

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*

Feature	Description
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

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

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

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

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

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

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

Notes on the Essay Data

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

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

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

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

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

In [1]:

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

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

```

import 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
import sys
from sklearn import tree

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

```

1.1 Reading Data

In [2]:

```

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

```

In [3]:

```

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

```

Number of data points in train data (109248, 17)

```

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

```

In [4]:

```

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

```

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

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

Out[4]:

```

id          description  quantity  price

```

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

In [5]:

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

1.2 preprocessing of project_subject_categories

In [6]:

```
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 [7]:

```
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('&', '_') # we are replacing the & value into
    sub_cat_list.append(temp.strip())
```

```

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]))

```

1.3 Text preprocessing

In [8]:

```

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

```

In [9]:

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

Out[9]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_c
0	160221 p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades
1	140945 p258326	897464ce9ddc600bcded1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Gra

In [10]:

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

In [11]:

```

# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)

```

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

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

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

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

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving docto

rs, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.

The cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.

nannan

=====

In [12]:

```
# 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 [13]:

```
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations.

The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills.

They also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget that they are doing work and just have the fun a 6 year old deserves.

nannan

=====

In [14]:

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

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations.

The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. The want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills.

They also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget the

our success. The number, color and shape makes can make that happen. My students will forget are
y are doing work and just have the fun a 6 year old deserves,nannan

In [15]:

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

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

In [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', 't
            heir', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these',
            'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'd
            o', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'whil
            e', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'bef
            ore', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'a
            gain', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each
            ', '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', "doesn
            '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]:

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


In [25]:

```
# Updating dataframe for clean project title and remove old project title
project_data.drop(['project_grade_category'], axis=1, inplace=True)
project_data['project_grade_category'] = preprocessed_grade
project_data.head(2)
```

Out[25]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_essay_1
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	My students are English learners that are work...
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Our students arrive to our school eager to lea...

In [26]:

```
# remove unnecessary column: https://cmdlinetips.com/2018/04/how-to-drop-one-or-more-columns-in-pandas-dataframe/
project_data = project_data.drop(['Unnamed: 0', 'id', 'teacher_id', 'project_submitted_datetime', \
                                   'project_essay_1', 'project_essay_2', 'project_essay_3', 'project_essay_4', \
                                   'project_resource_summary'], axis=1)
```

In [27]:

```
project_data.head()
```

Out[27]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_category
0	Mrs.	IN	0	0	154.60	23	Literacy_Language
1	Mr.	FL	7	1	299.00	1	History_Civics_Health_Sports
2	Ms.	AZ	1	0	516.85	22	Health_Sports
3	Mrs.	KY	4	1	232.90	4	Literacy_Language_Math_Science
4	Mrs.	TX	1	1	67.98	4	Math_Science

Check whether each column contain NaN or Not

In [28]:

```
project_data['teacher_prefix'].isnull().values.any()
```

Out[28]:

True

In [29]:

```
project_data['school_state'].isnull().values.any()
```

Out[29]:

False

In [30]:

```
project_data['teacher_number_of_previously_posted_projects'].isnull().values.any()
```

Out[30]:

False

In [31]:

```
project_data['project_is_approved'].isnull().values.any()
```

Out[31]:

False

In [32]:

```
project_data['price'].isnull().values.any()
```

Out[32]:

False

In [33]:

```
project_data['quantity'].isnull().values.any()
```

Out[33]:

False

In [34]:

```
project_data['clean_categories'].isnull().values.any()
```

Out[34]:

False

In [35]:

```
project_data['clean_subcategories'].isnull().values.any()
```

Out[35]:

False

In [36]:

```
project_data['clean_essay'].isnull().values.any()
```

Out[36]:

False

In [37]:

```
project_data['clean_project_title'].isnull().values.any()
```

Out[37]:

False

In [38]:

```
project_data['project_grade_category'].isnull().values.any()
```

Out[38]:

False

Since we got 'teacher prefix' attributes which contain NaN. Let check how many NaN are contain in this attributes

In [39]:

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

Out[39]:

3

1.5 Preparing data for models

In [40]:

```
project_data.columns
```

Out[40]:

```
Index(['teacher_prefix', 'school_state',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'price', 'quantity', 'clean_categories', 'clean_subcategories',  
      'clean_essay', 'clean_project_title', 'project_grade_category'],  
      dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

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

1.5.1 Vectorizing Categorical data

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

In [0]:

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

```
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
Shape of matrix after one hot encodig  (109248, 9)
```

In [0]:

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

```
['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Extracurricular', 'Civics_Government', 'ForeignLanguages', 'NutritionEducation', 'Warmth', 'Care_Hunger', 'SocialSciences', 'PerformingArts', 'CharacterEducation', 'TeamSports', 'Other', 'College_CareerPrep', 'Music', 'History_Geography', 'Health_LifeScience', 'EarlyDevelopment', 'ESL', 'Gym_Fitness', 'EnvironmentalScience', 'VisualArts', 'Health_Wellness', 'AppliedSciences', 'SpecialNeeds', 'Literature_Writing', 'Mathematics', 'Literacy']
Shape of matrix after one hot encodig  (109248, 30)
```

In [0]:

```
# you can do the similar thing with state, teacher_prefix and project_grade_category also
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [0]:

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

```
Shape of matrix after one hot encodig  (109248, 16623)
```

In [0]:

```
# you can vectorize the title also
# before you vectorize the title make sure you preprocess it
```

1.5.2.2 TFIDF vectorizer

In [0]:

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

Shape of matrix after one hot encoding (109248, 16623)

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

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

# =====
Output:

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

# =====

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

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

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

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

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

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

'''
```

Out[0]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef loadGloveModel (
```

```

gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\r', encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n        word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel = loadGloveModel('glove.42B.300d.txt')\n\n\n# =====\nOutput:\n    Loading Glove Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n# =====\nnwords = []\nfor i in preprocod_titles:\n    words.extend(i.split(' '))\n\nfor i in preprocod_titles:\n    words.extend(i.split(' '))\n\nprint("all the words in the coupus", len(words))\nnwords = set(words)\nprint("the unique words in the coupus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\nprint("The number of words that are present in both glove vectors and our coupus", len(inter_words), "(", np.round(len(inter_words)/len(words)*100,3), "%)")\n\nnwords_courpus = {}\nnwords_glove = set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n\nprint("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pickle\nwith open('glove_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n\n'
```

In [0]:

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

In [0]:

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

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

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248 [00:27<00:00, 3953.36it/s]
```

109248
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

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

In [0]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
```



```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248 [03:22<00  
:00, 539.44it/s]
```

In [0]:

1.5.3 Vectorizing Numerical features

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

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

# price_standardized = StandardScaler.fit(project_data['price'].values)
# this will rise the error
# ValueError: Expected 2D array, got 1D array instead: array=[725.05 213.03 329.    ... 399.    287.73
5.5 ].
# Reshape your data either using array.reshape(-1, 1)

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

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

price standardized

```
array([[0.00098843, 0.00191166, 0.00330448, ..., 0.00153418, 0.00046704,
        0.00070265]])
```

1.5.4 Merging all the above features

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

In [0]:

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

```
(109248, 9)
(109248, 30)
(109248, 16623)
(109248, 1)
```

In [0]:

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

Out[0]:

```
(109248, 16663)
```

In [0]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

Computing Sentiment Scores

In [0]:

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

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

```
sid = SentimentIntensityAnalyzer()
```

```
for sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students with t
he biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelligence
s i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of different bac
kgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a caring com
munity of successful \
learners which can be seen through collaborative student project based learning in and out of the class
room kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice a sk
ill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the kinderg
arten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role play in
our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food i wil
l take their idea \
and create common core cooking lessons where we learn important math and writing concepts while cooking
delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into making th
```

```
e food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project would exp
and our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce
make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create our own cookbooks
to be printed and \
shared with families students will gain math and literature skills as well as a life long enjoyment for
healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

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

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

D:\installed\Anaconda3\lib\site-packages\nltk\twitter__init__.py:20: UserWarning:

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

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

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](#)

2. Decision Tree

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

In [41]:

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

In [41]:

```
# Combine the train.csv and resource.csv
# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
from sklearn.model_selection import train_test_split

# https://www.geeksforgeeks.org/python-pandas-dataframe-sample/
# Take 50k dataset
project_data = project_data.sample(n=50000)
# Remove that row which contain NaN. We observed that only 3 rows that contain NaN
project_data = project_data[pd.notnull(project_data['teacher_prefix'])]
project_data.shape
```

Out[41]:

(49998, 11)

In [42]:

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

Out[42]:

	teacher_prefix	school_state	teacher_number_of_previously_posted_projects	project_is_approved	price	quantity	clean_category
51811	Ms.	IL	0	0	430.48	5	Math_Sci
49397	Ms.	MN	0	0	148.72	17	AppliedLea

In [127]:

```
# Split the data into train and test
```

```
# Split train and test
tr_X, ts_X, tr_y, ts_y, = train_test_split(project_data, project_data['project_is_approved'].values, te
st_size=0.33, random_state=1, stratify=project_data['project_is_approved'].values)
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)

# After train data, We are going to perform KFold Cross validation at the time of training model

# Reset index of df
tr_X = tr_X.reset_index(drop=True)
ts_X = ts_X.reset_index(drop=True)
tr_X.drop(['project_is_approved'], axis=1, inplace=True)
ts_X.drop(['project_is_approved'], axis=1, inplace=True)

print('Shape of train data:', tr_X.shape)
print('Shape of test data:', ts_X.shape)
```

Shape of train data: (33498, 10)
Shape of test data: (16500, 10)

In [128]:

```
print('Shape of Train Data',[tr_X.shape, tr_y.shape])
print('Shape of Test Data',[ts_X.shape, ts_y.shape])
```

Shape of Train Data [(33498, 10), (33498,)]
Shape of Test Data [(16500, 10), (16500,)]

2.2 Make Data Model Ready: encoding numerical, categorical features

In [129]:

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

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

# # For Numerical with train data
# ### 1) quantity

from sklearn.preprocessing import Normalizer
# # normalization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.Normalizer.html

quantity_scaler = Normalizer()
quantity_scaler.fit(tr_X['quantity'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
quantity_normalized = quantity_scaler.transform(tr_X['quantity'].values.reshape(1, -1))

# ### 2) price

# # the cost feature is already in numerical values, we are going to represent the money, as numerical values within the range 0-1

price_scaler = Normalizer()
price_scaler.fit(tr_X['price'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
price_normalized = price_scaler.transform(tr_X['price'].values.reshape(1, -1))

# ### 3) For teacher_number_of_previously_projects

# # We are going to represent the teacher_number_of_previously_posted_projects, as numerical values within the range 0-1
```

```

teacher_number_of_previously_posted_projects_scalar = Normalizer()
teacher_number_of_previously_posted_projects_scalar.fit(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1)) # finding the mean and standard deviation of this data
teacher_number_of_previously_posted_projects_normalized = teacher_number_of_previously_posted_projects_scalar.transform(tr_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))

```

In [130]:

```

print('Shape of quantity:', quantity_normalized.T.shape)
print('Shape of price:', price_normalized.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', teacher_number_of_previously_posted_projects_normalized.T.shape)

```

Shape of quantity: (33498, 1)

Shape of price: (33498, 1)

Shape of teacher_number_of_previously_posted_projects: (33498, 1)

In [131]:

```

# # Transform numerical attributes for test data
ts_price = price_scalar.transform(ts_X['price'].values.reshape(1,-1))
ts_quantity = quantity_scalar.transform(ts_X['quantity'].values.reshape(1,-1))
ts_teacher_number_of_previously_posted_projects = \
teacher_number_of_previously_posted_projects_scalar.transform(ts_X['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))

```

In [132]:

```

print('-----Test data-----')
print('Shape of quantity:', ts_quantity.T.shape)
print('Shape of price:', ts_price.T.shape)
print('Shape of teacher_number_of_previously_posted_projects:', ts_teacher_number_of_previously_posted_projects.T.shape)

```

-----Test data-----

Shape of quantity: (16500, 1)

Shape of price: (16500, 1)

Shape of teacher_number_of_previously_posted_projects: (16500, 1)

In [133]:

```

# For categorical with train data
# Please do the similar feature encoding with state, teacher_prefix and project_grade_category also
# One hot encoding for school state

### 1) school_state
print('=====\n')
# Count Vectorize with vocabulary contains unique code of school state and we are doing boolean BoW
vectorizer_school_state = CountVectorizer(vocabulary=tr_X['school_state'].unique(), lowercase=False, binary=True)
vectorizer_school_state.fit(tr_X['school_state'].values)
print('List of feature in school_state',vectorizer_school_state.get_feature_names())

# Transform train data
school_state_one_hot = vectorizer_school_state.transform(tr_X['school_state'].values)
print("\nShape of school_state matrix after one hot encoding ",school_state_one_hot.shape)

### 2) project_subject_categories
print('=====\n')
vectorizer_categories = CountVectorizer(vocabulary=list(sorted_cat_dict.keys()), lowercase=False, binary=True)
vectorizer_categories.fit(tr_X['clean_categories'].values)
print('List of features in project_subject_categories',vectorizer_categories.get_feature_names())

# Transform train data
categories_one_hot = vectorizer_categories.transform(tr_X['clean_categories'].values)
print("\nShape of project_subject_categories matrix after one hot encoding ",categories_one_hot.shape)

```

```

### 3) project_subject_subcategories
print('=====\\n')
vectorizer_subcategories = CountVectorizer(vocabulary=list(sorted_sub_cat_dict.keys()), lowercase=False
, binary=True)
vectorizer_subcategories.fit(tr_X['clean_subcategories'].values)
print('List of features in project_subject_categories',vectorizer_subcategories.get_feature_names())

# Transform train data
subcategories_one_hot = vectorizer_subcategories.transform(tr_X['clean_subcategories'].values)
print("\\nShape of project_subject_subcategories matrix after one hot encodig ",subcategories_one_hot.sh
ape)

### 4) project_grade_category
print('=====\\n')
# One hot encoding for project_grade_category

# Count Vectorize with vocuabulary contains unique code of project_grade_category and we are doing bool
en BoW
vectorizer_grade_category = CountVectorizer(vocabulary=tr_X['project_grade_category'].unique(), lowerca
se=False, binary=True)
vectorizer_grade_category.fit(tr_X['project_grade_category'].values)
print('List of features in project_grade_category',vectorizer_grade_category.get_feature_names())

# Transform train data
project_grade_category_one_hot = vectorizer_grade_category.transform(tr_X['project_grade_category'].val
ues)
print("\\nShape of project_grade_category matrix after one hot encodig ",project_grade_category_one_hot.
shape)

### 5) teacher_prefix
print('=====\\n')
# One hot encoding for teacher_prefix

# Count Vectorize with vocuabulary contains unique code of teacher_prefix and we are doing boolean BoW
# Since some of the data is filled with nan. So we update the nan to 'None' as a string
# tr_X['teacher_prefix'] = tr_X['teacher_prefix'].fillna('None')
vectorizer_teacher_prefix = CountVectorizer(vocabulary=tr_X['teacher_prefix'].unique(), lowercase=False
, binary=True)
vectorizer_teacher_prefix.fit(tr_X['teacher_prefix'].values)
print('List of features in teacher_prefix',vectorizer_teacher_prefix.get_feature_names())

# Transform train data
teacher_prefix_one_hot = vectorizer_teacher_prefix.transform(tr_X['teacher_prefix'].values)
print("\\nShape of teacher_prefix matrix after one hot encoding ",teacher_prefix_one_hot.shape)

```

```

=====

List of feature in school_state ['SC', 'CA', 'PA', 'IL', 'FL', 'MI', 'MO', 'OR', 'NY', 'KY', 'OH', 'NC'
, 'GA', 'VA', 'SD', 'IN', 'NJ', 'MD', 'CT', 'AZ', 'ME', 'TX', 'KS', 'WI', 'MT', 'TN', 'LA', 'NV', 'MA',
'HI', 'WV', 'OK', 'WA', 'CO', 'ID', 'IA', 'MS', 'UT', 'AL', 'AR', 'NH', 'DC', 'MN', 'NE', 'NM', 'DE', '
VT', 'RI', 'AK', 'WY', 'ND']

```

```

Shape of school_state matrix after one hot encoding (33498, 51)
=====

```

```

List of features in project_subject_categories ['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts'
, 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']

```

```

Shape of project_subject_categories matrix after one hot encodig (33498, 9)
=====

```

```

List of features in project_subject_categories ['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', 'SpecialNeed
s', 'Literature_Writing', 'Mathematics', 'Literacy']

```

```

Shape of project_subject_subcategories matrix after one hot encodig (33498, 30)
=====

```

```

List of features in project_grade_category ['grades_prek_2', 'grades_9_12', 'grades_6_8', 'grades_3_5']

```

```

Shape of project_grade_category matrix after one hot encodig (33498, 4)
=====

```

List of features in teacher_prefix ['Mrs.', 'Ms.', 'Mr.', 'Teacher', 'Dr.']

Shape of teacher_prefix matrix after one hot encoding (33498, 5)

In [134]:

```
# Transform categorical for test data
ts_school_state = vectorizer_school_state.transform(ts_X['school_state'].values)
ts_project_subject_category = vectorizer_categories.transform(ts_X['clean_categories'].values)
ts_project_subject_subcategory = vectorizer_subcategories.transform(ts_X['clean_subcategories'].values)
ts_project_grade_category = vectorizer_grade_category.transform(ts_X['project_grade_category'].values)
ts_teacher_prefix = vectorizer_teacher_prefix.transform(ts_X['teacher_prefix'].values)
```

In [135]:

```
print('-----Test data-----')
print('Shape of school_state:', ts_school_state.shape)
print('Shape of project_subject_categories:', ts_project_subject_category.shape)
print('Shape of project_subject_subcategories:', ts_project_subject_subcategory.shape)
print('Shape of project_grade_category:', ts_project_grade_category.shape)
print('Shape of teacher_prefix:', ts_teacher_prefix.shape)
```

```
-----Test data-----
Shape of school_state: (16500, 51)
Shape of project_subject_categories: (16500, 9)
Shape of project_subject_subcategories: (16500, 30)
Shape of project_grade_category: (16500, 4)
Shape of teacher_prefix: (16500, 5)
```

In [136]:

```
list_features = ['quantity', 'price', 'teacher_number_of_previously_posted_projects'] # storing all features names for further print feature importance
for i in range(len(vectorizer_school_state.get_feature_names())):
    list_features.append(vectorizer_school_state.get_feature_names()[i])
for i in range(len(vectorizer_categories.get_feature_names())):
    list_features.append(vectorizer_categories.get_feature_names()[i])
for i in range(len(vectorizer_subcategories.get_feature_names())):
    list_features.append(vectorizer_subcategories.get_feature_names()[i])
for i in range(len(vectorizer_grade_category.get_feature_names())):
    list_features.append(vectorizer_grade_category.get_feature_names()[i])
for i in range(len(vectorizer_teacher_prefix.get_feature_names())):
    list_features.append(vectorizer_teacher_prefix.get_feature_names()[i])

len(list_features)
```

Out[136]:

102

2.3 Make Data Model Ready: encoding eassay, and project_title

In [137]:

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

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


Note:

We already have preprocessed both essay and project_title in Text processing section (1.3 and 1.4) above

2.4 Applying Decision Tree on different kind of featurization as mentioned in the instructions

Apply Decision Tree on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instructions

In [256]:

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

BoW

In [257]:

```
### BoW in Essay and Title on Train

# # We are considering only the bigram words which appeared in at least 10 documents with max feature =
5000(rows or projects).
vectorizer_bow = CountVectorizer(min_df=10, max_features=5000)
tr_essay = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on train",tr_essay.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = CountVectorizer(min_df=10, max_features=5000)
tr_title = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig ",tr_title.shape)

### BoW in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encodig on test",ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encodig on test",ts_title.shape)
```

Shape of essay matrix after one hot encodig on train (33499, 5000)

Shape of title matrix after one hot encodig (33499, 1554)

=====

Shape of essay matrix after one hot encodig on test (16500, 5000)

Shape of title matrix after one hot encodig on test (16500, 1554)

In [258]:

```
print('Shape of normalized essay in train data', tr_essay.shape)
print('Shape of normalized title in train data', tr_title.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

Shape of normalized essay in train data (33499, 5000)

Shape of normalized title in train data (33499, 1554)

=====

Shape of normalized essay in test data (16500, 5000)
Shape of normalized title in test data (16500, 1554)

In [259]:

```
for i in tqdm(range(len(vectorizer_bow.get_feature_names()))):
    list_features.append(vectorizer_bow.get_feature_names()[i])
for i in tqdm(range(len(vectorizer_bowt.get_feature_names()))):
    list_features.append(vectorizer_bowt.get_feature_names()[i])

len(list_features)
```

100% |██| 5000/5000 [00:09<00:00, 502.64it/s]
100% |██| 1554/1554 [00:00<00:00, 2022.95it/s]

Out[259]:

6656

TFIDF

In [138]:

```
### BoW in Essay and Title on Train

# # We are considering only the bigram words which appeared in at least 10 documents with max feature = 5000 (rows or projects).
vectorizer_bow = TfidfVectorizer(min_df=10, max_features=5000)
tr_essay = vectorizer_bow.fit_transform(tr_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on train", tr_essay.shape)

# # Similarly you can vectorize for title also
vectorizer_bowt = TfidfVectorizer(min_df=10, max_features=5000)
tr_title = vectorizer_bowt.fit_transform(tr_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding ", tr_title.shape)

### BoW in Essay and Title on Test
print('=====\\n')
ts_essay = vectorizer_bow.transform(ts_X['clean_essay'].values)
print("Shape of essay matrix after one hot encoding on test", ts_essay.shape)

ts_title = vectorizer_bowt.transform(ts_X['clean_project_title'].values)
print("Shape of title matrix after one hot encoding on test", ts_title.shape)
```

Shape of essay matrix after one hot encoding on train (33498, 5000)
Shape of title matrix after one hot encoding (33498, 1543)
=====

Shape of essay matrix after one hot encoding on test (16500, 5000)
Shape of title matrix after one hot encoding on test (16500, 1543)

In [139]:

```
print('Shape of normalized essay in train data', tr_essay.shape)
print('Shape of normalized title in train data', tr_title.shape)
print('=====\\n')
print('Shape of normalized essay in test data', ts_essay.shape)
print('Shape of normalized title in test data', ts_title.shape)
```

Shape of normalized essay in train data (33498, 5000)
Shape of normalized title in train data (33498, 1543)
=====

Shape of normalized essay in test data (16500, 5000)
Shape of normalized title in test data (16500, 1543)

```
shape of normalized title in test data (16500, 1543)
```

In [140]:

[illegible]

Out[140]:

6645

AVG W2V

In [54]:

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

In [55]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words=0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay.append(vector)

avg_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words=0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title.append(vector)

tr_essay = np.array(avg_w2v_essay)
tr_title = np.array(avg_w2v_title)
print('===== Train Essay =====')
print(len(avg_w2v_essay))
print(len(avg_w2v_essay[0]))
print('===== Train Title =====')
print(len(avg_w2v_title))
print(len(avg_w2v_title[0]))
# print(avg_w2v_essay[0])
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 33498/33498 [00:07<00:00, 4642.08it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 33498/33498 [00:00<00:00, 95834.07it/s]
```

```
===== Train Essay =====
33498
300
===== Train Title =====
33498
300
```

In [56]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay.append(vector)

avg_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title.append(vector)

ts_essay = np.array(avg_w2v_essay)
ts_title = np.array(avg_w2v_title)
print('===== Test Essay =====')
print(len(avg_w2v_essay))
print(len(avg_w2v_essay[0]))
print('===== Test Title =====')
print(len(avg_w2v_title))
print(len(avg_w2v_title[0]))
# print(avg_w2v_essay[0])
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:03<00:00, 4475.00it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:00<00:00, 83120.66it/s]
```

```
===== Test Essay =====
16500
300
===== Test Title =====
16500
300
```

TFIDF W2V

In [96]:

```
# Tfidf weighted w2v on essay in train
tfidf_model = TfidfVectorizer()
tfidf_model.fit(tr_X['clean_essay'].values)
```

```

# 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())

# tfidf Word2Vec
# compute average word2vec for each essay
tfidf_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)
            )/len(sentence.split()))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf
            value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay.append(vector)

tr_essay = np.array(tfidf_w2v_essay)

# compute average word2vec for each essay for test data
tfidf_w2v_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)
            )/len(sentence.split()))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf
            value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_essay.append(vector)

ts_essay = np.array(tfidf_w2v_essay)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33498/33498 [00:51<00:00, 654.18it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:24<00:00, 670.24it/s]

```

In [97]:

```

# tfidf Word2Vec on title
# compute average word2vec for each title for train data
tfidf_model = TfidfVectorizer()
tfidf_model.fit(tr_X['clean_project_title'].values)
# 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())

tfidf_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(tr_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)
            )/len(sentence.split()))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf
            value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:

```

```

        vector /= tf_idf_weight
        tfidf_w2v_title.append(vector)

tr_title = np.array(tfidf_w2v_title)

# compute average word2vec for each title for test data
tfidf_w2v_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(ts_X['clean_project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)
            )/len(sentence.split()))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf
            value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
        tfidf_w2v_title.append(vector)

ts_title = np.array(tfidf_w2v_title)

```

```

100%|████████████████████████████████████████████████████████████████████████████████| 33498/33498 [00:00<00:00, 46195.61it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 16500/16500 [00:00<00:00, 48373.49it/s]

```

In [98]:

```

print('Train essay and title shape:',tr_essay.shape,tr_title.shape)
print('Test essay and title shape:',ts_essay.shape,ts_title.shape)

```

```

Train essay and title shape: (33498, 300) (33498, 300)
Test essay and title shape: (16500, 300) (16500, 300)

```

Merge Them

In [141]:

```

quantity_normalized.T.shape, price_normalized.T.shape, teacher_number_of_previously_posted_projects_normalized.T.shape

```

Out[141]:

```
((33498, 1), (33498, 1), (33498, 1))
```

In [142]:

```

school_state_one_hot.shape, categories_one_hot.shape, subcategories_one_hot.shape, project_grade_category_one_hot.shape, \
    teacher_prefix_one_hot.shape

```

Out[142]:

```
((33498, 51), (33498, 9), (33498, 30), (33498, 4), (33498, 5))
```

In [143]:

```
tr_essay.shape, tr_title.shape
```

Out[143]:

```
((33498, 5000), (33498, 1543))
```

In [144]:

```
# for train data
from scipy.sparse import hstack
tr_X = hstack((quantity_normalized.T, price_normalized.T, teacher_number_of_previously_posted_projects_
normalized.T, \
               school_state_one_hot, categories_one_hot, subcategories_one_hot, project_grade_category_o
ne_hot, \
               teacher_prefix_one_hot, tr_essay, tr_title))
tr_X.shape
```

Out[144]:

(33498, 6645)

In [145]:

```
tr_X = tr_X.toarray()
```

In [146]:

```
tr_X
```

Out[146]:

```
array([[0.00538104, 0.00010463, 0.00036885, ..., 0.          , 0.          ,
        0.          ],
       [0.00017937, 0.00504037, 0.02342177, ..., 0.          , 0.          ,
        0.          ],
       [0.0044842 , 0.00039572, 0.00129096, ..., 0.          , 0.          ,
        0.          ],
       ...,
       [0.00663661, 0.00161304, 0.00018442, ..., 0.          , 0.          ,
        0.          ],
       [0.00089684, 0.00230254, 0.00036885, ..., 0.          , 0.          ,
        0.          ],
       [0.00089684, 0.00058183, 0.00018442, ..., 0.          , 0.          ,
        0.          ]])
```

In [147]:

```
tr_X.shape, tr_y.shape
```

Out[147]:

((33498, 6645), (33498,))

In [148]:

```
# for test data
ts_text = ts_X
ts_X = hstack((ts_quantity.T, ts_price.T, ts_teacher_number_of_previously_posted_projects.T, ts_school_
state, \
               ts_project_subject_category, ts_project_subject_subcategory, ts_project_grade_category, \
               ts_teacher_prefix, ts_essay, ts_title))
ts_X.shape
```

Out[148]:

(16500, 6645)

In [149]:

```
ts_X = ts_X.toarray()
```

In [150]:

```
ts_X.shape, ts_y.shape
```

Out[150]:

```
((16500, 6645), (16500,))
```

In [151]:

```
# check whether data still contain NaN or infinity or not
np.any(np.isnan(tr_X)), np.any(np.isnan(ts_X))
```

Out[151]:

```
(False, False)
```

In [152]:

```
np.all(np.isfinite(tr_X)), np.all(np.isfinite(ts_X))
```

Out[152]:

```
(True, True)
```

In [153]:

```
len(list_features)
```

Out[153]:

```
6645
```

In [154]:

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

In [71]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

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

def plot_cm(feature_names, tr_thresholds, train_fpr, train_tpr, y_train, y_train_pred, y_test, y_test_p
red):
    """
    Parameters:
    feature_name - (string) Write feature to print the plot title
    tr_thresholds - train threshold value
    train_fpr = FPR for train data
    train_tpr - TPR for train data
    y_true - test class data
    y_pred - test prediction value
```



```

y_pred = test_prediction_value

Return:
Plot the confusion matrix for Train and Test Data
"""
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
cm = metrics.confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
plt.figure(figsize = (10,7))
sns.heatmap(cm, annot=True, fmt="d")
plt.xlabel('Predicted Class')
plt.ylabel('True Class')
plt.title('Confusion matrix for Train Data when DecisionTree with {0} features'.format(feature_name
s))

print("Test confusion matrix")
cm = metrics.confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
plt.figure(figsize = (10,7))
sns.heatmap(cm, annot=True, fmt="d")
plt.xlabel('Predicted Class')
plt.ylabel('True Class')
plt.title('Confusion matrix for Test Data when DecisionTree with {0} features'.format(feature_names
))

```

In [155]:

```
clf = DecisionTreeClassifier(class_weight='balanced', random_state=1)
```

In [156]:

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

2.4.1 Applying Decision Trees on BOW, SET 1

In [277]:

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

In [278]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 28 candidates, totalling 84 fits

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

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5775526369340258, total= 3.5s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 4.2s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5671512052355716, total= 3.3s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 8.2s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5701297917152118, total= 3.2s
[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5775526369340258, total= 3.4s
[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5671512052355716, total= 3.5s
```

```

[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5701297917152118, total= 3.4s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5775526369340258, total= 3.7s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5671512052355716, total= 3.4s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5701297917152118, total= 3.4s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5775526369340258, total= 3.4s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5671512052355716, total= 3.3s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5701297917152118, total= 3.4s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6509033875831284, total= 10.9s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6510476931362127, total= 10.9s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6702323364958912, total= 10.8s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6510731317269554, total= 10.7s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6510476931362127, total= 10.9s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6702323364958912, total= 10.9s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6504560592659553, total= 11.1s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6520887983046797, total= 11.1s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6702323364958912, total= 10.9s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.6504560592659553, total= 12.0s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.6518758144328247, total= 11.5s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.671009829197873, total= 10.2s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.6475658777272185, total= 17.2s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.6594027550956117, total= 16.5s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.6687423869594741, total= 16.5s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.6485848113240092, total= 17.2s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.6627373043486087, total= 16.9s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.6668109998309881, total= 17.4s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.6516274563527558, total= 18.6s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.6664867893120411, total= 18.0s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.6702785272754366, total= 17.4s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.6549851217632991, total= 18.4s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.6718727342141292, total= 17.0s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.6786414022895626, total= 16.8s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.586766819044743, total= 38.4s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.5969071616584778, total= 35.6s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.5900872086918172, total= 33.7s
[CV] max_depth=50, min_samples_split=10 .....
[CV] max_depth=50, min_samples_split=10, score=0.5903949438744363, total= 37.0s
[CV] max_depth=50, min_samples_split=10 .....
[CV] max_depth=50, min_samples_split=10, score=0.6000199451661069, total= 32.8s
[CV] max_depth=50, min_samples_split=10 .....
[CV] max_depth=50, min_samples_split=10, score=0.5893198729609802, total= 30.9s
[CV] max_depth=50, min_samples_split=100 .....
[CV] max_depth=50, min_samples_split=100, score=0.6221107559289205, total= 34.9s
[CV] max_depth=50, min_samples_split=100 .....

