

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care & Hunger</code> <code>Health & Sports</code> <code>History & Civics</code> <code>Literacy & Language</code> <code>Math & Science</code> <code>Music & The Arts</code> <code>Special Needs</code> <code>Warmth</code> Examples: <ul style="list-style-type: none">• <code>Music & The Arts</code>• <code>Literacy & Language, Math & Science</code>
<code>school_state</code>		State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: <code>Literacy</code> <code>Literature & Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

Feature	Description
project_essay_4	Fourth application essay
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

- __project_essay_1__: "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- __project_essay_2__: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [1]:

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

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

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

```

C:\Users\Santosh\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
 warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")

Assignment 5: Logistic Regression

1. [Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay ('BOW with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay ('TFIDF with bi-grams' with 'min_df=10' and 'max_features=5000')
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

4. [Task-2] Apply Logistic Regression on the below feature set **Set 5** by finding the best hyper parameter as suggested in step 2 and step 3.

5. Consider these set of features **Set 5** :

- [school_state](#) : categorical data
- [clean_categories](#) : categorical data
- [clean_subcategories](#) : categorical data
- [project_grade_category](#) :categorical data

- [teacher_prefix](#) : categorical data
- [quantity](#) : numerical data
- [teacher_number_of_previously_posted_projects](#) : numerical data
- [price](#) : numerical data
- [sentiment_score's of each of the essay](#) : numerical data
- [number of words in the title](#) : numerical data
- [number of words in the combine essays](#) : numerical data

And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3.

6. Conclusion

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

Note: Data Leakage

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

1.1 Reading Data

In [2]:

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

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]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

Number of data points in train data (1541272, 4)

['id' 'description' 'quantity' 'price']

Out[4]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

In [5]:

```
# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
```

```

one_step
price_data=resource_data.groupby('id').agg({'price':'sum','quantity':'sum'}).reset_index()
price_data.head(2)

```

Out[5]:

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21

In [6]:

```

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

```

In [7]:

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

Out[7]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade

In [8]:

```

# presence of the numerical digits in a strings with numeric :
https://stackoverflow.com/a/19859308/8089731

def hasNumbers(inputString):
    return any(i.isdigit() for i in inputString)

p1=project_data[['id','project_resource_summary']]
p1=pd.DataFrame(data=p1)
p1.columns=['id','digits_in_summary']
p1['digits_in_summary']=p1['digits_in_summary'].map(hasNumbers)

# https://stackoverflow.com/a/17383325/8089731
p1['digits_in_summary'] = p1['digits_in_summary'].astype(int)
project_data=pd.merge(project_data,p1,on='id',how='left')
project_data.head(5)

```

Out[8]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade
2	21895 p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Grade

3	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Grades P

1.2 preprocessing of project_subject_categories

In [9]:

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

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

project_data['clean_categories']=cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)
project_data.head(5)
```

Out[9]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Grade
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17	Grades P
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Grades P

In [10]:

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

```
my_counter.update(word.split())
my_counter
```

Out[10]:

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

In [11]:

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

1.3 preprocessing of project_subject_subcategories

In [12]:

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

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

sub_cat_list = []
for i in sub_catogories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the catogory based on space "Math & Science"
            e=> "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 placeing 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())
```

In [13]:

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

Out[13]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945 p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade

In [14]:

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

In [15]:

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

1.3 Text preprocessing

In [16]:

```
# 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 [17]:

```
# 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 [18]:

```
# 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', "dc
esn't", 'hadn', \
```


◀ ▶

```
# Combining all the above statemennts
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['essay'].values):
    sent = decontracted(sentence)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    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 [02:
31<00:00, 329.65it/s]

```
preprocessed_essays[2000]
```

'describing students not easy task many would say inspirational creative hard working they unique unique interests learning abilities much what common desire learn day despite difficulties encounter our classroom amazing understand everyone learns pace as teacher i pride making sure students always engaged motivated inspired create learning this project help students choose seating appropriate developmentally many students tire sitting chairs lessons different seats available helps keep engaged learning flexible seating important classroom many students struggle attention focus engagement we currently stability balls seating well regular chairs stools help students trouble balance find difficult sit stability ball long period time we excited try stools part engaging classroom community nannan'

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

100%|███| 50000/50000
[00:07<00:00, 6441.63it/s]

```
preprocessed_titles[2000]
```

'steady stools active learning'

```
project_grade_categories = list(project_data['project_grade_category'].values)
# remove special characters from list of strings within:
```

```
# 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  
  
project_grade_cat_list = []  
for i in tqdm(project_grade_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 &  
unger"]  
        if 'The' in j.split(): # this will split each of the category based on space "Math & Scienc  
e">"Math","&", "Science"  
            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i  
. 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('&', '_')  
    project_grade_cat_list.append(temp.strip())
```

100% | ██████████ | 50000/50000
[00:00<00:00, 131367.17it/s]

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

Unnamed: 0	id			teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_title	project_id
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc		Mrs.	IN	2016-12-05 13:43:57	Educational Support for English Learners at Home	160221
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a		Mr.	FL	2016-10-25 09:22:10	Wanted: Projector for Hungry Learners	140945
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0		Ms.	AZ	2016-08-31 12:03:56	Soccer Equipment for AWESOME Middle School Stu...	21895
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60		Mrs.	KY	2016-10-06 21:16:17	Techie Kindergarteners	45
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec		Mrs.	TX	2016-07-11 01:10:09	Interactive Math Tools	172407

In [25]:

Out [25]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_title	pr
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Educational Support for English Learners at Home
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Wanted: Projector for Hungry Learners
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Soccer Equipment for AWESOME Middle School Stu...
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17	Techie Kindergarteners
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Interactive Math Tools

In [26]:

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

In [27]:

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

1.5 Preparing data for models

2. Logistic Regression

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

In []:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [28]:

```
project_data.columns
```

Out[28]:

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

In [29]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_val_score
from collections import Counter
from sklearn.metrics import accuracy_score
from sklearn import model_selection
```

In [30]:

```
X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'],
                                                    test_size=0.33, stratify = project_data['project_is_approved'])
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)

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

Vectorizing Categorical data

one hot encoding

In [31]:

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

```
categories_one_hot_xtrain = vectorizer_cat.transform(X_train['clean_categories'].values)
categories_one_hot_xcv = vectorizer_cat.transform(X_cv['clean_categories'].values)
categories_one_hot_xtest = vectorizer_cat.transform(X_test['clean_categories'].values)
print("Shape of matrix after one hot encodig_xtrain ", categories_one_hot_xtrain.shape)
print("Shape of matrix after one hot encodig_xcv ", categories_one_hot_xcv.shape)
print("Shape of matrix after one hot encodig_xtest ", categories_one_hot_xtest.shape)
```

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig_xtrain (22445, 9)
Shape of matrix after one hot encodig_xcv (11055, 9)
Shape of matrix after one hot encodig_xtest (16500, 9)
```

In [32]:

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

```

sub_categories_one_hot_xtrain = vectorizer_sub_cat.transform(X_train['clean_subcategories'].values)
sub_categories_one_hot_xcv = vectorizer_sub_cat.transform(X_cv['clean_subcategories'].values)
sub_categories_one_hot_xtest = vectorizer_sub_cat.transform(X_test['clean_subcategories'].values)
print("Shape of matrix after one hot encoding_xtrain ",sub_categories_one_hot_xtrain.shape)
print("Shape of matrix after one hot encoding_xcv ",sub_categories_one_hot_xcv.shape)
print("Shape of matrix after one hot encoding_xtest ",sub_categories_one_hot_xtest.shape)

```

```

['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular',
'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger',
'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other',
'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL',
'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences',
'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encoding_xtrain (22445, 30)
Shape of matrix after one hot encoding_xcv (11055, 30)
Shape of matrix after one hot encoding_xtest (16500, 30)

```

In [33]:

```

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

school_state_one_hot_xtrain = vectorizer_state.transform(X_train['school_state'].values)
school_state_one_hot_xcv = vectorizer_state.transform(X_cv['school_state'].values)
school_state_one_hot_xtest = vectorizer_state.transform(X_test['school_state'].values)
print("Shape of matrix after one hot encoding_train ",school_state_one_hot_xtrain.shape)
print("Shape of matrix after one hot encoding_cv ",school_state_one_hot_xcv.shape)
print("Shape of matrix after one hot encoding_test ",school_state_one_hot_xtest.shape)

```

```

['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DC', 'DE', 'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS',
'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']
Shape of matrix after one hot encoding_train (22445, 51)
Shape of matrix after one hot encoding_cv (11055, 51)
Shape of matrix after one hot encoding_test (16500, 51)

```

In [34]:

```

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

#https://stackoverflow.com/a/39308809/8089731
teacher_prefix_one_hot_xtrain =
vectorizer_teacherprefix.transform(X_train['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_xcv =
vectorizer_teacherprefix.transform(X_cv['teacher_prefix'].values.astype('U'))
teacher_prefix_one_hot_xtest = vectorizer_teacherprefix.transform(X_test['teacher_prefix'].values.
astype('U'))
print("Shape of matrix after one hot encoding_xtrain ",teacher_prefix_one_hot_xtrain.shape)
print("Shape of matrix after one hot encoding_xcv ",teacher_prefix_one_hot_xcv.shape)
print("Shape of matrix after one hot encoding_xtest ",teacher_prefix_one_hot_xtest.shape)

```

```

['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
Shape of matrix after one hot encoding_xtrain (22445, 5)
Shape of matrix after one hot encoding_xcv (11055, 5)
Shape of matrix after one hot encoding_xtest (16500, 5)

```

In [35]:

```

# we use count vectorizer to convert the values into one hot encoded features
from sklearn.feature_extraction.text import CountVectorizer
# https://stackoverflow.com/a/38161028/8089731
pattern = "(?m)\\b(\\w-1+\\b)"

```

```

pattern = (u/\d{1,4}\/)/
vectorizer_projectgrade = CountVectorizer(token_pattern=pattern, lowercase=False, binary=True)
vectorizer_projectgrade.fit(X_train['clean_project_grade_category'].values)
print(vectorizer_projectgrade.get_feature_names())

#https://stackoverflow.com/a/39308809/8089731
project_grade_cat_one_hot_xtrain =
vectorizer_projectgrade.transform(X_train['clean_project_grade_category'].values)
project_grade_cat_one_hot_xcv =
vectorizer_projectgrade.transform(X_cv['clean_project_grade_category'].values)
project_grade_cat_one_hot_xtest =
vectorizer_projectgrade.transform(X_test['clean_project_grade_category'].values)
print("Shape of matrix after one hot encodig_xtrain ",project_grade_cat_one_hot_xtrain.shape)
print("Shape of matrix after one hot encodig_xcv ",project_grade_cat_one_hot_xcv.shape)
print("Shape of matrix after one hot encodig_xtest ",project_grade_cat_one_hot_xtest.shape)

['Grades3-5', 'Grades6-8', 'Grades9-12', 'GradesPreK-2']
Shape of matrix after one hot encodig_xtrain (22445, 4)
Shape of matrix after one hot encodig_xcv (11055, 4)
Shape of matrix after one hot encodig_xtest (16500, 4)

```

Vectorizing Numerical features

In [36]:

```

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

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.
73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

# Now standardize the data with above maen and variance.
price_standardized_xtrain = price_scalar.transform(X_train['price'].values.reshape(-1, 1))
price_standardized_xcv = price_scalar.transform(X_cv['price'].values.reshape(-1, 1))
price_standardized_xtest = price_scalar.transform(X_test['price'].values.reshape(-1, 1))
print("shape of price_standardized_xtrain",price_standardized_xtrain.shape)
print("shape of price_standardized_xcv",price_standardized_xcv.shape)
print("shape of price_standardized_xtest",price_standardized_xtest.shape)

shape of price_standardized_xtrain (22445, 1)
shape of price_standardized_xcv (11055, 1)
shape of price_standardized_xtest (16500, 1)

```

In [37]:

```

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

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.
73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

# Now standardize the data with above mean and variance

```

```
# Now standardize the data with above mean and variance.
quantity_standardized_xtrain = quantity_scalar.transform(X_train['quantity'].values.reshape(-1, 1))
quantity_standardized_xcv = quantity_scalar.transform(X_cv['quantity'].values.reshape(-1, 1))
quantity_standardized_xtest = quantity_scalar.transform(X_test['quantity'].values.reshape(-1, 1))
print("shape of quantity_standardized_xtrain", quantity_standardized_xtrain.shape)
print("shape of quantity_standardized_xcv", quantity_standardized_xcv.shape)
print("shape of quantity_standardized_xtest", quantity_standardized_xtest.shape)
```

C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

```
shape of quantity_standardized_xtrain (22445, 1)
shape of quantity_standardized_xcv (11055, 1)
shape of quantity_standardized_xtest (16500, 1)
```

In [38]:

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

# price_standardized = standardScalar.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329. ... 399. 287.
73 5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

# Now standardize the data with above mean and variance.
teacher_num_prev_projects_standardized_xtrain = teacher_num_prev_projects_scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
teacher_num_prev_projects_standardized_xcv = teacher_num_prev_projects_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
teacher_num_prev_projects_standardized_xtest = teacher_num_prev_projects_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
print(" shape of teacher_number_of_previously_posted_projects_standardized_xtrain", teacher_num_prev_projects_standardized_xtrain.shape)
print(" shape of teacher_number_of_previously_posted_projects_standardized_xcv", teacher_num_prev_projects_standardized_xcv.shape)
print(" shape of teacher_number_of_previously_posted_projects_standardized_xtest", teacher_num_prev_projects_standardized_xtest.shape)
```

C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:

Data with input dtype int64 was converted to float64 by StandardScaler.

```
C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
C:\Users\Santosh\Anaconda3\lib\site-packages\sklearn\utils\validation.py:595:
DataConversionWarning:
```

```
Data with input dtype int64 was converted to float64 by StandardScaler.
```

```
shape of teacher_number_of_previously_posted_projects_standardized_xtrain (22445, 1)
shape of teacher_number_of_previously_posted_projects_standardized_xcv (11055, 1)
shape of teacher_number_of_previously_posted_projects_standardized_xtest (16500, 1)
```

2.2 Make Data Model Ready: encoding numerical, categorical features

```
In [ ]:
```

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Vectorizing Text data

BOW on eassay

2.3 Make Data Model Ready: encoding eassay, and project_title

```
In [ ]:
```

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

```
In [46]:
```

```
# BOW on eassay
# We are considering only the words which appeared in at least 10 documents(rows or projects).

vectorizer_bow_essays = CountVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
vectorizer_bow_essays.fit(X_train['preprocessed_essays'])

essay_text_bow_xtrain = vectorizer_bow_essays.transform(X_train['preprocessed_essays'])
essay_text_bow_xcv = vectorizer_bow_essays.transform(X_cv['preprocessed_essays'])
essay_text_bow_xtest = vectorizer_bow_essays.transform(X_test['preprocessed_essays'])
```



```
print("Shape of matrix after BOW_text_essay X_train ",essay_text_bow_xtrain.shape)
print("Shape of matrix after BOW_text_essay X_cv ",essay_text_bow_xcv.shape)
print("Shape of matrix after BOW_text_essay X_test ",essay_text_bow_xtest.shape)
```

```
Shape of matrix after BOW_text_essay X_train (22445, 5000)
Shape of matrix after BOW_text_essay X_cv (11055, 5000)
Shape of matrix after BOW_text_essay X_test (16500, 5000)
```

BOW on project_title

In [40]:

```
# BOW on project_title
# We are considering only the words which appeared in at least 10 documents(rows or projects).

vectorizer_bow_titles = CountVectorizer(min_df=10)
vectorizer_bow_titles.fit(X_train['preprocessed_titles'])

proj_title_bow_xtrain = vectorizer_bow_titles.transform(X_train['preprocessed_titles'])
proj_title_bow_xcv = vectorizer_bow_titles.transform(X_cv['preprocessed_titles'])
proj_title_bow_xtest = vectorizer_bow_titles.transform(X_test['preprocessed_titles'])

print("Shape of matrix after BOW project_title_xtrain ",proj_title_bow_xtrain.shape)
print("Shape of matrix after BOW project_title_xcv ",proj_title_bow_xcv.shape)
print("Shape of matrix after BOW project_title_xtest ",proj_title_bow_xtest.shape)
```

```
Shape of matrix after BOW project_title_xtrain (22445, 1230)
Shape of matrix after BOW project_title_xcv (11055, 1230)
Shape of matrix after BOW project_title_xtest (16500, 1230)
```

TFIDF Vectorizer on Essay

In [48]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer_tfidf_essays = TfidfVectorizer(min_df=10,max_features=5000,ngram_range=(1,2))
vectorizer_tfidf_essays.fit(X_train['preprocessed_essays'])

essay_tfidf_xtrain = vectorizer_tfidf_essays.transform(X_train['preprocessed_essays'])
essay_tfidf_xcv = vectorizer_tfidf_essays.transform(X_cv['preprocessed_essays'])
essay_tfidf_xtest = vectorizer_tfidf_essays.transform(X_test['preprocessed_essays'])

print("Shape of matrix after tfidf eassay_xtrain ",essay_tfidf_xtrain.shape)
print("Shape of matrix after tfidf essay_xcv ",essay_tfidf_xcv.shape)
print("Shape of matrix after tfidf essay_xtest ",essay_tfidf_xtest.shape)
```

```
Shape of matrix after tfidf eassay_xtrain (22445, 5000)
Shape of matrix after tfidf essay_xcv (11055, 5000)
Shape of matrix after tfidf essay_xtest (16500, 5000)
```

TFIDF Vectorizer on Project Title

In [42]:

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

proj_title_tfidf_xtrain = vectorizer_tfidf_title.transform(X_train['preprocessed_titles'])
proj_title_tfidf_xcv = vectorizer_tfidf_title.transform(X_cv['preprocessed_titles'])
proj_title_tfidf_xtest = vectorizer_tfidf_title.transform(X_test['preprocessed_titles'])

print("Shape of matrix after tfidf proj_title_xtrain ",proj_title_tfidf_xtrain.shape)
print("Shape of matrix after tfidf proj_title_xcv ",proj_title_tfidf_xcv.shape)
print("Shape of matrix after tfidf proj_title_xtest ",proj_title_tfidf_xtest.shape)
```

```
Shape of matrix after tfidf proj_title_xtrain (22445, 1230)
Shape of matrix after tfidf proj_title_xcv (11055, 1230)
Shape of matrix after tfidf proj_title_xtest (16500, 1230)
```

In [43]:

```
# Using Pretrained Models: Avg W2V
```

In [44]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-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())
```

Average Word2Vec on Essay

In [45]:

```
# average Word2Vec
# compute average word2vec for each review.

