

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12

Feature	Description	
	One or more (comma-separated) subject categories for the project from the following enumerated list of values:	
project_subject_categories	<ul style="list-style-type: none"> • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth 	
	Examples:	
	<ul style="list-style-type: none"> • Music & The Arts • Literacy & Language, Math & Science 	
school_state	State where school is located (Two-letter U.S. postal code (https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). Example: WY	
project_subject_subcategories	One or more (comma-separated) subject subcategories for the project. Examples:	
	<ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Sciences 	
project_resource_summary	An explanation of the resources needed for the project. Example:	
	<ul style="list-style-type: none"> • My students need hands on literacy materials to manage sensory needs!</code 	
project_essay_1	First application essay*	
project_essay_2	Second application essay*	
project_essay_3	Third application essay*	
project_essay_4	Fourth application essay*	
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245	
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56	
teacher_prefix	Teacher's title. One of the following enumerated values:	
	<ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher. 	
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2	

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

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

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

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

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

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

Notes on the Essay Data

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

- **project_essay_1:** "Introduce us to your classroom"
- **project_essay_2:** "Tell us more about your students"
- **project_essay_3:** "Describe how your students will use the materials you're requesting"
- **project_essay_3:** "Close by sharing why your project will make a difference"

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

- **project_essay_1:** "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- **project_essay_2:** "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [132]:

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4 import sqlite3
5 import pandas as pd
6 import numpy as np
7 import nltk
8 import string
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 from sklearn.feature_extraction.text import TfidfTransformer
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn.metrics import confusion_matrix
15 from sklearn import metrics
16 from sklearn.metrics import roc_curve, auc
17 from nltk.stem.porter import PorterStemmer
18 import re
19 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
20 import string
21 from nltk.corpus import stopwords
22 from nltk.stem import PorterStemmer
23 from nltk.stem.wordnet import WordNetLemmatizer
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27 from tqdm import tqdm
28 import os
29 import chart_studio.plotly
30 # from plotly import plotly
31 import plotly.offline as offline
32 import plotly.graph_objs as go
33 offline.init_notebook_mode()
34 from collections import Counter
35 from scipy.sparse import hstack, vstack
36 from sklearn.model_selection import train_test_split
37 from sklearn.neighbors import KNeighborsClassifier
38 from sklearn.metrics import accuracy_score
39 from sklearn.model_selection import cross_val_score
40 from sklearn import model_selection
41 from sklearn.preprocessing import StandardScaler
42 from sklearn.model_selection import RandomizedSearchCV
43 #from sklearn.impute import SimpleImputer
44 from sklearn.datasets import load_digits
45 #from sklearn.feature_selection import SelectKBest, chi2
```

```

46 from sklearn.model_selection import GridSearchCV
47 from sklearn.feature_selection import SelectKBest,f_classif
48 from prettytable import PrettyTable
49 from sklearn.naive_bayes import MultinomialNB
50 from sklearn.preprocessing import Normalizer
51 from sklearn.metrics import confusion_matrix
52 #import math
53 import pdb

```

1.1 Reading Data

In [133]:

```

1 Project_data = pd.read_csv('train_data.csv')
2 Resource_data = pd.read_csv('resourceS.csv')
3 print(Project_data.shape)
4 print(Resource_data.shape)

```

```

(109248, 17)
(1541272, 4)

```

In [134]:

```

1 # how to replace elements in List python: https://stackoverflow.com/a/2582163/4084039
2 cols = ['Date' if x=='project_submitted_datetime' else x for x in list(Project_data.columns)]
3 #sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
4 Project_data['Date'] = pd.to_datetime(Project_data['project_submitted_datetime'])
5 Project_data.drop('project_submitted_datetime', axis=1, inplace=True)
6 Project_data.sort_values(by=['Date'], inplace=True)
7 # how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
8 Project_data = Project_data[cols]
9 Project_data.head(2)

```

Out[134]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject_categories	project
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	Math & Science	Applie
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	Special Needs	

1.2 preprocessing of project_subject_categories

```
In [135]: 1 y = Project_data['project_is_approved'].values
          2 Project_data.drop(['project_is_approved'], axis=1, inplace=True)
          3 lpd = len(Project_data)
          4 ys = np.zeros(lpd, dtype=np.int32)
          5 X = Project_data
```

```
In [136]: 1 #Splitting the Dataset into three Train,CV and Test
          2 X1, X_Test, Y1, Y_Test = train_test_split(X, y, test_size=0.33, random_state=0, stratify=ys)
          3 nx1 = len(X1)
          4 ys1 = np.zeros(nx1, dtype=np.int32)
          5 X_Train, X_CV, Y_Train, Y_CV = train_test_split(X1, Y1, test_size=0.33, random_state=0, stratify=ys1)
          6 print('Shape of the X_Train data is {0} and Y_Train data is: {1}'.format(X_Train.shape,Y_Train.shape[0]))
          7 print('Shape of the X_CV data is {0} and Y_CV data is : {1}'.format(X_CV.shape,Y_CV.shape[0]))
          8 print('Shape of the X_Test data is {0} and Y_Test data is : {1}'.format(X_Test.shape,Y_Test.shape[0]))
```

Shape of the X_Train data is (49041, 16) and Y_Train data is: 49041

Shape of the X_CV data is (24155, 16) and Y_CV data is : 24155

Shape of the X_Test data is (36052, 16) and Y_Test data is : 36052

In [137]:

```
1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6
7 #*****Train Data*****
8 categories = list(X_Train['project_subject_categories'].values)
9 cat_list = []
10 for i in categories:
11     temp = ""
12     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
13     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
14         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
15             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
16             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
17             temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
18             temp = temp.replace('&', '_') # we are replacing the & value into
19     cat_list.append(temp.strip())
20
21 X_Train['clean_categories'] = cat_list
22 X_Train.drop(['project_subject_categories'], axis=1, inplace=True)
23
24 from collections import Counter
25 my_counter = Counter()
26 for word in X_Train['clean_categories'].values:
27     my_counter.update(word.split())
28
29 cat_dict = dict(my_counter)
30 sorted_cat_dict_Train = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
31 print(len(sorted_cat_dict_Train))
32 #*****CV Data*****
33 categories = list(X_CV['project_subject_categories'].values)
34 cat_list = []
35 for i in categories:
36     temp = ""
37     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
38     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
39         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
40             j=j.replace('The', '') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
41             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
42             temp+=j.strip()+" " #" abc ".strip() will return "abc", remove the trailing spaces
43             temp = temp.replace('&', '_') # we are replacing the & value into
44     cat_list.append(temp.strip())
45
```

```

46 X_CV['clean_categories'] = cat_list
47 X_CV.drop(['project_subject_categories'], axis=1, inplace=True)
48
49 *****Test Data*****
50 categories = list(X_Test['project_subject_categories'].values)
51 cat_list = []
52 for i in categories:
53     temp = ""
54     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
55     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
56         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
57             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
58             j = j.replace(' ', '') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
59             temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
60             temp = temp.replace('&','_') # we are replacing the & value into
61         cat_list.append(temp.strip())
62
63 X_Test['clean_categories'] = cat_list
64 X_Test.drop(['project_subject_categories'], axis=1, inplace=True)
65

```

9

1.3 preprocessing of project_subject_subcategories

In [138]:

```
1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6 #*****Train Data*****
7 sub_categories = list(X_Train['project_subject_subcategories'].values)
8 sub_cat_list = []
9 for i in sub_categories:
10     temp = ""
11     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
14             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
17             temp = temp.replace('&','_')
18     sub_cat_list.append(temp.strip())
19
20 X_Train['clean_subcategories'] = sub_cat_list
21 X_Train.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in X_Train['clean_subcategories'].values:
26     my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict_Train = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
30 print(len(sorted_sub_cat_dict_Train))
31 #*****CV Data*****
32 sub_categories = list(X_CV['project_subject_subcategories'].values)
33 sub_cat_list = []
34 for i in sub_categories:
35     temp = ""
36     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
37     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
38         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
39             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
40             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
41             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
42             temp = temp.replace('&','_')
43     sub_cat_list.append(temp.strip())
44
45 X_CV['clean_subcategories'] = sub_cat_list
```

