

# Apply Logistic Regression

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

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

In [2]:

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

In [3]:

```
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']
```

In [4]:

```
# Let's check for any "null" or "missing" values
project_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 17 columns):
Unnamed: 0                50000 non-null int64
id                        50000 non-null object
teacher_id               50000 non-null object
teacher_prefix           49998 non-null object
school_state             50000 non-null object
project_submitted_datetime 50000 non-null object
project_grade_category    50000 non-null object
project_subject_categories 50000 non-null object
project_subject_subcategories 50000 non-null object
project_title            50000 non-null object
project_essay_1          50000 non-null object
project_essay_2          50000 non-null object
project_essay_3          1685 non-null object
project_essay_4          1685 non-null object
project_resource_summary  50000 non-null object
teacher_number_of_previously_posted_projects 50000 non-null int64
project_is_approved       50000 non-null int64
dtypes: int64(3), object(14)
memory usage: 6.5+ MB
```

In [5]:

```
project_data['teacher_prefix'].isna().sum()
```

Out[5]:

2

In [6]:

```
# "teacher_prefix" seems to contain 3 "missing" values, let's use mode replacement strategy to fill
1 those missing values
project_data['teacher_prefix'].mode()
```

Out[6]:

```
0    Mrs.
dtype: object
```

In [7]:

```
# Let's replace the missing values with "Mrs." , as it is the mode of the "teacher_prefix"
project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('Mrs.')
```

In [8]:

```
project_data['teacher_prefix'].value_counts()
```

Out[8]:

```
Mrs.      26142
Ms.       17936
Mr.       4859
Teacher   1061
Dr.         2
Name: teacher_prefix, dtype: int64
```

In [9]:

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

In [10]:

```
# Let's select only the selected features or columns, dropping "project_resource_summary" as it is optional
#
project_data.drop(['id', 'teacher_id', 'project_submitted_datetime', 'project_resource_summary'], axis=1, inplace=True)
project_data.columns
```

Out[10]:

```
Index(['Unnamed: 0', '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',
       'price', 'quantity'],
      dtype='object')
```

In [11]:

```
# Data seems to be highly imbalanced since the ratio of "class 1" to "class 0" is nearly 5.5
project_data['project_is_approved'].value_counts()
```

Out[11]:

```
1    42286
0     7714
Name: project_is_approved, dtype: int64
```

In [12]:

```
number_of_approved = project_data['project_is_approved'][project_data['project_is_approved'] == 1].count()
number_of_not_approved = project_data['project_is_approved'][project_data['project_is_approved'] == 0].count()

print("Ratio of Project approved to Not approved is:", number_of_approved/number_of_not_approved)
```

Ratio of Project approved to Not approved is: 5.481721545242416

Let's first merge all the project\_essays into single columns

In [13]:

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

In [14]:

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

Out[14]:

	Unnamed: 0	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories
0	160221	Mrs.	IN	Grades PreK-2	Literacy & Language	ESL, Literacy

	Unnamed: 0	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories
1	140945	Mr.	FL	Grades 6-8	History & Civics, Health & Sports	Civics & Government, Team Sports

In [15]:

```
# Let's drop the project essay columns from the dataset now, as we have captured the essay text data into single "essay" column
project_data.drop(['project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4'], axis=1, inplace=True)
```

In [16]:

```
y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
X.head(1)
```

Out[16]:

	Unnamed: 0	teacher_prefix	school_state	project_grade_category	project_subject_categories	project_subject_subcategories
0	160221	Mrs.	IN	Grades PreK-2	Literacy & Language	ESL, Literacy

In [17]:

```
# 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)
```

## 2) Make Data Model Ready: encoding numerical, categorical features

In [18]:

```
def cleaning_text_data(list_text_feature, df, old_col_name, new_col_name):

    # remove special characters from list of strings python:
    https://stackoverflow.com/a/47301924/4084039
    # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
    # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
    # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
    feature_list = []
    for i in list_text_feature:
        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('&', '_') # we are replacing the & value into _
        feature_list.append(temp.strip())
```

```

df[new_col_name] = feature_list
df.drop([old_col_name], axis=1, inplace=True)

from collections import Counter
my_counter = Counter()
for word in df[new_col_name].values:
    my_counter.update(word.split())

feature_dict = dict(my_counter)
sorted_feature_dict = dict(sorted(feature_dict.items(), key=lambda kv: kv[1]))
return sorted_feature_dict

```

In [19]:

```

def clean_project_grade(list_text_feature, df, old_col_name, new_col_name):

    # remove special characters from list of strings python:
    https://stackoverflow.com/a/47301924/4084039
    # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
    # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
    # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
    feature_list = []
    for i in list_text_feature:
        temp = i.split(' ')
        last_dig = temp[-1].split('-')
        fin = [temp[0]]
        fin.extend(last_dig)
        feature = ' '.join(fin)
        feature_list.append(feature.strip())

    df[new_col_name] = feature_list
    df.drop([old_col_name], axis=1, inplace=True)

    from collections import Counter
    my_counter = Counter()
    for word in df[new_col_name].values:
        my_counter.update(word.split())

    feature_dict = dict(my_counter)
    sorted_feature_dict = dict(sorted(feature_dict.items(), key=lambda kv: kv[1]))
    return sorted_feature_dict

```

## 2.1) Text Preprocessing: project\_subject\_categories

In [20]:

```

x_train_sorted_category_dict = cleaning_text_data(X_train['project_subject_categories'], X_train, 'p
project_subject_categories', 'clean_categories')
x_test_sorted_category_dict =
cleaning_text_data(X_test['project_subject_categories'], X_test, 'project_subject_categories', 'clean_
categories')

```

## 2.2) Text Preprocessing : project\_subject\_subcategories

In [21]:

```

x_train_sorted_subcategories = cleaning_text_data(X_train['project_subject_subcategories'], X_train
, 'project_subject_subcategories', 'clean_subcategories')
x_test_sorted_subcategories = cleaning_text_data(X_test['project_subject_subcategories'], X_test, 'p
project_subject_subcategories', 'clean_subcategories')

```

## 2.3) Text Preprocessing: project\_grade\_category

In [22]:

```

x_train_sorted_grade =
clean_project_grade(X_train['project_grade_category'], X_train, 'project_grade_category', 'clean_grade

```

```
)
x_test_sorted_grade =
clean_project_grade(X_test['project_grade_category'],X_test,'project_grade_category','clean_grade'
)
```

## 2.4) Text Preprocessing (stowords): project\_essay, project\_title

In [23]:

# <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 [24]:

# <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 [25]:

```
# Combining all the above students
```

```
from tqdm import tqdm
```

```
def process_text(df,col name):
```

```
preprocessed_feature = []
```

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

```
for sentence in tqdm(df[col_name].values):
```

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

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

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

[illegible]

```

sent = sent.replace('\n', ' ')
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
# https://gist.github.com/sebleier/554280
sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
preprocessed_feature.append(sent.lower().strip())
return preprocessed_feature

```

In [26]:

```

x_train_essay_preprocessed = process_text(X_train, 'essay')
x_test_essay_preprocessed = process_text(X_test, 'essay')

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33500/33500
[00:33<00:00, 1005.37it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:16<00:00, 1013.43it/s]

```

In [27]:

```

x_train_title_preprocessed = process_text(X_train, 'project_title')
x_test_title_preprocessed = process_text(X_test, 'project_title')

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33500/33500
[00:01<00:00, 20241.77it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:00<00:00, 21373.04it/s]

```

## 2.5) Vectorizing Categorical Data

### project\_subject\_categories (clean\_categories)

In [28]:

```

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
def cat_vectorizer(X_train, df, col_name):
    vectorizer = CountVectorizer()
    vectorizer.fit(X_train[col_name].values)
    feature_one_hot = vectorizer.transform(df[col_name].values)
    print(vectorizer.get_feature_names())
    return feature_one_hot, vectorizer.get_feature_names()

```

In [29]:

```

x_train_cat_one_hot, x_train_cat_feat_list = cat_vectorizer(X_train, X_train, 'clean_categories')
x_test_cat_one_hot, x_test_cat_feat_list = cat_vectorizer(X_train, X_test, 'clean_categories')

```

```

['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']

```

In [30]:

```

# shape after categorical one hot encoding
print(x_train_cat_one_hot.shape)
print(x_test_cat_one_hot.shape)

```

```

(33500, 9)
(16500, 9)

```

### project\_subject\_subcategory (clean\_subcategory)

In [31]:

```
x_train_subcat_one_hot, x_train_subcat_feat_list =
cat_vectorizer(X_train,X_train,'clean_subcategories')
x_test_subcat_one_hot, x_test_subcat_feat_list =
cat_vectorizer(X_train,X_test,'clean_subcategories')
```

```
['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']
['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']
```

In [32]:

```
# shape after categorical one hot encoding
print(x_train_subcat_one_hot.shape)
print(x_test_subcat_one_hot.shape)
```

```
(33500, 30)
(16500, 30)
```

