

# DonorsChoose

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

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

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

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

## About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> 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  <b>Examples:</b> Music & The Arts Literacy & Language, Math & Science
<code>school_state</code>	State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> Literacy Literature & Writing, Social Sciences
<code>project_resource_summary</code>	An explanation of the resources needed for the project. <b>Example:</b> My students need hands on literacy materials to manage sensory needs!
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*

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

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

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

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

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

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

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

## Notes on the Essay Data

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

- \_\_project\_essay\_1\_\_: "Introduce us to your classroom"
- \_\_project\_essay\_2\_\_: "Tell us more about your students"
- \_\_project\_essay\_3\_\_: "Describe how your students will use the materials you're requesting"
- \_\_project\_essay\_3\_\_: "Close by sharing why your project will make a difference"

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

- \_\_project\_essay\_1\_\_: "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- \_\_project\_essay\_2\_\_: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

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

In [1]:

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

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

```

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

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

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

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

from tqdm import tqdm
import os

from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
warnings.filterwarnings("ignore")

```

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

## 1.1 Reading Data

In [2]:

```

project_data = pd.read_csv('E:\\Machine Learning\\Dataset\\train_data.csv')
resource_data = pd.read_csv('E:\\Machine Learning\\Dataset\\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]:

```

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

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

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

```

```
project_data.head(2)
project_data.project_is_approved.value_counts()
```

Out[4]:

```
1    92706
0    16542
Name: project_is_approved, dtype: int64
```

In [5]:

```
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[5]:

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

## 1.2 Preprocessing of project\_subject\_categories

In [6]:

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

## 1.4 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	Date	project_grade_category	project_title
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2 Engineer STEAM the Print Classroom
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5 Sensory Tools for

In [10]:

```
# printing some random reviews
```

```
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
```

I have been fortunate enough to use the Fairy Tale STEM kits in my classroom as well as the STEM journals, which my students really enjoyed. I would love to implement more of the Lakeshore STEM kits in my classroom for the next school year as they provide excellent and engaging STEM lessons. My students come from a variety of backgrounds, including language and socioeconomic status. Many of them don't have a lot of experience in science and engineering and these kits give me the materials to provide these exciting opportunities for my students. Each month I try to do several science or STEM/STEAM projects. I would use the kits and robot to help guide my science instruction in engaging and meaningful ways. I can adapt the kits to my current language arts pacing guide where we already teach some of the material in the kits like tall tales (Paul Bunyan) or Johnny Appleseed. The following units will be taught in the next school year where I will implement these kits: magnets, motion, sink vs. float, robots. I often get to these units and don't know if I am teaching the right way or using the right materials. The kits will give me additional ideas, strategies, and lessons to prepare my students in science. It is challenging to develop high quality science activities. These kits give me the materials I need to provide my students with science activities that will go along with the curriculum in my classroom. Although I have some things (like magnets) in my classroom, I don't know how to use them effectively. The kits will provide me with the right amount of materials and show me how to use them in an appropriate way.

I teach high school English to students with learning and behavioral disabilities. My students all vary in their ability level. However, the ultimate goal is to increase all students literacy levels. This includes their reading, writing, and communication levels. I teach a really dynamic group of students. However, my students face a lot of challenges. My students all live in poverty and in a dangerous neighborhood. Despite these challenges, I have students who have the desire to defeat these challenges. My students all have learning disabilities and currently all are performing below grade level. My students are visual learners and will benefit from a classroom that fulfills their preferred learning style. The materials I am requesting will allow my students to be prepared for the classroom with the necessary supplies. Too often I am challenged with students who come to school unprepared for class due to economic challenges. I want my students to be able to focus on learning and not how they will be able to get school supplies. The supplies will last all year. Students will be able to complete written assignments and maintain a classroom journal. The chart paper will be used to make learning more visual in class and to create posters to aid students in their learning. The students have access to a classroom printer. The toner will be used to print student work that is completed on the classroom Chromebooks. I want to try and remove all barriers for the students learning and create opportunities for learning. One of the biggest barriers is the students not having the resources to get pens, paper, and folders. My students will be able to increase their literacy skills because of this project.

"Life moves pretty fast. If you don't stop and look around once in awhile, you could miss it." from the movie, Ferris Bueller's Day Off. Think back...what do you remember about your grandparents? How amazing would it be to be able to flip through a book to see a day in their lives? My second graders are voracious readers! They love to read both fiction and nonfiction books. Their favorite characters include Pete the Cat, Fly Guy, Piggie and Elephant, and Mercy Watson. They also love to read about insects, space and plants. My students are hungry bookworms! My students are eager to learn and read about the world around them. My kids love to be at school and are like little sponges absorbing everything around them. Their parents work long hours and usually do not see their children. My students are usually cared for by their grandparents or a family friend. Most of my students do not have someone who speaks English at home. Thus it is difficult for my students to acquire language. Now think forward... wouldn't it mean a lot to your kids, nieces or nephews or grandchildren, to be able to see a day in your life today 30 years from now? Memories are so precious to us and being able to share these memories with future generations will be a rewarding experience. As part of our social studies curriculum, students will be learning about changes over time. Students will be studying photos to learn about how their community has changed over time. In particular, we will look at photos to study how the land, buildings, clothing, and schools have changed over time. As a culminating activity, my students will capture a slice of their history and preserve it through scrap booking. Key important events in their young lives will be documented with the date, location, and names. Students will be using photos from home and from school to create their second grade memories. Their scrap books will preserve their unique stories for future generations to enjoy. Your donation to this project will provide my second graders with an opportunity to learn about social studies in a fun and creative manner. Through their scrapbooks, children will share their story with others and have a historical document for the rest of their lives.

In [11]:

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

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

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"'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 [12]:

```
sent = decontracted(project_data['essay'].values[2000])
print(sent)
print("="*50)
```

```
\ "Creativity is intelligence having fun.\" --Albert Einstein. Our elementary library at Greenville
Elementary is anything but a quiet, hushed space. It is a place for collaboration and research. It
is a place for incorporating technology. It is a place for innovation. And it is a place for creat
ing.Our school serves 350 third and fourth graders who primarily live in rural and poverty-stricke
n areas in our community. Being a Title I school, approximately 85% of them receive free or
reduced lunch. But they are inquisitive, creative, and eager to learn. They love visiting the libr
ary to check out books, hear \r\nstories, create digital stories, and use the computer lab for lea
rning and fun. We want to build our library is Makerspace with activities revolving around art and
literacy to provide more engaging, hands-on activities.We want to begin \"Makerspace Fridays!\"
Our school recently received a $1000 grant for books for our arts-integrated Makerspace. We have r
eceived titles such as \"Origami for Everyone,\" \"How to Make Stuff with Ducktape,\" and \"Cool E
ngineering Activities for Girls.\" We now need supplies to correlate with these new informational
texts. By adding these art and craft supplies, students will be able to design and create
masterpieces related to their coursework. \r\n\r\nFor example, while studying Native Americans, st
udents can use the looms and yarn to recreate Navajo and/or Pueblo weaving. Weaving can also be in
tegrated with literacy through Greek mythology and the story of Arachne.\r\n\r\nCreating art with
perler beads has many possibilities! Students can design their own animals after studying their ch
aracteristics. They can use symmetry and patterning to create one-of-a-kind originals.
\r\n\r\nOrigami reinforces geometry, thinking skills, fractions, problem-solving, and just fun sci
ence!Our students need to be able to apply what they read and learn. If they read a how-to book, t
hey will apply that reading through a hands-on art activity and actually create a product. This is
a crucial skill in the real world. By creating and designing their own masterpieces, they are usin
g many critical thinking skills. Students will become more analytical thinkers.
=====
```

In [13]:

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

Creativity is intelligence having fun. --Albert Einstein. Our elementary library at Greenville Elementary is anything but a quiet, hushed space. It is a place for collaboration and research. It is a place for incorporating technology. It is a place for innovation. And it is a place for creating. Our school serves 350 third and fourth graders who primarily live in rural and poverty-stricken areas in our community. Being a Title I school, approximately 85% of them receive free or reduced lunch. But they are inquisitive, creative, and eager to learn. They love visiting the library to check out books, hear stories, create digital stories, and use the computer lab for learning and fun. We want to build our library is Makerspace with activities revolving around art and literacy to provide more engaging, hands-on activities. We want to begin Makerspace Fridays! Our school recently received a \$1000 grant for books for our arts-integrated Makerspace. We have received titles such as Origami for Everyone, How to Make Stuff with Ducktape, and Cool Engineering Activities for Girls. We now need supplies to correlate with these new informational texts. By adding these art and craft supplies, students will be able to design and create masterpieces related to their coursework. For example, while studying Native Americans, students can use the looms and yarn to recreate Navajo and/or Pueblo weaving. Weaving can also be integrated with literacy through Greek mythology and the story of Arachne. Creating art with perler beads has many possibilities! Students can design their own animals after studying their characteristics. They can use symmetry and patterning to create one-of-a-kind originals. Origami reinforces geometry,

thinking skills, fractions, problem-solving, and just fun science! Our students need to be able to apply what they read and learn. If they read a how-to book, they will apply that reading through a hands-on art activity and actually create a product. This is a crucial skill in the real world. By creating and designing their own masterpieces, they are using many critical thinking skills. Students will become more analytical thinkers.

In [14]:

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

Creativity is intelligence having fun Albert Einstein Our elementary library at Greenville Elementary is anything but a quiet hushed space It is a place for collaboration and research It is a place for incorporating technology It is a place for innovation And it is a place for creating Our school serves 350 third and fourth graders who primarily live in rural and poverty stricken areas in our community Being a Title I school approximately 85 of them receive free or reduced lunch But they are inquisitive creative and eager to learn They love visiting the library to check out books hear stories create digital stories and use the computer lab for learning and fun We want to build our library is Makerspace with activities revolving around art and literacy to provide more engaging hands on activities We want to begin Makerspace Fridays Our school recently received a 1000 grant for books for our arts integrated Makerspace We have received titles such as Origami for Everyone How to Make Stuff with Ducktape and Cool Engineering Activities for Girls We now need supplies to correlate with these new informational texts By adding these art and craft supplies students will be able to design and create masterpieces related to their coursework For example while studying Native Americans students can use the looms and yarn to recreate Navajo and or Pueblo weaving Weaving can also be integrated with literacy through Greek mythology and the story of Arachne Creating art with perler beads has many possibilities Students can design their own animals after studying their characteristics They can use symmetry and patterning to create one of a kind originals Origami reinforces geometry thinking skills fractions problem solving and just fun science Our students need to be able to apply what they read and learn If they read a how to book they will apply that reading through a hands on art activity and actually create a product This is a crucial skill in the real world By creating and designing their own masterpieces they are using many critical thinking skills Students will become more analytical thinkers

In [15]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', \
'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', \
'their', \
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", \
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', \
'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', \
while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', \
'before', 'after', \
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under' \
, 'again', 'further', \
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', '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 [16]:

```
# Combining all the above students
from tqdm import tqdm
preprocessed_essays = []
# tqdm is for printing the status bar
```



```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[01:15<00:00, 1453.92it/s]
```

```
# after preprocessing
print(preprocessed_essays[2000])
project_data['preprocessed essays']=preprocessed_essays
```

◀ ▶

```
project_data['totalwords essay'] = project_data['preprocessed essays'].str.split().str.len()
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:03<00:00, 34312.32it/s]
```

```
print(preprocessed_project_title[2000])
print("="*50)
project_data['preprocessed_project_title']=preprocessed_project_title
project_data.head(5)
```

Out[24]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2 Enginee STEAM the Prin Classr
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5 Sens Tools Fo
51140	74477	p189804	4a97f3a390bfe21b99cf5e2b81981c73	Mrs.	CA	2016-04-27 00:46:53	Grades PreK-2 Mc Learr wi Mc Lister Ce
473	100660	p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2 Flex Seating Flex Learr
41558	33679	p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5 Going De The Al Ir Think

5 rows × 22 columns



## 1.5.1 Converting Title to Number of Words

In [25]:

```
project_data['totalwords_title'] =
project_data['preprocessed_project_title'].str.split().str.len()
```

## 1.6 Preprocessing Grades

In [26]:

```
#Preprocessing grades i.e removing all the spaces from the grades and replace - by _
```

```
pre_grades = []
```

```
for grade in tqdm(project_data['project_grade_category'].values):
    sent = re.sub('[^A-Za-z0-9]+', '_', grade)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    pre_grades.append(sent.lower().strip())
```

```
project_data['grade_category']=pre_grades
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 109248/109248
[00:00<00:00, 169501.68it/s]
```

## 1.7 Preparing data for models

## 1.7 Preparing data for models

In [27]:

```
project_data.columns
```

Out[27]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_grade_category', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay',
      'preprocessed_essays', 'totalwords_essay', 'preprocessed_project_title',
      'totalwords_title', '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

## 1.8 Using Pretrained Models: Avg W2V

In [46]:

```
# 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('E:\Machine Learning\glove.42B.300d.txt')
```

Loading Glove Model

1917495it [08:44, 3657.81it/s]

Done. 1917495 words loaded!

In [47]:

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

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

```

model = set(model.keys())
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)

```

all the words in the coupus 15568853  
 the unique words in the coupus 59501  
 The number of words that are present in both glove vectors and our coupus 51613 ( 86.743 %)  
 word 2 vec length 51613

In [48]:

```

# 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 [139]:

```

import nltk
nltk.download('vader_lexicon')

```

[nltk\_data] Downloading package vader\_lexicon to C:\Users\Utkarsh  
 [nltk\_data] Sri\AppData\Roaming\nltk\_data...

Out[139]:

True

## 1.9 Calculating Sentiment score(Taking Compound in Consideration)

In [48]:

```

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 the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelligences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of different backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a caring community of successful \
learners which can be seen through collaborative student project based learning in and out of the classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice a skill before it is\

```

0.9975  
neg  
neu  
pos  
compound

```
#https://programminghistorian.org/en/lessons/sentiment-analysis
#The "neg", "neu", and "pos" values describe the fraction of weighted scores that fall into each category. VADER also sums all weighted scores to calculate a "compound" value normalized between -1 and 1
#this value attempts to describe the overall affect of the entire text from strongly negative (-1) to strongly positive (1).

sscore=[]
sid = SentimentIntensityAnalyzer()

for essay in tqdm(project_data['preprocessed_essays']):

    for_sentiment = essay
    ss = sid.polarity_scores(for_sentiment)
    sscore.append(ss['neu'])

project_data['sscore']=sscore
```

100%|██| 109248/109248  
[03:55<00:00, 464.58it/s]

```
#Combining the data from the resources from the project data and resource file for quantity and price

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

In [29]:

```
#replacing all the nan values from the teacher prefix to blank_space  
project_data.teacher_prefix=project_data.teacher_prefix.fillna('')
```

In [30]:

```
#Seprating the values of approved projects from the whole data i.e removing the target value from  
the data  
X=project_data
```

In [31]:

```
y =X['project_is_approved'].values  
X.drop(['project_is_approved'], axis=1, inplace=True)  
X.head(5)  
y.shape
```

Out[31]:

(109248,)

## Assignment 8: Apply DT

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

- **Set 1:** categorical, numerical features + project\_title(BOW) + preprocessed\_eassay (BOW)
- **Set 2:** categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (TFIDF)
- **Set 3:** categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_eassay (TFIDF W2V)




### 2. Hyper paramter tuning (best `depth` in range [1, 5, 10, 50, 100, 500, 100], and the best `min\_samples\_split` in range [5, 10, 100, 500])

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

### 3. Graphviz

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

### 4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure 
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test. 
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points 
- Once after you plot the confusion matrix with the test data, get all the `false positive data points`
  - Plot the WordCloud [WordCloud](#)
  - Plot the box plot with the `price` of these `false positive data points`
  - Plot the pdf with the `teacher\_number\_of\_previously\_posted\_projects` of these `false positive data points`

### 5. [Task-2]

- Select 5k best features from features of **Set 2** using [feature importances](#), discard all the other remaining features and then apply any of the model of your choice i.e. (Decision tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3

## 6. Conclusion

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

### Note: Data Leakage

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

## 2 Decision Tree

### 2.1 Splitting of data

In [32]:

```
#splitting of data using train_test_split

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, stratify=y)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.25, stratify=y_train)

print(X_train.shape, y_train.shape)
print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)
print(X_train.columns)

(61452, 24) (61452,)
(20484, 24) (20484,)
(27312, 24) (27312,)
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
       'Date', 'project_grade_category', 'project_title', 'project_essay_1',
       'project_essay_2', 'project_essay_3', 'project_essay_4',
       'project_resource_summary',
       'teacher_number_of_previously_posted_projects', 'clean_categories',
       'clean_subcategories', 'essay', 'preprocessed_essays',
       'totalwords_essay', 'preprocessed_project_title', 'totalwords_title',
       'grade_category', 'price', 'quantity'],
      dtype='object')
```

### 2.2 Vectorizing Numericals Features

#### 2.2.1 Price Standarized

In [33]:

```
from sklearn.preprocessing import StandardScaler
price_scalar = StandardScaler()

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

X_train_price_norm = price_scalar.transform(X_train['price'].values.reshape(-1,1))
X_cv_price_norm = price_scalar.transform(X_cv['price'].values.reshape(-1,1))
X_test_price_norm = price_scalar.transform(X_test['price'].values.reshape(-1,1))

print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(61452, 1) (61452,)  
(20484, 1) (20484,)  
(27312, 1) (27312,)  
=====
```

## 2.2.2 Teacher\_number\_of\_previously\_posted\_projects standardized

In [34]:

```
from sklearn.preprocessing import StandardScaler  
price_scalar = StandardScaler()  
  
price_scalar.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))  
  
X_train_pp_norm = price_scalar.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))  
X_cv_pp_norm = price_scalar.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))  
X_test_pp_norm = price_scalar.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))  
  
print("After vectorizations")  
print(X_train_pp_norm.shape, y_train.shape)  
print(X_cv_pp_norm.shape, y_cv.shape)  
print(X_test_pp_norm.shape, y_test.shape)  
print("="*100)
```

After vectorizations

```
(61452, 1) (61452,)  
(20484, 1) (20484,)  
(27312, 1) (27312,)  
=====
```

## 2.2.3 Quantity Standardized

In [35]:

```
from sklearn.preprocessing import StandardScaler  
price_scalar = StandardScaler()  
  
price_scalar.fit(X_train['quantity'].values.reshape(-1,1))  
  
X_train_quantity_norm = price_scalar.transform(X_train['quantity'].values.reshape(-1,1))  
X_cv_quantity_norm = price_scalar.transform(X_cv['quantity'].values.reshape(-1,1))  
X_test_quantity_norm = price_scalar.transform(X_test['quantity'].values.reshape(-1,1))  
  
print(X_train_quantity_norm.shape, y_train.shape)  
print(X_cv_quantity_norm.shape, y_cv.shape)  
print(X_test_quantity_norm.shape, y_test.shape)  
print("="*100)
```

```
(61452, 1) (61452,)  
(20484, 1) (20484,)  
(27312, 1) (27312,)  
=====
```

## 2.3 Vectorizing Categorical features

### 2.3.1 Vectorizing School\_state

In [36]:



In [36]:

```
# 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

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

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

print("After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

school_fea=vectorizer.get_feature_names()
```

After vectorizations

```
(61452, 51) (61452,)
(20484, 51) (20484,)
(27312, 51) (27312,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'k',
s', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm',
'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv',
', 'wy']
```

## 2.3.2 Vectorizing teacher\_prefix

In [37]:

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

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

print("After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
prefix=vectorizer.get_feature_names()
```

After vectorizations

```
(61452, 5) (61452,)
(20484, 5) (20484,)
(27312, 5) (27312,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
```

## 2.3.3 Vectorizing grade\_category

In [38]:

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

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

print("After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
grade_fea=vectorizer.get_feature_names()
```

```
After vectorizations
(61452, 4) (61452,)
(20484, 4) (20484,)
(27312, 4) (27312,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
=====
```

### 2.3.4 Vectorizing clean\_categories

In [39]:

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

# we use the fitted CountVectorizer to convert the text to vector
X_train_cat_ohe = vectorizer.transform(X_train['clean_categories'].values)
X_cv_cat_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_cat_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("After vectorizations")
print(X_train_cat_ohe.shape, y_train.shape)
print(X_cv_cat_ohe.shape, y_cv.shape)
print(X_test_cat_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
clean_fea=vectorizer.get_feature_names()
```

```
After vectorizations
(61452, 9) (61452,)
(20484, 9) (20484,)
(27312, 9) (27312,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']
=====
```

### 2.3.5 Vectorizing clean\_subcategories

In [40]:

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

# we use the fitted CountVectorizer to convert the text to vector
X_train_scat_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
X_cv_scat_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_scat_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("After vectorizations")
print(X_train_scat_ohe.shape, y_train.shape)
print(X_cv_scat_ohe.shape, y_cv.shape)
```

```
print(X_test_cat_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
sub_fea=vectorizer.get_feature_names()
```

After vectorizations

```
(61452, 30) (61452,)
(20484, 30) (20484,)
(27312, 30) (27312,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience',
'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness',
'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'mathematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
=====
```

## 2.4 Encoding categorical & numerical features

In [41]:

```
#combining all the numerical and categorical values together

from scipy.sparse import hstack
X_tr_com =
hstack((X_train_cat_ohe,X_train_cat_ohe,X_train_state_ohe,X_train_teacher_ohe,X_train_grade_ohe,X
_train_pp_norm,X_train_price_norm,X_train_quantity_norm)).tocsr()
X_cr_com = hstack((X_cv_cat_ohe,X_cv_cat_ohe,X_cv_state_ohe,X_cv_teacher_ohe,X_cv_grade_ohe,X_cv_p
p_norm,X_cv_price_norm,X_cv_quantity_norm)).tocsr()
X_te_com =
hstack((X_test_cat_ohe,X_test_cat_ohe,X_test_state_ohe,X_test_teacher_ohe,X_test_grade_ohe,X_test
pp_norm,X_test_price_norm,X_test_quantity_norm)).tocsr()

print("Final Data matrix")
print(X_tr_com.shape, y_train.shape)
print(X_cr_com.shape, y_cv.shape)
print(X_te_com.shape, y_test.shape)
print("="*100)
fea=clean_fea+sub_fea+school_fea+prefix+grade_fea+['previously_posted_project','price','quantity']

print(len(fea))
```

Final Data matrix

```
(61452, 102) (61452,)
(20484, 102) (20484,)
(27312, 102) (27312,)
=====
```

102

In [42]:

```
# 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
```

## 2.4 Appling DT on different kind of featurization as mentioned in the instructions

Apply DT 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 [43]:

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

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

    return y_data_pred
```

## 2.4.1 Applying DT on BOW, SET 1

### 2.4.1.1 Converting project\_title & essay into BOW

In [44]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train data

X_train_essay_bow = vectorizer.transform(X_train['preprocessed_essays'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['preprocessed_essays'].values)
X_test_essay_bow = vectorizer.transform(X_test['preprocessed_essays'].values)

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("=="*100)

essay_fea=vectorizer.get_feature_names()
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['preprocessed_project_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['preprocessed_project_title'].values)
X_cv_title_bow = vectorizer.transform(X_cv['preprocessed_project_title'].values)
X_test_title_bow = vectorizer.transform(X_test['preprocessed_project_title'].values)

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
print("=="*100)
title_fea=vectorizer.get_feature_names()
```

```
After vectorizations
(61452, 5000) (61452,)
(20484, 5000) (20484,)
(27312, 5000) (27312,)
```

=====

```
After vectorizations
(61452, 5000) (61452,)
(20484, 5000) (20484,)
(27312, 5000) (27312,)
```

=====

### 2.4.1.2 Encoding numerical,categorical & BOW

In [45]:

```
from scipy.sparse import hstack
X_tr_bow = hstack((X_train_essay_bow, X_train_title_bow,X_tr_com)).tocsr()
X_cr_bow = hstack((X_cv_essay_bow, X_cv_title_bow,X_cr_com)).tocsr()
X_te_bow = hstack((X_test_essay_bow, X_test_title_bow,X_te_com)).tocsr()

print("Final Data matrix")
print(X_tr_bow.shape, y_train.shape)
print(X_cr_bow.shape, y_cv.shape)
print(X_te_bow.shape, y_test.shape)
print("="*100)
feat=essay_fea+title_fea+fea
len(feat)
```

```
Final Data matrix
(61452, 10102) (61452,)
(20484, 10102) (20484,)
(27312, 10102) (27312,)
```

Out[45]:

10102

### 2.4.1.3 Applying DT

In [46]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []

max_depth=[1,5,10,50,100]
for i in max_depth:

    clf= DecisionTreeClassifier(max_depth=i)
    clf.fit(X_tr_bow, y_train)

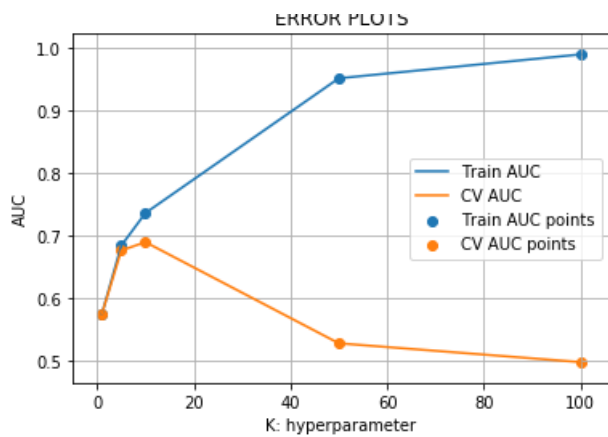
    y_train_pred = batch_predict(clf, X_tr_bow)
    y_cv_pred = batch_predict(clf, X_cr_bow)

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(max_depth, train_auc, label='Train AUC')
plt.plot(max_depth, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [43]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []
min_split=[5,10,100,500]
for i in min_split:

    clf= DecisionTreeClassifier(min_samples_split=i)
    clf.fit(X_tr_bow, y_train)

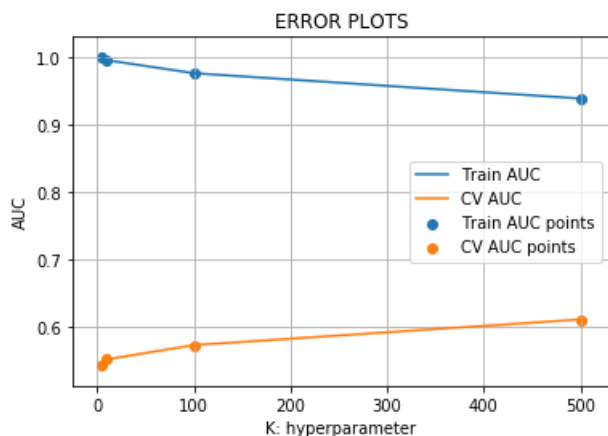
    y_train_pred = batch_predict(clf, X_tr_bow)
    y_cv_pred = batch_predict(clf, X_cr_bow)

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(min_split, train_auc, label='Train AUC')
plt.plot(min_split, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [43]:

```
from sklearn.metrics import roc_curve, auc

clf=DecisionTreeClassifier(max_depth=15,min_samples_split=100)
```

```

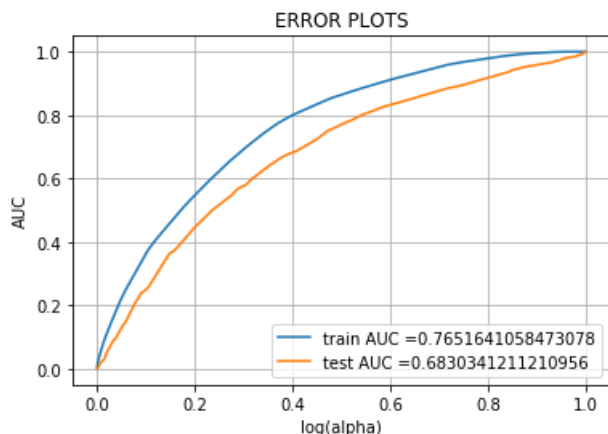
clf.fit(X_tr_bow, y_train)

y_train_pred =batch_predict(clf,X_tr_bow)
y_test_pred = batch_predict(clf,X_te_bow)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [72]:

```

def predict(proba, threshold, fpr, tpr):

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

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

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

```

In [46]:

```

#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn

print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
cmtr=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_tpr))
import seaborn as sns
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmtr, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

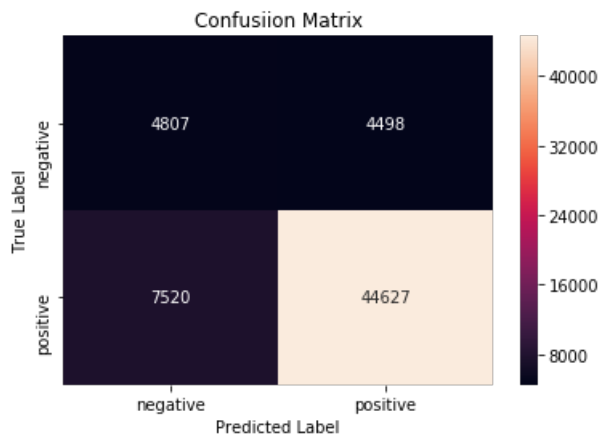
```

```

=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.2497243079691429 for threshold 0.727

```



In [47]:

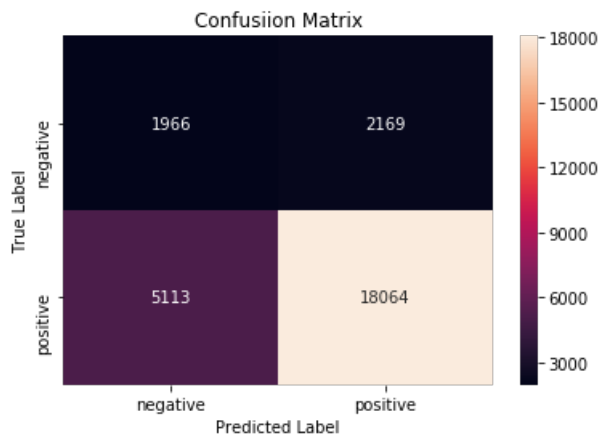
```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn

# plot confusion matrix to describe the performance of classifier.
import seaborn as sns
print("Test confusion matrix")

cmt=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmt, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix

the maximum value of  $tpr \cdot (1 - fpr)$  0.24982628313757718 for threshold 0.814



In [140]:

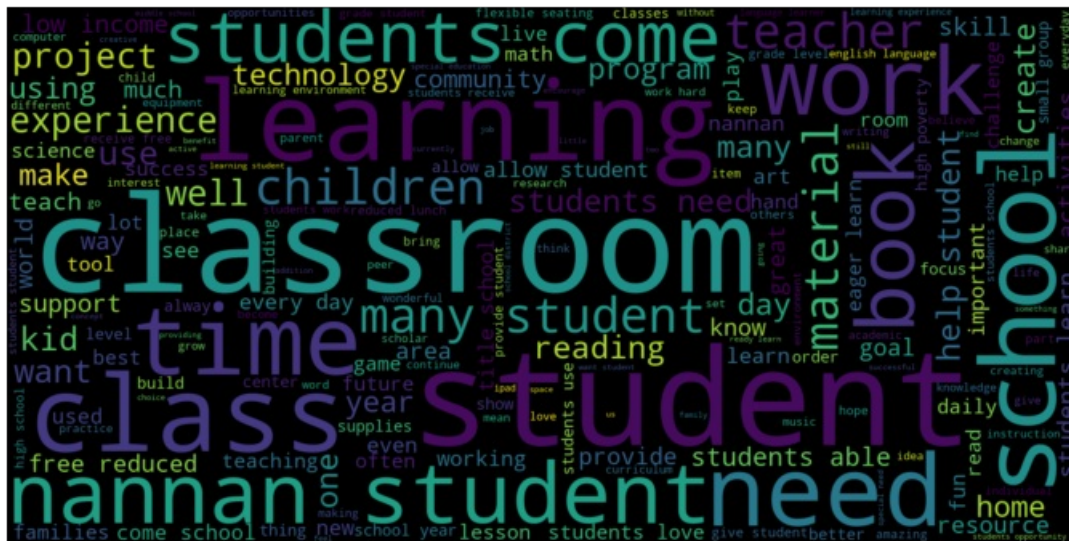
```
#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
y_test_pre = clf.predict(X_te_bow)
fp_phrases =list(X_test['preprocessed_essays'])
price_t= list(X_test['price'])
teach= list(X_test['teacher_number_of_previously_posted_projects'])
fp_words = []
fp_price=[]
fp_teacher=[]
for i in range(len(y_test_pre)):
    if y_test_pre[i]==1 and y_test[i]!=y_test_pre[i]:
        fp_words.append(neg_phrases[i])
        fp_price.append(price_t[i])
        fp_teacher.append(teach[i])
```



```
fp_text = pd.Series(neg_words).str.cat(sep=' ')
fp_text[:100]
```

'work low socio economic area students come us variety needs pertaining supplies food shelter love  
cl'

```
#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
from wordcloud import WordCloud
wordcloud = WordCloud(width=1600, height=800, max_font_size=200).generate(fp_text)
plt.figure(figsize=(12,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



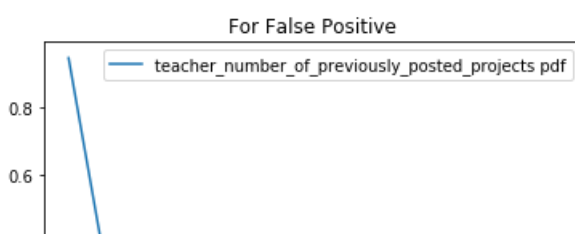
```
counts,bin_edges=np.histogram(fp_teacher,bins=10,density=True)
pdf=counts/sum(counts)
print("PDF Values ")
print(pdf)
print("Bin Edges ")
print(bin_edges)
plt.plot(bin_edges[1:],pdf,label="teacher_number_of_previously_posted_projects pdf")
plt.legend()
plt.title("For False Positive")
plt.xlabel("`Teacher_number_of_previously_posted_projects")
plt.show()
print("-----")
```

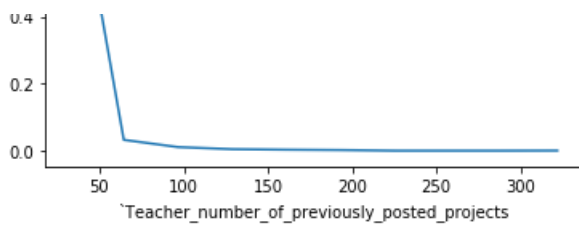
```
[9.47792208e-01 3.22077922e-02 1.06493506e-02 4.41558442e-03
 2.85714286e-03 1.81818182e-03 0.00000000e+00 0.00000000e+00
 0.00000000e+00 2.59740260e-04]
```

```

BIN Edges
[ 0.    32.2  64.4  96.6 128.8 161.   193.2 225.4 257.6 289.8 322. ]

```

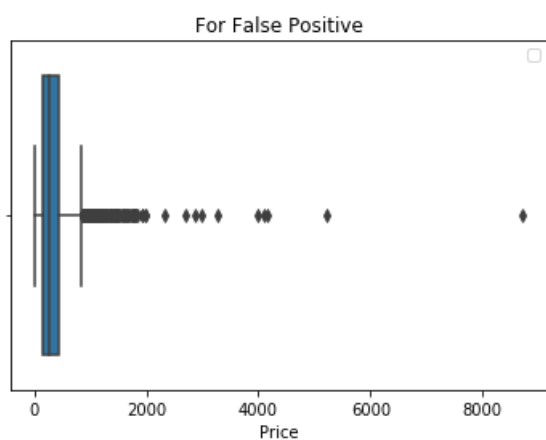




In [162]:

```
sns.boxplot(fp_price)
plt.legend()
plt.title("For False Positive")
plt.xlabel("Price")
plt.show()
```

No handles with labels found to put in legend.



In [79]:

```
import os
os.environ["PATH"] += os.pathsep + 'C:/Program Files (x86)/Graphviz2.38/bin'
```

In [93]:

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

clff.fit(X_tr_bow, y_train)

import graphviz
dot_data = tree.export_graphviz(clff, out_file=None)
graph = graphviz.Source(dot_data)
dot_data = tree.export_graphviz(clff, out_file=None, feature_names=feat, filled=True,
rounded=True, special_characters=True)
graph = graphviz.Source(dot_data)
graph.render(filename='D:\\graph.dot')
```

Out[93]:

'D:\\graph.dot.pdf'

## 2.4.2 Applying DT on TFIDF, SET 2

### 2.4.2.1 Converting Project\_title & essay into tf-idf

In [50]:

```

from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['preprocessed_essays'].values) # fit has to happen only on train data

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

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_cv_essay_tfidf.shape, y_cv.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)

vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['preprocessed_project_title'].values) # fit has to happen only on train
data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['preprocessed_project_title'].values)
X_cv_title_tfidf = vectorizer.transform(X_cv['preprocessed_project_title'].values)
X_test_title_tfidf = vectorizer.transform(X_test['preprocessed_project_title'].values)

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
print("="*100)

```

```

After vectorizations
(61452, 5000) (61452,)
(20484, 5000) (20484,)
(27312, 5000) (27312,)
=====

```

```

After vectorizations
(61452, 5000) (61452,)
(20484, 5000) (20484,)
(27312, 5000) (27312,)
=====

```

#### 2.4.2.2 Encoding numerical,categorical & Tf-idf

In [51]:

```

from scipy.sparse import hstack
X_tr_tfidf = hstack((X_train_essay_tfidf, X_train_title_tfidf,X_tr_com)).tocsr()
X_cr_tfidf = hstack((X_cv_essay_tfidf,X_cv_title_tfidf,X_cr_com)).tocsr()
X_te_tfidf = hstack((X_test_essay_tfidf, X_test_title_tfidf,X_te_com)).tocsr()

print("Final Data matrix")
print(X_tr_tfidf.shape, y_train.shape)
print(X_cr_tfidf.shape, y_cv.shape)
print(X_te_tfidf.shape, y_test.shape)
print("="*100)

```

```

Final Data matrix
(61452, 10102) (61452,)
(20484, 10102) (20484,)
(27312, 10102) (27312,)
=====

```

#### 2.4.2.3 Applying DT

In [92]:

```

import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

```

```

from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []
max_depth=[1,5,10,50,100]
for i in max_depth:

    clf= DecisionTreeClassifier(max_depth=i)
    clf.fit(X_tr_tfidf, y_train)

    y_train_pred = batch_predict(clf, X_tr_tfidf)
    y_cv_pred = batch_predict(clf, X_cr_tfidf)

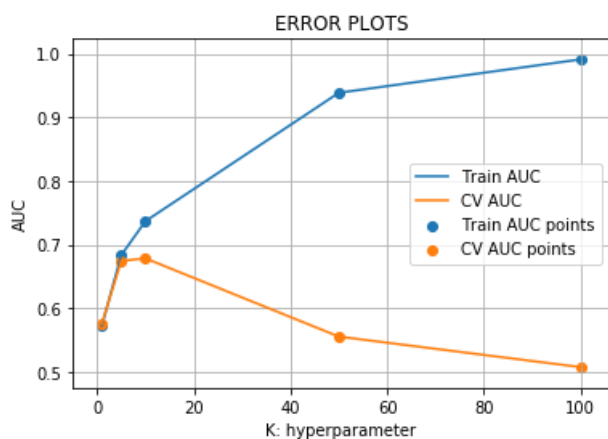
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(max_depth, train_auc, label='Train AUC')
plt.plot(max_depth, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [93]:

```

import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []
min_split=[5,10,100,500]
for i in min_split:

    clf= DecisionTreeClassifier(min_samples_split=i)
    clf.fit(X_tr_tfidf, y_train)

    y_train_pred = batch_predict(clf, X_tr_tfidf)
    y_cv_pred = batch_predict(clf, X_cr_tfidf)

    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

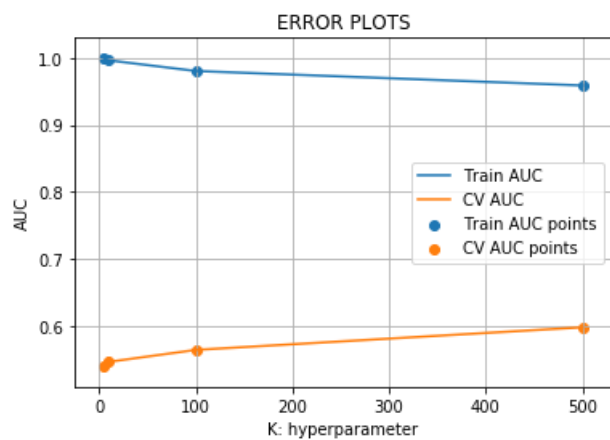
plt.plot(min_split, train_auc, label='Train AUC')
plt.plot(min_split, cv_auc, label='CV AUC')

plt.scatter(min_split, train_auc, label='Train AUC points')

```

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [52]:

```
from sklearn.metrics import roc_curve, auc

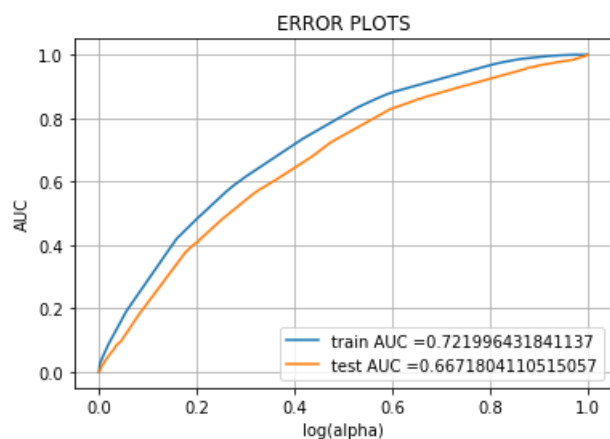
clf=DecisionTreeClassifier(max_depth=10,min_samples_split=100)

clf.fit(X_tr_tfidf, y_train)

y_train_pred =batch_predict(clf,X_tr_tfidf)
y_test_pred = batch_predict(clf,X_te_tfidf)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [95]:

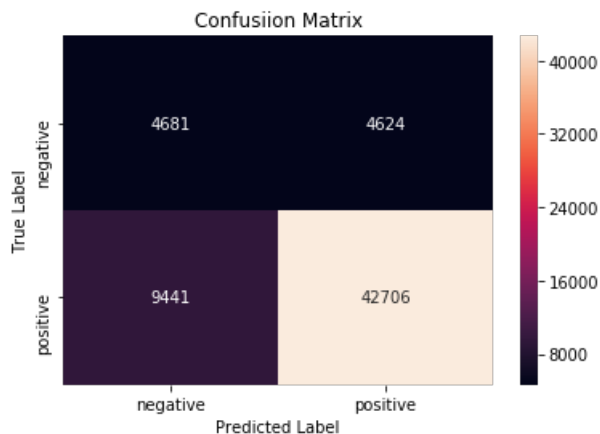
```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
print("="*100)
from sklearn.metrics import confusion matrix
```

```

print("Train confusion matrix")
#print(confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr)))
cmtr=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmtr, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24999061883088516 for threshold 0.821



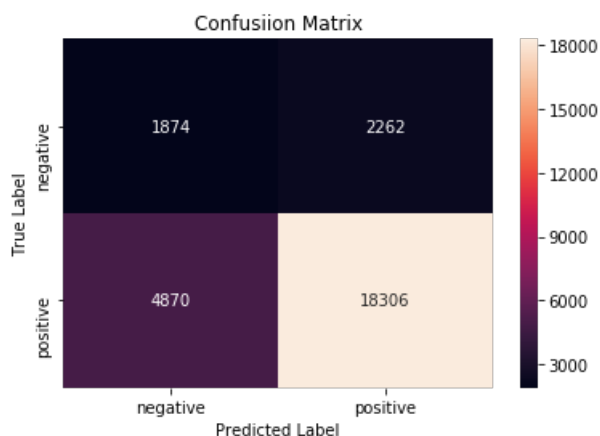
In [96]:

```

#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
print("Test confusion matrix")
cmt=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmt, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

```

Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24956764962269307 for threshold 0.826



#### 2.4.2.4 Word cloud for False Positive points

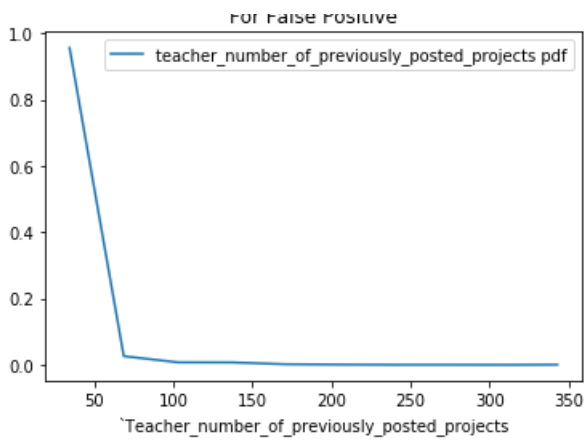
In [97]:

```

#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn

```



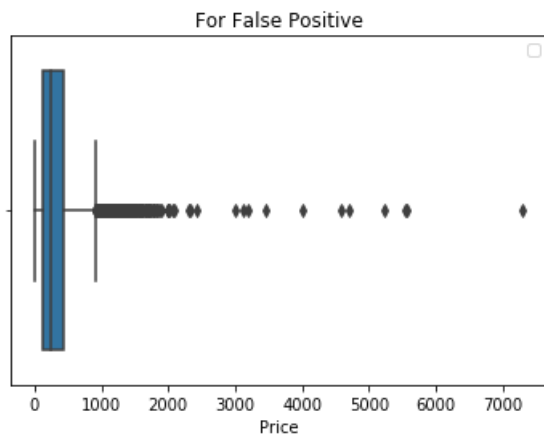


#### 2.4.2.6 Box-plot for 'Price'

In [99]:

```
sns.boxplot(fp_price)
plt.legend()
plt.title("For False Positive")
plt.xlabel("Price")
plt.show()
```

No handles with labels found to put in legend.



#### 2.4.2.7 Tree Representation

In [81]:

```
clff=tree.DecisionTreeClassifier(max_depth=3)

clff.fit(X_tr_tfidf, y_train)

import graphviz
dot_data = tree.export_graphviz(clff, out_file=None)
graph = graphviz.Source(dot_data)
dot_data = tree.export_graphviz(clff, out_file=None, feature_names=feat, filled=True,
rounded=True, special_characters=True)
graph = graphviz.Source(dot_data)
graph
```

Out[81]:



### 2.4.3 Applying DT on AVG W2V, SET 3

#### 2.4.3.1 Converting Project\_essay to Avg W2V

In [49]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_train_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train_essay.append(vector)

print(len(avg_w2v_vectors_train_essay))
print(len(avg_w2v_vectors_train_essay[0]))

avg_w2v_vectors_cv_essay = [];
for sentence in tqdm(X_cv['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv_essay.append(vector)

print(len(avg_w2v_vectors_cv_essay))
print(len(avg_w2v_vectors_cv_essay[0]))

avg_w2v_vectors_test_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_test_essay.append(vector)

print(len(avg_w2v_vectors_test_essay))
print(len(avg_w2v_vectors_test_essay[0]))
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 61452/61452  
[00:58<00:00, 1046.21it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

```
100%|██████████████████████████████████████████████████████████████████████████| 20484/20484  
[00:10<00:00, 2027.85it/s]
```

20484  
300

```
100%|██████████████████████████████████████████████████████████████████████████| 27312/27312  
[00:13<00:00, 2020.01it/s]
```

27312  
300

#### 2.4.3.2 Converting project title to Avg W2V

In [50]:

```

avg_w2v_vectors_train_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_project_title']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_train_title.append(vector)

print(len(avg_w2v_vectors_train_title))
print(len(avg_w2v_vectors_train_title[0]))

avg_w2v_vectors_cv_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['preprocessed_project_title']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_cv_title.append(vector)

print(len(avg_w2v_vectors_cv_title))
print(len(avg_w2v_vectors_cv_title[0]))

avg_w2v_vectors_test_title = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['preprocessed_project_title']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_test_title.append(vector)

print(len(avg_w2v_vectors_test_title))
print(len(avg_w2v_vectors_test_title[0]))

```

```
100%|██████████████████████████████████████████████████████████████████████████| 61452/61452  
[00:01<00:00, 36823.77it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

```
100%|██████████████████████████████████████████████████████████████████████████| 20484/20484  
[00:00<00:00, 30074.28it/s]
```

20484  
300

```
100%|██████████████████████████████████████████████████████████████████████████| 27312/27312  
[00:00<00:00, 39420.20it/s]
```

27312  
300

## 2.4.3.3 Combining numerical, categorical & Avg W2V

In [51]:

```
from scipy.sparse import hstack
X_tr_w2v = hstack((avg_w2v_vectors_train_essay ,avg_w2v_vectors_train_title, X_tr_com)).tocsr()
X_cr_w2v = hstack((avg_w2v_vectors_cv_essay, avg_w2v_vectors_cv_title,X_cr_com)).tocsr()
X_te_w2v = hstack((avg_w2v_vectors_test_essay, avg_w2v_vectors_test_title,X_te_com)).tocsr()

print("Final Data matrix")
print(X_tr_w2v.shape, y_train.shape)
print(X_cr_w2v.shape, y_cv.shape)
print(X_te_w2v.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(61452, 702) (61452,)
(20484, 702) (20484,)
(27312, 702) (27312,)
```

## 2.4.3.4 Applying DT

In [78]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []
neww=[]
max_depth=[1,5,10,50,100]
for i in max_depth:

    clf= DecisionTreeClassifier(max_depth=i)
    clf.fit(X_tr_w2v, y_train)

    y_train_pred_w2v = batch_predict(clf, X_tr_w2v)
    y_cv_pred_w2v = batch_predict(clf, X_cr_w2v)

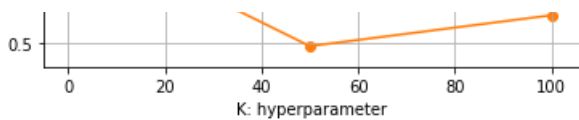
    train_auc.append(roc_auc_score(y_train,y_train_pred_w2v))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred_w2v))

plt.plot(max_depth, train_auc, label='Train AUC')
plt.plot(max_depth, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```





In [76]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier

train_auc = []
cv_auc = []
min_split=[5,10,100,500]
for i in min_split:

    clf= DecisionTreeClassifier(min_samples_split=i)
    clf.fit(X_tr_w2v, y_train)

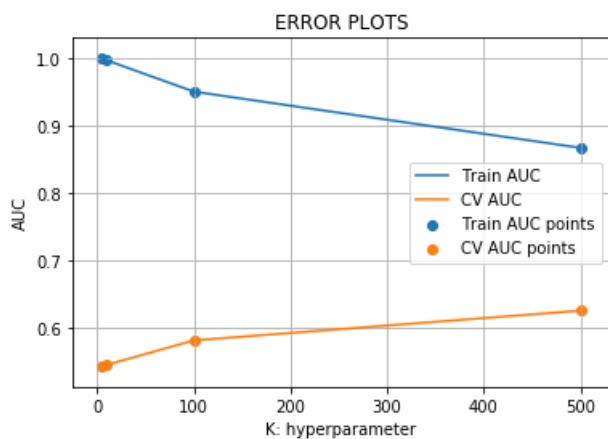
    y_train_pred_w2v = batch_predict(clf, X_tr_w2v)
    y_cv_pred_w2v = batch_predict(clf, X_cr_w2v)

    train_auc.append(roc_auc_score(y_train,y_train_pred_w2v))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred_w2v))

plt.plot(min_split, train_auc, label='Train AUC')
plt.plot(min_split, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [54]:

```
from sklearn.metrics import roc_curve, auc

clf=DecisionTreeClassifier(max_depth=15,min_samples_split=100)

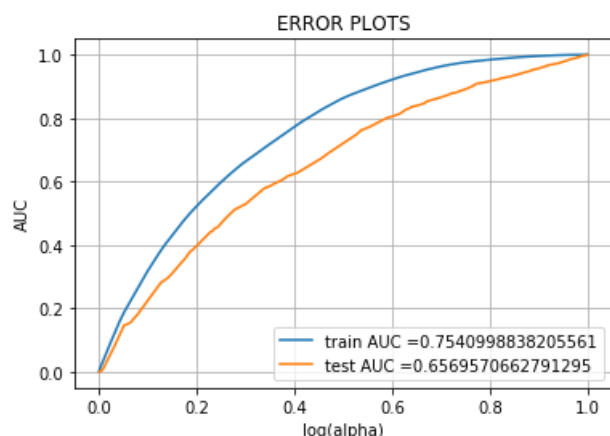
clf.fit(X_tr_w2v, y_train)

y_train_pred =batch_predict(clf,X_tr_w2v)
y_test_pred = batch_predict(clf,X_te_w2v)

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

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

```
plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



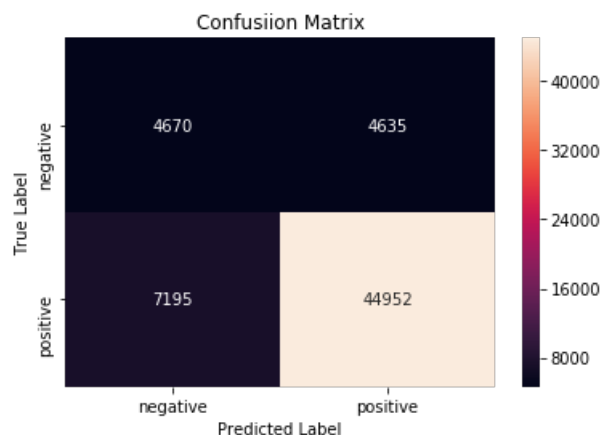
In [57]:

```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
print("="*100)
from sklearn.metrics import confusion_matrix

print("Train confusion matrix")

cmtr=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmtr, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Train confusion matrix  
the maximum value of tpr\*(1-fpr) 0.24999646293254363 for threshold 0.802



In [58]:

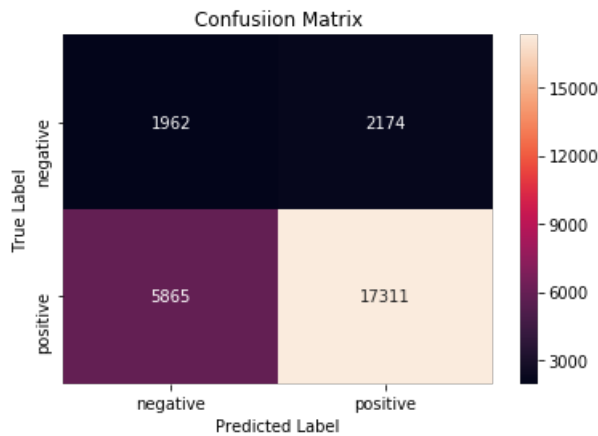
```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
import seaborn as sns
print("Test confusion matrix")

cmt=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmt, index = class_label, columns = class_label)
```

```
df_cm = pd.DataFrame(cm, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix

the maximum value of  $tpr \cdot (1 - fpr)$  0.24988682661837935 for threshold 0.858



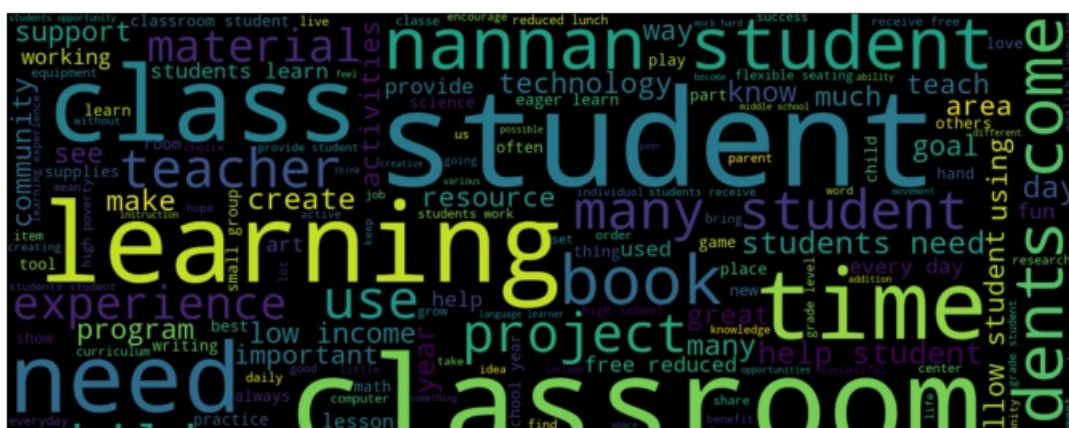
### 2.4.3.5 Wordcloud for 'False Positive' values

In [61]:

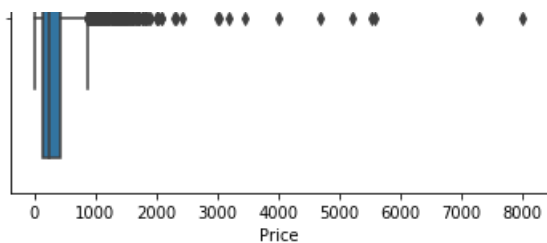
```
#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
y_test_pre = clf.predict(X_te_w2v)
fp_phrases = list(X_test['preprocessed_essays'])
price_t = list(X_test['price'])
teach = list(X_test['teacher_number_of_previously_posted_projects'])
fp_words = []
fp_price = []
fp_teacher = []
for i in range(len(y_test_pre)):
    if y_test_pre[i]==1 and y_test[i]!=y_test_pre[i]:
        fp_words.append(fp_phrases[i])
        fp_price.append(price_t[i])
        fp_teacher.append(teach[i])

fp_text = pd.Series(fp_words).str.cat(sep=' ')
fp_text[:100]

#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
from wordcloud import WordCloud
wordcloud = WordCloud(width=1600, height=800, max_font_size=200).generate(fp_text)
plt.figure(figsize=(12,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```







## 2.4.4 Applying DT on TFIDF W2V, SET 4

### 2.4.4.1 Converting project\_essay to TF-idf W2V

In [64]:

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

tfidf_w2v_vectors_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay.append(vector)

print(len(tfidf_w2v_vectors_essay))
print(len(tfidf_w2v_vectors_essay[0]))

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

tfidf_w2v_vectors_essay_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay_cv.append(vector)

print(len(tfidf_w2v_vectors_essay_cv))
print(len(tfidf_w2v_vectors_essay_cv[0]))

tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_essays'])
# we are converting a dictionary with word as a key, and the idf as a value
```



```
# we are converting a dictionary with word as a key, and the tf 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_vectors_essay_te = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split()))))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay_te.append(vector)

print(len(tfidf_w2v_vectors_essay_te))
print(len(tfidf_w2v_vectors_essay_te[0]))
```

```
100%|██████████████████████████████████████████████████████████████| 61452/61452 [02:  
39<00:00, 384.87it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

```
100%|██████████████████████████████████████████████████████████████████████████████| 61452/61452 [02:  
58<00:00, 344.30it/s]
```

61452  
300

```
100%|██████████████████████████████████████████████████████████| 61452/61452 [02:  
48<00:00, 365.13it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

#### 2.4.4.2 Converting project\_title to TF-idf W2V

In [65]:

```
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['preprocessed_project_title'])
# 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_vectors_essay = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay.append(vector)

print(len(tfidf_w2v_vectors_essay))
```

```
print(len(tfidf_w2v_vectors_essay))
print(len(tfidf_w2v_vectors_essay[0]))

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

tfidf_w2v_vectors_essay_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay_cv.append(vector)

print(len(tfidf_w2v_vectors_essay_cv))
print(len(tfidf_w2v_vectors_essay_cv[0]))

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

tfidf_w2v_vectors_essay_te = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['preprocessed_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_essay_te.append(vector)

print(len(tfidf_w2v_vectors_essay_te))
print(len(tfidf_w2v_vectors_essay_te[0]))
```

```
100%|███████████████████████████████████████████████████████████| 61452/61452 [02:  
45<00:00, 370.49it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

```
100% |██████████████████████████████████████████████████████| 61452/61452 [02:  
43<00:00, 376.67it/s]
```

$$\begin{array}{r} 61452 \\ 300 \end{array}$$

```
100% |██████████████████████████████████████████████████████| 61452/61452 [02:  
44<00:00, 373.45it/s]
```

61452

### 2.4.4.3 Combing numerical,categorical features & tf-idf W2V

In [66]:

```
from scipy.sparse import hstack
X_tr_w2v_tfidf = hstack((avg_w2v_vectors_train_essay ,avg_w2v_vectors_train_title,
X_tr_com)).tocsr()
X_cr_w2v_tfidf = hstack((avg_w2v_vectors_cv_essay, avg_w2v_vectors_cv_title,X_cr_com)).tocsr()
X_te_w2v_tfidf = hstack((avg_w2v_vectors_test_essay, avg_w2v_vectors_test_title,X_te_com)).tocsr()

print("Final Data matrix")
print(X_tr_w2v_tfidf.shape, y_train.shape)
print(X_cr_w2v_tfidf.shape, y_cv.shape)
print(X_te_w2v_tfidf.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(61452, 702) (61452,)
(20484, 702) (20484,)
(27312, 702) (27312,)
=====
```



### 2.4.4.4 Applying DT

In [77]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier
from mpl_toolkits.mplot3d import Axes3D

train_auc = []
cv_auc = []
neww=[]
max_depth=[1,5,10,50,100]
for i in max_depth:

    clf= DecisionTreeClassifier(max_depth=i)
    clf.fit(X_tr_w2v_tfidf, y_train)

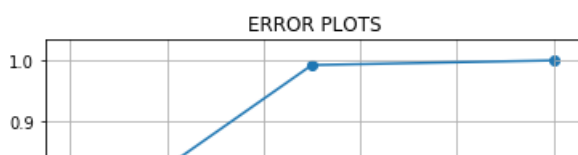
    y_train_pred = batch_predict(clf, X_tr_w2v_tfidf)
    y_cv_pred = batch_predict(clf, X_cr_w2v_tfidf)

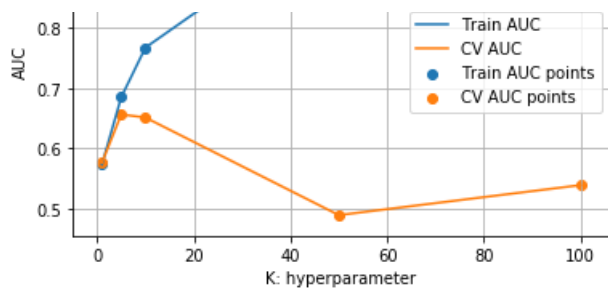
    train_auc.append(roc_auc_score(y_train,y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot(max_depth, train_auc, label='Train AUC')
plt.plot(max_depth, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```





In [75]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn.tree import DecisionTreeClassifier

train_auc = []
cv_auc = []
min_split=[5,10,100,500]
for i in min_split:

    clf= DecisionTreeClassifier(min_samples_split=i)
    clf.fit(X_tr_w2v_tfidf, y_train)

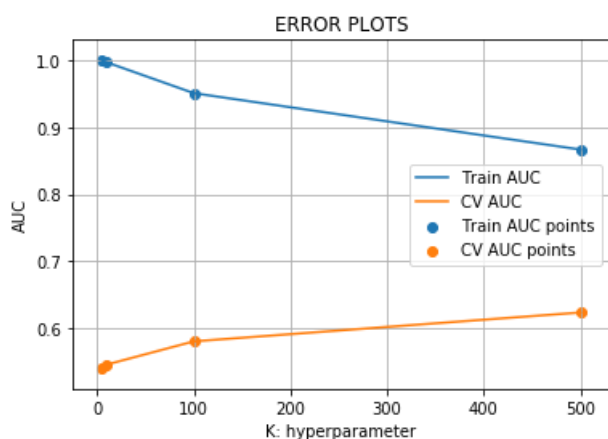
    y_train_pred_w2v = batch_predict(clf, X_tr_w2v_tfidf)
    y_cv_pred_w2v = batch_predict(clf, X_cr_w2v_tfidf)

    train_auc.append(roc_auc_score(y_train,y_train_pred_w2v))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred_w2v))

plt.plot(min_split, train_auc, label='Train AUC')
plt.plot(min_split, cv_auc, label='CV AUC')

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

plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



In [69]:

```
from sklearn.metrics import roc_curve, auc

clf=DecisionTreeClassifier(max_depth=5,min_samples_split=100)

clf.fit(X_tr_w2v_tfidf, y_train)

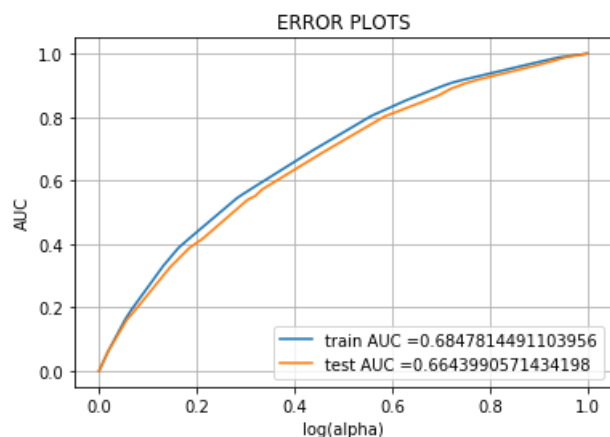
y_train_pred =batch_predict(clf,X_tr_w2v_tfidf)
y_test_pred = batch_predict(clf,X_te_w2v_tfidf)
```

```

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```



In [70]:

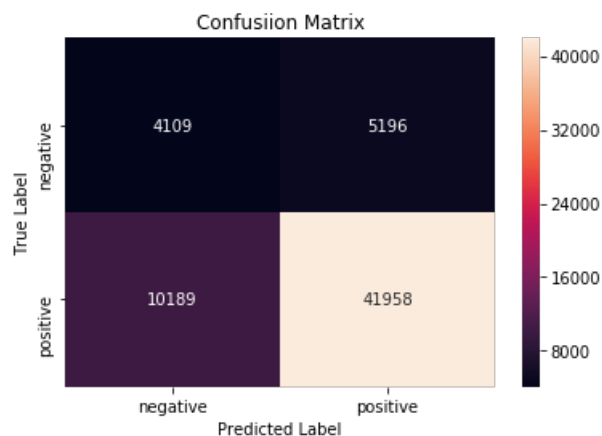
```

#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")

cmtr=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmtr, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()

```

Train confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24658833530013533 for threshold 0.835



In [71]:

```

#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn

```

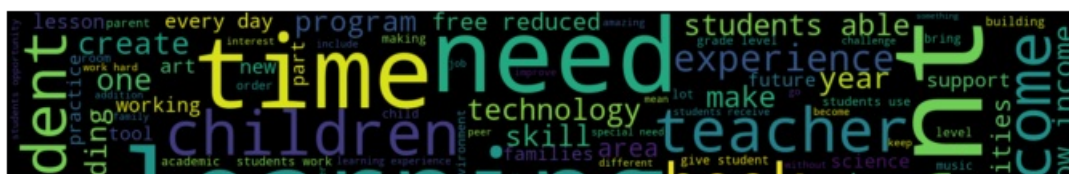
Test confusion matrix  
the maximum value of  $tpr \cdot (1 - fpr)$  0.24836968365701542 for threshold 0.843



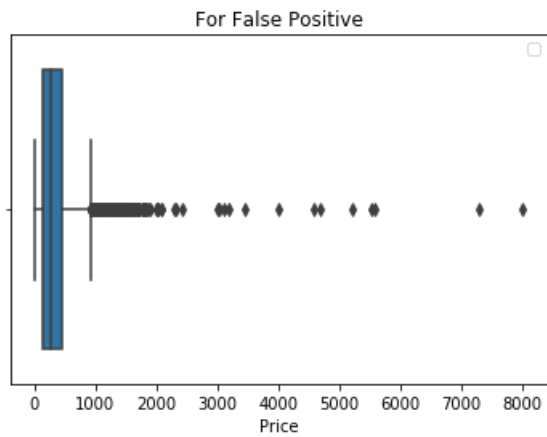
```
#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
y_test_pre = clf.predict(X_te_w2v_tfidf)
fp_phrases =list(X_test['preprocessed_essays'])
price_t= list(X_test['price'])
teach= list(X_test['teacher_number_of_previously_posted_projects'])
fp_words = []
fp_price=[]
fp_teacher=[]
for i in range(len(y_test_pre)):
    if y_test_pre[i]==1 and y_test[i]!=y_test_pre[i]:
        fp_words.append(fp_phrases[i])
        fp_price.append(price_t[i])
        fp_teacher.append(teach[i])

fp_text = pd.Series(fp_words).str.cat(sep=' ')
fp_text[:100]

#https://www.kaggle.com/divsinha/sentiment-analysis-countvectorizer-tf-idf
from wordcloud import WordCloud
wordcloud = WordCloud(width=1600, height=800, max_font_size=200).generate(fp_text)
plt.figure(figsize=(12,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```







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

### 2.5.1 Selecting top 5000 Features from set 2

In [59]:

```
#https://datascience.stackexchange.com/questions/6683/feature-selection-using-feature-importances-in-random-forests-with-scikit-learn
def selectKImportance(model, X, k=5):
    return X[:,model.feature_importances_.argsort()[::-1][:k]]
```

In [67]:

```
newX_tr = selectKImportance(clf,X_tr_tfidf,5000)
newX_cv = selectKImportance(clf,X_cv_tfidf,5000)
newX_te = selectKImportance(clf,X_te_tfidf,5000)
print(newX_tr.shape)
print(newX_cv.shape)
print(newX_te.shape)
```

```
(61452, 5000)
(20484, 5000)
(27312, 5000)
```

### 2.5.2 Applying Logistic Regression

In [69]:

```
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
from sklearn import linear_model
from math import log
from sklearn.model_selection import GridSearchCV

neigh=linear_model.SGDClassifier(loss='log')
parameters = {'alpha':[10**-5,10**-4,10**-3,10**-2,10**-1,1,10,100,1000,10000,100000]}
clf = GridSearchCV(neigh, parameters, cv=3, scoring='roc_auc')
clf.fit(newX_tr, y_train)

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

parameter=[log(y) for y in parameters['alpha']]

print(parameter)
plt.plot((parameter), train_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
```



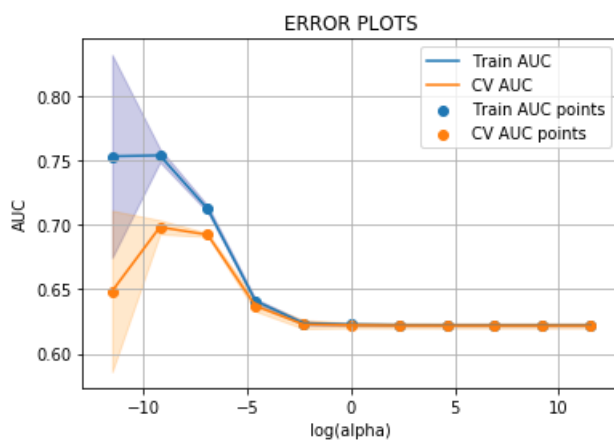
```
plt.gca().fill_between(parameter,train_auc - train_auc_std,train_auc +
train_auc_std,alpha=0.2,color='darkblue')

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

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

plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```

```
[-11.512925464970229, -9.210340371976182, -6.907755278982137, -4.605170185988091, -
2.3025850929940455, 0.0, 2.302585092994046, 4.605170185988092, 6.907755278982137,
9.210340371976184, 11.512925464970229]
```



In [70]:

```
k_5 = 10** -4

from sklearn.metrics import roc_curve, auc
from sklearn.calibration import CalibratedClassifierCV

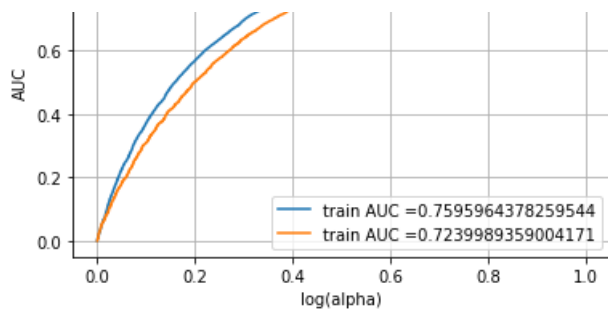
clf = linear_model.SGDClassifier(loss='log',alpha=k_5)
neigh=CalibratedClassifierCV(clf,method='sigmoid')
neigh.fit(newX_tr, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

y_train_pred = batch_predict(neigh, newX_tr)
y_test_pred = batch_predict(neigh, newX_te)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("log(alpha)")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()
```



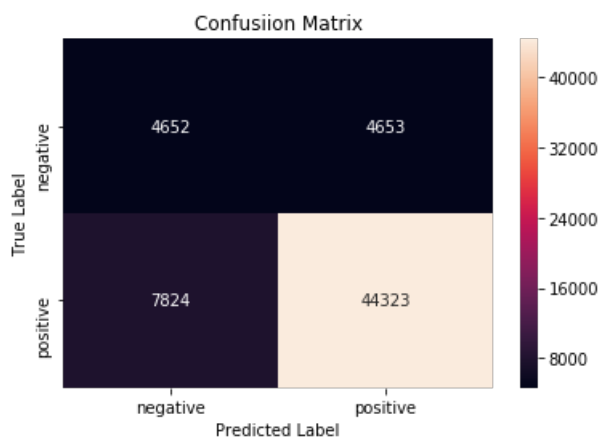


In [73]:

```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
print("Train confusion matrix")

cmt=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmt, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Train confusion matrix  
the maximum value of  $tpr \cdot (1-fpr)$  0.249999997112598 for threshold 0.787

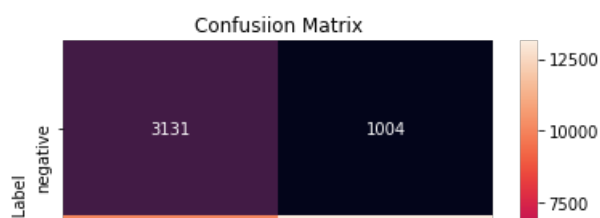


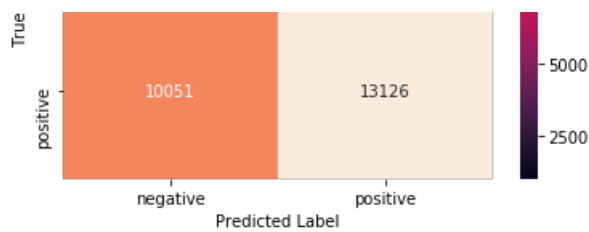
In [74]:

```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
import seaborn as sns
print("Test confusion matrix")

cmt=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
class_label = ["negative", "positive"]
df_cm = pd.DataFrame(cmt, index = class_label, columns = class_label)
sns.heatmap(df_cm, annot = True, fmt = "d")
plt.title("Confusiion Matrix")
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.show()
```

Test confusion matrix  
the maximum value of  $tpr \cdot (1-fpr)$  0.24999998537859927 for threshold 0.867





### 3. Conclusions

In [76]:

```
#https://www.kaggle.com/premvardhan/amazon-fine-food-reviews-analysis-using-knn
models = pd.DataFrame({'Model': ['DT with Bow', 'DT with TFIDF', 'DT with Avg_w2v', 'DT with tfidf_w2v'], 'Deapth': [15,10,15,5], 'Min_sample_split':[100,100,100,100], 'Train AUC': [.765,.721,.754,.684], 'Test AUC': [.683,.667,.656,.664]}, columns = ["Model", "Deapth", "Min_sample_split", "Train AUC", "Test AUC"])
models#.sort_values(by='Test AUC', ascending=False)
```

Out[76]:

	Model	Deapth	Min_sample_split	Train AUC	Test AUC
0	DT with Bow	15	100	0.765	0.683
1	DT with TFIDF	10	100	0.721	0.667
2	DT with Avg_w2v	15	100	0.754	0.656
3	DT with tfidf_w2v	5	100	0.684	0.664

In [ ]: