

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	Title of the project. Examples: <ul style="list-style-type: none">• 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: <ul style="list-style-type: none">• 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: <ul style="list-style-type: none">• Applied Learning• Care & Hunger• Health & Sports• History & Civics• Literacy & Language• Math & Science• Music & The Arts• Special Needs• Warmth Examples: <ul style="list-style-type: none">• Music & The Arts• Literacy & Language, Math & Science
<code>school_state</code>	State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. Examples: <ul style="list-style-type: none">• Literacy

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

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

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

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

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

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

Label	Description
<code>project_is_approved</code>	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"

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
```

```
/usr/local/lib/python3.6/dist-packages/smart_open/ssh.py:34: UserWarning: paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip install paramiko` to suppress
warnings.warn('paramiko missing, opening SSH/SCP/SFTP paths will be disabled. `pip install paramiko` to suppress')
```

1.1 Reading Data

In [0]:

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

In [3]:

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

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

```
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category']
```

```
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

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

```
Number of data points in train data (1541272, 4)
['id' 'description' 'quantity' 'price']
```

Out[4]:

	id	description	quantity	price
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 [0]:

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

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [0]:

```
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() + " #"
        temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

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

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

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

1.3 Text preprocessing

In [0]:

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

In [8]:

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

Out[8]:

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

In [0]:

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

In [10]:

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

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

=====

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

=====

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

my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nannan

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

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward\r\n\r\nMy school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but one smart, effective, efficient, and disciplined students with good character. In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.\r\n\r\nThe cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.nannan

In [0]:

```
# 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[20000])
print(sent)
print("="*50)
```

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

In [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)
```

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

In [14]:

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

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

In [0]:

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

```
100%|██████████| 109248/109248 [01:05<00:00, 1662.50it/s]
```

In [17]:

```
# after preprocessing
preprocessed_essays[20000]
```

Out[17]:

```
'my kindergarten students varied disabilities ranging speech language delays cognitive delays gross fine motor delays autism they eager beavers always strive work hardest working past limitations the materials ones i seek students i teach title i school students receive free reduced price lunch despite disabilities limitations students love coming school come eager learn explore have ever felt like ants pants needed groove move meeting this kids feel time the want able move learn say wobble chairs answer i love develop core enhances gross motor turn fine motor skills they also want learn games kids not want sit worksheets they want learn count jumping playing physical engagement key success the number toss color shape mats make happen my students forget work fun 6 year old deserves nannan'
```

1.4 Preprocessing of `project_title`

In [18]:

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

```
100%|██████████| 109248/109248 [00:03<00:00, 34963.94it/s]
```

1.5 Preparing data for models

In [19]:

```
project_data.columns
```

Out[19]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category', 'project_title',
      'project_essay_1', 'project_essay_2', 'project_essay_3',
```

```
'project_essay_4', 'project_resource_summary',
'teacher_number_of_previously_posted_projects', 'project_is_approved',
'clean_categories', 'clean_subcategories', 'essay'],
dtype='object')
```

we are going to consider

```
- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data

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

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

1.5.1 Vectorizing Categorical data

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

In [20]:

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

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

In [21]:

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

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

In [0]:

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

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [23]:

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

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

In [0]:

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

1.5.2.2 TFIDF vectorizer

In [25]:

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

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

1.5.2.3 Using Pretrained Models: Avg W2V

In [26]:

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

# =====
Output:

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

# =====

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

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

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

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

```

if i in words_glove:
    words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

```

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

```

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

```

```
'''
```

Out[26]:

```

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

```

In [0]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-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 [28]:

```

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

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

```

100%|██████████| 109248/109248 [00:38<00:00, 2831.04it/s]

109248
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

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

In [30]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors.append(vector)

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

```
100%|██████████| 109248/109248 [03:46<00:00, 483.36it/s]
```

```
109248
300
```

In [0]:

```
# Similarly you can vectorize for title also
```

1.5.3 Vectorizing Numerical features

In [0]:

```
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 [33]:

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

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

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

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

```
Mean : 298.1193425966608, Standard deviation : 367.49634838483496
```

In [34]:

```
price_standardized
```

Out[34]:

```
array([[ -0.3905327 ],
       [  0.00239637],
       [  0.59519138],
       ...,
       [-0.15825829],
       [-0.61243967],
       [-0.51216657]])
```

1.5.4 Merging all the above features

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

In [35]:

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

Out[35]:

```
(109248, 16663)
```

In [0]:

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

Computing Sentiment Scores

In [37]:

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

# import nltk
nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

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

```

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

```

```

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

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

```

```

/usr/local/lib/python3.6/dist-packages/nltk/twitter/__init__.py:20: UserWarning:

```

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

```

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,

```

Assignment 9: RF and GBDT

Response Coding: Example

The response label is built only on train dataset. For a category which is not there in train data and present in test data, we will encode them with default values Ex: in our test data if have State: D then we encode it as [0.5, 0.05]


1. Apply both Random Forrest and GBDT on these feature sets

- **Set 1:** categorical (instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- **Set 2:** categorical (instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(TFIDF) + preprocessed_eassay (TFIDF)
- **Set 3:** categorical (instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(AVG W2V) + preprocessed_eassay (AVG W2V)
- **Set 4:** categorical (instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(TFIDF W2V) + preprocessed_eassay (TFIDF W2V)

2. The hyper parameter tuning (Consider any two hyper parameters preferably n_estimators, max_depth)

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

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure

with X-axis as **n_estimators**, Y-axis as **max_depth**, and Z-axis as **AUC Score**, we have given the notebook which explains how to plot this 3d plot, you can find it in the same drive [3d_scatter_plot.ipynb](#)

or

- 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
[seaborn heat maps](#) with rows as **n_estimators**, columns as **max_depth**, and values inside the cell representing **AUC Score**
- You can choose either of the plotting techniques: 3d plot or heat map
- 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

4. 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. Random Forest and GBDT

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

In [38]:

```
#i'm making use of this code from sample assignment
data=project_data[['school_state','teacher_prefix','project_grade_category','teacher_number_of_previously_posted_projects','project_is_approved','clean_categories','clean_subcategories','price']]
data = data.replace(np.nan, '', regex=True)
data['essay'] =preprocessed_essays
data['title'] =preprocessed_titles
data['teacher_prefix'] = data['teacher_prefix'].str.replace(' ','No prefix')
data['project_grade_category'] = data['project_grade_category'].str.replace(' ','_')
data['project_grade_category'] = data['project_grade_category'].str.replace('-', '_')
print(data.shape)
```

(109248, 10)

In [0]:

```
y = data['project_is_approved'].values
```

In [0]:

```
X = data
# 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=0.33, stratify=y)
```

2.2 Make Data Model Ready: encoding numerical, categorical features

Response Encoding:School_state

In [0]:

```
#Response Encoding toward science,https://www.youtube.com/watch?v=zcsYYP4pXr8&list=PLC0PzjY99Q_V3b7GBh3gjmcCfe0l-408K&index=10
```



```
def get_gv_fea_dict(alpha, feature, train_df):
    value_count=train_df[feature].value_counts()
    gv_dict1=dict()
    gv_dict2=dict()
    for i,denominator in value_count.items():
        vec1=[]
        vec2=[]
        for k in range(0,2):
            cls_cnt = train_df.loc[(train_df['project_is_approved']==k) & (train_df[feature]==i)]
            vec1.append(cls_cnt.shape[0])
            vec2.append((cls_cnt.shape[0] + alpha*10)/(denominator+90*alpha)) # Here we are using alpha f
or laplase smoothning.
        gv_dict1[i]=vec1
        gv_dict2[i]=vec2
    return gv_dict2
```

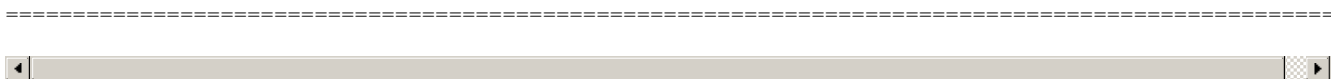
In [0]:

```
#Response Encodind towordatascience,https://www.youtube.com/watch?
v=zcsYYP4pXr8&list=PLC0PzjY99Q_V3b7GBh3gjmcCfe0l-408K&index=10
def get_gv_feature(feature,df,gv_dict):
    value_count=df[feature].value_counts()
    gv_fea=[]
    for index,row in df.iterrows():
        if row[feature] in dict(value_count).keys():
            gv_fea.append(gv_dict[row[feature]])
        else:
            gv_fea.append([1/2,1/2])
    return gv_fea
```

In [43]:

```
#Response Encodind towordatascience,https://www.youtube.com/watch?
v=zcsYYP4pXr8&list=PLC0PzjY99Q_V3b7GBh3gjmcCfe0l-408K&index=10
gv_dict=get_gv_fea_dict(1,"school_state",X_train)
X_train_state_response=np.array(get_gv_feature("school_state",X_train,gv_dict))
X_test_state_response=np.array(get_gv_feature("school_state",X_test,gv_dict))
print("After Response vectorizations:State")
print(X_train_state_response.shape, y_train.shape)
print(X_test_state_response.shape, y_test.shape)
print("=="*100)
```

After Response vectorizations:State
(73196, 2) (73196,)
(36052, 2) (36052,)



Response Encoding:teacher_prefix

In [44]:

```
#Response Encodind towordatascience,https://www.youtube.com/watch?
v=zcsYYP4pXr8&list=PLC0PzjY99Q_V3b7GBh3gjmcCfe0l-408K&index=10
gv_dict=get_gv_fea_dict(1,"teacher_prefix",X_train)
X_train_teacher_response=np.array(get_gv_feature("teacher_prefix",X_train,gv_dict))
X_test_teacher_response=np.array(get_gv_feature("teacher_prefix",X_test,gv_dict))
print("After Response vectorizations:State")
print(X_train_teacher_response.shape, y_train.shape)
print(X_test_teacher_response.shape, y_test.shape)
print("=="*100)
```

After Response vectorizations:State
(73196, 2) (73196,)
(36052, 2) (36052,)



Response Encoding:project_grade_category

In [45]:

```
#Response Encodind towordatascience,https://www.youtube.com/watch?v=zcsYYP4pXr8&list=PLC0PzjY99Q_V3b7GBh3gjmcCfe0l-408K&index=10
gv_dict=get_gv_fea_dict(1,"project_grade_category",X_train)
X_train_grade_response=np.array(get_gv_feature("project_grade_category",X_train,gv_dict))
X_test_grade_response=np.array(get_gv_feature("project_grade_category",X_test,gv_dict))
print("After Response vectorizations:State")
print(X_train_grade_response.shape, y_train.shape)
print(X_test_grade_response.shape, y_test.shape)
print("=="*100)
```

After Response vectorizations:State
(73196, 2) (73196,)
(36052, 2) (36052,)
=====



Normalizing:price

In [46]:

```
#Normalization of price
#I'm making use of the code from sample assignment
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(-1,1))

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

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

After vectorizations:price
(73196, 1) (73196,)
(36052, 1) (36052,)
=====



2.3 Make Data Model Ready: encoding eassay, and project_title

Bow Encoding:essay

In [47]:

```
#I'm making use of the code from the sample assignment.
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data

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

print("After Bow vectorizations:essay")
print(X_train_essay_bow.shape, y_train.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("=="*100)
```

After Bow vectorizations:essay
(73196, 5000) (73196,)
(36052, 5000) (36052,)

Bow Encoding:title

In [48]:

```
#I'm making use of the code from the sample assignment
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4),max_features=5000)
vectorizer.fit(X_train['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['title'].values)
X_test_title_bow = vectorizer.transform(X_test['title'].values)

print("After Bow vectorizations:title")
print(X_train_title_bow.shape, y_train.shape)
print(X_test_title_bow.shape, y_test.shape)
print("="*100)
```

```
After Bow vectorizations:title
(73196, 5000) (73196,)
(36052, 5000) (36052,)
=====
```

TFIDF Encoding:title

In [49]:

```
#I'm making use of the code from the sample assignment and second assignment
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,2),max_features=5000)
vectorizer.fit(X_train['title'].values)
# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['title'].values)
X_test_title_tfidf = vectorizer.transform(X_test['title'].values)
print("After tfidf vectorizations:title")
print(X_train_title_tfidf.shape, y_train.shape)
print(X_test_title_tfidf.shape, y_test.shape)
print("="*100)
```

```
After tfidf vectorizations:title
(73196, 5000) (73196,)
(36052, 5000) (36052,)
=====
```

TFIDF Encoding:essay

In [50]:

```
#I'm making use of the code from the sample assignment and second assignment
vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,2),max_features=5000)
vectorizer.fit(X_train['essay'].values)
# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['essay'].values)
X_test_essay_tfidf = vectorizer.transform(X_test['essay'].values)
print("After tfidf vectorizations:essay")
print(X_train_essay_tfidf.shape, y_train.shape)
print(X_test_essay_tfidf.shape, y_test.shape)
print("="*100)
```

```
After tfidf vectorizations:essay
(73196, 5000) (73196,)
(36052, 5000) (36052,)
=====
```

AVG W2V Encoding:title,essay

In [51]:

```
#I'm making use of the code from the second assignment
def fun(col):
    sample=[];
    for sentence in tqdm(col): # 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
        sample.append(vector)
    return sample
X_train_title_avg_w2v = fun(X_train['title'].values)
X_test_title_avg_w2v = fun(X_test['title'].values)# the avg-w2v for each sentence/review is stored
in this list
X_train_essay_avg_w2v = fun(X_train['essay'].values)
X_test_essay_avg_w2v = fun(X_test['essay'].values)
print("After AVG W2V vectorizations:essay,title")
print('\n',len(X_train_title_avg_w2v),len(X_train_title_avg_w2v[0]),y_train.shape)
print(len(X_test_title_avg_w2v),len(X_test_title_avg_w2v[0]),y_test.shape)
print("="*100)
print('\n',len(X_train_essay_avg_w2v),len(X_train_essay_avg_w2v[0]),y_train.shape)
print(len(X_test_essay_avg_w2v),len(X_test_essay_avg_w2v[0]),y_test.shape)

100%|██████████| 73196/73196 [00:01<00:00, 63815.38it/s]
100%|██████████| 36052/36052 [00:00<00:00, 63954.33it/s]
100%|██████████| 73196/73196 [00:23<00:00, 3078.37it/s]
100%|██████████| 36052/36052 [00:11<00:00, 3124.01it/s]
```

After AVG W2V vectorizations:essay,title

```
73196 300 (73196,)
36052 300 (36052,)
```

```
73196 300 (73196,)
36052 300 (36052,)
```

TFIDF W2V Encoding:essay

In [52]:

```
#I'm making use of the code from the second assignment
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['essay'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
def fun1(col):
    sample=[];
    for sentence in tqdm(col): # 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
        sample.append(vector)
```

```

    return sample
X_train_essay_tfidf_w2v = fun1(X_train['essay'].values)
X_test_essay_tfidf_w2v = fun1(X_test['essay'].values) # the avg-w2v for each sentence/review is
stored in this list
print("After tfidf W2V vectorizations:essay")
print('\n',len(X_train_essay_tfidf_w2v),len(X_train_essay_tfidf_w2v[0]),y_train.shape)
print(len(X_test_essay_tfidf_w2v),len(X_test_essay_tfidf_w2v[0]),y_test.shape)
print("=="*100)

```

```

100%|██████████| 73196/73196 [02:25<00:00, 501.39it/s]
100%|██████████| 36052/36052 [01:11<00:00, 504.84it/s]

```

After tfidf W2V vectorizations:essay

```

73196 300 (73196,)
36052 300 (36052,)
=====

```



TFIDF W2V Encoding:title

In [53]:

```

#I'm making use of the code from the second assignment
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['title'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
def fun2(col):
    sample=[];
    for sentence in tqdm(col): # 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
            sample.append(vector)
    return sample
X_train_title_tfidf_w2v = fun2(X_train['title'].values)
X_test_title_tfidf_w2v = fun2(X_test['title'].values) # the avg-w2v for each sentence/review is
stored in this list
print("After tfidf W2V vectorizations:title")
print('\n',len(X_train_title_tfidf_w2v),len(X_train_title_tfidf_w2v[0]),y_train.shape)
print(len(X_test_title_tfidf_w2v),len(X_test_title_tfidf_w2v[0]),y_test.shape)
print("=="*100)

```

```

100%|██████████| 73196/73196 [00:03<00:00, 24311.95it/s]
100%|██████████| 36052/36052 [00:01<00:00, 26462.43it/s]

```

After tfidf W2V vectorizations:title

```

73196 300 (73196,)
36052 300 (36052,)
=====

```



Concatinating all the features

In [54]:

```

#I'm making use of the code from the sample assignment
from scipy.sparse import hstack
from scipy import sparse
X_tr1 = hstack((X_train_essay_bow,X_train_title_bow,X_train_state_response,
X_train_teacher_response, X_train_grade_response, X_train_price_norm)).tocsr()
X_te1 = hstack((X_test_essay_bow,X_test_title_bow,X_test_state_response, X_test_teacher_response, X
_test_grade_response, X_test_price_norm)).tocsr()

print("Final Data matrix for Set-1")
print(X_tr1.shape, y_train.shape)
print(X_te1.shape, y_test.shape)
print("="*100)

X_tr2 = hstack((X_train_essay_tfidf,X_train_title_tfidf,X_train_state_response,
X_train_teacher_response, X_train_grade_response, X_train_price_norm)).tocsr()
X_te2 = hstack((X_test_essay_tfidf,X_test_title_tfidf,X_test_state_response,
X_test_teacher_response, X_test_grade_response, X_test_price_norm)).tocsr()

print("Final Data matrix for Set-2")
print(X_tr2.shape, y_train.shape)
print(X_te2.shape, y_test.shape)
print("="*100)

X_tr3 =
hstack((X_train_essay_avg_w2v,X_train_title_avg_w2v,sparse.csr_matrix(X_train_state_response), spa
rse.csr_matrix(X_train_teacher_response), sparse.csr_matrix(X_train_grade_response),
X_train_price_norm)).tocsr()
X_te3 = hstack((X_test_essay_avg_w2v,X_test_title_avg_w2v,sparse.csr_matrix(X_test_state_response)
, sparse.csr_matrix(X_test_teacher_response), sparse.csr_matrix(X_test_grade_response),
X_test_price_norm)).tocsr()

print("Final Data matrix for Set-3")
print(X_tr3.shape, y_train.shape)
print(X_te3.shape, y_test.shape)
print("="*100)

X_tr4 =
hstack((X_train_essay_tfidf_w2v,X_train_title_tfidf_w2v,sparse.csr_matrix(X_train_state_response),
sparse.csr_matrix(X_train_teacher_response), sparse.csr_matrix(X_train_grade_response),
X_train_price_norm)).tocsr()
X_te4 =
hstack((X_test_essay_tfidf_w2v,X_test_title_tfidf_w2v,sparse.csr_matrix(X_test_state_response), sp
arse.csr_matrix(X_test_teacher_response), sparse.csr_matrix(X_test_grade_response),
X_test_price_norm)).tocsr()

print("Final Data matrix for Set-4")
print(X_tr4.shape, y_train.shape)
print(X_te4.shape, y_test.shape)
print("="*100)

```

Final Data matrix for Set-1

(73196, 10007) (73196,)
(36052, 10007) (36052,)

Final Data matrix for Set-2

(73196, 10007) (73196,)
(36052, 10007) (36052,)

Final Data matrix for Set-3

(73196, 607) (73196,)
(36052, 607) (36052,)

Final Data matrix for Set-4

(73196, 607) (73196,)
(36052, 607) (36052,)

2.4 Applying Random Forest

Apply Random Forest 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

2.4.1 Applying Random Forests on BOW, SET 1

Finding Hyperparameter

In [0]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
nb= RandomForestClassifier();
parameters = {'n_estimators' : [1, 4,16, 32, 100, 200],
              'max_depth' : [ 1, 10 , 20 , 25 , 30, 35],
              'class_weight':['balanced']}
clf = GridSearchCV(nb, parameters, cv=3, scoring='roc_auc',verbose=15)
clf.fit(X_tr1, y_train)
```

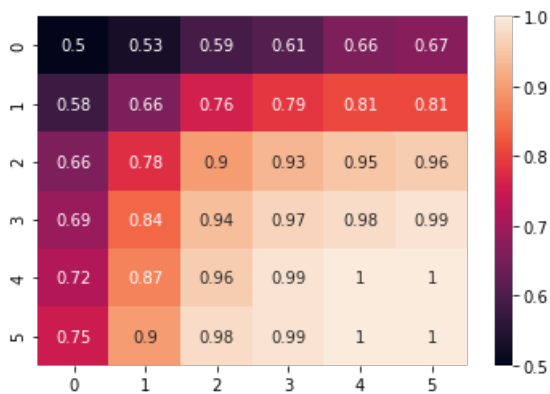
In [0]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc
train_auc=train_auc.reshape(6,6)
cv_auc=cv_auc.reshape(6,6)
import seaborn as sns
print("Train data Auc scores")
sns.heatmap(train_auc, annot=True)
```

Train data Auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca0914208>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35].
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more

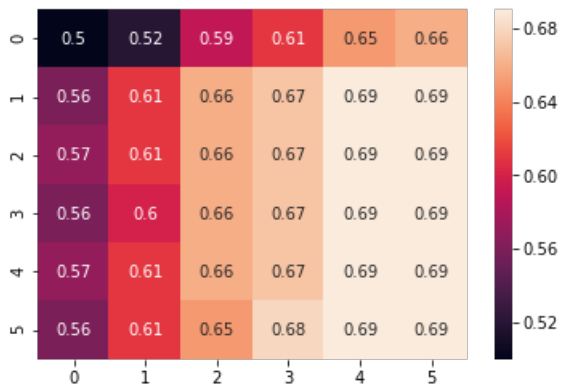
In [0]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
print("CV auc scores")
sns.heatmap(cv_auc, annot=True)
```

CV auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca0583400>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35].
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our cv AUC also Increase and it seems like its always better to take n_estimators to be more

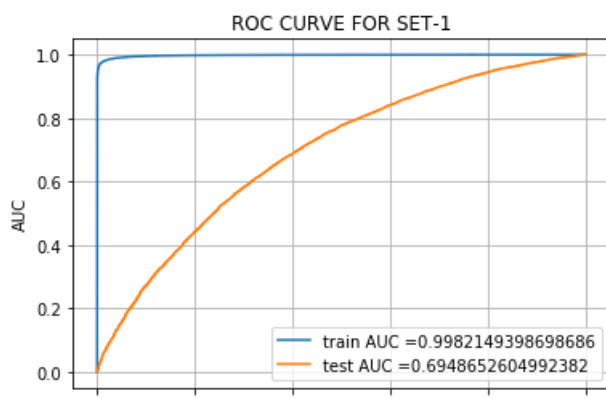
ROC Curve

In [0]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=35
best_n_estimators=200
neigh=RandomForestClassifier(n_estimators=best_n_estimators,max_depth=best_max_depth,class_weight=
'balanced');
neigh.fit(X_tr1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr1)[: ,1]
y_test_pred = neigh.predict_proba(X_te1)[: ,1]

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("alpha: hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-1")
plt.grid()
plt.show()
```



0.0 0.2 0.4 0.6 0.8 1.0
alpha: hyperparameter

Observation:

1. Here we took our Hyperparameter as max_depth=35 and n_estimators=200
2. The performance of our train-data was good with 99%
3. The performance of our train-data was good with 69%

confusion matrix

In [0]:

```
#I'm making use of the code from the sample assignment
# we will pick a threshold that will give the least fpr
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
```

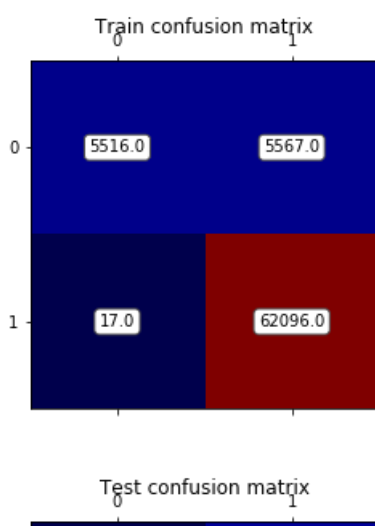
In [0]:

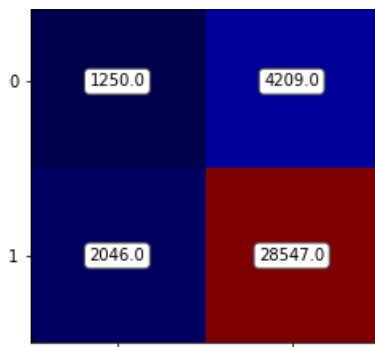
```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, neigh.predict(X_tel))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of tpr*(1-fpr) 0.24999470622237013 for threshold 0.408





2.4.2 Applying Random Forests on TFIDF, SET 2

Find Hyperparameter

In [0]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
nb= RandomForestClassifier();
parameters = {'n_estimators' : [1, 4,16, 32, 100, 200],
              'max_depth' : [ 1, 10 , 20 , 25 , 30, 35],
              'class_weight':['balanced']}
clf = GridSearchCV(nb, parameters, cv=3, scoring='roc_auc',verbose=15)
clf.fit(X_tr2, y_train)
```

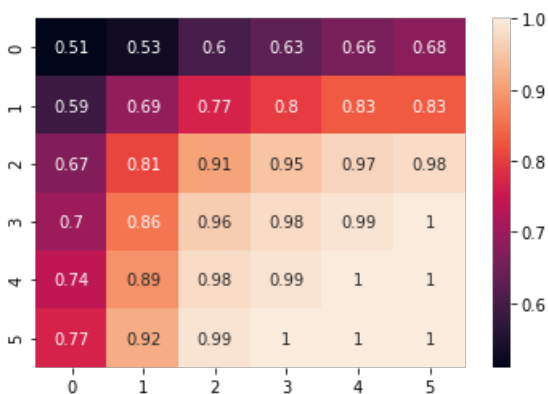
In [0]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc
train_auc=train_auc.reshape(6,6)
cv_auc=cv_auc.reshape(6,6)
import seaborn as sns
print("Train data Auc scores")
sns.heatmap(train_auc, annot=True)
```

Train data Auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca0667ef0>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35]
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because

seaborn heatmap cannot provide us the parameter.

3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more

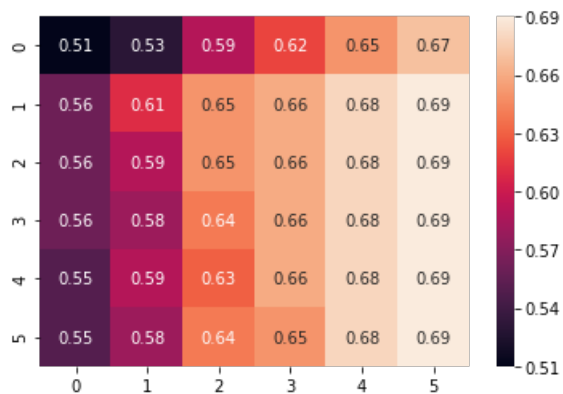
In [0]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
print("CV auc scores")
sns.heatmap(cv_auc, annot=True)
```

CV auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca0667ac8>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35].
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our cv Auc also Increase and it seems like its always better to take n_estimators to be more

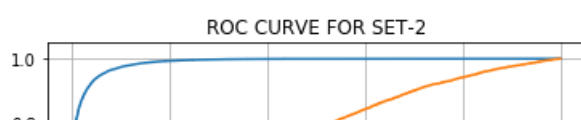
ROC Curve

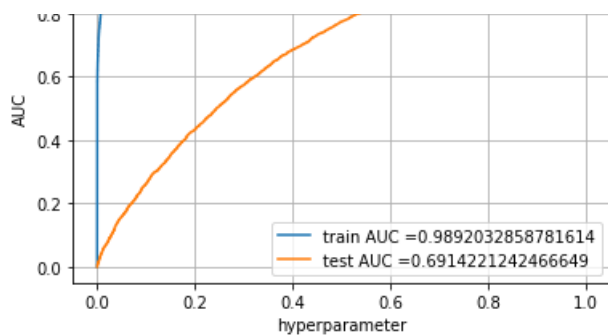
In [0]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=25
best_n_estimators=200
neigh=RandomForestClassifier(n_estimators=best_n_estimators,max_depth=best_max_depth,class_weight=
'balanced');
neigh.fit(X_tr2,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr2)[:,-1]
y_test_pred = neigh.predict_proba(X_te2)[:,-1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-2")
plt.grid()
plt.show()
```





Observation:

1. Here we took our Hyperparameter as max_depth=25 and n_estimators=200
2. The performance of our train-data was good with 99%
3. The performance of our train-data was good with 69%

confusion matrix

In [0]:

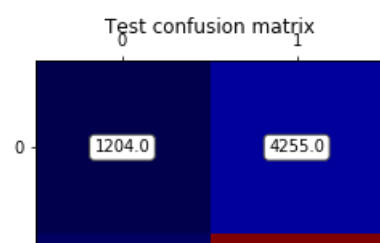
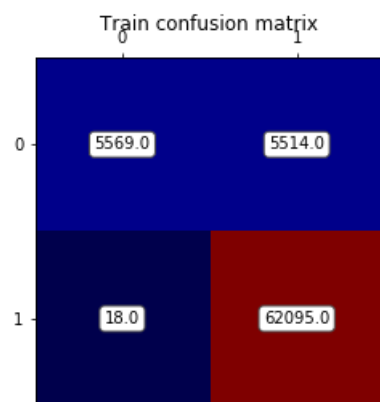
```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, neigh.predict(X_te2))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

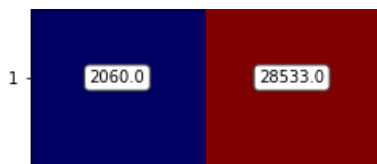
for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()

fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499938432613109 for threshold 0.438





2.4.3 Applying Random Forests on AVG W2V, SET 3

Finding Hyper Parameter

In [0]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
nb= RandomForestClassifier();
parameters = {'n_estimators': [1, 4,16, 32, 100, 200],
              'max_depth': [ 1, 10 , 20 , 25 , 30, 35],
              'class_weight':['balanced']}
clf = GridSearchCV(nb, parameters, cv=3, scoring='roc_auc',verbose=15)
clf.fit(X_tr3, y_train)
```

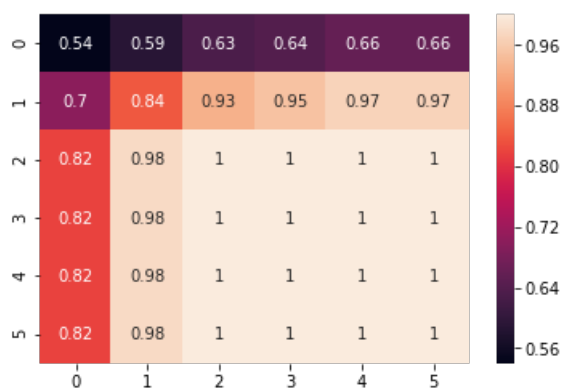
In [0]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc
train_auc=train_auc.reshape(6,6)
cv_auc=cv_auc.reshape(6,6)
import seaborn as sns
print("Train data Auc scores")
sns.heatmap(train_auc, annot=True)
```

Train data Auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca082ed30>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35].
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more

In [0]:

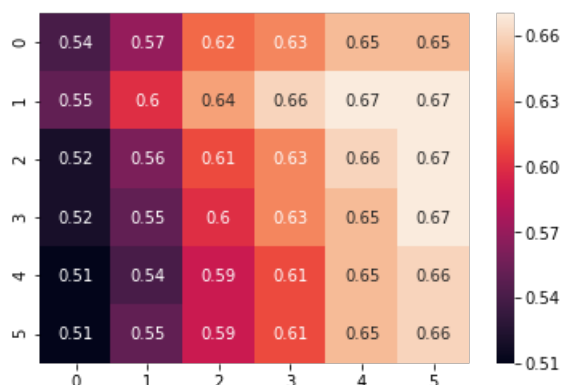
```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
print("CV auc scores")
```

```
print("CV auc scores")
sns.heatmap(cv_auc, annot=True)
```

CV auc scores

Out[0]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f4ca082e518>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35]
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more

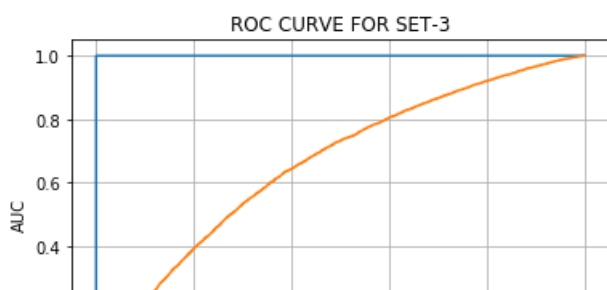
ROC Curve

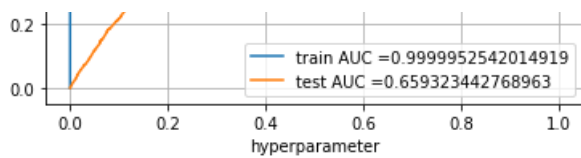
In [0]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=35
best_n_estimators=200
neigh=RandomForestClassifier(n_estimators=best_n_estimators,max_depth=best_max_depth,class_weight=
'balanced');
neigh.fit(X_tr3,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr3)[: ,1]
y_test_pred = neigh.predict_proba(X_te3)[: ,1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-3")
plt.grid()
plt.show()
```





Observation:

1. Here we took our Hyperparameter as max_depth=35 and n_estimators=200
2. The performance of our train-data was good with 99%
3. The performance of our train-data was good with 66%

confusion Matrix

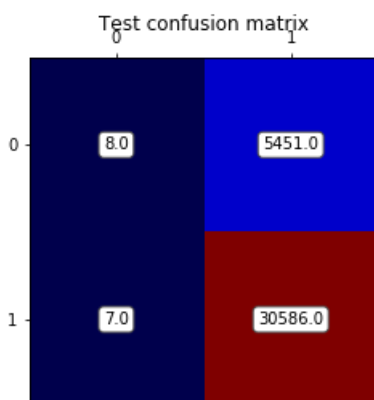
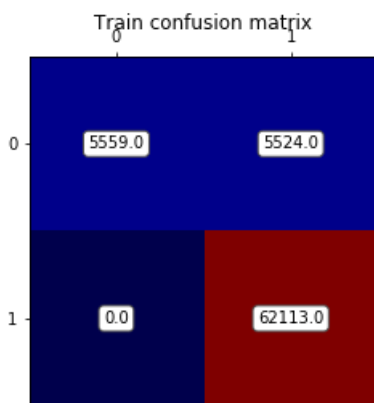
In [0]:

```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, neigh.predict(X_te3))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of $\text{tpr} \times (1 - \text{fpr})$ 0.2499975067752416 for threshold 0.305



2.4.4 Applying Random Forests on TFIDF W2V, SET 4

Finding Hyperparameter

In [0]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
nb= RandomForestClassifier();
parameters = {'n_estimators' : [1, 4,16, 32, 100, 200],
              'max_depth' : [ 1, 10 , 20 , 25 , 30, 35],
              'class_weight':['balanced']}
clf = GridSearchCV(nb, parameters, cv=3, scoring='roc_auc',verbose=15)
clf.fit(X_tr4, y_train)
```

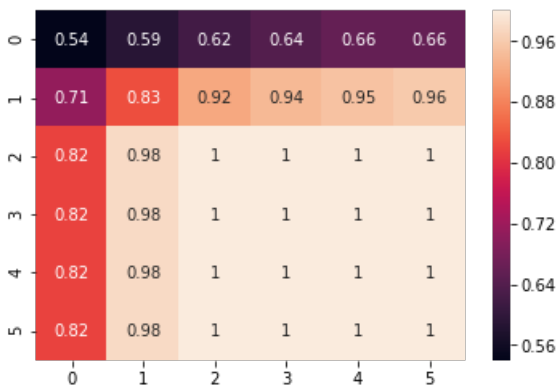
In [56]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc
train_auc=train_auc.reshape(6,6)
cv_auc=cv_auc.reshape(6,6)
import seaborn as sns
print("Train data Auc scores")
sns.heatmap(train_auc, annot=True)
```

Train data Auc scores

Out[56]:

<matplotlib.axes._subplots.AxesSubplot at 0x7fdb2fd53780>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35]
2. And The columns correspond to n_estimators [1, 4,16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more

In [57]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
print("CV auc scores")
sns.heatmap(cv_auc, annot=True)
```

CV auc scores

Out[57]:

<matplotlib.axes._subplots.AxesSubplot at 0x7fdb2fc25b00>



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1, 10 , 20 , 25 , 30, 35]
2. And The columns correspond to n_estimators [1, 4, 16, 32, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train AUC also Increase and it seems like its always better to take n_estimators to be more

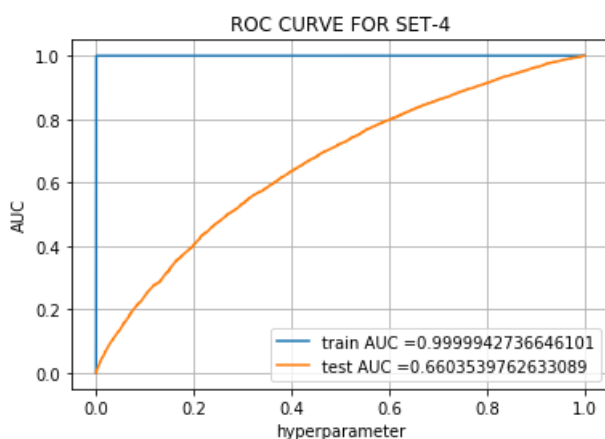
ROC Curve

In [58]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=35
best_n_estimators=200
neigh=RandomForestClassifier(n_estimators=best_n_estimators,max_depth=best_max_depth,class_weight=
'balanced');
neigh.fit(X_tr4,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr4)[: ,1]
y_test_pred = neigh.predict_proba(X_te4)[: ,1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-4")
plt.grid()
plt.show()
```



Observation:

1. Here we took our Hyperparameter as max_depth=35 and n_estimators=200
2. The performance of our train-data was good with 99%
3. The performance of our train-data was good with 66%

confusion Matrix

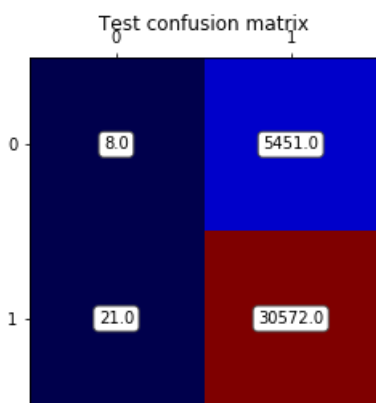
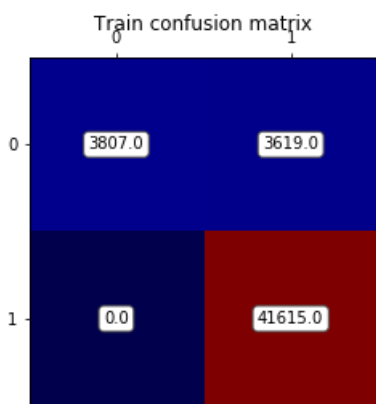
In [0]:

```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, neigh.predict(X_te4))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of tpr*(1-fpr) 0.2498397692677456 for threshold 0.305



2.5 Applying GBDT

Apply GBDT 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

2.5.1 Applying XGBOOST on BOW, SET 1

Find Hyperparameter

In [0]:

```
import xgboost as xgb
from sklearn.model_selection import GridSearchCV
parameters = {
    'max_depth': [1,3, 4, 5],
    'learning_rate': [0.1],
    'n_estimators': [1,16,64,100,200]
}
xgb_model = xgb.XGBClassifier()
clf = GridSearchCV(xgb_model, parameters, scoring = 'roc_auc', verbose=5)
clf.fit(X_train, y_train)
```

In [0]:

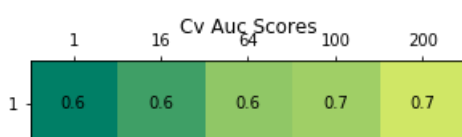
```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc=train_auc.reshape(4,5)
cv_auc=cv_auc.reshape(4,5)
#https://matplotlib.org/tutorials/colors/colormaps.html
#https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

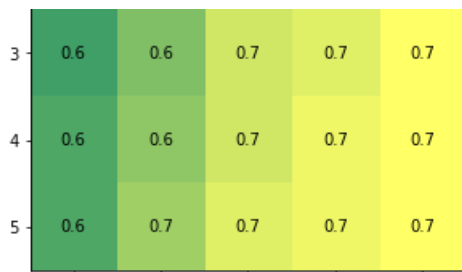
def showAucPlot(text,data):
    labels = [['1','16','64','100','200'],['1', '3','4','5']]
    fig = plt.figure()
    ax = fig.add_subplot(111)
    cax = ax.matshow(data,cmap="summer")
    #https://matplotlib.org/tutorials/colors/colormaps.html

    #https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

    for (i, j), z in np.ndenumerate(data):
        ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center')
    plt.title(text)
    ax.set_xticklabels([''] + labels[0])
    ax.set_yticklabels([''] + labels[1])
    plt.show()

showAucPlot("Train AUC Scores",train_auc)
showAucPlot("Cv AUC Scores",cv_auc)
```





Observation:

1. Here the Rows {0,4} correspond to the max_depth [1,3,4,5]
2. And The columns correspond to n_estimators [1,16,64, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train AUC also Increase and it seems like its always better to take n_estimators to be more
4. As the Max depth increases our Cv AUC also Increase and it seems like its always better to take n_estimators to be more

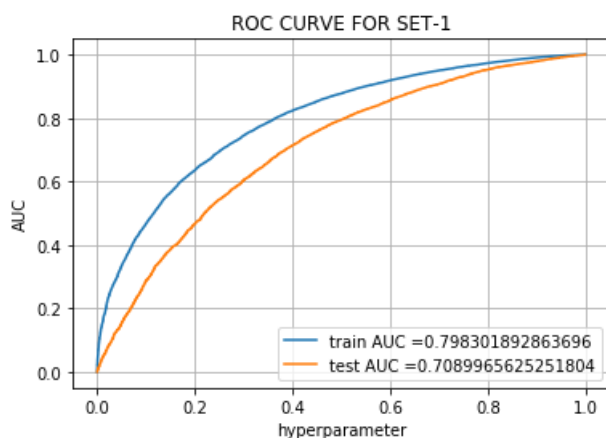
Roc Curve

In [0]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=4
best_n_estimators=200
neigh=xgb.XGBClassifier(max_depth=best_max_depth,n_estimators=best_n_estimators);
neigh.fit(X_tr1,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr1)[:,-1]
y_test_pred = neigh.predict_proba(X_tel)[:,-1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-1")
plt.grid()
plt.show()
```



Observation:

1. Here we took our Hyperparameter as max_depth=4 and n_estimators=100
2. The performance of our train-data was good with 79%
3. The performance of our train-data was good with 70%

Confusion matrix

In [0]:

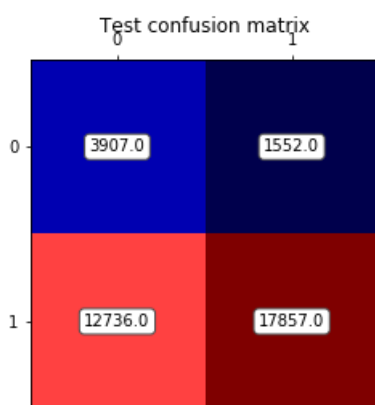
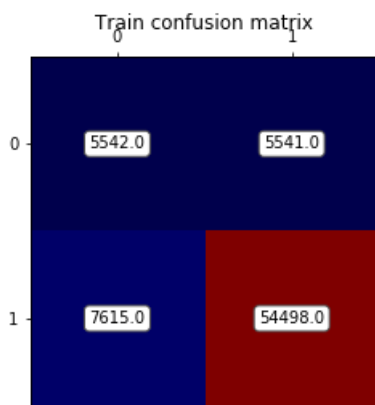
```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499999979647145 for threshold 0.783

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999999161092998 for threshold 0.854



2.5.2 Applying XGBOOST on TFIDF, SET 2

Find Hyperparameter

In [0]:

```
import xgboost as xgb
from sklearn.model_selection import GridSearchCV
parameters = {
    'max_depth': [1, 3, 4, 5]
```

```

        'max_depth': [1,3,4,5],
        'learning_rate': [0.1],
        'n_estimators': [1,16,64,100,200]
    }
xgb_model = xgb.XGBClassifier()
clf = GridSearchCV(xgb_model, parameters, scoring = 'roc_auc', verbose=5)
clf.fit(X_tr2, y_train)

```

In [0]:

```

#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc=train_auc.reshape(4,5)
cv_auc=cv_auc.reshape(4,5)
#https://matplotlib.org/tutorials/colors/colormaps.html
#https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

def showAucPlot(text,data):
    labels = [['1','16','64','100','200'], ['1', '3','4','5']]
    fig = plt.figure()
    ax = fig.add_subplot(111)
    cax = ax.matshow(data,cmap="summer")
    #https://matplotlib.org/tutorials/colors/colormaps.html

    #https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

    for (i, j), z in np.ndenumerate(data):
        ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center')
    plt.title(text)
    ax.set_xticklabels([''] + labels[0])
    ax.set_yticklabels([''] + labels[1])
    plt.show()

showAucPlot("Train Auc Scores",train_auc)
showAucPlot("Cv Auc Scores",cv_auc)

```



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1,3,4,5]
2. And The columns correspond to n_estimators [1,16,64, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train Auc also Increase and it seems like its always better to take n_estimators to be more
4. As the Max depth increases our Cv Auc also Increase and it seems like its always better to take n_estimators to be more

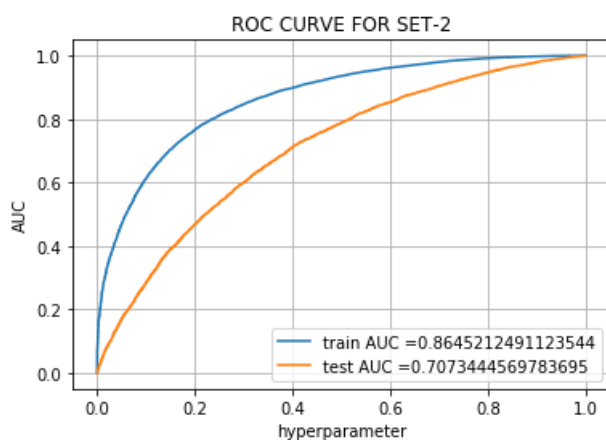
ROC Curve

In [0]:

```
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=5
best_n_estimators=200
neigh=xgb.XGBClassifier(max_depth=best_max_depth,n_estimators=best_n_estimators);
neigh.fit(X_tr2,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr2)[:,-1]
y_test_pred = neigh.predict_proba(X_te2)[:,-1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-2")
plt.grid()
plt.show()
```



Observation:

1. Here we took our Hyperparameter as max_depth=5 and n_estimators=200
2. The performance of our train-data was good with 86%
3. The performance of our train-data was good with 70%

confusion matrix

In [0]:

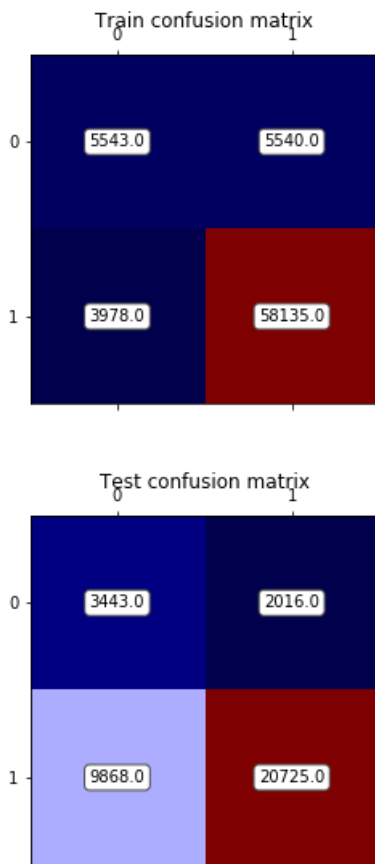
```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
```

```
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999998168243034 for threshold 0.75
the maximum value of $tpr \cdot (1 - fpr)$ 0.24999999161092998 for threshold 0.841



2.5.3 Applying XGBOOST on AVG W2V, SET 3

Finding Hyperparameter

In [0]:

```
import xgboost as xgb
from sklearn.model_selection import GridSearchCV
parameters = {
    'max_depth': [1, 3, 4, 5],
    'learning_rate': [0.1],
    'n_estimators': [1, 16, 64, 100, 200]
}
xgb_model = xgb.XGBClassifier()
clf = GridSearchCV(xgb_model, parameters, scoring = 'roc_auc', verbose=5)
clf.fit(X_tr3, y_train)
```

In [56]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
```



```

cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc=train_auc.reshape(4,5)
cv_auc=cv_auc.reshape(4,5)
#https://matplotlib.org/tutorials/colors/colormaps.html
#https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

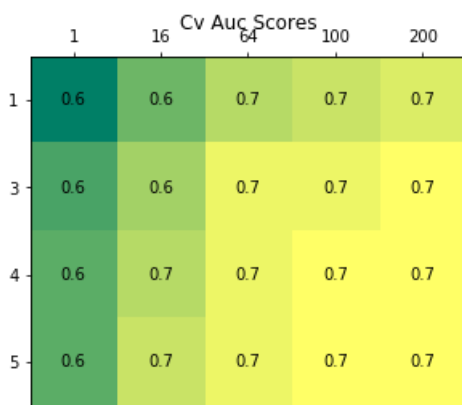
def showAucPlot(text,data):
    labels = [['1','16','64','100','200'], ['1', '3', '4', '5']]
    fig = plt.figure()
    ax = fig.add_subplot(111)
    cax = ax.matshow(data,cmap="summer")
    #https://matplotlib.org/tutorials/colors/colormaps.html

    #https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

    for (i, j), z in np.ndenumerate(data):
        ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center')
    plt.title(text)
    ax.set_xticklabels([''] + labels[0])
    ax.set_yticklabels([''] + labels[1])
    plt.show()

showAucPlot("Train Auc Scores",train_auc)
showAucPlot("Cv Auc Scores",cv_auc)

```



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1,3,4,5]
2. And The columns correspond to n_estimators [1,16,64, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train AUC also Increase and it seems like its always better to take n_estimators to be more
4. As the Max depth increases our Cv AUC also Increase and it seems like its always better to take n_estimators to be more

ROC Curve

In [64]:

```

#I'm making use of the code from the sample assignment
import xgboost as xgb
from sklearn.metrics import roc_curve, auc

```

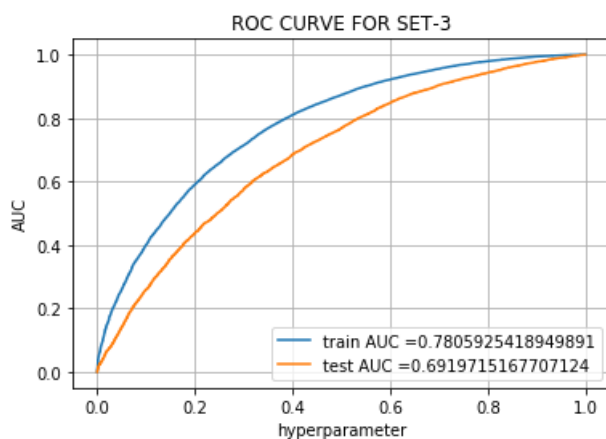
```

best_max_depth=4
best_n_estimators=100
neigh=xgb.XGBClassifier(max_depth=best_max_depth,n_estimators=best_n_estimators);
neigh.fit(X_tr3,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr3)[:,-1]
y_test_pred = neigh.predict_proba(X_te3)[:,-1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-3")
plt.grid()
plt.show()

```



Observation:

1. Here we took our Hyperparameter as max_depth=5 and n_estimators=200
2. The performance of our train-data was good with 78%
3. The performance of our train-data was good with 69%

confusion matrix

In [65]:

```

from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

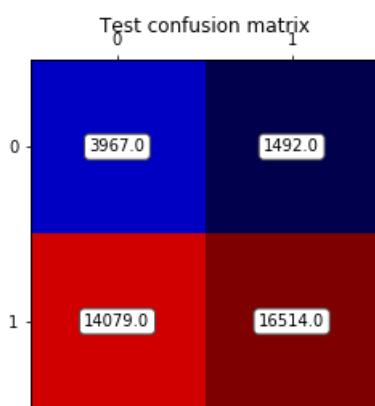
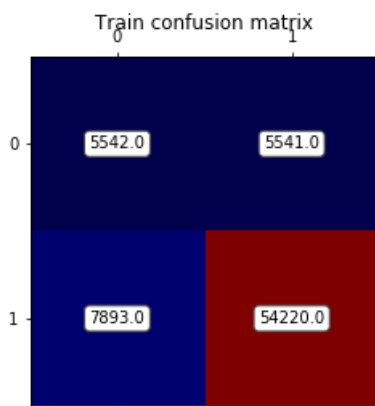
for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()

```

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499999979647145 for threshold 0.786

the maximum value of $\bar{tpr} \cdot (1 - \bar{fpr})$ 0.24999999161092998 for threshold 0.87



2.5.4 Applying XGBOOST on TFIDF W2V, SET 4

Finding Hyperparameter

In [0]:

```
import xgboost as xgb
from sklearn.model_selection import GridSearchCV
parameters = {
    'max_depth': [1, 3, 4, 5],
    'learning_rate': [0.1],
    'n_estimators': [1, 16, 64, 100, 200]
}
xgb_model = xgb.XGBClassifier()
clf = GridSearchCV(xgb_model, parameters, scoring = 'roc_auc', verbose=5)
clf.fit(X_tr4, y_train)
```

In [58]:

```
#https://seaborn.pydata.org/generated/seaborn.heatmap.html
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']
train_auc=np.around(train_auc, decimals=2, out=None)
cv_auc = np.around(cv_auc, decimals=2, out=None)
train_auc=train_auc.reshape(4,5)
cv_auc=cv_auc.reshape(4,5)
#https://matplotlib.org/tutorials/colors/colormaps.html
#https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib

def showAucPlot(text,data):
    labels = [['1', '16', '64', '100', '200'], ['1', '3', '4', '5']]
    fig = plt.figure()
    ax = fig.add_subplot(111)
    cax = ax.matshow(data, cmap="summer")
    #https://matplotlib.org/tutorials/colors/colormaps.html

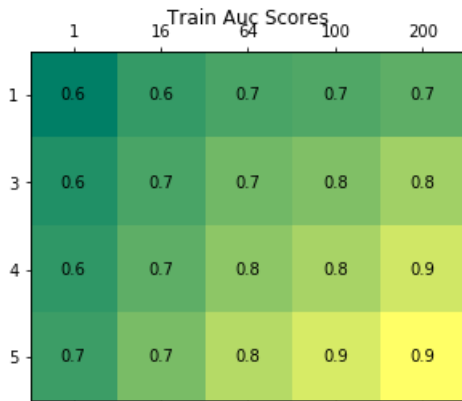
    #https://stackoverflow.com/questions/20998083/show-the-values-in-the-grid-using-matplotlib
```

```

for (i, j), z in np.ndenumerate(data):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center')
plt.title(text)
ax.set_xticklabels([''] + labels[0])
ax.set_yticklabels([''] + labels[1])
plt.show()

showAucPlot("Train AUC Scores", train_auc)
showAucPlot("Cv AUC Scores", cv_auc)

```



Observation:

1. Here the Rows {0,4} correspond to the max_depth [1,3,4,5]
2. And The columns correspond to n_estimators [1,16,64, 100, 200] as We can not show that directly there because seaborn heatmap cannot provide us the parameter.
3. As the Max depth increases our train AUC also Increase and it seems like its always better to take n_estimators to be more
4. As the Max depth increases our Cv AUC also Increase and it seems like its always better to take n_estimators to be more

ROC Curve

In [62]:

```

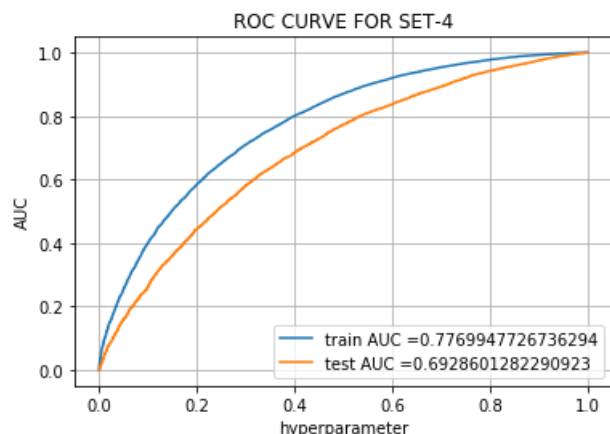
#I'm making use of the code from the sample assignment
from sklearn.metrics import roc_curve, auc
best_max_depth=4
best_n_estimators=100
neigh=xgb.XGBClassifier(max_depth=best_max_depth,n_estimators=best_n_estimators);
neigh.fit(X_tr4,y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
#https://github.com/scikit-learn/scikit-learn/blob/master/examples/model_selection/plot_roc.py
y_train_pred = neigh.predict_proba(X_tr4)[:,-1]
y_test_pred = neigh.predict_proba(X_te4)[:,-1]

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("hyperparameter")
plt.ylabel("AUC")
plt.title("ROC CURVE FOR SET-4")
plt.grid()
plt.show()
```



Observation:

1. Here we took our Hyperparameter as max_depth=4 and n_estimators=100
2. The performance of our train-data was good with 94.3%
3. The performance of our train-data was good with 69.2%

confusion matrix

In [63]:

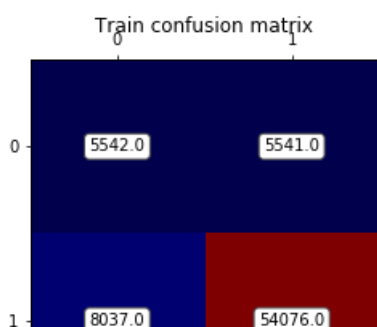
```
from sklearn.metrics import confusion_matrix
array1=confusion_matrix(y_train, predict(y_train_pred, tr_thresholds, train_fpr, train_fpr))
array2=confusion_matrix(y_test, predict(y_test_pred, tr_thresholds, test_fpr, test_fpr))
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array1, cmap='seismic')

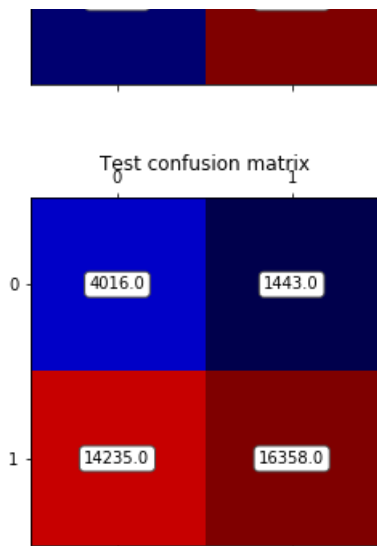
for (i, j), z in np.ndenumerate(array1):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Train confusion matrix')
plt.show()
fig, ax = plt.subplots()
# Using matshow here just because it sets the ticks up nicely. imshow is faster.
ax.matshow(array2, cmap='seismic')

for (i, j), z in np.ndenumerate(array2):
    ax.text(j, i, '{:0.1f}'.format(z), ha='center', va='center',
            bbox=dict(boxstyle='round', facecolor='white', edgecolor='0.3'))
plt.title('Test confusion matrix')
plt.show()
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.2499999979647145 for threshold 0.782

the maximum value of $tpr \cdot (1 - fpr)$ 0.24999999161092998 for threshold 0.869





3. Conclusion

In [66]:

```
# Please compare all your models using Prettytable library
# pretty table http://zetcode.com/python/prettytable/
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "max_depth", "n_estimators", "AUC"]
x.add_row(["BOW", "Random Forest", 35, 200, 70.0])
x.add_row(["TFIDF", "Random Forest", 25, 200, 69.6])
x.add_row(["AVGW2V", "Random Forest", 35, 200, 65.9])
x.add_row(["TFIDF WEIGHTED W2V", "Random Forest", 35, 200, 66.2])
x.add_row(["BOW", "XGBOOST", 4, 100, 70.8])
x.add_row(["TFIDF", "XGBOOST", 5, 200, 70.7])
x.add_row(["AVGW2V", "XGBOOST", 4, 100, 69.2])
x.add_row(["TFIDF WEIGHTED W2V", "XGBOOST", 4, 100, 69.2])
print(x)
```

Vectorizer	Model	max_depth	n_estimators	AUC
BOW	Random Forest	35	200	70.0
TFIDF	Random Forest	25	200	69.6
AVGW2V	Random Forest	35	200	65.9
TFIDF WEIGHTED W2V	Random Forest	35	200	66.2
BOW	XGBOOST	4	100	70.8
TFIDF	XGBOOST	5	200	70.7
AVGW2V	XGBOOST	4	100	69.2
TFIDF WEIGHTED W2V	XGBOOST	4	100	69.2

Final Observations:

1. Instead of encoding our Categorical features with oneHot Encoding we did it with Response encoding
2. We have 4 sets of data for which we have built both Random Forest and XGBoost on train data, find hyperparameters on cross validation data and test AUC on Test data
3. For all the 8 sets we did this process and plotted the AUC on a Seaborn HeatMap.
4. The rows indicate the max_depth and the columns indicate the n_estimators
5. And we noticed that as the Max depth increases our train AUC also increases and it seems like it's always better to take n_estimators should be more.
6. With Random Forest classifier set1(BOW), set2(TFIDF) seems to perform well.
7. Whereas with XGBoost classifier (GBDT) all of them performed almost the same but set2,4 did quite well with one percent