# average Word2Vec on X_train
essay_avg_w2v_vectors_xtrain = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    essay_avg_w2v_vectors_xtrain.append(vector)

print(len(essay_avg_w2v_vectors_xtrain))
print(len(essay_avg_w2v_vectors_xtrain[0]))

# average Word2Vec on X_cv
essay_avg_w2v_vectors_xcv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    essay_avg_w2v_vectors_xcv.append(vector)

print(len(essay_avg_w2v_vectors_xcv))
print(len(essay_avg_w2v_vectors_xcv[0]))

# average Word2Vec on X_test
essay_avg_w2v_vectors_xtest = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    essay_avg_w2v_vectors_xtest.append(vector)

print(len(essay_avg_w2v_vectors_xtest))
```

```
100%|███████████████████████████████████████████████████████████████| 22445/22445 [00:  
34<00:00, 655.93it/s]
```

[illegible][illegible]

Average Word2Vec on Project Title

```
# average Word2Vec

# compute average word2vec for each review.

# average Word2Vec on X_train

proj_title_avg_w2v_vectors_xtrain = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_titles']): # 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
    proj_title_avg_w2v_vectors_xtrain.append(vector)

print(len(proj_title_avg_w2v_vectors_xtrain))
print(len(proj_title_avg_w2v_vectors_xtrain[0]))

# average Word2Vec on X_cv

proj_title_avg_w2v_vectors_xcv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_titles']): # 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
    proj_title_avg_w2v_vectors_xcv.append(vector)

print(len(proj_title_avg_w2v_vectors_xcv))
print(len(proj_title_avg_w2v_vectors_xcv[0]))

# average Word2Vec on X_test

proj_title_avg_w2v_vectors_xtest = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_titles']): # 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
```

```
100%|██████████████████████████████████████████████████████████████████████████| 22445/22445  
[00:01<00:00, 11543.25it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:00<00:00, 12901.59it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:01<00:00, 13297.16it/s]
```

TFIDF weighted W2V on Essays

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

```
# average Word2Vec
# average Word2Vec on X_train

essay_tfidf_w2v_vectors_xtrain = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf 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
    essay_tfidf_w2v_vectors_xtrain.append(vector)

print(len(essay_tfidf_w2v_vectors_xtrain))
print(len(essay_tfidf_w2v_vectors_xtrain[0]))
```

```
# average Word2Vec on X_cv
essay_tfidf_w2v_vectors_xcv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf 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
    essay_tfidf_w2v_vectors_xcv.append(vector)

print(len(essay_tfidf_w2v_vectors_xcv))
print(len(essay_tfidf_w2v_vectors_xcv[0]))

# average Word2Vec on X_train
essay_tfidf_w2v_vectors_xtest = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf 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
    essay_tfidf_w2v_vectors_xtest.append(vector)

print(len(essay_tfidf_w2v_vectors_xtest))
print(len(essay_tfidf_w2v_vectors_xtest[0]))
```

```
100% |██████████████████████████████████████████████████████████████████████████| 22445/22445 [04  
:20<00:00, 86.26it/s]
```

22445
300

```
100%|███████████████████████████████████████████████████████████████████████| 11055/11055 [02  
:01<00:00, 88.33it/s]
```

11055
300

```
100%|███████████████████████████████████████████████████████████████████████████| 16500/16500 [03  
:07<00:00, 87.95it/s]
```

16500
300

TFIDF weighted W2V on Project Title

In [53]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_titles'])
# 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())
```

In [54]:

```
# TFIDF weighted W2V on Project Title
# compute average word2vec for each review.

# TFIDF weighted W2V on X_train

proj_title_tfidf_w2v_vectors_xtrain = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_titles']): # 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 tf
            idf 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
    proj_title_tfidf_w2v_vectors_xtrain.append(vector)

print(len(proj_title_tfidf_w2v_vectors_xtrain))
print(len(proj_title_tfidf_w2v_vectors_xtrain[0]))

# TFIDF weighted W2V on X_cv
proj_title_tfidf_w2v_vectors_xcv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_titles']): # 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 tf
            idf 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
    proj_title_tfidf_w2v_vectors_xcv.append(vector)

print(len(proj_title_tfidf_w2v_vectors_xcv))
print(len(proj_title_tfidf_w2v_vectors_xcv[0]))

# TFIDF weighted W2V on X_test
proj_title_tfidf_w2v_vectors_xtest = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_titles']): # 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 tf
            idf 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
    proj_title_tfidf_w2v_vectors_xtest.append(vector)

print(len(proj_title_tfidf_w2v_vectors_xtest))
print(len(proj_title_tfidf_w2v_vectors_xtest[0]))
```

22445
300

11055
300

16500
300

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

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train1=hstack((categories_one_hot_xtrain, sub_categories_one_hot_xtrain,
school_state_one_hot_xtrain,
                 teacher_prefix_one_hot_xtrain, project_grade_cat_one_hot_xtrain,
price_standardized_xtrain,
                 teacher_num_prev_projects_standardized_xtrain, quantity_standardized_xtrain,
                 essay_text_bow_xtrain, proj_title_bow_xtrain)).tocsr()

X_cv1=hstack((categories_one_hot_xcv, sub_categories_one_hot_xcv,
              school_state_one_hot_xcv, teacher_prefix_one_hot_xcv,
              project_grade_cat_one_hot_xcv, price_standardized_xcv,
              teacher_num_prev_projects_standardized_xcv, quantity_standardized_xcv,
              essay_text_bow_xcv, proj_title_bow_xcv)).tocsr()

X_test1=hstack((categories_one_hot_xtest, sub_categories_one_hot_xtest,
                school_state_one_hot_xtest, teacher_prefix_one_hot_xtest,
                project_grade_cat_one_hot_xtest, price_standardized_xtest,
                teacher_num_prev_projects_standardized_xtest, quantity_standardized_xtest,
                essay_text_bow_xtest, proj_title_bow_xtest)).tocsr()
```

```
print(X_train1.shape, y_train.shape)
print(X_cv1.shape, y_cv.shape)
print(X_test1.shape, y_test.shape)
```

```
(22445, 6332) (22445,)
(11055, 6332) (11055,)
(16500, 6332) (16500,)
```

GridSearchCV

In [58]:

```
import warnings
warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score

logireg = LogisticRegression(class_weight='balanced')
parameters = {'C':[0.0001, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 10, 100, 1000]}
clf = GridSearchCV(logireg, parameters, cv=10, scoring='roc_auc')
clf.fit(X_train1, 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']

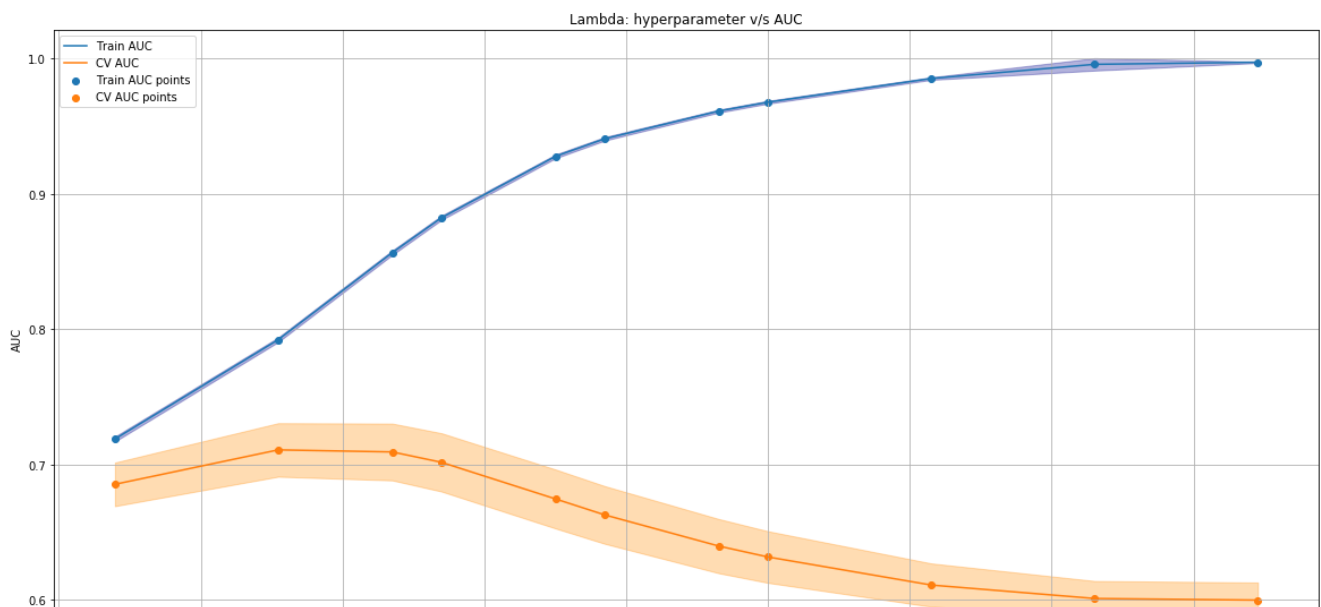
plt.figure(figsize=(20,10))
plt.plot(np.log(parameters['C']), train_auc, label='Train AUC')

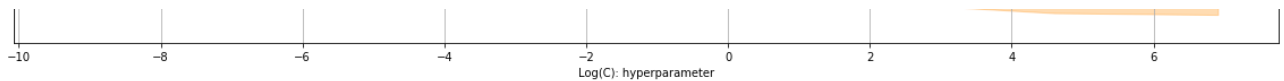
plt.gca().fill_between(np.log(parameters['C']), train_auc - train_auc_std, train_auc + train_auc_std,
alpha=0.3, color='darkblue')

plt.plot(np.log(parameters['C']), cv_auc, label='CV AUC')

plt.gca().fill_between(np.log(parameters['C']), cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')

# https://stackoverflow.com/a/48803361/4084039
plt.scatter(np.log(parameters['C']), train_auc, label='Train AUC points')
plt.scatter(np.log(parameters['C']), cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC")
plt.grid()
```





In [59]:

```
# Testing the performance of the model on test data, plotting ROC Curves
# Select best log(C) value
best_C_set_bow = clf.best_params_
print(best_C_set_bow)
```

```
{'C': 0.001}
```

Simple for loop (if you are having memory limitations use this)

In []:

```
# 4.23pm-6.07= 1.30hours
```

In [60]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive 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 [63]:

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

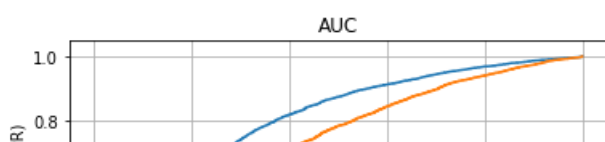
logireg = LogisticRegression(C=0.001, class_weight='balanced')
logireg.fit(X_train1[:, :], y_train[:])

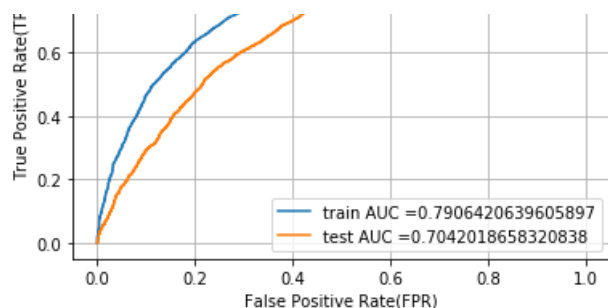
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(logireg, X_train1[:, :])
y_test_pred = batch_predict(logireg, X_test1[:])

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

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





In [64]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

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

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

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

In [65]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.42
[[ 1732  1731]
 [ 2410 16572]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.468
[[ 1337  1209]
 [ 3274 10680]]
```

In [66]:

```
# Confusion matrix for train data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

conf_matrix_xtrain = pd.DataFrame(confusion_matrix(y_train[:,], predict(y_train_pred,
tr_thresholds, train_fpr, train_tpr)))

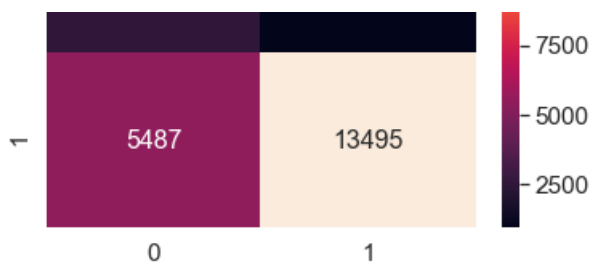
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matrix_xtrain, annot=True,annot_kws={"size": 16}, fmt='g')# font size
```

```
the maximum value of tpr*(1-fpr) 0.5193964171050934 for threshold 0.504
```

Out[66]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1f190c2f0f0>
```





In [67]:

```
# Confusion matrix for test data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

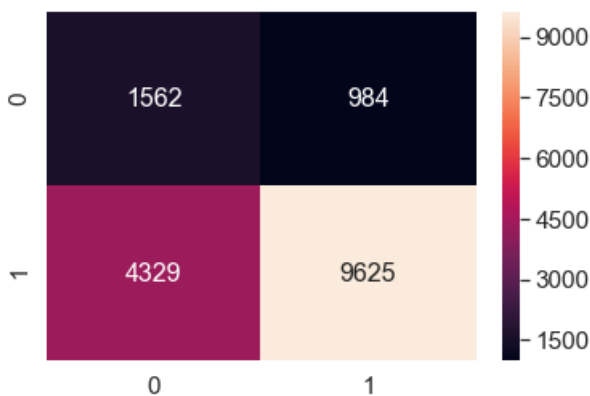
conf_matrix_xtest = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.4260662432427229 for threshold 0.501

Out [67]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f191734b70>



2.4.2 Applying Logistic Regression on TFIDF, SET 2

In [68]:

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

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train2=hstack((categories_one_hot_xtrain, sub_categories_one_hot_xtrain,
                school_state_one_hot_xtrain, teacher_prefix_one_hot_xtrain,
                project_grade_cat_one_hot_xtrain, price_standardized_xtrain,
                teacher_num_prev_projects_standardized_xtrain,
                quantity_standardized_xtrain,essay_tfidf_xtrain, proj_title_tfidf_xtrain)).tocsr()

X_cv2=hstack((categories_one_hot_xcv, sub_categories_one_hot_xcv,
              school_state_one_hot_xcv, teacher_prefix_one_hot_xcv,
              project_grade_cat_one_hot_xcv, price_standardized_xcv,
              teacher_num_prev_projects_standardized_xcv,quantity_standardized_xcv,
              essay_tfidf_xcv, proj_title_tfidf_xcv)).tocsr()

X_test2=hstack((categories_one_hot_xtest, sub_categories_one_hot_xtest,
                school_state_one_hot_xtest, teacher_prefix_one_hot_xtest,
                project_grade_cat_one_hot_xtest, price_standardized_xtest,
                teacher_num_prev_projects_standardized_xtest, quantity_standardized_xtest,
```

```

essay_tfidf_xtest, proj_title_tfidf_xtest)).tocsr()

print(X_train2.shape)
print(X_cv2.shape)
print(X_test2.shape)

```

(22445, 6332)
(11055, 6332)
(16500, 6332)

GridSearchCV

In [69]:

```

import warnings
warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score

logireg = LogisticRegression(class_weight='balanced')
parameters = {'C':[0.0001, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1]}
clf = GridSearchCV(logireg, parameters, cv=10, scoring='roc_auc')
clf.fit(X_train2, y_train)

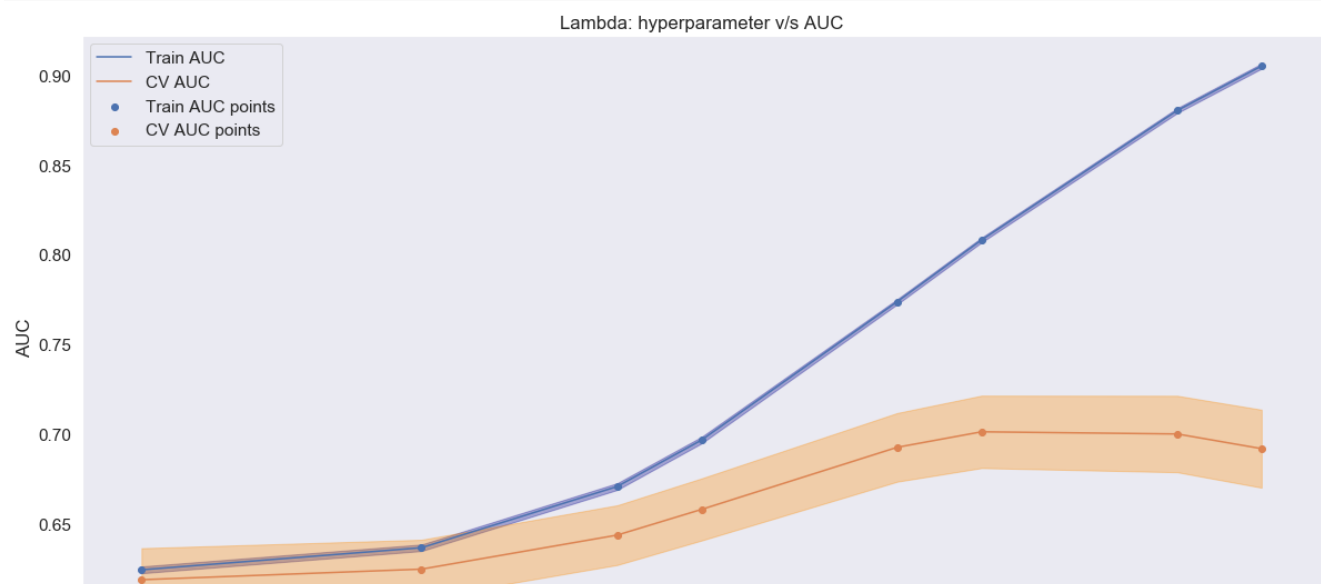
train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

plt.figure(figsize=(20,10))
plt.plot(np.log(parameters['C']), train_auc, label='Train AUC')

plt.gca().fill_between(np.log(parameters['C']), train_auc - train_auc_std, train_auc + train_auc_std,
alpha=0.3, color='darkblue')
plt.plot(np.log(parameters['C']), cv_auc, label='CV AUC')

# https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(np.log(parameters['C']), cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')
plt.scatter(np.log(parameters['C']), train_auc, label='Train AUC points')
plt.scatter(np.log(parameters['C']), cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC")
plt.grid()

```



0.60

-8

-6

-4

-2

0

Log(C): hyperparameter

In [70]:

```
# Testing the performance of the model on test data, plotting ROC Curves
# Select best log(C) value
best_C_set_tfidf = clf.best_params_
print(best_C_set_tfidf)
```

{'C': 0.1}

Simple for loop (if you are having memory limitations use this)

In [71]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive 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 [72]:

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

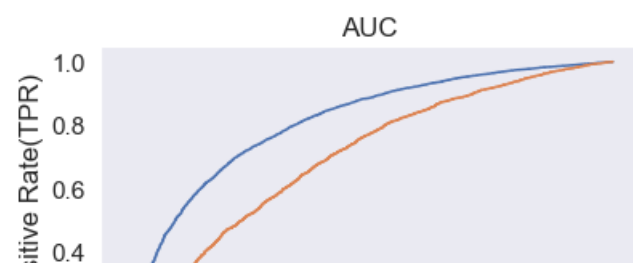
logireg = LogisticRegression(C=0.1, class_weight='balanced')
logireg.fit(X_train2[:, :], y_train[:])

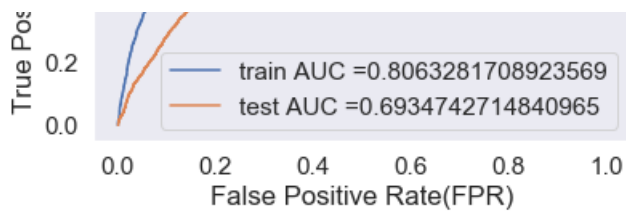
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(logireg, X_train2[:, :])
y_test_pred = batch_predict(logireg, X_test2[:])

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

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





In [73]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

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

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

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

In [74]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.5414888241768812 for threshold 0.502
[[ 2609   854]
 [ 5339 13643]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4101483822786148 for threshold 0.495
[[1473 1073]
 [4140 9814]]
```

In [75]:

```
# Confusion Matrix for Train Data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

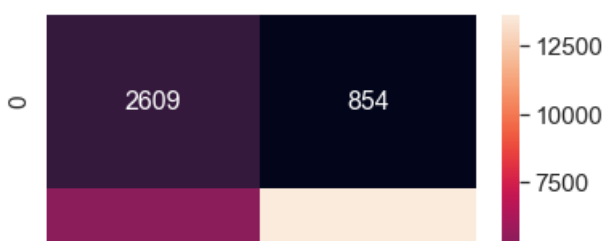
conf_matrix_xtrain = pd.DataFrame(confusion_matrix(y_train[:,], predict(y_train_pred,
tr_thresholds, train_fpr, train_tpr)))

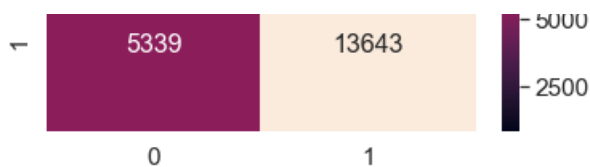
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matrix_xtrain, annot=True,annot_kws={"size": 16}, fmt='g')#font size
```

```
the maximum value of tpr*(1-fpr) 0.5414888241768812 for threshold 0.502
```

Out[75]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f19173acc0>





In [76]:

```
# Confusion matrix for test data
# Code for this segment from here -->> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

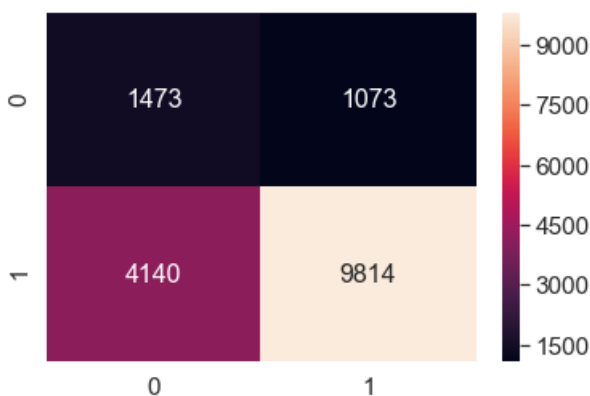
conf_matrix_xtest = pd.DataFrame(confusion_matrix(y_test[:], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.4101483822786148 for threshold 0.495

Out[76]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f1902ed978>



2.4.1 Applying Logistic Regression on AVG W2V, SET 3

In [77]:

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

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train3=hstack((categories_one_hot_xtrain, sub_categories_one_hot_xtrain,
                school_state_one_hot_xtrain, teacher_prefix_one_hot_xtrain,
                project_grade_cat_one_hot_xtrain, price_standardized_xtrain,
                teacher_num_prev_projects_standardized_xtrain,
                essay_avg_w2v_vectors_xtrain, proj_title_avg_w2v_vectors_xtrain)).tocsr()

X_cv3=hstack((categories_one_hot_xcv, sub_categories_one_hot_xcv,
              school_state_one_hot_xcv, teacher_prefix_one_hot_xcv,
              project_grade_cat_one_hot_xcv, price_standardized_xcv,
              teacher_num_prev_projects_standardized_xcv,
              essay_avg_w2v_vectors_xcv, proj_title_avg_w2v_vectors_xcv)).tocsr()

X_test3=hstack((categories_one_hot_xtest, sub_categories_one_hot_xtest,
                school_state_one_hot_xtest, teacher_prefix_one_hot_xtest,
                project_grade_cat_one_hot_xtest, price_standardized_xtest,
                teacher_num_prev_projects_standardized_xtest,
                essay_avg_w2v_vectors_xtest, proj_title_avg_w2v_vectors_xtest)).tocsr()

print(X_train3.shape, y_train.shape)
```

```
print(X_cv3.shape, y_cv.shape)
print(X_test3.shape, y_test.shape)
```

```
(22445, 701) (22445,)
(11055, 701) (11055,)
(16500, 701) (16500,)
```

GridSearchCV

In [78]:

```
import warnings
warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score

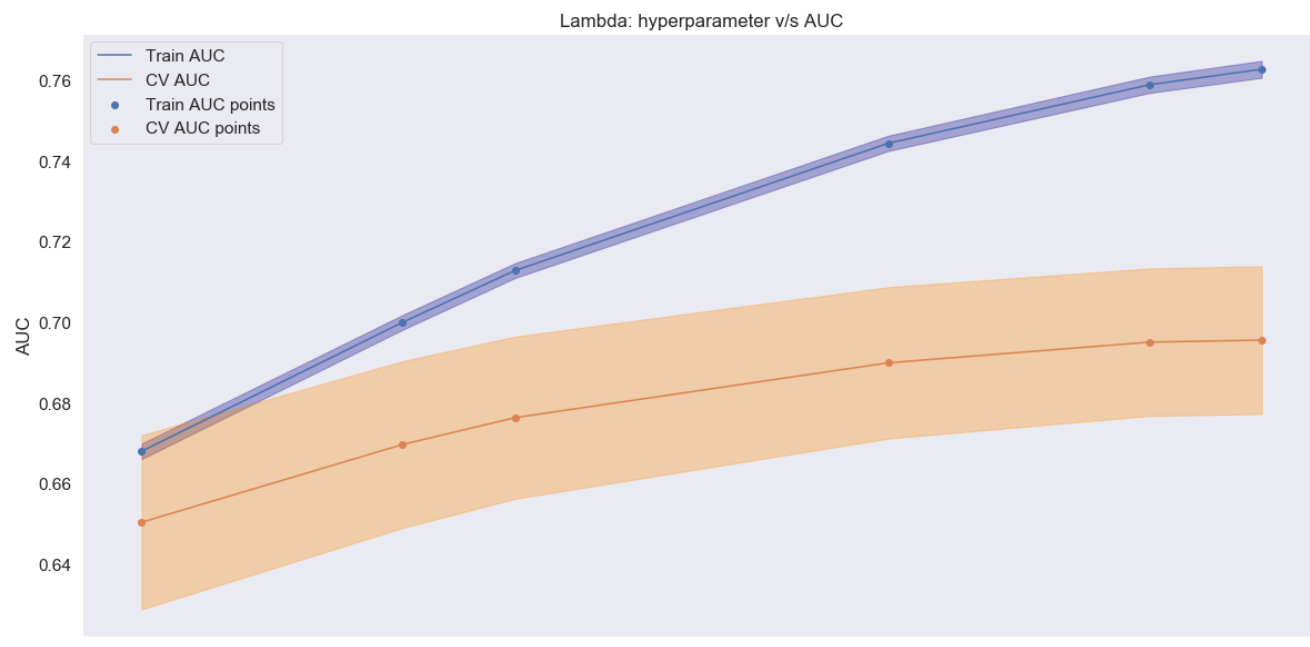
logireg = LogisticRegression(class_weight='balanced')
parameters = {'C': [0.001, 0.005, 0.01, 0.1, 0.5, 1]}
clf = GridSearchCV(logireg, parameters, cv=10, scoring='roc_auc')
clf.fit(X_train3, y_train)

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.figure(figsize=(20,10))
plt.plot(np.log(parameters['C']), train_auc, label='Train AUC')

plt.gca().fill_between(np.log(parameters['C']), train_auc - train_auc_std, train_auc + train_auc_std,
alpha=0.3, color='darkblue')
plt.plot(np.log(parameters['C']), cv_auc, label='CV AUC')

# https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(np.log(parameters['C']), cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')
plt.scatter(np.log(parameters['C']), train_auc, label='Train AUC points')
plt.scatter(np.log(parameters['C']), cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC")
plt.grid()
```



