

DonorsChoose

1. Importing All necessary Libs to Work

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
```

```
/anaconda3/lib/python3.7/site-packages/smart_open/ssh.py:34: UserWarning: paramiko missing, opening
SSH/SCP/SFTP paths will be disabled. `pip install paramiko` to suppress
  warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip install
paramiko` to suppress')
```

2. Reading Data

In [2]:

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

In [3]:

```
data = project_data[['id', 'teacher_prefix', 'school_state', 'project_grade_category',
                    'project_subject_categories', 'project_subject_subcategories', 'project_title',
                    'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4',
                    'teacher_number_of_previously_posted_projects', 'project_is_approved']].copy()
```

```
#data = resource_data[['quantity', 'price']].copy()
```

In [4]:

```
# I have to add to data frames to spilt it properly so i seached as below
# gSearch Key: df add column from other df
# https://stackoverflow.com/a/20603020/6000190
data = data.join(resource_data[['quantity', 'price']])
```

In [5]:

```
# Lets preprocess a little bit.. like creating one eassy and removing 4 parts of it.
# Dropping ID too.. Cause We might dont need id for anything.

data["essay"] = data["project_essay_1"].map(str) + \
    data["project_essay_2"].map(str) + \
    data["project_essay_3"].map(str) + \
    data["project_essay_4"].map(str)

data = data.drop(['id', 'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4'], axis=1)
```

In [6]:

```
data.head(2)
```

Out[6]:

	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories	project_title	teacher
0	Mrs.	IN	Grades PreK-2	Literacy & Language	ESL, Literacy	Educational Support for English Learners at Home	
1	Mr.	FL	Grades 6-8	History & Civics, Health & Sports	Civics & Government, Team Sports	Wanted: Projector for Hungry Learners	

In [7]:

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

```
Number of data points in train data (50000, 17)
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories'
 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
 'project_essay_4' 'project_resource_summary'
 'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

3. Text preprocessing

1.3.1 Essay Text

In [8]:

```
# 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 [9]:

```
sent = decontracted(data['essay'].values[100])
print(sent)
print("="*50)
```

I teach in a dual immersion 4th grade classroom. We teach 50% of the day in English and 50% in Spanish. My classroom is the English model for two classrooms of 30 students. Half of the 4th grade students at my school come to me to learn science, writing, and math in English in the morning, and the other 1/2 come to me in the afternoon to learn the same subjects. Most of my students are English Learners; however, many come to this school speaking English only. This school is a Title I school and is located in a high poverty area where many of the families are farm workers.

\r\n\r\nWe continually discuss how to implement the 4 Cs into our lessons: collaboration, communication, critical thinking, and creativity. We also never forget about the 5th and most important \nC,\" which is CARING! I have a great bunch of students who are enthusiastic about learning and reading. Most importantly, my students are so grateful for any help that we may receive. I really have a great group of kids who are so sweet. My students are becoming better researchers, project builders and writers each day. Currently, we do not have the ability for students to print and publish their work. All of my students have small student laptops, but they do not have the ability to print and our school printer is rather far away and is overworked! We currently have plenty of paper but nothing to do with it! I would like my students to freely print their work in order to display and share with other students and their families. We need printers, ink and computers that can connect to the printer. \r\nI would like my students to have a printing station within the classroom so that they can create colorful displays and start to love to use the resources available to them. This is a 21st-century skill and I would like my students to be able to learn a skill that is needed in the real world!nannan

=====

In [10]:

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

I teach in a dual immersion 4th grade classroom. We teach 50% of the day in English and 50% in Spanish. My classroom is the English model for two classrooms of 30 students. Half of the 4th grade students at my school come to me to learn science, writing, and math in English in the morning, and the other 1/2 come to me in the afternoon to learn the same subjects. Most of my students are English Learners; however, many come to this school speaking English only. This school is a Title I school and is located in a high poverty area where many of the families are farm workers. We continually discuss how to implement the 4 Cs into our lessons: collaboration, communication, critical thinking, and creativity. We also never forget about the 5th and most important C, which is CARING! I have a great bunch of students who are enthusiastic about learning and reading. Most importantly, my students are so grateful for any help that we may receive. I really have a great group of kids who are so sweet. My students are becoming better researchers, project builders and writers each day. Currently, we do not have the ability for students to print and publish their work. All of my students have small student laptops, but they do not have the ability to print and our school printer is rather far away and is overworked! We currently have plenty of paper but nothing to do with it! I would like my students to freely print their work in order to display and share with other students and their families. We need printers, ink and computers that can connect to the printer. I would like my students to have a printing station within the classroom so that they can create colorful displays and start to love to use the resources available to them. This is a 21st-century skill and I would like my students to be able to learn a skill that is needed in the real world!nannan

In [11]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

I teach in a dual immersion 4th grade classroom We teach 50 of the day in English and 50 in Spanish My classroom is the English model for two classrooms of 30 students Half of the 4th grade students at my school come to me to learn science writing and math in English in the morning and the other 12 come to me in the afternoon to learn the same subjects Most of my students are English Learners however many come to this school speaking English only This school is a Title I school and is located in a high poverty area where many of the families are farm workers We continually discuss how to implement the 4 Cs into our lessons collaboration communication critical thinking and creativity We also never forget about the 5th and most important C which is CARING I have a great bunch of students who are enthusiastic about learning and reading Most importantly my students are so grateful for any help that we may receive I really have a great group of kids who are so sweet My students are becoming better researchers project builders and writers each day Currently we do not have the ability for students to print and publish their work All of my students have small student laptops but they do not have the ability to print and our school printer is rather far away and is overworked We currently have plenty of paper but nothing to do with it I would like my students to freely print their work in order to display and share with other students and their families We need printers ink and computers that can connect to the printer I would like my students to have a printing station within the classroom so that they can create colorful displays and start to love to use the resources available to them This is a 21st century skill and I would like my students to be able to learn a skill that is needed in the real world nannan

In [12]:

```
# 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', "do \
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"]
```

In [13]:

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

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

```
100%|██████████| 50000/50000 [00:35<00:00, 1400.46it/s]
```

In [14]:

```
# after preprocessing
preprocessed_essays[191]
```

Out[14]:

```
'our school large elementary school warm inviting most students come variety different economic ba
ckgrounds 70 students receiving free reduced price lunches harvesters food bank provides backpack
snacks one hundred students many students arrive school little no school supplies my students exci
ted learning i commitment whatever takes ensure talents developed atmosphere high expectations per
sonal support students school responsible respectful always ready learn flexible classrooms give s
tudents choice kind learning space works best help work collaboratively communicate engage
critical thinking the couch table set give students opportunities i want classroom represent real
world seating arrangements classrooms 70 years ago students need seating confront gives
alternative desk i believe providing flexible soft alternative seating students would able move re
lease energy happier comfortable work nannan'
```

In [15]:

```
data['essay']=preprocessed_essays
```

1.3.2 Project title Text

In [16]:

```
# similarly you can preprocess the titles also
# Processing steps for project title
# 1. Clean pharase
# 2. Remove String patterns
# 3. Remove Special charcter
# 4. Remove Stop words

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

```
100%|██████████| 50000/50000 [00:01<00:00, 31309.66it/s]
```

In [17]:

```
# after preprocessing
preprocessed_project_title[0]
```

Out[17]:

```
'educational support english learners home'
```

In [18]:

```
data['project_title'][0]
```

Out[18]:

```
Out[18]:
```

```
'Educational Support for English Learners at Home'
```

```
In [19]:
```

```
# for i in tqdm(range(len(preprocessed_project_title))):
#     data['project_title'][i] = preprocessed_project_title[i]

data['project_title']=preprocessed_project_title
```

Preprocessing for project_grade_category data.

```
In [20]:
```

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

```
100%|██████████| 50000/50000 [00:01<00:00, 39564.72it/s]
```

```
In [21]:
```

```
print(data['project_grade_category'][0])
print(preprocessed_project_grade_category[0])
```

```
Grades PreK-2
grades prek 2
```

```
In [22]:
```

```
data['project_grade_category']=preprocessed_project_grade_category
```

Preprocessing for project_subject_categories data.

```
In [23]:
```

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

```
100%|██████████| 50000/50000 [00:01<00:00, 38009.69it/s]
```

```
In [24]:
```

```
print(data['project_subject_categories'][0])
print(preprocessed_project_sub_category[0])
```

```
Literacy & Language
literacy language
```

In [25]:

```
data['project_subject_categories']=preprocessed_project_sub_category
```

Preprocessing for project_subject_subcategories data.

In [26]:

```
# preprocessed_project_subject_sub_category = []
# # tqdm is for printing the status bar
# for sentence in tqdm(data['project_subject_subcategories'].values):
#     sent = decontracted(sentence)
#     sent = sent.replace('\r', ' ')
#     sent = sent.replace('\n', ' ')
#     sent = sent.replace('\n', ' ')
#     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
#     # https://gist.github.com/sebleier/554280
#     sent = ' '.join(e for e in sent.split() if e not in stopwords)
#     preprocessed_project_subject_sub_category.append(sent.lower().strip())
sub_categories = list(project_data['project_subject_subcategories'].values)
sub_cat_list = []
for i in tqdm(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 placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp +=j.strip()+" #" + abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&','_')
        sub_cat_list.append(temp.strip())
```

100%|██████████| 50000/50000 [00:00<00:00, 289802.78it/s]

In [27]:

```
project_data['project_subject_subcategories'] = sub_cat_list
```

In [28]:

```
print(data['project_subject_subcategories'][0])
print(sub_cat_list[0])
```

ESL, Literacy
ESL Literacy

In [29]:

```
data['project_subject_subcategories'] = sub_cat_list
# data.drop(['project_subject_subcategories'], axis=1, inplace=True)
data.head(2)
```

Out[29]:

	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories	project_title	teacher
0	Mrs.	IN	grades prek 2	literacy language	ESL Literacy	educational support english learners home	

1	Mrs.	IN	grades 3-5	literacy language	ESL Literacy	wanted projector	
---	------	----	------------	-------------------	--------------	------------------	--

In [30]:

```
# data['project_subject_subcategories']=preprocessed_project_subject_sub_category
```

4. Spliting Data

In [31]:

```
data.head(1)
```

Out[31]:

	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories	project_title	teacher
0	Mrs.	IN	grades prek 2	literacy language	ESL Literacy	educational support english learners home	

In [32]:

```
# Just Checking for Nan values .. before splitting them up
# Google search :- check which column has nan pandas
# https://stackoverflow.com/questions/36226083/how-to-find-which-columns-contain-any-nan-value-in-pandas-dataframe-python/36226137

data.isna().any()
```

Out[32]:

```
teacher_prefix      True
school_state        False
project_grade_category  False
project_subject_categories  False
project_subject_subcategories  False
project_title       False
teacher_number_of_previously_posted_projects  False
project_is_approved  False
quantity            False
price               False
essay               False
dtype: bool
```

In [33]:

```
data.fillna("", inplace=True)
print("Data Cleaned from nan")
```

Data Cleaned from nan

In [34]:

```
data.isna().any()
```

Out[34]:

```
teacher_prefix      False
school_state        False
project_grade_category  False
project_subject_categories  False
project_subject_subcategories  False
project_title       False
```



```

teacher_number_of_previously_posted_projects    False
project_is_approved                             False
quantity                                         False
price                                             False
essay                                             False
dtype: bool

```

In [35]:

```

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

```

Out[35]:

	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories	project_title	teacher
0	Mrs.	IN	grades prek 2	literacy language	ESL Literacy	educational support english learners home	

In [36]:

```
X = data
```

In [37]:

```

# train test split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)

```

5. Vectorizing Data

In [38]:

```

# For Selecting best feature.. lets create an empty array and lets add.

feature_names=[]

```

* Vectorizing Essay

In [39]:

```

from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4)) # , max_features=5000
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)
X_test_essay_bow = vectorizer.transform(X_test['essay'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("="*50)

feature_names.extend(vectorizer.get_feature_names())

```

```

After vectorizations
(22445, 71584) (22445,)
(11055, 71584) (11055,)
(16500, 71584) (16500,)

```

```
=====
```

In [40]:

```
print(X_train_essay_bow.shape)
# print(X_train_essay_bow[0:3,0:3])
print(X_train_essay_bow.toarray()[5,5])
print(X_cv_essay_bow.toarray()[5,5])
print(X_test_essay_bow.toarray()[5,5])
```

```
(22445, 71584)
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
```

* Vectorizing Title

In [41]:

```
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4)) # , max_features=5000
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['project_title'].values)
X_cv_title_bow = vectorizer.transform(X_cv['project_title'].values)
X_test_title_bow = vectorizer.transform(X_test['project_title'].values)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
print("="*50)

feature_names.extend(vectorizer.get_feature_names())
```

```
After vectorizations
(22445, 1993) (22445,)
(11055, 1993) (11055,)
(16500, 1993) (16500,)
=====
```

In [42]:

```
print(X_train_title_bow.shape)
# print(X_train_essay_bow[0:3,0:3])
print(X_train_title_bow.toarray()[5,5])
print(X_cv_title_bow.toarray()[5,5])
print(X_test_title_bow.toarray()[5,5])
```

```
(22445, 1993)
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
[[0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]
 [0 0 0 0 0]]
```

```
[0 0 0 0 0]]
[[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]]
```

* Vectorizing State

In [43]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train data

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

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

feature_names.extend(vectorizer.get_feature_names())
```

```
After vectorizations
(22445, 51) (22445,)
(11055, 51) (11055,)
(16500, 51) (16500,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'k',
s', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm',
'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv',
', 'wy']
=====
```

In [44]:

```
print(X_train_state_ohe.shape)
# print(X_train_essay_bow[0:3,0:3])
print(X_train_state_ohe.toarray()[0:5,:5])
print(X_cv_state_ohe.toarray()[0:5,:5])
print(X_test_state_ohe.toarray()[0:5,:5])
```

```
(22445, 51)
[[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]]
[[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 1]
[0 0 0 0 0]
[0 0 0 0 0]]
[[0 0 0 0 0]
[0 1 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]
[0 0 0 0 0]]
```

* Vectorizing Teacher_prefix

In [45]:

```
vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values) # fit has to happen only on train data
```

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

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

feature_names.extend(vectorizer.get_feature_names())
```

After vectorizations

```
(22445, 5) (22445,)
(11055, 5) (11055,)
(16500, 5) (16500,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
```

In [46]:

```
print(X_train_teacher_ohe.shape)
# print(X_train_essay_bow[0:3,0:3])
print(X_train_teacher_ohe.toarray()[:5,:5])
print(X_cv_teacher_ohe.toarray()[:5,:5])
print(X_test_teacher_ohe.toarray()[:5,:5])
```

```
(22445, 5)
[[0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 0 1 0]]
[[0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 1 0 0]
 [0 0 0 1 0]]
[[0 0 1 0 0]
 [0 0 1 0 0]
 [0 0 0 1 0]
 [0 0 0 1 0]
 [0 0 0 1 0]]
```

* Vectorizing Project Grade Category

In [47]:

```
my_counter_project_grade = Counter()
for word in X_train['project_grade_category'].values:
    my_counter_project_grade.update(word.split(","))

cat_dict_procat = dict(my_counter_project_grade)
sorted_procat = dict(sorted(cat_dict_procat.items(), key=lambda kv: kv[1]))

# we use count vectorizer to convert the values into one hot encoded features
vectorizer_teacher = CountVectorizer(vocabulary=list(sorted_procat.keys()), lowercase=False,
binary=True)
vectorizer_teacher.fit(X_train['project_grade_category'].values )

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

```

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer_teacher.get_feature_names())
print("="*100)

feature_names.extend(vectorizer.get_feature_names())

```

```

After vectorizations
(22445, 4) (22445,)
(11055, 4) (11055,)
(16500, 4) (16500,)
['grades 9 12', 'grades 6 8', 'grades 3 5', 'grades prek 2']
=====

```

* Vectorizing Project Subject Subcategories

In [48]:

```

vectorizer = CountVectorizer(min_df=10)
vectorizer.fit(X_train['project_subject_subcategories'].values) # fit has to happen only on train
data

# we use the fitted CountVectorizer to convert the text to vector
X_train_sub_cat = vectorizer.transform(X_train['project_subject_subcategories'].values)
X_cv_sub_cat = vectorizer.transform(X_cv['project_subject_subcategories'].values)
X_test_sub_cat = vectorizer.transform(X_test['project_subject_subcategories'].values)

print("After vectorizations")
print(X_train_sub_cat.shape, y_train.shape)
print(X_cv_sub_cat.shape, y_cv.shape)
print(X_test_sub_cat.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

feature_names.extend(vectorizer.get_feature_names())

```

```

After vectorizations
(22445, 30) (22445,)
(11055, 30) (11055,)
(16500, 30) (16500,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience',
'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness',
'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'm
athematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socia
lsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
=====

```

* Vectorizing Teacher Number Of Previously Posted Projects

In [49]:

```

# vectorizer = CountVectorizer()
# vectorizer.fit(X_train['teacher_number_of_previously_posted_projects'].values) # fit has to happ
en only on train data

from sklearn.preprocessing import Normalizer
normalizer = Normalizer()

normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_train_prPos_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].
values.reshape(-1,1))
X_cv_prPos_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values
.reshape(-1,1))
X_test_prPos_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].va
lues.reshape(-1,1))

```

```

print("After vectorizations")
print(X_train_prPos_norm.shape, y_train.shape)
print(X_cv_prPos_norm.shape, y_cv.shape)
print(X_test_prPos_norm.shape, y_test.shape)
print(X_cv_prPos_norm)
print("="*100)

feature_names.extend(X_train['teacher_number_of_previously_posted_projects'])

```

After vectorizations

```

(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
[[1.]
 [1.]
 [1.]
 ...
 [0.]
 [1.]
 [1.]]
=====

```

* Vectorizing Quantity

In [50]:

```

# vectorizer = CountVectorizer()
# vectorizer.fit(X_train['teacher_number_of_previously_posted_projects'].values) # fit has to happen only on train data

normalizer = Normalizer()

normalizer.fit(X_train['quantity'].values.reshape(-1,1))

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

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

feature_names.extend(X_train['quantity'])

```

After vectorizations

```

(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
[[1.]
 [1.]
 [1.]
 ...
 [1.]
 [1.]
 [1.]]
=====

```

* Standardising Price

In [51]:

```

# vectorizer = CountVectorizer()
# vectorizer.fit(X_train['teacher_number_of_previously_posted_projects'].values) # fit has to happen only on train data

```

```

normalizer = Normalizer()

normalizer.fit(X_train['price'].values.reshape(-1,1))

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

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

feature_names.extend(X_train['price'])

```

After vectorizations

```

(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
[[1.]
 [1.]
 [1.]
 ...
 [1.]
 [1.]
 [1.]]

```

In []:

In []:

In [52]:

```
X_train_essay_bow.shape
```

Out[52]:

```
(22445, 71584)
```

In [53]:

```
X_train_price_norm.shape
```

Out[53]:

```
(22445, 1)
```

In [54]:

```
X_cv_essay_bow.shape
```

Out[54]:

```
(11055, 71584)
```

In [55]:

```
X_cv_price_norm.shape
```

Out[55]:

```
(11055, 1)
```

In [56]:

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

```
from scipy.sparse import hstack
X_tr = hstack(
    (X_train_essay_bow,
     X_train_title_bow,
     X_train_state_ohe,
     X_train_teacher_ohe,
     X_train_sub_cat,
     X_train_prPos_norm,
     X_train_grade_ohe,
     X_train_quantity_norm,
     X_train_price_norm)).tocsr()
```

```
X_cr = hstack(
    (X_cv_essay_bow,
     X_cv_title_bow,
     X_cv_state_ohe,
     X_cv_teacher_ohe,
     X_cv_sub_cat,
     X_cv_prPos_norm,
     X_cv_grade_ohe,
     X_cv_quantity_norm,
     X_cv_price_norm)).tocsr()
```

```
X_te = hstack(
    (X_test_essay_bow,
     X_test_title_bow,
     X_test_state_ohe,
     X_test_teacher_ohe,
     X_test_sub_cat,
     X_test_prPos_norm,
     X_test_grade_ohe,
     X_test_quantity_norm,
     X_test_price_norm)).tocsr()
```

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

```
Final Data matrix
(22445, 73670) (22445,)
(11055, 73670) (11055,)
(16500, 73670) (16500,)
```

6. Let's Build a model with above data

Applying Naive Bayes on BOW, SET 1

In [57]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.naive_bayes import MultinomialNB
```



```

neigh = MultinomialNB()
parameters = {'alpha':[0.00001, 0.0001, 0.01, 0.1, 0.5, 1, 3, 10]}
clf = GridSearchCV(neigh, parameters, cv=3, scoring='roc_auc')
clf.fit(X_tr, y_train)

train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

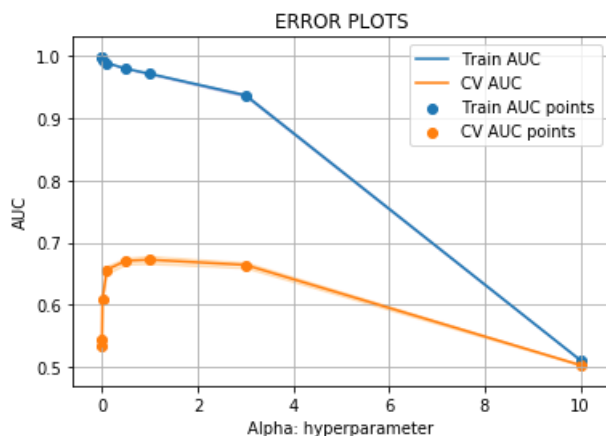
plt.plot(parameters['alpha'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['alpha'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.2,color='darkblue')

plt.plot(parameters['alpha'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['alpha'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color=
'darkorange')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



- As per the error Plot The minimum gap represents the best alpha value.
- And in this plot 71 is having less space between Train and Test AUC Graphs

In [58]:

```

def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    tive class
    # not the predicted outputs

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

    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred

```

In [59]:

```
X_tr.shape
```

Out[59]:

(22445, 73670)

In [60]:

X_te.shape

Out[60]:

(16500, 73670)

In [61]:

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

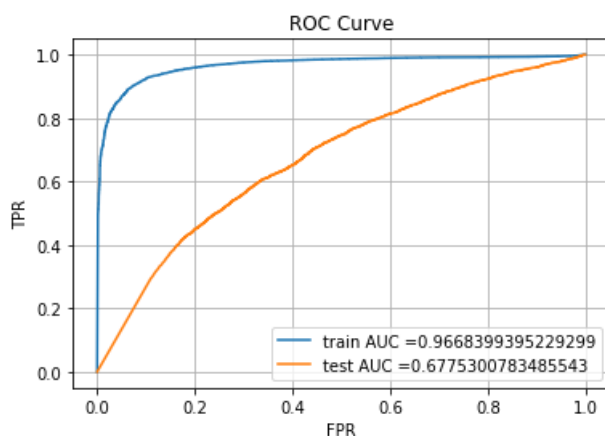
best_alpha_m1 = 0.1
# for i in tqdm(parameters):
naiveBayesClf_set1 = MultinomialNB(alpha=best_alpha_m1, class_prior=[0.5, 0.5] )
naiveBayesClf_set1.fit(X_tr, y_train)
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs

y_train_pred = batch_predict(naiveBayesClf_set1, X_tr)
y_test_pred = batch_predict(naiveBayesClf_set1, X_te)

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

m1_Auc = str(auc(train_fpr, train_tpr))

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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")
plt.grid()
plt.show()
```



- Train AUC values is 0.97
- Test AUC value is 0.64
- As i tried with Diffrent alpha Values it representing with diffrent values

In []:

In [62]:

```
def predict(proba, threshold, fpr, tpr):

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

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

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

In [63]:

```
from sklearn.metrics import confusion_matrix

print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915340996 for threshold 0.0
[[ 1731  1732]
 [  256 18726]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.24999984572938835 for threshold 0.0
[[  283  2263]
 [  583 13371]]
```

In [72]:

```
# https://stackoverflow.com/a/42265865/6000190
def plotConfusionMatrix(cm):
    # cm_normalized = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
    # cm_normalized = cm.astype('int') / 100
    sns.set(font_scale=1)#for label size
    # sns.heatmap(cm, annot=True,annot_kws={"size": 12})# font size
    sns.heatmap(cm, annot = True,annot_kws={"size": 16}, fmt='d')
```

In [73]:

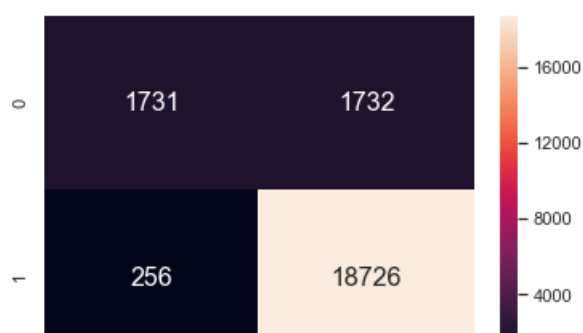
```
# Normalize the confusion matrix by row (i.e by the number of samples
# in each class)
# https://scikit-learn.org/0.16/auto_examples/model_selection/plot_confusion_matrix.html

cm_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))

plotConfusionMatrix(cm_train)

# sns.set(font_scale=1.4)#for label size
# sns.heatmap(cm_train, annot=True,annot_kws={"size": 16})# font size
```

```
the maximum value of tpr*(1-fpr) 0.24999997915340996 for threshold 0.0
```



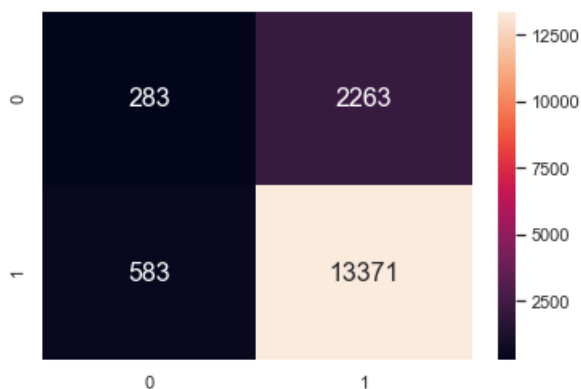


- TP = 1734
- TN = 1729
- FP = 203
- FN = 18779
- FN has a heigher value but TP Should have higher value as well.

In [74]:

```
cm_test = confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
plotConfusionMatrix(cm_test)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999984572938835 for threshold 0.0



- TP = 352
- TN = 2194
- FP = 668
- FN = 13286
- TP has higher values comparing to FN. and TN has heigher values which is not good model performance.

In []:

In [75]:

```
max_ind_pos=np.argsort((naiveBayesClf_set1.feature_log_prob_)[1])[::-1][0:10]
min_ind_pos=np.argsort((naiveBayesClf_set1.feature_log_prob_)[1])[::-1][0:10]
top_pos=np.take(feature_names,max_ind_pos)
last_pos=np.take(feature_names,min_ind_pos)

max_ind_neg=np.argsort((naiveBayesClf_set1.feature_log_prob_)[0])[::-1][0:10]
min_ind_neg=np.argsort((naiveBayesClf_set1.feature_log_prob_)[0])[::-1][0:10]
top_neg=np.take(feature_names,max_ind_neg)
last_neg=np.take(feature_names,min_ind_neg)

print("")
print("Positive hight cofe_ features: ")
print("")
print(top_pos)
print("="*50)
print("Positive low cofe_ features: ")
print("")
print(last_pos)

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

print("")
```

```

print("Negative hight cofe_ features: ")
print("")
print(top_neg)
print("="*50)
print("Negative low cofe_ features: ")
print("")
print(last_neg)

```

Positive hight cofe_ features:

```

['specialneeds' 'warmth' 'visualarts' 'teamsports' 'dr'
 'autism cerebral palsy' 'autism cerebral' 'ready play'
 'allowing students learn' 'the materials allow']
=====

```

Positive low cofe_ features:

```

['students' 'school' 'my' 'learning' 'classroom' 'the' 'they' 'not'
 'my students' 'learn']
=====
=====
=====

```

Negative hight cofe_ features:

```

['wobbly stools' 'space sit' 'but importantly' 'trying become'
 'multiple jobs make ends' 'multiple jobs make' 'multimedia projects'
 'multi step' 'busy bees' 'spanish hindi']
=====

```

Negative low cofe_ features:

```

['students' 'school' 'learning' 'my' 'classroom' 'learn' 'not' 'help'
 'they' 'the']

```

In []:

In []:

Applying Naive Bayes on TFIDF, SET 2 With Best 20 Features

In [76]:

```

# Making DataMatrix with eassay and title with TFIDF

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizerTfidf = TfidfVectorizer(min_df=10)

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizerTfidf.fit_transform(X_train['essay'].values)
X_cv_essay_tfidf = vectorizerTfidf.transform(X_cv['essay'].values)
X_test_essay_tfidf = vectorizerTfidf.transform(X_test['essay'].values)

print("Shape of matrix after one hot encodig ",X_train_essay_tfidf.shape)

```

Shape of matrix after one hot encodig (22445, 8882)

Similarly you can vectorize for title also with TFIDF

In [77]:

```

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizerTfidf.fit_transform(X_train['project_title'].values)
X_cv_title_tfidf = vectorizerTfidf.transform(X_cv['project_title'].values)

```

```
X_test_title_tfidf = vectorizerTfidf.transform(X_test['project_title'].values)
```

```
print("Shape of matrix after one hot encodig ",X_train_title_tfidf.shape)
```

Shape of matrix after one hot encodig (22445, 1225)

In [78]:

```
X_cv_essay_tfidf.shape
```

Out[78]:

(11055, 8882)

In [79]:

```
X_test_title_tfidf.shape
```

Out[79]:

(16500, 1225)

In [80]:

```
X_train_essay_tfidf.shape
```

Out[80]:

(22445, 8882)

In [81]:

```
X_cv_price_norm.shape
```

Out[81]:

(11055, 1)

In [82]:

```
# combining all features
```

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

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

```
X_tr = hstack(  
    (X_train_essay_tfidf,  
     X_train_title_tfidf,  
     X_train_state_ohe,  
     X_train_teacher_ohe,  
     X_train_sub_cat,  
     X_train_prPos_norm,  
     X_train_grade_ohe,  
     X_train_quantity_norm,  
     X_train_price_norm)).tocsr()
```

```
X_cr = hstack(  
    (X_cv_essay_tfidf,  
     X_cv_title_tfidf,  
     X_cv_state_ohe,  
     X_cv_teacher_ohe,  
     X_cv_sub_cat,  
     X_cv_prPos_norm,  
     X_cv_grade_ohe,  
     X_cv_quantity_norm,  
     X_cv_price_norm)).tocsr()
```

```
X_te = hstack(
    (X_test_essay_tfidf,
     X_test_title_tfidf,
     X_test_state_ohe,
     X_test_teacher_ohe,
     X_test_sub_cat,
     X_test_prPos_norm,
     X_test_grade_ohe,
     X_test_quantity_norm,
     X_test_price_norm)).tocsr()
```

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

```
Final Data matrix
(22445, 10200) (22445,)
(11055, 10200) (11055,)
(16500, 10200) (16500,)
=====
```

In [83]:

```
# Selecting Best 20 features for this set of Train Data
# from sklearn.datasets import load_digits
# from sklearn.feature_selection import SelectKBest, chi2

# X_new_train = SelectKBest(chi2, k=20).fit_transform(X_tr, y_train)
# X_new_train.shape
```

In [84]:

```
# # New best 20 feature for x test too
# X_new_test = SelectKBest(chi2, k=20).fit_transform(X_te, y_test)
# X_new_test.shape
```

In [85]:

```
neigh = MultinomialNB()
parameters = {'alpha':[0.00001,0.00005, 0.00008, 0.0001, 0.1, 0.2,0.3,0.4,0.5, 0.6, 0.8, 1]}
clf = GridSearchCV(neigh, parameters, cv=3, scoring='roc_auc')
clf.fit(X_tr, y_train)

train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

plt.plot(parameters['alpha'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['alpha'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.2,color='darkblue')

plt.plot(parameters['alpha'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['alpha'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color=
'darkorange')

plt.scatter(parameters['alpha'], train_auc, label='Train AUC points')
plt.scatter(parameters['alpha'], cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("Alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
```

```
plt.show()
```



- As per the error Plot The minimum gap represents the best alpha value.
- And in this plot 0.5 is having less space between Train and Test AUC Graphs

```
In [86]:
```

```
X_tr.shape
```

```
Out[86]:
```

```
(22445, 10200)
```

```
In [87]:
```

```
X_te.shape
```

```
Out[87]:
```

```
(16500, 10200)
```

```
In [88]:
```

```
y_test.shape
```

```
Out[88]:
```

```
(16500,)
```

```
In [89]:
```

```
best_alpha_m2 = 0.1
# for i in tqdm(parameters):
naiveBayesClf = MultinomialNB(alpha=best_alpha_m2, class_prior=[0.5, 0.5])
naiveBayesClf.fit(X_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(naiveBayesClf, X_tr)
y_test_pred = batch_predict(naiveBayesClf, X_te)

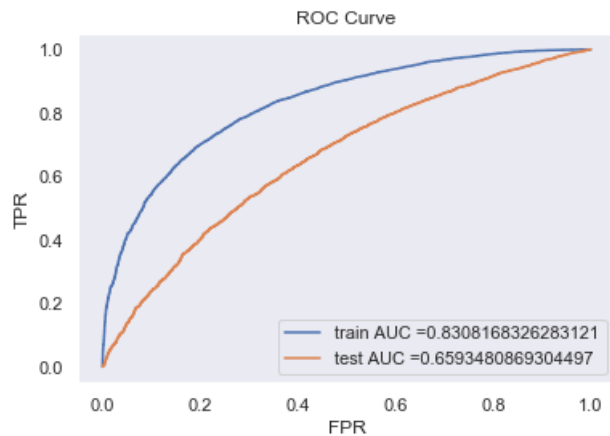
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)

m2_Auc = str(auc(train_fpr, train_tpr))

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("FPR")
plt.ylabel("TPR")
plt.title("ROC Curve")
```



```
plt.grid()
plt.show()
```



- Train AUC values is 0.77
- Test AUC value is 0.63
- Comparing to Knn Naive Bayes gives better results.

In [90]:

```
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Train confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.24999997915340996 for threshold 0.343

```
[[ 1731  1732]
 [ 1833 17149]]
```

Test confusion matrix
the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 0.408

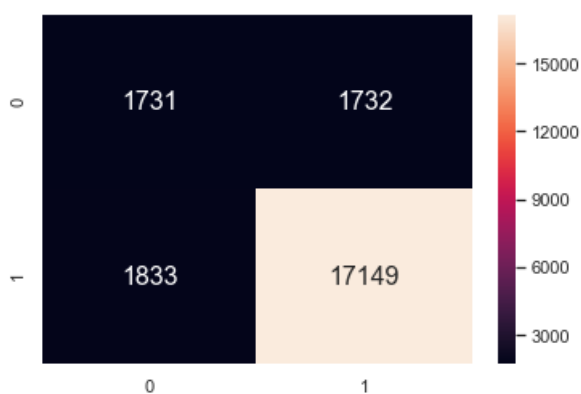
```
[[  968  1578]
 [ 2572 11382]]
```

In [91]:

```
cm_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))

plotConfusionMatrix(cm_train)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999997915340996 for threshold 0.343



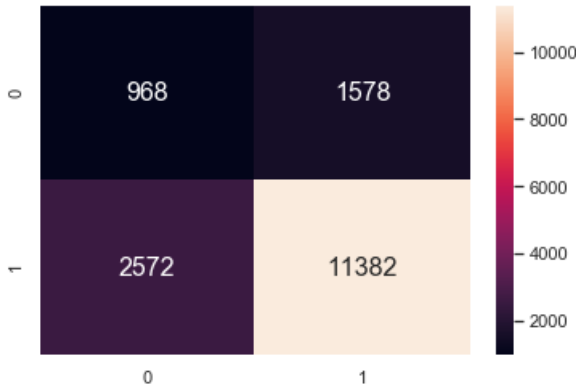
- TP = 1732
- TN = 1731
- FP = 2626
- FN = 16356
- FN has Higher value than other values which is good but with that.. TP should have higher value too.

FN has higher value then other values which is good but with that..

In [92]:

```
cm_test = confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
plotConfusionMatrix(cm_test)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.25 for threshold 0.408



- TP = 1108
- TN = 1438
- FP = 3479
- FN = 10475
- FN has Higher value than other values which is good but with that..

In [93]:

```
max_ind_pos=np.argsort((naiveBayesClf.feature_log_prob_)[1])[:,1][0:10]
min_ind_pos=np.argsort((naiveBayesClf.feature_log_prob_)[1])[:,1][0:10]
top_pos=np.take(feature_names,max_ind_pos)
last_pos=np.take(feature_names,min_ind_pos)

max_ind_neg=np.argsort((naiveBayesClf.feature_log_prob_)[0])[:,1][0:10]
min_ind_neg=np.argsort((naiveBayesClf.feature_log_prob_)[0])[:,1][0:10]
top_neg=np.take(feature_names,max_ind_neg)
last_neg=np.take(feature_names,min_ind_neg)

print("")
print("Positive high coe_ features: ")
print("")
print(top_pos)
print("="*50)
print("Positive low coe_ features: ")
print("")
print(last_pos)

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

print("")
print("Negative high coe_ features: ")
print("")
print(top_neg)
print("="*50)
print("Negative low coe_ features: ")
print("")
print(last_neg)
```

Positive high coe_ features:

```
['classroom need' 'classroom never' 'classroom needs' 'classroom needed'
 'benefit use' 'book these' 'art design' '98 students receive free' '2015'
 'bilingual households receive free']
```

=====

Positive low core_ features:

```
['classroom new' 'classroom next' 'classroom necessary' 'classroom meet'
 'classroom meet needs' 'classroom my first' 'classroom my school'
 'classroom my goal' 'classroom library needs' 'came']
```

Negative high core_ features:

```
['activity also' 'brooklyn' 'broaden' 'brings us'
 'access reading materials' 'access resources many' 'access students'
 'access technology homes' 'classroom environment the' 'bring stem']
```

Negative low core_ features:

```
['classroom new' 'classroom next' 'classroom necessary' 'classroom meet'
 'classroom meet needs' 'classroom my school' 'classroom my first'
 'classroom my goal' 'classroom library needs' 'classroom my students use']
```

In []:

In [94]:

```
# http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
```

In [95]:

```
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]

x.add_row(["BOW", "Naive Bayes", best_alpha_m1, m1_Auc])
x.add_row([" ", " ", " ", " "])
x.add_row(["TFIDF", "Naive Bayes", best_alpha_m2, m2_Auc])

print(x)
```

Vectorizer	Model	Hyper Parameter	AUC
BOW	Naive Bayes	0.1	0.9668399395229299
TFIDF	Naive Bayes	0.1	0.8308168326283121

In []:

Observations

- 1. I Have applied Naive Bayes For this Dataset.
- 1. I have taken 50k Data points for this computation
- 1. BOW vectorized data gave more good then TFIDF.
- 1. TFIDF vectorized Data set is not performing Good.
- 1. Finding top 20 features Listed Down, as they all are word.

In []:

