_
```

```
X_train.drop(['project_is_approved'], axis=1, inplace=True)
X_test.drop(['project_is_approved'], axis=1, inplace=True)
X_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

### ▼ 1.4.1 Vectorizing Categorical data


```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, bi
vectorizer_cat.fit(X_train['clean_categories'].values)
print(vectorizer_cat.get_feature_names())
```

```
categories_one_hot_train = vectorizer_cat.transform(X_train['clean_categories'].values)
categories_one_hot_cv = vectorizer_cat.transform(X_cv['clean_categories'].values)
categories_one_hot_test = vectorizer_cat.transform(X_test['clean_categories'].values)
print("Shape of matrix after one hot encodig_train ",categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",categories_one_hot_test.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNee
Shape of matrix after one hot encodig_train (22445, 9)
Shape of matrix after one hot encodig_cv (11055, 9)
Shape of matrix after one hot encodig_test (16500, 9)
```


```
# we use count vectorizer to convert the values into one hot encoded features
vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=F
vectorizer_sub_cat.fit(X_train['clean_subcategories'].values)
print(vectorizer_sub_cat.get_feature_names())
```

```
sub_categories_one_hot_train = vectorizer_sub_cat.transform(X_train['clean_subcategories'].va
sub_categories_one_hot_cv = vectorizer_sub_cat.transform(X_cv['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer_sub_cat.transform(X_test['clean_subcategories'].valu
print("Shape of matrix after one hot encodig_train ",sub_categories_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",sub_categories_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",sub_categories_one_hot_test.shape)
```

```
 ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurric
Shape of matrix after one hot encodig_train (22445, 30)
Shape of matrix after one hot encodig_cv (11055, 30)
Shape of matrix after one hot encodig_test (16500, 30)
```

```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_state = CountVectorizer( lowercase=False, binary=True)
vectorizer_state.fit(X_train['school_state'].values)
print(vectorizer_state.get_feature_names())
```


```
school_state_one_hot_train = vectorizer_state.transform(X_train['school_state'].values)
school_state_one_hot_cv = vectorizer_state.transform(X_cv['school_state'].values)
school_state_one_hot_test = vectorizer_state.transform(X_test['school_state'].values)
print("Shape of matrix after one hot encodig_train ",school_state_one_hot_train.shape)
print("Shape of matrix after one hot encodig_cv ",school_state_one_hot_cv.shape)
print("Shape of matrix after one hot encodig_test ",school_state_one_hot_test.shape)
```

```
 ['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL
Shape of matrix after one hot encodig_train (22445, 51)
Shape of matrix after one hot encodig_cv (11055, 51)
Shape of matrix after one hot encodig_test (16500, 51)
```

```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
vectorizer_teacherprefix = CountVectorizer( lowercase=False, binary=True)
vectorizer_teacherprefix.fit(X_train['teacher_prefix'].values.astype('U'))
print(vectorizer_teacherprefix.get_feature_names())
```


```
#https://stackoverflow.com/a/39308809/8089731
teacher_prefix_one_hot_train = vectorizer_teacherprefix.transform(X_train['teacher_prefix'].values)
teacher_prefix_one_hot_cv = vectorizer_teacherprefix.transform(X_cv['teacher_prefix'].values)
teacher_prefix_one_hot_test = vectorizer_teacherprefix.transform(X_test['teacher_prefix'].values)
print("Shape of matrix after one hot encoding_train ",teacher_prefix_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ",teacher_prefix_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ",teacher_prefix_one_hot_test[:5,:])
# print(X_train['teacher_prefix'].value_counts())
```

```

 ['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
Shape of matrix after one hot encoding_train (22445, 5)
Shape of matrix after one hot encoding_cv (11055, 5)
Shape of matrix after one hot encoding_test (0, 1) 1
(1, 2) 1
(2, 2) 1
(3, 3) 1
(4, 2) 1
```

```
print(project_data['clean_project_grade_category'].unique())
```


```

 ['GradesPreK-2' 'Grades6-8' 'Grades3-5' 'Grades9-12']
```

```
# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
# https://stackoverflow.com/a/38161028/8089731
pattern = "(?u)\\b[\\w-]+\\b"
vectorizer_projectgrade = CountVectorizer(token_pattern=pattern, lowercase=False, binary=True)
vectorizer_projectgrade.fit(X_train['clean_project_grade_category'].values)
print(vectorizer_projectgrade.get_feature_names())
```

```
#https://stackoverflow.com/a/39308809/8089731
project_grade_category_one_hot_train = vectorizer_projectgrade.transform(X_train['clean_project_grade_category'].values)
project_grade_category_one_hot_cv = vectorizer_projectgrade.transform(X_cv['clean_project_grade_category'].values)
project_grade_category_one_hot_test = vectorizer_projectgrade.transform(X_test['clean_project_grade_category'].values)
print("Shape of matrix after one hot encoding_train ",project_grade_category_one_hot_train.shape)
print("Shape of matrix after one hot encoding_cv ",project_grade_category_one_hot_cv.shape)
print("Shape of matrix after one hot encoding_test ",project_grade_category_one_hot_test[:5,:])
```

```

 ['Grades3-5', 'Grades6-8', 'Grades9-12', 'GradesPreK-2']
Shape of matrix after one hot encoding_train (22445, 4)
Shape of matrix after one hot encoding_cv (11055, 4)
Shape of matrix after one hot encoding_test (0, 0) 1
(1, 0) 1
(2, 0) 1
(3, 3) 1
(4, 3) 1
```

## ▼ Vectorizing Numerical features

# check this one: <https://www.youtube.com/watch?v=0H0q0c1n3Z4&t=530s>

# standardization sklearn: <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html>

```

from sklearn.preprocessing import StandardScaler

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

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

# Now standardize the data with above maen and variance.
price_standardized_train = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
price_standardized_cv = price_scalar.transform(X_cv['price'].values.reshape(-1, 1))
price_standardized_test = price_scalar.transform(X_test['price'].values.reshape(-1, 1))
print(price_standardized_train.shape)
print(price_standardized_cv.shape)
print(price_standardized_test.shape)

```



Mean : 300.23218311427934, Standard deviation : 401.43554589763  
 (22445, 1)  
 (11055, 1)  
 (16500, 1)

```

# check this one: https://www.youtube.com/watch?v=0H0q0cIn3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preproce
from sklearn.preprocessing import StandardScaler

```

```

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

quantity_scalar = StandardScaler()
quantity_scalar.fit(X_train['quantity'].values.reshape(-1,1)) # finding the mean and standard
# print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(quantity_scalar.va

# Now standardize the data with above maen and variance.
quantity_standardized_train = quantity_scalar.transform(X_train['quantity'].values.reshape(-1
quantity_standardized_cv = quantity_scalar.transform(X_cv['quantity'].values.reshape(-1, 1))
quantity_standardized_test = quantity_scalar.transform(X_test['quantity'].values.reshape(-1,
print(quantity_standardized_train.shape)
print(quantity_standardized_cv.shape)
print(quantity_standardized_test.shape)

```



```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
(22445, 1)
(11055, 1)
(16500, 1)
```

```
# check this one: https://www.youtube.com/watch?v=0H0qOcln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler
from sklearn.preprocessing import StandardScaler
```

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

```
teacher_number_of_previously_posted_projects_scalar = StandardScaler()
teacher_number_of_previously_posted_projects_scalar.fit(X_train['teacher_number_of_previously_posted_projects'])
# print(f"Mean : {teacher_number_of_previously_posted_projects_scalar.mean_[0]}, Standard dev
```

```
# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized_train = teacher_number_of_previously_posted_projects_scalar.transform(X_train['teacher_number_of_previously_posted_projects'])
teacher_number_of_previously_posted_projects_standardized_cv = teacher_number_of_previously_posted_projects_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'])
teacher_number_of_previously_posted_projects_standardized_test = teacher_number_of_previously_posted_projects_scalar.transform(X_test['teacher_number_of_previously_posted_projects'])
print(teacher_number_of_previously_posted_projects_standardized_train.shape)
print(teacher_number_of_previously_posted_projects_standardized_cv.shape)
print(teacher_number_of_previously_posted_projects_standardized_test.shape)
```



```
C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\LENOVO\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595: DataConversionWarning: Data with input dtype int64 was converted to float64 by StandardScaler.

(22445, 1)
(11055, 1)
(16500, 1)
```

```
# 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
```

## 2.3 Make Data Model Ready: encoding eassay, and project\_title

```
X_train.head(2)
```



Unnamed: 0		id	teacher_id	teacher_prefix	school_stat
6025	178512	p148474	b8274a422e6a205748d98b23a3c134e9	Mrs.	N'
26403	122685	p202047	df992790283e68a6338d9f9a86f14cb6	Mrs.	I

## ▼ Bag of Words(BOW) on project\_TEXT/ESSAYS (Train,Cv,Test)

```
# We are considering only the words which appeared in at least 10 documents(rows or projects)
vectorizer_bow_essays = CountVectorizer(min_df=10)
vectorizer_bow_essays.fit(X_train['preprocessed_essays'])
```

```
text_bow_train = vectorizer_bow_essays.transform(X_train['preprocessed_essays'])
text_bow_cv = vectorizer_bow_essays.transform(X_cv['preprocessed_essays'])
text_bow_test = vectorizer_bow_essays.transform(X_test['preprocessed_essays'])
print("Shape of matrix after BOW_text_train ",text_bow_train.shape)
print("Shape of matrix after BOW_text_cv ",text_bow_cv.shape)
print("Shape of matrix after BOW_text_test ",text_bow_test.shape)
```




```
Shape of matrix after BOW_text_train (22445, 8857)
Shape of matrix after BOW_text_cv (11055, 8857)
Shape of matrix after BOW_text_test (16500, 8857)
```

## ▼ Bag of Words(BOW) on project\_title (Train,Cv,Test)

```
# We are considering only the words which appeared in at least 10 documents(rows or projects)
vectorizer_bow_titles = CountVectorizer(min_df=10)
vectorizer_bow_titles.fit(X_train['preprocessed_titles'])
```

```
title_bow_train = vectorizer_bow_titles.transform(X_train['preprocessed_titles'])
title_bow_cv = vectorizer_bow_titles.transform(X_cv['preprocessed_titles'])
title_bow_test = vectorizer_bow_titles.transform(X_test['preprocessed_titles'])
print("Shape of matrix after BOW title train ".title bow train.shape)
```


```
print("Shape of matrix after BOW_title_cv ",title_bow_cv.shape)
print("Shape of matrix after BOW_title_test ",title_bow_test.shape)
```

 Shape of matrix after BOW\_title\_train (22445, 1250)  
 Shape of matrix after BOW\_title\_cv (11055, 1250)  
 Shape of matrix after BOW\_title\_test (16500, 1250)

## ▼ TFIDF Vectorizer on project\_TEXT/ESSAYS (Train,Cv,Test)

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essays = TfidfVectorizer(min_df=10)
vectorizer_tfidf_essays.fit(X_train['preprocessed_essays'])


text_tfidf_train = vectorizer_tfidf_essays.transform(X_train['preprocessed_essays'])
text_tfidf_cv = vectorizer_tfidf_essays.transform(X_cv['preprocessed_essays'])
text_tfidf_test = vectorizer_tfidf_essays.transform(X_test['preprocessed_essays'])
print("Shape of matrix after tfidf_text_train ",text_tfidf_train.shape)
print("Shape of matrix after tfidf_text_cv ",text_tfidf_cv.shape)
print("Shape of matrix after tfidf_text_test ",text_tfidf_test.shape)
```

 Shape of matrix after tfidf\_text\_train (22445, 8857)  
 Shape of matrix after tfidf\_text\_cv (11055, 8857)  
 Shape of matrix after tfidf\_text\_test (16500, 8857)

## ▼ TFIDF Vectorizer on project\_title (Train,Cv,Test)

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_title = TfidfVectorizer(min_df=10)
vectorizer_tfidf_title.fit(X_train['preprocessed_titles'])

title_tfidf_train = vectorizer_tfidf_title.transform(X_train['preprocessed_titles'])
title_tfidf_cv = vectorizer_tfidf_title.transform(X_cv['preprocessed_titles'])
title_tfidf_test = vectorizer_tfidf_title.transform(X_test['preprocessed_titles'])
print("Shape of matrix after tfidf_title_train ",title_tfidf_train.shape)
print("Shape of matrix after tfidf_title_cv ",title_tfidf_cv.shape)
print("Shape of matrix after tfidf_title_test ",title_tfidf_test.shape)
```

 Shape of matrix after tfidf\_title\_train (22445, 1250)  
 Shape of matrix after tfidf\_title\_cv (11055, 1250)  
 Shape of matrix after tfidf\_title\_test (16500, 1250)

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



```
# 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
```

```
import dill
# dill.dump_session('notebook_env.db')
dill.load_session('../notebook_env.db')
```



```
C:\Users\LENOVO\Anaconda3\lib\site-packages\smart_open\ssh.py:34: UserWarning: paramiko
warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip in
C:\Users\LENOVO\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected
warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

## 2.4 Appling NB() on different kind of featurization as mentione

Apply Naive Bayes on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instrucations

### ► 2.4.1 Applying Naive Bayes on BOW, SET 1

↳ 2 cells hidden

### ▼ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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 unti 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:])[0,1])
```

```
y_data_pred.extend(cir.predict_proba(data[tr_loop:])[1:,1])
return y_data_pred

import matplotlib.pyplot as plt
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import roc_auc_score
import math
from sklearn.model_selection import RandomizedSearchCV

train_auc = []
cv_auc = []
log_alphas = []
alphas = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50,

for i in tqdm(alphas):
    nb = MultinomialNB(alpha = i,class_prior=[0.5,0.5])
    nb.fit(X_tr, y_train)
    y_train_pred = batch_predict(nb, X_tr)
    y_cv_pred = batch_predict(nb, X_cr)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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))

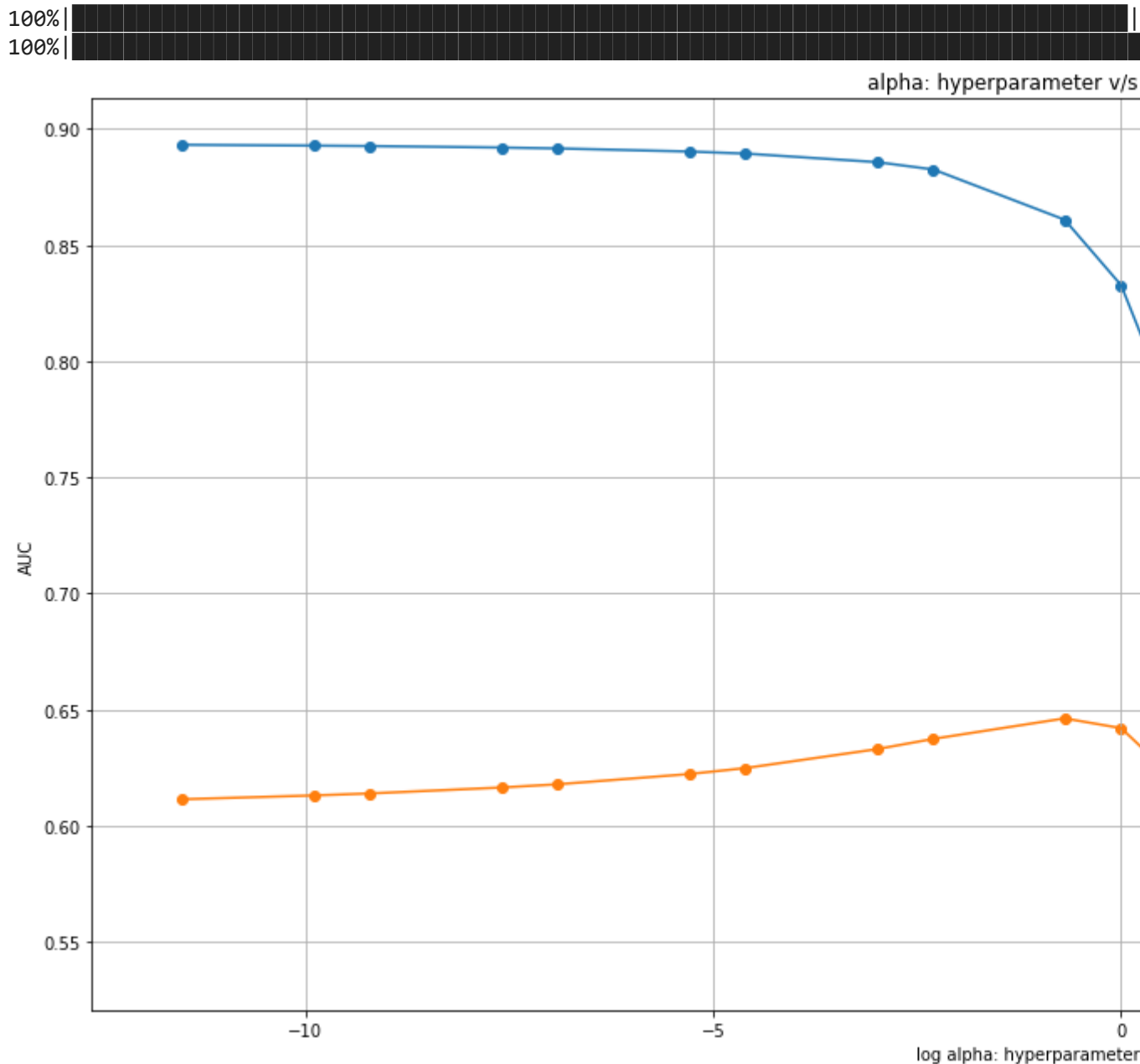
for a in tqdm(alphas):
    b = math.log(a)
    log_alphas.append(b)

log_alphas = np.array(log_alphas)
alphas = np.array(alphas)

plt.figure(figsize=(20,10))
plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')
plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("log alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC")
plt.grid()
plt.show()

# print(len(log_alphas))
# print(log_alphas.shape)
# print(train_auc.shape)
```





# [https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.metrics.roc\\_curve](https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve)  
 from sklearn.metrics import roc\_curve, auc

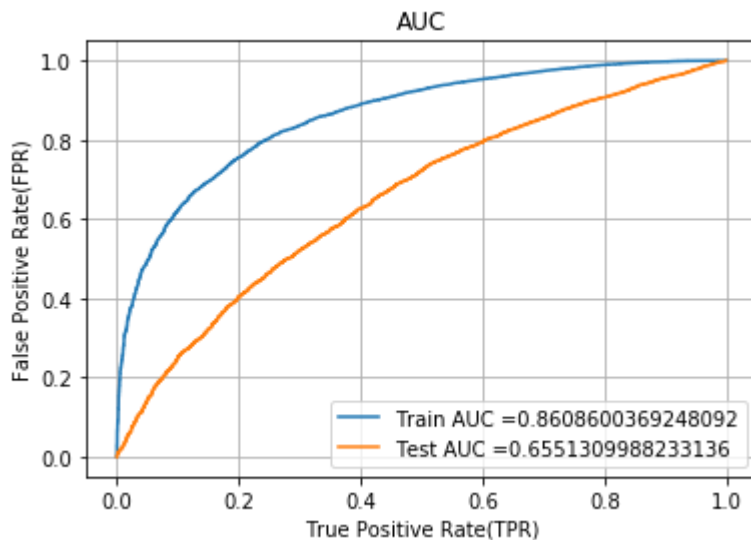
```
nb_bow = MultinomialNB(alpha = 0.5,class_prior=[0.5,0.5])
nb_bow.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs
```

```
y_train_pred = batch_predict(nb_bow, X_tr)
y_test_pred = batch_predict(nb_bow, X_te)
```

```
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("AUC")
plt.grid()
plt.show()
```



```
# 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

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, test_fpr, test_tpr)))
```



```
=====
```

Train confusion matrix

the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.6068681325618965 for threshold 0.49

```
[[ 2654   809]
 [ 3951 15031]]
```

Test confusion matrix

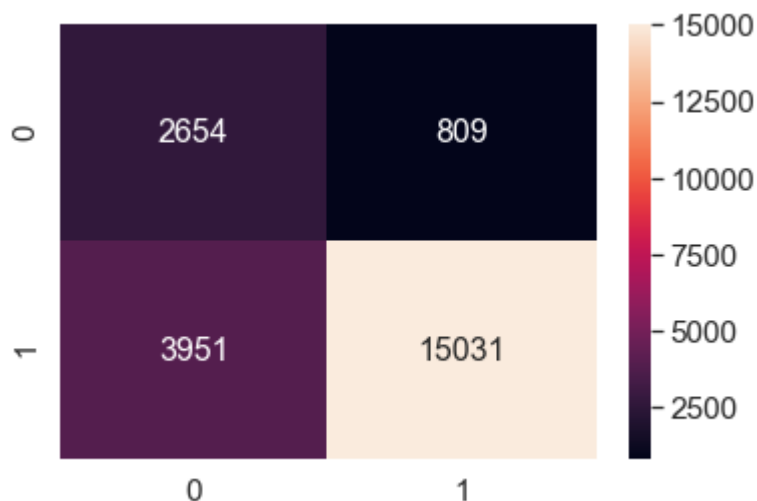
the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.3774556192431624 for threshold 0.422

```
[[ 1045   1501]
 [ 2968 10986]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.6068681325618965 for threshold 0.49  
<matplotlib.axes.\_subplots.AxesSubplot at 0x1f127fb4b38>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



# Please write all the code with proper documentation

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

# merge two sparse matrices: <https://stackoverflow.com/a/19710648/4084039>

```
from scipy.sparse import hstack
```

```
X_tr = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,project_grade_category_one_hot_train,price_standardized_train,quantity_standardized_train,teacher_number_of_previously_posted_projects_standardized_train,text_bow_train))
```

```
X_cr = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_number_of_previously_posted_projects_standardized_cv,project_grade_category_one_hot_cv,price_standardized_cv,quantity_standardized_cv,text_bow_cv,title_bow_cv))
```

```
X_te = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,project_grade_category_one_hot_test,price_standardized_test,quantity_standardized_test,teacher_number_of_previously_posted_projects_standardized_test,text_bow_test,title_bow_test))
```

```
print("Final Data matrix on BOW")
```

```
print(X_tr.shape, y_train.shape)
```

```
print(X_cr.shape, y_cv.shape)
```

```
print(X_te.shape, y_test.shape)
```

```
print("="*100)
```



Final Data matrix on BOW

```
(22445, 10209) (22445,)
```

```
(11055, 10209) (11055,)
```

```
(16500, 10209) (16500,)
```

```
=====
```

```
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
X_tr = scaler.fit_transform(X_tr,y_train)
```

```
X_cr = scaler.transform(X_cr)
```

```
X_te = scaler.transform(X_te)
```

```
print(X_tr.shape, y_train.shape)
```

```
print(X_cr.shape, y_cv.shape)
```

```
print(X_te.shape, y_test.shape)
```



```
(22445, 10209) (22445,)
```

```
(11055, 10209) (11055,)
```

```
(16500, 10209) (16500,)
```

```
from sklearn.naive_bayes import MultinomialNB
```

```
nb_bow = MultinomialNB(alpha = 0.5,class_prior=[0.5,0.5])
```

```
nb_bow.fit(X_tr, y_train)
```



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

```

# bow_features_probs1 = {}
# bow_features_probs0 = {}
# for a in range(10212) :
#     bow_features_probs1[a] = nb_bow.feature_log_prob_[1,a]
# bow_features_probs = {}
# for b in range(10212) :
#     bow_features_probs0[b] = nb_bow.feature_log_prob_[0,a]
# # print((bow_features_probs1.items[:3]))
# c1=0
# for k,v in bow_features_probs1.items():
#     print(k, v)
#     c1 = c1 + 1
#     if(c1==5):
#         break

# print("="*100)
# # print((bow_features_probs0))
# c2=0
# for k,v in bow_features_probs0.items():
#     print(k, v)
#     c2 = c2 + 1
#     if(c2==5):
#         break


```

```

bow_features_probs1 = []
for a in range(10209) :
    b = nb_bow.feature_log_prob_[1,a]
    bow_features_probs1.append(b)

```

```
len(bow_features_probs1)
```

 10209

```
bow_features_names = []
```

```

for a in vectorizer_cat.get_feature_names() :
    bow_features_names.append(a)

```

```

for a in vectorizer_sub_cat.get_feature_names() :
    bow_features_names.append(a)

```

```

for a in vectorizer_state.get_feature_names() :
    bow_features_names.append(a)

```

```

for a in vectorizer_teacherprefix.get_feature_names() :
    bow_features_names.append(a)

```

```

for a in vectorizer_projectgrade.get_feature_names() :
    bow_features_names.append(a)

bow_features_names.append("price")
bow_features_names.append("quantity")
bow_features_names.append("teacher_number_of_previously_posted")

for a in vectorizer_bow_essays.get_feature_names() :
    bow_features_names.append(a)

for a in vectorizer_bow_titles.get_feature_names() :
    bow_features_names.append(a)
len(bow_features_names)

```



10209

```

final_features_bow_df_pos = pd.DataFrame({'feature_prob_estimates' : bow_features_probs1, 'fea

final_features_bow_df_pos.sort_values(by = ['feature_prob_estimates'], ascending = False,inpl

final_features_bow_df_pos.head(10)

```



	feature_prob_estimates	feature_names
<b>92</b>	-4.215421	Mrs
<b>8</b>	-4.287486	Literacy_Language
<b>98</b>	-4.474116	GradesPreK-2
<b>7</b>	-4.574186	Math_Science
<b>93</b>	-4.601107	Ms
<b>95</b>	-4.653046	Grades3-5
<b>38</b>	-4.705384	Literacy
<b>7766</b>	-4.880016	students
<b>37</b>	-4.951733	Mathematics
<b>36</b>	-5.158836	Literature_Writing

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

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

```
bow_features_probs2 = []
```



```

for a in range(10209) :
    bb = nb_bow.feature_log_prob_[0,a]
    bow_features_probs2.append(bb)

# (bow_features_probs)

final_features_bow_df_neg = pd.DataFrame({'feature_prob_estimates' : bow_features_probs2,'fea

final_features_bow_df_neg.sort_values(by = ['feature_prob_estimates'], ascending = False,inpl

final_features_bow_df_neg.head(10)

```



	feature_prob_estimates	feature_names
92	-4.253758	Mrs
8	-4.402431	Literacy_Language
98	-4.458514	GradesPreK-2
7	-4.493592	Math_Science
93	-4.572391	Ms
95	-4.707219	Grades3-5
38	-4.874765	Literacy
37	-4.927322	Mathematics
7766	-4.941158	students
36	-5.226500	Literature_Writing

# Please write all the code with proper documentation

```

import dill
# dill.dump_session('notebook_env.db')
dill.load_session('../notebook_env.db')

```

## ▼ 2.4.2 Applying Naive Bayes on TFIDF, SET 2

```

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_tra
               ,project_grade_category_one_hot_train,price_standardized_train,quantity_standa
               ,teacher_number_of_previously_posted_projects_standardized_train,text_tfidf_tr
X_cr = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teache
               ,project_grade_category_one_hot_cv,price_standardized_cv,quantity_standardized
               ,teacher_number_of_previously_posted_projects_standardized_cv,text_tfidf_cv,ti

```

```
X_te = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,
               ,project_grade_category_one_hot_test,price_standardized_test,quantity_standard
               ,teacher_number_of_previously_posted_projects_standardized_test,text_tfidf_tes

print("Final Data matrix on TFIDF")
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
print("="*100)
```



Final Data matrix on TFIDF

(22445, 10209) (22445,)

(11055, 10209) (11055,)

(16500, 10209) (16500,)

=====

```
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_tr = scaler.fit_transform(X_tr,y_train)
X_cr = scaler.transform(X_cr)
X_te = scaler.transform(X_te)
print(X_tr.shape, y_train.shape)
print(X_cr.shape, y_cv.shape)
print(X_te.shape, y_test.shape)
```



(22445, 10209) (22445,)

(11055, 10209) (11055,)

(16500, 10209) (16500,)

## ➤ 1.1 Method 1: Simple for loop (if you are having memory limitations us

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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

import matplotlib.pyplot as plt
from sklearn.naive_bayes import MultinomialNB
```

```
from sklearn.metrics import roc_auc_score
import math
from sklearn.model_selection import RandomizedSearchCV

train_auc = []
cv_auc = []
log_alphas = []
alphas = [0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50,

for i in tqdm(alphas):
    nb = MultinomialNB(alpha = i,class_prior=[0.5,0.5])
    nb.fit(X_tr, y_train)
    y_train_pred = batch_predict(nb, X_tr)
    y_cv_pred = batch_predict(nb, X_cr)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the
    # 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))

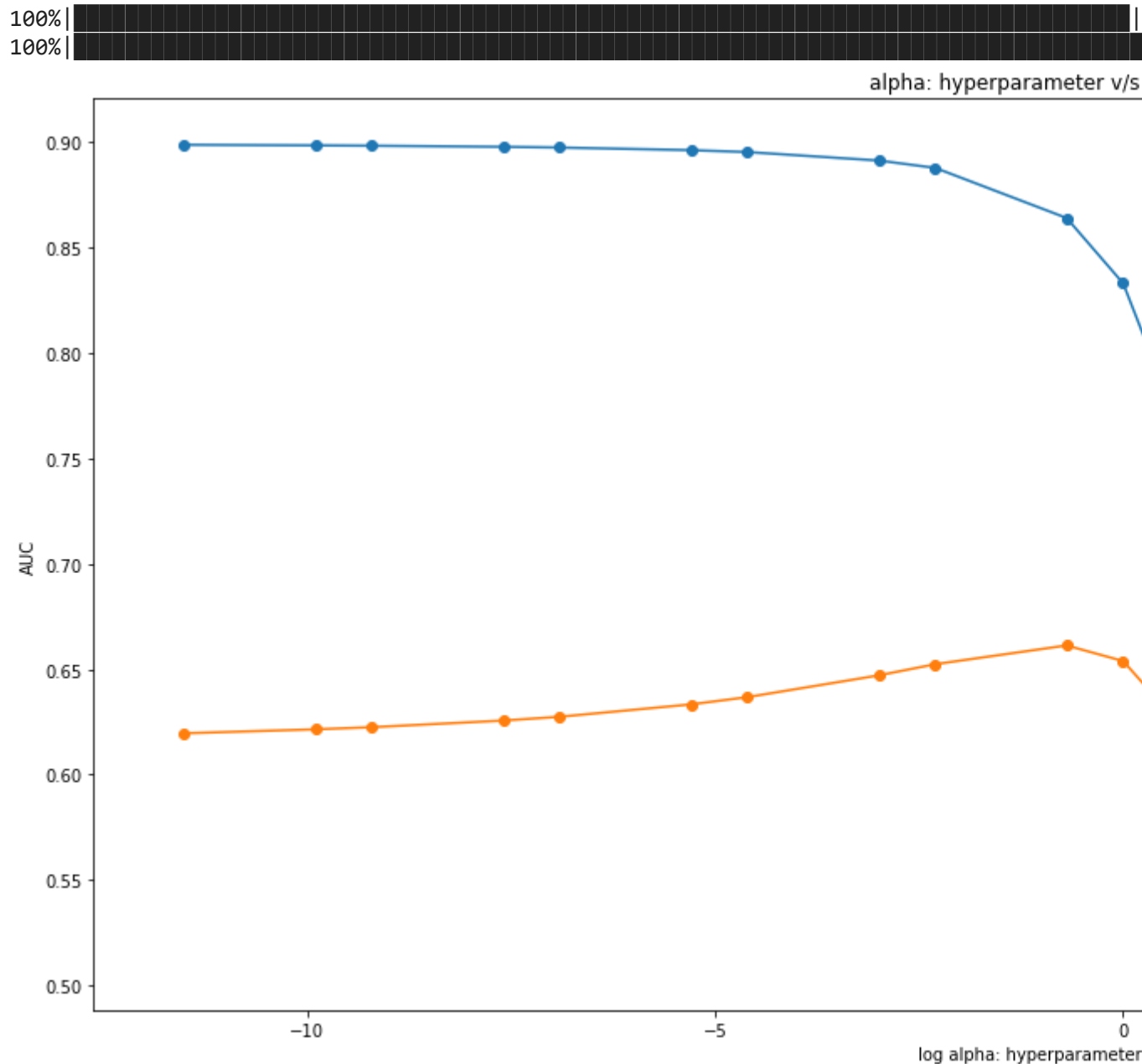
for a in tqdm(alphas):
    b = math.log(a)
    log_alphas.append(b)

log_alphas = np.array(log_alphas)
alphas = np.array(alphas)

plt.figure(figsize=(20,10))
plt.grid()
plt.plot(log_alphas, train_auc, label='Train AUC')
plt.plot(log_alphas, cv_auc, label='CV AUC')
plt.scatter(log_alphas, train_auc, label='Train AUC points')
plt.scatter(log_alphas, cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("log alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC")
plt.grid()
plt.show()

# print(len(log_alphas))
# print(log_alphas.shape)
# print(train_auc.shape)
```





# [https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc\\_curve.html#sklearn.metrics.roc\\_curve](https://scikitlearn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve)  
 from sklearn.metrics import roc\_curve, auc

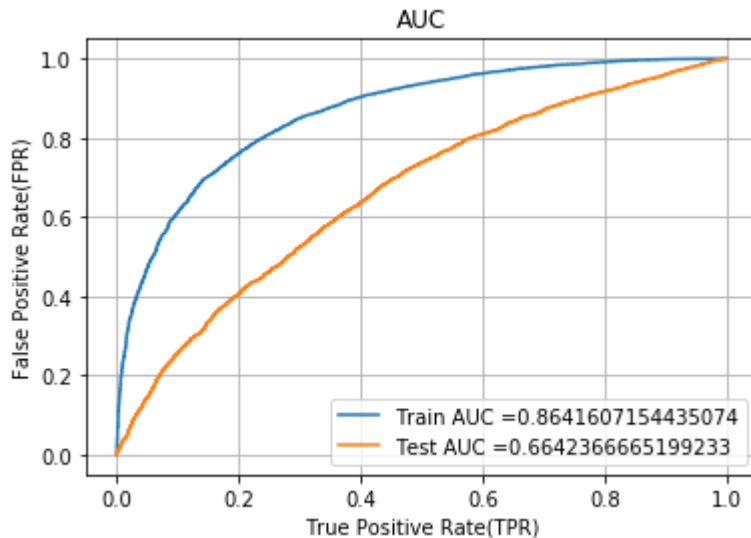
```
nb_tfidf = MultinomialNB(alpha = 0.5, class_prior=[0.5, 0.5])
nb_tfidf.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the pos
# not the predicted outputs
```

```
y_train_pred = batch_predict(nb_tfidf, X_tr)
y_test_pred = batch_predict(nb_tfidf, X_te)
```

```
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("AUC")
plt.grid()
plt.show()
```



```
# 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

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, test_fpr, test_tpr)))
```



```
=====
```

Train confusion matrix

the maximum value of  $\text{tpr} \cdot (1 - \text{fpr})$  0.6110093873451795 for threshold 0.485

```
[[ 2673   790]
 [ 3956 15026]]
```

Test confusion matrix

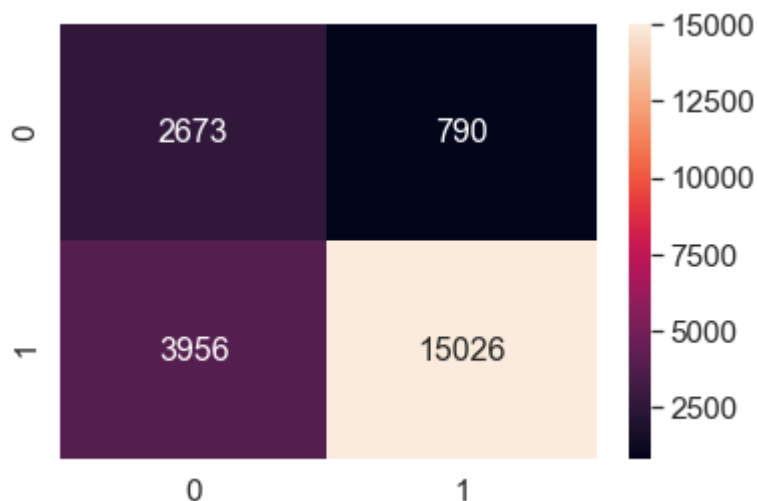
the maximum value of  $\text{tpr} \cdot (1 - \text{fpr})$  0.38422744871179804 for threshold 0.376

```
[[  984  1562]
 [ 2575 11379]]
```

```
conf_matr_df_train = pd.DataFrame(confusion_matrix(y_train[:], predict(y_train_pred, tr_thre
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train, annot=True,annot_kws={"size": 16}, fmt='g')
```



the maximum value of  $\text{tpr} \cdot (1 - \text{fpr})$  0.6110093873451795 for threshold 0.485  
<matplotlib.axes.\_subplots.AxesSubplot at 0x25de9172be0>



```
conf_matr_df_test = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresho
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test, annot=True,annot_kws={"size": 16}, fmt='g')
```



# Please write all the code with proper documentation

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

# merge two sparse matrices: <https://stackoverflow.com/a/19710648/4084039>

```
from scipy.sparse import hstack
```

```
X_tr = hstack((categories_one_hot_train,sub_categories_one_hot_train,school_state_one_hot_train,project_grade_category_one_hot_train,price_standardized_train,quantity_standardized_train,teacher_number_of_previously_posted_projects_standardized_train,text_tfidf_train))
```

```
X_cr = hstack((categories_one_hot_cv,sub_categories_one_hot_cv,school_state_one_hot_cv,teacher_number_of_previously_posted_projects_standardized_cv,text_tfidf_cv))
```

```
X_te = hstack((categories_one_hot_test,sub_categories_one_hot_test,school_state_one_hot_test,project_grade_category_one_hot_test,price_standardized_test,quantity_standardized_test,teacher_number_of_previously_posted_projects_standardized_test,text_tfidf_test))
```

```
print("Final Data matrix on TFIDF")
```

```
print(X_tr.shape, y_train.shape)
```

```
print(X_cr.shape, y_cv.shape)
```

```
print(X_te.shape, y_test.shape)
```

```
print("="*100)
```



Final Data matrix on TFIDF

```
(22445, 10209) (22445,)
```

```
(11055, 10209) (11055,)
```

```
(16500, 10209) (16500,)
```

```
=====
```

```
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
X_tr = scaler.fit_transform(X_tr,y_train)
```

```
X_cr = scaler.transform(X_cr)
```

```
X_te = scaler.transform(X_te)
```

```
print(X_tr.shape, y_train.shape)
```

```
print(X_cr.shape, y_cv.shape)
```

```
print(X_te.shape, y_test.shape)
```



```
(22445, 10209) (22445,)
```

```
(11055, 10209) (11055,)
```

```
(16500, 10209) (16500,)
```

```
from sklearn.naive_bayes import MultinomialNB
```

```
nb_tfidf = MultinomialNB(alpha = 0.5,class_prior=[0.5,0.5])
```


```
nb_tfidf.fit(X_tr, y_train)
```



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

```
tfidf_features_probs1 = []  
for a in range(10209) :  
    b = nb_tfidf.feature_log_prob_[1,a]  
    tfidf_features_probs1.append(b)
```

```
len(tfidf_features_probs1)
```

 10209

```
tfidf_features_names = []
```

```
for a in vectorizer_cat.get_feature_names() :  
    tfidf_features_names.append(a)
```

```
for a in vectorizer_sub_cat.get_feature_names() :  
    tfidf_features_names.append(a)
```

```
for a in vectorizer_state.get_feature_names() :  
    tfidf_features_names.append(a)
```


```
for a in vectorizer_teacherprefix.get_feature_names() :  
    tfidf_features_names.append(a)
```

```
for a in vectorizer_projectgrade.get_feature_names() :  
    tfidf_features_names.append(a)
```

```
tfidf_features_names.append("price")  
tfidf_features_names.append("quantity")  
tfidf_features_names.append("teacher_number_of_previously_posted")
```

```
for a in vectorizer_tfidf_essays.get_feature_names() :  
    tfidf_features_names.append(a)
```

```
for a in vectorizer_tfidf_title.get_feature_names() :  
    tfidf_features_names.append(a)  
len(tfidf_features_names)
```

 10209

```
final_features_tfidf_df_pos = pd.DataFrame({'feature_prob_estimates' : tfidf_features_probs1,
```

```
final_features_tfidf_df_pos.sort_values(by = ['feature_prob_estimates'], ascending = False,in
```

```
final_features_tfidf_df_pos.head(10)
```





	feature_prob_estimates	feature_names
<b>92</b>	-4.220127	Mrs
<b>8</b>	-4.292192	Literacy_Language
<b>98</b>	-4.478822	GradesPreK-2
<b>7766</b>	-4.573659	students
<b>7</b>	-4.578892	Math_Science
<b>93</b>	-4.605813	Ms
<b>95</b>	-4.657751	Grades3-5
<b>38</b>	-4.710090	Literacy
<b>37</b>	-4.956439	Mathematics
<b>36</b>	-5.163542	Literature_Writing

# Please write all the code with proper documentation

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

```
tfidf_features_probs2 = []
for a in range(10209) :
    bb = nb_tfidf.feature_log_prob_[0,a]
    tfidf_features_probs2.append(bb)
```

```
# (bow_features_probs)
```

```
final_features_tfidf_df_neg = pd.DataFrame({'feature_prob_estimates' : tfidf_features_probs2,
```

```
final_features_tfidf_df_neg.sort_values(by = ['feature_prob_estimates'], ascending = False,in
```

```
final_features_tfidf_df_neg.head(10)
```



	feature_prob_estimates	feature_names
92	-4.294053	Mrs
8	-4.442726	Literacy_Language
98	-4.498809	GradesPreK-2
7	-4.533887	Math_Science
93	-4.612686	Ms
7766	-4.616482	students
95	-4.747514	Grades3-5
38	-4.915060	Literacy
37	-4.967617	Mathematics
36	-5.266795	Literature_Writing

# Please write all the code with proper documentation

### 3. Conclusions

# 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", "AUC"]
x.add_row(["BOW", "Naive Bayes", 0.5, 0.65])
x.add_row(["TFIDF", "Naive Bayes", 0.5, 0.66])
print(x)
```



Vectorizer	Model	Hyper Parameter:Alpha	AUC
BOW	Naive Bayes	0.5	0.65
TFIDF	Naive Bayes	0.5	0.66