```

[CV] max_depth=50, min_samples_split=100, score=0.6299957609490274, total= 32.3s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.620631752354417, total= 30.1s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.63906679594554, total= 32.9s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6436538556737528, total= 30.4s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6309142949149389, total= 28.6s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.576881347593508, total= 45.3s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5861185519171861, total= 42.1s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5739863467806214, total= 39.0s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5813300691126156, total= 45.4s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5847951142031085, total= 43.1s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5761711894044477, total= 38.8s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.6137022960207874, total= 45.4s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.6057241626655493, total= 41.4s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.6078037478160949, total= 35.0s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6238996942105498, total= 38.4s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6277858227934183, total= 35.5s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6215679750762304, total= 31.1s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5680814324405808, total= 56.2s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5744675371951409, total= 57.1s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5634695063349497, total= 48.8s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5705678779082123, total= 57.4s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5734357764401273, total= 58.4s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.567297090593432, total= 45.6s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.59817048811816, total= 59.6s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5902326877708317, total= 57.0s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.584893308674916, total= 45.6s
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.5923298458703175, total= 1.0min
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.6013021299512264, total= 54.8s
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.601657873959051, total= 46.4s
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.5680814324405808, total= 1.2min
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.5744675371951409, total= 1.1min
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.5634695063349497, total= 57.1s
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.5705678779082123, total= 1.1min
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.5734357764401273, total= 1.1min
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.567297090593432, total= 49.6s
[CV] max_depth=1000, min_samples_split=100
[CV] max_depth=1000, min_samples_split=100, score=0.59817048811816, total= 1.1min
[CV] max_depth=1000, min_samples_split=100
[CV] max_depth=1000, min_samples_split=100, score=0.5902326877708317, total= 1.1min
[CV] max_depth=1000, min_samples_split=100
[CV] max_depth=1000, min_samples_split=100, score=0.584893308674916, total= 50.1s
[CV] max_depth=1000, min_samples_split=500
[CV] max_depth=1000, min_samples_split=500, score=0.5923298458703175, total= 1.0min

```
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6013021299512264, total= 57.5s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.601657873959051, total= 48.5s
```

```
[Parallel(n_jobs=1)]: Done 84 out of 84 | elapsed: 44.8min finished
```

Out[278]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=DecisionTreeClassifier(class_weight='balanced', 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=1,
                                              splitter='best'),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

Best param for BoW feature

In [279]:

```
best_d = clf.best_params_['max_depth']
best_split = clf.best_params_['min_samples_split']
best_d, best_split
```

Out[279]:

```
(10, 500)
```

Graphviz for BoW feature

In [280]:

```
# Taking max_depth=2 as per task
lr = DecisionTreeClassifier(max_depth=3, class_weight='balanced', min_samples_split=500)
lr.fit(tr_X, tr_y)
```

Out[280]:

```
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=3,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=500,
                       min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                       splitter='best')
```

In [281]:

```
from sklearn.externals.six import StringIO
dot_data = StringIO()
tree.export_graphviz(lr, out_file=dot_data, \
                     filled=True, rounded=True, special_characters=True, \
                     feature_names=list_features, class_names=['not approved', 'approved'])
```

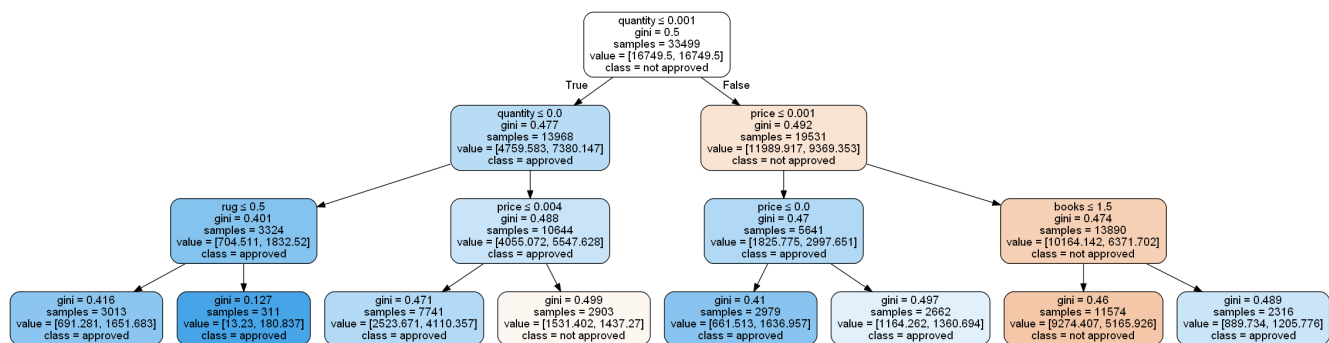
In [282]:

```
import pydotplus
from IPython.display import Image
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
```

In [283]:

```
Image(graph.create_png())
```

Out[283]:



In [284]:

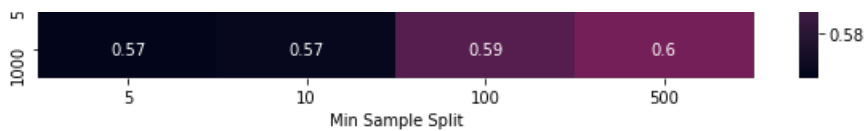
```
# Plot seaborn heatmap for gridsearchcv: https://www.kaggle.com/arindambanerjee/grid-search-simplified
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
samplesplit_list = list(clf.cv_results_['param_min_samples_split'].data)

plt.figure(1, figsize=(10,10))
plt.subplot(211)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                          'Min Sample Split':samplesplit_list, \
                          'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')

plt.subplot(212)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                          'Min Sample Split':samplesplit_list, \
                          'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')

plt.show()
```





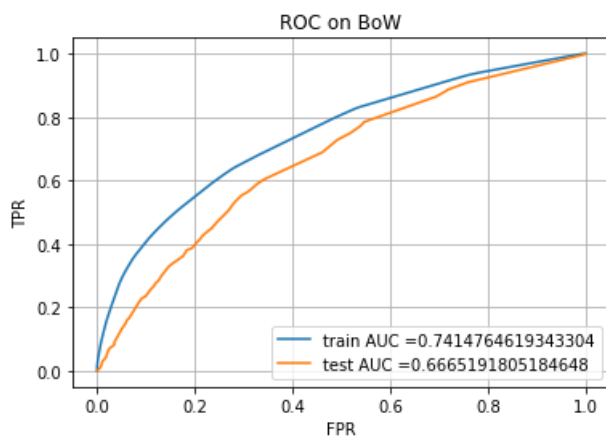
In [285]:

```
lr = DecisionTreeClassifier(max_depth=best_d, class_weight='balanced', min_samples_split=best_split, random_state=1)
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive classes
# not the predicted outputs

y_train_pred = lr.predict_proba(tr_X)[:,1]
y_test_pred = lr.predict_proba(ts_X)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("FPR")
plt.ylabel("TPR")
plt.title("ROC on BoW")
plt.grid()
plt.show()
```



In [286]:

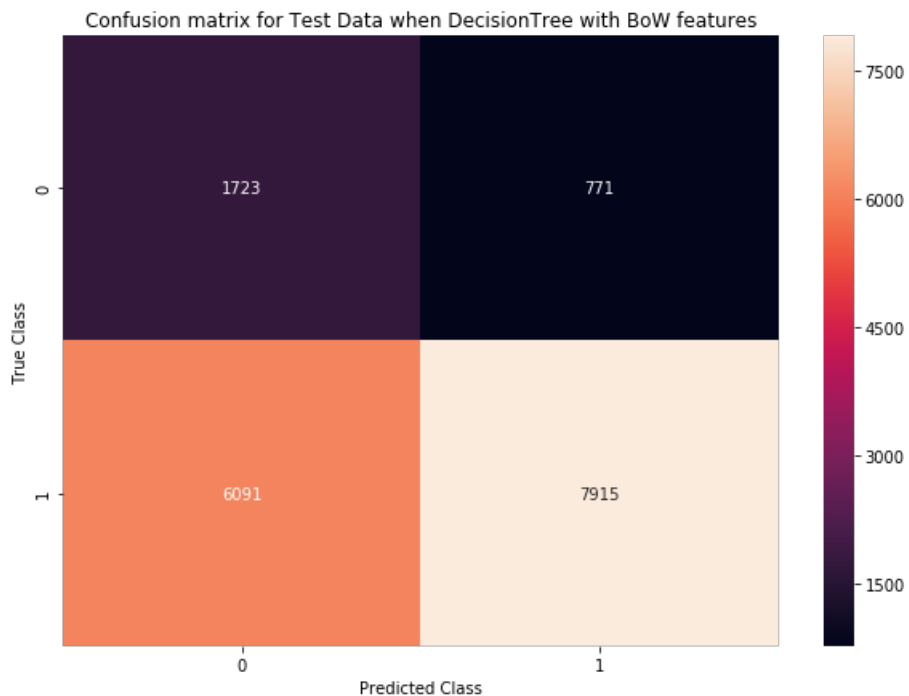
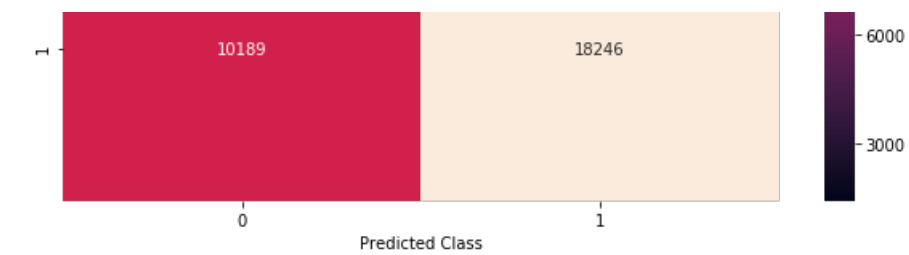
```
plot_cm('BoW', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr * (1 - fpr)$ 0.4607280094203376 for threshold 0.475

Train confusion matrix

Test confusion matrix





In [287]:

```
# Predict test data with best threshold value
ts_predict = predict_with_best_t(lr.predict_proba(ts_X[:,1]), 0.475)

false_datapoints = []

# Iterate over all data points in test data
for i in range(ts_X.shape[0]):
    # Select that data point where test true value is 0 and test predicted value is 1
    if (ts_y[i] == 0) and (ts_predict[i] == 1):
        # If it true, then put that datapoint into another array
        false_datapoints.append(ts_text.iloc[i].values)
```

In [288]:

```
false_datapoints = pd.DataFrame(data=false_datapoints, columns=ts_text.columns)
false_datapoints.shape
```

Out[288]:

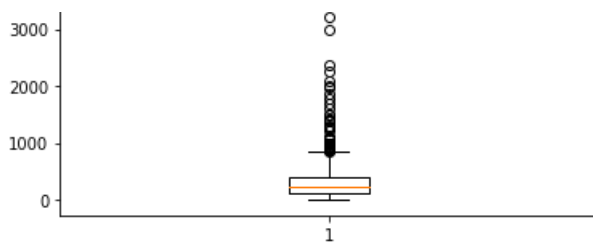
(771, 10)

In [289]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# WordCloud on Essay
from wordcloud import WordCloud

comment_words = ''
for val in false_datapoints['clean_essay']:

    # typecaste each val to string
    val = str(val)
```

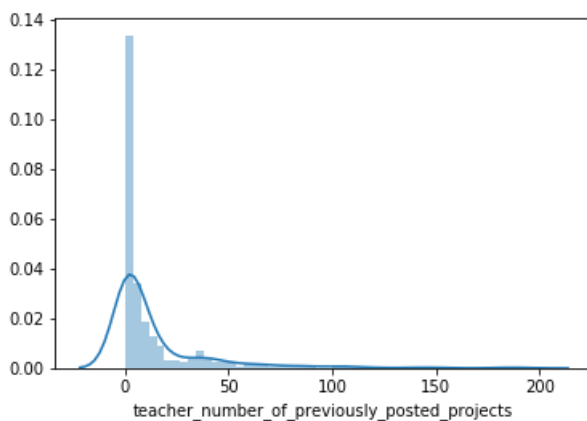



In [291]:

```
# PDF on teacher_number_of_previously_posted_projects
sns.distplot(false_datapoints['teacher_number_of_previously_posted_projects'], hist=True, kde=True)
```

Out[291]:

<matplotlib.axes._subplots.AxesSubplot at 0x1fca3362a20>



2.4.2 Applying Decision Trees on TFIDF, SET 2

In [74]:

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

In [157]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 28 candidates, totalling 84 fits

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

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5661502566606066, total= 3.2s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 3.8s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5616270486654009, total= 3.2s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 7.6s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5659347823809485, total= 3.2s
[CV] max_depth=1, min_samples_split=10 .....
```

[CV]	max_depth=1, min_samples_split=10, score=0.5661502566606066, total=	3.2s
[CV]	max_depth=1, min_samples_split=10	
[CV]	max_depth=1, min_samples_split=10, score=0.5616270486654009, total=	3.2s
[CV]	max_depth=1, min_samples_split=10	
[CV]	max_depth=1, min_samples_split=10, score=0.5659347823809485, total=	3.1s
[CV]	max_depth=1, min_samples_split=100	
[CV]	max_depth=1, min_samples_split=100, score=0.5661502566606066, total=	3.1s
[CV]	max_depth=1, min_samples_split=100	
[CV]	max_depth=1, min_samples_split=100, score=0.5616270486654009, total=	3.2s
[CV]	max_depth=1, min_samples_split=100	
[CV]	max_depth=1, min_samples_split=100, score=0.5659347823809485, total=	3.1s
[CV]	max_depth=1, min_samples_split=500	
[CV]	max_depth=1, min_samples_split=500, score=0.5661502566606066, total=	3.2s
[CV]	max_depth=1, min_samples_split=500	
[CV]	max_depth=1, min_samples_split=500, score=0.5616270486654009, total=	3.3s
[CV]	max_depth=1, min_samples_split=500	
[CV]	max_depth=1, min_samples_split=500, score=0.5659347823809485, total=	3.3s
[CV]	max_depth=5, min_samples_split=5	
[CV]	max_depth=5, min_samples_split=5, score=0.6503522532379667, total=	10.8s
[CV]	max_depth=5, min_samples_split=5	
[CV]	max_depth=5, min_samples_split=5, score=0.6508975250985194, total=	10.7s
[CV]	max_depth=5, min_samples_split=5	
[CV]	max_depth=5, min_samples_split=5, score=0.6507092187524943, total=	10.5s
[CV]	max_depth=5, min_samples_split=10	
[CV]	max_depth=5, min_samples_split=10, score=0.6497207436383853, total=	10.5s
[CV]	max_depth=5, min_samples_split=10	
[CV]	max_depth=5, min_samples_split=10, score=0.6506870619974311, total=	10.6s
[CV]	max_depth=5, min_samples_split=10	
[CV]	max_depth=5, min_samples_split=10, score=0.6507092187524943, total=	10.8s
[CV]	max_depth=5, min_samples_split=100	
[CV]	max_depth=5, min_samples_split=100, score=0.6497207436383853, total=	10.7s
[CV]	max_depth=5, min_samples_split=100	
[CV]	max_depth=5, min_samples_split=100, score=0.6507553363022113, total=	10.9s
[CV]	max_depth=5, min_samples_split=100	
[CV]	max_depth=5, min_samples_split=100, score=0.6509939202159927, total=	10.6s
[CV]	max_depth=5, min_samples_split=500	
[CV]	max_depth=5, min_samples_split=500, score=0.6506202117220995, total=	10.6s
[CV]	max_depth=5, min_samples_split=500	
[CV]	max_depth=5, min_samples_split=500, score=0.6510997927901336, total=	10.8s
[CV]	max_depth=5, min_samples_split=500	
[CV]	max_depth=5, min_samples_split=500, score=0.6509399794064592, total=	10.6s
[CV]	max_depth=10, min_samples_split=5	
[CV]	max_depth=10, min_samples_split=5, score=0.6513237896907272, total=	17.2s
[CV]	max_depth=10, min_samples_split=5	
[CV]	max_depth=10, min_samples_split=5, score=0.6559928573422893, total=	18.3s
[CV]	max_depth=10, min_samples_split=5	
[CV]	max_depth=10, min_samples_split=5, score=0.6505103548893696, total=	17.5s
[CV]	max_depth=10, min_samples_split=10	
[CV]	max_depth=10, min_samples_split=10, score=0.6514928839574002, total=	17.8s
[CV]	max_depth=10, min_samples_split=10	
[CV]	max_depth=10, min_samples_split=10, score=0.654672804353202, total=	18.4s
[CV]	max_depth=10, min_samples_split=10	
[CV]	max_depth=10, min_samples_split=10, score=0.6470103134827829, total=	17.6s
[CV]	max_depth=10, min_samples_split=100	
[CV]	max_depth=10, min_samples_split=100, score=0.6516338737412027, total=	17.4s
[CV]	max_depth=10, min_samples_split=100	
[CV]	max_depth=10, min_samples_split=100, score=0.6632887037716583, total=	18.3s
[CV]	max_depth=10, min_samples_split=100	
[CV]	max_depth=10, min_samples_split=100, score=0.6573942311146054, total=	17.5s
[CV]	max_depth=10, min_samples_split=500	
[CV]	max_depth=10, min_samples_split=500, score=0.6616236776093019, total=	16.7s
[CV]	max_depth=10, min_samples_split=500	
[CV]	max_depth=10, min_samples_split=500, score=0.669277491454338, total=	18.1s
[CV]	max_depth=10, min_samples_split=500	
[CV]	max_depth=10, min_samples_split=500, score=0.6719702850369565, total=	17.0s
[CV]	max_depth=50, min_samples_split=5	
[CV]	max_depth=50, min_samples_split=5, score=0.5676656553622668, total=	30.4s
[CV]	max_depth=50, min_samples_split=5	
[CV]	max_depth=50, min_samples_split=5, score=0.581789220549874, total=	35.1s
[CV]	max_depth=50, min_samples_split=5	
[CV]	max_depth=50, min_samples_split=5, score=0.582390720285037, total=	31.5s
[CV]	max_depth=50, min_samples_split=10	
[CV]	max_depth=50, min_samples_split=10, score=0.569767602722508, total=	30.6s
[CV]	max_depth=50, min_samples_split=10	
[CV]	max_depth=50, min_samples_split=10, score=0.5729840788434538, total=	35.2s
[CV]	max_depth=50, min_samples_split=10	
[CV]	max_depth=50, min_samples_split=10, score=0.5815812339861273, total=	31.5s

[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.5976236008360803, total= 29.8s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.6082884071165321, total= 34.8s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.6072701236011557, total= 31.0s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.636553949203297, total= 27.1s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6467071156346725, total= 32.2s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6235389835531043, total= 28.7s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5552285181110648, total= 33.6s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5556731615858037, total= 38.5s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5655779311872415, total= 34.7s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5564074731253219, total= 34.0s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5605698096820255, total= 38.5s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5656710648277486, total= 35.4s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5868498626120768, total= 33.4s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5897072948741983, total= 37.0s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5806572871601667, total= 35.1s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6272457569397196, total= 29.3s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6313652779384313, total= 34.4s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6228707722877189, total= 29.5s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5571854251522783, total= 37.6s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.550234307691694, total= 40.8s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5635129959311793, total= 36.9s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5539901760188124, total= 34.2s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5533073992768595, total= 41.2s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.559104379270366, total= 36.8s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5722326338287491, total= 34.2s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5730832027884585, total= 38.9s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5682411417342477, total= 37.5s
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.6171257127851468, total= 31.3s
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.6267537864661443, total= 34.7s
[CV] max_depth=500, min_samples_split=500
[CV] max_depth=500, min_samples_split=500, score=0.6228707722877189, total= 29.3s
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.5571854251522783, total= 37.8s
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.550234307691694, total= 41.1s
[CV] max_depth=1000, min_samples_split=5
[CV] max_depth=1000, min_samples_split=5, score=0.5635129959311793, total= 37.0s
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.5539901760188124, total= 34.5s
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.5533073992768595, total= 41.0s
[CV] max_depth=1000, min_samples_split=10
[CV] max_depth=1000, min_samples_split=10, score=0.559104379270366, total= 37.0s
[CV] max_depth=1000, min_samples_split=100
[CV] max_depth=1000, min_samples_split=100, score=0.5722326338287491, total= 34.1s
[CV] max_depth=1000, min_samples_split=100
[CV] max_depth=1000, min_samples_split=100, score=0.5730832027884585, total= 38.5s
[CV] max_depth=1000, min_samples_split=100

```
[CV] max_depth=1000, min_samples_split=100, score=0.5682411417342477, total= 37.7s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6171257127851468, total= 31.1s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6267537864661443, total= 34.4s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6228707722877189, total= 29.8s
```

```
[Parallel(n_jobs=1)]: Done 84 out of 84 | elapsed: 34.9min finished
```

Out[157]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=DecisionTreeClassifier(class_weight='balanced', 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=1,
                                              splitter='best'),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

Best Params for TFIDF Feature

In [158]:

```
best_d = clf.best_params_['max_depth']
best_split = clf.best_params_['min_samples_split']
best_d, best_split
```

Out[158]:

```
(10, 500)
```

Graphviz for TFIDF feature

In [159]:

```
# Taking max_depth=2 as per task
lr = DecisionTreeClassifier(max_depth=3, class_weight='balanced', min_samples_split=best_split)
lr.fit(tr_X, tr_y)
```

Out[159]:

```
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=3,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=500,
                       min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                       splitter='best')
```

In [160]:

```
from sklearn.externals.six import StringIO
dot_data = StringIO()
tree.export_graphviz(lr, out_file=dot_data, \
                     filled=True, rounded=True, special_characters=True, \
                     feature_names=list_features, class_names=['not approved', 'approved'])
```

In [161]:

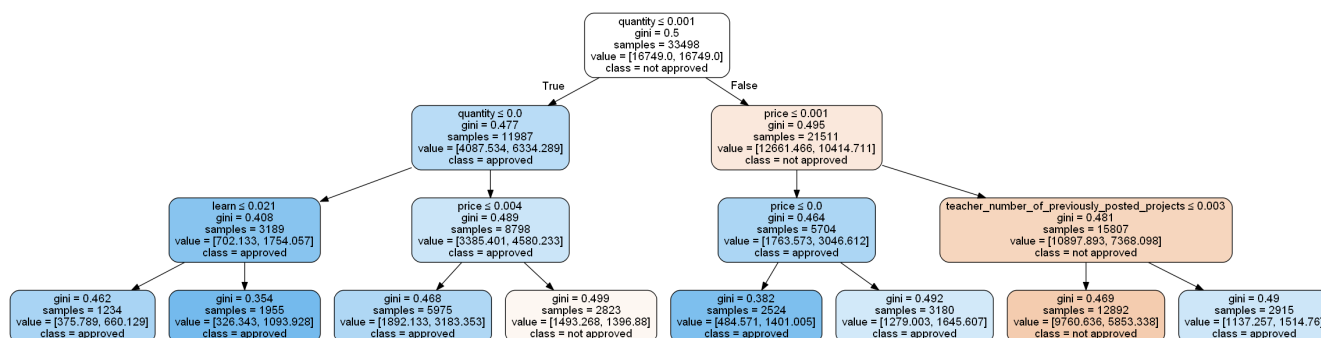
```
import pydotplus
from IPython.display import Image
```

```
from IPython.display import Image
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
```

In [162]:

```
Image(graph.create_png())
```

Out[162]:



In [163]:

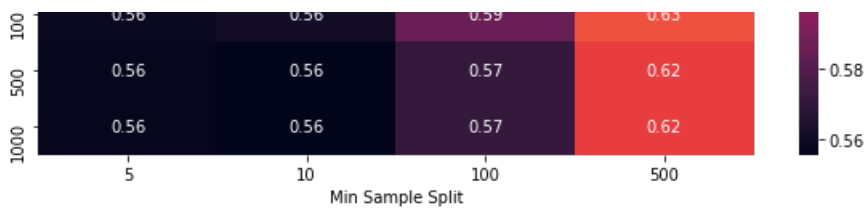
```
# Plot seaborn heatmap for gridsearchcv: https://www.kaggle.com/arindambanerjee/grid-search-simplified
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
samplesplit_list = list(clf.cv_results_['param_min_samples_split'].data)

plt.figure(1, figsize=(10,10))
plt.subplot(211)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                          'Min Sample Split':samplesplit_list, \
                          'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')

plt.subplot(212)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                          'Min Sample Split':samplesplit_list, \
                          'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')

plt.show()
```





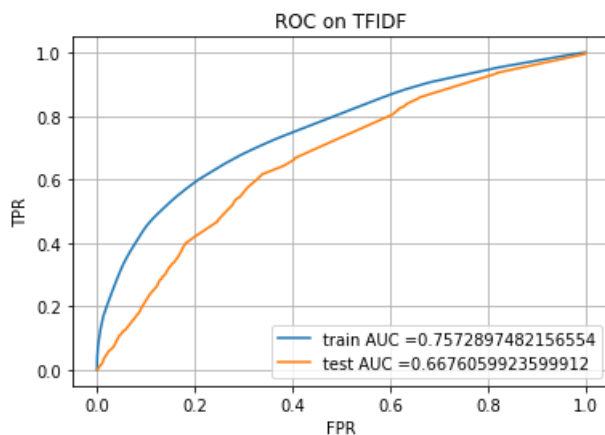
In [164]:

```
lr = DecisionTreeClassifier(max_depth=best_d, class_weight='balanced', min_samples_split=best_split, random_state=1)
lr.fit(tr_X, tr_y)
# 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 = lr.predict_proba(tr_X)[:,1]
y_test_pred = lr.predict_proba(ts_X)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("FPR")
plt.ylabel("TPR")
plt.title("ROC on TFIDF")
plt.grid()
plt.show()
```

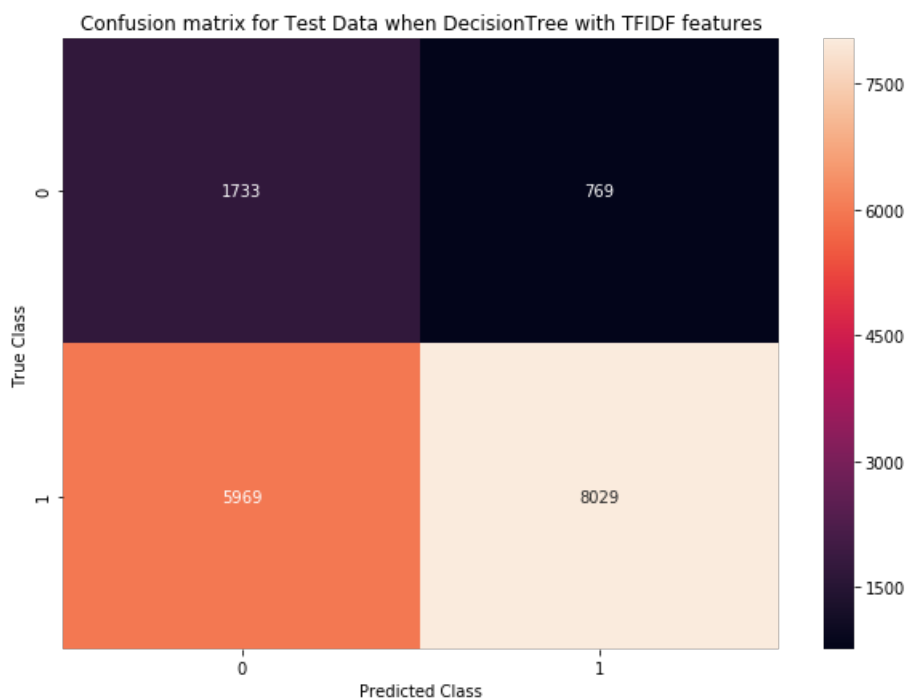
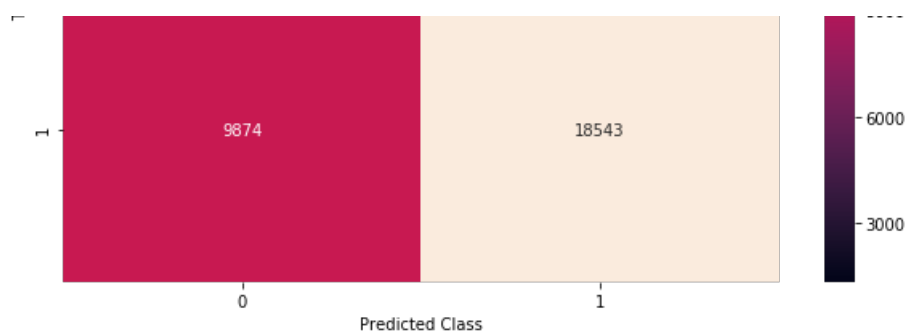


In [165]:

```
plot_cm('TFIDF', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4800559680059899 for threshold 0.469
Train confusion matrix
Test confusion matrix





In [166]:

```
# Predict test data with best threshold value
ts_predict = predict_with_best_t(lr.predict_proba(ts_X[:,1], 0.469)

false_datapoints = []

# Iterate over all data points in test data
for i in range(ts_X.shape[0]):
    # Select that data point where test true value is 0 and test predicted value is 1
    if (ts_y[i] == 0) and (ts_predict[i] == 1):
        # If it true, then put that datapoint into another array
        false_datapoints.append(ts_text.iloc[i].values)
```

In [167]:

```
false_datapoints = pd.DataFrame(data=false_datapoints, columns=ts_text.columns)
false_datapoints.shape
```

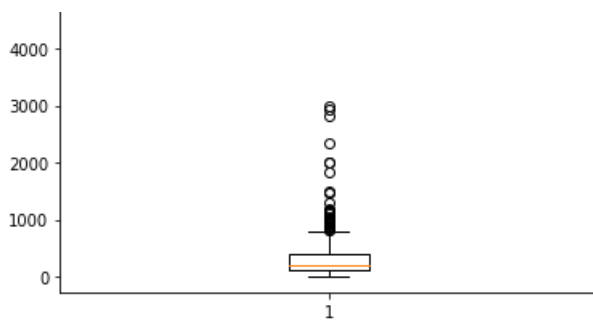
Out[167]:

(769, 10)

In [168]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# WordCloud on Essay
from wordcloud import WordCloud

comment_words = ' '
for val in false_datapoints['clean_essay']:
```

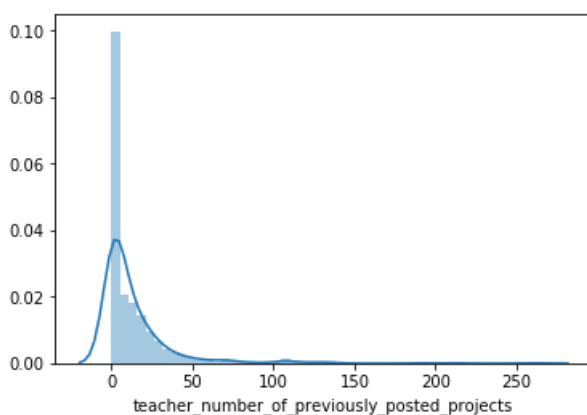



In [170]:

```
# PDF on teacher_number_of_previously_posted_projects
sns.distplot(false_datapoints['teacher_number_of_previously_posted_projects'], hist=True, kde=True)
```

Out[170]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f647028278>



2.4.3 Applying Decision Trees on AVG W2V, SET 3

In [0]:

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

In [74]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_X, tr_y)
```

Fitting 3 folds for each of 28 candidates, totalling 84 fits

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

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5661502566606066, total= 1.5s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 1.6s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5616270486654009, total= 1.5s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 3.2s remaining: 0.0s

```
[CV] max_depth=1, min_samples_split=5
```

```

[CV] max_depth=1, min_samples_split=5 .....
[CV] max_depth=1, min_samples_split=5, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5616270486654009, total= 1.5s
[CV] max_depth=1, min_samples_split=10 .....
[CV] max_depth=1, min_samples_split=10, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5616270486654009, total= 1.5s
[CV] max_depth=1, min_samples_split=100 .....
[CV] max_depth=1, min_samples_split=100, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5616270486654009, total= 1.5s
[CV] max_depth=1, min_samples_split=500 .....
[CV] max_depth=1, min_samples_split=500, score=0.5659347823809485, total= 1.5s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6546874787736376, total= 6.9s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6448343243442674, total= 6.9s
[CV] max_depth=5, min_samples_split=5 .....
[CV] max_depth=5, min_samples_split=5, score=0.6480399281720438, total= 6.9s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6546874787736376, total= 6.7s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6448343243442674, total= 6.9s
[CV] max_depth=5, min_samples_split=10 .....
[CV] max_depth=5, min_samples_split=10, score=0.6480399281720438, total= 6.5s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6546874787736376, total= 6.4s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6440967374916238, total= 6.5s
[CV] max_depth=5, min_samples_split=100 .....
[CV] max_depth=5, min_samples_split=100, score=0.6478287795233951, total= 6.5s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.6548224239125653, total= 6.7s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.6454046687533903, total= 6.9s
[CV] max_depth=5, min_samples_split=500 .....
[CV] max_depth=5, min_samples_split=500, score=0.6491657632880221, total= 6.9s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.5944759610704928, total= 11.8s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.6075583301258018, total= 12.1s
[CV] max_depth=10, min_samples_split=5 .....
[CV] max_depth=10, min_samples_split=5, score=0.5889428761214699, total= 11.8s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.5932629815603434, total= 11.4s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.6043339818197773, total= 11.7s
[CV] max_depth=10, min_samples_split=10 .....
[CV] max_depth=10, min_samples_split=10, score=0.5923265862214843, total= 11.8s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.6094826020787846, total= 11.5s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.617958242425572, total= 11.9s
[CV] max_depth=10, min_samples_split=100 .....
[CV] max_depth=10, min_samples_split=100, score=0.607610761372344, total= 11.6s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.648228993830848, total= 10.1s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.6389260573386993, total= 10.9s
[CV] max_depth=10, min_samples_split=500 .....
[CV] max_depth=10, min_samples_split=500, score=0.6329869813700293, total= 10.2s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.5408646970430222, total= 17.0s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.5355285023253773, total= 18.3s
[CV] max_depth=50, min_samples_split=5 .....
[CV] max_depth=50, min_samples_split=5, score=0.5357682443407672, total= 17.7s
[CV] max_depth=50, min_samples_split=10 .....
[CV] max_depth=50, min_samples_split=10, score=0.5449109947852787, total= 17.2s
[CV] max_depth=50, min_samples_split=10 .....