In [79]:

```
# Testing the performance of the model on test data, plotting ROC Curves
# Select best log(C) value
best_C_set_avgw2vec = clf.best_params_
print(best_C_set_avgw2vec)
```

```
{'C': 1}
```

In []:

```
# 6.34
```

Simple for loop (if you are having memory limitations use this)

In [80]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive 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 [81]:

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

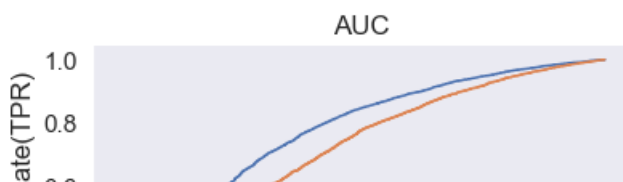
logireg = LogisticRegression(C=1, class_weight='balanced')
logireg.fit(X_train3[:, :], y_train[:])

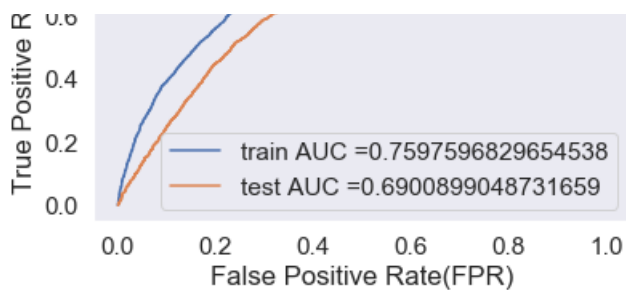
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(logireg, X_train3[:, :])
y_test_pred = batch_predict(logireg, X_test3[:, :])

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

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





In [82]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

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

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

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

In [83]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.4860366978969666 for threshold 0.487
[[ 2385  1078]
 [ 5586 13396]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4136171356880046 for threshold 0.522
[[1675  871]
 [5223 8731]]
```

In [84]:

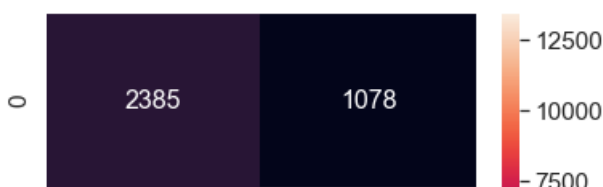
```
# Confusion Matrix for Train Data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

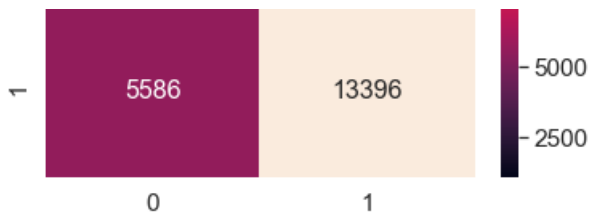
conf_matrix_xtrain = pd.DataFrame(confusion_matrix(y_train[:,], predict(y_train_pred,
tr_thresholds, train_fpr, train_tpr)))
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matrix_xtrain, annot=True,annot_kws={"size": 16}, fmt='g')
```

```
the maximum value of tpr*(1-fpr) 0.4860366978969666 for threshold 0.487
```

Out[84]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1f190ea6080>
```





In [85]:

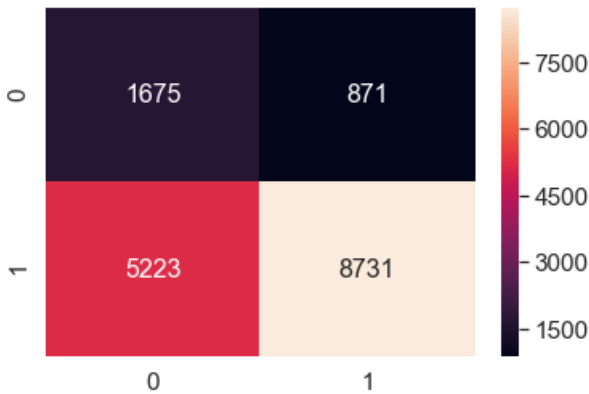
```
# Confusion matrix for test data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

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

the maximum value of $tpr \cdot (1 - fpr)$ 0.4136171356880046 for threshold 0.522

Out[85]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f191725f60>



2.4.1 Applying Logistic Regression on TFIDF Word2Vec, SET 4

In [86]:

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

# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack

X_train4=hstack((categories_one_hot_xtrain, sub_categories_one_hot_xtrain,
                 school_state_one_hot_xtrain, teacher_prefix_one_hot_xtrain,
                 project_grade_cat_one_hot_xtrain, price_standardized_xtrain,
                 teacher_num_prev_projects_standardized_xtrain, quantity_standardized_xtrain,
                 essay_tfidf_w2v_vectors_xtrain, proj_title_tfidf_w2v_vectors_xtrain)).tocsr()

X_cv4=hstack((categories_one_hot_xcv, sub_categories_one_hot_xcv,
              school_state_one_hot_xcv, teacher_prefix_one_hot_xcv,
              project_grade_cat_one_hot_xcv, price_standardized_xcv,
              teacher_num_prev_projects_standardized_xcv, quantity_standardized_xcv,
              essay_tfidf_w2v_vectors_xcv, proj_title_tfidf_w2v_vectors_xcv)).tocsr()

X_test4=hstack((categories_one_hot_xtest, sub_categories_one_hot_xtest,
                school_state_one_hot_xtest, teacher_prefix_one_hot_xtest,
                project_grade_cat_one_hot_xtest, price_standardized_xtest,
                teacher_num_prev_projects_standardized_xtest, quantity_standardized_xtest,
                essay_tfidf_w2v_vectors_xtest, proj_title_tfidf_w2v_vectors_xtest)).tocsr()
```

```
print(X_train4.shape, y_train.shape)
print(X_cv4.shape, y_cv.shape)
print(X_test4.shape, y_test.shape)
```

```
(22445, 702) (22445,)
(11055, 702) (11055,)
(16500, 702) (16500,)
```

GridSearchCV

In [87]:

```
import warnings
warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score

logireg = LogisticRegression(class_weight='balanced')
parameters = {'C': [0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1]}
clf = GridSearchCV(logireg, parameters, cv=10, scoring='roc_auc')
clf.fit(X_train4, 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']

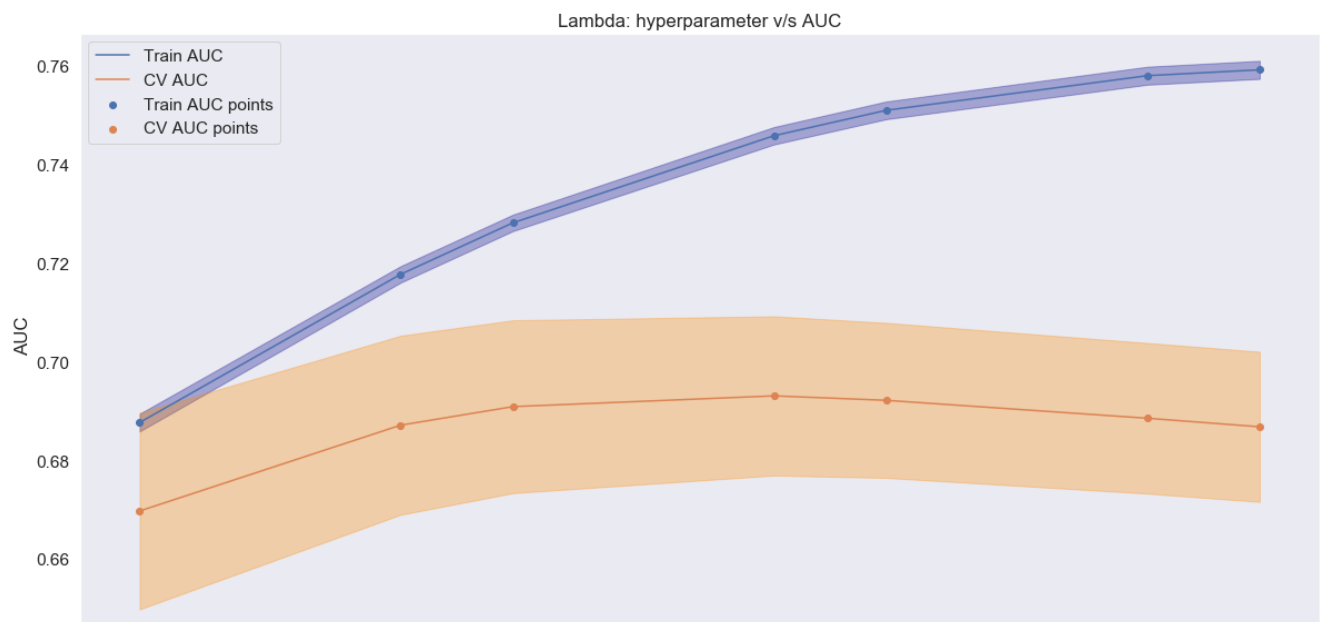
plt.figure(figsize=(20,10))
plt.plot(np.log(parameters['C']), train_auc, label='Train AUC')

plt.gca().fill_between(np.log(parameters['C']), train_auc - train_auc_std, train_auc + train_auc_std,
alpha=0.3, color='darkblue')

plt.plot(np.log(parameters['C']), cv_auc, label='CV AUC')

plt.gca().fill_between(np.log(parameters['C']), cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')

# https://stackoverflow.com/a/48803361/4084039
plt.scatter(np.log(parameters['C']), train_auc, label='Train AUC points')
plt.scatter(np.log(parameters['C']), cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC")
plt.grid()
```



In [88]:

```
# Testing the performance of the model on test data, plotting ROC Curves
# Select best log(C) value
best_C_set_tfidf2v = clf.best_params_
print(best_C_set_tfidf2v)
```

```
{'C': 0.05}
```

In [89]:

```
def batch_predict(clf, data):
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive 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 [90]:

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

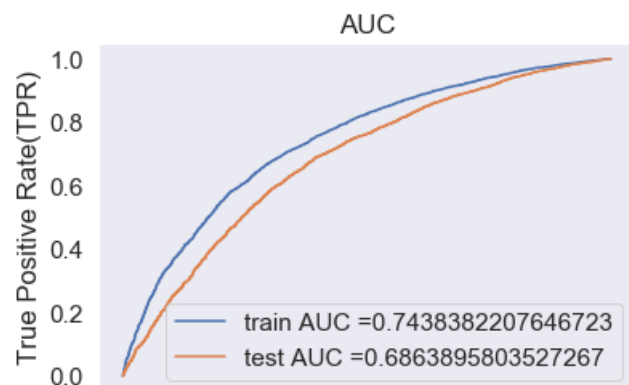
logireg = LogisticRegression(C=0.05, class_weight='balanced')
logireg.fit(X_train4[:, :], y_train[:])

# 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(logireg, X_train4[:, :])
y_test_pred = batch_predict(logireg, X_test4[:])

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

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



0.0 0.2 0.4 0.6 0.8 1.0
False Positive Rate(FPR)

In [91]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

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

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

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

In [92]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train[:,], predict(y_train_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, test_fpr, test_tpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.47083260451950876 for threshold 0.494
[[ 2408  1055]
 [ 6129 12853]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.4172781941698011 for threshold 0.504
[[1665   881]
 [5128  8826]]
```

2.5 Logistic Regression with added Features `Set 5`

New Features: Set 5

In [93]:

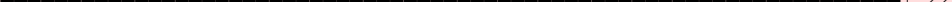
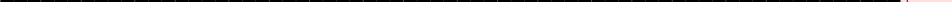
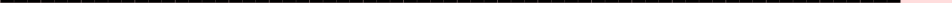
```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [131]:

```
print(X_train.shape)
print(X_cv.shape)
print(X_test.shape)
```

```
(22445, 18)
(11055, 18)
(16500, 18)
```

In [132]:

100%		22445/22445
[00:00<00:00, 102081.90it/s]		
100%		11055/11055
[00:00<00:00, 83798.00it/s]		
100%		16500/16500
[00:00<00:00, 83804.81it/s]		

In [133]:

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

```
# Essay Feature set

essay_train_wordcount = []
essay_train = list(X_train['preprocessed_essays'])
for i in tqdm(essay_train):
    b = len(str(i).split())
    essay_train_wordcount.append(b)
essay_train_wordcount = np.array(essay_train_wordcount)

essay_cv_wordcount = []
```

```

essay_cv = list(X_cv['preprocessed_essays'])
for i in tqdm(essay_cv):
    b = len(str(i).split())
    essay_cv_wordcount.append(b)
essay_cv_wordcount = np.array(essay_cv_wordcount)

essay_test_wordcount = []
essay_test = list(X_test['preprocessed_titles'])
for i in tqdm(essay_test):
    b = len(str(i).split())
    essay_test_wordcount.append(b)
essay_test_wordcount = np.array(essay_test_wordcount)

print(essay_train_wordcount.shape)
print(essay_cv_wordcount.shape)
print(essay_test_wordcount.shape)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 22445/22445
[00:03<00:00, 5806.06it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 11055/11055
[00:01<00:00, 10091.46it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:00<00:00, 92231.03it/s]

```

```

(22445,)
(11055,)
(16500,)

```

In [135]:

```

from sklearn.preprocessing import StandardScaler

essay_scalar_wordcount = StandardScaler()
essay_scalar_wordcount.fit(essay_train_wordcount.reshape(-1,1))

essay_wordcount_stand_train = essay_scalar_wordcount.transform(essay_train_wordcount.reshape(-1,1))
essay_wordcount_stand_cv = essay_scalar_wordcount.transform(essay_cv_wordcount.reshape(-1,1))
essay_wordcount_stand_test = essay_scalar_wordcount.transform(essay_test_wordcount.reshape(-1,1))

print(essay_wordcount_stand_train.shape)
print(essay_wordcount_stand_cv.shape)
print(essay_wordcount_stand_test.shape)

```

```

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

```

In [136]:

```

import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

import nltk
nltk.download('vader_lexicon')

sent = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students w
ith the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelli
gences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of differen
t backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a carin
g community of successful \
learners which can be seen through collaborative student project based learning in and out of the
classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice
a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the ki
ndergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role pla

```



```

y in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food
i will take their idea \
and create common core cooking lessons where we learn important math and writing concepts while co
oking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into maki
ng the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project woul
d expand our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make homemade apple
sauce make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create our own cook
books to be printed and \
shared with families students will gain math and literature skills as well as a life long enjoymen
t for healthy cooking \
nannan'
si = sent.polarity_scores(for_sentiment)

```

```
print(si)
```

```

lst = []
list(si.values())

```

```

[nltk_data] Downloading package vader_lexicon to
[nltk_data] C:\Users\Santosh\AppData\Roaming\nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!

```

```
{'neg': 0.01, 'neu': 0.745, 'pos': 0.245, 'compound': 0.9975}
```

```
Out[136]:
```

```
[0.01, 0.745, 0.245, 0.9975]
```

```
In [137]:
```

```

import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

import nltk
nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

# train data
essay_sentscore_xtrain = []
essay_train = list(X_train['preprocessed_essays'])
for i in tqdm(essay_train):
    ss = sid.polarity_scores(str(i))
    essay_sentscore_xtrain.append(ss)
essay_sentscore_xtrain = np.array(essay_sentscore_xtrain)

essay_negscore_xtrain = []
essay_neuscore_xtrain = []
essay_posscore_xtrain = []
essay_compoundscore_xtrain = []
for it in essay_sentscore_xtrain:
    a = it['neg']
    essay_negscore_xtrain.append(a)
    b = it['neu']
    essay_neuscore_xtrain.append(b)
    c = it['pos']
    essay_posscore_xtrain.append(c)
    d = it['compound']
    essay_compoundscore_xtrain.append(d)

essay_negscore_xtrain = np.array(essay_negscore_xtrain).reshape(-1,1)
essay_neuscore_xtrain = np.array(essay_neuscore_xtrain).reshape(-1,1)
essay_posscore_xtrain = np.array(essay_posscore_xtrain).reshape(-1,1)
essay_compoundscore_xtrain = np.array(essay_compoundscore_xtrain).reshape(-1,1)

print((essay_negscore_xtrain.shape))
print((essay_neuscore_xtrain.shape))
print((essay_posscore_xtrain.shape))
print((essay_compoundscore_xtrain.shape))

```



```
import warnings
warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt
```

```

from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import learning_curve, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import roc_auc_score

logireg = LogisticRegression(class_weight='balanced')
parameters = {'C':[ 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1]}
clf = GridSearchCV(logireg, parameters, cv= 10, scoring='roc_auc')
clf.fit(X_train5, y_train)

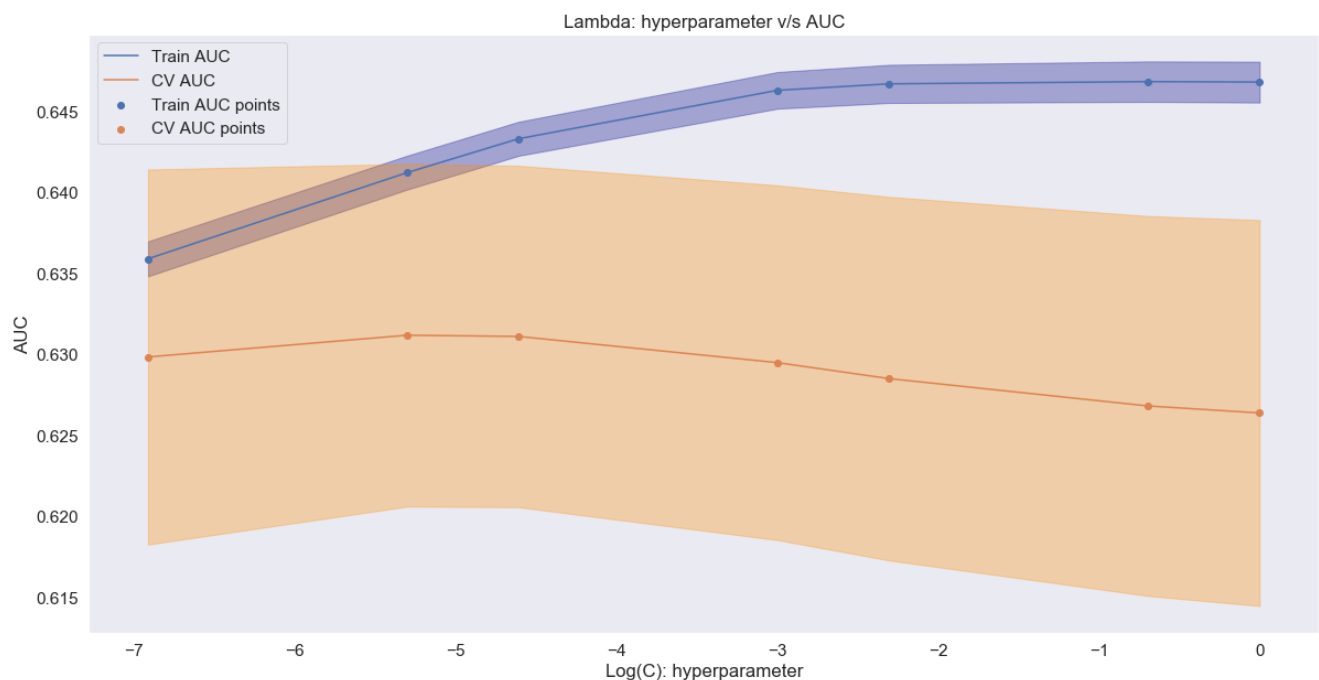
train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']
cv_auc = clf.cv_results_['mean_test_score']
cv_auc_std= clf.cv_results_['std_test_score']

plt.figure(figsize=(20,10))
plt.plot(np.log(parameters['C']), train_auc, label='Train AUC')

plt.gca().fill_between(np.log(parameters['C']),train_auc - train_auc_std,train_auc + train_auc_std,
alpha=0.3,color='darkblue')
plt.plot(np.log(parameters['C']), cv_auc, label='CV AUC')

# https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(np.log(parameters['C']),cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')
plt.scatter(np.log(parameters['C']), train_auc, label='Train AUC points')
plt.scatter(np.log(parameters['C']), cv_auc, label='CV AUC points')
plt.legend()
plt.xlabel("Log(C): hyperparameter")
plt.ylabel("AUC")
plt.title("Lambda: hyperparameter v/s AUC")
plt.grid()

```



In [142]:

```

# Testing the performance of the model on test data, plotting ROC Curves
# Select best log(C) value

best_C_set_wordcount = clf.best_params_
print(best_C_set_wordcount)

{'C': 0.005}

```

In [143]:

```

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 [144]:

```

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

logireg = LogisticRegression(C=0.005, class_weight='balanced')
logireg.fit(X_train5[:, :], y_train[:])

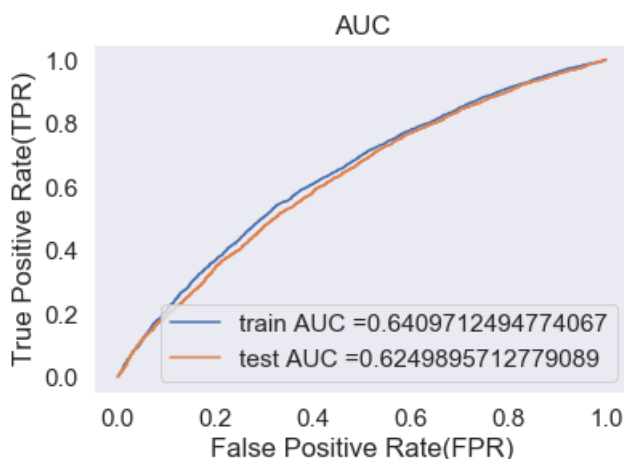
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(logireg, X_train5[:, :])
y_test_pred = batch_predict(logireg, X_test5[:])

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

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

```



In [145]:

```

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

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

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

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

```

```

        predictions.append(1)
    else:
        predictions.append(0)
    return predictions

```

In [146]:

```

from sklearn.metrics import confusion_matrix

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

```

```

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997915341 for threshold 0.476
[[ 1732  1731]
 [ 5710 13272]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 0.502
[[ 2513    33]
 [13541   413]]

```

In [147]:

```

# Confusion matrix for train data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

conf_matrix_xtrain = pd.DataFrame(confusion_matrix(y_train[:,], predict(y_train_pred,
tr_thresholds, train_fpr, train_tpr)))

sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matrix_xtrain, annot=True,annot_kws={"size": 16}, fmt='g')# font size

```

```

the maximum value of tpr*(1-fpr) 0.36887690278977003 for threshold 0.5

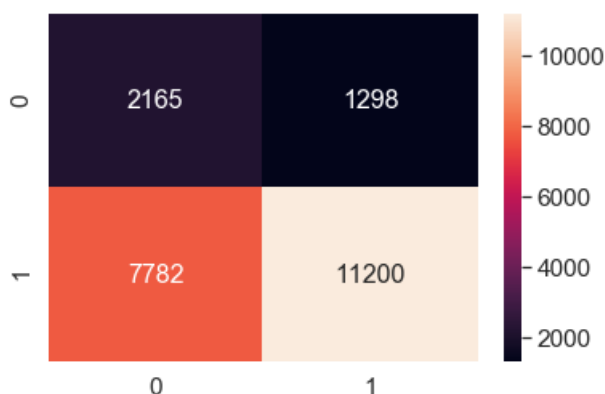
```

Out[147]:

```

<matplotlib.axes._subplots.AxesSubplot at 0x1f191716278>

```



In [149]:

```

# Confusion matrix for test data
# Code for this segment from here --> https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix

conf_matrix_xtest = pd.DataFrame(confusion_matrix(y_test[:,], predict(y_test_pred, tr_thresholds, t
est_fpr, test_tpr)))

sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matrix_xtest, annot=True,annot_kws={"size": 16}, fmt='g')#font size

```

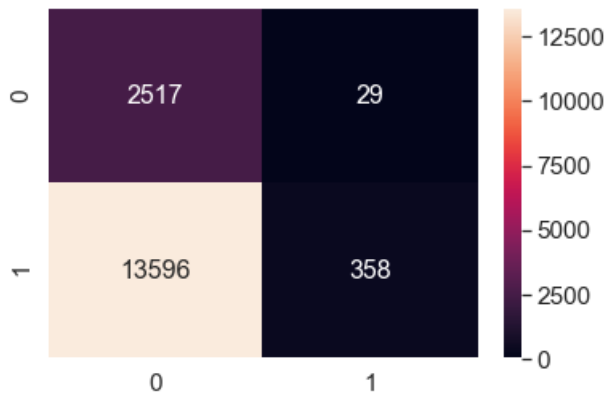
```

the maximum value of tpr*(1-fpr) 0.353005065121951 for threshold 0.52

```

Out[149]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f195d25390>



3. Conclusion

In []:

```
# Please compare all your models using Prettytable library
```

In [150]:

```
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]
x.add_row(["BOW", "Logistic Regression", 0.001, 0.70])
x.add_row(["TFIDF", "Logistic Regression", 0.1, 0.69])
x.add_row(["AVG W2V", "Logistic Regression", 1, 0.69])
x.add_row(["TFIDF W2V", "Logistic Regression", 0.05, 0.68])
x.add_row(["WORDCOUNT FEATURES", "Logistic Regression", 0.005, 0.62])
print(x)
```

Vectorizer	Model	Hyper Parameter	AUC
BOW	Logistic Regression	0.001	0.7
TFIDF	Logistic Regression	0.1	0.69
AVG W2V	Logistic Regression	1	0.69
TFIDF W2V	Logistic Regression	0.05	0.68
WORDCOUNT FEATURES	Logistic Regression	0.005	0.62

In []: