

# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

| Feature                                    |   | Description   |
|--|---|---|
| <code>project_id</code>                    |   | A unique identifier for the proposed project. <b>Example:</b> p036502   |
| <code>project_title</code>                 | <ul style="list-style-type: none"><li>•</li><li>•</li></ul>   | Title of the project. <b>Examples:</b><br><code>Art Will Make You Happy!</code><br><code>First Grade Fun</code>   |
| <code>project_grade_category</code>        | <ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li></ul>   | Grade level of students for which the project is targeted. One of the following enumerated values:<br><code>Grades PreK-2</code><br><code>Grades 3-5</code><br><code>Grades 6-8</code><br><code>Grades 9-12</code>  |
| <code>project_subject_categories</code>    | <ul style="list-style-type: none"><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li><li>•</li></ul> | One or more (comma-separated) subject categories for the project from the following enumerated list of values:<br><code>Applied Learning</code><br><code>Care &amp; Hunger</code><br><code>Health &amp; Sports</code><br><code>History &amp; Civics</code><br><code>Literacy &amp; Language</code><br><code>Math &amp; Science</code><br><code>Music &amp; The Arts</code><br><code>Special Needs</code><br><code>Warmth</code><br><br><b>Examples:</b><br><ul style="list-style-type: none"><li>• <code>Music &amp; The Arts</code></li><li>• <code>Literacy &amp; Language, Math &amp; Science</code></li></ul> |
| <code>school_state</code>                  |   | State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> WY   |
| <code>project_subject_subcategories</code> | <ul style="list-style-type: none"><li>•</li><li>•</li></ul>   | One or more (comma-separated) subject subcategories for the project. <b>Examples:</b><br><code>Literacy</code><br><code>Literature &amp; Writing, Social Sciences</code>  |
| <code>project_resource_summary</code>      | <ul style="list-style-type: none"><li>•</li></ul>   | An explanation of the resources needed for the project. <b>Example:</b><br><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. <b>Example:</b> 2016-04-28 12:43:56.245  |
| teacher_id                                   | A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56   |
| teacher_prefix                               | Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> <li>nan</li> <li>Dr.</li> <li>Mr.</li> <li>Mrs.</li> <li>Ms.</li> <li>Teacher.</li> </ul> |
| teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. <b>Example:</b> 2  |

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

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

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

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

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

| Label               | Description   |
|---------------------|---|
| 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

```

## 1.1 Reading Data

In [2]:

```

project_data = pd.read_csv('train_data.csv', nrows=50000)
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  |

## 1.2 preprocessing of project\_subject\_categories

In [5]:

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

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

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

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

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

## 1.3 preprocessing of project\_subject\_subcategories

In [6]:

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

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

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

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

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

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

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

## 1.3 Text preprocessing

In [7]:

```
# 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 [8]:

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

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             |

In [9]:

```
#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

In [10]:

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\"The limits of your language are the limits of your world.\"-Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English alongside of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills.\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\n\r\nannan

additional day or two for the year to come for each of these. (I know...)  
=====

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity. My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school. \r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still. nannan

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day. \r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas. \r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups. \r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one. \r\n\r\nIt costs a lot of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you! nannan

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. \r\n\r\nThe materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\nThey also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves. nannan

In [11]:

```
# 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"'\s're", " are", phrase)
    phrase = re.sub(r"'\s's", " is", phrase)
```

```

phrase = re.sub(r'\s', ' ', 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 [12]:

```

sent = decontracted(project_data['essay'].values[14000])
print(sent)
print("="*50)

```

Our school is located in a northern suburb of Atlanta in a diverse community. Our students come from various backgrounds, cultures, and countries. \r\n\r\nOur school is a Title I public school and my classroom alone serves around 120 students. \r\n\r\nMany of the students in my classroom are being served in one or more of the following programs: ESOL (English as a second language), Special Education, and TAG (our district is Gifted Program). The students that I have the privilege to teach are so sweet and work hard every day. My students need a class set of VR Pasonomi 3D VR Glasses to give them the opportunity to explore the world without leaving the classroom. \r\n\r\nWith the 3D VR classes, students will not just read about a location, they would experience it firsthand with a 360 degree panoramic view. Imagine virtual field-trips in Australia at Ayers Rocker, exploring the ancient Incan city of Machu Picchu or a walking through the Amazon Rainforest. When we are studying about Europe, I will have my students gaze up at and walk around cities and historical landmarks. While they are 'visiting' these wonders, students can hear commentary about what they are seeing (with certain applications) or even create their own commentary, without the thousands of dollars it would take to fly to these areas and be on time to math the next period. \r\n\r\nHaving this class set available to my students would definitely take my class to the next level and open students' eyes to the world around them without having to leave the classroom. Students will be able to dig deeper into any subject or location that we are studying in social studies. \r\n\r\nI can imagine conversations about the locations we visit- "is that really where the Aztecs would perform human sacrifice!?" I am excited at the uses of this technology and know that my students would thrive on our virtual trips! \r\n\r\nnnannan

=====

In [13]:

```

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

```

Our school is located in a northern suburb of Atlanta in a diverse community. Our students come from various backgrounds, cultures, and countries. Our school is a Title I public school and my classroom alone serves around 120 students. Many of the students in my classroom are being served in one or more of the following programs: ESOL (English as a second language), Special Education, and TAG (our district is Gifted Program). The students that I have the privilege to teach are so sweet and work hard every day. My students need a class set of VR Pasonomi 3D VR Glasses to give them the opportunity to explore the world without leaving the classroom. With the 3D VR classes, students will not just read about a location, they would experience it firsthand with a 360 degree panoramic view. Imagine virtual field-trips in Australia at Ayers Rocker, exploring the ancient Incan city of Machu Picchu or a walking through the Amazon Rainforest. When we are studying about Europe, I will have my students gaze up at and walk around cities and historical landmarks. While they are 'visiting' these wonders, students can hear commentary about what they are seeing (with certain applications) or even create their own commentary, without the thousands of dollars it would take to fly to these areas and be on time to math the next period. Having this class set available to my students would definitely take my class to the next level and open students' eyes to the world around them without having to leave the classroom. Students will be able to dig deeper into any subject or location that we are studying in social studies. I can imagine conversations about the locations we visit- "is that really where the Aztecs would perform human sacrifice!?" I am excited at the uses of this technology and know that my students would thrive on our virtual trips! nannan

In [14]:

```

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

```

Our school is located in a northern suburb of Atlanta in a diverse community Our students come from various backgrounds cultures and countries Our school is a Title I public school and my classroo

In various backgrounds cultures and countries our school is a little 1 public school and my classroom alone serves around 120 students Many of the students in my classroom are being served in one or more of the following programs ESOL English as a second language Special Education and TAG our district is Gifted Program The students that I have the privilege to teach are so sweet and work hard every day My students need a class set of VR Pasonomi 3D VR Glasses to give them the opportunity to explore the world without leaving the classroom With the 3D VR classes students will not just read about a location they would experience it firsthand with a 360 degree panoramic view Imagine virtual field trips in Australia at Ayers Rocker exploring the ancient Incan city of Machu Picchu or a walking through the Amazon Rainforest When we are studying about Europe I will have my students gaze up at and walk around cities and historical landmarks While they are visiting these wonders students can hear commentary about what they are seeing with certain applications or even create their own commentary without the thousands of dollars it would take to fly to these areas and be on time to math the next period Having this class set available to my students would definitely take my class to the next level and open students eyes to the world around them without having to leave the classroom Students will be able to dig deeper into any subject or location that we are studying in social studies I can imagine conversations about the locations we visit is that really where the Aztecs would perform human sacrifice I am excited at the uses of this technology and know that my students would thrive on our virtual trips nannan

In [15]:

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

In [16]:

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

```
100% |██████████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:40<00:00, 1219.57it/s]
```

In [17]:

```
# after preprocessing
```



```
preprocessed_essays[14000]
```

Out[17]:

```
'our school located northern suburb atlanta diverse community our students come various
backgrounds cultures countries our school title i public school classroom alone serves around 120
students many students classroom served one following programs esol english second language specia
l education tag district gifted program the students i privilege teach sweet work hard every day m
y students need class set vr pasonomi 3d vr glasses give opportunity explore world without leaving
classroom with 3d vr classes students not read location would experience firsthand 360 degree pano
ramic view imagine virtual field trips australia ayers rocker exploring ancient incan city machu
picchu walking though amazon rainforest when studying europe i students gaze walk around cities hi
storical landmarks while visiting wonders students hear commentary seeing certain applications eve
n create commentary without thousands dollars would take fly areas time math next period having cl
ass set available students would definitely take class next level open students eyes world around
without leave classroom students able dig deeper subject location studying social studies i imagin
e conversations locations visit really aztecs would perform human sacrifice i excited uses
technology know students would thrive virtual trips nannan'
```

## 1.4 Preprocessing of `project\_title`

In [18]:

```
# similarly you can preprocess the titles also
def preprocess_text_func(text_data):
    sent = decontracted(text_data)
    sent = sent.replace('\\r', ' ')
    sent = sent.replace('\\\"', ' ')
    sent = sent.replace('\\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    return sent.lower()
```

In [19]:

```
preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    preprocessed_titles.append(preprocess_text_func(sentence))
project_data['project_title']=preprocessed_titles
```

```
100% |████████████████████████████████████████████████████████████████████████████████| 50000/50000
[00:02<00:00, 24686.81it/s]
```

## 1.5 Preparing data for models

In [20]:

```
project_data.columns
```

Out[20]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category', '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',
      'clean_categories', 'clean_subcategories', 'essay'],
      dtype='object')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data

```

- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

```

## 1.5.1 Vectorizing Categorical data

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

In [21]:

```

# we use count vectorizer to convert the values into one
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
categories_one_hot = vectorizer.fit_transform(project_data['clean_categories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",categories_one_hot.shape)

```

```

['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig  (50000, 9)

```

In [22]:

```

# we use count vectorizer to convert the values into one
vectorizer = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
sub_categories_one_hot = vectorizer.fit_transform(project_data['clean_subcategories'].values)
print(vectorizer.get_feature_names())
print("Shape of matrix after one hot encodig ",sub_categories_one_hot.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 encodig  (50000, 30)

```

In [23]:

```

# you can do the similar thing with state, teacher_prefix and project_grade_category also
def perform_one_hot_encoding(listdata, category,fillnan_value=""):
    vectorizer = CountVectorizer(vocabulary=listdata, lowercase=False, binary=True)
    vectorizer.fit(project_data[category].fillna(fillnan_value).values)
    print(vectorizer.get_feature_names())
    print("="*50)
    return vectorizer.transform(project_data[category].fillna(fillnan_value).values)

```

In [24]:

```

# One hot encoding for school state
countries_list = sorted(project_data["school_state"].value_counts().keys())
school_state_one_hot = perform_one_hot_encoding(countries_list, "school_state")
print("Shape of matrix after one hot encodig ",school_state_one_hot.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 encodig  (50000, 51)

```

In [25]:

```
# Project_Grade_Category - replacing hyphens, spaces with Underscores
project_data['project_grade_category'] = project_data['project_grade_category'].map({'Grades PreK-2': 'Grades_PreK_2',
                                                                                       'Grades 6-8' :
grades_6_8',
                                                                                       'Grades 3-5' :
grades_3_5',
                                                                                       'Grades 9-12'
Grades_9_12'})
project_data['teacher_prefix'] = project_data['teacher_prefix'].map({'Mrs.': 'Mrs', 'Ms.': 'Ms', 'Mr.' : 'Mr',
                                                                                       'Teacher': 'Teacher', 'Dr.' :
'Dr'})
```

In [26]:

```
# Replacing Null values with most repititive values
project_data["teacher_prefix"].fillna("Mrs", inplace=True)
# One hot encoding for teacher_prefix
teacher_prefix_list = sorted(project_data["teacher_prefix"].value_counts().keys())
print (teacher_prefix_list)
teacher_prefix_one_hot = perform_one_hot_encoding(teacher_prefix_list, "teacher_prefix", "Mrs.")
print("Shape of matrix after one hot encodig ",teacher_prefix_one_hot.shape)
```

```
['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
['Dr', 'Mr', 'Mrs', 'Ms', 'Teacher']
=====
Shape of matrix after one hot encodig (50000, 5)
```

In [27]:

```
# we use count vectorizer to convert the values into one hot encoded features for
project_grade_category
grade_list = sorted(project_data["project_grade_category"].value_counts().keys())
vectorizer = CountVectorizer(vocabulary=list(grade_list), lowercase=False, binary=True)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())
```

```
grade_one_hot = vectorizer.transform(project_data['project_grade_category'].values)
print("Shape of matrix after one hot encodig ",grade_one_hot.shape)
```

```
['Grades_3_5', 'Grades_6_8', 'Grades_9_12', 'Grades_PreK_2']
Shape of matrix after one hot encodig (50000, 4)
```

## 1.5.2 Vectorizing Text data

### 1.5.2.1 Bag of words

In [28]:

```
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encodig ",text_bow.shape)
```

Shape of matrix after one hot encodig (50000, 12211)

In [29]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
vectorizer_titles = CountVectorizer(min_df=10)
text_bow_titles = vectorizer_titles.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encodig ",text bow titles.shape)
```

```
bow_titles_feature_names = vectorizer.get_feature_names()
```

Shape of matrix after one hot encoding (50000, 2135)

### 1.5.2.2 TFIDF vectorizer

In [30]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
text_tfidf = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot encoding ", text_tfidf.shape)
```

Shape of matrix after one hot encoding (50000, 12211)

In [31]:

```
# TFIDF Vectorizer for Preprocessed Title
vectorizer_title = TfidfVectorizer(min_df=10)
text_tfidf_title = vectorizer_title.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encoding ", text_tfidf_title.shape)
```

Shape of matrix after one hot encoding (50000, 2135)

### 1.5.2.3 Using Pretrained Models: Avg W2V

In [32]:

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

# =====
Output:

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

# =====

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

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

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

words_corpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_corpus[i] = model[i]
    else:
        words_corpus[i] = np.zeros(300)
'''
```

///

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\'r\',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\nword = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\nodel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\'glove.42B.300d.txt\')\n\n# =====\n\nOutput:\n\nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split(\'
\'))\n\nfor i in preprocod_titles:\n    words.extend(i.split(\' \'))\nprint("all the words in the
coupus", len(words))\nwords = set(words)\nprint("the unique words in the coupus",
len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words tha
t are present in both glove vectors and our coupus", len(inter_words),"
(",np.round(len(inter_words)/len(words)*100,3),"%")\n\nwords_courpus = {}\nwords_glove =
set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\nwith open(\'glove vectors\', \'wb\') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```

```
# 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
# compute average word2vec for each review.
avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(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
    avg_w2v_vectors.append(vector)

print(len(avg_w2v_vectors))
print(len(avg_w2v_vectors[0]))
```

50000  
300

### 1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(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
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(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
    tfidf_w2v_vectors.append(vector)

print(len(tfidf_w2v_vectors))
print(len(tfidf_w2v_vectors[0]))
```

50000  
300

```
# Similarly you can vectorize for title also
tfidf_model = TfidfVectorizer()
tfidf_model.fit(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())
```

```
# Similarly you can vectorize for title also
tfidf_w2v_vectors_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(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
    tfidf_w2v_vectors_title.append(vector)

print(len(tfidf_w2v_vectors_title))
print(len(tfidf_w2v_vectors_title[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:02<00:00. 23776.74it+ /s]
```

```
50000
300
```

### 1.5.3 Vectorizing Numerical features

In [39]:

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

In [40]:

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

# Now standardize the data with above mean and variance.
price_standardized = price_scaler.transform(project_data['price'].values.reshape(-1, 1))
```

Mean : 299.33367619999996, Standard deviation : 378.20927190421384

In [41]:

```
price_standardized
```

Out[41]:

```
array([[ -0.38268146],
       [ -0.00088225],
       [  0.57512161],
       ...,
       [-0.65382764],
       [-0.52109689],
       [  0.54492668]])
```

In [42]:

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

# Now standardize the data with above mean and variance.
teacher_number_of_previously_posted_projects_standardized =
teacher_number_of_previously_posted_projects_scaler.transform(project_data['teacher_number_of_previ
ously_posted_projects'].values.reshape(-1, 1))
```

Mean : 11.24934, Standard deviation : 28.15684303263418

### 1.5.4 Merging all the above features

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

In [43]:

```
print(categories_one_hot.shape)
print(sub_categories_one_hot.shape)
print(text_bow.shape)
print(price_standardized.shape)
```

```
(50000, 9)
(50000, 30)
(50000, 12211)
(50000, 1)
```

In [44]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
# with the same hstack function we are concatinating a sparse matrix and a dense matirx :)
X = hstack((school_state_one_hot, categories_one_hot, sub_categories_one_hot,
teacher_prefix_one_hot,
            grade_one_hot, text_bow, price_standardized))
X.shape
```

Out[44]:

```
(50000, 12311)
```

## Computing Sentiment Scores

In [45]:

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

# import nltk
# nltk.download('vader_lexicon')

sid = 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'
ss = sid.polarity_scores(for_sentiment)
```



```
for k in ss:
    print('{0}: {1}'.format(k, ss[k]), end='')

# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
```

neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,

In [131]:

```
from textblob import TextBlob
project_data['essay_sentiment'] = [ TextBlob(tb).sentiment.polarity for tb in project_data['essay']
]
```

## 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](#).

## 2. Logistic Regression

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

In [132]:

```
# 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
y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33)
```

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

In [133]:

```
# 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 [134]:

```
# School State
vectorizer = CountVectorizer()
vectorizer.fit(X_train['school_state'].values)

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

In [135]:

```
# teacher_prefix
vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values)

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

In [136]:

In [136]:

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

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

In [137]:

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

# we use the fitted CountVectorizer to convert the text to vector
X_train_category_ohe = vectorizer.transform(X_train['clean_categories'].values)
X_cv_category_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_category_ohe = vectorizer.transform(X_test['clean_categories'].values)
```

In [138]:

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

# we use the fitted CountVectorizer to convert the text to vector
X_train_subcategory_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
X_cv_subcategory_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_subcategory_ohe = vectorizer.transform(X_test['clean_subcategories'].values)
```

In [139]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
normalizer.fit(X_train['price'].values.reshape(-1,1))
X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(-1,1))
X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(-1,1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(-1,1))
```

In [140]:

```
# teacher previously posted projects
normalizer = Normalizer()
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_train_teach_prev_norm =
normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_cv_teach_prev_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_test_teach_prev_norm =
normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
```

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

In [141]:

```
# 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 [142]:

```
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['essay'].values)
```

Out[142]:

```
CountVectorizer(analyzer='word', binary=False, decode_error='strict',
                dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
                lowercase=True, max_df=1.0, max_features=5000, min_df=10,
                ngram_range=(1, 4), preprocessor=None, stop_words=None,
                strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
                tokenizer=None, vocabulary=None)
```

In [143]:

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

In [144]:

```
# Preprocessing project_title
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['project_title'].values)
```

Out[144]:

```
CountVectorizer(analyzer='word', binary=False, decode_error='strict',
                dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
                lowercase=True, max_df=1.0, max_features=5000, min_df=10,
                ngram_range=(1, 4), preprocessor=None, stop_words=None,
                strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
                tokenizer=None, vocabulary=None)
```

In [145]:

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_pj_title_bow = vectorizer.transform(X_train['project_title'].values)
X_cv_pj_title_bow = vectorizer.transform(X_cv['project_title'].values)
X_test_pj_title_bow = vectorizer.transform(X_test['project_title'].values)
```

In [146]:

```
# We are considering only the words which appeared in at least 10 documents (rows or projects).
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['project_resource_summary'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_resource_summary_bow = vectorizer.transform(X_train['project_resource_summary'].values)
X_cv_project_resource_summary_bow = vectorizer.transform(X_cv['project_resource_summary'].values)
X_test_project_resource_summary_bow = vectorizer.transform(X_test['project_resource_summary'].values)
```

## 2.4 Applying Logistic Regression on different kind of featurization as mentioned in the instructions

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

In [147]:

```
# 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
```

```

# Reading and understanding error messages will be very much helpful 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 [148]:

```
# set1:
```

In [149]:

```

#function to get heatmap confusion matrix
def get_confusion_matrix(clf,X_te,y_test):
    y_pred = clf.predict(X_te)
    df_cm = pd.DataFrame(confusion_matrix(y_test, y_pred), range(2),range(2))
    df_cm.columns = ['Predicted NO','Predicted YES']
    df_cm = df_cm.rename({0: 'Actual NO', 1: 'Actual YES'})
    sns.set(font_scale=1.4)#for label size
    sns.heatmap(df_cm, annot=True,annot_kws={"size": 16}, fmt='g')

```

In [150]:

```

from scipy.sparse import hstack
X_tr = hstack((X_train_essay_bow, X_train_state_ohe, X_train_teacher_ohe,
               X_train_grade_ohe, X_train_price_norm, X_train_category_ohe,
               X_train_subcategory_ohe, X_train_teach_prev_norm,
               X_train_pj_title_bow)).tocsr()

X_cr = hstack((X_cv_essay_bow, X_cv_state_ohe, X_cv_teacher_ohe,
               X_cv_grade_ohe, X_cv_category_ohe, X_cv_subcategory_ohe,
               X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_bow)).tocsr()

X_te = hstack((X_test_essay_bow, X_test_state_ohe, X_test_teacher_ohe,
               X_test_grade_ohe, X_test_category_ohe, X_test_subcategory_ohe,
               X_test_price_norm, X_test_teach_prev_norm,
               X_test_pj_title_bow)).tocsr()

```

In [151]:

```

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
    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
    if data.shape[0]%1000 !=0:
        y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])

    return y_data_pred

```

In [152]:

```

from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

```

In [153]:

```

lr = LogisticRegression(penalty='l1',class_weight='balanced')

parameters = {'C':[ 0.25,0.05,0.025, 0.01, 0.005,0.0025, 0.004, 0.003,0.001]}

clf = GridSearchCV(lr, parameters, cv= 3, scoring='roc_auc',return_train_score=True,verbose=2)

clf.fit(X_tr, y_train)

train_auc= clf.cv_results_['mean_train_score']

```

```

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(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("C: hyperparameter v/s AUC plot")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```

[CV] C=0.25 .....
[CV] ..... C=0.25, total= 9.4s

```

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 9.3s remaining: 0.0s

```

[CV] C=0.25 .....
[CV] ..... C=0.25, total= 11.1s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 4.9s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 1.7s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 1.4s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 1.5s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.8s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 1.0s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.8s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.6s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.6s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.6s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.5s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.4s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.4s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.4s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.4s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.4s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.4s
[CV] C=0.004 .....

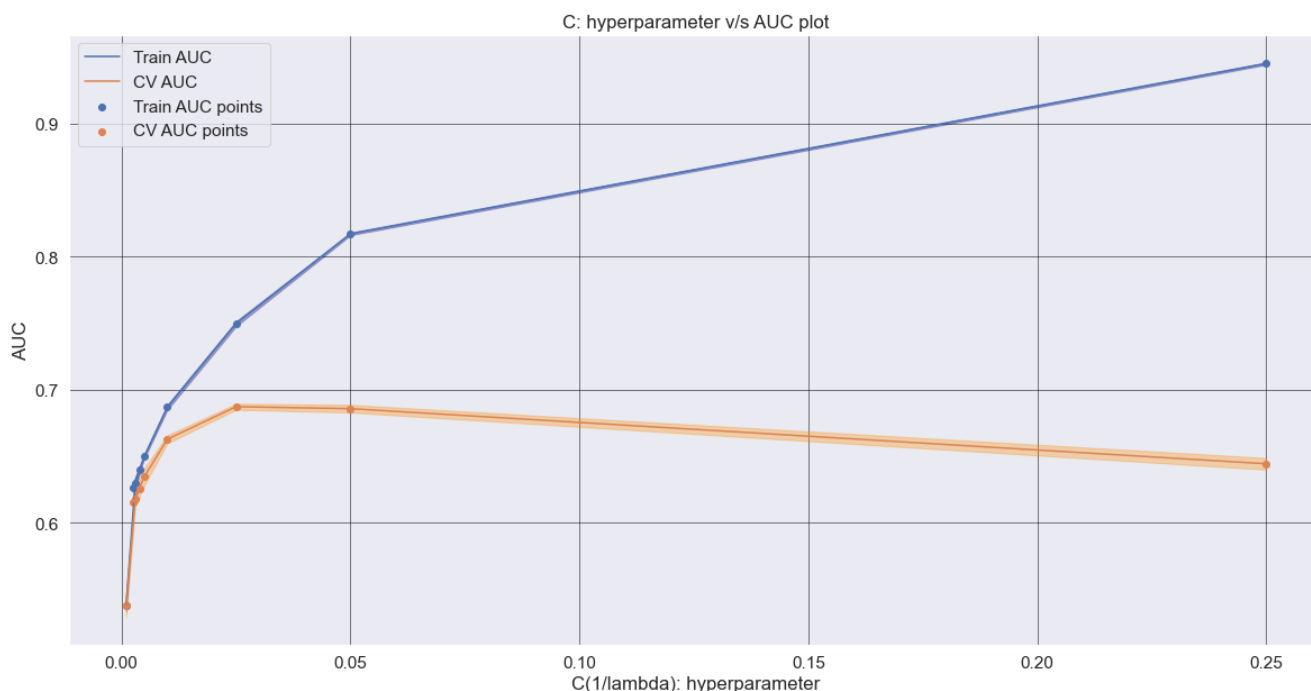
```

```

[CV] ..... C=0.004, total= 0.4s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.5s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.5s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.5s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.4s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.2s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.3s

```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 41.5s finished
```



In [154]:

```

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

model = LogisticRegression(C = 0.05,penalty='l1',class_weight='balanced')

model.fit(X_tr, y_train)

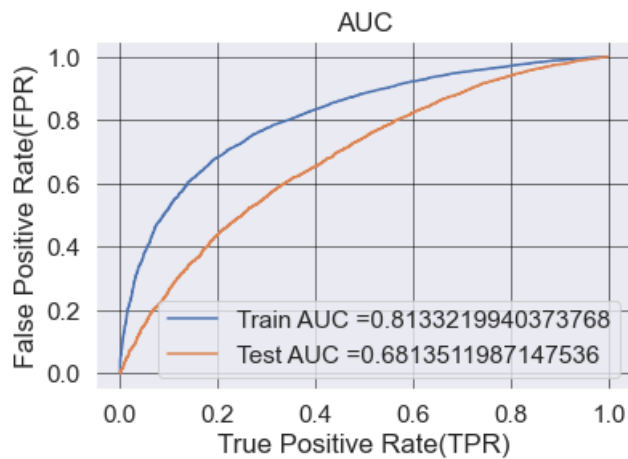
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

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

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```



In [155]:

```
def predict(proba, threshold, fpr, tpr):
    t = threshold[np.argmax(fpr*(1-tpr))]

    # (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 [156]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.373
[[ 1709  1710]
 [ 2206 16820]]
```

In [157]:

```
conf_matr_df_train_1 = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
train_fpr, train_fpr)), range(2),range(2))
```

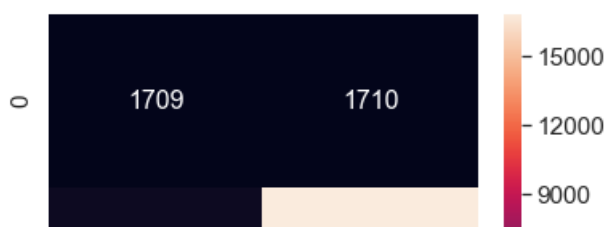
```
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.373
```

In [158]:

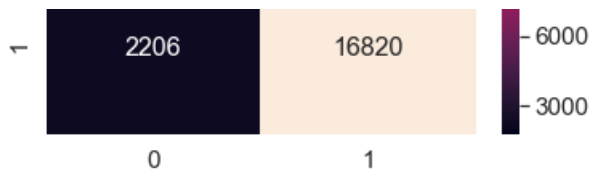
```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train_1, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[158]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1e73a492d48>
```







In [159]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999984572938835 for threshold 0.431  
[[ 744 1802]  
 [1508 12446]]

In [160]:

```
conf_matr_df_test_1 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2), range(2))
```

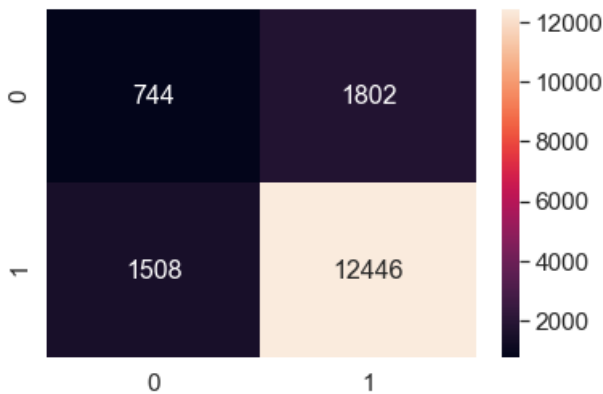
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999984572938835 for threshold 0.431

In [161]:

```
sns.set(font_scale=1.4) #for label size
sns.heatmap(conf_matr_df_test_1, annot=True, annot_kws={"size": 16}, fmt='g')
```

Out[161]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e731d92d88>



In [162]:

```
# set2:
```

In [163]:

```
# preprocessing TFIDF of Text Essays and Project Titles
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train["essay"].values)
X_train_essay_tfidf = vectorizer.transform(X_train['essay'].values)
X_cv_essay_tfidf = vectorizer.transform(X_cv['essay'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)

print("Shape of Datamatrix after TFIDF Vectorization")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

Shape of Datamatrix after TFIDF Vectorization

```
(22445, 5000) (22445,)
(11055, 5000) (11055,)
(16500, 5000) (16500,)
```

In [164]:

```
# Similarly you can vectorize for title also
vectorizer_titles = TfidfVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer_titles.fit(X_train["project_title"])

X_train_pj_title_tfidf = vectorizer.transform(X_train['project_title'].values)
X_cv_pj_title_tfidf = vectorizer.transform(X_cv['project_title'].values)
X_test_pj_title_tfidf = vectorizer.transform(X_test['project_title'].values)

print("Shape of Datamatrix after TFIDF Vectorization")
print(X_train_pj_title_tfidf.shape, y_train.shape)
print(X_cv_pj_title_tfidf.shape, y_cv.shape)
print(X_test_pj_title_tfidf.shape, y_test.shape)
print("="*100)
```

Shape of Datamatrix after TFIDF Vectorization

```
(22445, 5000) (22445,)
(11055, 5000) (11055,)
(16500, 5000) (16500,)
```

In [165]:

```
# Concatinating all the features for Set 2

X_tr = hstack((X_train_essay_tfidf, X_train_state_ohe, X_train_teacher_ohe,
               X_train_grade_ohe, X_train_price_norm, X_train_category_ohe,
               X_train_subcategory_ohe, X_train_teach_prev_norm,
               X_train_pj_title_tfidf)).tocsr()

X_cr = hstack((X_cv_essay_tfidf, X_cv_state_ohe, X_cv_teacher_ohe,
               X_cv_grade_ohe, X_cv_category_ohe, X_cv_subcategory_ohe,
               X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_tfidf)).tocsr()

X_te = hstack((X_test_essay_tfidf, X_test_state_ohe, X_test_teacher_ohe,
               X_test_grade_ohe, X_test_category_ohe, X_test_subcategory_ohe,
               X_test_price_norm, X_test_teach_prev_norm,
               X_test_pj_title_tfidf)).tocsr()
```

In [166]:

```
lr = LogisticRegression(penalty='l1', class_weight="balanced")

parameters = {'C':[1,0.5,0.25,0.05,0.025, 0.01, 0.005,0.0025, 0.004, 0.003,0.001]}

clf = GridSearchCV(lr, parameters, cv= 3, scoring='roc_auc',return_train_score=True,verbose=2)

clf.fit(X_tr, y_train)

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

plt.figure(figsize=(20,10))

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='dar
```

```

korange')

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

plt.legend()
plt.xlabel("C (1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("C: hyperparameter v/s AUC plot")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()

```

Fitting 3 folds for each of 11 candidates, totalling 33 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```

[CV] C=1 .....
[CV] ..... C=1, total= 6.0s

```

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 6.0s remaining: 0.0s

```

[CV] C=1 .....
[CV] ..... C=1, total= 5.3s
[CV] C=1 .....
[CV] ..... C=1, total= 5.3s
[CV] C=0.5 .....
[CV] ..... C=0.5, total= 2.0s
[CV] C=0.5 .....
[CV] ..... C=0.5, total= 2.0s
[CV] C=0.5 .....
[CV] ..... C=0.5, total= 2.3s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 0.9s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 0.9s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 0.9s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.5s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.5s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.6s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.5s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.4s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.4s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.4s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.4s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.2s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.3s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.2s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.2s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.2s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.2s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.2s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.2s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.2s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.2s

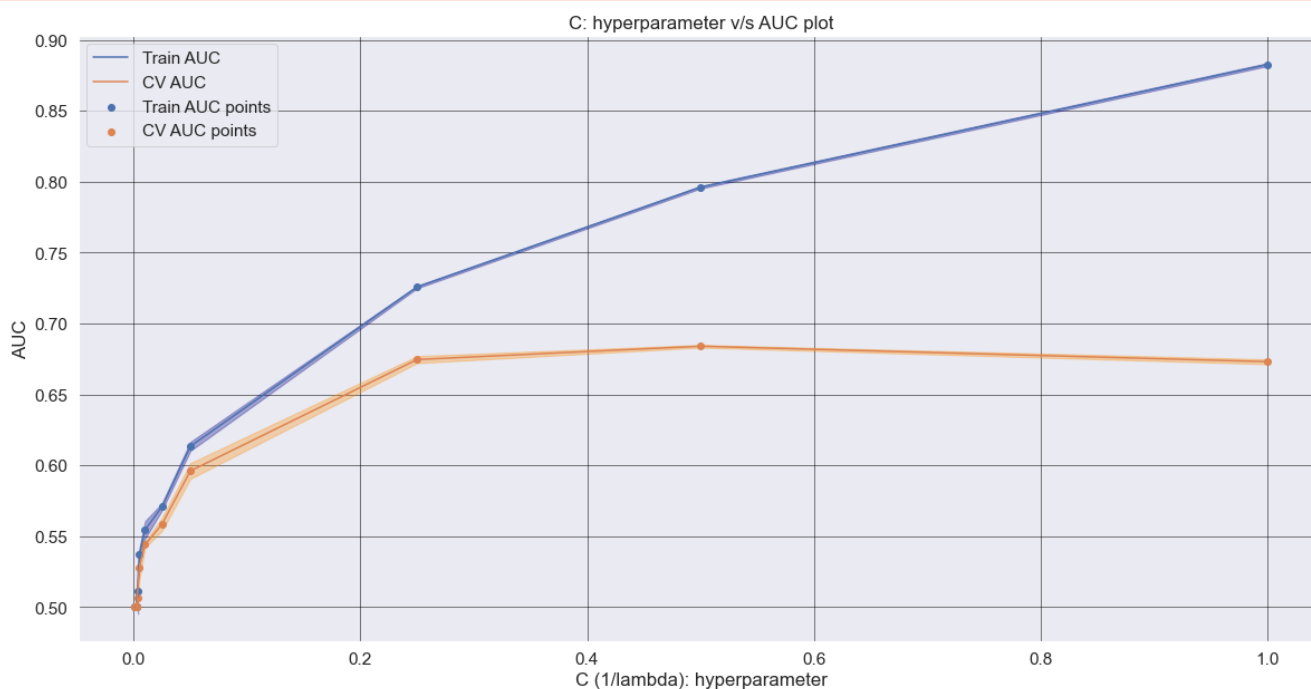
```

```

[CV] C=0.003 ..... C=0.003, total= 0.2s
[CV] ..... C=0.003, total= 0.2s
[CV] C=0.003 ..... C=0.003, total= 0.2s
[CV] ..... C=0.003, total= 0.2s
[CV] C=0.001 ..... C=0.001, total= 0.2s
[CV] ..... C=0.001, total= 0.2s
[CV] C=0.001 ..... C=0.001, total= 0.2s
[CV] ..... C=0.001, total= 0.2s
[CV] C=0.001 ..... C=0.001, total= 0.2s
[CV] ..... C=0.001, total= 0.2s

```

```
[Parallel(n_jobs=1)]: Done 33 out of 33 | elapsed: 33.8s finished
```



In [167]:

```
best_c2=clf.best_params_
print(best_c2)
```

```
{'C': 0.5}
```

In [168]:

```

model = LogisticRegression(C = 0.25,penalty='l1',class_weight='balanced')
model.fit(X_tr, y_train)

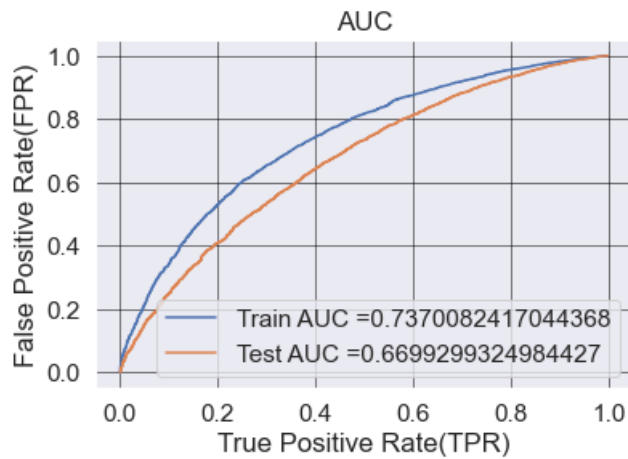
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

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

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='--', linewidth=0.5)
plt.show()

```



In [169]:

```
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999997861339704 for threshold 0.431  
[[ 1709 1710]  
 [ 3498 15528]]

In [170]:

```
conf_matr_df_train_2 = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
train_fpr, train_fpr)), range(2), range(2))
```

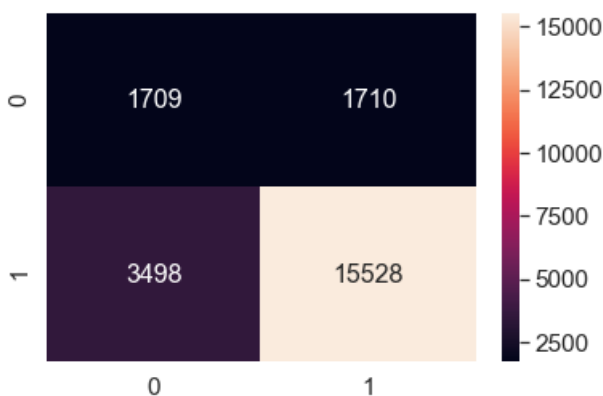
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999997861339704 for threshold 0.431

In [171]:

```
sns.set(font_scale=1.4) #for label size
sns.heatmap(conf_matr_df_train_2, annot=True, annot_kws={"size": 16}, fmt='g')
```

Out[171]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e7a6780548>



In [172]:

```
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.479  
[[ 196 2350]  
 [ 258 12322]]

```
[ 258 13696]]
```

```
In [173]:
```

```
conf_matr_df_test_2 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2), range(2))
```

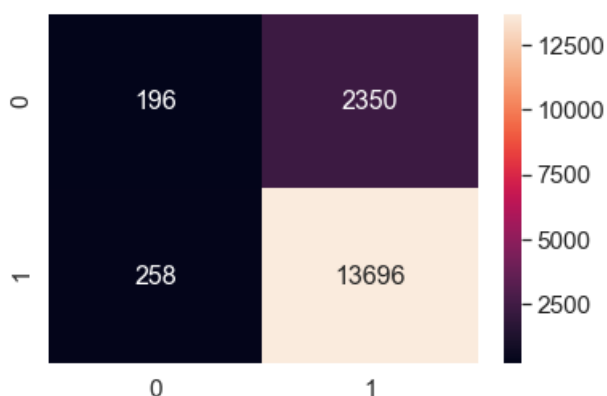
the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.479

```
In [174]:
```

```
sns.set(font_scale=1.4) # for label size
sns.heatmap(conf_matr_df_test_2, annot=True, annot_kws={"size": 16}, fmt='g')
```

```
Out[174]:
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e78e0866c8>



```
In [175]:
```

```
# set3
```

```
In [176]:
```

```
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

```
In [177]:
```

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train.append(vector)

print(len(avg_w2v_vectors_train))
print(len(avg_w2v_vectors_train[0]))
print(avg_w2v_vectors_train[0])
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 22445/22445
[00:11<00:00, 2022.00it/s]
```

```
22445
```

```
300
```

```
[-6.36130655e-02 -2.98671088e-02  1.44113993e-02 -1.65850472e-01
```

1.0.30130000e-02 -2.30071900e-02 1.44113333e-02 -1.00000000e-01  
2.43256556e-02 -2.18765736e-02 -3.67033792e+00 2.65864729e-01  
6.49939701e-02 -1.80636305e-01 1.59821424e-01 -1.93610935e-02  
1.76833326e-02 -1.59901647e-01 -7.70059131e-02 -1.29170147e-01  
-9.60884153e-02 -1.11404056e-01 5.66026465e-02 5.76478285e-02  
1.99858299e-02 7.02051958e-03 -5.50928901e-02 6.31444271e-02  
-4.13374374e-02 -6.04097757e-02 8.26678403e-02 -1.33313949e-01  
-1.47735663e-01 -1.09199348e-01 -2.87392524e-01 -4.63169326e-02  
6.41040479e-02 -8.82091313e-02 -1.65970039e-01 -2.77768625e-02  
-2.96260861e-02 -1.05177649e-01 6.59151689e-02 -9.73202803e-02  
-2.34277687e-02 8.91123083e-02 -9.93708472e-03 -2.16032828e-01  
-3.98970201e-02 -5.97881382e-02 4.48946958e-02 -1.88224151e-01  
-8.21719294e-02 -2.63194222e-02 2.43500537e-02 4.10525528e-02  
4.37783819e-03 7.22830208e-03 1.05519299e-01 -7.48155069e-02  
8.21422493e-02 2.02504014e-02 -7.49091799e-02 9.99924986e-02  
-2.53738340e-02 1.52611861e-02 7.16026097e-02 -1.30495044e-01  
-1.62361403e-02 2.06441751e-01 4.71444444e-02 2.10074868e-02  
1.78153912e-01 -1.27918056e-01 -1.57291954e-01 5.27378729e-02  
-4.67759306e-03 -9.50295354e-02 4.10353354e-02 -2.28342450e-01  
1.40315897e-01 5.49380840e-02 8.75941715e-02 -1.11539628e-01  
3.91888204e-02 -4.60752536e-01 -1.01493795e-01 -1.42145826e-01  
2.76480444e-02 6.76731972e-02 3.20082668e-02 -1.17512169e-01  
1.24697513e-01 -3.20006500e-02 7.84102576e-02 -4.06981493e-02  
-2.20253451e-02 9.41820000e-02 5.80388889e-02 -3.74516243e-01  
-2.52427521e+00 -5.42834799e-02 1.44900322e-01 4.11975236e-02  
-5.14347493e-02 8.37724604e-02 2.70724176e-01 5.76431257e-02  
-5.01261708e-02 -3.34370785e-02 8.14437472e-02 -1.07552562e-01  
-7.56728855e-02 -3.89855757e-02 1.55160725e-02 5.89680660e-02  
2.56501465e-02 2.22745988e-01 -7.30515569e-02 2.84419958e-02  
-4.45029160e-02 -2.39047701e-02 9.80545090e-02 7.91201549e-02  
-1.61729631e-01 1.77801403e-02 2.88151690e-02 -1.69847242e-01  
-4.00431944e-03 -1.19664931e-03 5.57457647e-02 -7.92253958e-02  
-2.13914333e-02 1.58984108e-01 6.27076902e-02 4.91708910e-02  
-6.63686847e-02 -1.55353248e-01 6.56521597e-02 1.84377618e-02  
1.29023467e-01 3.74861479e-02 2.57726804e-01 3.74205588e-01  
1.22936355e-01 4.44947681e-02 3.28010069e-03 6.95104444e-03  
3.07705972e-03 -2.71807854e-02 1.24882653e-01 -3.94405042e-02  
3.29752626e-01 1.58696653e-01 -2.27855528e-02 -3.52606681e-02  
4.76674859e-02 -6.21542691e-02 5.79896167e-02 -8.58536153e-02  
5.60179861e-02 3.08002493e-02 -9.38832333e-02 -3.76570340e-02  
6.21763368e-02 -5.95145833e-03 -5.15096042e-04 -4.24368653e-02  
-8.25822056e-02 2.25758125e-02 -1.83186500e-02 1.21015864e-01  
1.73836453e-01 -8.63686931e-02 -2.02193889e-02 1.62425743e-02  
-9.14413375e-02 -4.10978215e-02 -1.15233192e-01 2.13133384e-01  
-8.84506773e-02 -7.44204701e-02 -3.98922903e-02 -6.68608222e-02  
1.28370458e-01 1.10772359e-01 -1.85296799e-02 -3.68262479e-02  
1.31744097e-02 -2.50690256e-01 7.77075764e-03 1.77283958e-03  
1.90450258e-01 -2.14338403e-02 -8.51184479e-03 -5.33001736e-02  
-4.53764306e-02 -4.48636901e-02 3.25378118e-02 -1.12068207e-01  
9.47111174e-02 1.04112611e-01 8.07725208e-02 -2.79982931e-03  
1.70517945e-01 -4.00693241e-02 -2.07044375e-03 2.93203704e-02  
2.36346374e-02 1.18112264e-01 5.38878132e-02 -7.32100097e-02  
1.62453405e-01 -6.74826549e-02 1.22739669e-01 1.66126556e-02  
-8.59109374e-02 -6.44528208e-02 -2.35824522e-02 2.97553276e-02  
-8.35082944e-02 -1.16821694e-01 -4.78625674e-02 -5.91192361e-03  
-9.92166812e-02 -9.19806208e-02 -3.40431701e-02 3.42811257e-02  
-2.69738263e+00 1.24778683e-01 -3.43029097e-03 -5.90494167e-03  
-1.66246757e-02 -8.15915063e-02 1.21356901e-01 3.97240097e-02  
6.55093750e-04 -7.19937069e-02 -9.03269667e-02 1.14210451e-01  
2.52387333e-02 -7.68073153e-02 -8.26117201e-02 5.76407764e-02  
-1.76668065e-01 3.86877188e-02 -2.06318229e-01 4.70021389e-02  
5.20637153e-03 -2.14625743e-02 -4.52095826e-02 -1.20489826e-02  
-6.94916111e-02 4.47884231e-02 -3.39729826e-02 -6.88624701e-02  
1.59691319e-03 -7.08379722e-04 3.27276167e-02 -7.75372437e-02  
3.15919674e-02 -5.52375062e-02 5.94425597e-02 1.97262938e-02  
6.22628792e-02 -1.89401882e-02 1.33400243e-02 -2.34690653e-02  
4.23251819e-02 -1.54541648e-01 -2.76768412e-01 -8.25564694e-02  
1.08877513e-01 -4.95602708e-03 -2.04696847e-02 -1.42821148e-01  
-2.01050774e-01 9.51513653e-02 -1.54439583e-02 8.02743632e-02  
9.47042708e-02 6.11147708e-02 5.26816806e-03 -4.73692597e-02  
3.15961868e-01 1.26306417e-02 2.27548375e-02 1.82231715e-01  
-1.99326625e-02 1.04141790e-01 3.86473264e-02 6.39190340e-02  
1.93894368e-02 -2.94880833e-02 -1.14644382e-02 -9.39447319e-02  
-9.27013889e-02 -3.55651355e-02 -4.84626882e-02 2.08163111e-02  
-3.32993049e-02 -4.16279097e-03 1.15277511e-01 3.03436049e-02]

In [178]:

```
avg_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv.append(vector)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:05<00:00, 2028.86it/s]
```

In [179]:

```
avg_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_test.append(vector)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:08<00:00, 1921.06it/s]
```

In [180]:

```
# avg w2v for project_titles
avg_w2v_vectors_pj_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_pj_title_train.append(vector)

print(len(avg_w2v_vectors_pj_title_train))
print(len(avg_w2v_vectors_pj_title_train[0]))
print(avg_w2v_vectors_pj_title_train[0])
```

```
100%|██████████████████████████████████████████████████████████████████████████| 22445/22445  
[00:00<00:00, 39240.13it/s]
```

22445

300

```
[ 1.11699320e-01 -8.79170000e-02 -1.53535600e-02 -4.26230000e-01
 1.59202000e-01 -2.32649400e-01 -3.61468000e+00  5.16348000e-01
 1.98154000e-02 -2.10774800e-01  9.52240000e-02  7.98246000e-02
 2.61314000e-01  4.16078000e-02 -2.54906200e-01 -5.88602980e-02
 8.91600000e-03 -2.74082000e-01 -1.50301200e-02 -1.66438000e-01
-6.82560000e-02  7.53136000e-02 -1.17752600e-01  2.71002800e-01
-3.34654000e-02 -3.90990000e-01  1.38476000e-02 -3.15304000e-02
 1.20213000e-01 -1.01610000e-01 -2.66013820e-01  1.48666200e-01
 1.01532000e-01 -1.09025600e-01 -4.17345000e-02  3.24502000e-02
-3.82960000e-03  1.54252200e-01 -5.70732000e-02 -1.72186200e-01
 6.10400000e-03 -1.87940000e-02 -8.78790000e-02 -2.11696000e-01
-2.88234400e-01  2.36416000e-01 -2.49006000e-02 -2.84210000e-01
-2.70326400e-01  2.15745540e-01 -3.32757000e-01  2.71181400e-01
-1.90495120e-01  1.17212000e-01 -2.39188000e-02  3.56602000e-02
 1.04220400e-01 -5.91762000e-02 -2.05560400e-02  5.00864000e-02]
```



```

-9.43118000e-02 -1.52546000e-02 -8.59656000e-02 1.88516000e-01
1.96025000e-01 2.81256000e-01 1.86708600e-01 -2.15187400e-01
1.10492000e-02 -2.47324400e-01 -7.22104000e-02 3.08102600e-01
9.14994000e-02 1.70755000e-01 -1.81822000e-02 -6.71676000e-02
1.07398200e-01 -1.81245400e-01 -9.32060000e-02 1.54689600e-01
1.25239660e-01 -3.67030000e-01 -1.59842200e-01 -2.75040000e-02
-2.72288000e-01 8.32400000e-02 -3.85360000e-02 -1.36253200e-02
1.67356000e-01 -1.21447600e-01 2.40853000e-01 -5.48298000e-02
-2.84778000e-01 1.03124000e-01 2.13980000e-02 -1.65126000e-01
-2.36514000e+00 -2.65319800e-01 2.16966000e-01 3.29280000e-02
3.28620000e-02 3.99804000e-02 -9.85440000e-02 1.26317800e-01
-5.32422400e-02 1.92438200e-01 -3.96740000e-02 -2.17806800e-01
-1.26557800e-02 9.15252000e-02 -1.69097000e-01 -1.32732200e-01
3.01718600e-01 -8.05300000e-02 5.86140000e-02 -2.25324000e-01
1.35880600e-01 1.45590000e-01 -1.55858400e-01 1.46203200e-01
1.40157000e-02 6.09020000e-02 -3.54339000e-01 -2.34175600e-01
-8.32708000e-02 -7.06802000e-02 -2.56553400e-02 -1.46368000e-01
1.11478400e-01 -2.42980200e-01 3.74908000e-02 1.50360160e-01
1.02942600e-01 -9.38000000e-02 6.12340000e-03 2.22412600e-01
-5.91960000e-02 1.74040000e-02 9.98898000e-02 3.85549400e-01
-1.66880000e-02 -7.07206000e-02 -9.13318000e-02 -5.55466000e-02
1.66791600e-01 4.93060000e-02 -2.66220000e-03 -1.98870000e-01
2.42024000e-02 -4.07396000e-02 -4.09162000e-02 -4.85350000e-03
2.95816600e-01 3.53878000e-02 -1.45590000e-01 -1.54521620e-01
-6.72316000e-02 -1.16198888e-01 -5.08806000e-02 2.31204000e-02
-1.30304400e-01 3.61122000e-01 -1.20899580e-01 1.80122000e-02
1.70632000e-02 -4.98376000e-02 2.10686160e-01 2.01960000e-02
1.05433200e-01 -8.86694000e-02 1.67312000e-01 -1.00255400e-01
-2.93254000e-02 1.33303400e-01 1.53968000e-01 3.52944000e-01
1.49784800e-01 -2.40340000e-02 7.71000000e-02 -1.19030000e-02
1.42421840e-01 2.19328200e-01 2.66996400e-01 4.90206000e-02
-2.45366000e-02 1.33064200e-02 6.89080000e-03 8.83294000e-02
2.77948000e-01 3.02736000e-02 2.23984000e-01 -1.40974400e-01
1.05504600e-01 -2.59217600e-01 1.44482000e-02 -5.53758000e-02
-1.29432600e-01 2.55400000e-03 -2.88472000e-02 -1.45468000e-01
1.63878800e-01 -1.76732200e-01 -3.06820000e-02 1.89730000e-01
-1.32721000e-01 -5.14712000e-02 2.51472000e-02 -1.46345800e-01
-1.17862926e-01 -1.02354800e-01 -2.72900000e-01 2.18668000e-01
-1.27632800e-01 -8.74540000e-02 1.29516000e-01 -1.04500000e-01
-3.28492000e-02 3.12396000e-02 -2.10450000e-01 1.23800600e-02
-2.55718000e-01 -1.19636000e-01 1.37467200e-01 1.68883600e-01
-2.60378000e+00 3.31402000e-01 -4.60588000e-02 -1.50252000e-02
-3.65400000e-04 4.87828000e-02 4.52966000e-02 -1.43003200e-01
-1.36571600e-01 -1.44876400e-01 -8.44460000e-02 4.43616000e-02
1.39657880e-01 -2.01616000e-01 2.02671200e-01 -3.61978200e-01
-2.85010400e-02 -2.42614600e-01 -3.93906000e-01 -1.05038000e-01
1.69777600e-01 -5.68394000e-02 -1.15864480e-01 -2.07623200e-01
-2.04866600e-01 -6.93034000e-02 -1.32547200e-01 4.24738000e-02
-1.38944000e-01 -1.87889800e-01 3.00680000e-02 7.11320000e-02
2.69588400e-01 -5.74250000e-02 -3.54272800e-01 7.18992000e-02
1.00266200e-01 3.54644400e-02 1.89606620e-02 -7.71960000e-02
-5.48926000e-02 -2.60788000e-02 -3.44266000e-02 -2.67220000e-02
-1.48643400e-01 -5.62598000e-02 7.23542000e-02 -8.43545400e-02
-1.13358000e-02 -2.34902000e-01 6.87362000e-02 -1.04288400e-01
1.06281600e-01 1.94329800e-01 -4.28862000e-02 -2.32482000e-02
4.03778000e-01 -2.19486800e-01 2.46795400e-02 -6.70678000e-02
-2.76183600e-01 1.48226800e-01 -5.43020000e-03 2.87253000e-01
-1.62224000e-01 -3.11540000e-01 9.66960000e-03 1.21522340e-01
-1.45364600e-01 -2.47388000e-02 4.10900000e-02 -1.14907200e-01
-2.00677800e-01 -5.28853600e-02 1.54756000e-01 4.98604000e-02]

```

In [181]:

```

avg_w2v_vectors_pj_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_pj_title_cv.append(vector)

```

```
100%|██████████████████████████████████████████████████████████████████████████████| 11055/11055  
[00:00<00:00, 48736.17it/s]
```

In [182]:

```
avg_w2v_vectors_pj_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_pj_title_test.append(vector)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 53485.42it/s]
```

In [183]:

```
X_tr = hstack((avg_w2v_vectors_train, X_train_state_oh, X_train_teacher_oh,
               X_train_grade_oh, X_train_category_oh,
               X_train_subcategory_oh, X_train_price_norm,
               X_train_teach_prev_norm, avg_w2v_vectors_pj_title_train)).tocsr()

X_cr = hstack((avg_w2v_vectors_cv, X_cv_state_oh, X_cv_teacher_oh,
               X_cv_grade_oh, X_cv_category_oh,
               X_cv_subcategory_oh, X_cv_price_norm,
               X_cv_teach_prev_norm, avg_w2v_vectors_pj_title_cv)).tocsr()

X_te = hstack((avg_w2v_vectors_test, X_test_state_oh, X_test_teacher_oh,
               X_test_grade_oh, X_test_category_oh,
               X_test_subcategory_oh, X_test_price_norm,
               X_test_teach_prev_norm, avg_w2v_vectors_pj_title_test)).tocsr()
```

In [184]:

```
lr = LogisticRegression(penalty='l1',class_weight="balanced")

parameters = {'C':[0.25,0.05,0.025, 0.01, 0.005,0.0025, 0.004, 0.003,0.001]}

clf = GridSearchCV(lr, parameters, cv= 3, scoring='roc_auc',return_train_score=True,verbose=2)

clf.fit(X_tr, y_train)

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

plt.figure(figsize=(20,10))

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("C(1/lambda) : hyperparameter")
plt.ylabel("AUC")
plt.title("C: hyperparameter v/s AUC plot")
plt.grid(color='black', linestyle='-', linewidth=0.5)
```

```
plt.show()
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

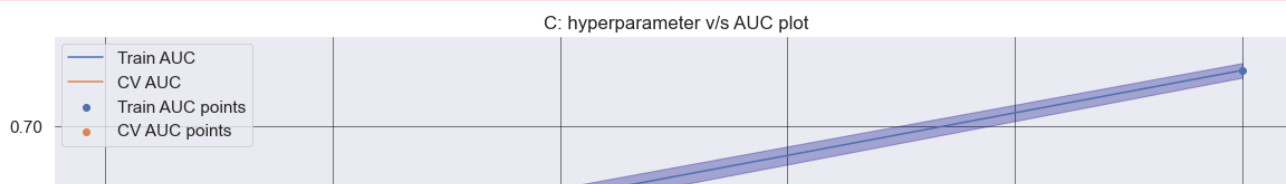
```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
```

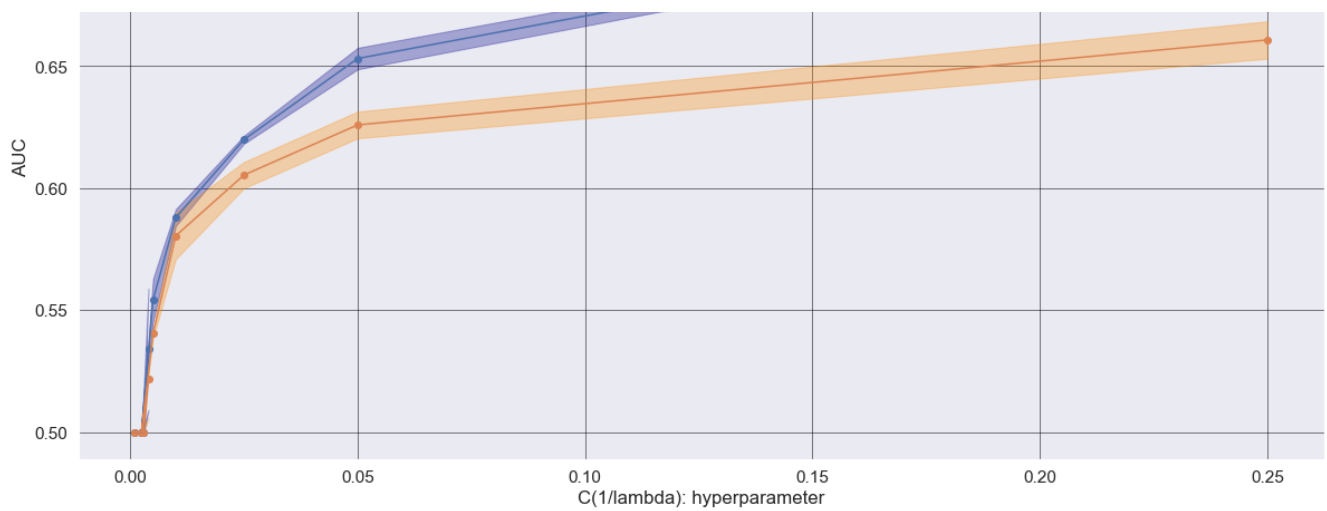
```
[CV] C=0.25 .....  
[CV] ..... C=0.25, total= 7.2min
```

```
[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 7.2min remaining: 0.0s
```

```
[CV] C=0.25 .....  
[CV] ..... C=0.25, total= 7.6min  
[CV] C=0.25 .....  
[CV] ..... C=0.25, total= 6.7min  
[CV] C=0.05 .....  
[CV] ..... C=0.05, total= 26.2s  
[CV] C=0.05 .....  
[CV] ..... C=0.05, total= 21.1s  
[CV] C=0.05 .....  
[CV] ..... C=0.05, total= 23.5s  
[CV] C=0.025 .....  
[CV] ..... C=0.025, total= 1.3s  
[CV] C=0.025 .....  
[CV] ..... C=0.025, total= 4.6s  
[CV] C=0.025 .....  
[CV] ..... C=0.025, total= 5.7s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, total= 0.7s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, total= 0.8s  
[CV] C=0.01 .....  
[CV] ..... C=0.01, total= 0.8s  
[CV] C=0.005 .....  
[CV] ..... C=0.005, total= 0.6s  
[CV] C=0.005 .....  
[CV] ..... C=0.005, total= 0.5s  
[CV] C=0.005 .....  
[CV] ..... C=0.005, total= 0.7s  
[CV] C=0.0025 .....  
[CV] ..... C=0.0025, total= 0.4s  
[CV] C=0.0025 .....  
[CV] ..... C=0.0025, total= 0.4s  
[CV] C=0.0025 .....  
[CV] ..... C=0.0025, total= 0.4s  
[CV] C=0.004 .....  
[CV] ..... C=0.004, total= 0.4s  
[CV] C=0.004 .....  
[CV] ..... C=0.004, total= 0.4s  
[CV] C=0.004 .....  
[CV] ..... C=0.004, total= 0.6s  
[CV] C=0.003 .....  
[CV] ..... C=0.003, total= 0.4s  
[CV] C=0.003 .....  
[CV] ..... C=0.003, total= 0.4s  
[CV] C=0.003 .....  
[CV] ..... C=0.003, total= 0.4s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, total= 0.4s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, total= 0.4s  
[CV] C=0.001 .....  
[CV] ..... C=0.001, total= 0.4s
```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 23.0min finished
```





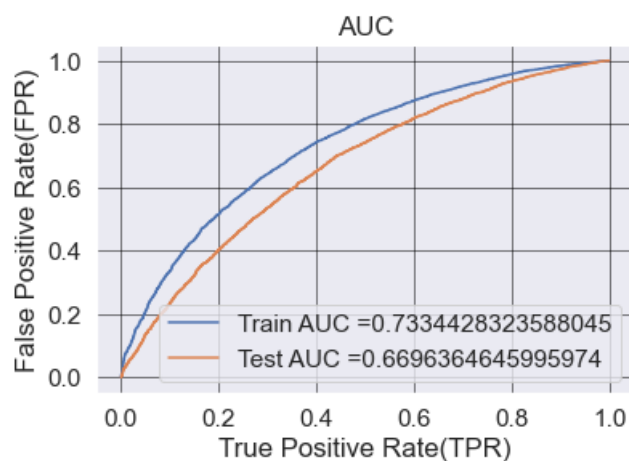
In [185]:

```
best_c3=clf.best_params_  
print(best_c3)
```

```
{'C': 0.25}
```

In [186]:

```
model = LogisticRegression(C=0.5,penalty='l1',class_weight='balanced')  
  
model.fit(X_tr, y_train)  
  
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive  
# class  
# not the predicted outputs  
  
y_train_pred = batch_predict(model, X_tr)  
y_test_pred = batch_predict(model, X_te)  
  
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)  
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)  
  
plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))  
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))  
plt.legend()  
plt.xlabel("True Positive Rate(TPR)")  
plt.ylabel("False Positive Rate(FPR)")  
plt.title("AUC")  
plt.grid(color='black', linestyle='--', linewidth=0.5)  
plt.show()
```



In [187]:

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

the maximum value of  $tpr*(1-fpr)$  0.24999997861339704 for threshold 0.414

Out[187]:

```
array([[ 1709,  1710],
       [ 3501, 15525]], dtype=int64)
```

In [188]:

```
conf_matr_df_train_3 = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
train_fpr, train_fpr)), range(2), range(2))
```

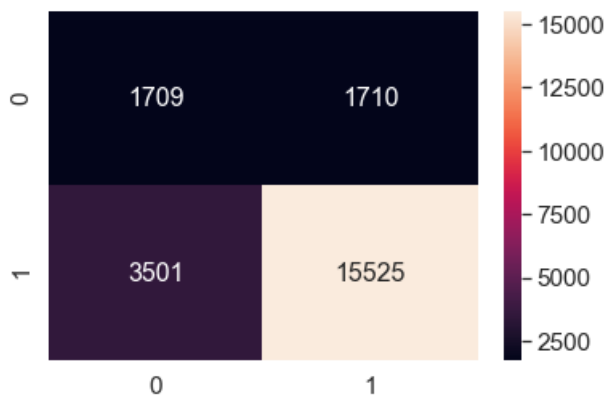
the maximum value of  $tpr*(1-fpr)$  0.24999997861339704 for threshold 0.414

In [189]:

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train_3, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[189]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e74a22b908>



In [190]:

```
confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
```

the maximum value of  $tpr*(1-fpr)$  0.25 for threshold 0.477

Out[190]:

```
array([[1428, 1118],
       [4239, 9715]], dtype=int64)
```

In [191]:

```
conf_matr_df_test_4 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, tes
t_fpr, test_fpr)), range(2), range(2))
```

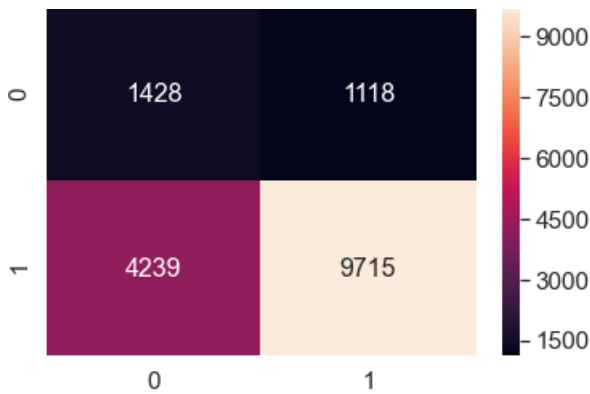
the maximum value of  $tpr*(1-fpr)$  0.25 for threshold 0.477

In [192]:

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test_4, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[192]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e7349d8188>



In [193]:

```
# Please write all the code with proper documentation
# preprocessing project_title and essay with TFIDF W2V Vectorization
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['essay'])
# 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 [194]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay']): # 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):
            # 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
    tfidf_w2v_vectors_train.append(vector)

print(len(tfidf_w2v_vectors_train))
print(len(tfidf_w2v_vectors_train[0]))
```

```
100%|██████████████████████████████████████████████████████████| 22445/22445 [02:  
22<00:00, 157.87it/s]
```

22445  
300

In [195]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essay']): # 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):
            # 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
```

```
100%|███████████████████████████████████████████████████████████| 11055/11055 [01:  
08<00:00, 160.8lit/s]
```

In [196]:

```
100%|███████████████████████████████████████████████████████████████████████████████| 16500/16500 [01:  
45<00:00, 156.04it/s]
```

In [197]:

In [198]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_pj_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title']): # 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):
            # 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
    tfidf_w2v_vectors_pj_title_train.append(vector)

print(len(tfidf_w2v_vectors_pj_title_train))
print(len(tfidf_w2v_vectors_pj_title_train[0]))
```

```
printf("len(%c%c%c_wz_v_vectors_p) = %d\n", c1, c2, c3, len(c1));
```

22445  
300

In [199]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_pj_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title']): # 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):
            # 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
    tfidf_w2v_vectors_pj_title_cv.append(vector)

print(len(tfidf_w2v_vectors_pj_title_cv))
print(len(tfidf_w2v_vectors_pj_title_cv[0]))
```

11055  
300

In [200]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_pj_title_test = []; # the avg-w2v for each sentence/review is stored in this list

for sentence in tqdm(X_test['project_title']): # 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):
            # 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
    tfidf_w2v_vectors_pj_title_test.append(vector)

print(len(tfidf_w2v_vectors_pj_title_test))
print(len(tfidf_w2v_vectors_pj_title_test[0]))
```

16500  
300

In [201]:

```
# Concatinating all the features
X_tr = hstack((tfidf_w2v_vectors_train, X_train_state_one, X_train_teacher_one,
```



```

X_cr = hstack((tfidf_w2v_vectors_cv, X_cv_state_oh, X_cv_teacher_oh,
              X_cv_grade_oh, X_cv_category_oh,
              X_cv_subcategory_oh, X_cv_price_norm,
              X_cv_teach_prev_norm, tfidf_w2v_vectors_pj_title_cv)).tocsr()

X_tr = hstack((tfidf_w2v_vectors_tr, X_tr_state_oh, X_tr_teacher_oh,
              X_tr_grade_oh, X_tr_category_oh,
              X_tr_subcategory_oh, X_tr_price_norm,
              X_tr_teach_prev_norm, tfidf_w2v_vectors_pj_title_tr)).tocsr()

X_te = hstack((tfidf_w2v_vectors_test, X_test_state_oh, X_test_teacher_oh,
              X_test_grade_oh, X_test_category_oh,
              X_test_subcategory_oh, X_test_price_norm,
              X_test_teach_prev_norm, tfidf_w2v_vectors_pj_title_test)).tocsr()

```

In [202]:

```

lr = LogisticRegression(penalty='l1', class_weight="balanced")

parameters = {'C':[0.25,0.05,0.025, 0.01, 0.005,0.0025, 0.004, 0.003,0.001]}

clf = GridSearchCV(lr, parameters, cv= 3, scoring='roc_auc',return_train_score=True,verbose=2)

clf.fit(X_tr, y_train)

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

plt.figure(figsize=(20,10))

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.3,color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'],cv_auc - cv_auc_std,cv_auc + cv_auc_std,alpha=0.3,color='darkorange')

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

plt.legend()
plt.xlabel("C(1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("C: hyperparameter v/s AUC plot")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()

```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```

[CV] C=0.25 .....
[CV] ..... C=0.25, total= 2.3s

```

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 2.2s remaining: 0.0s

```

[CV] C=0.25 .....
[CV] ..... C=0.25, total= 1.8s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 2.0s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 1.5s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.9s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 1.1s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.8s

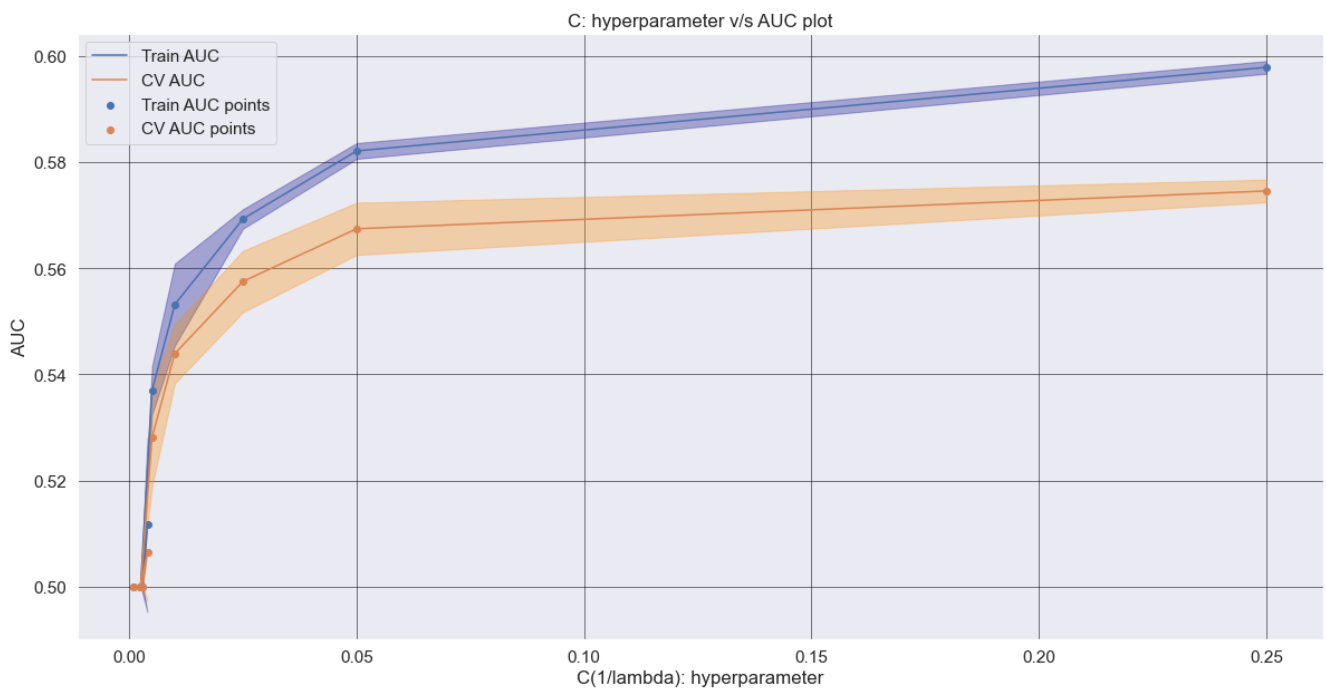
```

```

[CV] ..... C=0.025, total= 2.2s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 1.0s
[CV] C=0.025 .....
[CV] ..... C=0.025, total= 0.9s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.8s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.8s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.9s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.4s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.5s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.5s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.4s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.5s
[CV] C=0.0025 .....
[CV] ..... C=0.0025, total= 0.3s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.4s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.4s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.4s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.4s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.4s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.4s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.4s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.4s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.4s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.3s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.4s

```

```
[Parallel(n_jobs=1)]: Done 27 out of 27 | elapsed: 24.0s finished
```



In [203]:

```

best_c4=clf.best_params_
print(best_c4)

```

```
{'C': 0.25}
```

In [204]:

```
model = LogisticRegression(C = 0.25,penalty="l1",class_weight="balanced")

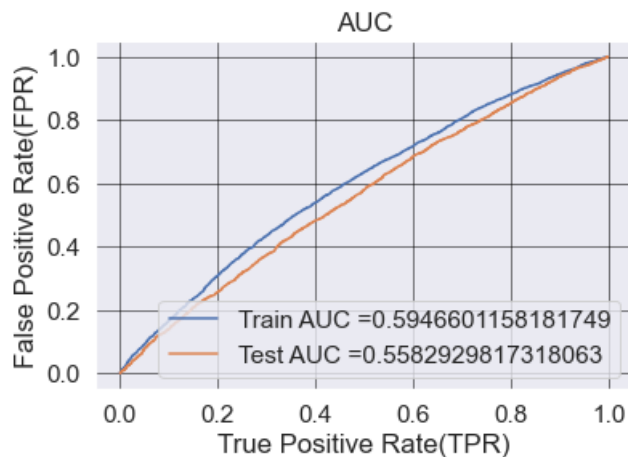
model.fit(X_tr, y_train)

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

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

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```



In [205]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.487
[[ 1709  1710]
 [ 6935 12091]]
```

In [206]:

```
conf_matr_df_train_4 = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
train_fpr, train_fpr)), range(2),range(2))
```

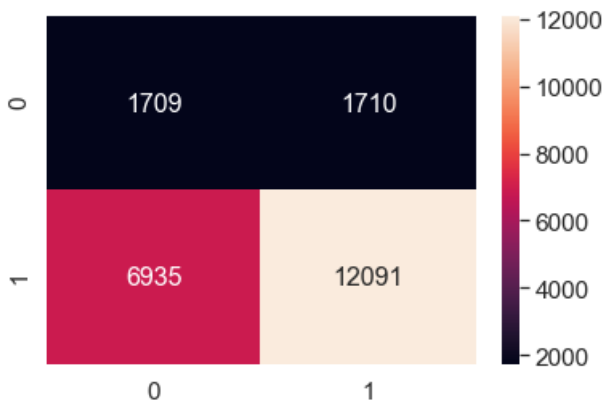
```
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.487
```

In [207]:

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train_4, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[207]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e747e05f88>



In [208]:

```
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.518

```
[[1549  997]
 [7342 6612]]
```

In [209]:

```
conf_matr_df_test_4 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2), range(2))
```

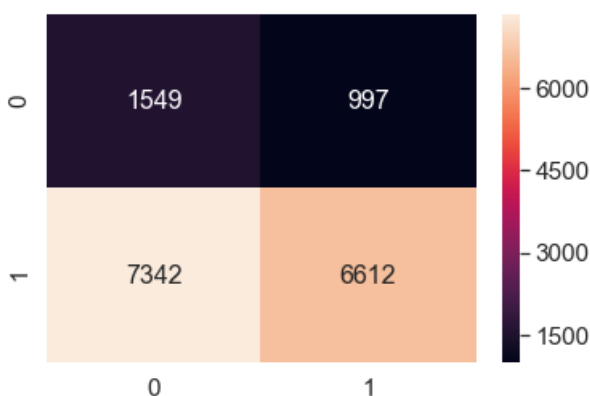
the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.518

In [210]:

```
sns.set(font_scale=1.4) # for label size
sns.heatmap(conf_matr_df_test_4, annot=True, annot_kws={"size": 16}, fmt='g')
```

Out[210]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e74659dc08>



## 2.5 Logistic Regression with added Features `Set 5`

In [211]:

```
# 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
```

```
# reading and understanding error messages will be very much helpful 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 [212]:

```
X_train_essay_sentiment = X_train['essay_sentiment'][:,np.newaxis]
X_test_essay_sentiment = X_test['essay_sentiment'][:,np.newaxis]
X_cv_essay_sentiment = X_cv['essay_sentiment'][:,np.newaxis]
print(X_train_essay_sentiment.shape)
```

(22445, 1)

In [213]:

```
normalizer = Normalizer()

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

X_train_essay_sentiment = normalizer.transform(X_train['essay_sentiment'].values.reshape(1,-1))
X_test_essay_sentiment = normalizer.transform(X_test['essay_sentiment'].values.reshape(1,-1))
X_cv_essay_sentiment = normalizer.transform(X_cv['essay_sentiment'].values.reshape(1,-1))

X_train_essay_sentiment = X_train_essay_sentiment.reshape(-1,1)
X_test_essay_sentiment = X_test_essay_sentiment.reshape(-1,1)
X_cv_essay_sentiment = X_cv_essay_sentiment.reshape(-1,1)
```

In [214]:

```
normalizer = Normalizer()

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

X_train_essay_sentiment = normalizer.transform(X_train['essay_sentiment'].values.reshape(1,-1))
X_test_essay_sentiment = normalizer.transform(X_test['essay_sentiment'].values.reshape(1,-1))
X_cv_essay_sentiment = normalizer.transform(X_cv['essay_sentiment'].values.reshape(1,-1))

X_train_essay_sentiment = X_train_essay_sentiment.reshape(-1,1)
X_test_essay_sentiment = X_test_essay_sentiment.reshape(-1,1)
X_cv_essay_sentiment = X_cv_essay_sentiment.reshape(-1,1)
```

In [215]:

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

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

In [216]:

```
vectorizer = TfidfVectorizer(min_df=5)
vectorizer.fit(X_train['project_resource_summary'].values) # fit has to happen only on train
datadata

# we use the fitted CountVectorizer to convert the text to vector
X_train_summary_tfidf = vectorizer.transform(X_train['project_resource_summary'].values)
X_cv_summary_tfidf = vectorizer.transform(X_cv['project_resource_summary'].values)
X_test_summary_tfidf = vectorizer.transform(X_test['project_resource_summary'].values)
```

In [217]:

```
X_tr = hstack((X_train_essay_tfidf, X_train_state_oh, X_train_teacher_oh,
               X_train_grade_oh, X_train_price_norm, X_train_category_oh,
               X_train_subcategory_oh, X_train_teach_prev_norm,
               X_train_pj_title_tfidf, X_train_summary_tfidf, X_train_essay_sentiment)).tocsr()

X_cr = hstack((X_cv_essay_tfidf, X_cv_state_oh, X_cv_teacher_oh,
               X_cv_grade_oh, X_cv_category_oh, X_cv_subcategory_oh,
               X_cv_price_norm, X_cv_teach_prev_norm, X_cv_pj_title_tfidf, X_cv_summary_tfidf, X_cv_essay_sentiment)).tocsr()

X_te = hstack((X_test_essay_tfidf, X_test_state_oh, X_test_teacher_oh,
               X_test_grade_oh, X_test_category_oh, X_test_subcategory_oh,
               X_test_price_norm, X_test_teach_prev_norm,
               X_test_pj_title_tfidf, X_test_summary_tfidf, X_test_essay_sentiment)).tocsr()
```

In [218]:

```
lr = LogisticRegression(penalty='l1', class_weight="balanced")

parameters = {'C': [0.25, 0.1, 0.05, 0.01, 0.005, 0.004, 0.003, 0.002, 0.001]}

clf = GridSearchCV(lr, parameters, cv=3, scoring='roc_auc', return_train_score=True, verbose=2)

clf.fit(X_tr, y_train)

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

plt.figure(figsize=(20, 10))

plt.plot(parameters['C'], train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], train_auc - train_auc_std, train_auc +
train_auc_std, alpha=0.3, color='darkblue')

plt.plot(parameters['C'], cv_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(parameters['C'], cv_auc - cv_auc_std, cv_auc + cv_auc_std, alpha=0.3, color='darkorange')

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

plt.legend()
plt.xlabel("C (1/lambda): hyperparameter")
plt.ylabel("AUC")
plt.title("C: hyperparameter v/s AUC plot")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```

Fitting 3 folds for each of 9 candidates, totalling 27 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

```
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 0.6s
```

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 0.5s remaining: 0.0s

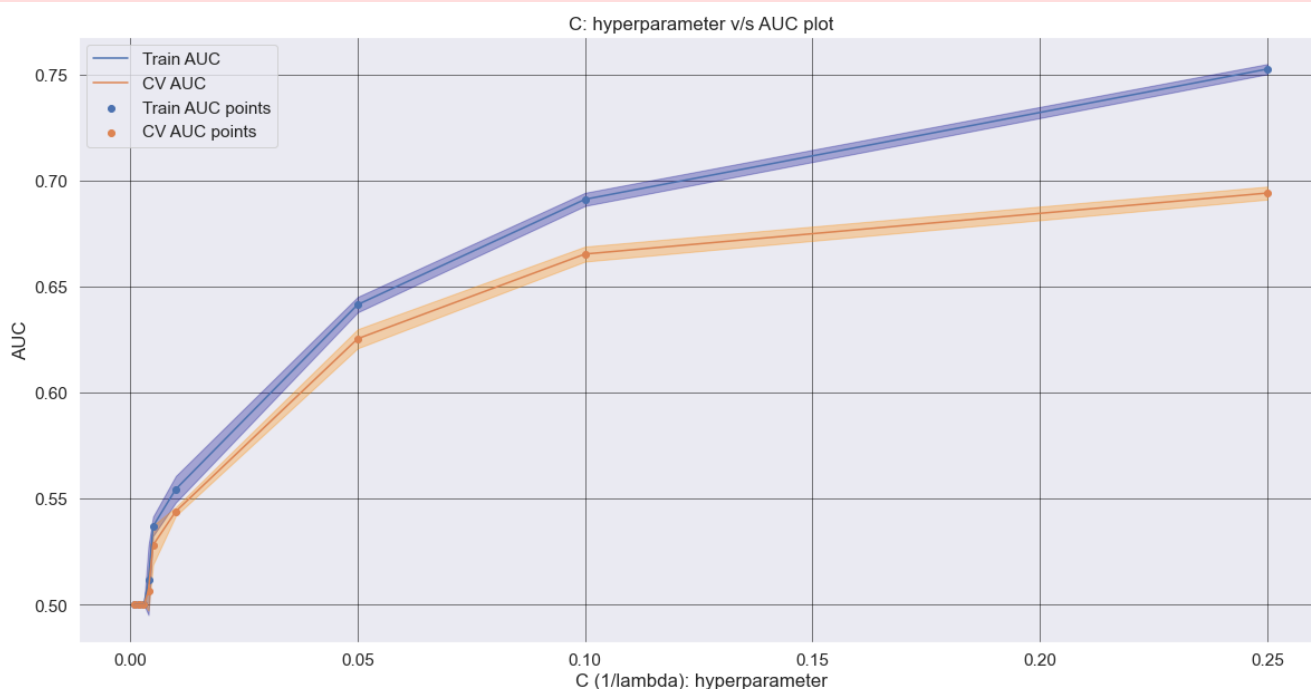
```
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 0.7s
[CV] C=0.25 .....
[CV] ..... C=0.25, total= 1.0s
[CV] C=0.1 .....
[CV] ..... C=0.1, total= 0.4s
[CV] C=0.1 .....
[CV] ..... C=0.1, total= 0.5s
[CV] C=0.1 .....
[CV] ..... C=0.1, total= 0.4s
[CV] C=0.05 .....
```

```

[CV] ..... C=0.05, total= 0.4s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.4s
[CV] C=0.05 .....
[CV] ..... C=0.05, total= 0.4s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.3s
[CV] C=0.01 .....
[CV] ..... C=0.01, total= 0.4s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.2s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.3s
[CV] C=0.005 .....
[CV] ..... C=0.005, total= 0.2s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.1s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.1s
[CV] C=0.004 .....
[CV] ..... C=0.004, total= 0.2s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.1s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.1s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.1s
[CV] C=0.003 .....
[CV] ..... C=0.003, total= 0.1s
[CV] C=0.002 .....
[CV] ..... C=0.002, total= 0.1s
[CV] C=0.002 .....
[CV] ..... C=0.002, total= 0.1s
[CV] C=0.002 .....
[CV] ..... C=0.002, total= 0.1s
[CV] C=0.002 .....
[CV] ..... C=0.002, total= 0.1s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.1s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.1s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.1s
[CV] C=0.001 .....
[CV] ..... C=0.001, total= 0.1s

```

[Parallel(n\_jobs=1)]: Done 27 out of 27 | elapsed: 8.6s finished



In [219]:

```

best_c5=clf.best_params_
print(best_c5)

```

```
{'C': 0.25}
```

In [220]:

```
model = LogisticRegression(C = 0.1,penalty='l1',class_weight="balanced")

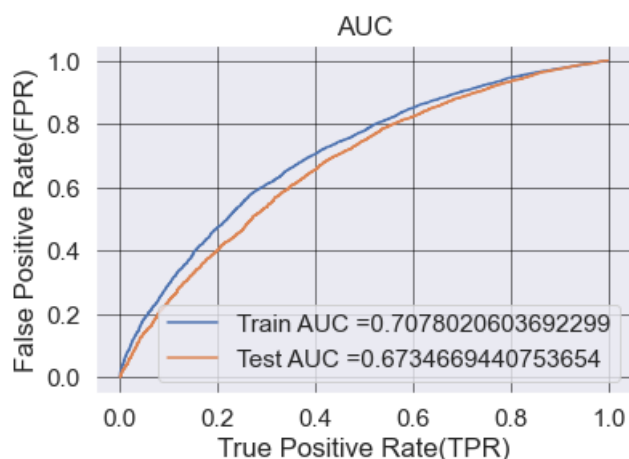
model.fit(X_tr, y_train)

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

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

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

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("True Positive Rate(TPR)")
plt.ylabel("False Positive Rate(FPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```



In [221]:

```
print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
```

```
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.469
[[ 1709  1710]
 [ 4212 14814]]
```

In [222]:

```
conf_matr_df_train_5 = pd.DataFrame(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
train_fpr, train_fpr)), range(2),range(2))
```

```
the maximum value of tpr*(1-fpr) 0.24999997861339704 for threshold 0.469
```

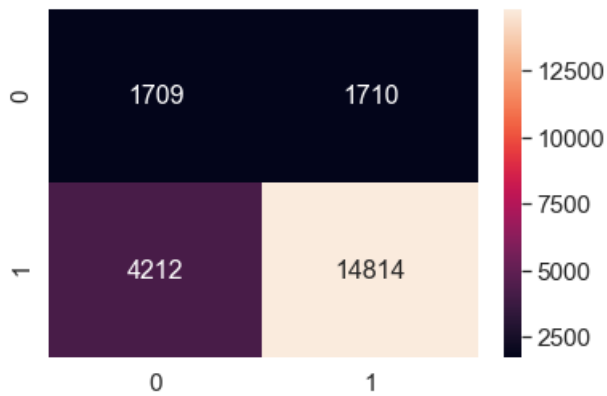
In [223]:

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_train_5, annot=True,annot_kws={"size": 16}, fmt='g')
```

Out[223]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1e731d60788>
```





In [224]:

```
print(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)))
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.505  
[[ 575 1971]  
[ 1074 12880]]

In [225]:

```
conf_matr_df_test_5 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2), range(2))
```

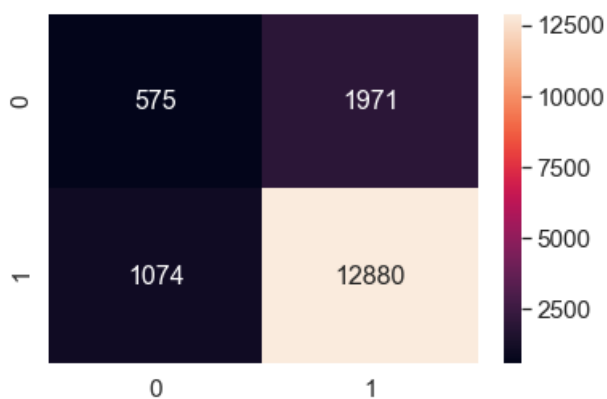
the maximum value of  $tpr \cdot (1 - fpr)$  0.25 for threshold 0.505

In [226]:

```
sns.set(font_scale=1.4) #for label size
sns.heatmap(conf_matr_df_test_5, annot=True, annot_kws={"size": 16}, fmt='g')
```

Out[226]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x1e73682bc88>



### 3. Conclusion

In [227]:

```
# Please compare all your models using Prettytable library
# http://zetcode.com/python/prettytable/

from prettytable import PrettyTable

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

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

```
x.add_row(["BOW", "Logistic Regression", 0.05, 0.69])
x.add_row(["TFIDF", "Logistic Regression", 0.25, 0.69])
x.add_row(["AVG W2V", "Logistic Regression", 0.5, 0.71])
x.add_row(["TFIDF W2V", "Logistic Regression", 0.25, 0.71])
x.add_row(["WITHOUT TEXT", "Logistic Regression", 0.1, 0.64])

print(x)
```

| Vectorizer   | Model               | Alpha:Hyper Parameter | AUC  |
|--------------|---------------------|-----------------------|------|
| BOW          | Logistic Regression | 0.05                  | 0.69 |
| TFIDF        | Logistic Regression | 0.25                  | 0.69 |
| AVG W2V      | Logistic Regression | 0.5                   | 0.71 |
| TFIDF W2V    | Logistic Regression | 0.25                  | 0.71 |
| WITHOUT TEXT | Logistic Regression | 0.1                   | 0.64 |

In [ ]: