

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth Examples: Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*

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

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

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

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

In [5]:

```
%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 [6]:

```

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

```

In [7]:

```

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

```

Number of data points in train data (50000, 17)

```

-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']

```

In [8]:

```

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[8]:

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

In [9]:

```
# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]

#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
project_data.drop('project_submitted_datetime', axis=1, inplace=True)
project_data.sort_values(by=['Date'], inplace=True)

# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
project_data = project_data[cols]

project_data.head(2)
```

Out[9]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_s
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5 L

In [10]:

```
project_grade_category = []

for i in range(len(project_data)):
    a = project_data["project_grade_category"][i].replace(" ", "_").replace("-", "_")
    project_grade_category.append(a)
```

In [11]:

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

In [12]:

```
project_data["project_grade_category"] = project_grade_category
```

In [13]:

```
project_data.head(5)
```

Out[13]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_subject_categories	proje
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Applied Learning
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Literacy & Language
29891	146723	p099708	c0a28c79fe8ad5810da49de47b3fb491	Mrs.	CA	2016-04-27 01:10:09	Math & Science, History & Civics M
23374	72317	n087808	598621c141cda5fb184ee7e8ccdd3fcc	Ms	CA	2016-04-27	Literacy & Language

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_subject_categories	proj
49228	57854	p099430	4000cfe0c8b2df75a218347c1765e283	Ms.	IL	2016-04-27 07:19:44	Literacy & Language

1.2 preprocessing of project_subject_categories

In [14]:

```
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 placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_') # we are replacing the & value into
    cat_list.append(temp.strip())

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

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 [15]:

```
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 placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())
```

```

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]))

```

Clean Titles

In [16]:

```

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

```

In [17]:

```

# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"\'re", " are", phrase)
    phrase = re.sub(r"\s", " is", phrase)
    phrase = re.sub(r"\d", " would", phrase)
    phrase = re.sub(r"\ll", " will", phrase)
    phrase = re.sub(r"\t", " not", phrase)
    phrase = re.sub(r"\ve", " have", phrase)
    phrase = re.sub(r"\m", " am", phrase)
    return phrase

```

In [18]:

```

clean_titles = []

```

```
100%|██████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:01<00:00, 36777.18it/s]
```

```
project_data["clean_titles"] = clean_titles
```

```
project_data.drop(['project title'], axis=1, inplace=True)
```

```
title word count = []
```

```
project_data["title word count"] = title_word_count
```

```
project data.head(5)
```

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	Date	project_essay_1	project_essay_2
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	I recently read an article about giving studen...	I teach at a low-income (Title 1) school. Ever...
41558	33679	p137682	06f6e62e17de34cf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	My students crave challenge, they eat obstacle...	We are an urban, public k-5 elementary school....
29891	146723	p099708	c0a28c79fe8ad5810da49de47b3fb491	Mrs.	CA	2016-04-27 01:10:09	It's the end of the school year. Routines have...	My students desire challenges, movement, and c...
23374	72317	p087808	598621c141cda5fb184ee7e8ccdd3fcc	Ms.	CA	2016-04-27 02:04:15	Never has society so rapidly changed. Technolo...	Our Language Arts and Social Justice Magnet Sc...
49228	57854	p099430	4000cfe0c8b2df75a218347c1765e283	Ms.	IL	2016-04-27 07:19:44	My students yearn for a classroom environment ...	I have the privilege of teaching an incredible...

1.3 Text preprocessing

In [25]:

```
# 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)
```

Clean Essays (Text preprocessing)

In [26]:

```
clean_essay = []

for ess in tqdm(project_data["essay"]):
    ess = decontracted(ess)
    ess = ess.replace('\\r', ' ')
    ess = ess.replace('\\n', ' ')
    ess = ess.replace('\\n', ' ')
    ess = re.sub('[^A-Za-z0-9]+', ' ', ess)
    ess = ' '.join(f for f in ess.split() if f not in stopwords)
    clean_essay.append(ess.lower().strip())
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 50000/50000  
[00:44<00:00, 1134.92it/s]
```

In [27]:

```
project data["clean essays"] = clean_essay
```

In [28]:

```
project_data.drop(['essay'], axis=1, inplace=True)
```

Number of Words in Essay

In [29]:

```
essay word count = []
```

In [30]:

```
for ess in project_data["clean_essays"] :
    c = len(ess.split())
    essay word count.append(c)
```

In [31]:

```
project_data["essay word count"] = essay_word_count
```

In [32]:

```
project data.head(5)
```

Out[32]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_essay_1	project_essay_2
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	I recently read an article about giving studen	I teach at a low-income (Title 1) school. Ever

In [40]:

```
project_data.head(5)
```

Out[40]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_essay_1	project_essay_2
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	I recently read an article about giving studen... I teach at a low-income (Title 1) school. Ever...
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	My students crave challenge, they eat obstacle... We are an urban, public k-5 elementary school....
29891	146723	p099708	c0a28c79fe8ad5810da49de47b3fb491	Mrs.	CA	2016-04-27 01:10:09	It's the end of the school year. Routines have... My students desire challenges, movement, and c...
23374	72317	p087808	598621c141cda5fb184ee7e8ccdd3fcc	Ms.	CA	2016-04-27 02:04:15	Never has society so rapidly changed. Technolo... Our Language Arts and Social Justice Magnet Sc...
49228	57854	p099430	4000cfe0c8b2df75a218347c1765e283	Ms.	IL	2016-04-27 07:19:44	My students yearn for a classroom environment ... I have the privilege of teaching an incredible...

5 rows × 24 columns



Test - Train Split

In [41]:

```
# train test split

from sklearn.model_selection import train_test_split

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

In [42]:

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

Preparing data for models

In [43]:

```
project_data.columns
```

Out[43]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_essay_1', 'project_essay_2', 'project_essay_3',
      'project_essay_4', 'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'project_grade_category', 'clean_categories', 'clean_subcategories',
      'clean_titles', 'title_word_count', 'clean_essays', 'essay_word_count',
      'pos', 'neg', 'neu', 'compound'],
      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

Vectorizing Categorical data

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

One Hot Encode - Clean Categories of Projects

In [44]:

```
# we use count vectorizer to convert the values into one

from sklearn.feature_extraction.text import CountVectorizer

vectorizer_proj = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_proj.fit(X_train['clean_categories'].values)

categories_one_hot_train = vectorizer_proj.transform(X_train['clean_categories'].values)
categories_one_hot_test = vectorizer_proj.transform(X_test['clean_categories'].values)

print(vectorizer_proj.get_feature_names())

print("Shape of matrix of Train data after one hot encoding ",categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",categories_one_hot_test.shape)
```

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

```
Shape of matrix of Train data after one hot encoding (33500, 9)
```

```
Shape of matrix of Test data after one hot encoding (16500, 9)
```

One Hot Encode - Clean Sub-Categories of Projects

In [45]:

```
# we use count vectorizer to convert the values into one

vectorizer_sub_proj = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_sub_proj.fit(X_train['clean_subcategories'].values)

sub_categories_one_hot_train = vectorizer_sub_proj.transform(X_train['clean_subcategories'].values)
sub_categories_one_hot_test = vectorizer_sub_proj.transform(X_test['clean_subcategories'].values)

print(vectorizer_sub_proj.get_feature_names())
print("Shape of matrix of Train data after one hot encoding ",sub_categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding ",sub_categories_one_hot_test.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 of Train data after one hot encoding (33500, 30)
Shape of matrix of Test data after one hot encoding (16500, 30)
```

One Hot Encode - School States

In [46]:

```
my_counter = Counter()
for state in project_data['school_state'].values:
    my_counter.update(state.split())
```

In [47]:

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

In [48]:

```
## we use count vectorizer to convert the values into one hot encoded features

vectorizer_states = CountVectorizer(vocabulary=list(sorted_school_state_cat_dict.keys()),
lowercase=False, binary=True)
vectorizer_states.fit(X_train['school_state'].values)

school_state_categories_one_hot_train = vectorizer_states.transform(X_train['school_state'].values)
school_state_categories_one_hot_test = vectorizer_states.transform(X_test['school_state'].values)

print(vectorizer_states.get_feature_names())

print("Shape of matrix of Train data after one hot encoding", school_state_categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding", school_state_categories_one_hot_test.shape)

['VT', 'WY', 'ND', 'MT', 'RI', 'NH', 'SD', 'NE', 'AK', 'DE', 'WV', 'ME', 'NM', 'HI', 'DC', 'KS', 'ID', 'IA', 'AR', 'CO', 'MN', 'OR', 'MS', 'KY', 'NV', 'MD', 'TN', 'CT', 'AL', 'UT', 'WI', 'VA', 'AZ', 'NJ', 'OK', 'MA', 'LA', 'WA', 'MO', 'IN', 'OH', 'PA', 'MI', 'GA', 'SC', 'IL', 'NC', 'FL', 'TX', 'NY', 'CA']
Shape of matrix of Train data after one hot encoding (33500, 51)
Shape of matrix of Test data after one hot encoding (16500, 51)
```

One Hot Encode - Project Grade Category

In [49]:

```
my_counter = Counter()
for project_grade in project_data['project_grade_category'].values:
    my_counter.update(project_grade.split())
```

In [50]:

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

In [51]:

```
## we use count vectorizer to convert the values into one hot encoded features

vectorizer_grade = CountVectorizer(vocabulary=list(sorted_project_grade_cat_dict.keys()),
lowercase=False, binary=True)
vectorizer_grade.fit(X_train['project_grade_category'].values)

project_grade_categories_one_hot_train =
```

```

vectorizer_grade.transform(X_train['project_grade_category'].values)
project_grade_categories_one_hot_test = vectorizer_grade.transform(X_test['project_grade_category'].values)

print(vectorizer_grade.get_feature_names())

print("Shape of matrix of Train data after one hot encoding", project_grade_categories_one_hot_train.shape)
print("Shape of matrix of Test data after one hot encoding", project_grade_categories_one_hot_test.shape)

['Grades_9_12', 'Grades_6_8', 'Grades_3_5', 'Grades_PreK_2']
Shape of matrix of Train data after one hot encoding (33500, 4)
Shape of matrix of Test data after one hot encoding (16500, 4)

```

One Hot Encode - Teacher Prefix

In [52]:

```
project_data["teacher_prefix"].fillna(" ", inplace = True)
```

In [53]:

```

teacher_prefix = []

for i in range(len(project_data)):
    a = project_data["teacher_prefix"][i].replace('.', ' ')
    teacher_prefix.append(a)

```

In [54]:

```
project_data.drop(['teacher_prefix'], axis=1, inplace=True)
```

In [55]:

```
project_data["teacher_prefix"] = teacher_prefix
```

In [56]:

```

my_counter = Counter()
for teacher_prefix in project_data['teacher_prefix'].values:
    teacher_prefix = str(teacher_prefix)
    my_counter.update(teacher_prefix.split())

```

In [57]:

```

teacher_prefix_cat_dict = dict(my_counter)
sorted_teacher_prefix_cat_dict = dict(sorted(teacher_prefix_cat_dict.items(), key=lambda kv: kv[1]))

```

In [58]:

```

## we use count vectorizer to convert the values into one hot encoded features
## Unlike the previous Categories this category returns a
## ValueError: np.nan is an invalid document, expected byte or unicode string.
## The link below explains how to tackle such discrepancies.
## https://stackoverflow.com/questions/39303912/tfidfvectorizer-in-scikit-learn-valueerror-np-nan-is-an-invalid-document/39308809#39308809

vectorizer_teacher = CountVectorizer(vocabulary=list(sorted_teacher_prefix_cat_dict.keys()), lower
case=False, binary=True)
vectorizer_teacher.fit(X_train['teacher_prefix'].values.astype("U"))

teacher_prefix_categories_one_hot_train = vectorizer_teacher.transform(X_train['teacher_prefix'].v
alues.astype("U"))
teacher_prefix_categories_one_hot_test =
vectorizer_teacher.transform(X_test['teacher_prefix'].values.astype("U"))

print(vectorizer_teacher.get_feature_names())

```

```
print(vectorizer_teacher.get_feature_names())

print("Shape of matrix after one hot encoding ",teacher_prefix_categories_one_hot_train.shape)
print("Shape of matrix after one hot encoding ",teacher_prefix_categories_one_hot_test.shape)
```

```
['Dr', 'Teacher', 'Mr', 'Ms', 'Mrs']
Shape of matrix after one hot encoding (33500, 5)
Shape of matrix after one hot encoding (16500, 5)
```

Vectorizing Text data

Bag of words

Train Data

In [59]:

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

vectorizer_bow_essay = CountVectorizer(min_df=10)

vectorizer_bow_essay.fit(X_train["clean_essays"])

text_bow_train = vectorizer_bow_essay.transform(X_train["clean_essays"])

print("Shape of matrix after one hot encoding ",text_bow_train.shape)
```

```
Shape of matrix after one hot encoding (33500, 10463)
```

Test Data

In [60]:

```
text_bow_test = vectorizer_bow_essay.transform(X_test["clean_essays"])
print("Shape of matrix after one hot encoding ",text_bow_test.shape)
```

```
Shape of matrix after one hot encoding (16500, 10463)
```

Train Data - Titles

In [61]:

```
vectorizer_bow_title = CountVectorizer(min_df=10)

vectorizer_bow_title.fit(X_train["clean_titles"])

title_bow_train = vectorizer_bow_title.transform(X_train["clean_titles"])
print("Shape of matrix after one hot encoding ",title_bow_train.shape)
```

```
Shape of matrix after one hot encoding (33500, 1631)
```

Test Data - Titles

In [62]:

```
title_bow_test = vectorizer_bow_title.transform(X_test["clean_titles"])
print("Shape of matrix after one hot encoding ",title_bow_test.shape)
```

```
Shape of matrix after one hot encoding (16500, 1631)
```

TFIDF vectorizer with min_df=10

Train Data - Essays

In [63]:

```
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer_tfidf_essay = TfidfVectorizer(min_df=10)
vectorizer_tfidf_essay.fit(X_train["clean_essays"])

text_tfidf_train = vectorizer_tfidf_essay.transform(X_train["clean_essays"])
print("Shape of matrix after one hot encoding ",text_tfidf_train.shape)
```

Shape of matrix after one hot encoding (33500, 10463)

Test Data - Essays

In [64]:

```
text_tfidf_test = vectorizer_tfidf_essay.transform(X_test["clean_essays"])
print("Shape of matrix after one hot encoding ",text_tfidf_test.shape)
```

Shape of matrix after one hot encoding (16500, 10463)

TFIDF-Train Data - Titles

In [65]:

```
vectorizer_tfidf_titles = TfidfVectorizer(min_df=10)

vectorizer_tfidf_titles.fit(X_train["clean_titles"])
title_tfidf_train = vectorizer_tfidf_titles.transform(X_train["clean_titles"])
print("Shape of matrix after one hot encoding ",title_tfidf_train.shape)
```

Shape of matrix after one hot encoding (33500, 1631)

TFIDF - Test Data - Titles

In [66]:

```
title_tfidf_test = vectorizer_tfidf_titles.transform(X_test["clean_titles"])
print("Shape of matrix after one hot encoding ",title_tfidf_test.shape)
```

Shape of matrix after one hot encoding (16500, 1631)

Using Pretrained Models : AVG W2V

In [67]:

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

Train - Essays

In [68]:

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

avg_w2v_vectors_train = [];

for sentence in tqdm(X_train["clean_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
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 33500/33500  
[00:18<00:00, 1825.67it/s]
```

Test - Essays

```
for sentence in tqdm(X_test["clean_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_test.append(vector)

print(len(avg_w2v_vectors_test))
print(len(avg_w2v_vectors_test[0]))
```

```
100% |██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:03<00:00, 4673.59it/s]
```

Train - Titles

```
avg_w2v_vectors_titles_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train["clean_titles"]): # for each title
    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_titles_train.append(vector)

print(len(avg_w2v_vectors_titles_train))
print(len(avg_w2v_vectors_titles_train[0]))
```

```
100% |██████████████████████████████████████████████████████████████████████████| 33500/33500  
[00:00<00:00, 84951.29it/s]
```


33500
300

Test - Titles

In [71]:

```
# Similarly you can vectorize for title also

avg_w2v_vectors_titles_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test["clean_titles"]): # for each title
    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_titles_test.append(vector)

print(len(avg_w2v_vectors_titles_test))
print(len(avg_w2v_vectors_titles_test[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 78827.64it/s]
```

16500
300

Using Pretrained Models: TFIDF weighted W2V

Train - Essays

In [72]:

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

In [73]:

```
# 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["clean_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_train.append(vector)

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

```
100%|███████████████████████████████████████████████████████████████| 33500/33500 [02:  
12<00:00, 252.64it/s]
```

33500
300

Test - Essays

In [74]:

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

```
tfidf_w2v_vectors_test = []; # the avg-w2v for each sentence/review is stored in this list  
for sentence in tqdm(X_test["clean_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_test.append(vector)
```

```
print(len(tfidf_w2v_vectors_test))  
print(len(tfidf_w2v_vectors_test[0]))
```

100%|██████████████████████████████████████████████████████████████████████████████| 16500/16500 [01:

13<00:00, 223.90it/s]

16500
300

Train - Titles

In [75]:

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

In [76]:

```
# compute average word2vec for each review.

tfidf_w2v_vectors_titles_train = [];

for sentence in tqdm(X_train["clean_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_titles_train.append(vector)

print(len(tfidf_w2v_vectors_titles_train))
print(len(tfidf_w2v_vectors_titles_train[0]))
```

33500
300

16500
300

In [80]:

```

from sklearn.preprocessing import Normalizer

normalizer = Normalizer()

# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.

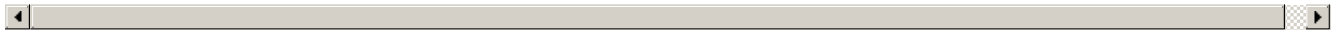
normalizer.fit(X_train['price'].values.reshape(-1,1))

price_train = normalizer.transform(X_train['price'].values.reshape(-1,1))
price_test = normalizer.transform(X_test['price'].values.reshape(-1,1))

print("After vectorizations")
print(price_train.shape, y_train.shape)
print(price_test.shape, y_test.shape)
print("="*100)

```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====



B) Quantity

In [81]:

```

normalizer = Normalizer()

# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.

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

quantity_train = normalizer.transform(X_train['quantity'].values.reshape(-1,1))
quantity_test = normalizer.transform(X_test['quantity'].values.reshape(-1,1))

print("After vectorizations")
print(quantity_train.shape, y_train.shape)
print(quantity_test.shape, y_test.shape)
print("="*100)

```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====



C) Number of Projects previously proposed by Teacher

In [82]:

```

normalizer = Normalizer()

# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.

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

prev_projects_train = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

```

```

.values.reshape(-1,1))
prev_projects_test = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

print("After vectorizations")
print(prev_projects_train.shape, y_train.shape)
print(prev_projects_test.shape, y_test.shape)
print("="*100)

```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====



D) Title word Count

In [83]:

```

normalizer = Normalizer()

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

title_word_count_train = normalizer.transform(X_train['title_word_count'].values.reshape(-1,1))
title_word_count_test = normalizer.transform(X_test['title_word_count'].values.reshape(-1,1))

print("After vectorizations")
print(title_word_count_train.shape, y_train.shape)
print(title_word_count_test.shape, y_test.shape)
print("="*100)

```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====



E) Essay word Count

In [84]:

```

normalizer = Normalizer()

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

essay_word_count_train = normalizer.transform(X_train['essay_word_count'].values.reshape(-1,1))
essay_word_count_test = normalizer.transform(X_test['essay_word_count'].values.reshape(-1,1))

print("After vectorizations")
print(essay_word_count_train.shape, y_train.shape)
print(essay_word_count_test.shape, y_test.shape)
print("="*100)

```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====



F) Essay Sentiments - pos

In [85]:

```

normalizer = Normalizer()

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

essay_sent_pos_train = normalizer.transform(X_train['pos'].values.reshape(-1,1))
essay_sent_pos_test = normalizer.transform(X_test['pos'].values.reshape(-1,1))

print("After vectorizations")

```

```
print("After vectorizations")
print(essay_sent_pos_train.shape, y_train.shape)
print(essay_sent_pos_test.shape, y_test.shape)
print("="*100)
```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====

G) Essay Sentiments - neg

In [86]:

```
normalizer = Normalizer()

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

essay_sent_neg_train = normalizer.transform(X_train['neg'].values.reshape(-1,1))
essay_sent_neg_test = normalizer.transform(X_test['neg'].values.reshape(-1,1))

print("After vectorizations")
print(essay_sent_neg_train.shape, y_train.shape)
print(essay_sent_neg_test.shape, y_test.shape)
print("="*100)
```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====

H) Essay Sentiments - neu

In [87]:

```
normalizer = Normalizer()

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

essay_sent_neu_train = normalizer.transform(X_train['neu'].values.reshape(-1,1))
essay_sent_neu_test = normalizer.transform(X_test['neu'].values.reshape(-1,1))

print("After vectorizations")
print(essay_sent_neu_train.shape, y_train.shape)
print(essay_sent_neu_test.shape, y_test.shape)
print("="*100)
```

After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====

I) Essay Sentiments - compound

In [88]:

```
normalizer = Normalizer()

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

essay_sent_comp_train = normalizer.transform(X_train['compound'].values.reshape(-1,1))
essay_sent_comp_test = normalizer.transform(X_test['compound'].values.reshape(-1,1))

print("After vectorizations")
print(essay_sent_comp_train.shape, y_train.shape)
print(essay_sent_comp_test.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
```

Assignment 8: DT

1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V)




2. Hyper paramter tuning (best `depth` in range [1, 5, 10, 50, 100, 500, 100], and the best `min_samples_split` in range [5, 10, 100, 500])

- 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. Graphviz

- Visualize your decision tree with Graphviz. It helps you to understand how a decision is being made, given a new vector.
- Since feature names are not obtained from word2vec related models, visualize only BOW & TFIDF decision trees using Graphviz
- Make sure to print the words in each node of the decision tree instead of printing its index.
- Just for visualization purpose, limit max_depth to 2 or 3 and either embed the generated images of graphviz in your notebook, or directly upload them as .png files.

4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure 
- 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 
- Once after you plot the confusion matrix with the test data, get all the `false positive data points`
 - Plot the WordCloud [WordCloud](#)
 - Plot the box plot with the `price` of these `false positive data points`
 - Plot the pdf with the `teacher_number_of_previously_posted_projects` of these `false positive data points`

5. [Task-2]

- Select 5k best features from features of **Set 2** using `feature_importances_`, discard all the other remaining features and then apply any of the model of your choice i.e. (Decision tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure 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](#)
-

2. Decision Tree

Set 1: Categorical, Numerical features + Project_title(BOW) + Preprocessed_essay (BOW with min_df=10)

In [89]:

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

X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train,
school_state_categories_one_hot_train, project_grade_categories_one_hot_train,
teacher_prefix_categories_one_hot_train, price_train, quantity_train, prev_projects_train, title_word_count_train, essay_word_count_train, essay_sent_pos_train, essay_sent_neg_train,
essay_sent_neu_train, essay_sent_comp_train, title_bow_train, text_bow_train)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test,
school_state_categories_one_hot_test, project_grade_categories_one_hot_test,
teacher_prefix_categories_one_hot_test, price_test, quantity_test, prev_projects_test,
title_word_count_test, essay_word_count_test, essay_sent_pos_test, essay_sent_neg_test, essay_sent_neu_test, essay_sent_comp_test, title_bow_test, text_bow_test)).tocsr()
```

In [90]:

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

```
Final Data matrix
(33500, 12202) (33500,)
(16500, 12202) (16500,)
```

=====



In [91]:

```
bow_features_names = []
```

In [92]:

```
## Obtain Feature names for Project titles

for a in vectorizer_proj.get_feature_names() :
    bow_features_names.append(a)
```

In [93]:

```
## Obtain Feature names for Project Sub-titles
for a in vectorizer_sub_proj.get_feature_names() :
    bow_features_names.append(a)
```

In [94]:

```
for a in vectorizer_states.get_feature_names() :
    bow_features_names.append(a)
```

In [95]:

```
## Obtain Feature names for Project Grade Category
for a in vectorizer_grade.get_feature_names() :
    bow_features_names.append(a)
```

In [96]:

```
## Obtain Feature names for Teacher Title
for a in vectorizer_teacher.get_feature_names() :
    bow_features_names.append(a)
```

In [97]:

```
len(bow_features_names)
bow_features_names.append("price")
bow_features_names.append("quantity")
bow_features_names.append("prev_projects")
```



```
bow_features_names.append("prev_proposed_projects")
bow_features_names.append("title_word_count")
bow_features_names.append("essay_word_count")
bow_features_names.append("pos")
bow_features_names.append("neg")
bow_features_names.append("neu")
bow_features_names.append("compound")
for a in vectorizer_bow_title.get_feature_names():
    bow_features_names.append(a)
for a in vectorizer_bow_essay.get_feature_names():
    bow_features_names.append(a)
```

In [98]:

```
len(bow_features_names)
```

Out[98]:

12202

In [99]:

```
from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(max_depth=3)
```

In [100]:

```
clf = dtree.fit(X_tr, y_train)
```

In [101]:

```
# Visualize data
import graphviz
from sklearn import tree
from graphviz import Source

dot_data = tree.export_graphviz(dtree, out_file=None, feature_names=bow_features_names)
graph = graphviz.Source(dot_data)
graph.render("Bow tree",view = True)
```

Out[101]:

'Bow tree.pdf'

GridSearchCV (K fold Cross Validation)

In [102]:

```
from sklearn.model_selection import GridSearchCV
from sklearn.tree import DecisionTreeClassifier
```

In [103]:

```
dt = DecisionTreeClassifier()

parameters = {'max_depth':[1, 5, 10, 50, 100, 500], 'min_samples_split': [5, 10, 100, 500]}

clf = GridSearchCV(dt, parameters, cv= 3, scoring='roc_auc')

clf.fit(X_tr, y_train)

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

In [104]:

```
train_auc
```

Out[104]:

```
array([0.55239772, 0.55239772, 0.55239772, 0.55239772, 0.62143325,
       0.6212722 , 0.62076759, 0.61991603, 0.67250569, 0.67096022,
       0.66556902, 0.65599991, 0.89613622, 0.88928822, 0.86001558,
       0.79319391, 0.97143957, 0.96371831, 0.93397666, 0.87971578,
       0.99866876, 0.99418318, 0.97597894, 0.93643729])
```

In [105]:

```
cv_auc
```

Out[105]:

```
array([0.55239776, 0.55239776, 0.55239776, 0.55239776, 0.59770303,
       0.59816312, 0.59766831, 0.59761585, 0.59863186, 0.59825058,
       0.59806453, 0.59862196, 0.49965871, 0.50495851, 0.52408317,
       0.56552297, 0.48258742, 0.48756936, 0.50551413, 0.53851181,
       0.50973769, 0.51760409, 0.52400558, 0.54606139])
```

Plot for Train Data

In [106]:

```
import plotly.plotly as py
import plotly.graph_objs as go
```

In [107]:

```
import plotly
plotly.tools.set_credentials_file(username='chandana_saribala', api_key='Z5CHrYM8sPITAbv8u8MS')
```

In [108]:

```
x1 = [[0.55120541, 0.55120541, 0.55120541, 0.55120541, 0.62279835,
       0.62257113, 0.62243708, 0.62212488, 0.66744193, 0.6654276 ,
       0.65909786, 0.65098396, 0.89711643, 0.88350259, 0.85558414,
       0.79684316, 0.97391624, 0.96529621, 0.94181771, 0.89162551,
       0.99880609, 0.99385652, 0.97615667, 0.94150511, 0.9698162 ,
       0.96413439, 0.9403805 , 0.88949701]]
```

In [109]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,10,50,50,50,50,50,100,100,100,100,100,500,500,500,500,500], index = x1)
```

In [110]:

```
y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x1)
```

In [111]:

```
trace = go.Scatter3d(
    x=x1, y=y1, z=z1,
    marker=dict(
        size=4,
        color=z1,
        colorscale='Viridis',
    ),
    line=dict(
        color='#1f77b4',
        width=1
    )
)
```

In [112]:

```
data = [trace]
```

In [113]:

```
layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- TRAIN Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        ),
        aspectratio = dict( x=1, y=1, z=0.7 ),
        aspectmode = 'manual'
    ),
)
```

In [114]:

```
fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-a', height=700)
```

High five! You successfully sent some data to your account on plotly. View your plot in your browser at https://plot.ly/~chandana_saribala/0 or inside your plot.ly account where it is named 'Decision-trees-a'

Out[114]:

Plot for Cross Validation Data

In [115]:

```
x2 = [[0.55120507, 0.55120507, 0.55120507, 0.55120507, 0.60567636,  
        0.60607942, 0.60566638, 0.60527114, 0.60029542, 0.60164221,  
        0.60510206, 0.60950152, 0.52622504, 0.52801446, 0.54987881,  
        0.58466911, 0.50442353, 0.50606574, 0.51938558, 0.55616451,  
        0.51926658, 0.5202411 , 0.53729777, 0.55517838, 0.49877934,  
        0.50904987, 0.52035284, 0.54911636]]
```

In [116]:

```
y2 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x2)
```

In [117]:

```
z2 = pd.Series([1,1,1,1,1,5,5,5,5,5,10,10,10,10,10,50,50,50,50,50,100,100,100,100,500,500,500,500], index = x2)
```

In [118]:

```
trace = go.Scatter3d(  
    x=x2, y=y2, z=z2,  
    marker=dict(  
        size=4,  
        color=z2,  
        colorscale='Viridis',  
    ),  
    line=dict(  
        color='#1f77b4',  
        width=1  
    )  
)
```

In [119]:

```
data = [trace]
```

In [120]:

```

layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- Cross Validation Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        ),
        aspectratio = dict( x=1, y=1, z=0.7 ),
        aspectmode = 'manual'
    ),
)

```

In [121]:

```

fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-b', height=900)

```

Out[121]:

Observations : 1) We understand from the 2 plots that the Decision Tree with a depth of 100 & 500 performs great on Training Data but performs pretty bad on unseen data (cross validation data) . => Probably a case of Overfitting. 2) Decision trees with depth 1 & 5 performs poor on both Train data as well as Cross Validation Data. =>Probably a case of Underfitting. 3) Decision Tree with maximum depth 10, performs decently on both Train as well as Cross Validation Data. 4) 100 as the value for Minimum samples per split is considered.

In [122]:

```
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
import seaborn as sea
```

In [127]:

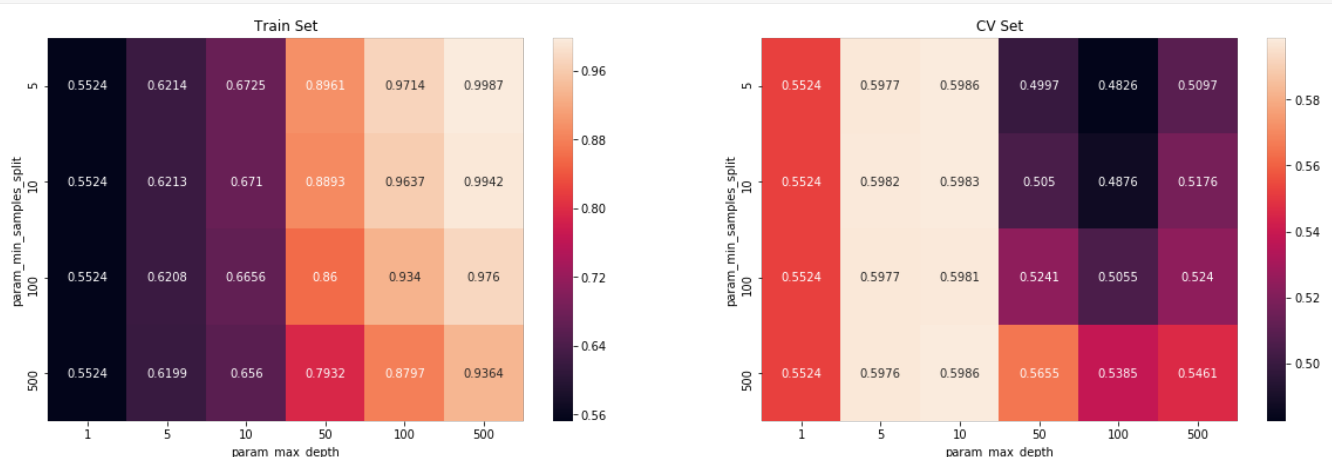
```
max_scores = pd.DataFrame(clf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth'])
.max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))

sea.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sea.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Train the model using the best hyper parameter value

In [128]:

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

```

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

return y_data_pred

```

In [129]:

```

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

model = DecisionTreeClassifier(max_depth = 10, min_samples_split = 100)

model.fit(X_tr, y_train)

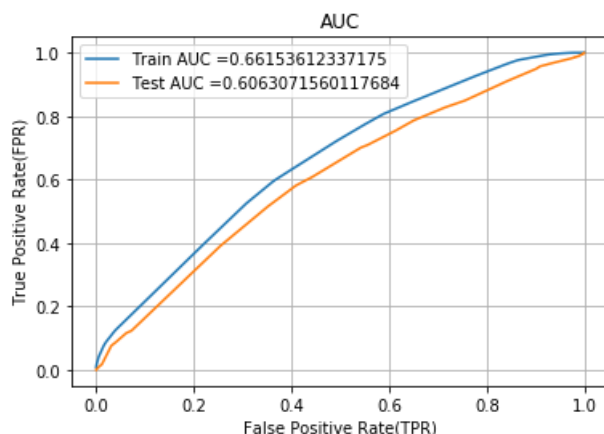
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

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("False Positive Rate(TPR)")
plt.ylabel("True Positive Rate(FPR)")
plt.title("AUC")
plt.grid()
plt.show()

```



Confusion Matrix

In [130]:

```

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

```

Train Data

In [131]:

```

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) 0.24999996255834908 for threshold 0.835
[[ 2583  2585]
 [ 7687 20645]]

```

In [132]:

```

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) 0.24999996255834908 for threshold 0.835

```

In [133]:

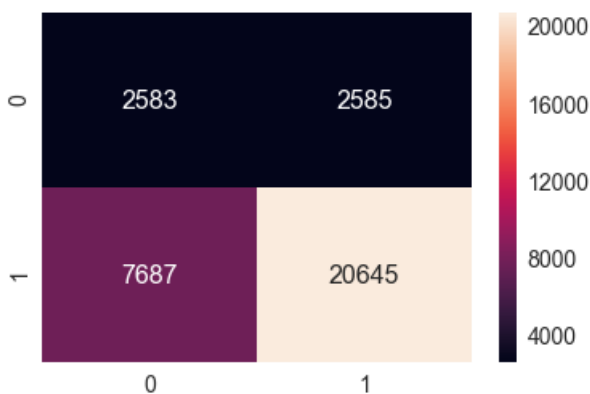
```

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[133]:

<matplotlib.axes._subplots.AxesSubplot at 0x2b3b57f5a90>



Test Data

In [134]:

```

print("="*100)
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*(1-fpr) 0.24820058758590557 for threshold 0.835
[[1130 1416]
 [4058 9896]]

```

In [135]:


```
In [135]:
```

```
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))
```

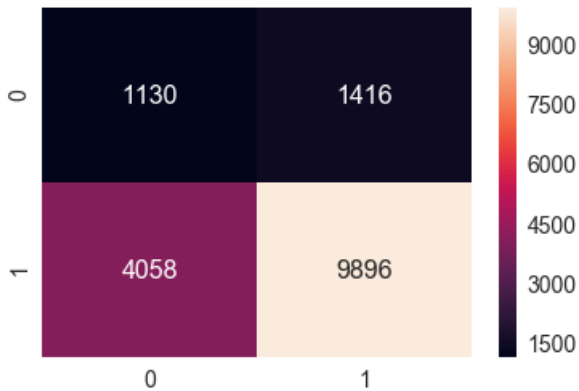
the maximum value of $tpr \cdot (1 - fpr)$ 0.24820058758590557 for threshold 0.835

```
In [136]:
```

```
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[136]:
```

<matplotlib.axes._subplots.AxesSubplot at 0x2b3b4944fd0>



Applying False positives from BOW encoded essays

```
In [137]:
```

```
bow_test = text_bow_test.todense()
```

```
In [138]:
```

```
bow_test.shape
```

```
Out[138]:
```

(16500, 10463)

```
In [139]:
```

```
vectorizer_bow_essay = CountVectorizer(min_df=10)
av = vectorizer_bow_essay.fit(X_train["clean_essays"])
```

```
In [140]:
```

```
bow_features = av.get_feature_names()
```

```
In [141]:
```

```
len(bow_features)
```

```
Out[141]:
```

10463

```
In [142]:
```

```
y_test_converted = list(y_test[:, :])
```

In [143]:

```
false_positives_index_a = []
fp_count = 0

for i in tqdm(range(len(y_test_pred))):
    if y_test_converted[i] == 0 and y_test_pred[i] <= 0.839:
        false_positives_index_a.append(i)
        fp_count = fp_count + 1
    else :
        continue
```

```
100%|██████████████████████████████████████████████████████████████████████████| 16500/16500  
[00:00<00:00, 2462409.39it/s]
```

In [144]:

fp count

Out [144] :

1165

In [145]:

```
false_positives_index_a[0:5]
```

Out[145]:

[1, 6, 13, 29, 32]

In [146]:

```
df1 = pd.DataFrame(bow_test)
```

In [147]:

```
df1_final = df1.iloc[false_positives_index_a,:]
```

In [148]:

```
df1_final.shape
```

Out[148]:

(1165, 10463)

In [149]:

```
df1 final[0].sum()
```

Out [149] :

3

In [150]:

```
best_indices = []

for j in range(10453):

    s = df1_final[j].sum()

    if s >= 100 :
        best_indices.append(j)
    else :
        continue
```

```
In [151]:
```

```
len(best_indices)
```

```
Out[151]:
```

```
294
```

```
In [152]:
```

```
best_indices[0:10]
```

```
Out[152]:
```

```
[181, 182, 207, 220, 291, 298, 482, 488, 502, 512]
```

```
In [153]:
```

```
bow_features[0:10]
```

```
Out[153]:
```

```
['00', '000', '10', '100', '1000', '101', '104', '105', '10th', '11']
```

```
In [154]:
```

```
fp_words = []
```

```
for a in best_indices :  
    fp_words.append(str(bow_features[a]))
```

```
In [155]:
```

```
fp_words[0:10]
```

```
Out[155]:
```

```
['ability',  
 'able',  
 'academic',  
 'access',  
 'active',  
 'activities',  
 'all',  
 'allow',  
 'also',  
 'always']
```

Word Cloud for False Positives words

```
In [158]:
```

```
from wordcloud import WordCloud
```

```
In [159]:
```

```
#convert list to string and generate  
unique_string=(" ").join(fp_words)  
wordcloud = WordCloud(width = 1000, height = 500).generate(unique_string)  
plt.figure(figsize=(25,10))  
plt.imshow(wordcloud)  
plt.axis("off")  
plt.savefig("your_file_name"+" .png", bbox_inches='tight')  
plt.show()  
plt.close()
```


them are Extremely costs costing more than 3000 Dollars.

PDF with the Teacher_number_of_previously_posted_projects of these False Positive data points

In [164]:

```
df3 = pd.DataFrame(X_test['teacher_number_of_previously_posted_projects'])
```

In [165]:

```
df3_final = df3.iloc[false_positives_index_a,:]
```

In [166]:

```
df3_final.shape
```

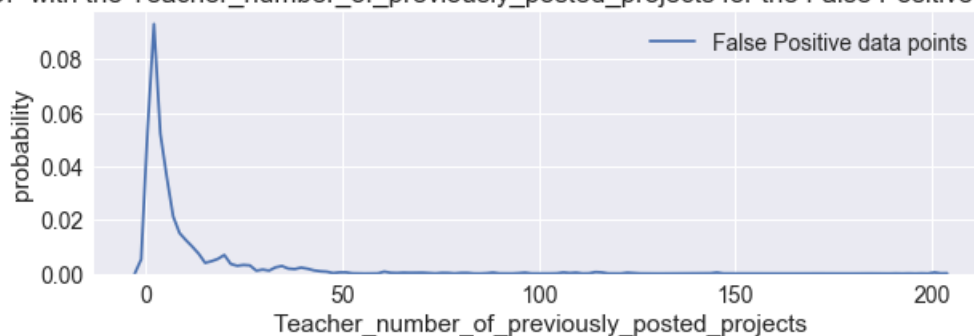
Out[166]:

```
(1165, 1)
```

In [167]:

```
plt.figure(figsize=(10,3))
sns.distplot(df3_final.values, hist=False, label="False Positive data points")
plt.title('PDF with the Teacher_number_of_previously_posted_projects for the False Positive data points')
plt.xlabel('Teacher_number_of_previously_posted_projects')
plt.ylabel('probability')
plt.legend()
plt.show()
```

PDF with the Teacher_number_of_previously_posted_projects for the False Positive data points



Inference:

Majority of the cases have Teachers with previously posted projects as 0 (which is nearly 10% of the total data)

Set 2 : Categorical, Numerical features + Project_title(TFIDF) + Preprocessed_essay (TFIDF min_df=10)

In [168]:

```
X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train,
school_state_categories_one_hot_train, project_grade_categories_one_hot_train,
teacher_prefix_categories_one_hot_train, price_train, quantity_train, prev_projects_train, title_word_count_train,
essay_word_count_train, essay_sent_pos_train, essay_sent_neg_train,
essay_sent_neu_train, essay_sent_comp_train, title_tfidf_train, text_tfidf_train)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test,
school_state_categories_one_hot_test, project_grade_categories_one_hot_test,
teacher_prefix_categories_one_hot_test, price_test, quantity_test, prev_projects_test,
title_word_count_test, essay_word_count_test, essay_sent_pos_test, essay_sent_neg_test, essay_sent_neu_test,
essay_sent_comp_test, title_tfidf_test, text_tfidf_test)).tocsr()
```

In [169]:

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

```
Final Data matrix
(33500, 12202) (33500,)
(16500, 12202) (16500,)
```

=====

In [170]:

```
tfidf_features_names = []
## Obtain Feature names for Project titles

for a in vectorizer_proj.get_feature_names() :
    tfidf_features_names.append(a)

## Obtain Feature names for Project Sub-titles
for a in vectorizer_sub_proj.get_feature_names() :
    tfidf_features_names.append(a)

## Obtain Feature names for states
for a in vectorizer_states.get_feature_names() :
    tfidf_features_names.append(a)

## Obtain Feature names for Project Grade Category
for a in vectorizer_grade.get_feature_names() :
    tfidf_features_names.append(a)

## Obtain Feature names for Teacher Title
for a in vectorizer_teacher.get_feature_names() :
    tfidf_features_names.append(a)

tfidf_features_names.append("price")

tfidf_features_names.append("quantity")

tfidf_features_names.append("prev_proposed_projects")

tfidf_features_names.append("title_word_count")

tfidf_features_names.append("essay_word_count")

tfidf_features_names.append("positive_sentiment")

tfidf_features_names.append("negative_sentiment")

tfidf_features_names.append("neutral_sentiment")

tfidf_features_names.append("compound")

for a in vectorizer_tfidf_titles.get_feature_names() :
    tfidf_features_names.append(a)

for a in vectorizer_tfidf_essay.get_feature_names() :
    tfidf_features_names.append(a)
```

In [171]:

```
len(tfidf_features_names)
```

Out[171]:

12202

In [172]:

```
from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(max_depth=3)
```

In [173]:

```
clf = dtree.fit(X_tr, y_train)
```

In [174]:

```
# Visualize data
import graphviz
from sklearn import tree
from graphviz import Source

dot_data = tree.export_graphviz(dtree, out_file=None, feature_names=tfidf_features_names)
graph = graphviz.Source(dot_data)
graph.render("Tfidf tree", view = True)
```

Out[174]:

'Tfidf tree.pdf'

GridSearchCV (K fold Cross Validation)

In [175]:

```
dt_2 = DecisionTreeClassifier()

parameters = {'max_depth':[1, 5, 10, 50, 100, 500], 'min_samples_split': [5, 10, 100, 500]}

clf = GridSearchCV(dt_2, parameters, cv= 3, scoring='roc_auc')

clf.fit(X_tr, y_train)

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

In [176]:

train_auc

Out[176]:

```
array([0.55426048, 0.55426048, 0.55426048, 0.55426048, 0.62139423,
        0.62130717, 0.62050275, 0.62034926, 0.67463554, 0.6734489 ,
        0.66371938, 0.65805757, 0.87979657, 0.87301819, 0.83494096,
        0.79601948, 0.95725685, 0.95344406, 0.92289302, 0.88957186,
        0.99932173, 0.99595624, 0.97773142, 0.96173474, 0.9993614 ,
        0.99568348, 0.97743846, 0.96067759])
```

In [177]:

cv_auc

Out[177]:

```
array([0.55401884, 0.55401884, 0.55401884, 0.55401884, 0.60051756,
        0.6001876 , 0.60079485, 0.60048695, 0.59315466, 0.59300666,
        0.6008911 , 0.60514658, 0.50072369, 0.50797618, 0.54699355,
        0.56864204, 0.48792224, 0.49317578, 0.5284145 , 0.53963103,
        0.52677872, 0.53287878, 0.55017869, 0.55427039, 0.52921551,
        0.53412304, 0.54876724, 0.55741219])
```

Plot for Train Data

In [178]:

```
x1 = [[0.55278549, 0.55278549, 0.55278549, 0.55278549, 0.62550369,
        0.62535175, 0.62471044, 0.62301915, 0.67253364, 0.67200378,
        0.66955318, 0.65998752, 0.87169988, 0.86733981, 0.83749587,
```

In [179]:

In [180]:

In [181]:

In [182]:

In [183]:

```

layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- TRAIN Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        )
    ),
)

```



```

        aspectratio = dict( x=1, y=1, z=0. / ),
        aspectmode = 'manual'
    ),
)

```

In [184]:

```

fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-c', height=900)

```

Out[184]:

Plot for Cross Validation Data

In [185]:

```

x1 = [0.55210814, 0.55210814, 0.55210814, 0.55210814, 0.60625189,
      0.6058133 , 0.60636619, 0.60776402, 0.60264809, 0.60200219,
      0.60876201, 0.61338102, 0.50811386, 0.51411963, 0.54883589,
      0.56141945, 0.50528828, 0.50864596, 0.53671867, 0.55859065,
      0.52998987, 0.53038417, 0.54293448, 0.55622729, 0.52577628,
      0.52917346, 0.54701332, 0.55743405]

```

In [186]:

```

y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x1)

```

In [187]:

In [187]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,50,50,50,50,100,100,100,100,500,500,500,500,1000,1000,1000,1000],index = x1)
```

In [188]:

```
trace = go.Scatter3d(
    x=x1, y=y1, z=z1,
    marker=dict(
        size=4,
        color=z1,
        colorscale='Viridis',
    ),
    line=dict(
        color='#1f77b4',
        width=1
    )
)
```

In [189]:

```
data = [trace]
```

In [190]:

```
layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- Cross Validation Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        ),
        aspectratio = dict( x=1, y=1, z=0.7 ),
        aspectmode = 'manual'
    ),
)
```

In [191]:

```
fig = dict(data=data, layout=layout)

py.ipplot(fig, filename='Decision-trees-d', height=900)
```

In [191]:

Out[191]:

Observations : 1) We understand from the 2 plots that the Decision Tree with a depth of 100, 500, 1000 performs great on Training Data but performs pretty bad on unseen data (cross validation data) . => Probably a case of Overfitting. 2) Decision trees with depth 1 & 5 performs poor on both Train data as well as Cross Validation Data. =>Probably a case of Underfitting. 3) Decision Tree with maximum depth 10, performs fairly well on both Train as well as Cross Validation Data. 4) 500 as the value for Minimum samples per split is considered.

Train the model using the best hyper parameter value

In [192]:

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

model = DecisionTreeClassifier(max_depth = 10, min_samples_split = 500)

model.fit(X_tr, y_train)

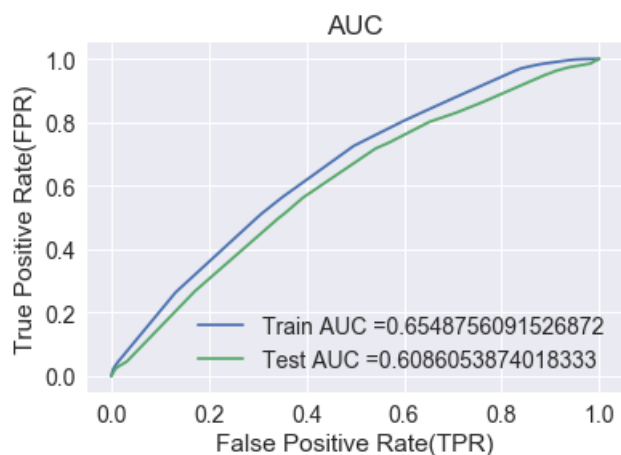
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

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("False Positive Rate(TPR)")
plt.ylabel("True Positive Rate(FPR)")
plt.title("AUC")
plt.show()
```



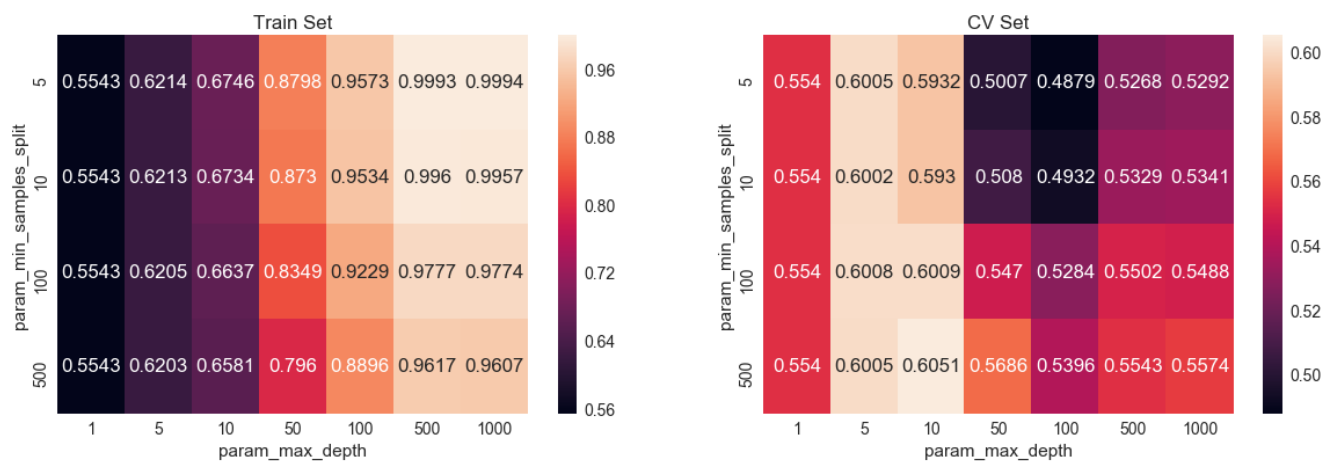
In [193]:

```
import seaborn as sea
max_scores = pd.DataFrame(clf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth'])
.max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))
sea.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sea.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Confusion Matrix

Train data

In [194]:

```
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 $tpr*(1-fpr)$ 0.2499891793628809 for threshold 0.858

```
[[ 2601  2567]
 [ 7778 20554]]
```

In [195]:

```
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))
```

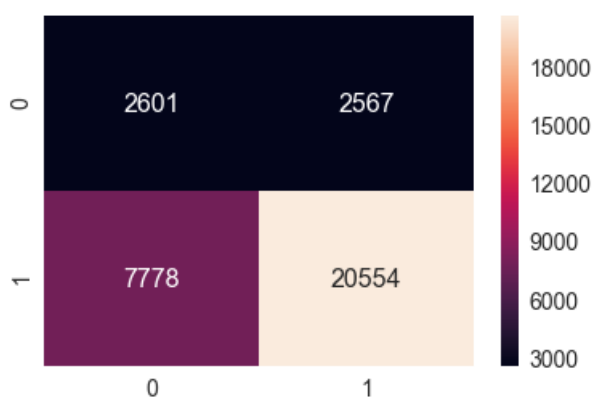
the maximum value of $tpr*(1-fpr)$ 0.2499891793628809 for threshold 0.858

In [196]:

```
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[196]:

<matplotlib.axes._subplots.AxesSubplot at 0x2b3bcaef5f8>



Test Data

In [197]:

```
print("="*100)
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 $tpr*(1-fpr)$ 0.24826661540768477 for threshold 0.858

```
[[ 1167  1379]
 [ 3942 10012]]
```

In [198]:

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

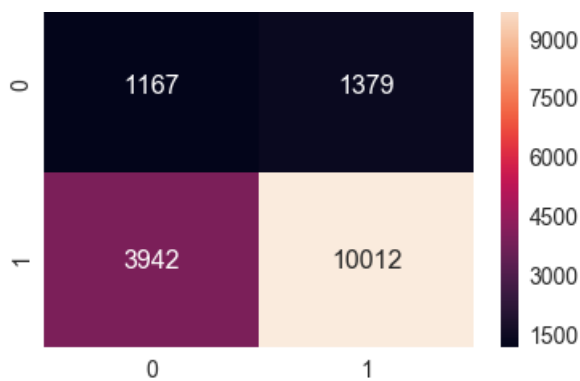
the maximum value of $tpr*(1-fpr)$ 0.24826661540768477 for threshold 0.858

In [199]:

```
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[199]:

<matplotlib.axes._subplots.AxesSubplot at 0x2b3bd4f89b0>



Obtaining the False Positive words from TFIDF encoded Essays

In [200]:

```
tfidf_test = text_tfidf_test.todense()

tfidf_test.shape
```

Out[200]:

```
(16500, 10463)
```

In [201]:

```
vectorizer_tfidf_essay = TfidfVectorizer(min_df=10)
bv = vectorizer_tfidf_essay.fit(X_train["clean_essays"])
```

In [202]:

```
tfidf_features = bv.get_feature_names()

len(tfidf_features)
```

Out[202]:

```
10463
```

In [203]:

```
y_test_converted = list(y_test[:, :])
```

In [204]:

```
false_positives_index_b = []
fp_count = 0

for i in tqdm(range(len(y_test_pred))):
    if y_test_converted[i] == 0 and y_test_pred[i] <= 0.84:
        false_positives_index_b.append(i)
        fp_count = fp_count + 1
    else :
        continue
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500
[00:00<00:00, 1942570.48it/s]
```

In [205]:

```
fp_count
```

Out[205]:

1167

In [206]:

```
false_positives_index_b[0:5]
```

Out[206]:

```
[29, 51, 81, 104, 119]
```

In [207]:

```
df2 = pd.DataFrame(tfidf_test)
df2_final = df2.iloc[false_positives_index_b,:]
```

In [208]:

```
df2_final.shape
```

Out[208]:

```
(1167, 10463)
```

In [209]:

```
best_indices_b = []

for j in range(10453):

    s = df2_final[j].sum()

    if s >= 10 :
        best_indices_b.append(j)
    else :
        continue
```

In [210]:

```
len(best_indices_b)
```

Out[210]:

```
170
```

In [211]:

```
best_indices_b[0:10]
```

Out[211]:

```
[182, 207, 220, 291, 298, 488, 502, 512, 518, 670]
```

In [212]:

```
tfidf_features[0:10]
```

Out[212]:

```
['00', '000', '10', '100', '1000', '101', '104', '105', '10th', '11']
```

In [213]:

```
fp_words_b = []

for a in best_indices_b :
    fp_words_b.append(str(tfidf_features[a]))
```

```
fp_words_b[0:10]
```

```
['able',
 'academic',
 'access',
 'active',
 'activities',
 'allow',
 'also',
 'always',
 'amazing',
 'area']
```

Word Cloud for False Positives words

```
#convert list to string and generate
unique_string=(" ").join(fp_words_b)
wordcloud = WordCloud(width = 1000, height = 500, background_color = 'white').generate(unique_string
)
plt.figure(figsize=(25,10))
plt.imshow(wordcloud)
plt.axis("off")
plt.savefig("Word_Cloud_tfidf"+" .png", bbox_inches='tight')
plt.show()
plt.close()
```



Box - Plot with the price of these False positive data points

```
df2 b = pd.DataFrame(X test['price'])
```

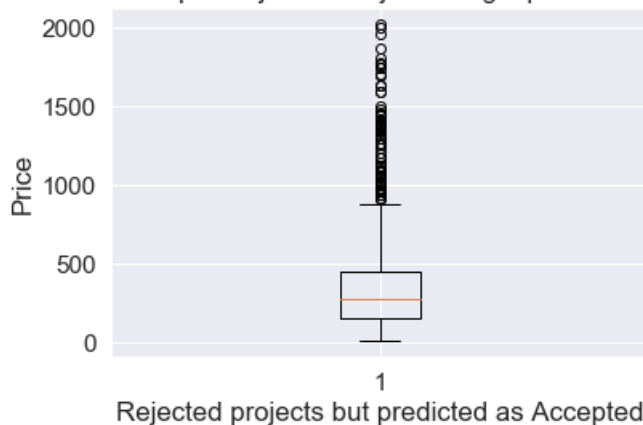
```
df2 b final = df2 b.iloc[false positives index b,:]
```

T. 10141.


```
In [214]:
```

```
plt.boxplot(df2_b_final.values)
plt.title('Box Plots of Cost per Rejected Project that got predicted as Accepted')
plt.xlabel('Rejected projects but predicted as Accepted')
plt.ylabel('Price')
plt.show()
```

Box Plots of Cost per Rejected Project that got predicted as Accepted



Inference 1) Majority of the projects that were rejected but predicted as accepted Costs almost less than 500 Dollars. 2) Good number of incorrectly prediced projects cost around 1000 - 2000 Dollars. 3) A Few of them are Extremely costly costing more than 2000 Dollars.

PDF with the Teacher_number_of_previously_posted_projects of these False Positive data points

```
In [215]:
```

```
df3_b = pd.DataFrame(X_test['teacher_number_of_previously_posted_projects'])
```

```
In [216]:
```

```
df3_b_final = df3_b.iloc[false_positives_index_b,:]
```

```
In [217]:
```

```
df3_b_final.shape
```

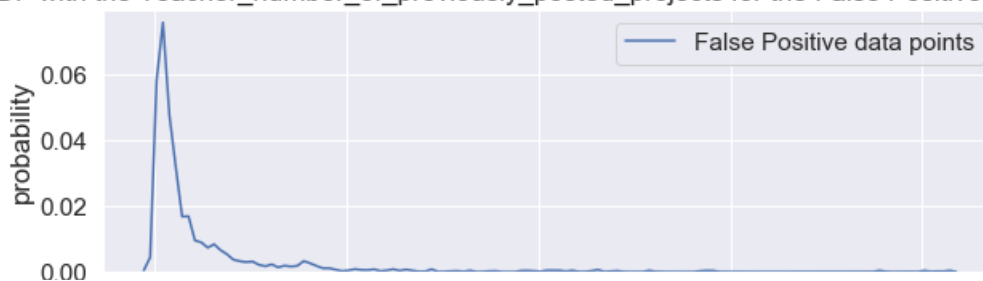
```
Out[217]:
```

```
(1268, 1)
```

```
In [218]:
```

```
plt.figure(figsize=(10,3))
sns.distplot(df3_b_final.values, hist=False, label="False Positive data points")
plt.title('PDF with the Teacher_number_of_previously_posted_projects for the False Positive data p
oints')
plt.xlabel('Teacher_number_of_previously_posted_projects')
plt.ylabel('probability')
plt.legend()
plt.show()
```

PDF with the Teacher_number_of_previously_posted_projects for the False Positive data points



0 50 100 150 200
Teacher_number_of_previously_posted_projects

Inference 1) Majority of the cases have Teachers with previously posted projects as 0. 2) Good percent of Teachers have 10 or fewer projects previously posted projects.

Set 3 : Categorical, Numerical features + Project_title(AVG W2V) + Preprocessed_essay (AVG W2V)

In [219]:

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

X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train,
school_state_categories_one_hot_train, project_grade_categories_one_hot_train,
teacher_prefix_categories_one_hot_train, price_train, quantity_train, prev_projects_train, title_w
rd_count_train, essay_word_count_train, essay_sent_pos_train, essay_sent_neg_train,
essay_sent_neu_train, essay_sent_comp_train, avg_w2v_vectors_train, avg_w2v_vectors_titles_train))
.tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test,
school_state_categories_one_hot_test, project_grade_categories_one_hot_test,
teacher_prefix_categories_one_hot_test, price_test, quantity_test, prev_projects_test,
title_word_count_test, essay_word_count_test, essay_sent_pos_test, essay_sent_neg_test, essay_sent_
neu_test, essay_sent_comp_test, avg_w2v_vectors_test, avg_w2v_vectors_titles_test)).tocsr()
```

In [220]:

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

```
Final Data matrix
(33500, 708) (33500,)
(16500, 708) (16500,)
```



GridSearchCV (K fold Cross Validation)

In [221]:

```
dt_3 = DecisionTreeClassifier()

parameters = {'max_depth':[1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}

clf = GridSearchCV(dt_3, parameters, cv= 3, scoring='roc_auc')

clf.fit(X_tr, y_train)

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

In [223]:

```
train_auc
```

Out[223]:

```
array([0.56261079, 0.56261079, 0.56261079, 0.56261079, 0.65106466,
0.65102953, 0.65095451, 0.64944063, 0.75917609, 0.75486241,
0.73047794, 0.69831248, 0.99278022, 0.98905981, 0.93318717,
0.81835938, 0.99939538, 0.99592675, 0.94364995, 0.83725025,
0.99949096, 0.99619551, 0.94276764, 0.83757524, 0.99942112,
0.9960447, 0.94295993, 0.83635325])
```

```
0.5500417 , 0.5429999, 0.5500925],
```

In [224]:

```
cv_auc
```

Out[224]:

```
array([0.56189007, 0.56189007, 0.56189007, 0.56189007, 0.59669849,
       0.59666333, 0.59661517, 0.59657631, 0.55781516, 0.56017799,
       0.57626127, 0.58390923, 0.49458698, 0.49984367, 0.53367696,
       0.5748384 , 0.52427531, 0.52552411, 0.55331475, 0.57280215,
       0.52522253, 0.531215 , 0.54791693, 0.57236022, 0.53079113,
       0.53025665, 0.55244107, 0.57358791])
```

Plot for Train Data

In [225]:

```
x1 = [0.56261079, 0.56261079, 0.56261079, 0.56261079, 0.65106466,
      0.65102953, 0.65095451, 0.64944063, 0.75917609, 0.75486241,
      0.73047794, 0.69831248, 0.99278022, 0.98905981, 0.93318717,
      0.81835938, 0.99939538, 0.99592675, 0.94364995, 0.83725025,
      0.99949096, 0.99619551, 0.94276764, 0.83757524, 0.99942112,
      0.9960447 , 0.94295993, 0.83635325]
```

In [226]:

```
y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x1)
```

In [227]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,50,50,50,50,100,100,100,100,500,500,500,500,1000,1000,1000,1000], index = x1)
```

In [228]:

```
trace = go.Scatter3d(
    x=x1, y=y1, z=z1,
    marker=dict(
        size=4,
        color=z1,
        colorscale='Viridis',
    ),
    line=dict(
        color='#1f77b4',
        width=1
    )
)
```

In [229]:

```
data = [trace]
```

In [230]:

```
layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- TRAIN Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        yaxis=dict(
```

```

yaxis=dict(
    gridcolor='rgb(255, 255, 255)',
    zerolinecolor='rgb(255, 255, 255)',
    showbackground=True,
    backgroundcolor='rgb(230, 230,230)'
),
zaxis=dict(
    gridcolor='rgb(255, 255, 255)',
    zerolinecolor='rgb(255, 255, 255)',
    showbackground=True,
    backgroundcolor='rgb(230, 230,230)'
),
camera=dict(
    up=dict(
        x=0,
        y=0,
        z=1
    ),
    eye=dict(
        x=-1.7428,
        y=1.0707,
        z=0.7100,
    )
),
aspectratio = dict( x=1, y=1, z=0.7 ),
aspectmode = 'manual'
),
)

```

In [231]:

```

fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-e', height=900)

```

Out[231]:

Plot for Cross Validation Data

In [232]:

```
x1 = [0.56189007, 0.56189007, 0.56189007, 0.56189007, 0.59669849,  
      0.59666333, 0.59661517, 0.59657631, 0.55781516, 0.56017799,  
      0.57626127, 0.58390923, 0.49458698, 0.49984367, 0.53367696,  
      0.5748384 , 0.52427531, 0.52552411, 0.55331475, 0.57280215,  
      0.52522253, 0.531215 , 0.54791693, 0.57236022, 0.53079113,  
      0.53025665, 0.55244107, 0.57358791]
```

In [233]:

```
y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,1  
00,500], index = x1)
```

In [234]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,50,50,50,50,100,100,100,100,500,500,500,500,1000,1000,10  
00,1000], index = x1)
```

In [235]:

```
trace = go.Scatter3d(  
    x=x1, y=y1, z=z1,  
    marker=dict(  
        size=4,  
        color=z1,  
        colorscale='Viridis',  
    ),  
    line=dict(  
        color='#1f77b4',  
        width=1  
    )  
)
```

In [236]:

```
data = [trace]
```

In [237]:

```
layout = dict(  
    width=800,  
    height=700,  
    autosize=False,  
    title='Hyper Parameter Tuning -- Cross Validation Data',  
    scene=dict(  
        xaxis=dict(  
            gridcolor='rgb(255, 255, 255)',  
            zerolinecolor='rgb(255, 255, 255)',  
            showbackground=True,  
            backgroundcolor='rgb(230, 230,230)'  
        ),  
        yaxis=dict(  
            gridcolor='rgb(255, 255, 255)',  
            zerolinecolor='rgb(255, 255, 255)',  
            showbackground=True,  
            backgroundcolor='rgb(230, 230,230)'  
        ),  
        zaxis=dict(  
            gridcolor='rgb(255, 255, 255)',  
            zerolinecolor='rgb(255, 255, 255)',  
            showbackground=True,  
            backgroundcolor='rgb(230, 230,230)'  
        ),  
    ),  
)
```

```

camera=dict(
    up=dict(
        x=0,
        y=0,
        z=1
    ),
    eye=dict(
        x=-1.7428,
        y=1.0707,
        z=0.7100,
    )
),
aspectratio = dict( x=1, y=1, z=0.7 ),
aspectmode = 'manual'
),
)

```

In [238]:

```

fig = dict(data=data, layout=layout)
py.iplot(fig, filename='Decision-trees-f', height=900)

```

Out[238]:

Observations : 1) We understand from the 2 plots that the Decision Tree with a depth of 500, 1000 performs great on Training Data but performs pretty bad on unseen data (cross validation data) . => Probably a case of Overfitting. 2) Decision trees with depth 1 performs poor on both Train data as well as Cross Validation Data. =>Probably a case of Underfitting. 3) Decision Tree with maximum depth 5 or 10 performs fairly well on both Train as well as Cross Validation Data. 4) 500 as the value for Minimum samples per split is considered.

Train the model using the best hyper parameter value

In [239]:

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

model = DecisionTreeClassifier(max_depth = 5, min_samples_split = 500)

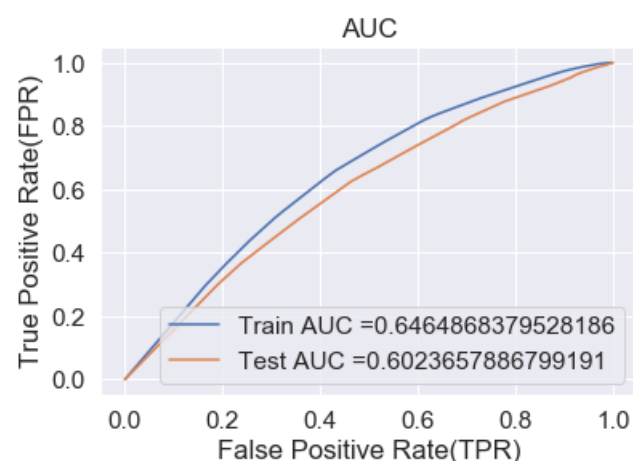
model.fit(X_tr, y_train)

# 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("False Positive Rate(TPR)")
plt.ylabel("True Positive Rate(FPR)")
plt.title("AUC")
plt.show()
```



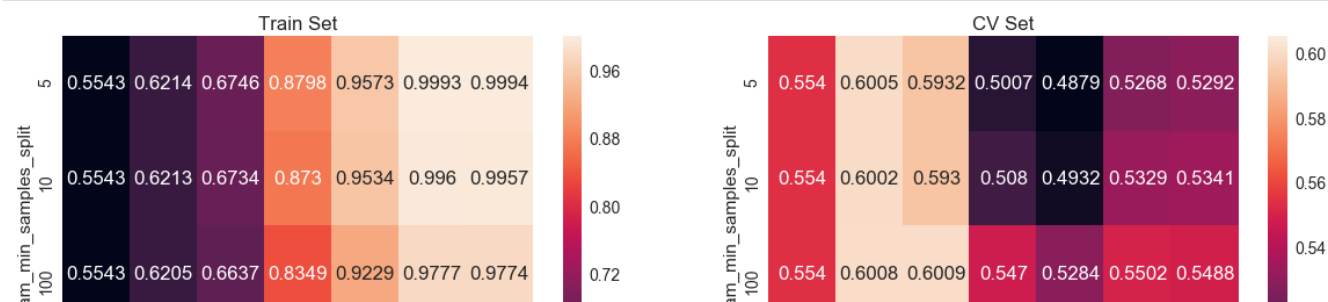
In [216]:

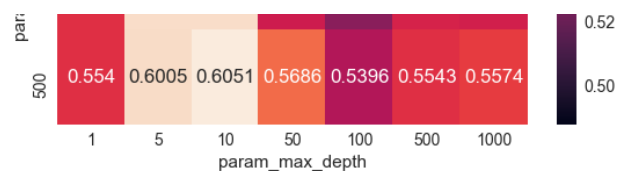
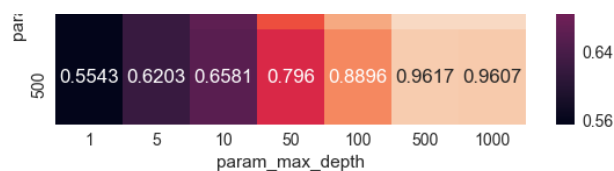
```
import seaborn as sea
max_scores = pd.DataFrame(clf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth'])
.max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))
sea.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sea.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```





Confusion Matrix

Train matrix

In [240]:

```
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) 0.2499976037343404 for threshold 0.833
[[ 2592  2576]
 [ 7937 20395]]
```

In [241]:

```
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))
```

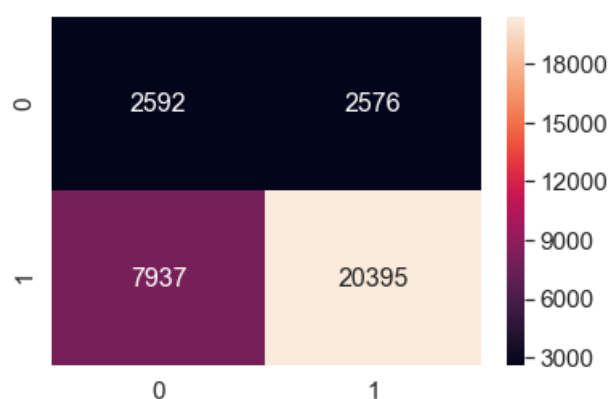
```
the maximum value of tpr*(1-fpr) 0.2499976037343404 for threshold 0.833
```

In [242]:

```
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[242]:

<matplotlib.axes._subplots.AxesSubplot at 0x1e51833f4e0>



Test data

In [243]:

```
print("="*100)
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*(1-fpr) 0.2498875367241191 for threshold 0.861
[[1300 1246]
 [4931 9023]]
```

In [244]:

```
conf_matr_df_test_3 = pd.DataFrame(confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr)), range(2), range(2))
```

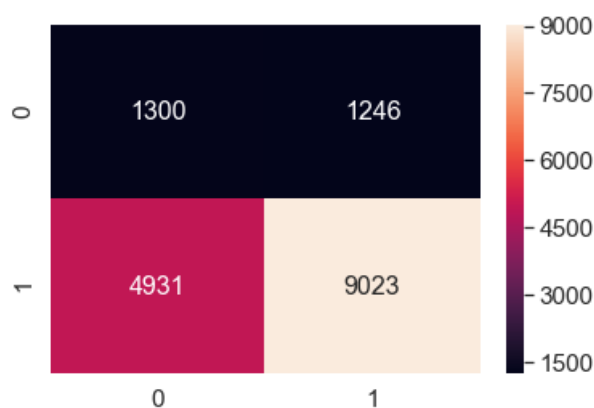
the maximum value of tpr*(1-fpr) 0.2498875367241191 for threshold 0.861

In [245]:

```
sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_df_test_3, annot=True, annot_kws={"size": 16}, fmt='g')
```

Out[245]:

<matplotlib.axes._subplots.AxesSubplot at 0x1e5174cd0b8>



Set 4 : Categorical, Numerical features + Project_title(TFIDF W2V) + Preprocessed_essay (TFIDF W2V)

In [246]:

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

X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train,
school_state_categories_one_hot_train, project_grade_categories_one_hot_train,
teacher_prefix_categories_one_hot_train, price_train, quantity_train, prev_projects_train, title_word_count_train,
essay_word_count_train, essay_sent_pos_train, essay_sent_neg_train, essay_sent_neu_train, essay_sent_comp_train,
tfidf_w2v_vectors_train, tfidf_w2v_vectors_titles_train)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test,
school_state_categories_one_hot_test, project_grade_categories_one_hot_test,
teacher_prefix_categories_one_hot_test, price_test, quantity_test, prev_projects_test, title_word_count_test,
essay_word_count_test, essay_sent_pos_test, essay_sent_neg_test, essay_sent_neu_test, essay_sent_comp_test,
tfidf_w2v_vectors_test, tfidf_w2v_vectors_titles_test)).tocsr()
```

In [247]:

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

```
Final Data matrix
(33500, 708) (33500,)
(16500, 708) (16500,)
```

In [248]:

```
dt_4 = DecisionTreeClassifier()

parameters = {'max_depth':[1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}

clf = GridSearchCV(dt_4, parameters, cv= 3, scoring='roc_auc')

clf.fit(X_tr, y_train)

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

In [250]:

```
train_auc
```

Out[250]:

```
array([0.56134309, 0.56134309, 0.56134309, 0.56134309, 0.66012691,
        0.66012691, 0.65997374, 0.65868173, 0.78427688, 0.7816141 ,
        0.75295522, 0.71396313, 0.99891471, 0.99512885, 0.94212355,
        0.81307748, 0.99947637, 0.99620729, 0.94239419, 0.81442112,
        0.99952158, 0.99620628, 0.94216665, 0.81236292, 0.99946323,
        0.9961006 , 0.94270562, 0.81367683])
```

In [252]:

```
cv_auc
```

Out[252]:

```
array([0.55731383, 0.55731383, 0.55731383, 0.55731383, 0.6151335 ,
        0.61557978, 0.61501655, 0.61523553, 0.56921624, 0.56693834,
        0.58887011, 0.60691274, 0.51475315, 0.51617444, 0.55683711,
        0.59144841, 0.53165678, 0.53479783, 0.56194992, 0.59047246,
        0.53071935, 0.53103767, 0.56298964, 0.59255993, 0.53242729,
        0.53397968, 0.56096015, 0.59157814])
```

Plot for Train Data

In [253]:

```
x1 = [0.56134309, 0.56134309, 0.56134309, 0.56134309, 0.66012691,
        0.66012691, 0.65997374, 0.65868173, 0.78427688, 0.7816141 ,
        0.75295522, 0.71396313, 0.99891471, 0.99512885, 0.94212355,
        0.81307748, 0.99947637, 0.99620729, 0.94239419, 0.81442112,
        0.99952158, 0.99620628, 0.94216665, 0.81236292, 0.99946323,
        0.9961006 , 0.94270562, 0.81367683]
```

In [254]:

```
y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x1)
```

In [255]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,50,50,50,50,100,100,100,100,500,500,500,500,1000,1000,1000,1000], index = x1)
```

In [256]:

```
trace = go.Scatter3d(
    x=x1, y=y1, z=z1,
    marker=dict(
        size=4,
```

```

        color=z1,
        colorscale='Viridis',
    ),
    line=dict(
        color='#1f77b4',
        width=1
    )
)

```

In [257]:

```
data = [trace]
```

In [258]:

```

layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- TRAIN Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230,230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        ),
        aspectratio = dict( x=1, y=1, z=0.7 ),
        aspectmode = 'manual'
    ),
)

```

In [259]:

```

fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-g', height=900)

```

Out[259]:

Plot for Cross Validation Data

In [260]:

```
x1 = [0.55731383, 0.55731383, 0.55731383, 0.55731383, 0.6151335 ,
      0.61557978, 0.61501655, 0.61523553, 0.56921624, 0.56693834,
      0.58887011, 0.60691274, 0.51475315, 0.51617444, 0.55683711,
      0.59144841, 0.53165678, 0.53479783, 0.56194992, 0.59047246,
      0.53071935, 0.53103767, 0.56298964, 0.59255993, 0.53242729,
      0.53397968, 0.56096015, 0.59157814]
```

In [261]:

```
y1 = pd.Series([5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500,5,10,100,500], index = x1)
```

In [262]:

```
z1 = pd.Series([1,1,1,1,5,5,5,5,10,10,10,10,50,50,50,50,100,100,100,100,500,500,500,500,1000,1000,1000,1000], index = x1)
```

In [263]:

```
trace = go.Scatter3d(
    x=x1, y=y1, z=z1,
    marker=dict(
        size=4,
        color=z1,
        colorscale='Viridis',
    ),
    line=dict(
        color='#1f77b4',
        width=1
    )
)
```

In [264]:

$$-1 = 1 = \dots \quad \vdash \quad 1 = \dots = \dots = 1$$

```
data = [trace]
```

In [265]:

```
layout = dict(
    width=800,
    height=700,
    autosize=False,
    title='Hyper Parameter Tuning -- Cross Validation Data',
    scene=dict(
        xaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        yaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        zaxis=dict(
            gridcolor='rgb(255, 255, 255)',
            zerolinecolor='rgb(255, 255, 255)',
            showbackground=True,
            backgroundcolor='rgb(230, 230, 230)'
        ),
        camera=dict(
            up=dict(
                x=0,
                y=0,
                z=1
            ),
            eye=dict(
                x=-1.7428,
                y=1.0707,
                z=0.7100,
            )
        ),
        aspectratio = dict( x=1, y=1, z=0.7 ),
        aspectmode = 'manual'
    ),
)
```

In [266]:

```
fig = dict(data=data, layout=layout)

py.iplot(fig, filename='Decision-trees-h ', height=900)
```

Out[266]:

Observations : 1) We understand from the 2 plots that the Decision Tree with a depth of 500, 1000 performs great on Training Data but performs pretty bad on unseen data (cross validation data) . => Probably a case of Overfitting. 2) Decision trees with depth 1 performs poor on both Train data as well as Cross Validation Data. =>Probably a case of Underfitting. 3) Decision Tree with maximum depth 10 performs fairly well on both Train as well as Cross Validation Data. 4) 100 as the value for Minimum samples per split is considered.

Train the model using the best hyper parameter value

In [267]:

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

model = DecisionTreeClassifier(max_depth = 10, min_samples_split = 100)

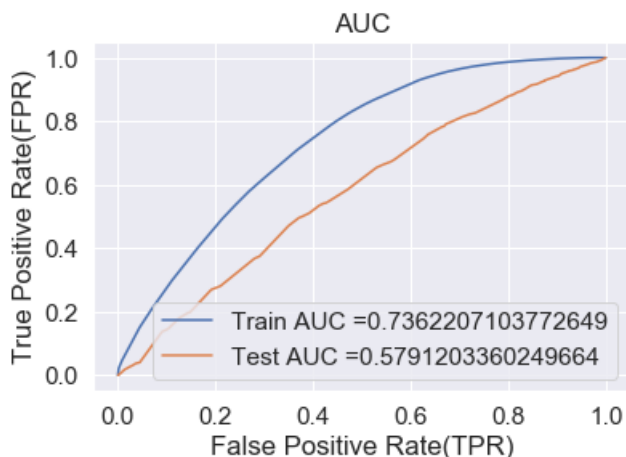
model.fit(X_tr, y_train)

# 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("False Positive Rate(TPR)")
plt.ylabel("True Positive Rate(FPR)")
plt.title("AUC")
plt.show()
```



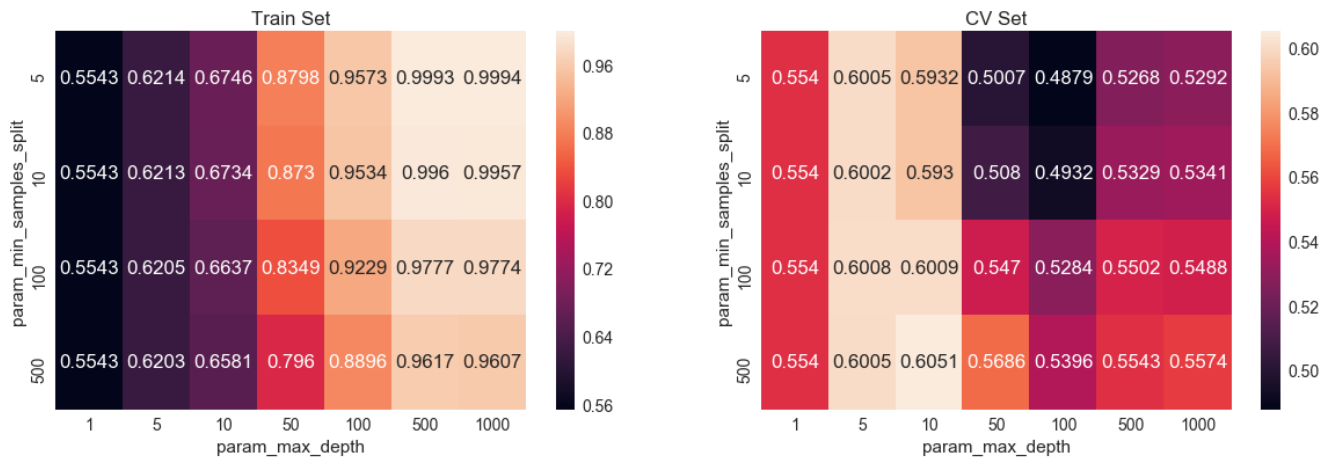
In [217]:

```
import seaborn as sea
max_scores = pd.DataFrame(clf.cv_results_).groupby(['param_min_samples_split', 'param_max_depth'])
.max().unstack()[['mean_test_score', 'mean_train_score']]

fig, ax = plt.subplots(1,2, figsize=(20,6))
sea.heatmap(max_scores.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sea.heatmap(max_scores.mean_test_score, annot = True, fmt='.4g', ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('CV Set')

plt.show()
```



Confusion Matrix

Train Data

In [268]:

```
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) 0.24998348823193933 for threshold 0.823
[[ 2563  2605]
 [ 4203 24129]]
```

In [269]:

```
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))
```

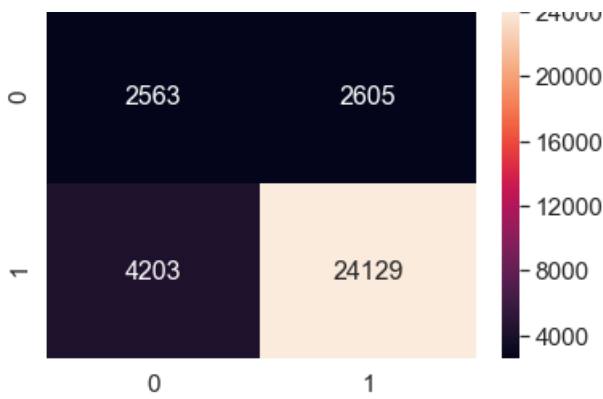
the maximum value of tpr*(1-fpr) 0.24998348823193933 for threshold 0.823

In [270]:

```
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[270]:

<matplotlib.axes._subplots.AxesSubplot at 0x1e5183ad080>



Test data

In [271]:

```
print("="*100)
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*(1-fpr) 0.24917789191060447 for threshold 0.881
[[1346 1200]
 [5787 8167]]
```

In [272]:

```
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))
```

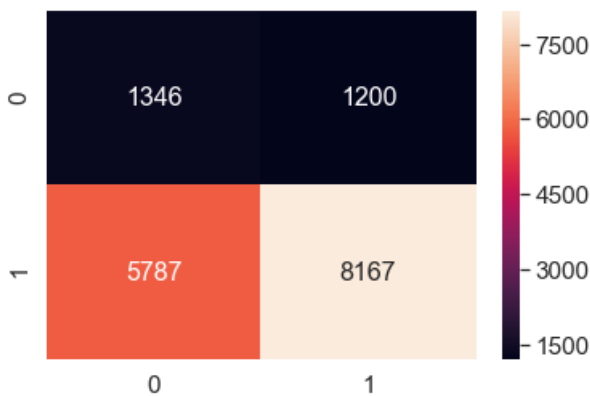
the maximum value of tpr*(1-fpr) 0.24917789191060447 for threshold 0.881

In [273]:

```
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[273]:

<matplotlib.axes._subplots.AxesSubplot at 0x1e51812dcf8>



2.5 [Task-2]Getting top 5k features using `feature_importances_`

In [274]:

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



```

X_tr = hstack((categories_one_hot_train, sub_categories_one_hot_train,
school_state_categories_one_hot_train, project_grade_categories_one_hot_train,
teacher_prefix_categories_one_hot_train, price_train, quantity_train, prev_projects_train, title_wor
rd_count_train, essay_word_count_train, essay_sent_pos_train, essay_sent_neg_train,
essay_sent_neu_train, essay_sent_comp_train, title_tfidf_train, text_tfidf_train)).tocsr()
X_te = hstack((categories_one_hot_test, sub_categories_one_hot_test,
school_state_categories_one_hot_test, project_grade_categories_one_hot_test,
teacher_prefix_categories_one_hot_test, price_test, quantity_test, prev_projects_test,
title_word_count_test, essay_word_count_test, essay_sent_pos_test, essay_sent_neg_test, essay_sent_
neu_test, essay_sent_comp_test, title_tfidf_test, text_tfidf_test)).tocsr()

```

In [275]:

```

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

```

```

Final Data matrix
(33500, 12207) (33500,)
(16500, 12207) (16500,)
=====

```



In [276]:

```

## Fit the Model to obtain the best 5k features

model = DecisionTreeClassifier()

model.fit(X_tr, y_train)

```

Out[276]:

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')

```

In [277]:

```

## Compute the Feature importances for our Train Features

a=model.tree_.compute_feature_importances(normalize=False)

```

In [278]:

```
df9 = pd.DataFrame(a)
```

In [279]:

```
df9 = np.transpose(df9)
```

In [281]:

```

## Store the indexes of the features with atleast some importance. Lets ignore the features with 0
## as the feature importance value and instead consider all the values other than these

best_ind = []

for j in range(12207):

    s = df9[j].sum()

    if s > 0 :
        best_ind.append(j)
    else :

```

continue

In [282]:

```
a[0:2]
```

Out[282]:

```
array([0., 0.])
```

In [283]:

```
## Identify number of Features after feature importance step  
len(best_ind)
```

Out[283]:

```
2055
```

Inference While calculating the feature importances of every Column using the Decision Trees classifier, we hardly get 2755 columns with some importance. The remaining columns of the the total 14355 columns contribute 0 importance.

In [284]:

```
best_ind[0:10]
```

Out[284]:

```
[3, 6, 7, 8, 15, 22, 23, 31, 33, 34]
```

In [285]:

```
## Convert the sparse matrix to a dense matrix to separate the important features and its corresponding values.
```

```
a = X_tr.todense()
```

In [286]:

```
a.shape
```

Out[286]:

```
(33500, 12207)
```

In [287]:

```
df_set_5_x_train = pd.DataFrame(a)
```

In [288]:

```
final_df = df_set_5_x_train.iloc[:, best_ind]
```

In [289]:

```
final_df.shape
```

Out[289]:

```
(33500, 2055)
```

In [290]:

```
y_train.shape
```

Out[290]:

(33500,)

In [291]:

```
b = X_te.todense()
```

In [293]:

```
b.shape
```

Out[293]:

(16500, 12207)

In [294]:

```
df_set_5_x_test = pd.DataFrame(b)
```

In [295]:

```
final_df_test = df_set_5_x_test.iloc[:, best_ind]
```

In [296]:

```
final_df_test.shape
```

Out[296]:

(16500, 2055)

GridSearchCV - L2 regularization

In [297]:

```
from sklearn.linear_model import SGDClassifier

sv = SGDClassifier(loss='hinge', penalty='l2')

parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(sv, parameters, cv= 3, scoring='roc_auc')

clf.fit(final_df, 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']
```

In [298]:

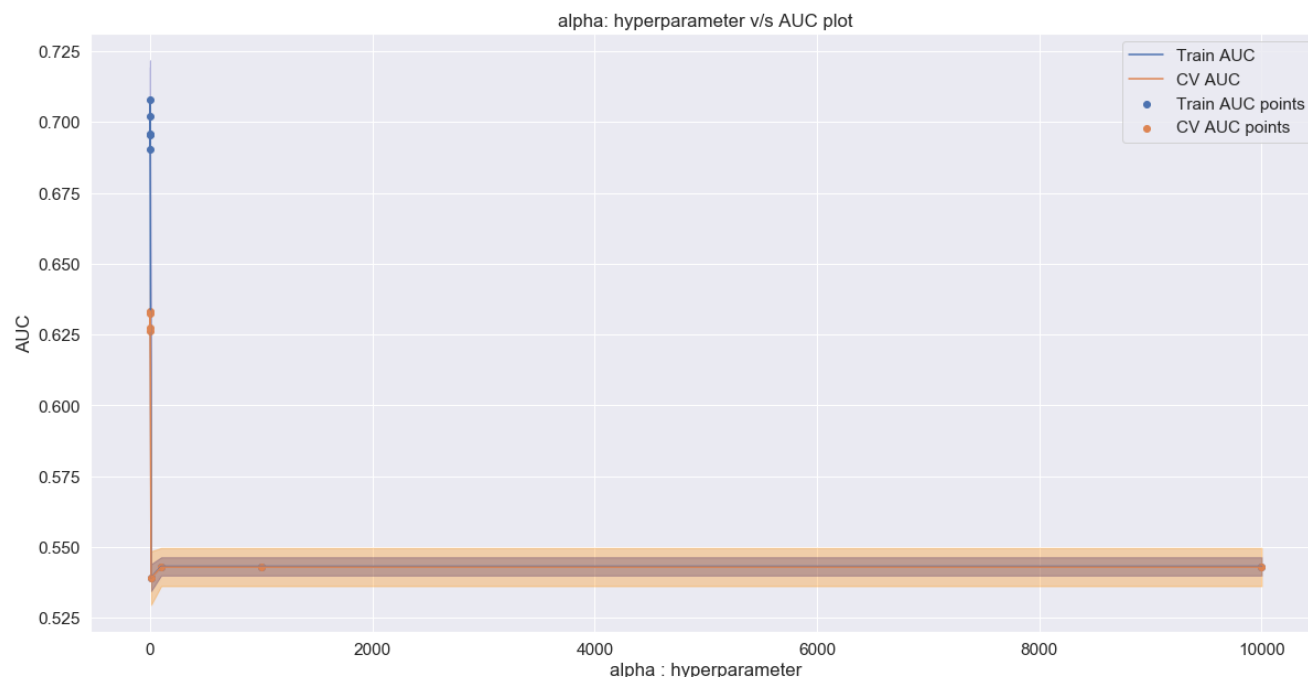
```
plt.figure(figsize=(20,10))

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

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

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

```
plt.legend()
plt.xlabel("alpha : hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC plot")
plt.show()
```



In [299]:

```
from sklearn.linear_model import SGDClassifier

sv = SGDClassifier(loss='hinge', penalty='l2')

parameters = {'alpha':[0.01, 0.05, 0.1, 0.5, 0.9, 1.5, 2.0]}

clf = GridSearchCV(sv, parameters, cv= 3, scoring='roc_auc')

clf.fit(final_df, 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']
```

In [300]:

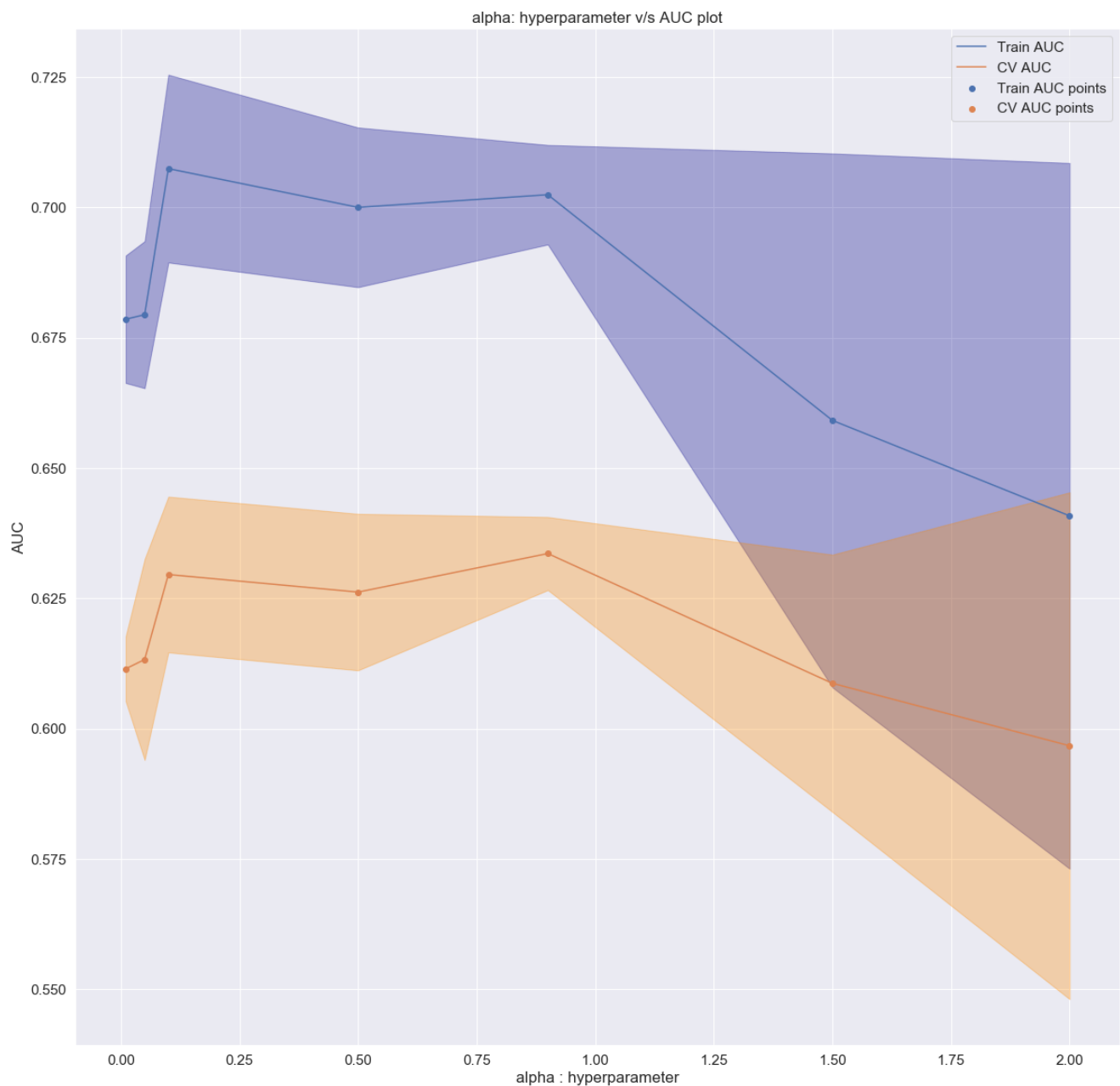
```
plt.figure(figsize=(20,20))

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

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

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

plt.legend()
plt.xlabel("alpha : hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC plot")
plt.show()
```



Inference For L2 regularization i have considered 1.5 to be the best alpha value. But, there is a huge difference in the performance of the model between the Cross Validation and Train datasets.

GridSearchCV - L1 regularization

In [301]:

```
sv = SGDClassifier(loss='hinge', penalty='l1')

parameters = {'alpha':[10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2, 10**3, 10**4]}

clf = GridSearchCV(sv, parameters, cv= 3, scoring='roc_auc')

clf.fit(final_df, 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']
```

In [302]:

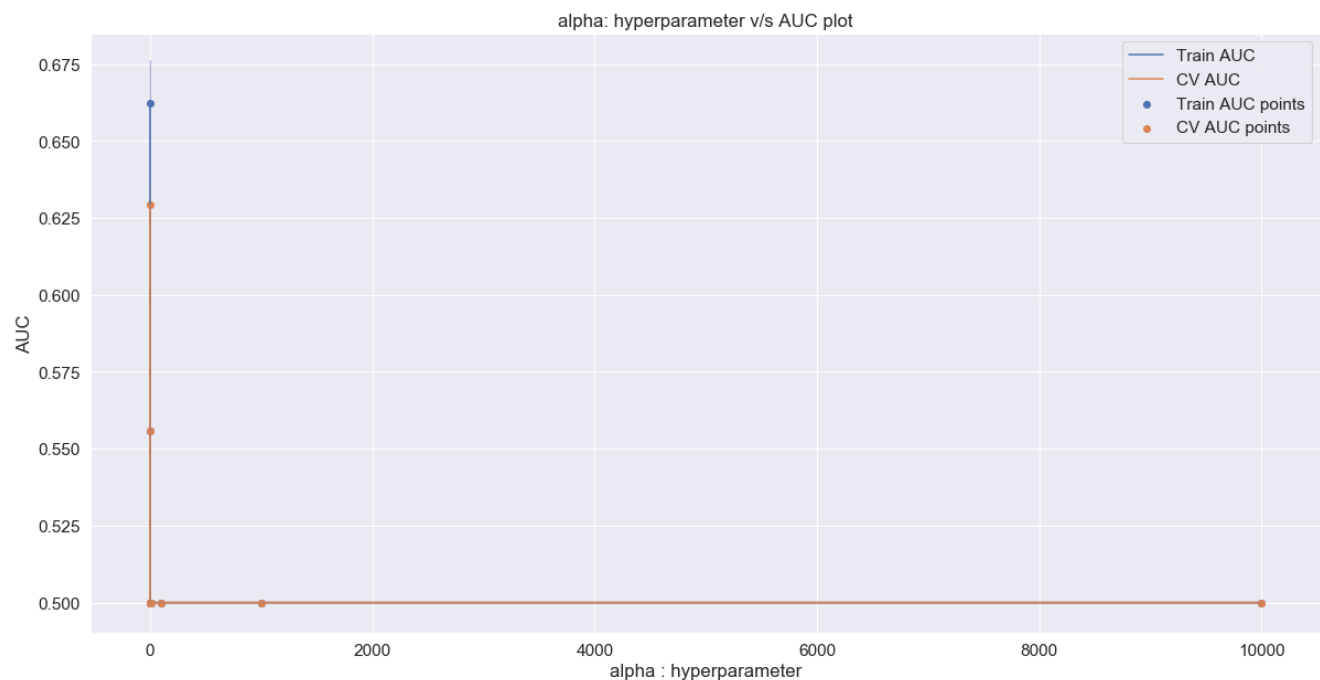
```
plt.figure(figsize=(20,10))
```

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

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

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

plt.legend()
plt.xlabel("alpha : hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC plot")
plt.show()
```



In [303]:

```
sv = SGDClassifier(loss='hinge', penalty='l1')

parameters = {'alpha':[0.00003, 0.00006, 0.00008, 0.0001, 0.0004]}

clf = GridSearchCV(sv, parameters, cv= 3, scoring='roc_auc')

clf.fit(final_df, 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']
```

In [304]:

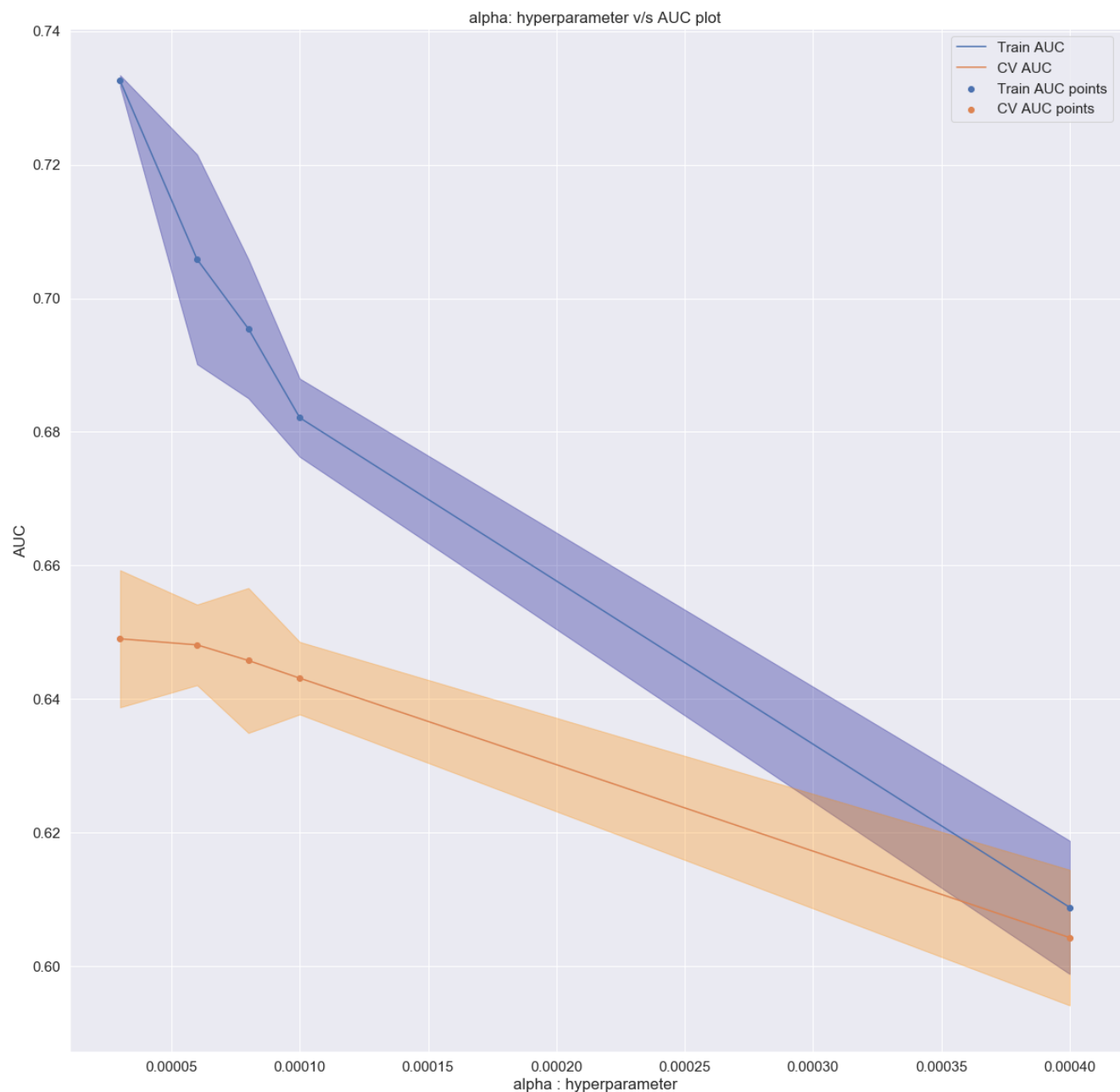
```
plt.figure(figsize=(20,20))

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

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

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

plt.legend()
plt.xlabel("alpha : hyperparameter")
plt.ylabel("AUC")
plt.title("alpha: hyperparameter v/s AUC plot")
plt.show()
```



Inference 0.00009 is considered as the best hyperparameter value for Alpha as it yields a better result both on training and test data. L1 regularization performs better and yields a better result compared to the L2 regularization, this is based on the AUC score.

Train the model using the best hyper parameter value

In [305]:

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

model = SGDClassifier(loss='hinge', penalty='l1', alpha=0.00009)

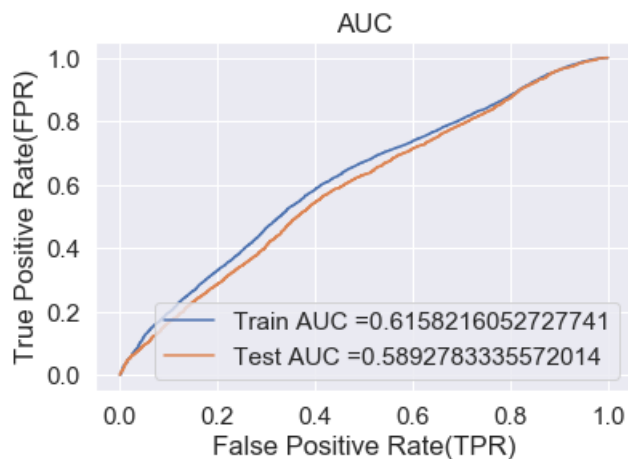
model.fit(final_df, y_train)
```

```
# 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 = model.decision_function(final_df)
y_test_pred = model.decision_function(final_df_test)

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

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



Confusion Matrix

Train Data

In [306]:

```
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_tpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.25 for threshold 1.04
[[ 2584  2584]
 [ 9281 19051]]
```

In [307]:

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

```
the maximum value of tpr*(1-fpr) 0.25 for threshold 1.04
```

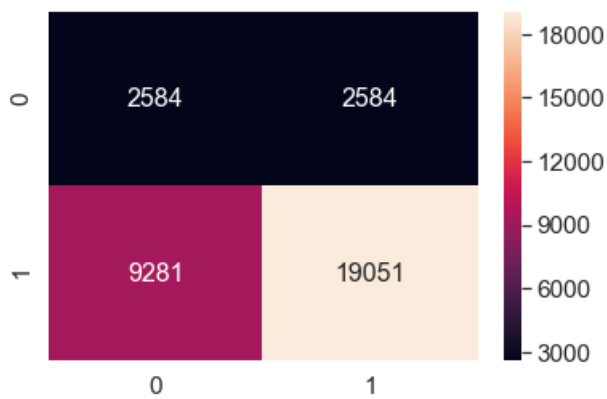
In [308]:

```
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[308]:

```
[[ 2584  2584]
 [ 9281 19051]]
```


<matplotlib.axes._subplots.AxesSubplot at 0x1e518390a58>



Test Data

In [309]:

```
print("="*100)
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*(1-fpr) 0.25 for threshold 1.104
[[1846  700]
 [8741 5213]]
```

In [310]:

```
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))
```

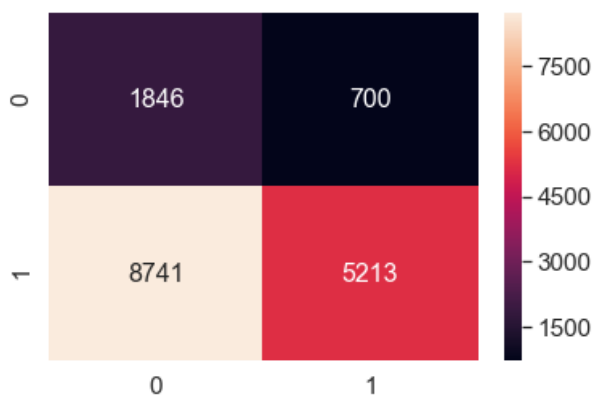
the maximum value of tpr*(1-fpr) 0.25 for threshold 1.104

In [311]:

```
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[311]:

<matplotlib.axes._subplots.AxesSubplot at 0x1e51a370940>



Conclusion

In [312]:

```
# 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", "Hyperparameters(max depth,min samples split)", "Train AUC",  
                , "Test AUC"]
```

```
x.add_row(["BOW", "Decision Trees", "(10, 100)", 0.665, 0.615])
```

```
x.add_row(["TFIDF", "Decision Trees", "(10, 500)", 0.653, 0.615])
```

```
x.add_row(["AVG W2V", "Decision Trees", "(5, 500)", 0.637, 0.602])
```

```
x.add_row(["TFIDF W2V", "Decision Trees", "(10, 100)", 0.586, 0.555])
```

```
x.add_row(["TFIDF-5k Features", "Linear SVM", "[L1 penalty, alpha = 0.00009]", 0.648, 0.635])
```

```
print(x)
```

```
+-----+-----+-----+-----+-----+
-----+
| Vectorizer | Model | Hyperparameters(max depth,min samples split) | Train AUC | Test AUC |
+-----+-----+-----+-----+-----+
-----+
| BOW | Decision Trees | (10, 100) | 0.665 | 0.615 |
| TFIDF | Decision Trees | (10, 500) | 0.653 | 0.615 |
| AVG W2V | Decision Trees | (5, 500) | 0.637 | 0.602 |
| TFIDF W2V | Decision Trees | (10, 100) | 0.586 | 0.555 |
| TFIDF-5k Features | Linear SVM | [L1 penalty, alpha = 0.00009] | 0.648 | 0.635 |
+-----+-----+-----+-----+-----+
-----+
```