```

46 X_CV.drop(['project_subject_subcategories'], axis=1, inplace=True)
47
48 *****Test Data*****
49 sub_categories = list(X_Test['project_subject_subcategories'].values)
50 sub_cat_list = []
51 for i in sub_categories:
52     temp = ""
53     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
54     for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
55         if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
56             j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
57             j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
58             temp +=j.strip()+" #" abc ".strip() will return "abc", remove the trailing spaces
59             temp = temp.replace('&','_')
60     sub_cat_list.append(temp.strip())
61
62 X_Test['clean_subcategories'] = sub_cat_list
63 X_Test.drop(['project_subject_subcategories'], axis=1, inplace=True)

```

30

1.3 Text preprocessing

In [139]:

```

1 # merge two column text dataframe:
2 X_Train["essay"] = X_Train["project_essay_1"].map(str) +\
3                     X_Train["project_essay_2"].map(str) + \
4                     X_Train["project_essay_3"].map(str) + \
5                     X_Train["project_essay_4"].map(str)
6
7 X_CV["essay"] = X_CV["project_essay_1"].map(str) +\
8                 X_CV["project_essay_2"].map(str) + \
9                 X_CV["project_essay_3"].map(str) + \
10                X_CV["project_essay_4"].map(str)
11
12 X_Test["essay"] = X_Test["project_essay_1"].map(str) +\
13                  X_Test["project_essay_2"].map(str) + \
14                  X_Test["project_essay_3"].map(str) + \
15                  X_Test["project_essay_4"].map(str)

```

In [140]:

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

In [141]:

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

In [142]:

```

1 # Combining all the above students
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Essays in Train data set-----
5 preprocessed_essays_Train = []
6 for sentence in tqdm(X_Train['essay'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\r', ' ')
9     sent = sent.replace('\n', ' ')
10    sent = sent.replace('\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14    preprocessed_essays_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Essays in CV data set-----
18 preprocessed_essays_CV = []
19 for sentence in tqdm(X_CV['essay'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\r', ' ')
22     sent = sent.replace('\n', ' ')
23     sent = sent.replace('\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
27     preprocessed_essays_CV.append(sent.lower().strip())
28 # pdb.set_trace()
29
30 #-----PreProcessing of Essays in Test data set-----
31 preprocessed_essays_Test = []
32 for sentence in tqdm(X_Test['essay'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\r', ' ')
35     sent = sent.replace('\n', ' ')
36     sent = sent.replace('\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
40     preprocessed_essays_Test.append(sent.lower().strip())
41 # pdb.set_trace()

```

100%		49041/49041	[00:44<00:00, 1111.52it/s]
100%		24155/24155	[00:20<00:00, 1158.57it/s]

[illegible]

In [143]:

```

1 word_count_essay_Train = []
2 for a in tqdm(X_Train["essay"]) :
3     b = len(a.split())
4     word_count_essay_Train.append(b)
5
6 X_Train["word_count_essay_Train"] = word_count_essay_Train
7
8 word_count_essay_CV = []
9 for a in tqdm(X_CV["essay"]) :
10     b = len(a.split())
11     word_count_essay_CV.append(b)
12
13 X_CV["word_count_essay_CV"] = word_count_essay_CV
14
15 word_count_essay_Test = []
16 for a in tqdm(X_Test["essay"]) :
17     b = len(a.split())
18     word_count_essay_Test.append(b)
19
20 X_Test["word_count_essay_Test"] = word_count_essay_Test

```

[illegible]

1.4 Preprocessing of project_title

In [144]:

```
1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Project Title in Train data set-----
5 preprocessed_titles_Train = []
6 for sentence in tqdm(X_Train['project_title'].values):
7     sent = decontracted(sentence)
8     sent = sent.replace('\\r', ' ')
9     sent = sent.replace('\\\"', ' ')
10    sent = sent.replace('\\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280
13    sent = ' '.join(e for e in sent.split() if e not in stopwords)
14    preprocessed_titles_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Project Title in CV data set-----
18 preprocessed_titles_CV = []
19 for sentence in tqdm(X_CV['project_title'].values):
20     sent = decontracted(sentence)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\\"', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e not in stopwords)
27     preprocessed_titles_CV.append(sent.lower().strip())
28 # pdb.set_trace()
29
30 #-----PreProcessing of Project Title in Test data set-----
31 preprocessed_titles_Test = []
32 for sentence in tqdm(X_Test['project_title'].values):
33     sent = decontracted(sentence)
34     sent = sent.replace('\\r', ' ')
35     sent = sent.replace('\\\"', ' ')
36     sent = sent.replace('\\n', ' ')
37     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
38     # https://gist.github.com/sebleier/554280
39     sent = ' '.join(e for e in sent.split() if e not in stopwords)
40     preprocessed_titles_Test.append(sent.lower().strip())
41 # pdb.set_trace()
```

100%	<div></div>	49041/49041	[00:01<00:00, 24818.33it/s]
100%	<div></div>	24155/24155	[00:00<00:00, 24648.09it/s]
100%	<div></div>	36052/36052	[00:01<00:00, 24359.48it/s]

In [145]:

```

1 word_count_title_Train = []
2 for a in tqdm(X_Train["project_title"]) :
3     b = len(a.split())
4     word_count_title_Train.append(b)
5
6 X_Train["word_count_title_Train"] = word_count_title_Train
7
8 word_count_title_CV = []
9 for a in tqdm(X_CV["project_title"]) :
10     b = len(a.split())
11     word_count_title_CV.append(b)
12
13 X_CV["word_count_title_CV"] = word_count_title_CV
14
15 word_count_title_Test = []
16 for a in tqdm(X_Test["project_title"]) :
17     b = len(a.split())
18     word_count_title_Test.append(b)
19
20 X_Test["word_count_title_Test"] = word_count_title_Test

```

[illegible]

1.5 Preparing data for models

1.5.1 Vectorizing Categorical data

- <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>
(<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>)

In [146]:

```
1  #-----Vectorizing categorical data for Train,CV and Test-----
2
3  # we use count vectorizer to convert the values into one hot encoding
4  vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict_Train.keys()), lowercase=False, binary=True)
5  vectorizer_cat.fit(X_Train['clean_categories'].values)
6  categories_one_hot_Train = vectorizer_cat.transform(X_Train['clean_categories'].values)
7  categories_one_hot_CV = vectorizer_cat.transform(X_CV['clean_categories'].values)
8  categories_one_hot_Test = vectorizer_cat.transform(X_Test['clean_categories'].values)
9  print(vectorizer_cat.get_feature_names())
10 print("-"*120)
11 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Train.shape))
12 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(categories_one_hot_CV.shape))
13 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Test.shape))
```

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

```
Shape of Train dataset matrix after one hot encoding is: (49041, 9)
```

```
Shape of CV dataset matrix after one hot encoding is: (24155, 9)
```

```
Shape of Test dataset matrix after one hot encoding is: (36052, 9)
```



```
In [147]: 1 #-----Vectorizing Sub categorical data for Train,CV and Test-----
2
3 # we use count vectorizer to convert the values into one hot encoding
4 vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict_Train.keys()), lowercase=False, binary=True)
5 vectorizer_sub_cat.fit(X_Train['clean_categories'].values)
6 sub_categories_one_hot_Train = vectorizer_sub_cat.transform(X_Train['clean_subcategories'].values)
7 sub_categories_one_hot_CV = vectorizer_sub_cat.transform(X_CV['clean_subcategories'].values)
8 sub_categories_one_hot_Test = vectorizer_sub_cat.transform(X_Test['clean_subcategories'].values)
9 print(vectorizer_sub_cat.get_feature_names())
10 print("-"*120)
11 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Train.shape))
12 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_CV.shape))
13 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Test.shape))
```

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

```
-----
Shape of Train dataset matrix after one hot encoding is: (49041, 30)
```

```
Shape of CV dataset matrix after one hot encoding is: (24155, 30)
```

```
Shape of Test dataset matrix after one hot encoding is: (36052, 30)
```

School State

In [148]:

```
1  #-----Vectorizing categorical data of School state for Train dataset-----
2
3  school_catogories_Train = list(X_Train['school_state'].values)
4  school_list_Train = []
5  for sent in school_catogories_Train:
6      school_list_Train.append(sent.lower().strip())
7  X_Train['school_categories'] = school_list_Train
8  X_Train.drop(['school_state'], axis=1, inplace=True)
9
10 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
11 my_counter_school_Train = Counter()
12 for word in X_Train['school_categories'].values:
13     my_counter_school_Train.update(word.split())
14
15 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
16 school_dict_Train = dict(my_counter_school_Train)
17 sorted_school_dict_Train = dict(sorted(school_dict_Train.items(), key=lambda kv: kv[1]))
18
19 vectorizer_school = CountVectorizer(vocabulary=list(sorted_school_dict_Train.keys()), lowercase=False, binary=True)
20 vectorizer_school.fit(X_Train['school_categories'].values)
21 #print(vectorizer.get_feature_names())
22
23 school_one_hot_Train = vectorizer_school.transform(X_Train['school_categories'].values)
24
25 #-----Vectorizing categorical data of School state for CV dataset-----
26
27 school_catogories_CV = list(X_CV['school_state'].values)
28 school_list_CV = []
29 for sent in school_catogories_CV:
30     school_list_CV.append(sent.lower().strip())
31 X_CV['school_categories'] = school_list_CV
32 X_CV.drop(['school_state'], axis=1, inplace=True)
33
34 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
35 my_counter_school_CV = Counter()
36 for word in X_CV['school_categories'].values:
37     my_counter_school_CV.update(word.split())
38
39 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
40 school_dict_CV = dict(my_counter_school_CV)
41 sorted_school_dict_CV = dict(sorted(school_dict_CV.items(), key=lambda kv: kv[1]))
42 school_one_hot_CV = vectorizer_school.transform(X_CV['school_categories'].values)
43
44 #-----Vectorizing categorical data of School state for Test dataset-----
45
```

```

46 school_catogories_Test = list(X_Test['school_state'].values)
47 school_list_Test = []
48 for sent in school_catogories_Test:
49     school_list_Test.append(sent.lower().strip())
50 X_Test['school_categories'] = school_list_Test
51 X_Test.drop(['school_state'], axis=1, inplace=True)
52
53 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
54 my_counter_school_Test = Counter()
55 for word in X_Test['school_categories'].values:
56     my_counter_school_Test.update(word.split())
57
58 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
59 school_dict_Test = dict(my_counter_school_Test)
60 sorted_school_dict_Test = dict(sorted(school_dict_Test.items(), key=lambda kv: kv[1]))
61 school_one_hot_Test = vectorizer_school.transform(X_Test['school_categories'].values)
62 print("-"*120)
63 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(school_one_hot_Train.shape))
64 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(school_one_hot_CV.shape))
65 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(school_one_hot_Test.shape))

```

```

Shape of Train dataset matrix after one hot encoding is: (49041, 51)
Shape of CV dataset matrix after one hot encoding is: (24155, 51)
Shape of Test dataset matrix after one hot encoding is: (36052, 51)

```

Prefix

In [149]:

```
1 #-----Vectorizing categorical data of Teacher Prefix for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 prefix_catogories_Train = list(X_Train['teacher_prefix'].values)
8 prefix_list_Train = []
9 for sent in prefix_catogories_Train:
10     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
11     # https://gist.github.com/sebleier/554280
12     sent = ' '.join(e for e in sent.split())
13     prefix_list_Train.append(sent.lower().strip())
14 X_Train['prefix_catogories'] = prefix_list_Train
15 X_Train.drop(['teacher_prefix'], axis=1, inplace=True)
16
17 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
18 my_counter_prefix_Train = Counter()
19 for word in X_Train['prefix_catogories'].values:
20     my_counter_prefix_Train.update(word.split())
21
22 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
23 prefix_dict_Train = dict(my_counter_prefix_Train)
24 sorted_prefix_dict_Train = dict(sorted(prefix_dict_Train.items(), key=lambda kv: kv[1]))
25
26
27 vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict_Train.keys()), lowercase=False, binary=True)
28 vectorizer_prefix.fit(X_Train['prefix_catogories'].values)
29 #print(vectorizer.get_feature_names())
30
31 prefix_one_hot_Train = vectorizer_prefix.transform(X_Train['prefix_catogories'].values)
32 #print("Shape of matrix after one hot encodig ",prefix_one_hot.shape)
33
34 #-----Vectorizing categorical data of Teacher Prefix for CV dataset-----
35
36 prefix_catogories_CV = list(X_CV['teacher_prefix'].values)
37 prefix_list_CV = []
38 for sent in prefix_catogories_CV:
39     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
40     # https://gist.github.com/sebleier/554280
41     sent = ' '.join(e for e in sent.split())
42     prefix_list_CV.append(sent.lower().strip())
43 X_CV['prefix_catogories'] = prefix_list_CV
44 X_CV.drop(['teacher_prefix'], axis=1, inplace=True)
45
```

```

46 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
47 my_counter_prefix_CV = Counter()
48 for word in X_CV['prefix_catogories'].values:
49     my_counter_prefix_CV.update(word.split())
50
51 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
52 prefix_dict_CV = dict(my_counter_prefix_CV)
53 sorted_prefix_dict_CV = dict(sorted(prefix_dict_CV.items(), key=lambda kv: kv[1]))
54 prefix_one_hot_CV = vectorizer_prefix.transform(X_CV['prefix_catogories'].values)
55
56 #-----Vectorizing categorical data of Teacher Prefix for Test dataset-----
57
58 prefix_catogories_Test = list(X_Test['teacher_prefix'].values)
59 prefix_list_Test = []
60 for sent in prefix_catogories_Test:
61     sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
62     # https://gist.github.com/sebleier/554280
63     sent = ' '.join(e for e in sent.split())
64     prefix_list_Test.append(sent.lower().strip())
65 X_Test['prefix_catogories'] = prefix_list_Test
66 X_Test.drop(['teacher_prefix'], axis=1, inplace=True)
67
68 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
69 my_counter_prefix_Test = Counter()
70 for word in X_Test['prefix_catogories'].values:
71     my_counter_prefix_Test.update(word.split())
72
73 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
74 prefix_dict_Test = dict(my_counter_prefix_Test)
75 sorted_prefix_dict_Test = dict(sorted(prefix_dict_Test.items(), key=lambda kv: kv[1]))
76 prefix_one_hot_Test = vectorizer_prefix.transform(X_Test['prefix_catogories'].values)
77 print("-"*120)
78 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Train.shape))
79 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_CV.shape))
80 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Test.shape))

```

```

-----
Shape of Train dataset matrix after one hot encoding is: (49041, 6)
Shape of CV dataset matrix after one hot encoding is: (24155, 6)
Shape of Test dataset matrix after one hot encoding is: (36052, 6)

```

project_grade_category

In [150]:

```
1 #-----Vectorizing categorical data of Project Grade for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 grade_catogories_Train = list(X_Train['project_grade_category'].values)
8 grade_list_Train = []
9 for sent in grade_catogories_Train:
10     sent = sent.replace('-', '_')
11     sent = sent.replace(' ', '_')
12     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
13     # https://gist.github.com/sebleier/554280
14     sent = ' '.join(e for e in sent.split())
15     grade_list_Train.append(sent.lower().strip())
16
17 # temp = temp.replace('-', '_')
18 X_Train['new_grade_category'] = grade_list_Train
19 X_Train.drop(['project_grade_category'], axis=1, inplace=True)
20
21 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
22 my_counter_grade_Train = Counter()
23 for word in X_Train['new_grade_category'].values:
24     my_counter_grade_Train.update(word.split())
25
26 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
27 grade_dict_Train = dict(my_counter_grade_Train)
28 sorted_grade_dict_Train = dict(sorted(grade_dict_Train.items(), key=lambda kv: kv[1]))
29
30 vectorizer_grade = CountVectorizer(vocabulary=list(sorted_grade_dict_Train.keys()), lowercase=False, binary=True)
31 vectorizer_grade.fit(X_Train['new_grade_category'].values)
32 #print(vectorizer.get_feature_names())
33
34 grade_one_hot_Train = vectorizer_grade.transform(X_Train['new_grade_category'].values)
35
36 #-----Vectorizing categorical data of Project Grade for CV dataset-----
37
38 grade_catogories_CV = list(X_CV['project_grade_category'].values)
39 grade_list_CV = []
40 for sent in grade_catogories_CV:
41     sent = sent.replace('-', '_')
42     sent = sent.replace(' ', '_')
43     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
44     # https://gist.github.com/sebleier/554280
45     sent = ' '.join(e for e in sent.split())
```

```

46     grade_list_CV.append(sent.lower().strip())
47
48     # temp = temp.replace('-', '_')
49     X_CV['new_grade_category'] = grade_list_CV
50     X_CV.drop(['project_grade_category'], axis=1, inplace=True)
51
52     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
53     my_counter_grade_CV = Counter()
54     for word in X_CV['new_grade_category'].values:
55         my_counter_grade_CV.update(word.split())
56
57     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
58     grade_dict_CV = dict(my_counter_grade_CV)
59     sorted_grade_dict_CV = dict(sorted(grade_dict_CV.items(), key=lambda kv: kv[1]))
60
61     grade_one_hot_CV = vectorizer_grade.transform(X_CV['new_grade_category'].values)
62
63     #-----Vectorizing categorical data of Project Grade for Train dataset-----
64
65     grade_categories_Test = list(X_Test['project_grade_category'].values)
66     grade_list_Test = []
67     for sent in grade_categories_Test:
68         sent = sent.replace('-', '_')
69         sent = sent.replace(' ', '_')
70         # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
71         # https://gist.github.com/sebleier/554280
72         sent = ' '.join(e for e in sent.split())
73         grade_list_Test.append(sent.lower().strip())
74
75     # temp = temp.replace('-', '_')
76     X_Test['new_grade_category'] = grade_list_Test
77     X_Test.drop(['project_grade_category'], axis=1, inplace=True)
78
79     # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
80     my_counter_grade_Test = Counter()
81     for word in X_Test['new_grade_category'].values:
82         my_counter_grade_Test.update(word.split())
83
84     # dict sort by value python: https://stackoverflow.com/a/613218/4084039
85     grade_dict_Test = dict(my_counter_grade_Test)
86     sorted_grade_dict_Test = dict(sorted(grade_dict_Test.items(), key=lambda kv: kv[1]))
87
88     grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)
89     print("-"*120)
90     print('Shape of Train dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Train.shape))
91     print('Shape of CV dataset matrix after one hot encoding is: {}'.format(grade_one_hot_CV.shape))

```

```
92 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Test.shape))
```

Shape of Train dataset matrix after one hot encoding is: (49041, 4)
Shape of CV dataset matrix after one hot encoding is: (24155, 4)
Shape of Test dataset matrix after one hot encoding is: (36052, 4)

1.5.2 Vectorizing Numerical features

```
In [151]: 1 price_data = Resource_data.groupby('id').agg({'price': 'sum', 'quantity': 'sum'}).reset_index()
          2 X_Train = pd.merge(X_Train, price_data, on='id', how='left')
          3 X_CV = pd.merge(X_CV, price_data, on='id', how='left')
          4 X_Test = pd.merge(X_Test, price_data, on='id', how='left')
```

```
In [152]: 1 price_norm = Normalizer(norm='l2', copy=False)
          2 price_norm.fit(X_Train['price'].values.reshape(1,-1))
          3
          4 p=price_norm.transform(X_Train['price'].values.reshape(1,-1))
          5 price_norm.transform(X_CV['price'].values.reshape(1,-1))
          6 price_norm.transform(X_Test['price'].values.reshape(1,-1))
          7 price_norm_Train = (X_Train['price'].values.reshape(-1,1))
          8 price_norm_CV = (X_CV['price'].values.reshape(-1,1))
          9 price_norm_Test = (X_Test['price'].values.reshape(-1,1))
         10 print("-"*120)
         11 print('Shape of Train normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Train.shape))
         12 print('Shape of CV normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_CV.shape))
         13 print('Shape of Test normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Test.shape))
```

Shape of Train normalized price dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized price dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized price dataset matrix after one hot encoding is: (36052, 1)


```
In [153]: 1 quantity_norm = Normalizer(norm='l2', copy=False)
2 quantity_norm.fit(X_Train['quantity'].values.reshape(1,-1))
3
4 quantity_norm.transform(X_Train['quantity'].values.reshape(1,-1))
5 quantity_norm.transform(X_CV['quantity'].values.reshape(1,-1))
6 quantity_norm.transform(X_Test['quantity'].values.reshape(1,-1))
7 quantity_norm_Train = quantity_norm.transform(X_Train['quantity'].values.reshape(-1,1))
8 quantity_norm_CV = quantity_norm.transform(X_CV['quantity'].values.reshape(-1,1))
9 quantity_norm_Test = quantity_norm.transform(X_Test['quantity'].values.reshape(-1,1))
10 print("-"*120)
11 print('Shape of Train normalized quantity dataset matrix after one hot encoding is: {0}'.format(quantity_norm_Train.shape))
12 print('Shape of CV normalized quantity dataset matrix after one hot encoding is: {0}'.format(quantity_norm_CV.shape))
13 print('Shape of Test normalized quantity dataset matrix after one hot encoding is: {0}'.format(quantity_norm_Test.shape))
```

```
-----
Shape of Train normalized quantity dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized quantity dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized quantity dataset matrix after one hot encoding is: (36052, 1)
```

```
In [154]: 1 teacher_prev_post_norm = Normalizer(norm='l2', copy=False)
2 teacher_prev_post_norm.fit(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
3
4 teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
5 teacher_prev_post_norm.transform(X_CV['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
6 teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
7 teacher_prev_post_norm_Train = teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
8 teacher_prev_post_norm_CV = teacher_prev_post_norm.transform(X_CV['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
9 teacher_prev_post_norm_Test = teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
10 print("-"*120)
11 print('Shape of Train normalized previously posted project dataset matrix after one hot encoding is: {0}'.format(teacher_prev_post_norm_Train.shape))
12 print('Shape of CV normalized previously posted project dataset matrix after one hot encoding is: {0}'.format(teacher_prev_post_norm_CV.shape))
13 print('Shape of Test normalized previously posted project dataset matrix after one hot encoding is: {0}'.format(teacher_prev_post_norm_Test.shape))
```

```
-----
Shape of Train normalized previously posted project dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized previously posted project dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized previously posted project dataset matrix after one hot encoding is: (36052, 1)
```

```
In [155]: 1 title_norm = Normalizer(norm='l2', copy=False)
2 title_norm.fit(X_Train['word_count_title_Train'].values.reshape(1,-1))
3 title_norm.transform(X_Train['word_count_title_Train'].values.reshape(1,-1))
4 title_norm.transform(X_CV['word_count_title_CV'].values.reshape(1,-1))
5 title_norm.transform(X_Test['word_count_title_Test'].values.reshape(1,-1))
6 word_count_title_Train = title_norm.transform(X_Train['word_count_title_Train'].values.reshape(-1,1))
7 word_count_title_CV = title_norm.transform(X_CV['word_count_title_CV'].values.reshape(-1,1))
8 word_count_title_Test = title_norm.transform(X_Test['word_count_title_Test'].values.reshape(-1,1))
9 print("-"*120)
10 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Train.shape))
11 print('Shape of CV normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_CV.shape))
12 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Test.shape))
```

```
-----
Shape of Train normalized title dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized title dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)
```

```
In [156]: 1 essay_norm = Normalizer(norm='l2', copy=False)
2 essay_norm.fit(X_Train['word_count_essay_Train'].values.reshape(1,-1))
3 essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(1,-1))
4 essay_norm.transform(X_CV['word_count_essay_CV'].values.reshape(1,-1))
5 essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(1,-1))
6 word_count_essay_Train = essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(-1,1))
7 word_count_essay_CV = essay_norm.transform(X_CV['word_count_essay_CV'].values.reshape(-1,1))
8 word_count_essay_Test = essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(-1,1))
9 print("-"*120)
10 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Train.shape))
11 print('Shape of CV normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_CV.shape))
12 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Test.shape))
```

```
-----
Shape of Train normalized title dataset matrix after one hot encoding is: (49041, 1)
Shape of CV normalized title dataset matrix after one hot encoding is: (24155, 1)
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)
```

1.5.3 Vectorizing Text data

1.5.3.1 Bag of words

```
In [157]: 1 # We are considering only the words which appeared in at least 10 documents(rows or projects).
2 vectorizer_essays_bow = CountVectorizer(min_df=10)
3 text_bow_Train = vectorizer_essays_bow.fit_transform(preprocessed_essays_Train)
4 text_bow_CV = vectorizer_essays_bow.transform(preprocessed_essays_CV)
5 text_bow_Test = vectorizer_essays_bow.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying Bag Of Words for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(text_bow_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(text_bow_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(text_bow_Test.shape))
```

Applying Bag Of Words for Text Data

Shape of Train dataset matrix after one hot encoding is: (49041, 12097)
Shape of CV dataset matrix after one hot encoding is: (24155, 12097)
Shape of Test dataset matrix after one hot encoding is: (36052, 12097)

Bag of Words for Project Title

```
In [158]: 1 # you can vectorize the title also
2 # before you vectorize the title make sure you preprocess it
3 vectorizer_titles_bow = CountVectorizer(min_df=10)
4 title_bow_Train = vectorizer_titles_bow.fit_transform(preprocessed_titles_Train)
5 title_bow_CV = vectorizer_titles_bow.transform(preprocessed_titles_CV)
6 title_bow_Test = vectorizer_titles_bow.transform(preprocessed_titles_Test)
7 print("-"*120)
8 print("Applying Bag Of Words for Project Title Data")
9 print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(title_bow_Train.shape))
11 print('Shape of CV dataset matrix after one hot encoding is: {}'.format(title_bow_CV.shape))
12 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(title_bow_Test.shape))
13
```

Applying Bag Of Words for Project Title Data

Shape of Train dataset matrix after one hot encoding is: (49041, 2083)
Shape of CV dataset matrix after one hot encoding is: (24155, 2083)
Shape of Test dataset matrix after one hot encoding is: (36052, 2083)

1.5.2.2 TFIDF vectorizer

TFIDF Vectorizer

```
In [159]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer_essays_tfidf = TfidfVectorizer(min_df=10)
3 text_tfidf_Train = vectorizer_essays_tfidf.fit_transform(preprocessed_essays_Train)
4 text_tfidf_CV = vectorizer_essays_tfidf.transform(preprocessed_essays_CV)
5 text_tfidf_Test = vectorizer_essays_tfidf.transform(preprocessed_essays_Test)
6 print("-"*120)
7 print("Applying TFIDF for Text Data")
8 print("-"*120)
9 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Train.shape))
10 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(text_tfidf_CV.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(text_tfidf_Test.shape))
```

Applying TFIDF for Text Data

Shape of Train dataset matrix after one hot encoding is: (49041, 12097)
Shape of CV dataset matrix after one hot encoding is: (24155, 12097)
Shape of Test dataset matrix after one hot encoding is: (36052, 12097)

TFIDF vectorizer for Project Title

```
In [160]: 1 vectorizer_titles_tfidf = TfidfVectorizer(min_df=10)
2 title_tfidf_Train = vectorizer_titles_tfidf.fit_transform(preprocessed_titles_Train)
3 title_tfidf_CV = vectorizer_titles_tfidf.transform(preprocessed_titles_CV)
4 title_tfidf_Test = vectorizer_titles_tfidf.transform(preprocessed_titles_Test)
5 print("-"*120)
6 print("Applying TFIDF for Project Title")
7 print("-"*120)
8 print('Shape of Train dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Train.shape))
9 print('Shape of CV dataset matrix after one hot encoding is: {0}'.format(title_tfidf_CV.shape))
10 print('Shape of Test dataset matrix after one hot encoding is: {0}'.format(title_tfidf_Test.shape))
```

Applying TFIDF for Project Title

Shape of Train dataset matrix after one hot encoding is: (49041, 2083)
Shape of CV dataset matrix after one hot encoding is: (24155, 2083)
Shape of Test dataset matrix after one hot encoding is: (36052, 2083)

Assignment 4: Naive Bayes

1. Apply Multinomial NaiveBayes on these feature sets

- **Set 1**: categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- **Set 2**: categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)

2. The hyper paramter tuning(find best Alpha)

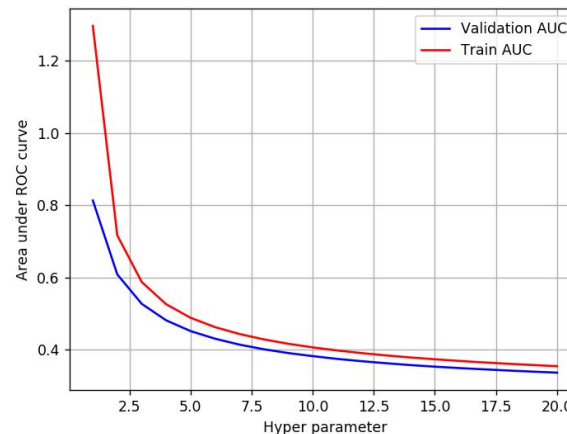
- Find the best hyper parameter which will give the maximum [AUC \(https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/\)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Consider a wide range of alpha values for hyperparameter tuning, start as low as 0.00001
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Feature importance

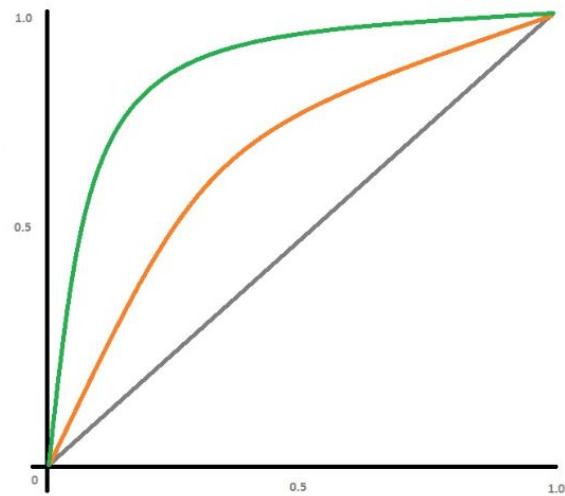
- Find the top 10 features of positive class and top 10 features of negative class for both feature sets **Set 1** and **Set 2** using values of feature_log_prob_ parameter of [MultinomialNB \(https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html) and print their corresponding feature names

4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure. Here on X-axis you will have alpha values, since they have a wide range, just to represent those alpha values on the graph, apply log function on those alpha values.



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



- Along with plotting ROC curve, you need to print the [confusion matrix](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/>) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

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

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

5. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link (<http://zetcode.com/python/prettytable/>).

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

1.5.4 Merging all the above features

- we need to merge all the numerical vectors i.e catogorical, text, numerical vectors

```
In [161]: 1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 BOW_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one_h
3 print(BOW_Train.shape)
4 TFIDF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one
5 print(TFIDF_Train.shape)
```

(49041, 14285)

(49041, 14285)

```
In [162]: 1 BOW_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,text_bow
2 print(BOW_CV.shape)
3 TFIDF_CV = hstack((categories_one_hot_CV,sub_categories_one_hot_CV,school_one_hot_CV,grade_one_hot_CV,prefix_one_hot_CV,text_t
4 print(TFIDF_CV.shape)
```

(24155, 14285)

(24155, 14285)

```
In [163]: 1 BOW_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Te
2 print(BOW_Test.shape)
3 TFIDF_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_
4 print(TFIDF_Test.shape)
```

(36052, 14285)

(36052, 14285)

Loading Pickle files

2. Naive Bayes

2.4 Applying NB() on different kind of featurization as mentioned in the instructions

Apply Naive Bayes on different kind of featurization as mentioned in the instructions
For Every model that you work on make sure you do the step 2 and step 3 of instructions

2.4.1 Applying Naive Bayes on BOW, SET 1

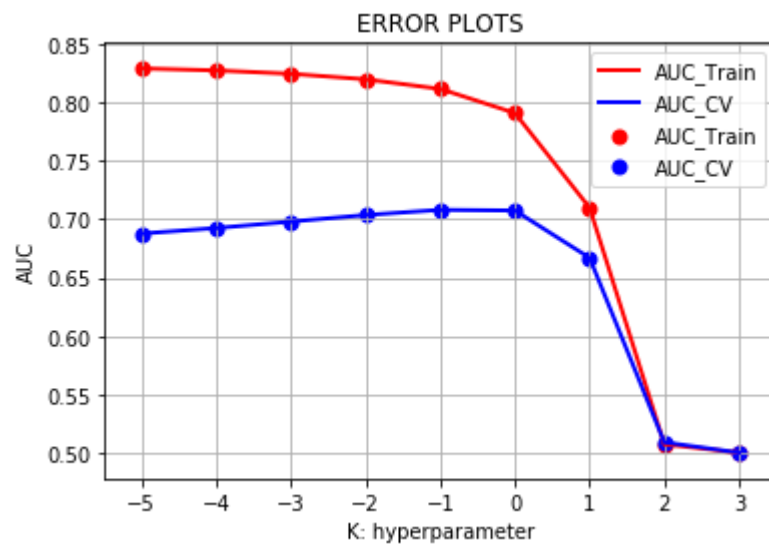
In [164]:

```

1 %%time
2
3 BOW_TR_CSR = BOW_Train.tocsr()
4 BOW_CV_CSR = BOW_CV.tocsr()
5 BOW_TS_CSR = BOW_Test.tocsr()
6
7 alpha = [0.00001, 0.0001, 0.001, 0.01,0.1, 1, 10,100,1000]
8 L_ALPHA= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(alpha):
14     nb = MultinomialNB(alpha = i,fit_prior=True,class_prior=[0.5,0.5])
15     # pdb.set_trace()
16     nb.fit(BOW_Train, Y_Train)
17     pred = nb.predict(BOW_CV)
18     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
19     ACCV.append(acc)
20     a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, nb.predict_proba(BOW_TR_CSR)[:,1])
21     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
22
23     a_fpr_cv, a_tpr_cv, thresholds = roc_curve(Y_CV, nb.predict_proba(BOW_CV_CSR)[:,1])
24     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
25
26
27 for av in tqdm(alpha):
28     b = np.log10(av)
29     L_ALPHA.append(b)
30
31 # Performance of model on Train data and Test data for each hyper parameter.
32 plt.plot(L_ALPHA, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.scatter(L_ALPHA, AUC_TR, label='AUC_Train',color='red',linewidth=2)
34 plt.gca()
35 plt.plot(L_ALPHA, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.scatter(L_ALPHA, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
37 plt.gca()
38 plt.legend()
39 plt.xlabel("K: hyperparameter")
40 plt.ylabel("AUC")
41 plt.title("ERROR PLOTS")
42 plt.grid()
43 plt.show()

```

[illegible][illegible]



Wall time: 21.4 s
 Compiler : 376 ms
 Parser : 291 ms

In [165]:

```
1 print(AUC_TR)
2 print(AUC_CV)
```

```
[0.8294012820005231, 0.8276958244116213, 0.8248167281856434, 0.8200770452977412, 0.811853158302271, 0.7911488255143547, 0.7099457
736235792, 0.5071686477789238, 0.5]
[0.6878082577138354, 0.692571394489643, 0.6980880457054459, 0.7037202471954482, 0.7081541941967424, 0.7075945397845345, 0.6669896
607278614, 0.5088766907673291, 0.5]
```

OBSERVATION:

1. Optimal alpha has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha(\text{alpha})=0.1$ is having the highest AUC CV score.
3. Hence $\alpha(\text{alpha})=0.1$ is choosen as the optimal alpha value.

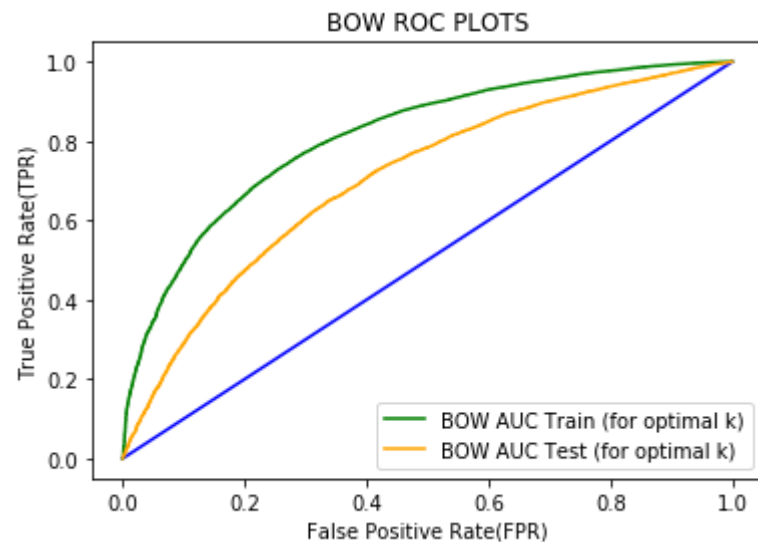
```
In [166]: 1 A_OPT=0.1
2 BOW_opt = MultinomialNB(alpha=A_OPT,fit_prior=True,class_prior=[0.5,0.5])
3 BOW_opt.fit(BOW_Train, Y_Train)
4 pred = BOW_opt.predict(BOW_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha(a) = {0} is {1}%'.format(A_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, BOW_opt.predict_proba(BOW_TR_CSR)[:,1])
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, BOW_opt.predict_proba(BOW_TS_CSR)[:,1])
```

Test accuracy for alpha(a) = 0.1 is 69.73815599689338%

BOW ROC PLOT

In [167]:

```
1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="BOW AUC Train (for optimal k)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="BOW AUC Test (for optimal k)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("BOW ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal k) =", auc(a_fpr_Test, a_tpr_Test))
16 BOW_AOPT=A_OPT
17 BOW_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred1 = BOW_opt.predict(BOW_Train)
19 pred2 = BOW_opt.predict(BOW_Test)
20
```



AUC Train (for optimal k) = 0.811853158302271

AUC Test (for optimal k) = 0.7047236730122236

Wall time: 1.21 s

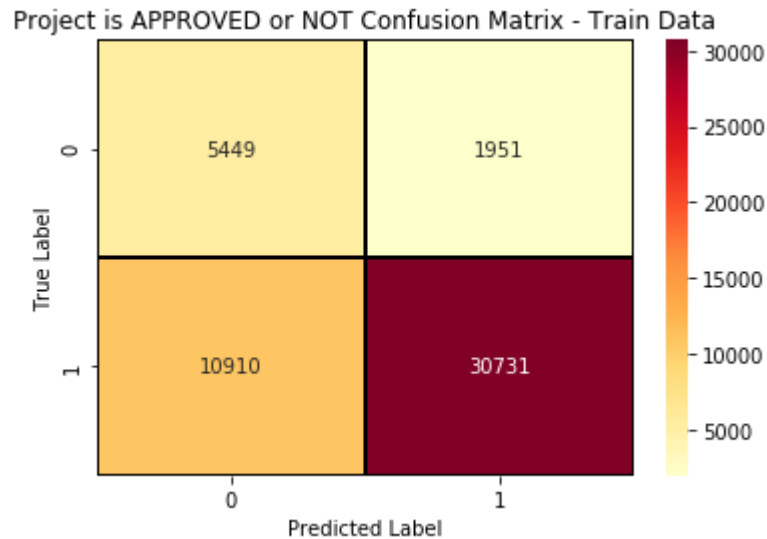
BOW CONFUSION MATRIX

In [168]:

```
1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred1)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Wall time: 1.24 s

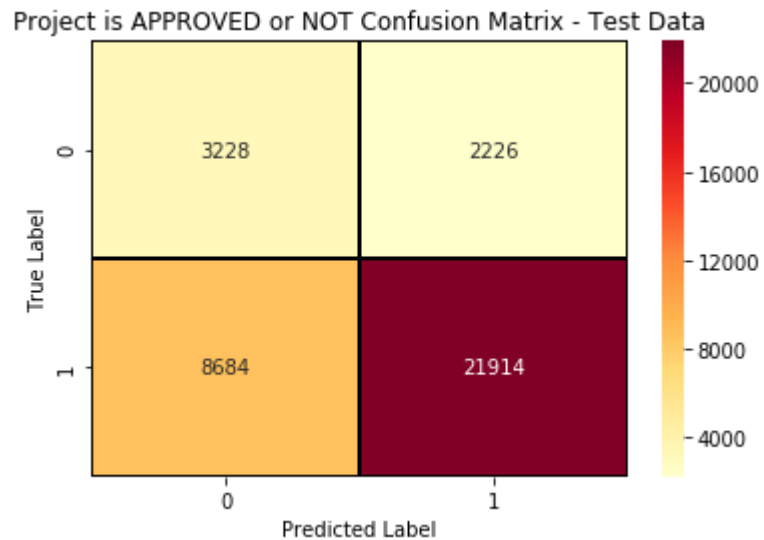
Out[168]: Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



```
In [169]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred2)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

Wall time: 240 ms

Out[169]: Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')



2.4.1.1 Top 10 important features of positive class from SET 1

```
In [170]: 1 B_NB = MultinomialNB(alpha = A_OPT, class_prior=[0.5,0.5])
2 B_NB.fit(BOW_TR_CSR, Y_Train)
3
4 BOW_F = []
5
6 for BOW in vectorizer_cat.get_feature_names() :
7     BOW_F.append(BOW)
8 for BOW1 in vectorizer_sub_cat.get_feature_names() :
9     BOW_F.append(BOW1)
10 for BOW2 in vectorizer_school.get_feature_names() :
11     BOW_F.append(BOW2)
12 for BOW3 in vectorizer_grade.get_feature_names() :
13     BOW_F.append(BOW3)
14 for BOW4 in vectorizer_prefix.get_feature_names() :
15     BOW_F.append(BOW4)
16 for BOW5 in vectorizer_essays_bow.get_feature_names() :
17     BOW_F.append(BOW5)
18 for BOW6 in vectorizer_titles_bow.get_feature_names() :
19     BOW_F.append(BOW6)
20 BOW_F.append("price")
21 BOW_F.append("quantity")
22 BOW_F.append("prev_proposed_projects")
23 BOW_F.append("title_word_count")
24 BOW_F.append("essay_word_count")
25 len(BOW_F)
```

Out[170]: 14285

```
In [171]: 1 POS_F_BOW = B_NB.feature_log_prob_[1, :].argsort()[::-1]
2 for i in POS_F_BOW[:10]:
3     print(BOW_F[i])
```

```
students
school
learning
classroom
not
learn
help
essay_word_count
title_word_count
quantity
```

2.4.1.2 Top 10 important features of negative class from SET 1

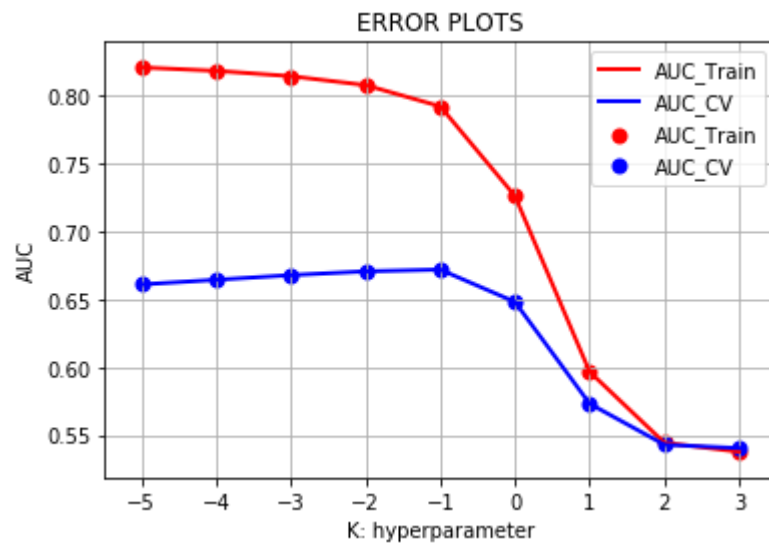
```
In [172]: 1 NEG_F_BOW = B_NB.feature_log_prob_[0, :].argsort()[::-1]
2         for i in NEG_F_BOW[:10]:
3             print(BOW_F[i])
```

```
students
school
learning
classroom
not
learn
help
essay_word_count
title_word_count
quantity
```

2.4.2 Applying Naive Bayes on TFIDF, SET 2

In [173]:

```
1 %%time
2
3 TFIDF_TR_CSR = TFIDF_Train.tocsr()
4 TFIDF_CV_CSR = TFIDF_CV.tocsr()
5 TFIDF_TS_CSR = TFIDF_Test.tocsr()
6
7 alpha = [0.00001, 0.0001, 0.001, 0.01,0.1,1,10,100,1000]
8 L_ALPHA= []
9 ACCV = []
10 AUC_TR = []
11 AUC_CV = []
12
13 for i in tqdm(alpha):
14     nb = MultinomialNB(alpha = i,fit_prior=True,class_prior=[0.5,0.5])
15     # pdb.set_trace()
16     nb.fit(TFIDF_Train, Y_Train)
17     pred = nb.predict(TFIDF_CV)
18     acc = accuracy_score(Y_CV, pred, normalize=True) * float(100)
19     ACCV.append(acc)
20     a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, nb.predict_proba(TFIDF_TR_CSR)[:,1])
21     AUC_TR.append(auc(a_fpr_train, a_tpr_train))
22
23     a_fpr_cv, a_tpr_cv, thresholds = roc_curve(Y_CV, nb.predict_proba(TFIDF_CV_CSR)[:,1])
24     AUC_CV.append(auc(a_fpr_cv, a_tpr_cv))
25
26
27 for av in tqdm(alpha):
28     b = np.log10(av)
29     L_ALPHA.append(b)
30
31 # Performance of model on Train data and Test data for each hyper parameter.
32 plt.plot(L_ALPHA, AUC_TR, label='AUC_Train',color='red',linewidth=2)
33 plt.scatter(L_ALPHA, AUC_TR, label='AUC_Train',color='red',linewidth=2)
34 plt.gca()
35 plt.plot(L_ALPHA, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
36 plt.scatter(L_ALPHA, AUC_CV, label='AUC_CV',color='blue',linewidth=2)
37 plt.gca()
38 plt.legend()
39 plt.xlabel("K: hyperparameter")
40 plt.ylabel("AUC")
41 plt.title("ERROR PLOTS")
42 plt.grid()
43 plt.show()
```



Wall time: 7.12 s

In [174]:

```
1 print(AUC_TR)
2 print(AUC_CV)
```

```
[0.8209130343210337, 0.8184105792952243, 0.8144660748859134, 0.8078640480373748, 0.792222775175454, 0.7261600491848925, 0.5965382854216575, 0.5442793176813133, 0.5375339793096331]
[0.6609192836423523, 0.6643669795100033, 0.6677887752645999, 0.6705764620620442, 0.6718734284394317, 0.6480679522520089, 0.5730040988684286, 0.5428300498437408, 0.540232275128462]
```

OBSERVATION:

1. Optimal alpha has been choose based on the alpha which maximize the AUC CV score.
2. In this case $\alpha(\text{alpha})=0.1$ is having the highest AUC CV score.
3. Hence $\alpha(\text{alpha})=0.1$ is choosen as the optimal alpha value.

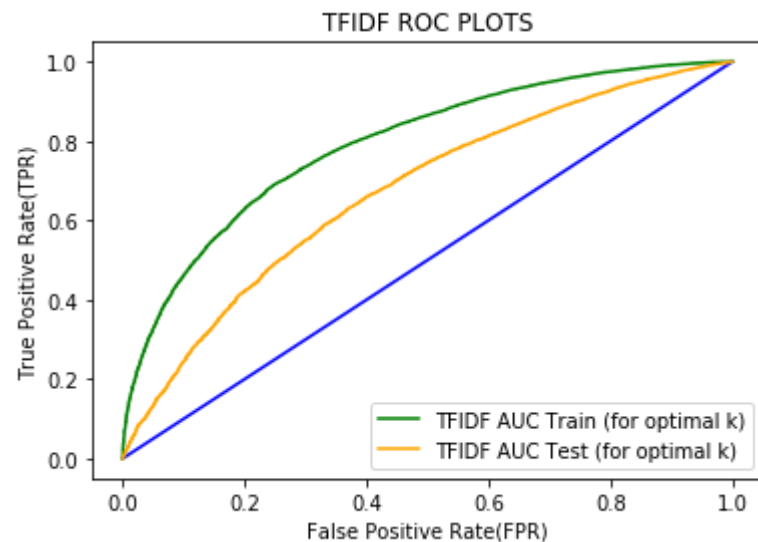
```
In [175]: 1 A_OPT=0.1
2 TFIDF_opt = MultinomialNB(alpha=A_OPT,fit_prior=True,class_prior=[0.5,0.5])
3 TFIDF_opt.fit(TFIDF_Train, Y_Train)
4 pred = TFIDF_opt.predict(TFIDF_Test)
5 acc = accuracy_score(Y_Test, pred, normalize=True) * float(100)
6 print('\nTest accuracy for alpha(a) = {0} is {1}%'.format(A_OPT,acc))
7
8 a_fpr_train, a_tpr_train, thresholds = roc_curve(Y_Train, TFIDF_opt.predict_proba(TFIDF_TR_CSR)[:,1])
9 a_fpr_Test, a_tpr_Test, thresholds = roc_curve(Y_Test, TFIDF_opt.predict_proba(TFIDF_TS_CSR)[:,1])
```

Test accuracy for alpha(a) = 0.1 is 67.11139465216908%

TFIDF ROC PLOT

In [176]:

```
1 %%time
2
3 #https://stackoverflow.com/questions/52910061/implementing-roc-curves-for-k-nn-machine-learning-algorithm-using-python-and-sci
4
5 plt.plot([0,1],[0,1], 'k-', color='blue')
6 plt.plot(a_fpr_train, a_tpr_train, label="TFIDF AUC Train (for optimal k)", color='green')
7 plt.plot(a_fpr_Test, a_tpr_Test, label="TFIDF AUC Test (for optimal k)", color='orange')
8 plt.legend()
9 plt.ylabel("True Positive Rate(TPR)")
10 plt.xlabel("False Positive Rate(FPR)")
11 plt.title("TFIDF ROC PLOTS")
12 plt.show()
13 print("-"*120)
14 print("AUC Train (for optimal k) =", auc(a_fpr_train, a_tpr_train))
15 print("AUC Test (for optimal k) =", auc(a_fpr_Test, a_tpr_Test))
16 TFIDF_AOPT=A_OPT
17 TFIDF_AUC=round(auc(a_fpr_Test, a_tpr_Test)*100)
18 pred3 = TFIDF_opt.predict(TFIDF_Train)
19 pred4 = TFIDF_opt.predict(TFIDF_Test)
20
```



AUC Train (for optimal k) = 0.792222775175454

AUC Test (for optimal k) = 0.6727470054018932

Wall time: 1.28 s

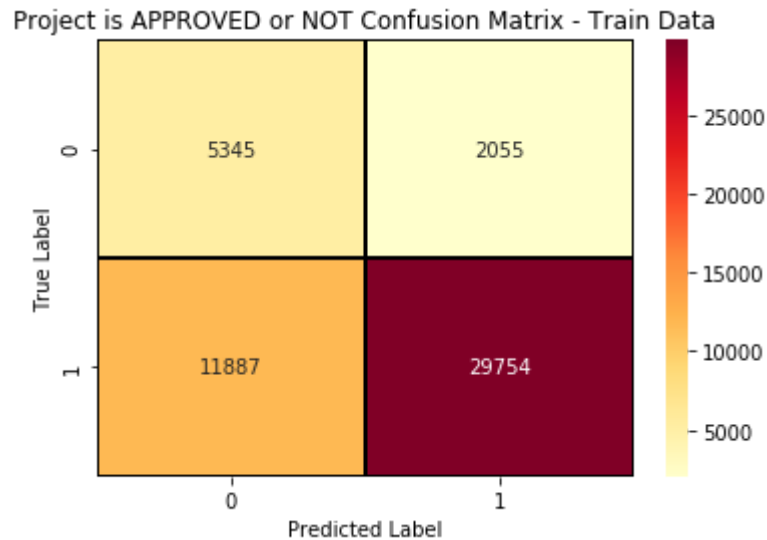
TFIDF CONFUSION MATRIX

In [177]:

```
1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 %matplotlib inline
5 from sklearn.metrics import confusion_matrix
6 Train = confusion_matrix(Y_Train, pred3)
7 sns.heatmap(Train,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
8 plt.ylabel('True Label')
9 plt.xlabel('Predicted Label')
10 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

Wall time: 388 ms

Out[177]: Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



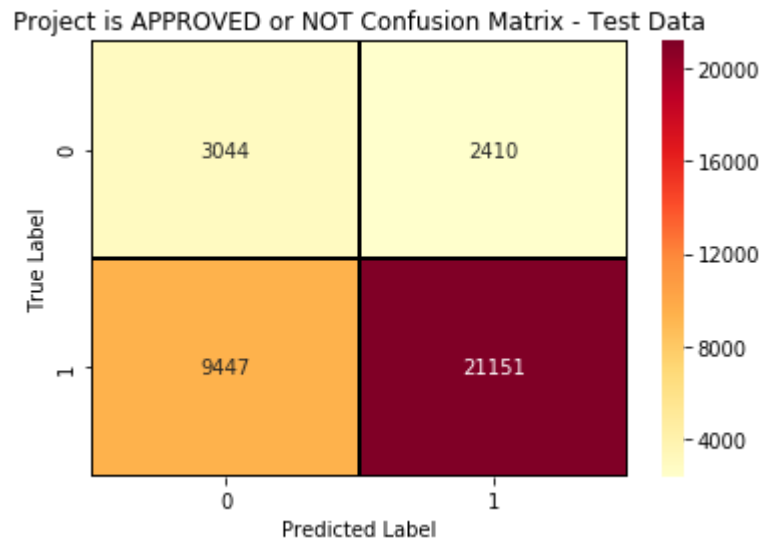
```

In [178]: 1 %%time
2 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
3 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
4 Test = confusion_matrix(Y_Test, pred4)
5 sns.heatmap(Test,annot=True,cbar=True,fmt='d',cmap='YlOrRd',linewidths=1,linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')

```

Wall time: 156 ms

Out[178]: Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')



2.4.2.1 Top 10 important features of positive class from SET 2

```
In [179]: 1 T_NB = MultinomialNB(alpha = A_OPT, class_prior=[0.5,0.5])
2 T_NB.fit(TFIDF_TR_CSR, Y_Train)
3
4 TFIDF_F = []
5
6 for TFIDF in vectorizer_cat.get_feature_names() :
7     TFIDF_F.append(TFIDF)
8 for TFIDF1 in vectorizer_sub_cat.get_feature_names() :
9     TFIDF_F.append(TFIDF1)
10 for TFIDF2 in vectorizer_school.get_feature_names() :
11     TFIDF_F.append(TFIDF2)
12 for TFIDF3 in vectorizer_grade.get_feature_names() :
13     TFIDF_F.append(TFIDF3)
14 for TFIDF4 in vectorizer_prefix.get_feature_names() :
15     TFIDF_F.append(TFIDF4)
16 for TFIDF5 in vectorizer_essays_tfidf.get_feature_names() :
17     TFIDF_F.append(TFIDF5)
18 for TFIDF6 in vectorizer_titles_tfidf.get_feature_names() :
19     TFIDF_F.append(TFIDF6)
20 TFIDF_F.append("price")
21 TFIDF_F.append("quantity")
22 TFIDF_F.append("prev_proposed_projects")
23 TFIDF_F.append("title_word_count")
24 TFIDF_F.append("essay_word_count")
25 len(TFIDF_F)
```

Out[179]: 14285

```
In [180]: 1 POS_F_TFIDF = T_NB.feature_log_prob_[1, :].argsort()[::-1]
2 for i in POS_F_TFIDF[:10]:
3     print(TFIDF_F[i])
```

```
essay_word_count
title_word_count
quantity
prev_proposed_projects
mrs
Literacy_Language
grades_prek_2
Math_Science
ms
grades_3_5
```

2.4.2.2 Top 10 important features of negative class from SET 2

```
In [181]: 1 NEG_F_TFIDF = T_NB.feature_log_prob_[0, :].argsort()[::-1]
2         for i in NEG_F_TFIDF[:10]:
3             print(TFIDF_F[i])
```

```
essay_word_count
quantity
title_word_count
prev_proposed_projects
mrs
Literacy_Language
Math_Science
grades_prek_2
ms
grades_3_5
```

3. Conclusions

```
In [182]: 1 # Please compare all your models using Prettytable Library
2
3 pt = PrettyTable()
4 pt.field_names= ("Vectorizer", "Model", "HyperParameter", "AUC")
5 pt.add_row(["BOW", "Naive Bayes", BOW_AOPT, BOW_AUC])
6 pt.add_row(["Tf-Idf", "Naive Bayes", TFIDF_AOPT, TFIDF_AUC])
7 print(pt)
```

Vectorizer	Model	HyperParameter	AUC
BOW	Naive Bayes	0.1	70.0
Tf-Idf	Naive Bayes	0.1	67.0

SUMMARY:

1. Compare to KNN "Naive Bayes" is giving better AUC value and good interpretable.
2. Also latency time is very low when compare to KNN.
3. Feature importance can be find very easily from the given data.
4. Alpha value needs to be choose correctly else there is a chance for high variance and low bias model(i.e.,Overfitting).