```

```

[CV] max_depth=50, min_samples_split=10, score=0.5335133664324004, total= 18.3s
[CV] max_depth=50, min_samples_split=10 .....
[CV] max_depth=50, min_samples_split=10, score=0.5409420098258328, total= 17.3s
[CV] max_depth=50, min_samples_split=100 .....
[CV] max_depth=50, min_samples_split=100, score=0.5743079411303935, total= 15.9s
[CV] max_depth=50, min_samples_split=100 .....
[CV] max_depth=50, min_samples_split=100, score=0.5674650782072657, total= 16.9s
[CV] max_depth=50, min_samples_split=100 .....
[CV] max_depth=50, min_samples_split=100, score=0.5659443545361664, total= 16.0s
[CV] max_depth=50, min_samples_split=500 .....
[CV] max_depth=50, min_samples_split=500, score=0.6352298603961224, total= 11.8s
[CV] max_depth=50, min_samples_split=500 .....
[CV] max_depth=50, min_samples_split=500, score=0.6257025616045503, total= 12.2s
[CV] max_depth=50, min_samples_split=500 .....
[CV] max_depth=50, min_samples_split=500, score=0.6271312543900959, total= 10.7s
[CV] max_depth=100, min_samples_split=5 .....
[CV] max_depth=100, min_samples_split=5, score=0.5346249098444333, total= 18.0s
[CV] max_depth=100, min_samples_split=5 .....
[CV] max_depth=100, min_samples_split=5, score=0.5322112623248987, total= 19.0s
[CV] max_depth=100, min_samples_split=5 .....
[CV] max_depth=100, min_samples_split=5, score=0.5352628844327199, total= 17.9s
[CV] max_depth=100, min_samples_split=10 .....
[CV] max_depth=100, min_samples_split=10, score=0.5417444421359855, total= 18.1s
[CV] max_depth=100, min_samples_split=10 .....
[CV] max_depth=100, min_samples_split=10, score=0.5383047518168257, total= 18.0s
[CV] max_depth=100, min_samples_split=10 .....
[CV] max_depth=100, min_samples_split=10, score=0.5363705418076818, total= 17.3s
[CV] max_depth=100, min_samples_split=100 .....
[CV] max_depth=100, min_samples_split=100, score=0.5666579756721115, total= 15.4s
[CV] max_depth=100, min_samples_split=100 .....
[CV] max_depth=100, min_samples_split=100, score=0.5613772289020869, total= 16.6s
[CV] max_depth=100, min_samples_split=100 .....
[CV] max_depth=100, min_samples_split=100, score=0.564716031882074, total= 15.7s
[CV] max_depth=100, min_samples_split=500 .....
[CV] max_depth=100, min_samples_split=500, score=0.6352298603961224, total= 11.3s
[CV] max_depth=100, min_samples_split=500 .....
[CV] max_depth=100, min_samples_split=500, score=0.6257025616045503, total= 11.6s
[CV] max_depth=100, min_samples_split=500 .....
[CV] max_depth=100, min_samples_split=500, score=0.6271312543900959, total= 10.4s
[CV] max_depth=500, min_samples_split=5 .....
[CV] max_depth=500, min_samples_split=5, score=0.5346249098444333, total= 16.9s
[CV] max_depth=500, min_samples_split=5 .....
[CV] max_depth=500, min_samples_split=5, score=0.5322112623248987, total= 18.4s
[CV] max_depth=500, min_samples_split=5 .....
[CV] max_depth=500, min_samples_split=5, score=0.5352628844327199, total= 18.4s
[CV] max_depth=500, min_samples_split=10 .....
[CV] max_depth=500, min_samples_split=10, score=0.5417444421359855, total= 18.1s
[CV] max_depth=500, min_samples_split=10 .....
[CV] max_depth=500, min_samples_split=10, score=0.5383047518168257, total= 19.2s
[CV] max_depth=500, min_samples_split=10 .....
[CV] max_depth=500, min_samples_split=10, score=0.5363705418076818, total= 18.4s
[CV] max_depth=500, min_samples_split=100 .....
[CV] max_depth=500, min_samples_split=100, score=0.5666579756721115, total= 16.6s
[CV] max_depth=500, min_samples_split=100 .....
[CV] max_depth=500, min_samples_split=100, score=0.5613772289020869, total= 17.6s
[CV] max_depth=500, min_samples_split=100 .....
[CV] max_depth=500, min_samples_split=100, score=0.564716031882074, total= 16.9s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6352298603961224, total= 11.8s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6257025616045503, total= 12.5s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6271312543900959, total= 11.0s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5346249098444333, total= 17.9s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5322112623248987, total= 19.2s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5352628844327199, total= 18.6s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.5417444421359855, total= 18.0s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.5383047518168257, total= 19.0s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.5363705418076818, total= 18.2s
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.5666579756721115, total= 16.3s

```

```
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.5613772289020869, total= 17.7s
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.564716031882074, total= 16.0s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6352298603961224, total= 11.8s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6257025616045503, total= 12.4s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6271312543900959, total= 11.0s
```

```
[Parallel(n_jobs=1)]: Done 84 out of 84 | elapsed: 17.0min finished
```

Out[74]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=DecisionTreeClassifier(class_weight='balanced', 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=1,
                                              splitter='best'),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

Best Params for AVGW2V Feature

In [75]:

```
best_d = clf.best_params_['max_depth']
best_split = clf.best_params_['min_samples_split']
best_d, best_split
```

Out[75]:

```
(5, 500)
```

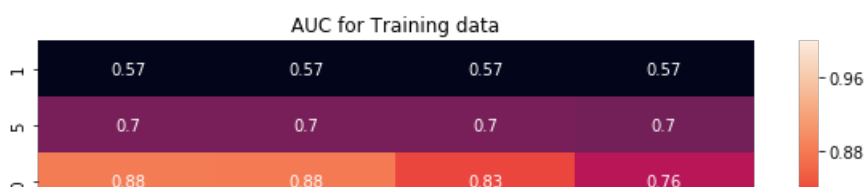
In [76]:

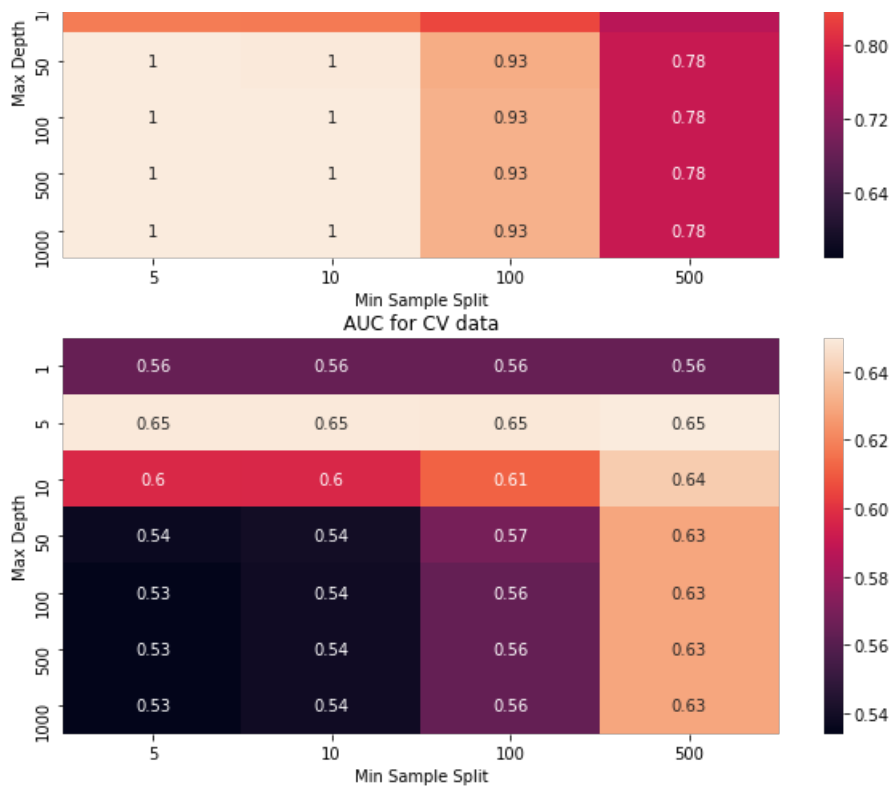
```
# Plot seaborn heatmap for gridsearchcv: https://www.kaggle.com/arindambanerjee/grid-search-simplified
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
samplesplit_list = list(clf.cv_results_['param_min_samples_split'].data)

plt.figure(1, figsize=(10,10))
plt.subplot(211)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                        'Min Sample Split':samplesplit_list, \
                        'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')

plt.subplot(212)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                        'Min Sample Split':samplesplit_list, \
                        'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')

plt.show()
```





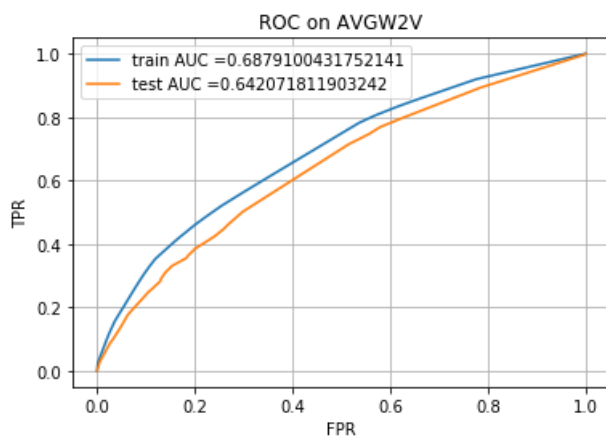
In [77]:

```
lr = DecisionTreeClassifier(max_depth=best_d, class_weight='balanced', min_samples_split=best_split, random_state=1)
lr.fit(tr_X, tr_y)
# 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 = lr.predict_proba(tr_X)[:,1]
y_test_pred = lr.predict_proba(ts_X)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("FPR")
plt.ylabel("TPR")
plt.title("ROC on AVGW2V")
plt.grid()
plt.show()
```



In [78]:

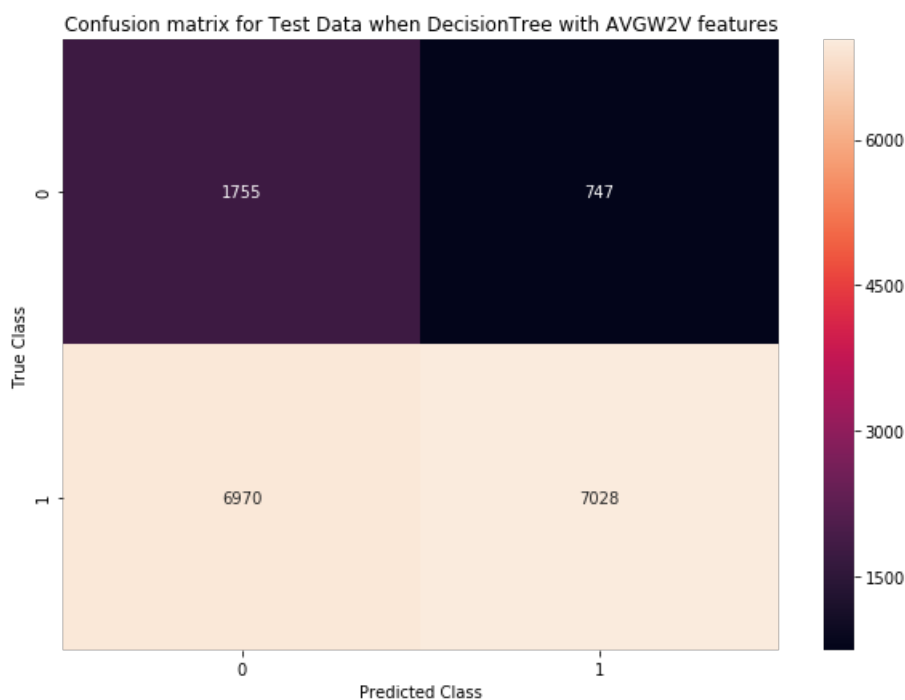
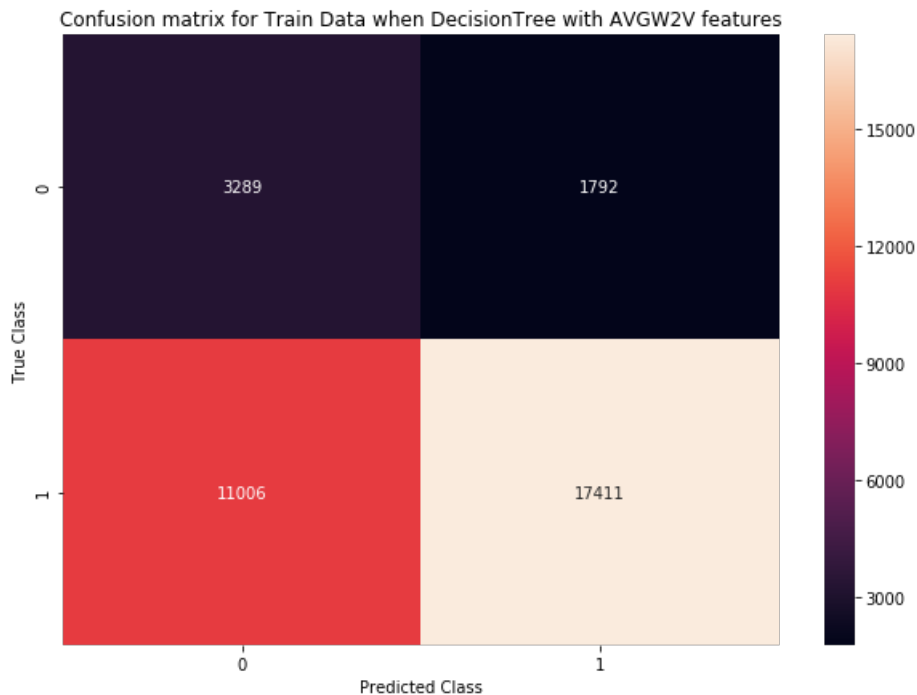
```
plot_cm('AVGW2V', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

```
plot_cm(AVGW2V, cf_thresholds, train_tpr, train_fpr, cf_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.39660680977732465 for threshold 0.482

Train confusion matrix

Test confusion matrix



In [80]:

```
# Predict test data with best threshold value
ts_predict = predict_with_best_t(y_test_pred, 0.482)

false_datapoints = []

# Iterate over all data points in test data
for i in range(ts_X.shape[0]):
    # Select that data point where test true value is 0 and test predicted value is 1
    if (ts_y[i] == 0) and (ts_predict[i] == 1):
        # If it true, then put that datapoint into another array
        false_datapoints.append(ts_text.iloc[i].values)
```

```
false_datapoints = pd.DataFrame(data=false_datapoints, columns=ts_text.columns)
false_datapoints.shape
```

 $(747, 10)$

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# WordCloud on Essay
from wordcloud import WordCloud
```

```
comment_words = ' '
for val in false_datapoints['clean_essay']:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()

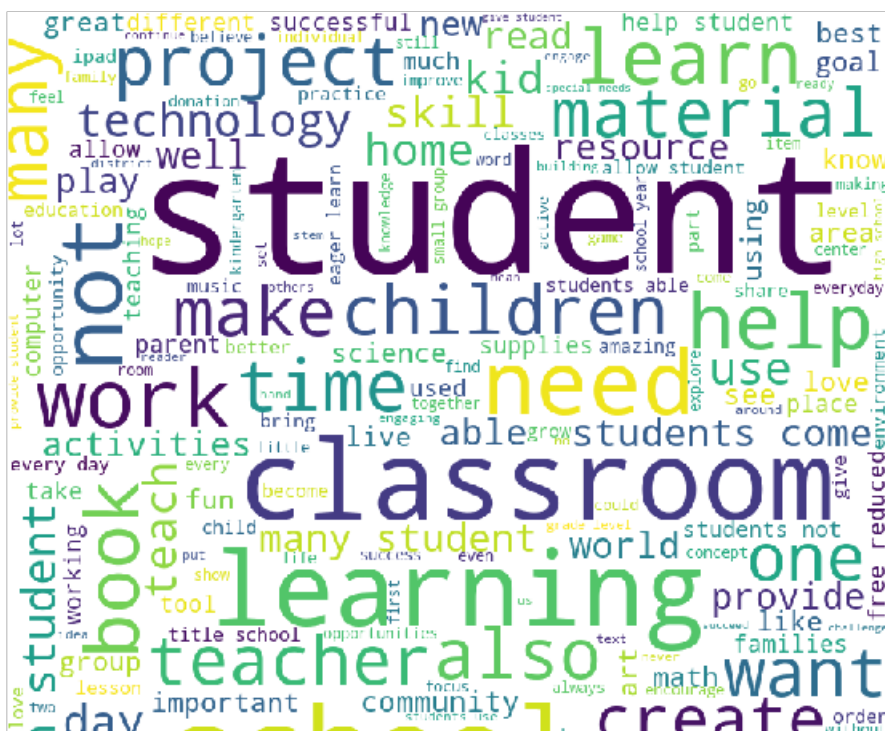
    # Converts each token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

    for words in tokens:
        comment_words = comment_words + words + ' '
```

```
wordcloud = WordCloud(width = 800, height = 800,  
                        background_color='white',  
                        stopwords = stopwords,  
                        min font size = 10).generate(comment words)
```

```
# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)

plt.show()
```



Fitting 3 folds for each of 28 candidates, totalling 84 fits

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

[CV] max_depth=1, min_samples_split=5
[CV] max_depth=1, min_samples_split=5, score=0.5661502566606066, total= 1.5s

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 1.5s remaining: 0.0s

[CV] max_depth=1, min_samples_split=5
[CV] max_depth=1, min_samples_split=5, score=0.5616270486654009, total= 1.5s

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 3.2s remaining: 0.0s

[CV] max_depth=1, min_samples_split=5
[CV] max_depth=1, min_samples_split=5, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=10
[CV] max_depth=1, min_samples_split=10, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=10
[CV] max_depth=1, min_samples_split=10, score=0.5616270486654009, total= 1.5s
[CV] max_depth=1, min_samples_split=10
[CV] max_depth=1, min_samples_split=10, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=100
[CV] max_depth=1, min_samples_split=100, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=100
[CV] max_depth=1, min_samples_split=100, score=0.5616270486654009, total= 1.5s
[CV] max_depth=1, min_samples_split=100
[CV] max_depth=1, min_samples_split=100, score=0.5659347823809485, total= 1.5s
[CV] max_depth=1, min_samples_split=500
[CV] max_depth=1, min_samples_split=500, score=0.5661502566606066, total= 1.5s
[CV] max_depth=1, min_samples_split=500
[CV] max_depth=1, min_samples_split=500, score=0.5616270486654009, total= 1.6s
[CV] max_depth=1, min_samples_split=500
[CV] max_depth=1, min_samples_split=500, score=0.5659347823809485, total= 1.6s
[CV] max_depth=5, min_samples_split=5
[CV] max_depth=5, min_samples_split=5, score=0.6458938602734846, total= 7.0s
[CV] max_depth=5, min_samples_split=5
[CV] max_depth=5, min_samples_split=5, score=0.6427663701278759, total= 7.0s
[CV] max_depth=5, min_samples_split=5
[CV] max_depth=5, min_samples_split=5, score=0.6407767202191855, total= 7.2s
[CV] max_depth=5, min_samples_split=10
[CV] max_depth=5, min_samples_split=10, score=0.6458938602734846, total= 7.2s
[CV] max_depth=5, min_samples_split=10
[CV] max_depth=5, min_samples_split=10, score=0.6427663701278759, total= 7.0s
[CV] max_depth=5, min_samples_split=10
[CV] max_depth=5, min_samples_split=10, score=0.6407767202191855, total= 7.0s
[CV] max_depth=5, min_samples_split=100
[CV] max_depth=5, min_samples_split=100, score=0.6454716075552328, total= 6.9s
[CV] max_depth=5, min_samples_split=100
[CV] max_depth=5, min_samples_split=100, score=0.6425567483806119, total= 7.0s
[CV] max_depth=5, min_samples_split=100
[CV] max_depth=5, min_samples_split=100, score=0.6407767202191855, total= 7.0s
[CV] max_depth=5, min_samples_split=500
[CV] max_depth=5, min_samples_split=500, score=0.6487388627418185, total= 6.9s
[CV] max_depth=5, min_samples_split=500
[CV] max_depth=5, min_samples_split=500, score=0.6423604324882733, total= 7.1s
[CV] max_depth=5, min_samples_split=500
[CV] max_depth=5, min_samples_split=500, score=0.6426904029509427, total= 6.9s
[CV] max_depth=10, min_samples_split=5
[CV] max_depth=10, min_samples_split=5, score=0.59343029982311, total= 11.8s
[CV] max_depth=10, min_samples_split=5
[CV] max_depth=10, min_samples_split=5, score=0.6174409967911388, total= 11.8s
[CV] max_depth=10, min_samples_split=5
[CV] max_depth=10, min_samples_split=5, score=0.6079518356587538, total= 11.8s
[CV] max_depth=10, min_samples_split=10
[CV] max_depth=10, min_samples_split=10, score=0.5953621558618536, total= 11.3s
[CV] max_depth=10, min_samples_split=10
[CV] max_depth=10, min_samples_split=10, score=0.6166696623017646, total= 11.8s
[CV] max_depth=10, min_samples_split=10
[CV] max_depth=10, min_samples_split=10, score=0.6074256789183601, total= 11.8s

[CV] max_depth=10, min_samples_split=100
[CV] max_depth=10, min_samples_split=100, score=0.6133191132543359, total= 11.5s
[CV] max_depth=10, min_samples_split=100
[CV] max_depth=10, min_samples_split=100, score=0.6243918569912889, total= 11.6s
[CV] max_depth=10, min_samples_split=100
[CV] max_depth=10, min_samples_split=100, score=0.6185499575457767, total= 11.8s
[CV] max_depth=10, min_samples_split=500
[CV] max_depth=10, min_samples_split=500, score=0.6233156160845382, total= 10.1s
[CV] max_depth=10, min_samples_split=500
[CV] max_depth=10, min_samples_split=500, score=0.6392726951143144, total= 10.1s
[CV] max_depth=10, min_samples_split=500
[CV] max_depth=10, min_samples_split=500, score=0.6388037025969413, total= 10.1s
[CV] max_depth=50, min_samples_split=5
[CV] max_depth=50, min_samples_split=5, score=0.5369379835638006, total= 16.9s
[CV] max_depth=50, min_samples_split=5
[CV] max_depth=50, min_samples_split=5, score=0.5354430270090781, total= 16.3s
[CV] max_depth=50, min_samples_split=5
[CV] max_depth=50, min_samples_split=5, score=0.5304583172862023, total= 17.2s
[CV] max_depth=50, min_samples_split=10
[CV] max_depth=50, min_samples_split=10, score=0.5333242206676753, total= 17.0s
[CV] max_depth=50, min_samples_split=10
[CV] max_depth=50, min_samples_split=10, score=0.5395232814444462, total= 16.6s
[CV] max_depth=50, min_samples_split=10
[CV] max_depth=50, min_samples_split=10, score=0.532944333832873, total= 17.2s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.5768635172779006, total= 15.6s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.5638113901608219, total= 14.4s
[CV] max_depth=50, min_samples_split=100
[CV] max_depth=50, min_samples_split=100, score=0.5686729176477866, total= 15.3s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6096686151195139, total= 10.7s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6381750399861196, total= 10.1s
[CV] max_depth=50, min_samples_split=500
[CV] max_depth=50, min_samples_split=500, score=0.6345424721827557, total= 10.4s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5318054880639452, total= 17.7s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5405477076286735, total= 16.0s
[CV] max_depth=100, min_samples_split=5
[CV] max_depth=100, min_samples_split=5, score=0.5302629143651922, total= 16.3s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5334821603834972, total= 17.0s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.540635395393918, total= 15.6s
[CV] max_depth=100, min_samples_split=10
[CV] max_depth=100, min_samples_split=10, score=0.5343180472354367, total= 16.0s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5713994387329128, total= 15.2s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5664217994651233, total= 14.3s
[CV] max_depth=100, min_samples_split=100
[CV] max_depth=100, min_samples_split=100, score=0.5640484442098626, total= 15.8s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6096686151195139, total= 11.3s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6381750399861196, total= 10.2s
[CV] max_depth=100, min_samples_split=500
[CV] max_depth=100, min_samples_split=500, score=0.6345424721827557, total= 10.8s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5318054880639452, total= 17.5s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5405477076286735, total= 16.7s
[CV] max_depth=500, min_samples_split=5
[CV] max_depth=500, min_samples_split=5, score=0.5302629143651922, total= 17.1s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5334821603834972, total= 17.0s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.540635395393918, total= 15.7s
[CV] max_depth=500, min_samples_split=10
[CV] max_depth=500, min_samples_split=10, score=0.5343180472354367, total= 16.7s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5713994387329128, total= 16.2s
[CV] max_depth=500, min_samples_split=100
[CV] max_depth=500, min_samples_split=100, score=0.5664217994651233, total= 14.7s
[CV] max_depth=500, min_samples_split=100

```
[CV] max_depth=500, min_samples_split=100, score=0.5640484442098626, total= 14.7s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6096686151195139, total= 11.1s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6381750399861196, total= 9.8s
[CV] max_depth=500, min_samples_split=500 .....
[CV] max_depth=500, min_samples_split=500, score=0.6345424721827557, total= 10.5s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5318054880639452, total= 17.3s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5405477076286735, total= 16.4s
[CV] max_depth=1000, min_samples_split=5 .....
[CV] max_depth=1000, min_samples_split=5, score=0.5302629143651922, total= 15.9s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.5334821603834972, total= 16.4s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.540635395393918, total= 15.4s
[CV] max_depth=1000, min_samples_split=10 .....
[CV] max_depth=1000, min_samples_split=10, score=0.5343180472354367, total= 15.7s
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.5713994387329128, total= 14.9s
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.5664217994651233, total= 14.0s
[CV] max_depth=1000, min_samples_split=100 .....
[CV] max_depth=1000, min_samples_split=100, score=0.5640484442098626, total= 14.4s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6096686151195139, total= 10.2s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6381750399861196, total= 9.3s
[CV] max_depth=1000, min_samples_split=500 .....
[CV] max_depth=1000, min_samples_split=500, score=0.6345424721827557, total= 10.0s
```

```
[Parallel(n_jobs=1)]: Done 84 out of 84 | elapsed: 15.9min finished
```

Out[115]:

```
GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=DecisionTreeClassifier(class_weight='balanced', 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=1,
                                              splitter='best'),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='roc_auc', verbose=3)
```

Best Params for TFIDFW2V Feature

In [116]:

```
best_d = clf.best_params_['max_depth']
best_split = clf.best_params_['min_samples_split']
best_d, best_split
```

Out[116]:

```
(5, 500)
```

In [117]:

```
# Plot seaborn heatmap for gridsearchcv: https://www.kaggle.com/arindambanerjee/grid-search-simplified
max_depth_list = list(clf.cv_results_['param_max_depth'].data)
samplesplit_list = list(clf.cv_results_['param_min_samples_split'].data)

plt.figure(1, figsize=(10,10))
plt.subplot(211)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
                          'Min Sample Split':samplesplit_list, \
```

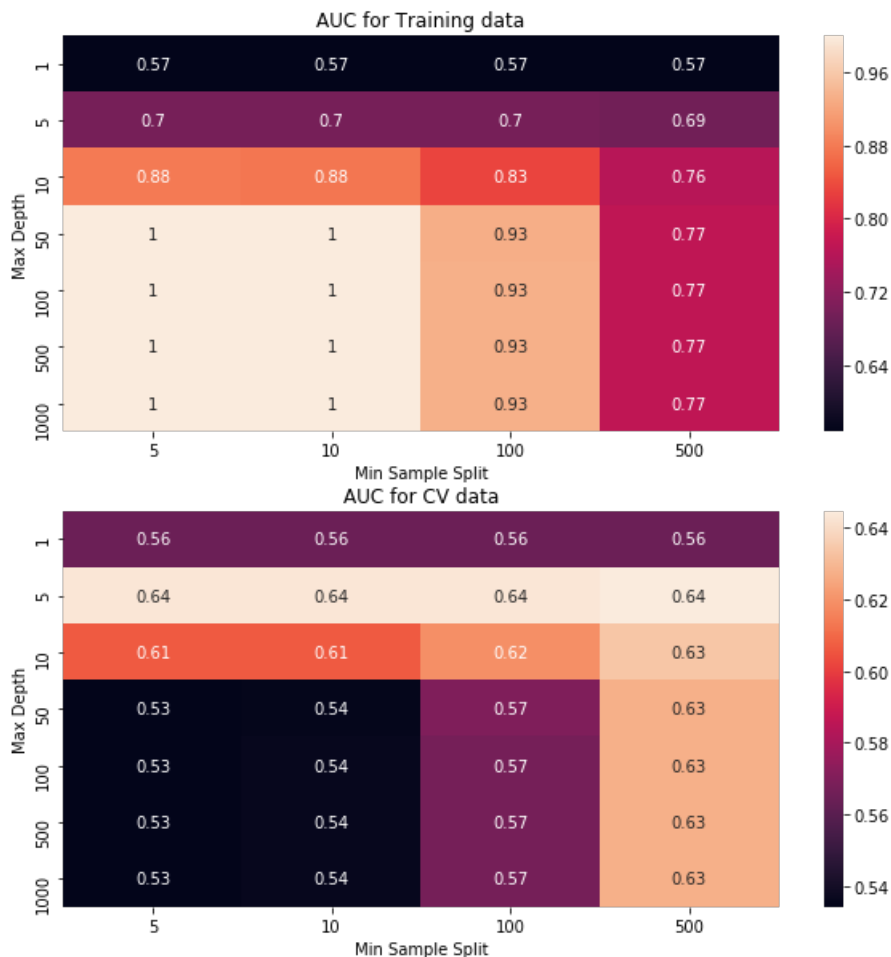
```

        'Min Sample Split':samplesplit_list, \
        'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')

plt.subplot(212)
data = pd.DataFrame(data={'Max Depth':max_depth_list, \
        'Min Sample Split':samplesplit_list, \
        'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='Max Depth', columns='Min Sample Split', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')

plt.show()

```



In [118]:

```

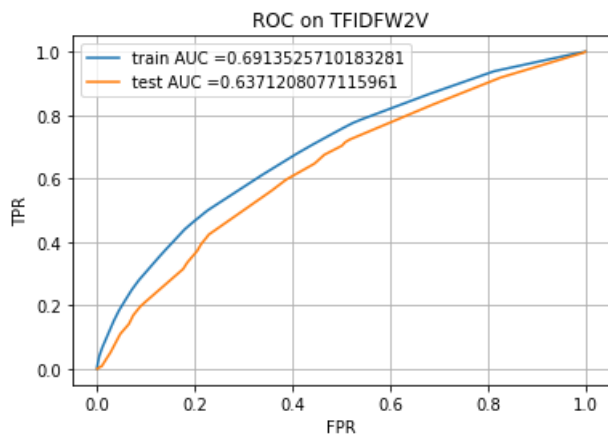
lr = DecisionTreeClassifier(max_depth=best_d, class_weight='balanced', min_samples_split=best_split, random_state=1)
lr.fit(tr_X, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive classes
# not the predicted outputs

y_train_pred = lr.predict_proba(tr_X)[:,1]
y_test_pred = lr.predict_proba(ts_X)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("FPR")
plt.ylabel("TPR")
plt.title("ROC on TFIDFW2V")
plt.grid()
plt.show()

```



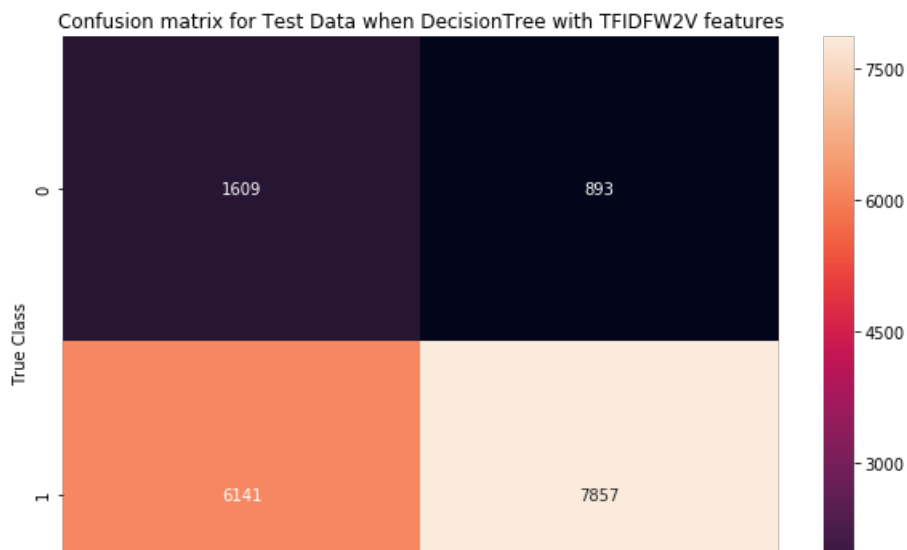
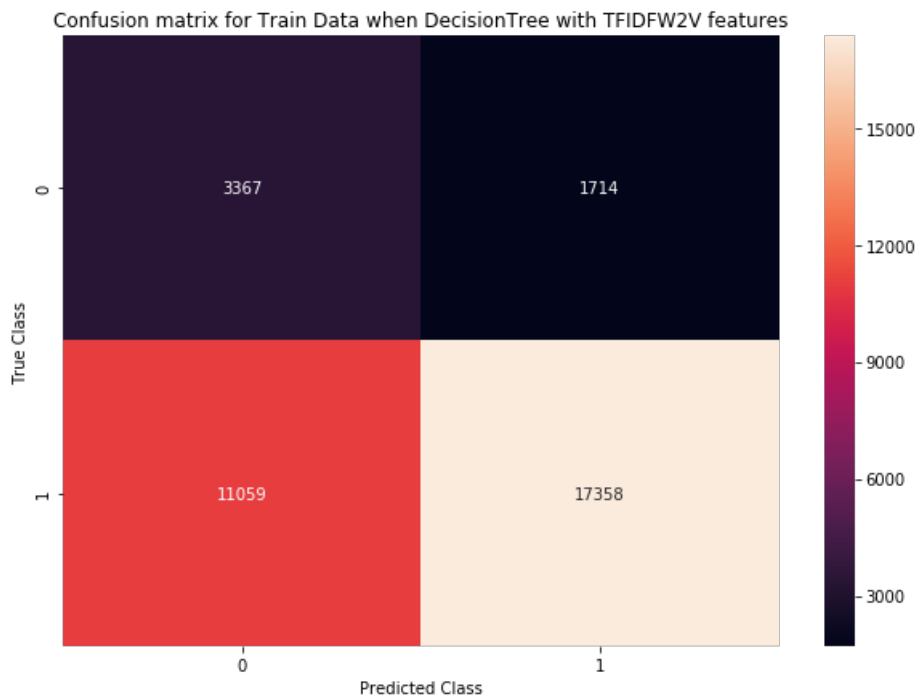
In [119]:

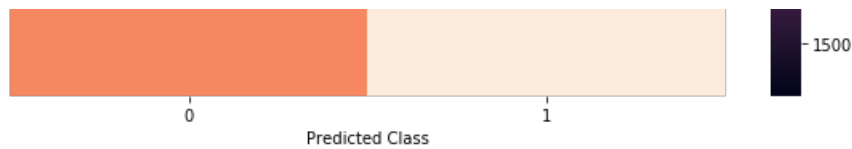
```
plot_cm('TFIDFW2V', tr_thresholds, train_fpr, train_tpr, tr_y, y_train_pred, ts_y, y_test_pred)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4047765814455433 for threshold 0.5

Train confusion matrix

Test confusion matrix





In [121]:

```
# Predict test data with best threshold value
ts_predict = predict_with_best_t(lr.predict_proba(ts_X)[: ,1], 0.5)

false_datapoints = []

# Iterate over all data points in test data
for i in range(ts_X.shape[0]):
    # Select that data point where test true value is 0 and test predicted value is 1
    if (ts_y[i] == 0) and (ts_predict[i] == 1):
        # If it true, then put that datapoint into another array
        false_datapoints.append(ts_text.iloc[i].values)
```

In [122]:

```
false_datapoints = pd.DataFrame(data=false_datapoints, columns=ts_text.columns)
false_datapoints.shape
```

Out[122]:

 $(893, 10)$

In [123]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# WordCloud on Essay
from wordcloud import WordCloud

comment_words = ' '

for val in false_datapoints['clean_essay']:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()

    # Converts each token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

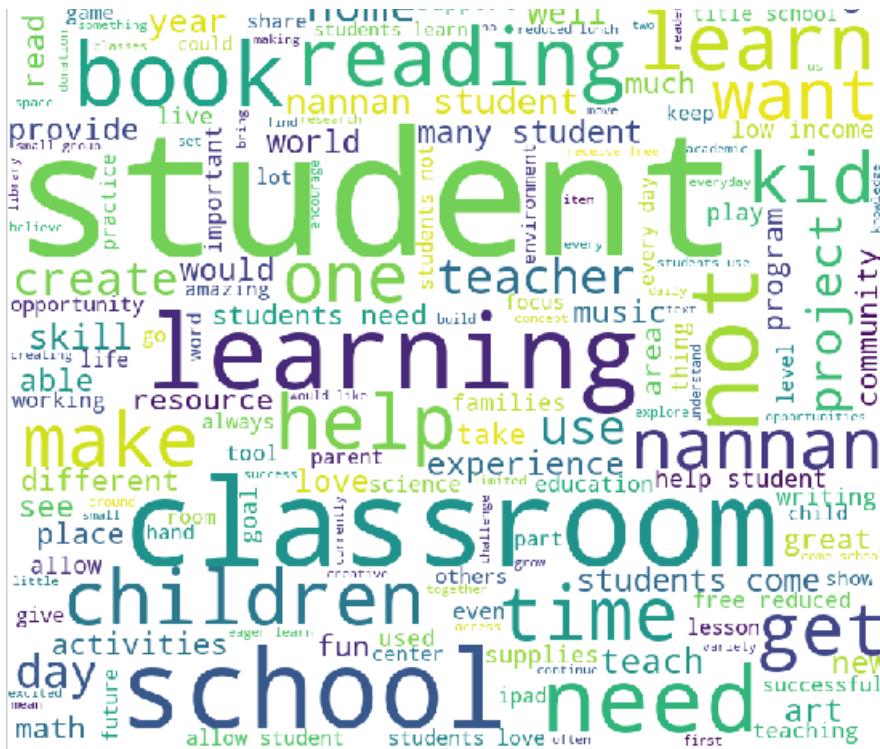
    for words in tokens:
        comment_words = comment_words + words + ' '

wordcloud = WordCloud(width = 800, height = 800,
                        background_color='white',
                        stopwords = stopwords,
                        min_font_size = 10).generate(comment_words)

# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)

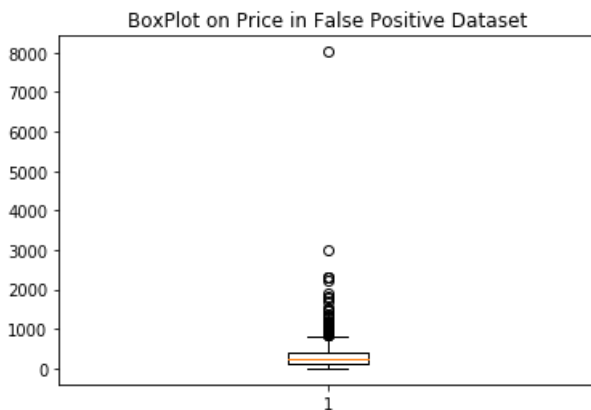
plt.show()
```





In [124]:

```
# BoxPlot on test price
plt.boxplot(false_datapoints['price'])
plt.title('BoxPlot on Price in False Positive Dataset')
plt.show()
```



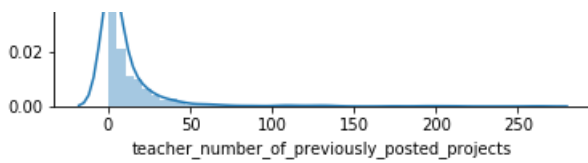
In [125]:

```
# PDF on teacher_number_of_previously_posted_projects
sns.distplot(false_datapoints['teacher number of previously posted projects'], hist=True, kde=True)
```

Out[125]:

```
<matplotlib.axes. subplots.AxesSubplot at 0x1f6481ac0b8>
```





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

In [0]:

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

In [324]:

```
# Before proceed here, get compile and run again in [SET 2] on TFIDF Feature
```

In [171]:

```
# taking top 5k feature importance
feature_imp5k = np.flip(np.argsort(lr.feature_importances_))[:5000]
```

In [172]:

```
# Dimension of this variable
feature_imp5k.shape
```

Out[172]:

```
(5000,)
```

In [173]:

```
# print the feature importance column number
feature_imp5k
```

Out[173]:

```
array([ 1,  0,  2, ..., 3857, 3858, 3859], dtype=int64)
```

In [174]:

```
tr_5k = tr_X[:,feature_imp5k]
ts_5k = ts_X[:,feature_imp5k]
```

In [175]:

```
tr_5k.shape, ts_5k.shape
```

Out[175]:

```
((33498, 5000), (16500, 5000))
```

In [177]:

```
# Training the Logistic Regression Model for 5k feature importance
```



```
# Training the Logistic Regression Model for Feature Importance
from sklearn.linear_model import LogisticRegression
```

In [178]:

```
clf = LogisticRegression(class_weight='balanced', random_state=1)
parameters = {'C':[10**-4,10**-3,10**-2,10**-1,1,10,100,10**3,10**4],
              'penalty':['l1','l2']}
```

In [179]:

```
clf = GridSearchCV(clf, parameters, cv=3, scoring='roc_auc', verbose=3)
clf.fit(tr_5k, tr_y)
```

Fitting 3 folds for each of 18 candidates, totalling 54 fits

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

```
[CV] C=0.0001, penalty=l1 .....
[CV] ..... C=0.0001, penalty=l1, score=0.5, total= 5.0s
```

[Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 5.3s remaining: 0.0s

```
[CV] C=0.0001, penalty=l1 .....
[CV] ..... C=0.0001, penalty=l1, score=0.5, total= 4.8s
```

[Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 10.5s remaining: 0.0s

```
[CV] C=0.0001, penalty=l1 .....
[CV] ..... C=0.0001, penalty=l1, score=0.5, total= 4.8s
[CV] C=0.0001, penalty=l2 .....
[CV] ... C=0.0001, penalty=l2, score=0.5450458776082798, total= 4.8s
[CV] C=0.0001, penalty=l2 .....
[CV] ... C=0.0001, penalty=l2, score=0.5579748251978365, total= 4.9s
[CV] C=0.0001, penalty=l2 .....
[CV] ... C=0.0001, penalty=l2, score=0.5698773566833224, total= 4.8s
[CV] C=0.001, penalty=l1 .....
[CV] ..... C=0.001, penalty=l1, score=0.5, total= 4.8s
[CV] C=0.001, penalty=l1 .....
[CV] ..... C=0.001, penalty=l1, score=0.5, total= 4.7s
[CV] C=0.001, penalty=l1 .....
[CV] ..... C=0.001, penalty=l1, score=0.5, total= 4.8s
[CV] C=0.001, penalty=l2 .....
[CV] .... C=0.001, penalty=l2, score=0.557766178429691, total= 4.8s
[CV] C=0.001, penalty=l2 .....
[CV] .... C=0.001, penalty=l2, score=0.5734346082357447, total= 4.9s
[CV] C=0.001, penalty=l2 .....
[CV] .... C=0.001, penalty=l2, score=0.5858172088767739, total= 4.9s
[CV] C=0.01, penalty=l1 .....
[CV] .... C=0.01, penalty=l1, score=0.5267530062137703, total= 4.9s
[CV] C=0.01, penalty=l1 .....
[CV] .... C=0.01, penalty=l1, score=0.5402202651847539, total= 4.8s
[CV] C=0.01, penalty=l1 .....
[CV] .... C=0.01, penalty=l1, score=0.5419946662828659, total= 4.9s
[CV] C=0.01, penalty=l2 .....
[CV] .... C=0.01, penalty=l2, score=0.6163265110272395, total= 5.0s
[CV] C=0.01, penalty=l2 .....
[CV] .... C=0.01, penalty=l2, score=0.6305110545167363, total= 5.0s
[CV] C=0.01, penalty=l2 .....
[CV] ..... C=0.01, penalty=l2, score=0.638104810547405, total= 5.1s
[CV] C=0.1, penalty=l1 .....
[CV] .... C=0.1, penalty=l1, score=0.6514270471810082, total= 4.8s
[CV] C=0.1, penalty=l1 .....
[CV] .... C=0.1, penalty=l1, score=0.6528906299858004, total= 4.8s
[CV] C=0.1, penalty=l1 .....
[CV] .... C=0.1, penalty=l1, score=0.660906869103303, total= 4.9s
[CV] C=0.1, penalty=l2 .....
[CV] .... C=0.1, penalty=l2, score=0.674788633724557, total= 5.4s
[CV] C=0.1, penalty=l2 .....
```

```

[CV] C=0.1, penalty=l2 .....
[CV] ..... C=0.1, penalty=l2, score=0.6768596163127414, total= 5.3s
[CV] C=0.1, penalty=l2 .....
[CV] ..... C=0.1, penalty=l2, score=0.6824521379767245, total= 5.5s
[CV] C=1, penalty=l1 .....
[CV] ..... C=1, penalty=l1, score=0.6776913718988324, total= 8.6s
[CV] C=1, penalty=l1 .....
[CV] ..... C=1, penalty=l1, score=0.6841549641620664, total= 6.7s
[CV] C=1, penalty=l1 .....
[CV] ..... C=1, penalty=l1, score=0.6854099651186922, total= 5.7s
[CV] C=1, penalty=l2 .....
[CV] ..... C=1, penalty=l2, score=0.6678840913795762, total= 6.1s
[CV] C=1, penalty=l2 .....
[CV] ..... C=1, penalty=l2, score=0.6725758788969016, total= 5.9s
[CV] C=1, penalty=l2 .....
[CV] ..... C=1, penalty=l2, score=0.6694181676138631, total= 6.2s
[CV] C=10, penalty=l1 .....
[CV] ..... C=10, penalty=l1, score=0.6335105639828152, total= 7.0s
[CV] C=10, penalty=l1 .....
[CV] ..... C=10, penalty=l1, score=0.6391467724919428, total= 13.9s
[CV] C=10, penalty=l1 .....
[CV] ..... C=10, penalty=l1, score=0.6260060428672914, total= 7.8s
[CV] C=10, penalty=l2 .....
[CV] ..... C=10, penalty=l2, score=0.6425674984305734, total= 7.3s
[CV] C=10, penalty=l2 .....
[CV] ..... C=10, penalty=l2, score=0.6463737525527299, total= 7.9s
[CV] C=10, penalty=l2 .....
[CV] ..... C=10, penalty=l2, score=0.6396465199509905, total= 7.9s
[CV] C=100, penalty=l1 .....
[CV] ..... C=100, penalty=l1, score=0.6190741448603506, total= 7.4s
[CV] C=100, penalty=l1 .....
[CV] ..... C=100, penalty=l1, score=0.6221203886331728, total= 17.4s
[CV] C=100, penalty=l1 .....
[CV] ..... C=100, penalty=l1, score=0.6108723095696109, total= 7.1s
[CV] C=100, penalty=l2 .....
[CV] ..... C=100, penalty=l2, score=0.6259674080226272, total= 11.2s
[CV] C=100, penalty=l2 .....
[CV] ..... C=100, penalty=l2, score=0.6288653041138836, total= 11.1s
[CV] C=100, penalty=l2 .....
[CV] ..... C=100, penalty=l2, score=0.6188877891476829, total= 11.1s
[CV] C=1000, penalty=l1 .....
[CV] ..... C=1000, penalty=l1, score=0.6162641951006969, total= 10.9s
[CV] C=1000, penalty=l1 .....
[CV] ..... C=1000, penalty=l1, score=0.6180372673625515, total= 7.1s
[CV] C=1000, penalty=l1 .....
[CV] ..... C=1000, penalty=l1, score=0.6076238880086524, total= 7.5s
[CV] C=1000, penalty=l2 .....
[CV] ..... C=1000, penalty=l2, score=0.6186892816980243, total= 19.4s
[CV] C=1000, penalty=l2 .....
[CV] ..... C=1000, penalty=l2, score=0.6212660717277194, total= 15.8s
[CV] C=1000, penalty=l2 .....
[CV] ..... C=1000, penalty=l2, score=0.6102528321107581, total= 20.6s
[CV] C=10000, penalty=l1 .....
[CV] .... C=10000, penalty=l1, score=0.6152875175839966, total= 11.4s
[CV] C=10000, penalty=l1 .....
[CV] .... C=10000, penalty=l1, score=0.6161522608610677, total= 7.3s
[CV] C=10000, penalty=l1 .....
[CV] .... C=10000, penalty=l1, score=0.6062017837758018, total= 8.0s
[CV] C=10000, penalty=l2 .....
[CV] .... C=10000, penalty=l2, score=0.6162376485159899, total= 29.5s
[CV] C=10000, penalty=l2 .....
[CV] ..... C=10000, penalty=l2, score=0.618389888098695, total= 21.5s
[CV] C=10000, penalty=l2 .....
[CV] .... C=10000, penalty=l2, score=0.6076276919270126, total= 29.6s

```

```
[Parallel(n_jobs=1)]: Done 54 out of 54 | elapsed: 8.1min finished
```

```
Out[179]:
```

```

GridSearchCV(cv=3, error_score='raise-deprecating',
             estimator=LogisticRegression(C=1.0, class_weight='balanced', dual=False,
             fit_intercept=True, intercept_scaling=1, max_iter=100,
             multi_class='warn', n_jobs=None, penalty='l2', random_state=1,
             solver='warn', tol=0.0001, verbose=0, warm_start=False),
             fit_params=None, iid='warn', n_jobs=None,
             param_grid={'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000], 'penalty': ['l1', 'l2']},

```

```
pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
scoring='roc_auc', verbose=3)
```

In [180]:

```
best_C = clf.best_params_['C']
best_penalty = clf.best_params_['penalty']
best_C, best_penalty
```

Out[180]:

```
(1, 'l1')
```

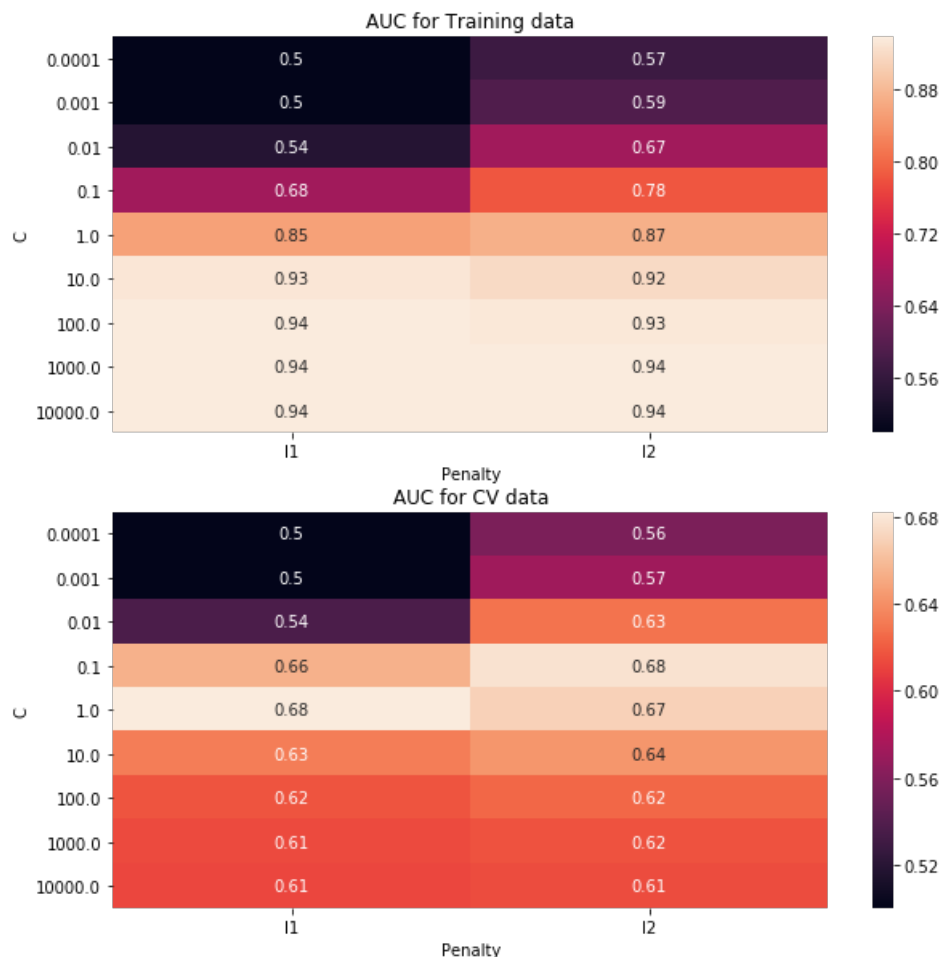
In [181]:

```
# Plot seaborn heatmap for gridsearchcv: https://www.kaggle.com/arindambanerjee/grid-search-simplified
C_list = list(clf.cv_results_['param_C'].data)
penalty_list = list(clf.cv_results_['param_penalty'].data)

plt.figure(1, figsize=(10,10))
plt.subplot(211)
data = pd.DataFrame(data={'C':C_list, \
                          'Penalty':penalty_list, \
                          'AUC':clf.cv_results_['mean_train_score']})
data = data.pivot(index='C', columns='Penalty', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for Training data')

plt.subplot(212)
data = pd.DataFrame(data={'C':C_list, \
                          'Penalty':penalty_list, \
                          'AUC':clf.cv_results_['mean_test_score']})
data = data.pivot(index='C', columns='Penalty', values='AUC')
sns.heatmap(data, annot=True).set_title('AUC for CV data')

plt.show()
```



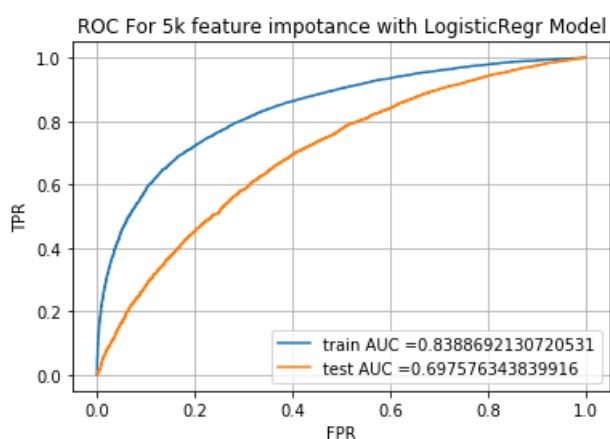
In [182]:

```
lr = LogisticRegression(C=best_C, class_weight='balanced', penalty=best_penalty, random_state=1)
lr.fit(tr_5k, tr_y)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

# Since LinearSVC doesnt have predict_proba attribute: https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html
y_train_pred = lr.predict_proba(tr_5k)[:,1]
y_test_pred = lr.predict_proba(ts_5k)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(tr_y, y_train_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(ts_y, 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("FPR")
plt.ylabel("TPR")
plt.title("ROC For 5k feature impotence with LogisticRegr Model")
plt.grid()
plt.show()
```



In [191]:

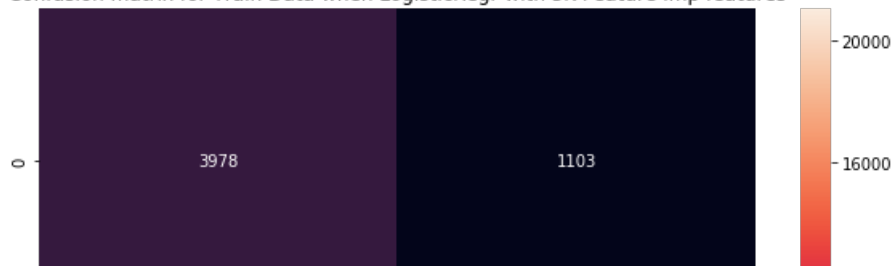
```
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
cm = metrics.confusion_matrix(tr_y, predict_with_best_t(y_train_pred, best_t))
plt.figure(figsize = (10,7))
sns.heatmap(cm, annot=True, fmt="d")
plt.xlabel('Predicted Class')
plt.ylabel('True Class')
plt.title('Confusion matrix for Train Data when LogisticRegr with 5K Feature Imp features')
```

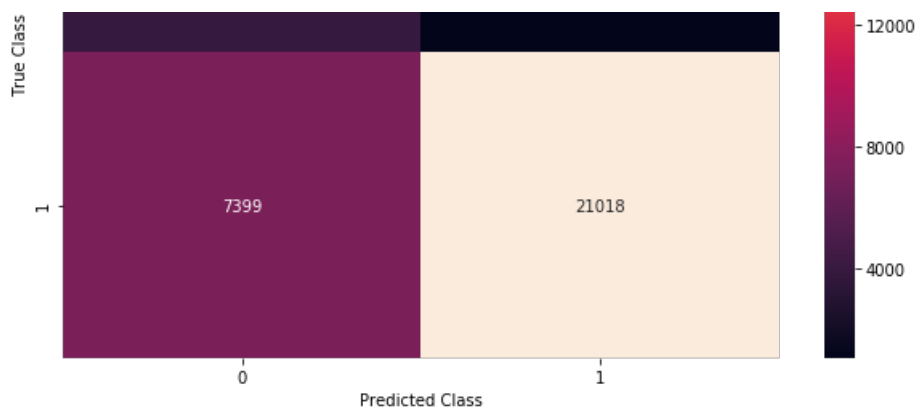
the maximum value of $tpr \cdot (1 - fpr)$ 0.5790669044437498 for threshold 0.498
Train confusion matrix

Out[191]:

Text(0.5, 1.0, 'Confusion matrix for Train Data when LogisticRegr with 5K Feature Imp features')

Confusion matrix for Train Data when LogisticRegr with 5K Feature Imp features



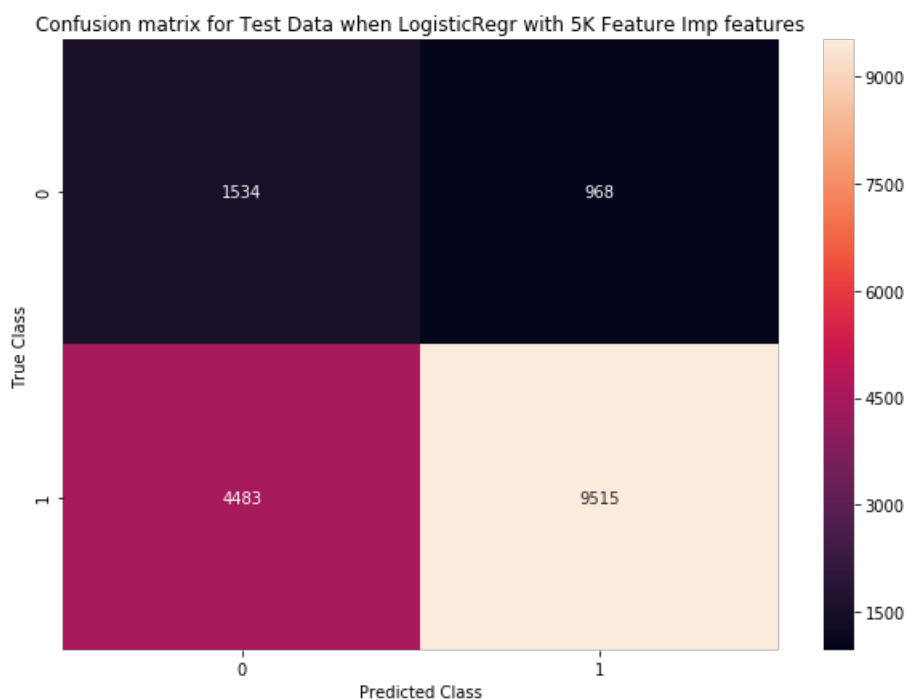


In [192]:

```
cm = metrics.confusion_matrix(ts_y, predict_with_best_t(y_test_pred, 0.498))
plt.figure(figsize = (10,7))
sns.heatmap(cm, annot=True, fmt="d")
plt.xlabel('Predicted Class')
plt.ylabel('True Class')
plt.title('Confusion matrix for Test Data when LogisticRegr with 5K Feature Imp features')
```

Out[192]:

Text(0.5, 1.0, 'Confusion matrix for Test Data when LogisticRegr with 5K Feature Imp features')



In [184]:

```
# Predict test data with best threshold value
ts_predict = predict_with_best_t(lr.predict_proba(ts_5k[:,1], 0.498)

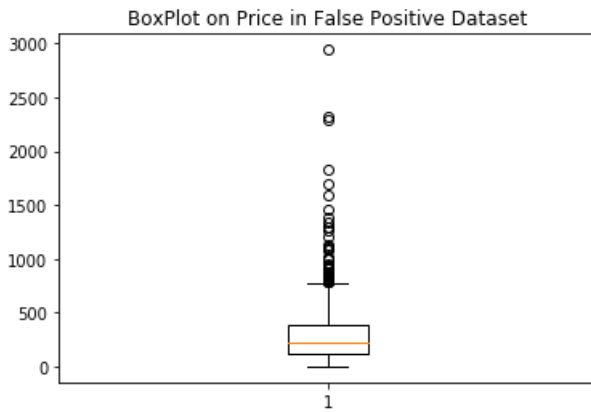
false_datapoints = []

# Iterate over all data points in test data
for i in range(ts_5k.shape[0]):
    # Select that data point where test true value is 0 and test predicted value is 1
    if (ts_y[i] == 0) and (ts_predict[i] == 1):
        # If it true, then put that datapoint into another array
        false_datapoints.append(ts_text.iloc[i].values)
```

In [185]:

In [194]:

```
# BoxPlot on test price
plt.boxplot(false_datapoints['price'])
plt.title('BoxPlot on Price in False Positive Dataset')
plt.show()
```

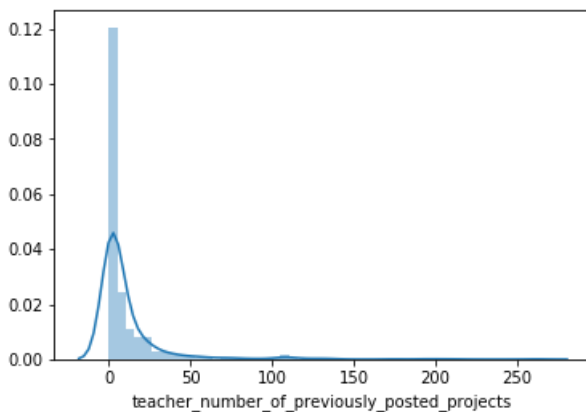


In [195]:

```
# PDF on teacher_number_of_previously_posted_projects
sns.distplot(false_datapoints['teacher_number_of_previously_posted_projects'], hist=True, kde=True)
```

Out[195]:

<matplotlib.axes._subplots.AxesSubplot at 0x1f647c99ba8>



3. Conclusion

In [196]:

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

x = PrettyTable()

x.field_names = ["Features", "Model", "max_depth", "min_sample_split", "Maximum AUC score",]

x.add_row(["BoW", "DT", 10, 500, 0.66651])
x.add_row(["TFIDF", "DT", 10, 500, 0.6676])
x.add_row(["AVG W2V", "DT", 5, 500, 0.64207])
```

```
x.add_row(["TFIDF W2V","DT", 5,500, 0.63712])
print(x)

y = PrettyTable()

y.field_names = ["Features", "Model", "C", "penalty", "Maximum AUC score",]
y.add_row(["5k Feature Importance from DT","LogisticRegr", 1,"l1", 0.69757])
print(y)
```

```
+-----+-----+-----+-----+-----+
| Features | Model | max_depth | min_sample_split | Maximum AUC score |
+-----+-----+-----+-----+-----+
| BoW      | DT    | 10        | 500               | 0.66651           |
| TFIDF    | DT    | 10        | 500               | 0.6676            |
| AVG W2V  | DT    | 5         | 500               | 0.64207           |
| TFIDF W2V | DT    | 5         | 500               | 0.63712           |
+-----+-----+-----+-----+-----+
```

```
+-----+-----+-----+-----+-----+
| Features | Model | C | penalty | Maximum AUC score |
+-----+-----+-----+-----+-----+
| 5k Feature Importance from DT | LogisticRegr | 1 | l1 | 0.69757 |
+-----+-----+-----+-----+-----+
```

In []: