

# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> <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 enumerated values: <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>
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <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> <b>Examples:</b> <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> WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Literacy</li></ul>

Feature	Description
<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 materials to manage sensory needs!</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> 2016-04-28 12:43:56.245
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> <li>nan</li> <li>Dr.</li> <li>Mr.</li> <li>Mrs.</li> <li>Ms.</li> <li>Teacher.</li> </ul>
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. <b>Example:</b> 2

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

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

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502
<code>description</code>	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. <b>Example:</b> 3
<code>price</code>	Price of the resource required. <b>Example:</b> 9.95

**Note:** Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project was not approved, and a value of 1 indicates the project was approved.

## Notes on the Essay Data

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

- `__project_essay_1__` "Introduce us to your classroom"
- `__project_essay_2__` "Tell us more about your students"
- `__project_essay_3__` "Describe how your students will use the materials you're requesting"
- `__project_essay_3__` "Close by sharing why your project will make a difference"

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

- `__project_essay_1__` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful"

your neighborhood, and your school are all helpful.

- \_\_project\_essay\_2\_\_: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with project\_submitted\_datetime of 2016-05-17 and later, the values of project\_essay\_3 and project\_essay\_4 will be NaN.

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import 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
print('done')

!pip install -U -q PyDrive
```

paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip install paramiko` to suppress

done

In [0]:

```
from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials
# Authenticate and create the PyDrive client.
auth.authenticate_user()
gauth = GoogleAuth()
gauth.credentials = GoogleCredentials.get_application_default()
drive = GoogleDrive(gauth)

# get links to drive to access data
link='https://drive.google.com/open?id=18VAiuw3vfETGcuJOdicvkgQT0pSxF7Wy'
link3='https://drive.google.com/open?id=1Z6bjXmyCaoEzXYo_tRDwLTsfeA2F3K3j'
fluff, id2 = link3.split('=')
#print (id2) # Verify that you have everything after '='
downloaded = drive.CreateFile({'id':id2})
```

```
downloaded.GetContentFile('glove_vectors')
```

## 1.1 Reading Data

In [3]:

```
#Project_data

fluff, id = link.split('=')
print (id) # Verify that you have everything after '='

downloaded = drive.CreateFile({'id':id})
downloaded.GetContentFile('train_data.csv')
project_data = pd.read_csv('train_data.csv',nrows=100000)

print(project_data.shape)
link1='https://drive.google.com/open?id=1luHEj9KOGWD9SU-CPgKyb6VrWqVos4uV'
```

```
18VAiuw3vfETGcuJOdicvkgQT0pSxF7Wy
(100000, 17)
```

In [4]:

```
#Resource_data

fluff1, idi = link1.split('=')
print (idi) # Verify that you have everything after '='
downloaded = drive.CreateFile({'id':idi})
downloaded.GetContentFile('resources .csv')
resource_data = pd.read_csv('resources .csv')

print(resource_data .head(3))
```

```
11uHEj9KOGWD9SU-CPgKyb6VrWqVos4uV
   id                description  quantity \
0  p233245  LC652 - Lakeshore Double-Space Mobile Drying Rack          1
1  p069063      Bouncy Bands for Desks (Blue support pipes)           3
2  p069063  Cory Stories: A Kid's Book About Living With Adhd          1

   price
0  149.00
1   14.95
2    8.45
```

In [5]:

```
print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)

print(resource_data.shape)
print(resource_data.columns.values)
```

```
Number of data points in train data (100000, 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']
(1541272, 4)
['id' 'description' 'quantity' 'price']
```

In [6]:

```
#Sort the datapoints by date <-
```

```
# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
```

```
# how to replace elements in list python: https://stackoverflow.com/a/2002103/5004039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]

#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
project_data.drop('project_submitted_datetime', axis=1, inplace=True)# we drop the col
project_data.sort_values(by=['Date'], inplace=True)# sort the values y date

# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
project_data = project_data[cols]

project_data.head(2)
```

Out[6]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_cate
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5

### 1.3 Text preprocessing

In [8]:

```
# merge two column text dataframe:
project_data["essay"] = project_data["project_essay_1"].map(str) + \
    project_data["project_essay_2"].map(str) + \
    project_data["project_essay_3"].map(str) + \
    project_data["project_essay_4"].map(str)

project_data.head(2)
```

Out[8]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_cate
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5

In [0]:

```
# https://stackoverflow.com/a/47091490/4084039
import re
```

```
def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\ 're", " are", phrase)
    phrase = re.sub(r"\ 's", " is", phrase)
    phrase = re.sub(r"\ 'd", " would", phrase)
    phrase = re.sub(r"\ 'll", " will", phrase)
    phrase = re.sub(r"\ 't", " not", phrase)
    phrase = re.sub(r"\ 've", " have", phrase)
    phrase = re.sub(r"\ 'm", " am", phrase)
    return phrase
```

In [0]:

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

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

**Note:** Firstly i clean the project\_subject\_categories and project\_subject\_subcategories or preprocess it. then after cleaing i convert to train,test and cv.

## Preprocessing of project\_subject\_categories

In [0]:

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

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

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

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

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

## Preprocessing of project\_subject\_subcategories

In [0]:

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

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

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

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
mv_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]))

```

## Preprocessing of project\_grade\_category

In [13]:

```
print(project_data['project_grade_category'][:20])# we have to remove the grades from every row
```

```

55660    Grades PreK-2
76127    Grades 3-5
51140    Grades PreK-2
473      Grades PreK-2
41558    Grades 3-5
29891    Grades 3-5
81565    Grades 3-5
79026    Grades 3-5
23374    Grades PreK-2
86551    Grades 3-5
49228    Grades PreK-2
72638    Grades 9-12
7176     Grades PreK-2
70898    Grades 3-5
72593    Grades PreK-2
35006    Grades 3-5
5145     Grades 3-5
48237    Grades 9-12
64637    Grades PreK-2
98973    Grades 3-5
Name: project_grade_category, dtype: object

```

In [0]:

```

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

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

grade_cat_list = []
for i in d:
    # consider we have text like this:
    for j in i.split(' '): # # split by space
        j=j.replace('Grades','')# clean grades from the row
    grade_cat_list.append(j.strip())

project_data['clean_grade'] = grade_cat_list
project_data.drop(['project_grade_category'], 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_grade'].values:
    my_counter.update(word.split())

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

```

## Assignment 4: Apply Naive Bayes

1. [Task-1] Apply Naive Bayes(MultinomialNB) on these feature sets



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

## 2. Hyper paramter tuning to find best Alpha

- Find the best hyper parameter which results in the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation (or) simple cross validation data
- Use gridsearch-cv or randomsearch-cv or write your own for loops to do this task

## 3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, as shown in the figure
- Once you find the best hyper parameter, you need to train your model-M using the best hyper-param. Now, find the AUC on test data and plot the ROC curve on both train and test using model-M.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points

## 4. [Task-2]

- Select the top 20 features from set one and set two by using the absolute values of the coeff\_paramter of the multinomial.
- Repeat the steps 2 and 3 on the data matrix after feature selection

## 5. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library [link](#)

### Note: Data Leakage

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakage, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method fit\_transform() on you train data, and apply the method transform() on cv/test data.
4. For more details please go through this [link](#).

## 2. Preparing our data for the models

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

In [0]:

```
#Splitting Data into train and Test sklearn https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(project_data,
                                                    project_data['project_is_approved'],
                                                    stratify= project_data['project_is_approved'],
                                                    test_size = 0.33
                                                    )
```

In [0]:

```
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, stratify= y_train,
                                                test_size = 0.33)
```

In [18]:

```
print(y_train.value_counts())
print(y_test.value_counts())
```

```
print(y_test.value_counts())
print(y_cv.value_counts())
# huge imbalance
```

```
1    38087
0     6803
Name: project_is_approved, dtype: int64
1    27999
0     5001
Name: project_is_approved, dtype: int64
1    18760
0     3350
Name: project_is_approved, dtype: int64
```

In [0]:

```
#dropping the y labels
#https://stackoverflow.com/questions/13411544/delete-column-from-pandas-dataframe-by-column-name
#x_train =
X_train.drop(["project_is_approved"], axis = 1, inplace = True)
#x_test =
X_test.drop(["project_is_approved"], axis = 1, inplace = True)
#x_cv =
X_cv.drop(["project_is_approved"], axis = 1, inplace = True)
#print(X_train)
```

## Text preprocessing of train,test and cv

In [20]:

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

```
100%|██████████| 44890/44890 [00:50<00:00, 884.68it/s]
```

In [21]:

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

```
100%|██████████| 33000/33000 [00:36<00:00, 901.47it/s]
```

In [22]:

```
#Proprocessing for essay
# Combining all the above stundents
```

```

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

```

100%|██████████| 22110/22110 [00:24<00:00, 892.40it/s]

In [23]:

```

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

```

100%|██████████| 22110/22110 [00:01<00:00, 21314.25it/s]

In [24]:

```

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

```

100%|██████████| 44890/44890 [00:02<00:00, 21182.28it/s]

In [25]:

```

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

```

100%|██████████| 33000/33000 [00:01<00:00, 21210.84it/s]

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

### 1. vectorize categorical data

1.project\_subject\_categories convert categorical to vectors\*

In [26]:

```
# convert train,cv and test data of clean_categories into vectors

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer1 = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False,
binary=True)
vectorizer1.fit(X_train['clean_categories'].values)

# firstly convert fit the train data into the vectoriaer then it learn hte vocablery

# we use the fitted CountVectorizer to convert the text to vector
X_train_cat = vectorizer1.transform(X_train['clean_categories'].values)
X_cv_cat = vectorizer1.transform(X_cv['clean_categories'].values)
X_test_cat = vectorizer1.transform(X_test['clean_categories'].values)

print(vectorizer1.get_feature_names())

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
```

In [27]:

```
print("After vectorizations")
print(X_train_cat.shape, y_train.shape)
print(X_cv_cat.shape, y_cv.shape)
print(X_test_cat.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(44890, 9) (44890,)
(22110, 9) (22110,)
(33000, 9) (33000,)
```



2.project\_subject\_subcategories convert categorical to vectors\*

In [28]:

```
# convert train,cv and test data of clean_categories into vectors

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer2 = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary
=True)
vectorizer2.fit(X_train['clean_subcategories'].values)

# firstly convert fit the train data into the vectoriaer then it learn hte vocablery

# we use the fitted CountVectorizer to convert the text to vector
X_train_subcat = vectorizer2.transform(X_train['clean_subcategories'].values)
X_cv_subcat = vectorizer2.transform(X_cv['clean_subcategories'].values)
X_test_subcat = vectorizer2.transform(X_test['clean_subcategories'].values)
```

```
x_train_subcat = vectorizer2.transform(x_train['clean_subcategory'].values)

print(vectorizer2.get_feature_names())
```

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

In [29]:

```
print("After vectorizations")
print(X_train_subcat.shape, y_train.shape)
print(X_cv_subcat.shape, y_cv.shape)
print(X_test_subcat.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(44890, 30) (44890,)
(22110, 30) (22110,)
(33000, 30) (33000,)
```

**\*3 school\_state convert categorical to vectors\*\***

In [31]:

```
#first convert to dict.
# now time to cont the each words
from collections import Counter
my_counter = Counter()
for word in project_data['school_state'].values:
    my_counter.update(word.split()) # count the words

school_state_dict = dict(my_counter) # store in dictionary
sorted_school_state_dict = dict(sorted(school_state_dict.items(), key=lambda kv: kv[1])) # sor it
print(sorted_school_state_dict)

# convert train,cv and test data of clean_categories into vectors

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer3 = CountVectorizer(vocabulary=list(sorted_school_state_dict.keys()), lowercase=False, b
inary=True)
vectorizer3.fit(project_data['school_state'].values)

# firstly convert fit the train data into the vectoriaer then it learn hte vocablery

# we use the fitted CountVectorizer to convert the text to vector
X_train_school_state = vectorizer3.transform(X_train['school_state'].values)
X_cv_school_state = vectorizer3.transform(X_cv['school_state'].values)
X_test_school_state = vectorizer3.transform(X_test['school_state'].values)

print(vectorizer3.get_feature_names())
```

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

In [32]:

```
print("After vectorizations")
print(X_train_school_state .shape, y_train.shape)
print(X_cv_school_state .shape, y_cv.shape)
print(X_test_school_state .shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(44890, 51) (44890,)
(22110, 51) (22110,)
(33000, 51) (33000,)
```

#### \*4. project\_grade\_category categorical to vectors\*

In [35]:

```
# convert train,cv and test data of clean_categories into vectors

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer4 = CountVectorizer(vocabulary=list(sorted_project_grade_category_dict.keys()),
lowercase=False, binary=True)
vectorizer4.fit(project_data['clean_grade'].values)

# firstly convert fit the train data into the vectoriaer then it learn hte vocablery

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_grade_category = vectorizer4.transform(X_train['clean_grade'].values)
X_cv_project_grade_category = vectorizer4.transform(X_cv['clean_grade'].values)
X_test_project_grade_category = vectorizer4.transform(X_test['clean_grade'].values)

print(vectorizer4.get_feature_names())
```

```
['9-12', '6-8', '3-5', 'PreK-2']
```

In [36]:

```
print("After vectorizations")
print(X_train_project_grade_category .shape, y_train.shape)
print(X_cv_project_grade_category .shape, y_cv.shape)
print(X_test_project_grade_category .shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(44890, 4) (44890,)
(22110, 4) (22110,)
(33000, 4) (33000,)
```

#### 5. teacher\_prefix categorical to vectors\*

In [37]:

```
#https://stackoverflow.com/questions/42224700/attributeerror-float-object-has-no-attribute-split
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna(" ")# filll the null values
with space
X_train['teacher_prefix'][:3]# dots is the problme for us
```

Out[37]:

```
22903    Ms.
38384    Mr.
40869    Ms.
Name: teacher_prefix, dtype: object
```

```
name: teacher_prefix, dtype: object
```

In [38]:

```
my_counter = Counter()
for word in project_data['teacher_prefix'].values:
    my_counter.update(word.split())

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

# convert train,cv and test data of clean_categories into vectors

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer5 = CountVectorizer(vocabulary=list(sorted_teacher_prefix_dict.keys()), lowercase=False,
binary=True)
vectorizer5.fit(project_data['teacher_prefix'].values.astype('U'))

# firstly convert fit the train data into the vectoriaer then it learn hte vocablery

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_prefix = vectorizer5.transform(X_train['teacher_prefix'].values.astype('U'))
X_cv_teacher_prefix= vectorizer5.transform(X_cv['teacher_prefix'].values.astype('U'))
X_test_teacher_prefix = vectorizer5.transform(X_test['teacher_prefix'].values.astype('U'))

print(vectorizer5.get_feature_names())

# when i executeed this error comes
#np.nan is an invalid document, expected byte or unicode string.
# then iconvert to unicode      just writ .astype('U') after the .values in fit and trainform
#https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is
-an-invalid-document

['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
```

In [39]:

```
print("After vectorizations")
print(X_train_teacher_prefix .shape, y_train.shape)
print(X_cv_teacher_prefix .shape, y_cv.shape)
print(X_test_teacher_prefix .shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(44890, 5) (44890,)
(22110, 5) (22110,)
(33000, 5) (33000,)
```



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

-> preprocess essay and title for gts ready for apply featureization

Apply Baw featureization essay

In [40]:

```

X_train_essay=preprocessed_essays_train
X_cv_essay=preprocessed_essays_cv
X_test_essay=preprocessed_essays_test

X_train_title=preprocessed_titles_train
X_cv_title=preprocessed_titles_cv
X_test_title=preprocessed_titles_test

# We are considering only the words which appeared in at least 10 documents (rows or projects).
vectorizer6 = CountVectorizer(min_df=10) # its a countvectors used for convert text to vectors
vectorizer6.fit(X_train_essay) # that is learned from trained data

# we use the fitted CountVectorizer to convert the text to vector
X_train_bow = vectorizer6.transform(X_train_essay)
X_cv_bow = vectorizer6.transform(X_cv_essay)
X_test_bow = vectorizer6.transform(X_test_essay)

print("After vectorizations")
print(X_train_bow.shape, y_train.shape)
print(X_cv_bow.shape, y_cv.shape)
print(X_test_bow.shape, y_test.shape)
print("=="*100)
# so the dimension of all are the same by using first fit and then transform

```

```

After vectorizations
(44890, 11636) (44890,)
(22110, 11636) (22110,)
(33000, 11636) (33000,)
=====

```



### Apply Bow featureization Title

In [41]:

```

vectorizer7 = CountVectorizer(min_df=10)
vectorizer7.fit(X_train_title) # that is learned from trained data

# we use the fitted CountVectorizer to convert the text to vector
X_train_bow_title = vectorizer7.transform(X_train_title)
X_cv_bow_title = vectorizer7.transform(X_cv_title)
X_test_bow_title = vectorizer7.transform(X_test_title)

print("After vectorizations")
print(X_train_bow_title.shape, y_train.shape)
print(X_cv_bow_title.shape, y_cv.shape)
print(X_test_bow_title.shape, y_test.shape)
print("=="*100)
# so the dimension of all are the same by using first fit and then transform

```

```

After vectorizations
(44890, 1888) (44890,)
(22110, 1888) (22110,)
(33000, 1888) (33000,)
=====

```



### Apply tf-idf featureization titles



In [42]:

```
#for titles
from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer8 = TfidfVectorizer(min_df=10)# its a countvectors used for convert text to vectors
vectorizer8.fit(X_train_title)# that is learned from trained data

# we use the fitted CountVectorizer to convert the text to vector
X_train_tf_title = vectorizer8.transform(X_train_title)
X_cv_tf_title= vectorizer8.transform(X_cv_title)
X_test_tf_title = vectorizer8.transform(X_test_title)

print("After vectorizations")
print(X_train_tf_title.shape, y_train.shape)
print(X_cv_tf_title.shape, y_cv.shape)
print(X_test_tf_title.shape, y_test.shape)
print("="*100)
# so the dimension of all are the same by using first fit and then transform
```

```
After vectorizations
(44890, 1888) (44890,)
(22110, 1888) (22110,)
(33000, 1888) (33000,)
```

### Apply tf-idf featureization Essays

In [43]:

```
#for essay
from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer9 = TfidfVectorizer(min_df=10)# its a countvectors used for convert text to vectors
vectorizer9.fit(X_train_essay)# that is learned from trained data

# we use the fitted CountVectorizer to convert the text to vector
X_train_tf_essay = vectorizer9.transform(X_train_essay)
X_cv_tf_essay= vectorizer9.transform(X_cv_essay)
X_test_tf_essay = vectorizer9.transform(X_test_essay)

print("After vectorizations")
print(X_train_tf_essay.shape, y_train.shape)
print(X_cv_tf_essay.shape, y_cv.shape)
print(X_test_tf_essay.shape, y_test.shape)
print("="*100)
# so the dimension of all are the same by using first fit and then transform
```

```
After vectorizations
(44890, 11636) (44890,)
(22110, 11636) (22110,)
(33000, 11636) (33000,)
```

## 1.5.3 Vectorizing Numerical features¶

In [44]:

```

price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
print(price_data.head(2))

# we also have to do this in train, test and cv
# so also merge the resource data with the train, cv and test

X_train = pd.merge(X_train, price_data, on = "id", how = "left")
#print(x_train.columns)
X_test = pd.merge(X_test, price_data, on = "id", how = "left")
X_cv = pd.merge(X_cv, price_data, on = "id", how = "left")

```

```

      id  price  quantity
0  p000001  459.56         7
1  p000002  515.89        21

```

### Standardized price for the train, test and cv

In [48]:

```

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

price_scalar = MinMaxScaler()

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

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

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

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

```

Out[48]:

```

array([[0.01842206],
       [0.01966626],
       [0.03824835],
       ...,
       [0.04143488],
       [0.02000132],
       [0.01980629]])

```

### Shapes"

In [49]:

```

print(train_price_standar.shape, y_train.shape)
print(test_price_standar.shape, y_test.shape)
print(cv_price_standar.shape, y_cv.shape)

```

```

(44890, 1) (44890,)
(33000, 1) (33000,)
(22110, 1) (22110,)

```

## Standardized Previous\_year\_teacher\_projects train,test and cv

In [50]:

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

# Now standardize the data with above mean and variance.
train_prev_proj_standar =
price_scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
train_prev_proj_standar

# Now standardize the data with above mean and variance.
test_prev_proj_standar =
price_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
test_prev_proj_standar

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

Out[50]:

```
array([[0.00461894],
       [0.         ],
       [0.00692841],
       ...,
       [0.06928406],
       [0.00461894],
       [0.01847575]])
```

## Shapes of all

In [51]:

```
print(train_prev_proj_standar.shape, y_train.shape)
print(test_prev_proj_standar.shape, y_test.shape)
print(cv_prev_proj_standar.shape, y_cv.shape)
```

```
(44890, 1) (44890,)
(33000, 1) (33000,)
(22110, 1) (22110,)
```

## Standardized the Quantity column of the train,test and cv

In [52]:

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

# Now standardize the data with above mean and variance.
train_qnty_standar = price_scalar.transform(X_train['quantity'].values.reshape(-1, 1))
train_qnty_standar

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

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

Out[52]:

```
array([[0.00861141],
       [0.0452099 ],
       [0.          ],
       ...,
       [0.00968784],
       [0.01399354],
       [0.03336921]])
```

## Shapes

In [53]:

```
print(train_qnty_standar.shape, y_train.shape)
print(test_qnty_standar.shape, y_test.shape)
print(cv_qnty_standar.shape, y_cv.shape)
```

```
(44890, 1) (44890,)
(33000, 1) (33000,)
(22110, 1) (22110,)
```

## Merge all features which we clean till now\*\*

### Prepare for set 1:

In [59]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set1_train = hstack((X_train_bow_title, X_train_bow,
                       X_train_teacher_prefix, X_train_cat, X_train_subcat,
                       X_train_project_grade_category, X_train_school_state))

print(X_set1_train.shape, y_train.shape)
```

```
(44890, 13623) (44890,)
```

In [60]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set1_cv = hstack((X_cv_bow_title, X_cv_bow,
                    X_cv_teacher_prefix, X_cv_cat, X_cv_subcat,
                    X_cv_project_grade_category, X_cv_school_state))

print(X_set1_cv.shape, y_cv.shape)
```

```
(22110, 13623) (22110,)
```

In [61]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set1_test = hstack((X_test_bow_title, X_test_bow,
                      X_test_teacher_prefix, X_test_cat, X_test_subcat,
                      X_test_project_grade_category, X_test_school_state))
```

```
print(X_set1_test.shape, y_test.shape)
```

```
(33000, 13623) (33000,)
```

### Prepare for set 2:

In [62]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set2_train = hstack((X_train_tf_essay,X_train_tf_title,
                       X_train_teacher_prefix,X_train_cat,X_train_subcat,
                       X_train_project_grade_category,X_train_school_state))

print(X_set2_train.shape, y_train.shape)
```

```
(44890, 13623) (44890,)
```

In [63]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set2_cv = hstack((X_cv_tf_essay,X_cv_tf_title,
                   X_cv_teacher_prefix,X_cv_cat,X_cv_subcat,
                   X_cv_project_grade_category,X_cv_school_state))

print(X_set2_cv.shape, y_cv.shape)
```

```
(22110, 13623) (22110,)
```

In [64]:

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_set2_test = hstack((X_test_tf_essay,X_test_tf_title,
                     X_test_teacher_prefix,X_test_cat,X_test_subcat,
                     X_test_project_grade_category,X_test_school_state))

print(X_set2_test.shape, y_test.shape)
```

```
(33000, 13623) (33000,)
```

## Applying Naive bayes section

### 2.4.1 Applying Naive Bayes(MultinomialNB) on BOW, SET 1

In [65]:

```
#http://localhost:8888/notebooks/Assignment_SAMPLE_SOLUTION%20(1).ipynb (for reference) Which you
provided

#
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score
from sklearn.naive_bayes import MultinomialNB
import matplotlib.pyplot as plt

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.
```

THE SAME, BUT FOR BINARY CLASS PROBLEMS:

`y_score` : array, shape = `[n_samples]` or `[n_samples, n_classes]`  
Target scores, can either be probability estimates of the positive class, confidence values, or no `n`-thresholded measure of decisions (as returned by "decision\_function" on some classifiers).  
For binary `y_true`, `y_score` is supposed to be the score of the class with greater label.

```
"""

train_auc = []
cv_auc = []
alpha =[0.0001,0.001,0.01,0.1,1,10,100,1000]

for i in tqdm(alpha):

    neigh = MultinomialNB(alpha=i)# takes the alpha from the i th list value
    neigh.fit(X_set1_train, y_train)# fit the model

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    y_train_pred = neigh.predict_proba(X_set1_train)[:,-1]#Return probability estimates for the set1x ,for the class label 1 or +ve.
    y_cv_pred = neigh.predict_proba(X_set1_cv)[:,-1]#Return probability estimates for the setcvx,for the class label 1 or +ve .

    # roc curve
    #Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores.
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(alpha, train_auc, label='Train AUC')
plt.plot(alpha, cv_auc, label='CV AUC')
plt.xscale('log')# we take the log in the x axis

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

plt.legend()
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

100%|██████████| 8/8 [00:02<00:00, 3.27it/s]



In [66]:

```

score_t_cv = [x for x in cv_auc]
opt_t_cv = alpha[score_t_cv.index(max(score_t_cv))]
print("Maximum AUC score of cv is:" + ' ' + str(max(score_t_cv)))
print("Corresponding alpha value of cv is:",opt_t_cv, '\n')
best_alp=opt_t_cv
print(best_alp)

```

Maximum AUC score of cv is: 0.7024476577666041  
Corresponding alpha value of cv is: 1

1

## Fitting Model to Hyper-Parameter Curve

In [68]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

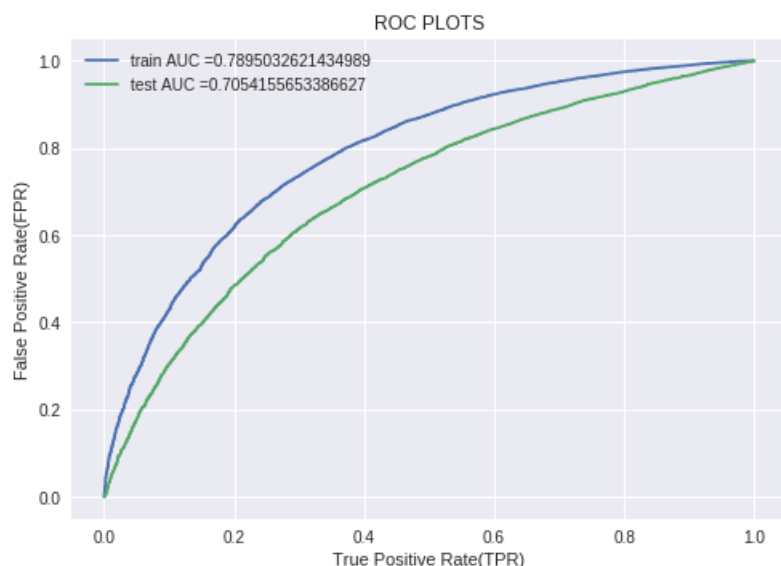
neigh = MultinomialNB(alpha=1)
neigh.fit(X_set1_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set1_train)[: ,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set1_test)[: ,1])

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

print("="*100)

```



**OBSERVATIONS:** As we seen form the roc plot ,Model work good on the train data , also model works good on the test data, only a little bit overfitting

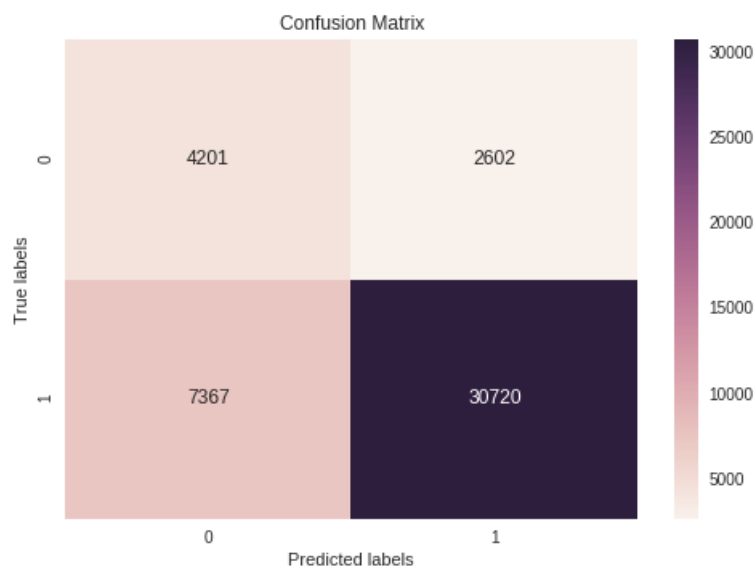
**Confusion matrix :**

In [69]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set1_train )), annot=True, ax = ax,fmt='g');
#annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```

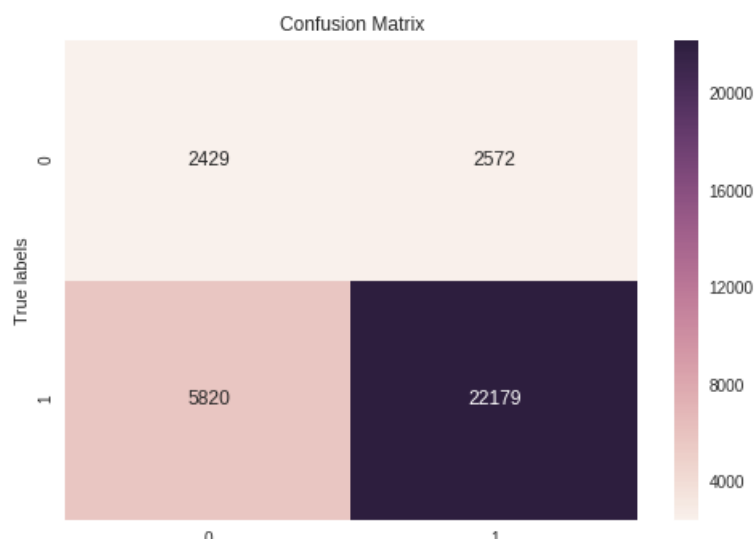


In [70]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set1_test )), annot=True, ax = ax,fmt='g'); #a
nnot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```





## 2.4.2 Applying KNN brute force on TFIDF, SET 2

In [98]:

```
#http://localhost:8888/notebooks/Assignment_SAMPLE_SOLUTION%20(1).ipynb (for reference) Which you
provided

from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import roc_auc_score

import matplotlib.pyplot as plt

"""
y_true : array, shape = [n_samples] or [n_samples, n_classes]
True binary labels or binary label indicators.

y_score : array, shape = [n_samples] or [n_samples, n_classes]
Target scores, can either be probability estimates of the positive class, confidence values, or no
n-thresholded measure of
decisions (as returned by "decision_function" on some classifiers).
For binary y_true, y_score is supposed to be the score of the class with greater label.
"""

train_auc = []
cv_auc = []
alpha = [0.0001,0.001,0.01,0.1,1,10,100,1000]
for i in tqdm(alpha):

    neigh = MultinomialNB(alpha=i)# takes the k from the i th list value
    neigh.fit(X_set2_train, y_train)# fit the model

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    tive class
    # not the predicted outputs
    y_train_pred = neigh.predict_proba(X_set2_train)[:,-1]#Return probability estimates for the
    set1x ,for the class label 1 or +ve.
    y_cv_pred = neigh.predict_proba(X_set2_cv)[:,-1]#Return probability estimates for the
    setcvx,for the class label 1 or +ve .

    # roc curve
    #Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scor
    es.
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

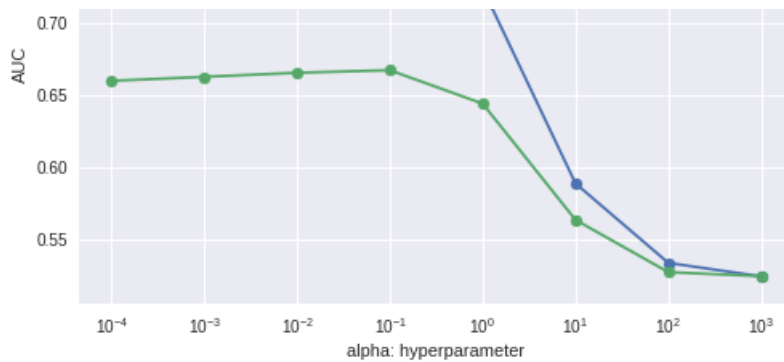
plt.plot(alpha, train_auc, label='Train AUC')
plt.plot(alpha, cv_auc, label='CV AUC')

plt.scatter(alpha, train_auc, label='Train AUC points')
plt.scatter(alpha, cv_auc, label='CV AUC points')
plt.xscale('log')
plt.legend()
plt.xlabel("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

100%|██████████| 8/8 [00:06<00:00, 1.32it/s]

ERROR PLOTS





### Fitting Model to Hyper-Parameter Curve:

In [100]:

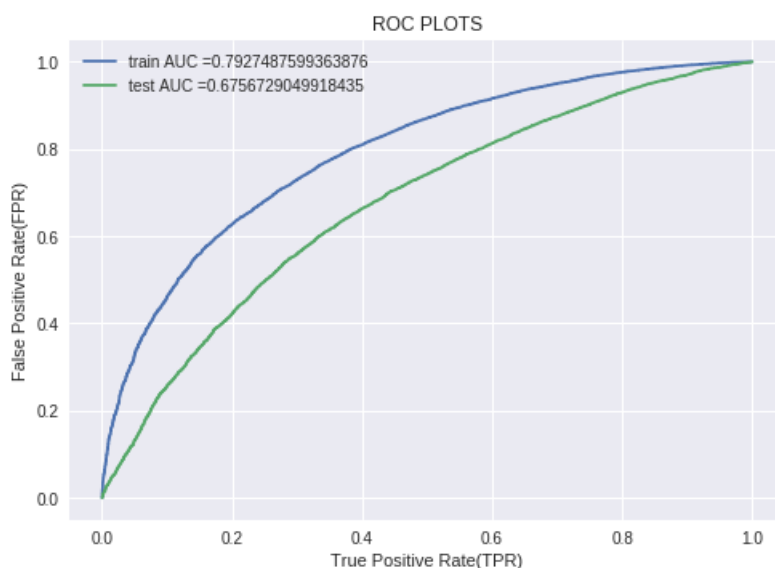
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

neigh = MultinomialNB(alpha=0.1)
neigh.fit(X_set2_train, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

train_fpr, train_tpr, thresholds = roc_curve(y_train, neigh.predict_proba(X_set2_train)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(y_test, neigh.predict_proba(X_set2_test)[:,1])

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

print("="*100)
```



**OBSERVATIONS:** As we seen form the roc plot ,, only a little bit overfitting, but roc curve are not so good only 65 score.

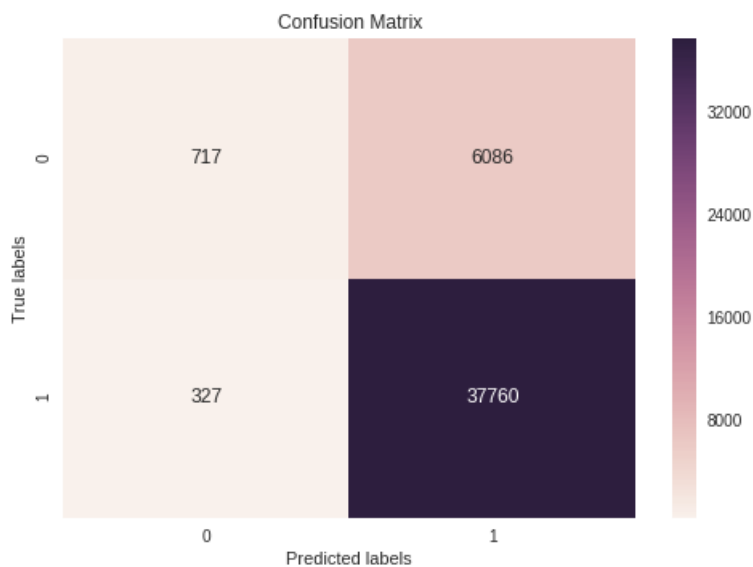
**Confusion matrix**

In [101]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, neigh.predict(X_set2_train )), annot=True, ax = ax,fmt='g');
#annot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```

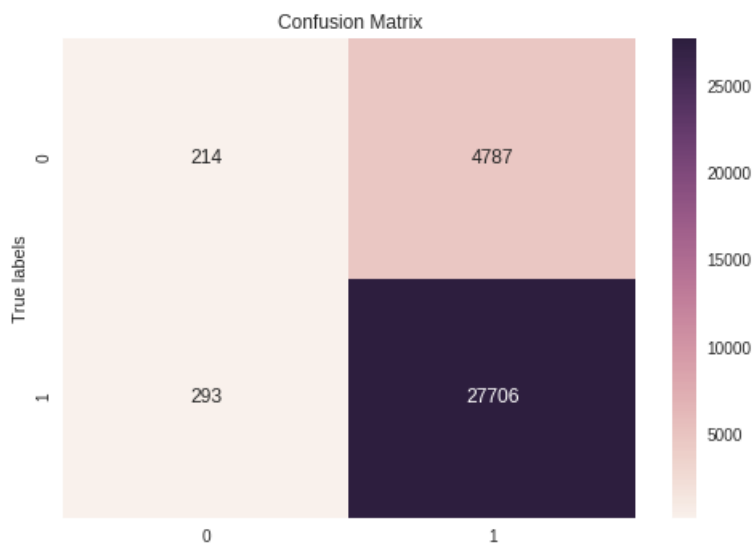


In [102]:

```
#https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
import seaborn as sns
import matplotlib.pyplot as plt

ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, neigh.predict(X_set2_test )), annot=True, ax = ax,fmt='g'); #a
nnot=True to annotate cells

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('True labels');
ax.set_title('Confusion Matrix');
#ax.xaxis.set_ticklabels(['business', 'health']); ax.yaxis.set_ticklabels(['health', 'business']);
```



**OBSERVATIONS:** As we see from this confusion matrix, True negatives are very less in this case because also in the original data it is very less, so bcz of this imbalance this work not good, dominating the negatives

## Top 10 features (negatives and positives)

For BOW

In [103]:

```
nb = MultinomialNB(alpha=0.1) # takes the k from the i th list value
nb.fit(X_set1_train, y_train) # fit the model

# now make a dictionary of all the probabilities for the weights
bow_features_probs = []
for a in range(13623): # loop till the (X_set1_train.shape)
    bow_features_probs.append(nb.feature_log_prob_[0,a] ) # negative feature probabilities

print(len(bow_features_probs))

bow_features_names = []
for a in vectorizer1.get_feature_names(): # clean categories
    bow_features_names.append(a)

for a in vectorizer2.get_feature_names(): # sub categories
    bow_features_names.append(a)
for a in vectorizer3.get_feature_names(): # school state
    bow_features_names.append(a)
for a in vectorizer4.get_feature_names(): # grade categories
    bow_features_names.append(a)
for a in vectorizer5.get_feature_names(): # teacher prefix
    bow_features_names.append(a)
for a in vectorizer6.get_feature_names(): # titles bow
    bow_features_names.append(a)
for a in vectorizer7.get_feature_names(): # essays bow
    bow_features_names.append(a)
print( len(bow_features_names))
```

13623  
13623

In [104]:

```
#top 10 negatives

final_bow_features = pd.DataFrame({'feature_prob_estimates' : bow_features_probs, 'feature_names'
: bow_features_names})
a = final_bow_features.sort_values(by = ['feature_prob_estimates'], ascending = False)
#print(final_bow_features.head(6))
a.head(10)
```

Out[104]:

	feature_prob_estimates	feature_names
11931	-2.988128	breaking
10992	-4.088530	unaware
7898	-4.403996	ponder
3830	-4.561124	ergonomic
7894	-4.755719	polynesian
6000	-4.756070	..

8899	-4.758973	reusable
6839	-4.776171	morning
8736	-4.961146	renovated
8253	-4.989949	protectors
8783	-5.104457	reputable

In [105]:

```
#top 10 Positives

# now make a dictionary of all the probabilityies fo the weights
bow_features_probs_pos = []
for a in range(13623):
    bow_features_probs_pos.append(nb.feature_log_prob_[1,a] )# negative feature probabilities

#len(bow_features_probs)
final_bow_features = pd.DataFrame({'feature_prob_estimates_pos' : bow_features_probs_pos,
'feature_names' : bow_features_names})

a =final_bow_features.sort_values(by = ['feature_prob_estimates_pos'], ascending = False)
#print(final_bow_features.head(6))
a.head(10)
```

Out[105]:

	feature_prob_estimates_pos	feature_names
11931	-2.980569	breaking
10992	-4.122557	unaware
7898	-4.485766	ponder
3830	-4.512538	ergonomic
8899	-4.787055	reusable
7894	-4.833971	polynesian
6839	-4.863757	morning
8253	-4.994075	protectors
8736	-5.017349	renovated
8783	-5.122240	reputable

For set 2 tf\_idf

In [106]:

```
nb = MultinomialNB(alpha=1)# takes the k from the i th list value
nb.fit(X_set2_train, y_train)# fit the model

# now make a dictionary of all the probabilityies fo the weights
tf_features_probs = []
for a in range(13623):# loop till (shape of data)
    tf_features_probs.append(nb.feature_log_prob_[0,a] )# negative feature probabilities

#len(bow_features_probs)
tf_features_names = []
for a in vectorizer1.get_feature_names() :# clean categories
    tf_features_names.append(a)

for a in vectorizer2.get_feature_names() :# sub categories
```

```

for a in vectorizer2.get_feature_names() :# sub categories
    tf_features_names.append(a)
for a in vectorizer3.get_feature_names() :#school state
    tf_features_names.append(a)
for a in vectorizer4.get_feature_names() :# grade categories
    tf_features_names.append(a)
for a in vectorizer5.get_feature_names() :# teacher prefix
    tf_features_names.append(a)
len(tf_features_names)

for a in vectorizer8.get_feature_names() :#titles tf_idf
    tf_features_names.append(a)
for a in vectorizer9.get_feature_names() :# essays tf_idf
    tf_features_names.append(a)

# top 10 -ves
final_tf_features = pd.DataFrame({'feature_prob_estimates' : tf_features_probs, 'feature_names' :
tf_features_names})
a =final_tf_features.sort_values(by = ['feature_prob_estimates'], ascending = False)
#print(final_bow_features.head(6))

a.head(10)

```

Out[106]:

	feature_prob_estimates	feature_names
13537	-3.631837	worth
13536	-3.713545	worst
13567	-4.103081	xylophones
13566	-4.128766	xylophone
13565	-4.453319	xtramath
13564	-4.743612	xtra
13534	-4.743612	worrying
10043	-4.808164	productivity
13622	-4.827585	zumba
13535	-4.836728	worse

In [107]:

```

#top 10 Positives

# now make a dictionary of all the probabilityies fo the weights
bow_features_probs_pos = []
for a in range(13623):
    bow_features_probs_pos.append(nb.feature_log_prob_[1,a] )# negative feature probabilities

#len(bow_features_probs)
final_bow_features = pd.DataFrame({'feature_prob_estimates_pos' : bow_features_probs_pos,
'feature_names' : bow_features_names})

a =final_bow_features.sort_values(by = ['feature_prob_estimates_pos'], ascending = False)
#print(final_bow_features.head(6))
a.head(10)

```

Out[107]:

	feature_prob_estimates_pos	feature_names
--	----------------------------	---------------

id	feature_pos	feature_names
13537	-3.661805	walking
13536	-3.830088	wiggle
13566	-4.045229	wide
13565	-4.264523	whole
13622	-4.635902	zone
10043	-4.705550	stood
13535	-4.741242	walk
13564	-4.777131	whiteboards
13534	-4.777131	wait

## 3. Conclusions

In [110]:

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

tb = PrettyTable()
tb.field_names= ("Vectorizer", "Model", "HyperParameter", "AUC")
tb.add_row(["BOW", "Auto",1, 70])
tb.add_row(["Tf-Idf", "Auto", 0.1, 67])
print(tb.get_string(titles = "KNN - Observations"))
#print(tb)
```

```
+-----+-----+-----+-----+
| Vectorizer | Model | HyperParameter | AUC |
+-----+-----+-----+-----+
|    BOW    |  Auto |         1       |  70 |
|   Tf-Idf  |  Auto |        0.1      |  67 |
+-----+-----+-----+-----+
```