## school\_state

In [33]:

```
# we use count vectorizer to convert the values into one hot encoding
# CountVectorizer for "school_state"
x_train_state_one_hot, x_train_state_feat_list = cat_vectorizer(X_train,X_train,'school_state')
x_test_state_one_hot, x_test_state_feat_list = cat_vectorizer(X_train,X_test,'school_state')
```

```
['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']
['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 [34]:

```
# shape after categorical one hot encoding
print(x_train_state_one_hot.shape)
print(x_test_state_one_hot.shape)
```

```
(33500, 51)
(16500, 51)
```

## teacher\_prefix

In [35]:

```
# we use count vectorizer to convert the values into one hot encoding
# CountVectorizer for teacher_prefix
x_train_teacher_prefix_one_hot,x_train_teacher_prefix_feat_list = cat_vectorizer(X_train,X_train,'
teacher_prefix')
x_test_teacher_prefix_one_hot,x_test_teacher_prefix_feat_list =
cat_vectorizer(X_train,X_test,'teacher_prefix')
```

```
['mr', 'mrs', 'ms', 'teacher']
['mr', 'mrs', 'ms', 'teacher']
```



In [36]:

```
# shape after categorical one hot encoding
print(x_train_teacher_prefix_one_hot.shape)
print(x_test_teacher_prefix_one_hot.shape)
```

```
(33500, 4)
(16500, 4)
```

## project\_grade\_category

In [37]:

```
# using count vectorizer for one-hot encoding of project_grade_category
x_train_grade_one_hot, x_train_grade_feat_list = cat_vectorizer(X_train,X_train,'clean_grade')
x_test_grade_one_hot, x_test_grade_feat_list = cat_vectorizer(X_train,X_test,'clean_grade')
```

```
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
```

In [38]:

```
# shape after categorical one hot encoding
print(x_train_grade_one_hot.shape)
print(x_test_grade_one_hot.shape)
```

```
(33500, 4)
(16500, 4)
```

## 2.6) Vectorizing Text Data

### 2.6.1) Bag of Words (essay)

In [39]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
def bow_vectorizer(X_train,col_name,df):
    vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
    vectorizer.fit(X_train[col_name].values)
    df_bow = vectorizer.transform(df[col_name].values)
    return df_bow, vectorizer.get_feature_names()
```

In [40]:

```
x_train_essay_bow, x_train_essay_feat = bow_vectorizer(X_train,'essay',X_train)
x_test_essay_bow, x_test_essay_feat = bow_vectorizer(X_train,'essay',X_test)
```

In [41]:

```
print(x_train_essay_bow.shape)
print(x_test_essay_bow.shape)
```

```
(33500, 5000)
(16500, 5000)
```

### 2.7.2) Bag of Words (title)

In [42]:

```
def bow_vectorizer_title(X_train,col_name,df):
    vectorizer = CountVectorizer()
```

```
vectorizer.fit(X_train[col_name].values)
df_bow = vectorizer.transform(df[col_name].values)
return df_bow, vectorizer.get_feature_names()
```

In [43]:

```
x_train_title_bow, x_train_title_feat = bow_vectorizer_title(X_train, 'project_title', X_train)
x_test_title_bow, x_test_title_feat = bow_vectorizer_title(X_train, 'project_title', X_test)
```

In [44]:

```
print(x_train_title_bow.shape)
print(x_test_title_bow.shape)
```

```
(33500, 9968)
(16500, 9968)
```

### 2.7.3) TFIDF (essay)

In [45]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents (rows or projects).
def tfidf_vectorizer(X_train, col_name, df):
    vectorizer = TfidfVectorizer(min_df=10, ngram_range = (2,2), max_features = 5000)
    vectorizer.fit(X_train[col_name].values)
    df_tfidf = vectorizer.transform(df[col_name].values)
    return df_tfidf, vectorizer.get_feature_names()
```

In [46]:

```
# Lets vectorize essay
x_train_essay_tfidf, x_train_essay_tfidf_feat = tfidf_vectorizer(X_train, 'essay', X_train)
x_test_essay_tfidf, x_test_essay_tfidf_feat = tfidf_vectorizer(X_train, 'essay', X_test)
```

In [47]:

```
print(x_train_essay_tfidf.shape)
print(x_test_essay_tfidf.shape)
```

```
(33500, 5000)
(16500, 5000)
```

### 2.7.4) TFIDF (title)

In [48]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
def tfidf_vectorizer_title(X_train, col_name, df):
    vectorizer = TfidfVectorizer()
    vectorizer.fit(X_train[col_name].values)
    df_tfidf = vectorizer.transform(df[col_name].values)
    return df_tfidf, vectorizer.get_feature_names()
```

In [49]:

```
# Lets vectorize essay
x_train_title_tfidf, x_train_title_tfidf_feat =
tfidf_vectorizer_title(X_train, 'project_title', X_train)
x_test_title_tfidf, x_test_title_tfidf_feat =
tfidf_vectorizer_title(X_train, 'project_title', X_test)
```

In [50]:

```
print(x_train_title_tfidf.shape)
```

```
print(x_test_title_tfidf.shape)
```

```
(33500, 9968)
(16500, 9968)
```

## 2.7.5) Using Pretrained Models: Avg W2V

In [51]:

```
'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

words = []
for i in preprocod_texts:
    words.extend(i.split(' '))

for i in preprocod_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(", np.round(len(inter_words)/len(words)*100,3), "%")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

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

'''
```

Out[51]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\r',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n
word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\glove.42B.300d.txt')\n\n# =====\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
```

```

=====
\n\nwords = []\nfor i in preproced_texts:\n    words.extend(i.split('\n\n'))\n\nfor i in preproced_titles:\n    words.extend(i.split('\n\n'))\nprint("all the words in the coupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha t are present in both glove vectors and our coupus", len(inter_words), " ("np.round(len(inter_words)/len(words)*100,3),"%")\n\nwords_courpus = {}\nwords_glove = set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\nprint("word 2 vec length", len(words_courpus))\n\n# stronging variables into pickle files python : http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pickle\nwith open('glove_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n
```

In [52]:

```

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

```

In [53]:

```

# Combining all the above stundents
from tqdm import tqdm
def preprocess_essay(df,col_name):
    preprocessed_essays = []
    # tqdm is for printing the status bar
    for sentence in tqdm(df[col_name].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_essays.append(sent.lower().strip())
    return preprocessed_essays

```

In [54]:

```

# average Word2Vec
# compute average word2vec for each review.
def compute_avg_W2V(preprocessed_feature):
    avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(preprocessed_feature): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        cnt_words =0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if word in glove_words:
                vector += model[word]
                cnt_words += 1
        if cnt_words != 0:
            vector /= cnt_words
        avg_w2v_vectors.append(vector)
    return avg_w2v_vectors

```

In [55]:

```

x_train_preprocessed_essay = preprocess_essay(X_train,'essay')
x_test_preprocessed_essay = preprocess_essay(X_test,'essay')

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33500/33500 [00:
33<00:00, 998.33it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:15<00:00, 1068.75it/s]

```

In [56]:

```

x_train_preprocessed_title = preprocess_essay(X_train,'project_title')
x_test_preprocessed_title = preprocess_essay(X_test,'project_title')

```

```
100%|██████████████████████████████████████████████████████████████████████████| 33500/33500  
[00:01<00:00, 19420.45it/s]  
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 20939.10it/s]
```

In [57]:

```
x_train_avg_w2v_essay = compute_avg_W2V(x_train_preprocessed_essay)
x_test_avg_w2v_essay = compute_avg_W2V(x_test_preprocessed_essay)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 33500/33500  
[00:24<00:00, 1366.46it/s]  
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:11<00:00, 1383.07it/s]
```

In [58]:

```
x_train_avg_w2v_title = compute_avg_W2V(x_train_preprocessed_title)
x_test_avg_w2v_title = compute_avg_W2V(x_test_preprocessed_title)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 33500/33500  
[00:01<00:00, 28438.27it/s]  
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 28846.32it/s]
```

### 2.7.6) Using Pretrained Models: TFIDF Weighted W2V

In [59]:

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

In [60]:

```
# average Word2Vec
# compute average word2vec for each review.
def compute_tfidf_w2v_vectors(preprocessed_feature):
    tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    dictionary, tfidf_words = get_tfidf_dict(preprocessed_feature)
    for sentence in tqdm(preprocessed_feature): # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
        for word in sentence.split(): # for each word in a review/sentence
            if (word in glove_words) and (word in tfidf_words):
                vec = model[word] # getting the vector for each word
                # here we are multiplying idf value(dictionary[word]) and the tf
                value((sentence.count(word)/len(sentence.split())))
                tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting
                the tfidf value for each word
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
            tfidf_w2v_vectors.append(vector)
    return tfidf_w2v_vectors
```

In [61]:

```
x_train_weighted_w2v_essay = compute_tfidf_w2v_vectors(x_train_essay_preprocessed)
x_test_weighted_w2v_essay = compute_tfidf_w2v_vectors(x_test_essay_preprocessed)
```

```
100% |██████████████████████████████████████████████████████████████| 33500/33500 [02:  
17<00:00, 243.50it/s]
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [01:08<00:00, 240.56it/s]
```

In [62]:

```
x_train_weighted_w2v_title = compute_tfidf_w2v_vectors(x_train_title_preprocessed)
x_test_weighted_w2v_title= compute_tfidf_w2v_vectors(x_test_title_preprocessed)
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 33500/33500
[00:02<00:00, 12869.89it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:01<00:00, 12547.54it/s]
```

## 2.7.7) Vectorizing Numerical Features

We have 2 numerical features left, "price" and "teacher\_number\_of\_previously\_posted\_projects". Let's check for the "missing" or "NaN" values present in those numerical features and use "Mean Replacement" for "price" and "Mode Replacement" for "teacher\_number\_of\_previously\_posted\_projects".

In [63]:

```
print("Total number of \"Missing\" Values present in X_train price:",X_train['price'].isna().sum()
)
print("Total number of \"Missing\" Values present in X_test price:",X_test['price'].isna().sum())
```

```
Total number of "Missing" Values present in X_train price: 32451
Total number of "Missing" Values present in X_test price: 15975
```

In [64]:

```
print("Total number of \"Missing\" Values present in X_train previous teacher number:",X_train['teacher_number_of_previously_posted_projects'].isna().sum())
print("Total number of \"Missing\" Values present in X_test previous teacher number:",X_test['teacher_number_of_previously_posted_projects'].isna().sum())
```

```
Total number of "Missing" Values present in X_train previous teacher number: 0
Total number of "Missing" Values present in X_test previous teacher number: 0
```

In [65]:

```
print("Total number of \"Missing\" Values present in X_train quantity:",X_train['quantity'].isna().sum())
print("Total number of \"Missing\" Values present in X_test quantity:",X_test['quantity'].isna().sum())
```

```
Total number of "Missing" Values present in X_train quantity: 32451
Total number of "Missing" Values present in X_test quantity: 15975
```

"teacher\_number\_of\_previously\_posted\_projects" does not have any "missing" values.

In [66]:

```
X_train['price'].mean()
```

Out[66]:

```
287.1204575786462
```

In [67]:

```
X_train['price'] = X_train['price'].fillna(287.1204)
```

In [68]:

```
X_test['price'] = X_test['price'].fillna(287.1204)
```

```
X_test['price'].mean()
```

Out[68]:

291.04784761904796

In [69]:

```
X_test['price'] = X_test['price'].fillna(291.0478)
```

In [70]:

```
print(X_train['quantity'].mean())
print(X_test['quantity'].mean())
```

19.142993326978075

17.588571428571427

In [71]:

```
X_train['quantity'] = X_train['quantity'].fillna(19.1429)
X_test['quantity'] = X_test['quantity'].fillna(17.5885)
```

In [72]:

```
# check this one: https://www.youtube.com/watch?v=0HOqOcln3Z4&t=530s
# https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html
from sklearn.preprocessing import StandardScaler
def scaler_function(df,col_name):

    scaler = StandardScaler()
    scaler.fit(df[col_name].values.reshape(-1,1)) # finding the mean and standard deviation of this
data

    # Now standardize the data with above maen and variance.
    print(f"Mean : {scaler.mean_[0]}, Standard deviation : {np.sqrt(scaler.var_[0])}")
    scaled = scaler.transform(df[col_name].values.reshape(-1, 1))
    return scaled
```

**teacher\_number\_of\_previously\_posted\_projects**

In [73]:

```
x_train_teacher_number = scaler_function(X_train,'teacher_number_of_previously_posted_projects')
x_test_teacher_number = scaler_function(X_test,'teacher_number_of_previously_posted_projects')
```

Mean : 11.331701492537313, Standard deviation : 28.327010543978243

Mean : 11.082121212121212, Standard deviation : 27.807397888307367

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:  
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:  
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:  
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:  
Data with input dtype int64 was converted to float64 by StandardScaler.

**price**

In [74]:

```
x_train_price = scaler_function(X_train, 'price')
x_test_price = scaler_function(X_test, 'price')
```

Mean : 287.120401802985, Standard deviation : 57.04108764222917  
Mean : 291.04780151515143, Standard deviation : 72.65831828074468

In [75]:

```
x_train_quantity = scaler_function(X_train, 'quantity')
x_test_quantity = scaler_function(X_test, 'quantity')
```

Mean : 19.14290292238806, Standard deviation : 6.69524792716915  
Mean : 17.58850227272726, Standard deviation : 4.106322777547319

## 2.8) Merging all the features and building the sets

In [76]:

```
# train dataset
print("After Vectorization and One hot encoding train dataset shape becomes:")
print(x_train_cat_one_hot.shape)
print(x_train_subcat_one_hot.shape)
print(x_train_state_one_hot.shape)
print(x_train_teacher_prefix_one_hot.shape)
print(x_train_grade_one_hot.shape)
print(x_train_essay_bow.shape)
print(x_train_title_bow.shape)
print(x_train_essay_tfidf.shape)
print(x_train_title_tfidf.shape)
print(np.asarray(x_train_avg_w2v_essay).shape)
print(np.asarray(x_train_avg_w2v_title).shape)
print(np.asarray(x_train_weighted_w2v_essay).shape)
print(np.asarray(x_train_weighted_w2v_title).shape)
print(x_train_teacher_number.shape)
print(x_train_price.shape)
print(x_train_quantity.shape)
print("="*50)

# test dataset
print("After Vectorization and One hot encoding test dataset shape becomes:")
print(x_test_cat_one_hot.shape)
print(x_test_subcat_one_hot.shape)
print(x_test_state_one_hot.shape)
print(x_test_teacher_prefix_one_hot.shape)
print(x_test_grade_one_hot.shape)
print(x_test_essay_bow.shape)
print(x_test_title_bow.shape)
print(x_test_essay_tfidf.shape)
print(x_test_title_tfidf.shape)
print(np.asarray(x_test_avg_w2v_essay).shape)
print(np.asarray(x_test_avg_w2v_title).shape)
print(np.asarray(x_test_weighted_w2v_essay).shape)
print(np.asarray(x_test_weighted_w2v_title).shape)
print(x_test_teacher_number.shape)
print(x_test_price.shape)
print(x_test_quantity.shape)
print("="*50)
```

After Vectorization and One hot encoding train dataset shape becomes:

```
(33500, 9)
(33500, 30)
(33500, 51)
(33500, 4)
(33500, 4)
(33500, 5000)
(33500, 9968)
(33500, 5000)
(33500, 9968)
(33500, 300)
```



```
(33500, 300)
(33500, 300)
(33500, 300)
(33500, 1)
(33500, 1)
(33500, 1)
```

=====

After Vectorization and One hot encoding test dataset shape becomes:

```
(16500, 9)
(16500, 30)
(16500, 51)
(16500, 4)
(16500, 4)
(16500, 5000)
(16500, 9968)
(16500, 5000)
(16500, 9968)
(16500, 300)
(16500, 300)
(16500, 300)
(16500, 300)
(16500, 1)
(16500, 1)
(16500, 1)
```

=====

### 2.8.1) Set 1: categorical, numerical features + project\_title(BOW) + preprocessed\_essay (BOW with bi-grams with min\_df=10 and max\_features=5000)

In [77]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X_train_set_1 =
hstack((x_train_cat_one_hot,x_train_subcat_one_hot,x_train_state_one_hot,x_train_teacher_prefix_one_hot,\
x_train_grade_one_hot,x_train_teacher_number,x_train_price,x_train_quantity,x_train_title_bow,x_train_essay_bow)).tocsr()
X_test_set_1 =
hstack((x_test_cat_one_hot,x_test_subcat_one_hot,x_test_state_one_hot,x_test_teacher_prefix_one_hot,\
x_test_grade_one_hot,x_test_teacher_number,x_test_price,x_test_quantity,x_test_title_bow,x_test_essay_bow)).tocsr()
```

In [78]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.GridSearchCV.html
# for class prior i referred https://scikit-learn.org/stable/modules/generated/sklearn.naive\_bayes.MultinomialNB.html,\
# and https://stackoverflow.com/questions/42498208/setting-prior-probabilities-in-naive-bayes-multinomialnb
from sklearn.model_selection import GridSearchCV
from scipy.stats import randint as sp_randint
from sklearn.model_selection import RandomizedSearchCV
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
import math

log_reg = LogisticRegression(class_weight = "balanced")
parameters = {'C':[10**x for x in range(-4,5)]}
clf = RandomizedSearchCV(log_reg, parameters,n_iter = 9, scoring='roc_auc')
clf.fit(X_train_set_1, y_train)

results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['param_C'])

train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
```

```

cv_auc_std= results['std_test_score']
C_ = results['param_C'].apply(lambda x: math.log10(x))

plt.plot(C_, train_auc, label='Train AUC')

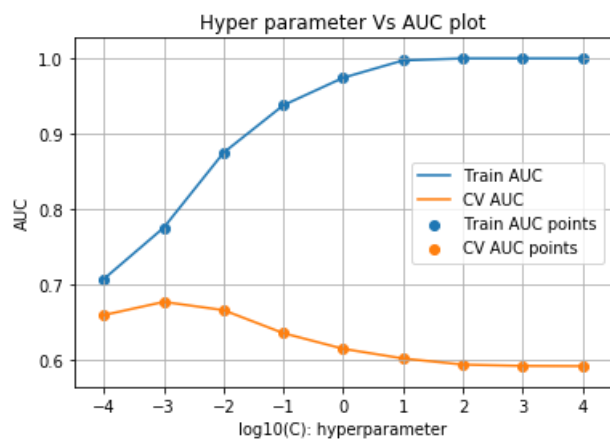
plt.plot(C_, cv_auc, label='CV AUC')

plt.scatter(C_, train_auc, label='Train AUC points')
plt.scatter(C_, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("log10(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

results.head()

```



Out[78]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	s
0	0.811543	0.044782	0.014673	0.000942	0.0001	{'C': 0.0001}	0.679561	0.642758	0
1	1.277167	0.109940	0.012008	0.000009	0.001	{'C': 0.001}	0.691917	0.664490	0
2	2.487216	0.194052	0.012671	0.001877	0.01	{'C': 0.01}	0.672369	0.657613	0
3	5.141678	0.089234	0.010669	0.002878	0.1	{'C': 0.1}	0.634602	0.632223	0
4	11.795676	0.715600	0.009663	0.000955	1	{'C': 1}	0.608311	0.616749	0

In [79]:

```

# From the AUC plot, we find that the best value for "C" - "Inverse of Regularization Strength" for the LogisticRegression is 0.01
best_C = 0.001

```

In [80]:

```

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

log_reg = LogisticRegression(C=best_C, class_weight = "balanced")
log_reg.fit(X_train_set_1, 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 = log_reg.predict_proba(X_train_set_1)
y_test_pred = log_reg.predict_proba(X_test_set_1)

y_train_pred_prob = []
y_test_pred_prob = []

for index in range(len(y_train_pred)):
    y_train_pred_prob.append(y_train_pred[index][1])

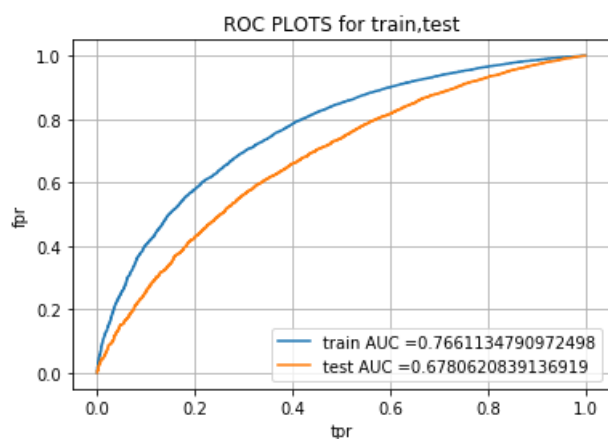
for index in range(len(y_test_pred)):
    y_test_pred_prob.append(y_test_pred[index][1])

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_prob)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_prob)

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("tpr")
plt.ylabel("fpr")
plt.title("ROC PLOTS for train,test")
plt.grid()
plt.show()

```



In [81]:

```

def compute_auc_with_hyper_para(x_tr,y_tr,x_cv,y_cv,para_list):
    auc_tr = []
    auc_cv = []
    y_train_pred_prob = []
    y_cv_pred_prob = []

    for para in para_list:
        log_reg = LogisticRegression(C = para, class_weight = "balanced")
        log_reg.fit(x_tr,y_tr)

        y_train_pred = log_reg.predict_proba(x_tr)
        y_cv_pred = log_reg.predict_proba(x_cv)

        for index in range(len(y_train_pred)):
            y_train_pred_prob.append(y_train_pred[index][1])

        for index in range(len(y_cv_pred)):
            y_cv_pred_prob.append(y_cv_pred[index][1])

        train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, y_train_pred_prob)
        cv_fpr, cv_tpr, tc_thresholds = roc_curve(y_cv, y_cv_pred_prob)

        y_train_pred_prob = []
        y_cv_pred_prob = []

```

```

    auc_tr.append(auc(train_fpr,train_tpr))
    auc_cv.append(auc(cv_fpr,cv_tpr))
plt.plot(para_list,auc_tr,label="train_auc")
plt.plot(para_list,auc_cv,label="cv_auc")
plt.legend()
plt.xlabel("Hyper Parameter:C")
plt.ylabel("Area Under ROC Curve")
plt.title("auc v/s hyper parameters")
plt.grid()
plt.show()

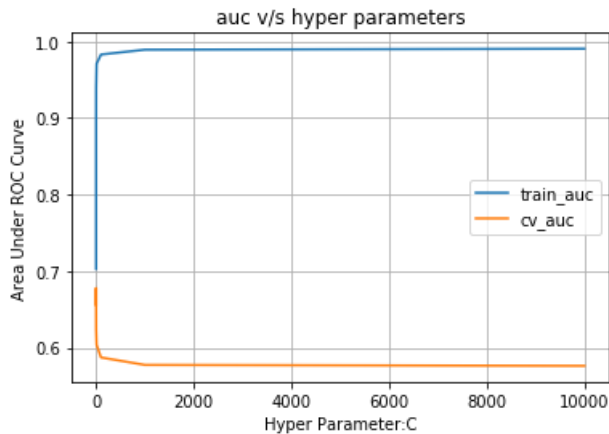
```

In [82]:

```

para_list = [10**x for x in range(-4,5)]
compute_auc_with_hyper_para(X_train_set_1,y_train,X_test_set_1,y_test,para_list)

```



In [83]:

```

# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(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))
    return t

def predict_with_best_t(proba, threshold):
    predictions = []
    for i in proba:
        if i>=threshold:
            predictions.append(1)
        else:
            predictions.append(0)
    return predictions

```

In [84]:

```

print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, predict_with_best_t(y_train_pred_prob, best_t)),annot = True,
fmt = "d", cbar=False)

```

```

=====
the maximum value of tpr*(1-fpr) 0.4883840281247254 for threshold 0.5
Train confusion matrix

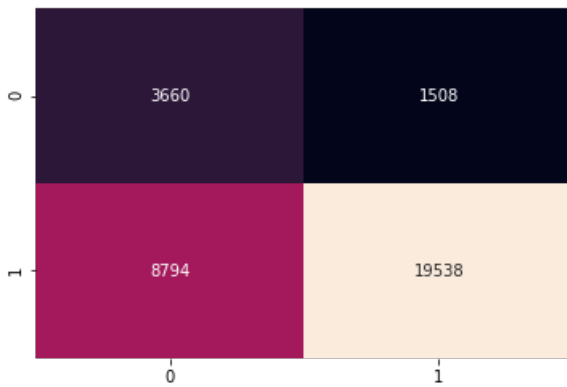
```

Out[84]:

```

<matplotlib.axes._subplots.AxesSubplot at 0x23dda3200f0>

```



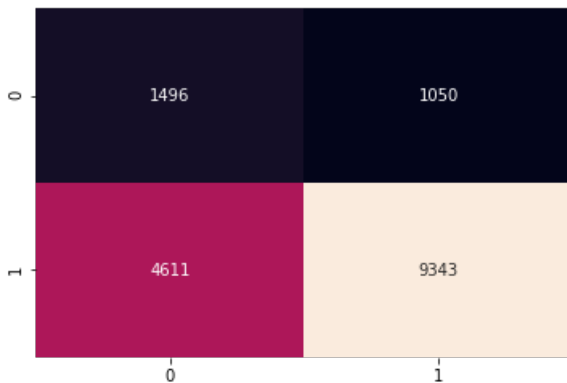
In [85]:

```
print("Test confusion matrix")
sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

Test confusion matrix

Out[85]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23dd9f3df98>



## 2.8.2) Set 2: categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (TFIDF with bi-grams with min\_df=10 and max\_features=5000)

In [86]:

```
X_train_set_2 =
hstack((x_train_cat_one_hot,x_train_subcat_one_hot,x_train_state_one_hot,x_train_teacher_prefix_one
_hot,\
x_train_grade_one_hot,x_train_teacher_number,x_train_price,x_train_quantity,x_train_title_tfidf,x_
train_essay_tfidf)).tocsr()
X_test_set_2 =
hstack((x_test_cat_one_hot,x_test_subcat_one_hot,x_test_state_one_hot,x_test_teacher_prefix_one_hot
,\
x_test_grade_one_hot,x_test_teacher_number,x_test_price,x_test_quantity,x_
test_title_tfidf,x_test_essay_tfidf)).tocsr()
```

In [87]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
# for class prior i referred https://scikit-
learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html,\
# and https://stackoverflow.com/questions/42498208/setting-prior-probabilities-in-naive-bayes-mult
inomialnb
from sklearn.model_selection import GridSearchCV
from scipy.stats import randint as sp_randint
from sklearn.model_selection import RandomizedSearchCV
import matplotlib.pyplot as plt
```

```

from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
import math

log_reg = LogisticRegression(class_weight = "balanced")
parameters = {'C': [10**x for x in range(-4,5)]}
clf = RandomizedSearchCV(log_reg, parameters, n_iter = 9, scoring='roc_auc')
clf.fit(X_train_set_2, y_train)

results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['param_C'])

train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
C_ = results['param_C'].apply(lambda x: math.log10(x))

plt.plot(C_, train_auc, label='Train AUC')

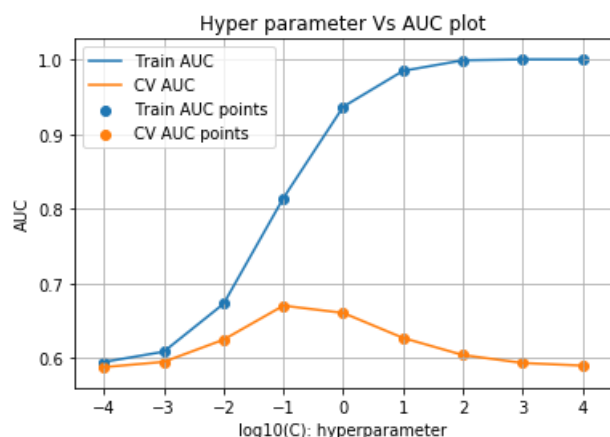
plt.plot(C_, cv_auc, label='CV AUC')

plt.scatter(C_, train_auc, label='Train AUC points')
plt.scatter(C_, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("log10(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

results.head()

```



Out[87]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	s
0	0.547330	0.023224	0.011667	0.000472	0.0001	{'C': 0.0001}	0.608939	0.578229	0
1	0.829996	0.063717	0.012667	0.000941	0.001	{'C': 0.001}	0.616226	0.588160	0
2	1.655000	0.258317	0.015672	0.001252	0.01	{'C': 0.01}	0.642466	0.622090	0
3	2.867334	0.221158	0.009666	0.001248	0.1	{'C': 0.1}	0.680244	0.666542	0
4	6.709006	0.136877	0.011993	0.001424	1	{'C': 1}	0.663853	0.657622	0

In [88]:

```
# From the AUC plot, we find that the best value for "C" - "Inverse of Regularization Strength" for the LogisticRegression is 0.01
best_C = 0.1
```

In [89]:

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

log_reg = LogisticRegression(C=best_C, class_weight = "balanced")
log_reg.fit(X_train_set_2, y_train)

y_train_pred = log_reg.predict_proba(X_train_set_2)
y_test_pred = log_reg.predict_proba(X_test_set_2)

y_train_pred_prob = []
y_test_pred_prob = []

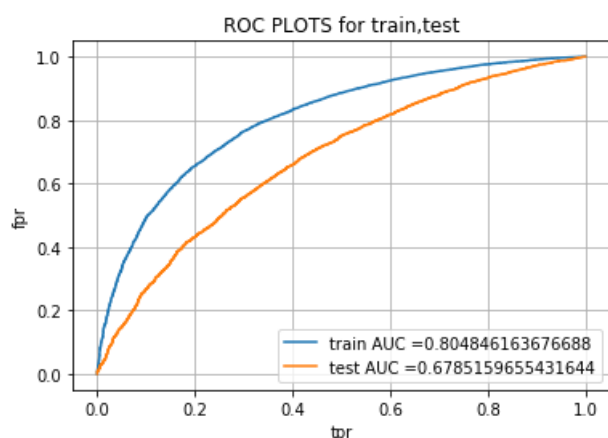
for index in range(len(y_train_pred)):
    y_train_pred_prob.append(y_train_pred[index][1])

for index in range(len(y_test_pred)):
    y_test_pred_prob.append(y_test_pred[index][1])

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_prob)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_prob)

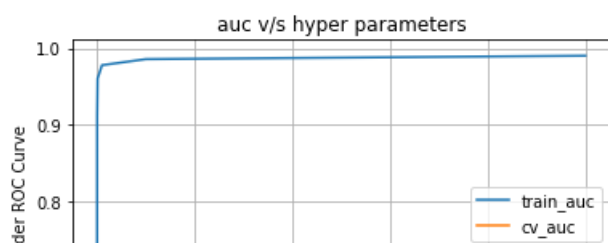
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))

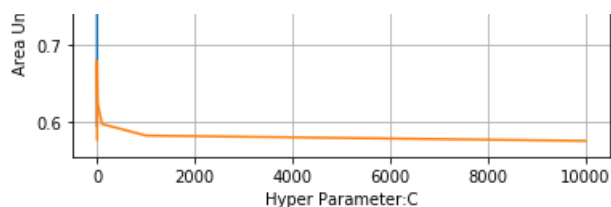
plt.legend()
plt.xlabel("tpr")
plt.ylabel("fpr")
plt.title("ROC PLOTS for train,test")
plt.grid()
plt.show()
```



In [90]:

```
para_list = [10**x for x in range(-4,5)]
compute_auc_with_hyper_para(X_train_set_2,y_train,X_test_set_2,y_test,para_list)
```





In [91]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, predict_with_best_t(y_train_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

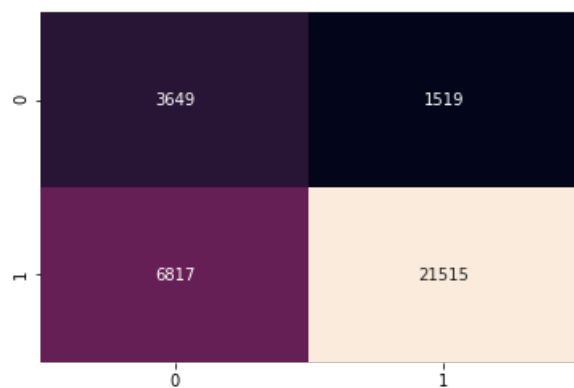
=====

the maximum value of tpr\*(1-fpr) 0.5361860067317682 for threshold 0.482

Train confusion matrix

Out[91]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23df36640b8>



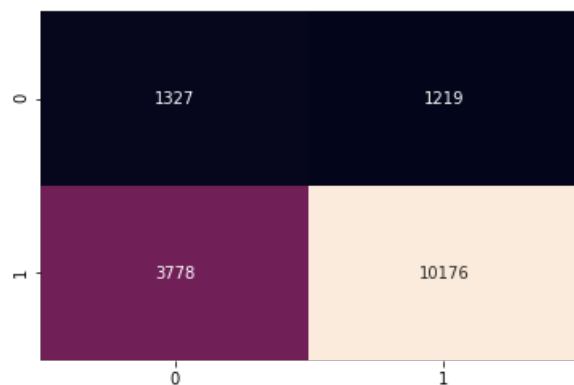
In [92]:

```
print("Test confusion matrix")
sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

Test confusion matrix

Out[92]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23dda4f8940>





### 2.8.3) Set 3: categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)

In [93]:

```
X_train_set_3 =
hstack((x_train_cat_one_hot,x_train_subcat_one_hot,x_train_state_one_hot,x_train_teacher_prefix_one_hot,\
x_train_grade_one_hot,x_train_teacher_number,x_train_price,x_train_quantity,x_train_avg_w2v_title,\
x_train_avg_w2v_essay)).tocsr()
X_test_set_3 =
hstack((x_test_cat_one_hot,x_test_subcat_one_hot,x_test_state_one_hot,x_test_teacher_prefix_one_hot,\
x_test_grade_one_hot,x_test_teacher_number,x_test_price,x_test_quantity,x_test_avg_w2v_title,x_test_avg_w2v_essay)).tocsr()
```

In [94]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
# for class prior i referred https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html,\
# and https://stackoverflow.com/questions/42498208/setting-prior-probabilities-in-naive-bayes-multinomialnb
from sklearn.model_selection import GridSearchCV
from scipy.stats import randint as sp_randint
from sklearn.model_selection import RandomizedSearchCV
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
import math

log_reg = LogisticRegression(class_weight = "balanced")
parameters = {'C':[10**x for x in range(-4,5)]}
clf = RandomizedSearchCV(log_reg, parameters,n_iter = 9, scoring='roc_auc')
clf.fit(X_train_set_3, y_train)

results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['param_C'])

train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
C_ = results['param_C'].apply(lambda x: math.log10(x))

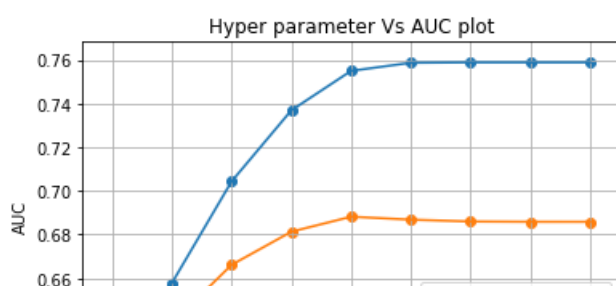
plt.plot(C_, train_auc, label='Train AUC')

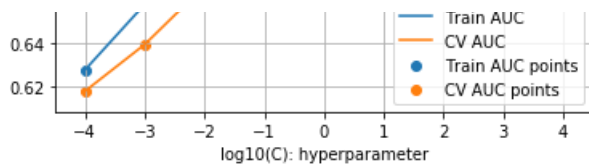
plt.plot(C_, cv_auc, label='CV AUC')

plt.scatter(C_, train_auc, label='Train AUC points')
plt.scatter(C_, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("log10(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

results.head()
```





Out [94]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	s
0	1.147349	0.031332	0.020652	0.001689	0.0001	{'C': 0.0001}	0.633438	0.612602	0
1	2.190009	0.094049	0.021341	0.002064	0.001	{'C': 0.001}	0.655723	0.635413	0
2	4.128012	0.229492	0.020663	0.001248	0.01	{'C': 0.01}	0.679335	0.663050	0
3	7.940000	0.327600	0.019999	0.000002	0.1	{'C': 0.1}	0.691738	0.677693	0
4	24.082243	3.012708	0.044334	0.003399	1	{'C': 1}	0.697724	0.684370	0

In [95]:

```
# From the AUC plot, we find that the best value for "C" - "Inverse of Regularization Strength" for the LogisticRegression is 0.01
best_C = 1
```

In [96]:

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

log_reg = LogisticRegression(C=best_C, class_weight = "balanced")
log_reg.fit(X_train_set_3, y_train)

y_train_pred = log_reg.predict_proba(X_train_set_3)
y_test_pred = log_reg.predict_proba(X_test_set_3)

y_train_pred_prob = []
y_test_pred_prob = []

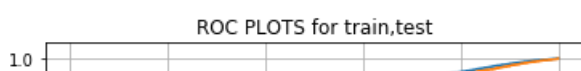
for index in range(len(y_train_pred)):
    y_train_pred_prob.append(y_train_pred[index][1])

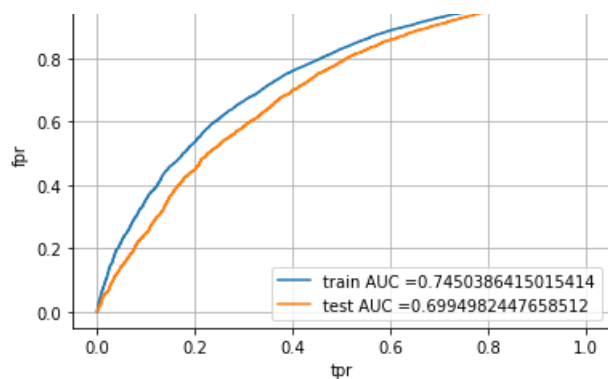
for index in range(len(y_test_pred)):
    y_test_pred_prob.append(y_test_pred[index][1])

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_prob)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_prob)

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)))

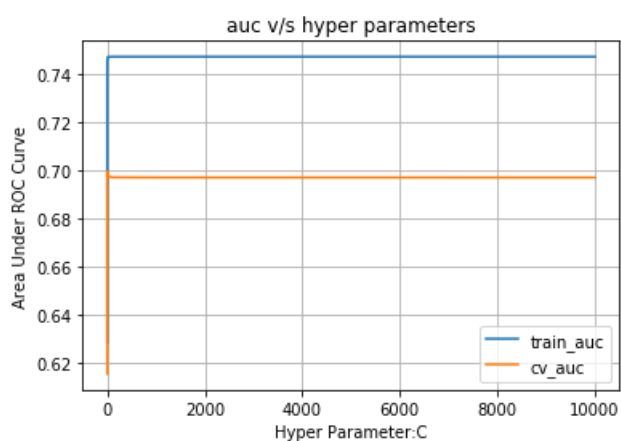
plt.legend()
plt.xlabel("tpr")
plt.ylabel("fpr")
plt.title("ROC PLOTS for train,test")
plt.grid()
plt.show()
```





In [98]:

```
para_list = [10**x for x in range(-4,5)]
compute_auc_with_hyper_para(X_train_set_3,y_train,X_test_set_3,y_test,para_list)
```



In [99]:

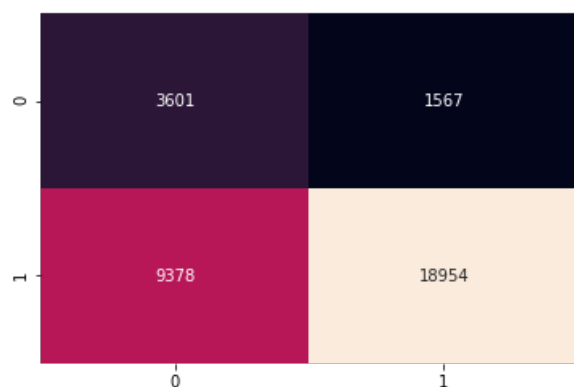
```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, predict_with_best_t(y_train_pred_prob, best_t)),annot = True,
fmt = "d", cbar=False)
```

=====

the maximum value of  $tpr \cdot (1 - fpr)$  0.4661484661744053 for threshold 0.5  
Train confusion matrix

Out[99]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23dd9c40208>



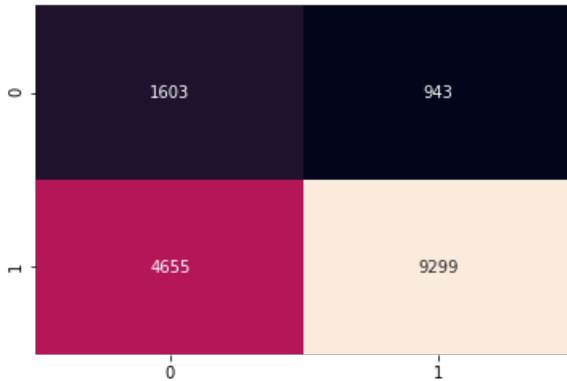
In [100]:

```
print("Test confusion matrix")
sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

Test confusion matrix

Out[100]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23dd9f4ac50>



## 2.8.4) Set 4: categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_essay (TFIDF W2V)

In [101]:

```
X_train_set_4 =
hstack((x_train_cat_one_hot,x_train_subcat_one_hot,x_train_state_one_hot,x_train_teacher_prefix_one
_hot,\

x_train_grade_one_hot,x_train_teacher_number,x_train_price,x_train_quantity,x_train_weighted_w2v_title,x_train_weighted_w2v_essay)).tocsr()
X_test_set_4 =
hstack((x_test_cat_one_hot,x_test_subcat_one_hot,x_test_state_one_hot,x_test_teacher_prefix_one_hot
,\

x_test_grade_one_hot,x_test_teacher_number,x_test_price,x_test_quantity,x_test_weighted_w2v_title,x_test_weighted_w2v_essay)).tocsr()
```

In [102]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
# for class prior i referred https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html,\
# and https://stackoverflow.com/questions/42498208/setting-prior-probabilities-in-naive-bayes-multinomialnb
from sklearn.model_selection import GridSearchCV
from scipy.stats import randint as sp_randint
from sklearn.model_selection import RandomizedSearchCV
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
import math

log_reg = LogisticRegression()
parameters = {'C':[10**x for x in range(-4,5)]}
clf = RandomizedSearchCV(log_reg, parameters,n_iter = 9, scoring='roc_auc')
clf.fit(X_train_set_4, y_train)

results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['param_C'])

train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
C_ = results['param_C'].apply(lambda x: math.log10(x))
```

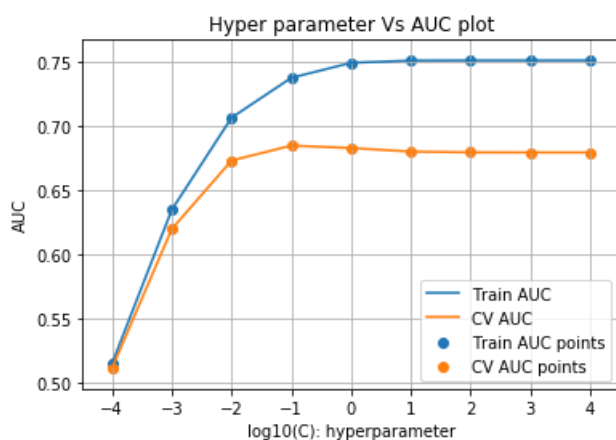
```
plt.plot(C_, train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# plt.gca().fill_between(K, train_auc - train_auc_std, train_auc +
train_auc_std,alpha=0.2,color='darkblue')

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

plt.scatter(C_, train_auc, label='Train AUC points')
plt.scatter(C_, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("log10(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

results.head()
```



Out[102]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	s
0	1.236678	0.049048	0.021324	0.001253	0.0001	{'C': 0.0001}	0.512632	0.510574	0
1	1.702676	0.048663	0.023995	0.000820	0.001	{'C': 0.001}	0.632200	0.615017	0
2	3.008675	0.170859	0.020334	0.000470	0.01	{'C': 0.01}	0.685462	0.668724	0
3	5.669356	0.489519	0.020338	0.000469	0.1	{'C': 0.1}	0.695624	0.680045	0
4	12.793331	0.559630	0.021334	0.001887	1	{'C': 1}	0.692743	0.679128	0

In [103]:

```
# From the AUC plot, we find that the best value for "C" - "Inverse of Regularization Strength" for the LogisticRegression is 0.01
best_C = 0.1
```

In [104]:

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

log_reg = LogisticRegression(C=best_C, class_weight = "balanced")
```

```

log_reg.fit(X_train_set_4, y_train)

y_train_pred = log_reg.predict_proba(X_train_set_4)
y_test_pred = log_reg.predict_proba(X_test_set_4)

y_train_pred_prob = []
y_test_pred_prob = []

for index in range(len(y_train_pred)):
    y_train_pred_prob.append(y_train_pred[index][1])

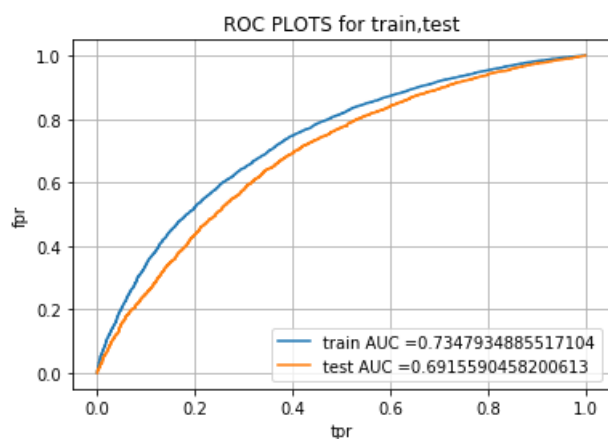
for index in range(len(y_test_pred)):
    y_test_pred_prob.append(y_test_pred[index][1])

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_prob)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_prob)

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("tpr")
plt.ylabel("fpr")
plt.title("ROC PLOTS for train,test")
plt.grid()
plt.show()

```

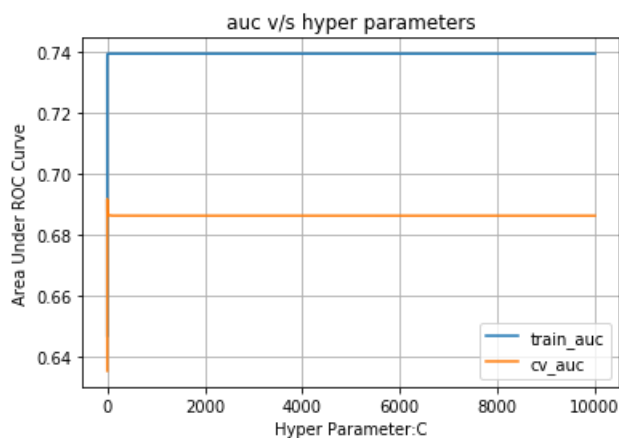


In [105]:

```

para_list = [10**x for x in range(-4,5)]
compute_auc_with_hyper_para(X_train_set_4,y_train,X_test_set_4,y_test,para_list)

```



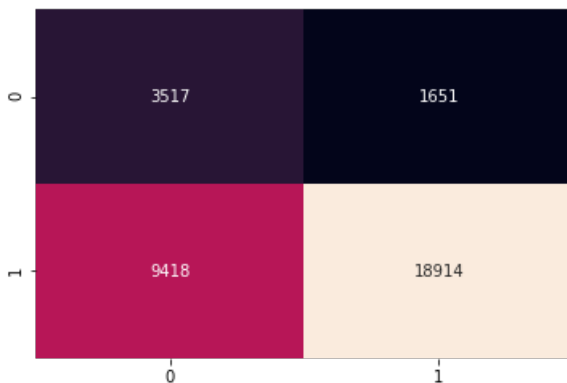
In [106]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, predict_with_best_t(y_train_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

the maximum value of tpr\*(1-fpr) 0.4543138899488549 for threshold 0.491  
Train confusion matrix

Out[106]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23dda39c6a0>



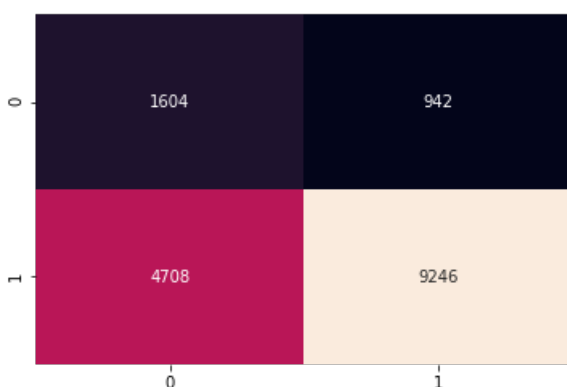
In [107]:

```
print("Test confusion matrix")
sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

Test confusion matrix

Out[107]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23df3673860>



## 2.8.5) Calculate Sentiment Score for each essay (combined)

In [112]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# import nltk
# nltk.download('vader_lexicon')
def compute_sentiment_score(df):
    score_list = []
    sid = SentimentIntensityAnalyzer()
```

```

for essay in df['essay']:
    ss = sid.polarity_scores(essay)
    score_list.append(ss)
return score_list
# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93

```

C:\ProgramData\Anaconda3\lib\site-packages\nltk\twitter\\_\_init\_\_.py:20: UserWarning:

The twython library has not been installed. Some functionality from the twitter package will not be available.

In [113]:

```

x_train_score = compute_sentiment_score(X_train)
x_test_score = compute_sentiment_score(X_test)

```

In [114]:

```

def populate_list(score_dicts):

    neg_score = []
    neu_score = []
    pos_score = []
    compound_score = []
    for dict_ in score_dicts:
        neg_score.append(dict_['neg'])
        neu_score.append(dict_['neu'])
        pos_score.append(dict_['pos'])
        compound_score.append(dict_['compound'])
    return neg_score, neu_score, pos_score, compound_score

```

In [115]:

```

x_train_neg, x_train_neu, x_train_pos, x_train_compound = populate_list(x_train_score)
x_test_neg, x_test_neu, x_test_pos, x_test_compound = populate_list(x_test_score)

```

In [116]:

```

X_train['words_project_title'] = X_train['project_title'].apply(lambda x: len(x.split()))
X_train['words_essay'] = X_train['essay'].apply(lambda x: len(x.split()))

```

In [117]:

```

X_test['words_project_title'] = X_test['project_title'].apply(lambda x: len(x.split()))
X_test['words_essay'] = X_test['essay'].apply(lambda x: len(x.split()))

```

In [119]:

```

# X_cv['words_project_title'] = X_cv['project_title'].apply(lambda x: len(x.split()))
# X_cv['words_essay'] = X_cv['essay'].apply(lambda x: len(x.split()))

```

Let's join the all the sentiment scores to the respective dataframes

In [120]:

```

# for training set
X_train['neg'] = x_train_neg
X_train['neu'] = x_train_neu
X_train['pos'] = x_train_pos
X_train['compound'] = x_train_compound

# for testing set
X_test['neg'] = x_test_neg
X_test['neu'] = x_test_neu
X_test['pos'] = x_test_pos
X_test['compound'] = x_test_compound

```



```
# # for cv set
# X_cv['neg'] = x_cv_neg
# X_cv['neu'] = x_cv_neu
# X_cv['pos'] = x_cv_pos
# X_cv['compound'] = x_cv_compound
```

In [121]:

```
x_train_neg_resaped = X_train.values.reshape(-1,1)
```

## 2.8.6) Vectorizing newly added numerical features: words\_project\_title, words\_essay

In [122]:

```
#Let's do for all the title words
x_train_words_title_scaled = scaler_function(X_train,'words_project_title')
x_test_words_title_scaled = scaler_function(X_test,'words_project_title')
# x_cv_words_title_scaled = scaler_function(X_cv,'words_project_title')

# let's do for all the essay words
x_train_words_essay_scaled = scaler_function(X_train,'words_essay')
x_test_words_essay_scaled = scaler_function(X_test,'words_essay')
# x_cv_words_essay_scaled = scaler_function(X_cv,'words_essay')
```

```
Mean : 5.1657014925373135, Standard deviation : 2.0973842345434
Mean : 5.149030303030303, Standard deviation : 2.101262744231156
Mean : 255.32268656716417, Standard deviation : 65.3401665586849
Mean : 254.10315151515152, Standard deviation : 64.80778593148109
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

## 2.8.6) Prepare dataset 5

In [123]:

```
X_train_set_5 =
hstack((x_train_cat_one_hot,x_train_subcat_one_hot,x_train_state_one_hot,x_train_teacher_prefix_one
_hot,\
```

```

x_train_grade_one_hot,x_train_teacher_number,x_train_price,x_train_quantity,x_train_words_title_scaled,x_train_words_essay_scaled, X_train['neg'].values.reshape(-1,1),X_train['neu'].values.reshape(-1,1),X_train['pos'].values.reshape(-1,1),X_train['compound'].values.reshape(-1,1)).tocsr()
X_test_set_5 =
hstack((x_test_cat_one_hot,x_test_subcat_one_hot,x_test_state_one_hot,x_test_teacher_prefix_one_hot,\
        x_test_grade_one_hot,x_test_teacher_number,x_test_price,x_test_quantity,x_test_words_title_scaled,x_test_words_essay_scaled,X_test['neg'].values.reshape(-1,1),X_test['neu'].values.reshape(-1,1),X_test['pos'].values.reshape(-1,1),X_test['compound'].values.reshape(-1,1)).tocsr()

```

In [124]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
# for class prior i referred https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html,\
# and https://stackoverflow.com/questions/42498208/setting-prior-probabilities-in-naive-bayes-multinomialnb
from sklearn.model_selection import GridSearchCV
from scipy.stats import randint as sp_randint
from sklearn.model_selection import RandomizedSearchCV
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.linear_model import LogisticRegression
import math

log_reg = LogisticRegression(class_weight = "balanced")
parameters = {'C':[10**x for x in range(-4,5)]}
clf = RandomizedSearchCV(log_reg, parameters,n_iter = 9, scoring='roc_auc')
clf.fit(X_train_set_5, y_train)

results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['param_C'])

train_auc= results['mean_train_score']
train_auc_std= results['std_train_score']
cv_auc = results['mean_test_score']
cv_auc_std= results['std_test_score']
C_ = results['param_C'].apply(lambda x: math.log10(x))

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

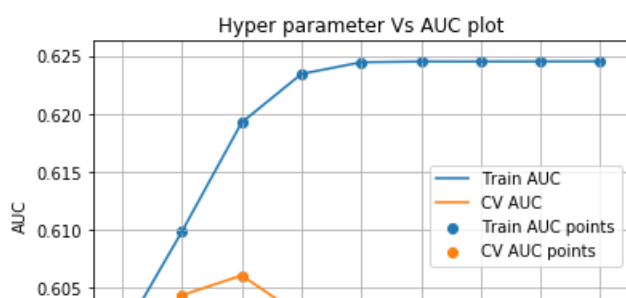
plt.plot(C_, cv_auc, label='CV AUC')
# # this code is copied from here: https://stackoverflow.com/a/48803361/4084039
# # plt.gca().fill_between(K, cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.2,color='darkorange')

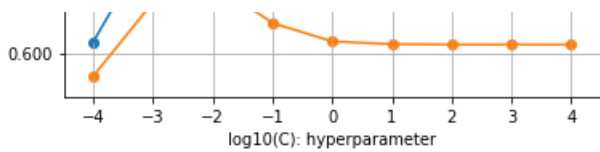
plt.scatter(C_, train_auc, label='Train AUC points')
plt.scatter(C_, cv_auc, label='CV AUC points')

plt.legend()
plt.xlabel("log10(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Hyper parameter Vs AUC plot")
plt.grid()
plt.show()

results.head()

```





Out [124]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_C	params	split0_test_score	split1_test_score	s
0	0.090675	0.011902	0.003999	0.000007	0.0001	{'C': 0.0001}	0.610298	0.592021	0
1	0.133999	0.002946	0.003662	0.000467	0.001	{'C': 0.001}	0.619420	0.600561	0
2	0.225341	0.005438	0.003332	0.000470	0.01	{'C': 0.01}	0.620285	0.606057	0
3	0.446016	0.034758	0.003337	0.000471	0.1	{'C': 0.1}	0.615055	0.605570	0
4	0.638688	0.014291	0.003991	0.000011	1	{'C': 1}	0.612012	0.605105	0

In [125]:

```
# From the AUC plot, we find that the best value for "C" - "Inverse of Regularization Strength" for the LogisticRegression is 0.01
best_C = 0.01
```

In [126]:

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

log_reg = LogisticRegression(C=best_C, class_weight = "balanced")
log_reg.fit(X_train_set_5, y_train)

y_train_pred = log_reg.predict_proba(X_train_set_5)
y_test_pred = log_reg.predict_proba(X_test_set_5)

y_train_pred_prob = []
y_test_pred_prob = []

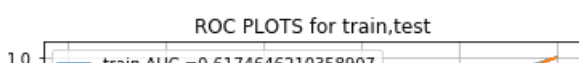
for index in range(len(y_train_pred)):
    y_train_pred_prob.append(y_train_pred[index][1])

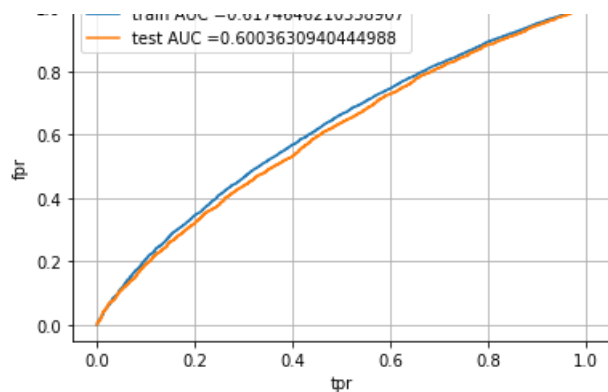
for index in range(len(y_test_pred)):
    y_test_pred_prob.append(y_test_pred[index][1])

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred_prob)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred_prob)

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)))

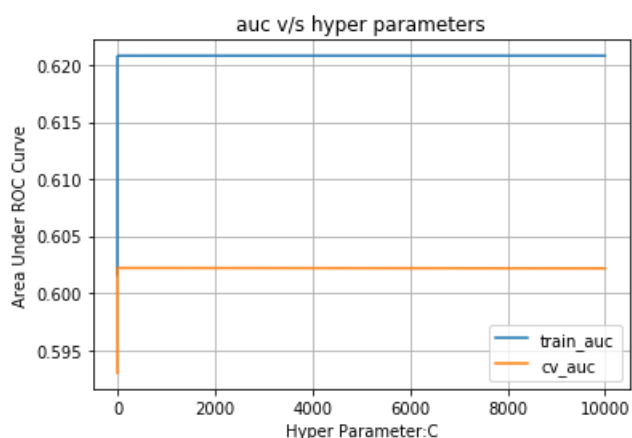
plt.legend()
plt.xlabel("tpr")
plt.ylabel("fpr")
plt.title("ROC PLOTS for train,test")
plt.grid()
plt.show()
```





In [127]:

```
para_list = [10**x for x in range(-4,5)]
compute_auc_with_hyper_para(X_train_set_5,y_train,X_test_set_5,y_test,para_list)
```



In [128]:

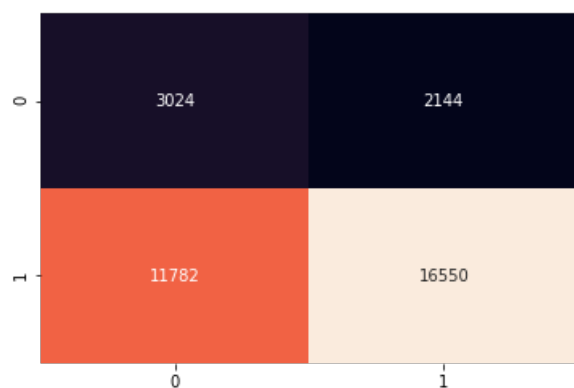
```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
sns.heatmap(confusion_matrix(y_train, predict_with_best_t(y_train_pred_prob, best_t)),annot = True,
fmt = "d", cbar=False)
```

=====

the maximum value of  $tpr \cdot (1 - fpr)$  0.3418062871507193 for threshold 0.491  
Train confusion matrix

Out[128]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23d87ad70f0>



In [129]:

```
In [129]:
```

```
print("Test confusion matrix")
sns.heatmap(confusion_matrix(y_test, predict_with_best_t(y_test_pred_prob, best_t)), annot = True,
fmt = "d", cbar=False)
```

Test confusion matrix

Out[129]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x23d826a9208>

