

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>	<ul style="list-style-type: none">••	Title of the project. Examples: <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care & Hunger</code> <code>Health & Sports</code> <code>History & Civics</code> <code>Literacy & Language</code> <code>Math & Science</code> <code>Music & The Arts</code> <code>Special Needs</code> <code>Warmth</code> Examples: <ul style="list-style-type: none">• <code>Music & The Arts</code>• <code>Literacy & Language, Math & Science</code>
<code>school_state</code>		State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: <code>Literacy</code> <code>Literature & Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

Feature	Description
<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."
- `__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 [0]:

```
#from google.colab import drive
#drive.mount('/content/drive')
```

In [0]:

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

```

import squites
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

!pip install chart_studio

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

```

```

Requirement already satisfied: chart_studio in /usr/local/lib/python3.6/dist-packages (1.0.0)
Requirement already satisfied: retrying>=1.3.3 in /usr/local/lib/python3.6/dist-packages (from
chart_studio) (1.3.3)
Requirement already satisfied: plotly in /usr/local/lib/python3.6/dist-packages (from
chart_studio) (4.1.1)
Requirement already satisfied: six in /usr/local/lib/python3.6/dist-packages (from chart_studio)
(1.12.0)
Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-packages (from
chart_studio) (2.21.0)
Requirement already satisfied: chardet<3.1.0,>=3.0.2 in /usr/local/lib/python3.6/dist-packages
(from requests->chart_studio) (3.0.4)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.6/dist-packages (from
requests->chart_studio) (2019.9.11)
Requirement already satisfied: urllib3<1.25,>=1.21.1 in /usr/local/lib/python3.6/dist-packages
(from requests->chart_studio) (1.24.3)
Requirement already satisfied: idna<2.9,>=2.5 in /usr/local/lib/python3.6/dist-packages (from
requests->chart_studio) (2.8)

```

1.1 Reading Data

In [0]:

```

project_data = pd.read_csv('/content/drive/My
Drive/Mass3/Assignments_DonorsChoose_2018/train_data.csv')
resource_data = pd.read_csv('/content/drive/My
Drive/Mass3/Assignments_DonorsChoose_2018/resources.csv')

```

In [0]:

```

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 [0]:

```
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[0]:

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

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

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

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

1.3 preprocessing of project_subject_subcategories

In [0]:

```
sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python:
```

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039

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

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

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

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

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

1.3 Text preprocessing

In [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 [0]:

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

Out[0]:

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

In [0]:

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

In [0]:

```
# 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\n\r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity. My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still.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\n\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an \"open classroom\" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.\r\n\r\nYour generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.\r\n\r\nIt costs a lot of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!nnnn

=====

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\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 [0]:

```
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\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 [0]:

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

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

In [0]:

```
# remove special 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 The want to be able to move as they learn or so they say Wobble chairs are the answer and I love them because they develop their core which enhances gross motor and in turn fine motor skills They also want to learn through games my kids do not want to sit and do worksheets They want to learn to count by jumping and playing Physical engagement is the key to our success The number toss and color and shape mats can make that happen My students will forget 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 [0]:

```
# Combining all the above stundents
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 [00:54<00:00, 1987.74it/s]

In [0]:

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

Out[0]:

'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 de serves nannan'

In [0]:

```
project_data["essays"] = preprocessed_essays

project_data.drop(['essay'], axis=1, inplace=True)
project_data.drop(['project_essay_1'], axis=1, inplace=True)
project_data.drop(['project_essay_2'], axis=1, inplace=True)
project_data.drop(['project_essay_3'], axis=1, inplace=True)
project_data.drop(['project_essay_4'], axis=1, inplace=True)
```

In [0]:

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

Out[0]:

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

In [0]:

In [0]:

```
#Count of number of words in essay column
#https://stackoverflow.com/questions/49984905/count-number-of-words-per-row

project_data['count_project_essay']=project_data.essays.apply(lambda x: len(x.split()))
project_data.head()
```

Out[0]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Grade
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17	Grades P
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Grades P

1.4 Preprocessing of `project_title`

In [0]:

```
# Combining all the above statemennts
preprocessed_title = []
# 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('\n', ' ')
    sent = sent.replace('\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_title.append(sent.lower().strip())
```

100%|██████████| 109248/109248 [00:02<00:00, 41847.38it/s]

In [0]:

```
# after preprocesing
preprocessed_title[0:4]
```

Out[0]:

```
['educational support english learners home',
 'wanted projector hungry learners',
 'soccer equipment awesome middle school students',
 ...]
```

```
'techie kindergarteners']
```

In [0]:

```
project_data["project_title"] = preprocessed_title
```

In [0]:

```
#Count of number of words in title column
#https://stackoverflow.com/questions/49984905/count-number-of-words-per-row

project_data['count_project_title']=project_data.project_title.apply(lambda x: len(x.split()))
project_data.head()
```

Out[0]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades P
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grade
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Grade
3	45	p246581	f3cb9bffba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17	Grades P
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Grades P

Preprocessing of 'Project_grade_category'

In [0]:

```
project_data['project_grade_category']=project_data['project_grade_category'].replace('Grades PreK-2','Grades_PreK_2')
project_data['project_grade_category']=project_data['project_grade_category'].replace('Grades 3-5','Grades_3_5')
project_data['project_grade_category']=project_data['project_grade_category'].replace('Grades 6-8','Grades_6_8')
project_data['project_grade_category']=project_data['project_grade_category'].replace('Grades 9-12','Grades_9_12')
```

1.5 Preparing data for models

In [0]:

```
project_data.columns
```

Out[0]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'project_submitted_datetime', 'project_grade_category', 'project_title',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essays',
      'count_project_essay', 'count_project_title'],
      dtype='object')
```

we are going to consider

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

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

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

1.5.1 Vectorizing Categorical data

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

In [0]:

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

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

In [0]:

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

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

State Dictionary

In [0]:

```
from collections import Counter
my_counter1 = Counter()
```

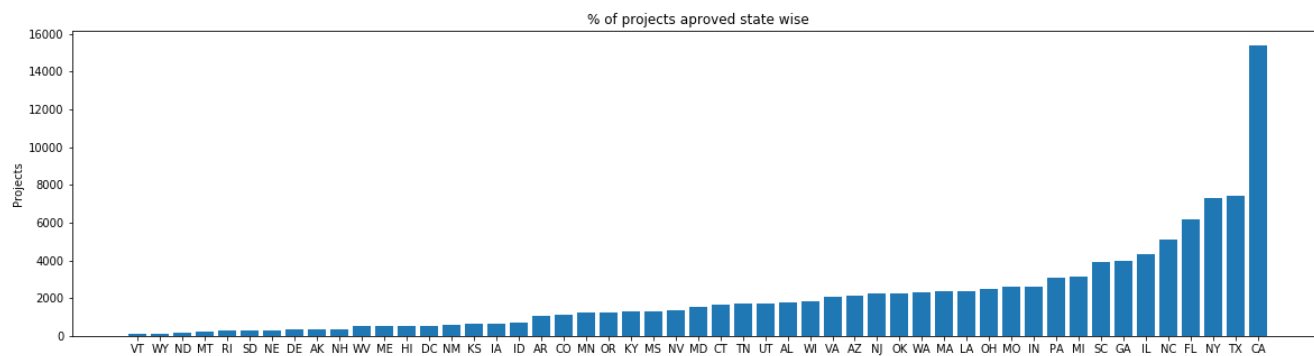
```

for word in project_data['school_state'].values:
    my_counter1.update(word.split())

state_dict = dict(my_counter1)
sorted_state_dict = dict(sorted(state_dict.items(), key=lambda kv: kv[1]))
ind1 = np.arange(len(sorted_state_dict))
plt.figure(figsize=(20,5))
p2 = plt.bar(ind1, list(sorted_state_dict.values()))
plt.ylabel('Projects')
plt.title('% of projects aproved state wise')
plt.xticks(ind1, list(sorted_state_dict.keys()))
plt.show()

for i, j in sorted_state_dict.items():
    print("{:20} {:10}".format(i,j))

```



VT	:	80
WY	:	98
ND	:	143
MT	:	245
RI	:	285
SD	:	300
NE	:	309
DE	:	343
AK	:	345
NH	:	348
WV	:	503
ME	:	505
HI	:	507
DC	:	516
NM	:	557
KS	:	634
IA	:	666
ID	:	693
AR	:	1049
CO	:	1111
MN	:	1208
OR	:	1242
KY	:	1304
MS	:	1323
NV	:	1367
MD	:	1514
CT	:	1663
TN	:	1688
UT	:	1731
AL	:	1762
WI	:	1827
VA	:	2045
AZ	:	2147
NJ	:	2237
OK	:	2276
WA	:	2334
MA	:	2389
LA	:	2394
OH	:	2467
MO	:	2576
IN	:	2620
PA	:	3109
MI	:	3161
SC	:	3936
GA	:	3963

```
IL          :      4350
NC          :      5091
FL          :      6185
NY          :      7318
TX          :      7396
CA          :     15388
```

State

In [0]:

```
vectorizer = CountVectorizer(vocabulary=list(sorted_state_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(project_data['school_state'].values)
print(vectorizer.get_feature_names())
```

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

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

```
Shape of matrix after one hot encodig (109248, 51)
```

Teacher_prefix Dictionary

In [0]:

```
project_data['teacher_prefix']=project_data['teacher_prefix'].fillna(method='ffill')
```

In [0]:

```
project_data['teacher_prefix'][30368:30372]
```

Out[0]:

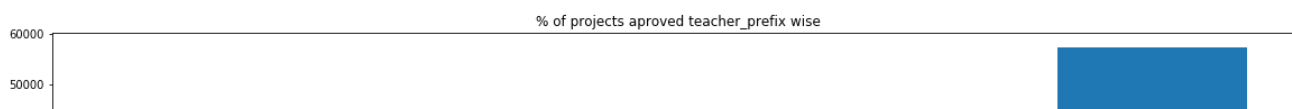
```
30368    Mrs.
30369    Mrs.
30370     Ms.
30371     Ms.
Name: teacher_prefix, dtype: object
```

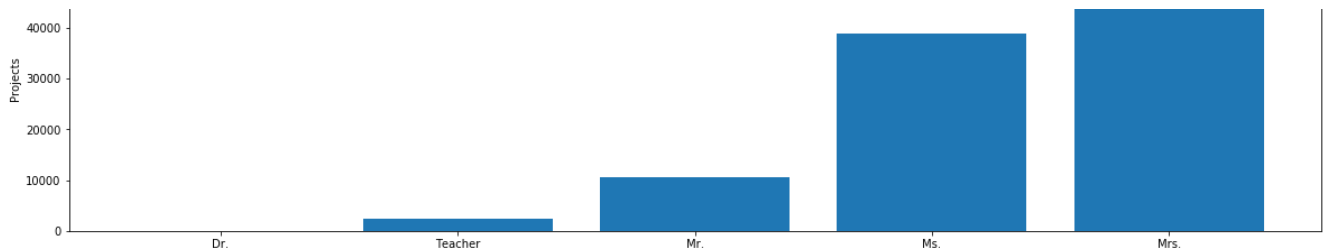
In [0]:

```
from collections import Counter
my_counter2 = Counter()
for word in project_data['teacher_prefix'].values:
    my_counter2.update(word.split())

prefix_dict = dict(my_counter2)
sorted_prefix_dict = dict(sorted(prefix_dict.items(), key=lambda kv: kv[1]))
ind2 = np.arange(len(sorted_prefix_dict))
plt.figure(figsize=(20,5))
p2 = plt.bar(ind2, list(sorted_prefix_dict.values()))
plt.ylabel('Projects')
plt.title('% of projects aproved teacher_prefix wise')
plt.xticks(ind2, list(sorted_prefix_dict.keys()))
plt.show()

for i, j in sorted_prefix_dict.items():
    print("{:20} {:10}".format(i,j))
```





```
Dr.          :      13
Teacher      :     2360
Mr.          :    10648
Ms.          :    38956
Mrs.         :   57271
```

In [0]:

```
#Teacher_prefix

vectorizer = CountVectorizer(vocabulary=list(sorted_prefix_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(project_data['teacher_prefix'].values)
print(vectorizer.get_feature_names())

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

['Dr.', 'Teacher', 'Mr.', 'Ms.', 'Mrs.']
Shape of matrix after one hot encodig  (109248, 5)
```

In [0]:

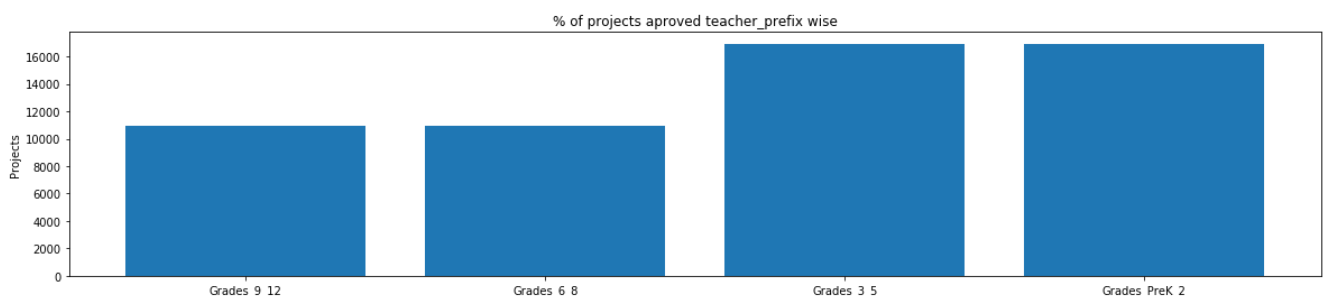
```
#project_grade_category Dictionary

#How_to_update_keys: https://www.geeksforgeeks.org/python-ways-to-change-keys-in-dictionary/

from collections import Counter
my_counter3 = Counter()
for word in project_data['project_grade_category'].values:
    my_counter3.update(word.split('_'))

grade_dict = dict(my_counter3)
sorted_grade_dict = dict(sorted(grade_dict.items(), key=lambda kv: kv[1]))
dict((sorted_grade_dict[key], value) for (key, value) in sorted_grade_dict.items())
updt_keys = ['Grades_9_12', 'Grades_6_8', 'Grades_3_5', 'Grades_PreK_2']
sorted_grade_dict = dict(zip(updt_keys, list(sorted_grade_dict.values())))
ind3 = np.arange(len(sorted_grade_dict))
plt.figure(figsize=(20,4))
p3 = plt.bar(ind3, list(sorted_grade_dict.values()))
plt.ylabel('Projects')
plt.title('% of projects aproved teacher_prefix wise')
plt.xticks(ind3, list(sorted_grade_dict.keys()))
plt.show()

for i, j in sorted_grade_dict.items():
    print("{:20} {:10}".format(i,j))
```



```
Grades_9_12      :    10963
Grades_6_8       :    10963
Grades_3_5       :    16923
Grades_PreK_2    :    16923
```

In [0]:

```
#project_grade_category
vectorizer = CountVectorizer(vocabulary=list(sorted_grade_dict.keys()), lowercase=False, binary=True)
vectorizer.fit(project_data['project_grade_category'].values)
print(vectorizer.get_feature_names())

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

```
['Grades_9_12', 'Grades_6_8', 'Grades_3_5', 'Grades_PreK_2']
Shape of matrix after one hot encoding  (109248, 4)
```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [0]:

```
# We are considering only the words which appeared in at least 10 documents (rows or projects).
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_essays)
print("Shape of matrix after one hot 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
preprocessed_title = []
# 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('\n', ' ')
    sent = sent.replace('\n', ' ')
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
    # https://gist.github.com/sebleier/554280
    sent = ' '.join(e for e in sent.split() if e not in stopwords)
    preprocessed_title.append(sent.lower().strip())
```

100%|██████████| 109248/109248 [00:02<00:00, 46775.49it/s]

In [0]:

```
# Similarly you can vectorize for title also
vectorizer = CountVectorizer(min_df=10)
text_bow = vectorizer.fit_transform(preprocessed_title)
print("Shape of matrix after one hot encoding ", text_bow.shape)
```

Shape of matrix after one hot encoding (109248, 3222)

1.5.2.2 TFIDF vectorizer

In [0]:


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

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

1.5.2.3 Using Pretrained Models: Avg W2V

In [0]:

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

# =====
Output:

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

# =====

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

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

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

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

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

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

'''
```

Out[0]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\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        model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
```

```

loadGloveModel('\\glove.42B.300d.txt\\')\n\n# =====\nOutput:\n    \nLoading Glove Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n# =====\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split('\\' \\'))\n\nfor i in preprocod_titles:\n    words.extend(i.split('\\' \\'))\n\nprint("all the words in the corpus", len(words))\n\nwords = set(words)\n\nprint("the unique words in the corpus", len(words))\n\ninter_words = set(model.keys()).intersection(words)\n\nprint("The number of words that are present in both glove vectors and our corpus", len(inter_words), "\n", np.round(len(inter_words)/len(words)*100,3), "%")\n\nwords_corpus = {}\n\nwords_glove = set(model.keys())\n\nfor i in words:\n    if i in words_glove:\n        words_corpus[i] = model[i]\n\nprint("word 2 vec length", len(words_corpus))\n\n\n# stronging variables into pickle files python : http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pickle\n\nwith open('\\glove_vectors\\', '\\wb\\') as f:\n    pickle.dump(words_corpus, f)\n\n\n
```

In [0]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('/content/drive/My Drive/Mass3/Assignments_DonorsChoose_2018/glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

In [0]:

```

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

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

```

100%|██████████| 109248/109248 [00:37<00:00, 2938.28it/s]

109248
300

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [0]:

```

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

```

In [0]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    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 [03:56<00:00, 462.02it/s]

109248
300

Using Pretrained Models: TFIDF weighted W2V on project_title

In [0]:

```

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

```

In [0]:

```

# Similarly you can vectorize for title also

tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_title): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            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 [00:04<00:00, 26541.31it/s]

109248
300

In [0]:

```

# Conversion of a list to sparse matrix

from scipy.sparse import coo_matrix
tfidf_w2v_matrix=np.reshape(np.asarray(tfidf_w2v_vectors), (109248,300))
sparse_tfidf_w2v_matrix=coo_matrix(tfidf_w2v_matrix).tocsr()

print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix.shape)

```

Shape of matrix after one hot encoding (109248, 300)

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 [0]:

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

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

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

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

Mean : 298.1193425966608, Standard deviation : 367.49634838483496

In [0]:

```
price_standardized
```

Out[0]:

```
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 categorical, text, numerical vectors

In [0]:

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

```
(109248, 9)
(109248, 30)
(109248, 3222)
(109248, 1)
```

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
```

```
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
X.shape
```

Out[0]:

```
(109248, 3262)
```

Computing Sentiment Scores

In [0]:

```
import nltk
nltk.download('vader_lexicon')
from nltk.sentiment.vader import SentimentIntensityAnalyzer

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

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students with the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelligences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of different backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a caring community of successful \
learners which can be seen through collaborative student project based learning in and out of the classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role play in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food i will take their idea \
and create common core cooking lessons where we learn important math and writing concepts while cooking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into making the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project would 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 enjoyment for healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

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

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

```
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
neg: 0.01, neu: 0.745, pos: 0.245, compound: 0.9975,
```

Assignment 5: Logistic Regression

1. [Task-1] Logistic Regression(either SGDClassifier with log loss, or LogisticRegression) on these feature sets

- **Set 1:** categorical, numerical features + project_title(BOW) + preprocessed_essay ('BOW with bi-grams' with 'min_df=10' and 'max_features=5000')

- **Set 2:** categorical, numerical features + project_title(TFIDF)+ preprocessed_eassay (`TFIDF with bi-grams` with `min_df=10` and `max_features=5000`)
- **Set 3:** categorical, numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project_title(TFIDF W2V)+ preprocessed_essay (TFIDF W2V)

2. Hyper paramter tuning (find best hyper parameters corresponding the algorithm that you choose)

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

3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.
- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points. Please visualize your confusion matrices using [seaborn heatmaps](#).

4. [\[Task-2\] Apply Logistic Regression on the below feature set Set 5 by finding the best hyper parameter as suggested in step 2 and step 3.](#)

5. [Consider these set of features Set 5 :](#)

- [school_state](#) : categorical data
- [clean_categories](#) : categorical data
- [clean_subcategories](#) : categorical data
- [project_grade_category](#) :categorical data
- [teacher_prefix](#) : categorical data
- [quantity](#) : numerical data
- [teacher_number_of_previously_posted_projects](#) : numerical data
- [price](#) : numerical data
- [sentiment score's of each of the essay](#) : numerical data
- [number of words in the title](#) : numerical data
- [number of words in the combine essays](#) : numerical data

[And apply the Logistic regression on these features by finding the best hyper paramter as suggested in step 2 and step 3.](#)

6. [Conclusion](#)

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

Note: Data Leakage

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

2. Logistic Regression

In [0]:

```
#Choosing top 50k datapoints due to memory issue.

project_data=project_data.loc[:49999,:]
project_data.shape
```

Out [0]:

```
(50000, 18)
```

In [0]:

```

sid=SentimentIntensityAnalyzer()

positive_sid= []
negative_sid = []
neutral_sid = []
compound_sid = []
for i in tqdm(project_data['essays']):
    positive_sid.append(sid.polarity_scores(i)['pos'])
    negative_sid.append(sid.polarity_scores(i)['neg'])
    neutral_sid.append(sid.polarity_scores(i)['neu'])
    compound_sid.append(sid.polarity_scores(i)['compound'])

```

100%|██████████| 50000/50000 [06:08<00:00, 135.80it/s]

In [0]:

```

#Converting List to dataframe
#https://stackoverflow.com/questions/42049147/convert-list-to-pandas-dataframe-column

project_data['essay_pos']=positive_sid
project_data['essay_neg']=negative_sid
project_data['essay_neu']=neutral_sid
project_data['essay_comp']=compound_sid
project_data.head()

```

Out[0]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_cat
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57	Grades_Pr
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr.	FL	2016-10-25 09:22:10	Grades
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56	Grades
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17	Grades_Pr
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09	Grades_Pr

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

In [0]:

```

y = project_data['project_is_approved'].values
X = project_data.drop(['project_is_approved'], axis=1)
X.head(1)

```

Out[0]:

Unnamed: 0 id teacher_id teacher_prefix school_state project_submitted_datetime project_grade_category

0 160221 p253737 c90749f5d961ff158d4b4d1e7dc665fc Mrs. IN 2016-12-05 13:43:57 Grades_Prel

In [0]:

```
# 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)
X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train)
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

2.2 Make Data Model Ready: encoding numerical, categorical features

Encoding categorical features: School State

In [0]:

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

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

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

state_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(state_pref_features))
print(state_pref_features)
```

```
After vectorizations
(22445, 51) (22445,)
(11055, 51) (11055,)
(16500, 51) (16500,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'k
s', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm',
'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv
', 'wy']
```

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

Encoding categorical features: teacher_prefix

In [0]:

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

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

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

teacher_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(teacher_pref_features))
print(teacher_pref_features)
```

```
After vectorizations
(22445, 5) (22445,)
(11055, 5) (11055,)
(16500, 5) (16500,)
['dr', 'mr', 'mrs', 'ms', 'teacher']
=====
```

```
size of important feature list: 5
['dr', 'mr', 'mrs', 'ms', 'teacher']
```

Encoding categorical features: project_grade_category

In [0]:

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

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

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

grade_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(grade_pref_features))
print(grade_pref_features)
```

```
After vectorizations
(22445, 4) (22445,)
(11055, 4) (11055,)
(16500, 4) (16500,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
=====
```

```
size of important feature list: 4
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
```

Encoding categorical features: clean_categories

In [0]:

```

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

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

print("After vectorizations")
print(X_train_clean_cat_ohe.shape, y_train.shape)
print(X_cv_clean_cat_ohe.shape, y_cv.shape)
print(X_test_clean_cat_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

clean_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(clean_pref_features))
print(clean_pref_features)

```

```

After vectorizations
(22445, 9) (22445,)
(11055, 9) (11055,)
(16500, 9) (16500,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']
=====

```

```

size of important feature list: 9
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']

```

Encoding categorical features: clean_subcategories

In [0]:

```

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

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

print("After vectorizations")
print(X_train_clean_subcat_ohe.shape, y_train.shape)
print(X_cv_clean_subcat_ohe.shape, y_cv.shape)
print(X_test_clean_subcat_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)

clean_sub_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(clean_sub_pref_features))
print(clean_sub_pref_features)

```

```

After vectorizations
(22445, 30) (22445,)
(11055, 30) (11055,)
(16500, 30) (16500,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience',
'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness',
'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'm
athematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socia
lsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
=====

```

```

size of important feature list: 30
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience',
'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness',
'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'm
athematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socia
lsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']

```

```
atnemantics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
```

Encoding numerical features: Price

In [0]:

```
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_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))

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

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_price_norm=X_train_price_norm.reshape(-1,1)
X_cv_price_norm=X_cv_price_norm.reshape(-1,1)
X_test_price_norm=X_test_price_norm.reshape(-1,1)

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

```
After vectorizations
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: teacher_number_of_previously_posted_projects

In [0]:

```
normalizer = Normalizer()

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

X_train_prev_posted = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_cv_prev_posted = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
X_test_prev_posted = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_prev_posted.shape, y_train.shape)
```

```
print(X_cv_prev_posted.shape, y_cv.shape)
print(X_test_prev_posted.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_prev_posted=X_train_prev_posted.reshape(-1,1)
X_cv_prev_posted=X_cv_prev_posted.reshape(-1,1)
X_test_prev_posted=X_test_prev_posted.reshape(-1,1)
```

```
print("After vectorizations")
print(X_train_prev_posted.shape, y_train.shape)
print(X_cv_prev_posted.shape, y_cv.shape)
print(X_test_prev_posted.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: quantity

In [0]:

```
normalizer = Normalizer()

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

X_train_quantity = normalizer.transform(X_train['quantity'].values.reshape(1,-1))
X_cv_quantity = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
X_test_quantity = normalizer.transform(X_test['quantity'].values.reshape(1,-1))
```

```
print("After vectorizations")
print(X_train_quantity.shape, y_train.shape)
print(X_cv_quantity.shape, y_cv.shape)
print(X_test_quantity.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_quantity=X_train_quantity.reshape(-1,1)
X_cv_quantity=X_cv_quantity.reshape(-1,1)
X_test_quantity=X_test_quantity.reshape(-1,1)
```

```
print("After vectorizations")
print(X_train_quantity.shape, y_train.shape)
print(X_cv_quantity.shape, y_cv.shape)
print(X_test_quantity.shape, y_test.shape)
print("="*100)
```

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

```
(16500, 1) (16500,)
```

Encoding numerical features: Essay_pos

In [0]:

```
normalizer = Normalizer()

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

X_train_pos = normalizer.transform(X_train['essay_pos'].values.reshape(1,-1))
X_cv_pos = normalizer.transform(X_cv['essay_pos'].values.reshape(1,-1))
X_test_pos = normalizer.transform(X_test['essay_pos'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_pos.shape, y_train.shape)
print(X_cv_pos.shape, y_cv.shape)
print(X_test_pos.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_pos=X_train_pos.reshape(-1,1)
X_cv_pos=X_cv_pos.reshape(-1,1)
X_test_pos=X_test_pos.reshape(-1,1)

print("After vectorizations")
print(X_train_pos.shape, y_train.shape)
print(X_cv_pos.shape, y_cv.shape)
print(X_test_pos.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: Essay_neg

In [0]:

```
normalizer = Normalizer()

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

X_train_neg = normalizer.transform(X_train['essay_neg'].values.reshape(1,-1))
X_cv_neg = normalizer.transform(X_cv['essay_neg'].values.reshape(1,-1))
X_test_neg = normalizer.transform(X_test['essay_neg'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_neg.shape, y_train.shape)
print(X_cv_neg.shape, y_cv.shape)
print(X_test_neg.shape, y_test.shape)
print("="*100)
```

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

```
17, 10000, 10000,)
```

In [0]:

```
X_train_neg=X_train_neg.reshape(-1,1)
X_cv_neg=X_cv_neg.reshape(-1,1)
X_test_neg=X_test_neg.reshape(-1,1)

print("After vectorizations")
print(X_train_neg.shape, y_train.shape)
print(X_cv_neg.shape, y_cv.shape)
print(X_test_neg.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: Essay_neu

In [0]:

```
normalizer = Normalizer()

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

X_train_neu = normalizer.transform(X_train['essay_neu'].values.reshape(1,-1))
X_cv_neu = normalizer.transform(X_cv['essay_neu'].values.reshape(1,-1))
X_test_neu = normalizer.transform(X_test['essay_neu'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_neu.shape, y_train.shape)
print(X_cv_neu.shape, y_cv.shape)
print(X_test_neu.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_neu=X_train_neu.reshape(-1,1)
X_cv_neu=X_cv_neu.reshape(-1,1)
X_test_neu=X_test_neu.reshape(-1,1)

print("After vectorizations")
print(X_train_neu.shape, y_train.shape)
print(X_cv_neu.shape, y_cv.shape)
print(X_test_neu.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: Essay_comp

In [0]:

```
normalizer = Normalizer()

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

X_train_comp = normalizer.transform(X_train['essay_comp'].values.reshape(1,-1))
X_cv_comp = normalizer.transform(X_cv['essay_comp'].values.reshape(1,-1))
X_test_comp = normalizer.transform(X_test['essay_comp'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_comp.shape, y_train.shape)
print(X_cv_comp.shape, y_cv.shape)
print(X_test_comp.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
=====
```

In [0]:

```
X_train_comp=X_train_comp.reshape(-1,1)
X_cv_comp=X_cv_comp.reshape(-1,1)
X_test_comp=X_test_comp.reshape(-1,1)

print("After vectorizations")
print(X_train_comp.shape, y_train.shape)
print(X_cv_comp.shape, y_cv.shape)
print(X_test_comp.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
=====
```

Encoding numerical features: Count words essay

In [0]:

```
normalizer = Normalizer()

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

X_train_count_essay = normalizer.transform(X_train['count_project_essay'].values.reshape(1,-1))
X_cv_count_essay = normalizer.transform(X_cv['count_project_essay'].values.reshape(1,-1))
X_test_count_essay = normalizer.transform(X_test['count_project_essay'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_count_essay.shape, y_train.shape)
print(X_cv_count_essay.shape, y_cv.shape)
print(X_test_count_essay.shape, y_test.shape)
print("=="*100)
```

```
After vectorizations
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
=====
```

In [0]:

```
X_train_count_essay=X_train_count_essay.reshape(-1,1)
X_cv_count_essay=X_cv_count_essay.reshape(-1,1)
```

```
X_cv_count_essay=X_cv_count_essay.reshape(-1,1)
X_test_count_essay=X_test_count_essay.reshape(-1,1)

print("After vectorizations")
print(X_train_count_essay.shape, y_train.shape)
print(X_cv_count_essay.shape, y_cv.shape)
print(X_test_count_essay.shape, y_test.shape)
print("=="*100)
```

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

Encoding numerical features: Count words title

In [0]:

```
normalizer = Normalizer()

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

X_train_count_title = normalizer.transform(X_train['count_project_title'].values.reshape(1,-1))
X_cv_count_title = normalizer.transform(X_cv['count_project_title'].values.reshape(1,-1))
X_test_count_title = normalizer.transform(X_test['count_project_title'].values.reshape(1,-1))

print("After vectorizations")
print(X_train_count_title.shape, y_train.shape)
print(X_cv_count_title.shape, y_cv.shape)
print(X_test_count_title.shape, y_test.shape)
print("=="*100)
```

After vectorizations

```
(1, 22445) (22445,)
(1, 11055) (11055,)
(1, 16500) (16500,)
```

In [0]:

```
X_train_count_title=X_train_count_title.reshape(-1,1)
X_cv_count_title=X_cv_count_title.reshape(-1,1)
X_test_count_title=X_test_count_title.reshape(-1,1)

print("After vectorizations")
print(X_train_count_title.shape, y_train.shape)
print(X_cv_count_title.shape, y_cv.shape)
print(X_test_count_title.shape, y_test.shape)
print("=="*100)
```

After vectorizations

```
(22445, 1) (22445,)
(11055, 1) (11055,)
(16500, 1) (16500,)
```

2.3 Make Data Model Ready: encoding eassay, and project_title

2.3.1 Bag of Words

In [0]:

```
#PROJECT_TITLE
```



```

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

print("="*100)

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

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

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

bow_title_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(bow_title_pref_features))

```

```

(22445, 21) (22445,)
(11055, 21) (11055,)
(16500, 21) (16500,)
=====

```

```

After vectorizations
(22445, 2009) (22445,)
(11055, 2009) (11055,)
(16500, 2009) (16500,)
=====

```

```
size of important feature list: 2009
```

In [0]:

```

#ESSAY

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

print("="*100)

vectorizer = CountVectorizer(min_df=10,ngram_range=(2,2), max_features=5000)
vectorizer.fit(X_train['essays'].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['essays'].values)
X_cv_essay_bow = vectorizer.transform(X_cv['essays'].values)
X_test_essay_bow = vectorizer.transform(X_test['essays'].values)

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

bow_essay_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(bow_essay_pref_features))

```

```

(22445, 21) (22445,)
(11055, 21) (11055,)
(16500, 21) (16500,)
=====

```

```

After vectorizations
(22445, 5000) (22445,)
(11055, 5000) (11055,)

```

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

```
size of important feature list: 5000
```

2.3.2 TFIDF

```
In [0]:
```

```
#ESSAY
```

```
In [0]:
```

```
vectorizer = TfidfVectorizer(min_df=10, ngram_range=(2,2), max_features=5000)

vectorizer.fit(X_train['essays'])
essay_tfidf_train=vectorizer.transform(X_train['essays'])

print("Shape of matrix after one hot encoding ",essay_tfidf_train.shape)
print("="*100)

essay_tfidf_cv=vectorizer.transform(X_cv['essays'])

print("Shape of matrix after one hot encoding ",essay_tfidf_cv.shape)
print("="*100)

essay_tfidf_test=vectorizer.transform(X_test['essays'])

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

tfidf_essay_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(tfidf_essay_pref_features))
```

```
Shape of matrix after one hot encoding (22445, 5000)
```

```
Shape of matrix after one hot encoding (11055, 5000)
```

```
Shape of matrix after one hot encoding (16500, 5000)
size of important feature list: 5000
```

```
In [0]:
```

```
#PROJECT_TITLE
```

```
In [0]:
```

```
vectorizer = TfidfVectorizer(min_df=10)

vectorizer.fit(X_train['project_title'])
title_tfidf_train=vectorizer.transform(X_train['project_title'])

print("Shape of matrix after one hot encoding ",title_tfidf_train.shape)
print("="*100)

title_tfidf_cv=vectorizer.transform(X_cv['project_title'])

print("Shape of matrix after one hot encoding ",title_tfidf_cv.shape)
print("="*100)

title_tfidf_test=vectorizer.transform(X_test['project_title'])
```

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

tfidf_title_pref_features= vectorizer.get_feature_names()

print('size of important feature list:', len(tfidf_title_pref_features))
```

Shape of matrix after one hot encoding (22445, 1249)

Shape of matrix after one hot encoding (11055, 1249)

Shape of matrix after one hot encoding (16500, 1249)
size of important feature list: 1249

2.3.3 AVG W2V

In [0]:

```
#ESSAY

# 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('/content/drive/My Drive/Mass3/Assignments_DonorsChoose_2018/glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

In [0]:

```
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_essay_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay_train.append(vector)

print(len(avg_w2v_essay_train))
print(len(avg_w2v_essay_train[0]))
```

100%|██████████| 22445/22445 [00:07<00:00, 3062.75it/s]

22445

300

In [0]:

```
avg_w2v_essay_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay_cv.append(vector)
```

100%|██████████| 11055/11055 [00:03<00:00, 3017.08it/s]

In [0]:

```
avg_w2v_essay_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_essay_test.append(vector)
```

100%|██████████| 16500/16500 [00:05<00:00, 2997.67it/s]

In [0]:

```
#PROJECT_TITLE

# average Word2Vec
# compute average word2vec for each review.
avg_w2v_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title_train.append(vector)

print(len(avg_w2v_title_train))
print(len(avg_w2v_title_train[0]))
```

100%|██████████| 22445/22445 [00:00<00:00, 58491.19it/s]

22445
300

In [0]:

```
avg_w2v_title_cv = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_cv['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_title_cv.append(vector)
```

100%|██████████| 11055/11055 [00:00<00:00, 59306.14it/s]

In [0]:

```
avg_w2v_title_test = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
```

```

if cnt_words != 0:
    vector /= cnt_words
avg_w2v_title_test.append(vector)

```

100%|██████████| 16500/16500 [00:00<00:00, 58281.84it/s]

2.3.4 TFIDF W2V

In [0]:

```

#ESSAY

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

# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())
# average Word2Vec
# compute average word2vec for each review.
def tf_idf_done(word_list):
    train_title_tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
    for sentence in tqdm(word_list): # 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.wv[word]
                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()))
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf
            if tf_idf_weight != 0:
                vector /= tf_idf_weight
                train_title_tfidf_w2v_vectors.append(vector)
        print(len(train_title_tfidf_w2v_vectors))
        print(len(train_title_tfidf_w2v_vectors[0]))
    return train_title_tfidf_w2v_vectors

```

In [0]:

```
tfidf_w2v_matrix_essay_train=tf_idf_done(X_train['essays'])
```

100%|██████████| 22445/22445 [00:46<00:00, 482.45it/s]

22445
300

In [0]:

```

tfidf_w2v_matrix_essay_train=np.reshape(np.asarray(tf_idf_done(X_train['essays'])), (22445,300))
sparse_tfidf_w2v_matrix_essay_train=coo_matrix(tfidf_w2v_matrix_essay_train).tocsr()

print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_essay_train.shape)
print("="*100)

```

100%|██████████| 22445/22445 [00:46<00:00, 484.24it/s]

22445
300
Shape of matrix after one hot encoding (22445, 300)

=====

In [0]:

```
tfidf_w2v_matrix_essay_cv=tf_idf_done(X_cv['essays'])
```

100%|██████████| 11055/11055 [00:22<00:00, 486.49it/s]

11055
300

In [0]:

```
tfidf_w2v_matrix_essay_cv=np.reshape(np.asarray(tf_idf_done(X_cv['essays'])), (11055,300))  
sparse_tfidf_w2v_matrix_essay_cv=coo_matrix(tfidf_w2v_matrix_essay_cv).tocsr()
```

```
print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_essay_cv.shape)  
print("="*100)
```

100%|██████████| 11055/11055 [00:22<00:00, 486.81it/s]

11055
300
Shape of matrix after one hot encoding (11055, 300)
=====

In [0]:

```
tfidf_w2v_matrix_essay_test=tf_idf_done(X_test['essays'])
```

100%|██████████| 16500/16500 [00:34<00:00, 481.43it/s]

16500
300

In [0]:

```
tfidf_w2v_matrix_essay_test=np.reshape(np.asarray(tf_idf_done(X_test['essays'])), (16500,300))  
sparse_tfidf_w2v_matrix_essay_test=coo_matrix(tfidf_w2v_matrix_essay_test).tocsr()
```

```
print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_essay_test.shape)
```

100%|██████████| 16500/16500 [00:33<00:00, 487.07it/s]

16500
300
Shape of matrix after one hot encoding (16500, 300)

In [0]:

```
#PROJECT_TITLE
```

In [0]:

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

```
# we are converting a dictionary with word as a key, and the idf as a value  
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))  
tfidf_words = set(tfidf_model.get_feature_names())  
# average Word2Vec  
# compute average word2vec for each review.
```

```
tfidf_w2v_matrix_title_train=tf_idf_done(X_train['project_title'])
```

```
100%|██████████| 22445/22445 [00:01<00:00, 19918.22it/s]
```

```
22445  
300
```

In [0]:

```
tfidf_w2v_matrix_title_train=np.reshape(np.asarray(tf_idf_done(X_train['project_title'])), (22445,300))  
sparse_tfidf_w2v_matrix_title_train=coo_matrix(tfidf_w2v_matrix_title_train).tocsr()  
  
print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_title_train.shape)  
print("="*100)
```

```
100%|██████████| 22445/22445 [00:01<00:00, 20835.29it/s]
```

```
22445  
300  
Shape of matrix after one hot encoding (22445, 300)  
=====
```



In [0]:

```
tfidf_w2v_matrix_title_cv=tf_idf_done(X_cv['project_title'])
```

```
100%|██████████| 11055/11055 [00:00<00:00, 32644.05it/s]
```

```
11055  
300
```

In [0]:

```
tfidf_w2v_matrix_title_cv=np.reshape(np.asarray(tf_idf_done(X_cv['project_title'])), (11055,300))  
sparse_tfidf_w2v_matrix_title_cv=coo_matrix(tfidf_w2v_matrix_title_cv).tocsr()  
  
print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_title_cv.shape)  
print("="*100)
```

```
100%|██████████| 11055/11055 [00:00<00:00, 28147.28it/s]
```

```
11055  
300  
Shape of matrix after one hot encoding (11055, 300)  
=====
```



In [0]:

```
tfidf_w2v_matrix_title_test=tf_idf_done(X_test['project_title'])
```

```
100%|██████████| 16500/16500 [00:00<00:00, 33030.85it/s]
```

```
16500  
300
```

In [0]:

```
tfidf_w2v_matrix_title_test=np.reshape(np.asarray(tf_idf_done(X_test['project_title'])), (16500,300))
```

```
sparse_tfidf_w2v_matrix_title_test=coo_matrix(tfidf_w2v_matrix_title_test).tocsr()

print("Shape of matrix after one hot encoding ",sparse_tfidf_w2v_matrix_title_test.shape)
```

```
100%|██████████| 16500/16500 [00:00<00:00, 28161.60it/s]
```

```
16500
300
Shape of matrix after one hot encoding (16500, 300)
```

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

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

Concatinating the features

Set 1:

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr1 = hstack((X_train_essay_bow,
X_train_title_bow,X_train_clean_cat_ohe,X_train_clean_subcat_ohe, X_train_state_ohe,
X_train_teacher_ohe, X_train_grade_ohe, X_train_price_norm,X_train_prev_posted)).tocsr()
X_cv1 = hstack((X_cv_essay_bow, X_cv_title_bow,X_cv_clean_cat_ohe,X_cv_clean_subcat_ohe, X_cv_state
_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_price_norm,X_cv_prev_posted)).tocsr()
X_te1 = hstack((X_test_essay_bow, X_test_title_bow,X_test_clean_cat_ohe,X_test_clean_subcat_ohe,
X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe,
X_test_price_norm,X_test_prev_posted)).tocsr()

print("Final Data matrix")
print(X_tr1.shape, y_train.shape)
print(X_cv1.shape, y_cv.shape)
print(X_te1.shape, y_test.shape)
print("=="*100)
```

```
Final Data matrix
(22445, 7110) (22445,)
(11055, 7110) (11055,)
(16500, 7110) (16500,)
```



Set 2:

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr2 = hstack((essay_tfidf_train,
title_tfidf_train,X_train_clean_cat_ohe,X_train_clean_subcat_ohe, X_train_state_ohe,
X_train_teacher_ohe, X_train_grade_ohe, X_train_price_norm,X_train_prev_posted)).tocsr()
X_cv2 = hstack((essay_tfidf_cv, title_tfidf_cv,X_cv_clean_cat_ohe,X_cv_clean_subcat_ohe, X_cv_state
_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_price_norm,X_cv_prev_posted)).tocsr()
X_te2 = hstack((essay_tfidf_test, title_tfidf_test,X_test_clean_cat_ohe,X_test_clean_subcat_ohe,
X_test_state_ohe, X_test_teacher_ohe, X_test_grade_ohe,
X_test_price_norm,X_test_prev_posted)).tocsr()

print("Final Data matrix")
print(X_tr2.shape, y_train.shape)
print(X_cv2.shape, y_cv.shape)
print(X_te2.shape, y_test.shape)
print("=="*100)
```



```
Final Data matrix
(22445, 6350) (22445,)
(11055, 6350) (11055,)
(16500, 6350) (16500,)
=====
```

Set 3:

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
X_tr3 =
hstack((avg_w2v_essay_train,avg_w2v_title_train,X_train_clean_cat_oh, X_train_clean_subcat_oh,
X_train_state_oh, X_train_teacher_oh, X_train_grade_oh, X_train_price_norm,X_train_prev_posted)
).tocsr()
X_cv3 = hstack((avg_w2v_essay_cv,avg_w2v_title_cv,X_cv_clean_cat_oh,X_cv_clean_subcat_oh,
X_cv_state_oh, X_cv_teacher_oh, X_cv_grade_oh, X_cv_price_norm,X_cv_prev_posted)).tocsr()
X_te3 = hstack((avg_w2v_essay_test,avg_w2v_title_test,X_test_clean_cat_oh,X_test_clean_subcat_oh
, X_test_state_oh, X_test_teacher_oh, X_test_grade_oh,
X_test_price_norm,X_test_prev_posted)).tocsr()

print("Final Data matrix")
print(X_tr3.shape, y_train.shape)
print(X_cv3.shape, y_cv.shape)
print(X_te3.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(22445, 701) (22445,)
(11055, 701) (11055,)
(16500, 701) (16500,)
=====
```

Set 4:

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr4 =
hstack((sparse_tfidf_w2v_matrix_essay_train,sparse_tfidf_w2v_matrix_title_train,X_train_clean_cat_oh,X_train_clean_subcat_oh, X_train_state_oh, X_train_teacher_oh, X_train_grade_oh,
X_train_price_norm,X_train_prev_posted)).tocsr()
X_cv4 =
hstack((sparse_tfidf_w2v_matrix_essay_cv,sparse_tfidf_w2v_matrix_title_cv,X_cv_clean_cat_oh,X_cv_clean_subcat_oh, X_cv_state_oh, X_cv_teacher_oh, X_cv_grade_oh, X_cv_price_norm,X_cv_prev_posted)
).tocsr()
X_te4 =
hstack((sparse_tfidf_w2v_matrix_essay_test,sparse_tfidf_w2v_matrix_title_test,X_test_clean_cat_oh,X_test_clean_subcat_oh, X_test_state_oh, X_test_teacher_oh, X_test_grade_oh,
X_test_price_norm,X_test_prev_posted)).tocsr()

print("Final Data matrix")
print(X_tr4.shape, y_train.shape)
print(X_cv4.shape, y_cv.shape)
print(X_te4.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(22445, 701) (22445,)
(11055, 701) (11055,)
(16500, 701) (16500,)
=====
```

Set 5:

In [0]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr5 = hstack((X_train_count_title,X_train_count_essay,X_train_pos,X_train_neg,X_train_neu,X_train_comp,X_train_quantity,X_train_clean_cat_ohe,X_train_clean_subcat_ohe, X_train_state_ohe,
X_train_teacher_ohe, X_train_grade_ohe, X_train_price_norm,X_train_prev_posted)).tocsr()
X_cv5 =
hstack((X_cv_count_title,X_cv_count_essay,X_cv_pos,X_cv_neg,X_cv_neu,X_cv_comp,X_cv_quantity,X_cv_clean_cat_ohe,X_cv_clean_subcat_ohe, X_cv_state_ohe, X_cv_teacher_ohe, X_cv_grade_ohe, X_cv_price_norm,X_cv_prev_posted)).tocsr()
X_te5 = hstack((X_test_count_title,X_test_count_essay,X_test_pos,X_test_neg,X_test_neu,X_test_comp,X_test_quantity,X_test_clean_cat_ohe,X_test_clean_subcat_ohe, X_test_state_ohe,
X_test_teacher_ohe, X_test_grade_ohe, X_test_price_norm,X_test_prev_posted)).tocsr()

print("Final Data matrix")
print(X_tr5.shape, y_train.shape)
print(X_cv5.shape, y_cv.shape)
print(X_te5.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix
(22445, 108) (22445,)
(11055, 108) (11055,)
(16500, 108) (16500,)
=====
```

2.4.1 Applying Logistic Regression on BOW, SET 1

In [0]:

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

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

    return y_data_pred
```

In [0]:

```
from sklearn.linear_model import LogisticRegression
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score
import math

train_auc = []
cv_auc = []

K = [10**-5, 10**-4, 10**-3, 10**-2, 10**-1,10**0,10**1,10**2]
for i in tqdm(K):
    neigh = LogisticRegression(penalty='l1',C=i, class_weight='balanced')
    neigh.fit(X_tr1, y_train)

    y_train_pred = batch_predict(neigh, X_tr1)
    y_cv_pred = batch_predict(neigh, X_cv1)

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train,y_train_pred))
```

```

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

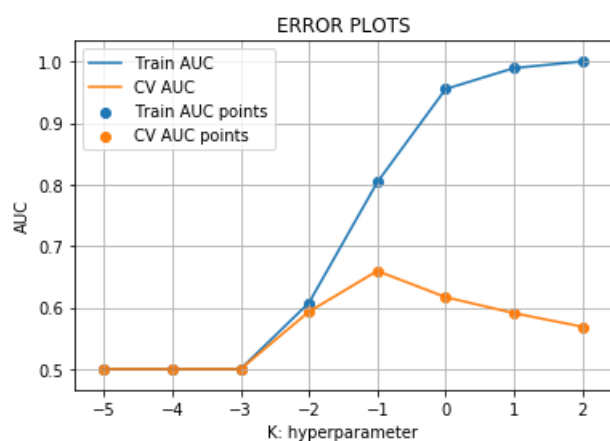
plt.plot([np.log10(i) for i in K], train_auc, label='Train AUC')
plt.plot([np.log10(i) for i in K], cv_auc, label='CV AUC')

plt.scatter([np.log10(i) for i in K], train_auc, label='Train AUC points')
plt.scatter([np.log10(i) for i in K], cv_auc, label='CV AUC points')

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

```

100%|██████████| 8/8 [02:04<00:00, 33.56s/it]



Testing the performance of the model on test data, plotting ROC Curves

In [0]:

```
# from the error plot we choose C such that, we will have maximum AUC on cv data.
```

```
#here we are choosing the best_c based on forloop results
best_c = 0.1
```

In [0]:

```

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

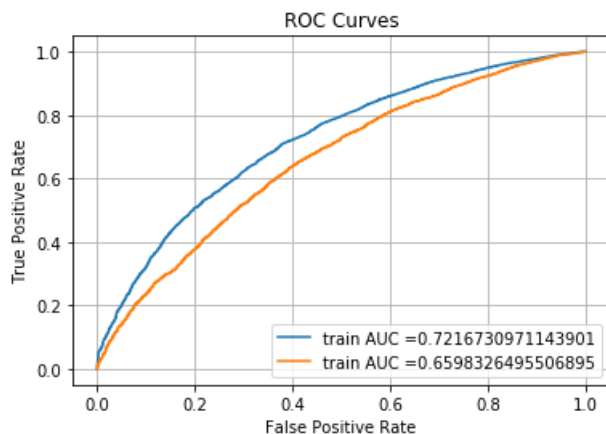
neigh = LogisticRegression(penalty='l1', C=best_c)
neigh.fit(X_tr1, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

y_train_pred = batch_predict(neigh, X_tr1)
y_test_pred = batch_predict(neigh, X_te1)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves")
plt.grid()
plt.show()

```



In [0]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr

def find_best_threshold(threshold, fpr, tpr):
    t = threshold[np.argmax(tpr*(1-fpr))]
    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    return t

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

In [0]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

```
=====

the maximum value of tpr*(1-fpr) 0.44046600921346435 for threshold 0.827
Train confusion matrix
[[ 2158  1305]
 [ 5565 13417]]
Test confusion matrix
[[1385 1161]
 [4340 9614]]
```

In [0]:

```
#Confusion matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt='.1f')
```

```

plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

print('='*100)

print("Test Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt= '.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

```

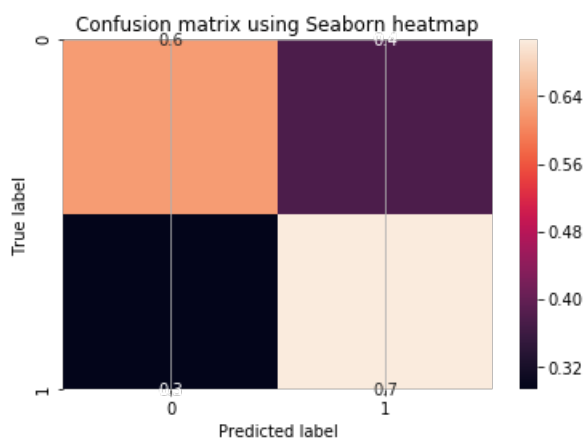
the maximum value of $tpr*(1-fpr)$ 0.44046600921346435 for threshold 0.827

Train Confusion matrix

```

[[ 2158  1305]
 [ 5565 13417]]

```

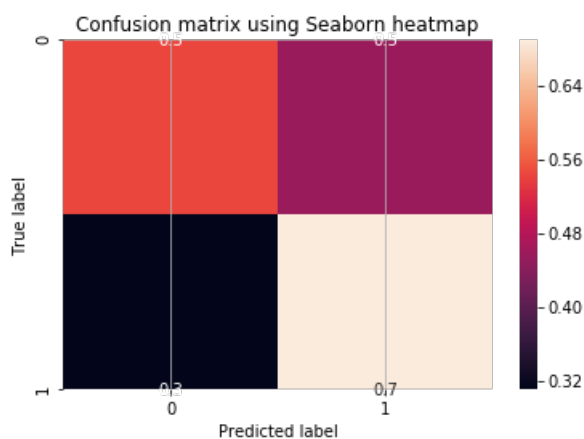


Test Confusion matrix

```

[[1385 1161]
 [4340 9614]]

```



2.4.2 Applying Logistic Regression on TFIDF, SET 2

In [0]:

```
from sklearn.linear_model import LogisticRegression
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
K = [10**-5, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]
for i in tqdm(K):
    neigh = LogisticRegression(penalty='l1', C=i, class_weight='balanced')
    neigh.fit(X_tr2, y_train)

    y_train_pred = batch_predict(neigh, X_tr2)
    y_cv_pred = batch_predict(neigh, X_cv2)

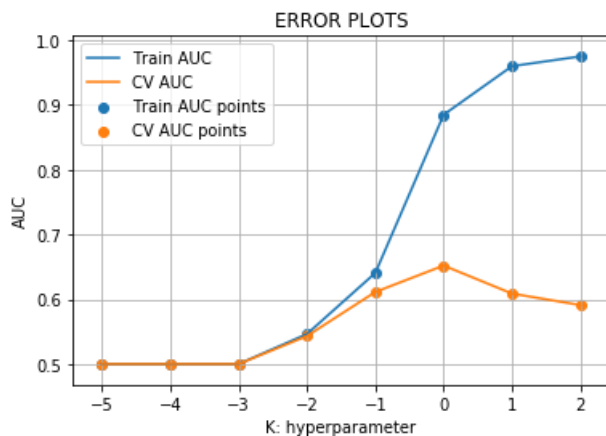
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    tive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot([np.log10(i) for i in K], train_auc, label='Train AUC')
plt.plot([np.log10(i) for i in K], cv_auc, label='CV AUC')

plt.scatter([np.log10(i) for i in K], train_auc, label='Train AUC points')
plt.scatter([np.log10(i) for i in K], cv_auc, label='CV AUC points')

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

100%|██████████| 8/8 [00:43<00:00, 11.86s/it]



Testing the performance of the model on test data, plotting ROC Curves

In [0]:

```
# from the error plot we choose C such that, we will have maximum AUC on cv data.

#here we are choosing the best_c based on forloop results
best_c = 1
```

In [0]:

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

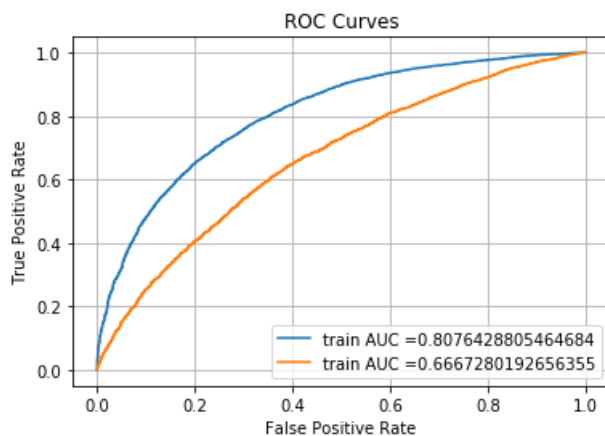
neigh = LogisticRegression(penalty='l1', C=best_c)
neigh.fit(X_tr2, y_train)
```

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

y_train_pred = batch_predict(neigh, X_tr2)
y_test_pred = batch_predict(neigh, X_te2)

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

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves")
plt.grid()
plt.show()
```



In [0]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t)))
```

=====

the maximum value of tpr*(1-fpr) 0.5334064981786018 for threshold 0.831

Train confusion matrix

```
[[ 2571   892]
 [ 5344 13638]]
```

Test confusion matrix

```
[[1429 1117]
 [4433 9521]]
```

In [0]:

```
#Confusion matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt='.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
```

```

plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

print('='*100)

print("Test Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt= '.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

```

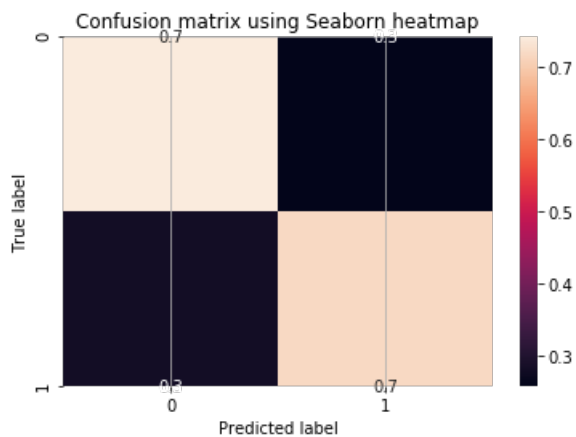
the maximum value of $tpr*(1-fpr)$ 0.5334064981786018 for threshold 0.831

Train Confusion matrix

```

[[ 2571   892]
 [ 5344 13638]]

```

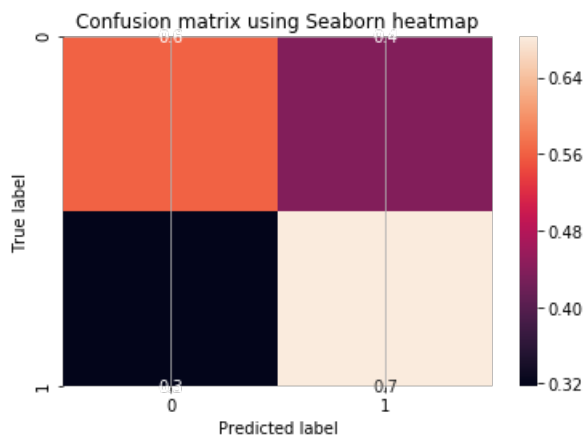


Test Confusion matrix

```

[[1429 1117]
 [4433 9521]]

```



2.4.3 Applying KNN brute force on AVG W2V, SET 3

In [0]:


```

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

train_auc = []
cv_auc = []
K = [10**-5, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1]
for i in tqdm(K):
    neigh = LogisticRegression(penalty='l1', C=i, class_weight='balanced')
    neigh.fit(X_tr3, y_train)

    y_train_pred = batch_predict(neigh, X_tr3)
    y_cv_pred = batch_predict(neigh, X_cv3)

    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot([np.log10(i) for i in K], train_auc, label='Train AUC')
plt.plot([np.log10(i) for i in K], cv_auc, label='CV AUC')

plt.scatter([np.log10(i) for i in K], train_auc, label='Train AUC points')
plt.scatter([np.log10(i) for i in K], cv_auc, label='CV AUC points')

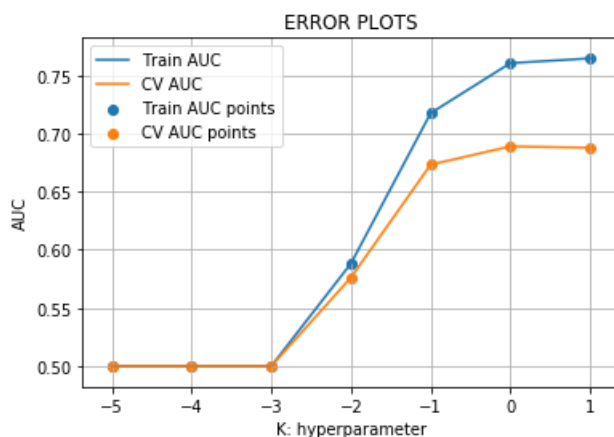
plt.legend()
plt.xlabel("K: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid()
plt.show()

```

```

0%|          | 0/7 [00:00<?, ?it/s]
14%|█        | 1/7 [00:01<00:06, 1.07s/it]
29%|██       | 2/7 [00:02<00:05, 1.07s/it]
43%|███      | 3/7 [00:03<00:04, 1.06s/it]
57%|████     | 4/7 [00:05<00:04, 1.37s/it]
71%|█████    | 5/7 [03:28<02:03, 61.86s/it]
86%|██████   | 6/7 [36:57<10:46, 646.16s/it]
100%|████████| 7/7 [1:31:18<00:00, 1430.47s/it]

```



Testing the performance of the model on test data, plotting ROC Curves

In [0]:

```

# from the error plot we choose K such that, we will have maximum AUC on cv data and gap between t
he train and cv is less
# Note: based on the method you use you might get different hyperparameter values as best one
# so, you choose according to the method you choose, you use gridsearch if you are having more com
puting power and note it will take more time
# if you increase the cv values in the GridSearchCV you will get more robust results.

#here we are choosing the best_k based on forloop results

best_c = 1

```

In [0]:

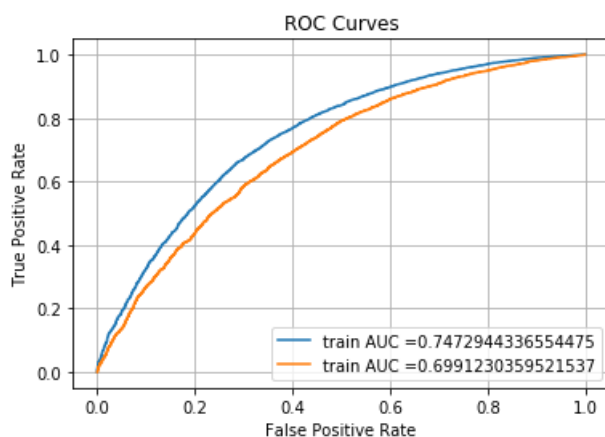
```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc
from sklearn.linear_model import LogisticRegression
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

neigh = LogisticRegression(penalty='l1', C=best_c)
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

y_train_pred2 = batch_predict(neigh, X_tr3)
y_test_pred2 = batch_predict(neigh, X_te3)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred2)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred2)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves")
plt.grid()
plt.show()
```



In [0]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred2, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred2, best_t)))
```

```
=====

the maximum value of tpr*(1-fpr) 0.47241531583959057 for threshold 0.827
Train confusion matrix
[[ 2257  1206]
 [ 5223 13759]]
Test confusion matrix
[[1478 1068]
 [4014 9940]]
```

In [0]:

```
#Confusion matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
```

```

print("Train Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred2, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred2, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt='.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

print('='*100)

print("Test Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred2, best_t))
print(conf_mat)

#For heatmap

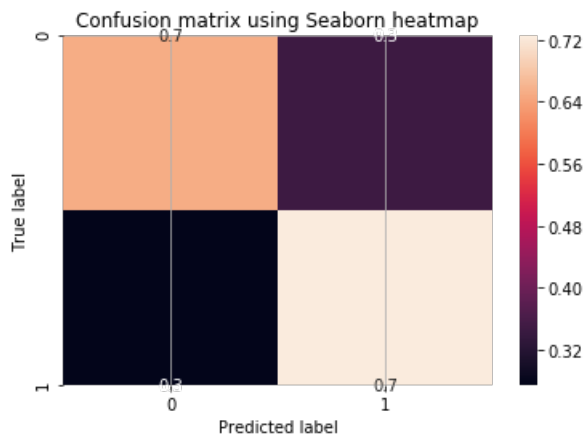
import seaborn as sns
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred2, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt= '.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

```

the maximum value of $tpr*(1-fpr)$ 0.47241531583959057 for threshold 0.827

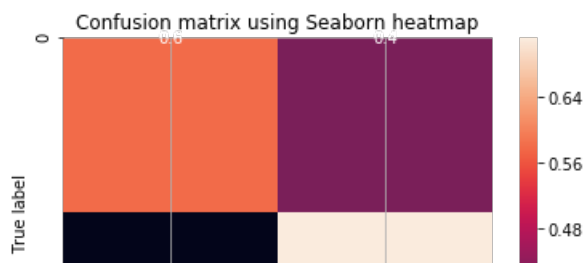
Train Confusion matrix

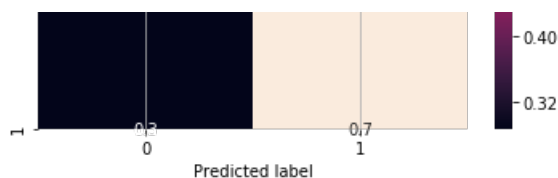
```
[[ 2257  1206]
 [ 5223 13759]]
```



Test Confusion matrix

```
[[1478 1068]
 [4014 9940]]
```





2.4.4 Applying KNN brute force on TFIDF W2V, SET 4

In [0]:

```
from sklearn.linear_model import LogisticRegression
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
K = [10**-5, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1]
for i in tqdm(K):
    neigh = LogisticRegression(penalty='l1', C=i, class_weight='balanced')
    neigh.fit(X_tr4, y_train)

    y_train_pred = batch_predict(neigh, X_tr4)
    y_cv_pred = batch_predict(neigh, X_cv4)

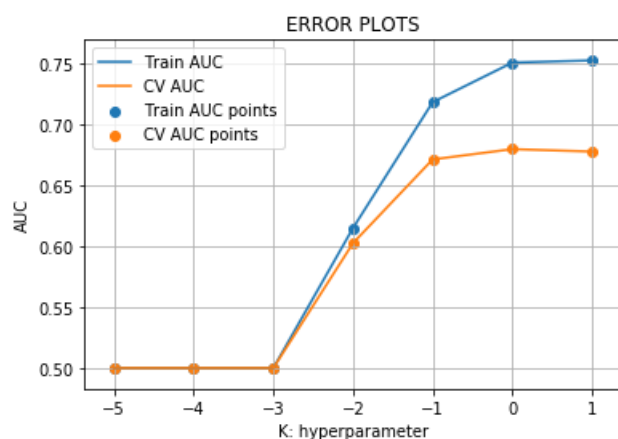
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the posi
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot([np.log10(i) for i in K], train_auc, label='Train AUC')
plt.plot([np.log10(i) for i in K], cv_auc, label='CV AUC')

plt.scatter([np.log10(i) for i in K], train_auc, label='Train AUC points')
plt.scatter([np.log10(i) for i in K], cv_auc, label='CV AUC points')

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

100% |██████████| 7/7 [41:59<00:00, 700.20s/it]



Testing the performance of the model on test data, plotting ROC Curves

In [0]:

```
# from the error plot we choose C such that, we will have maximum AUC on cv data

#here we are choosing the best_k based on forloop results

best_c = 1
```

```
best_c = 1
```

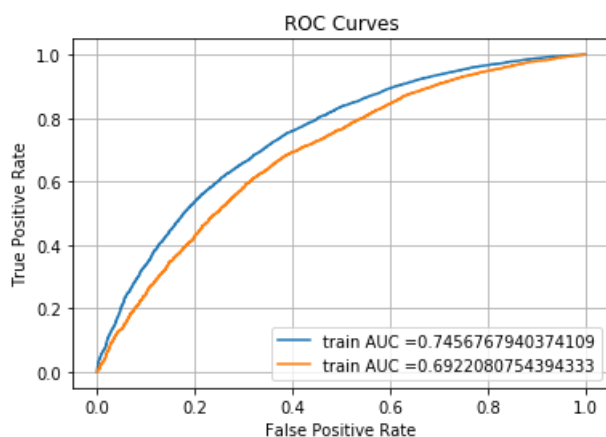
In [0]:

```
from sklearn.metrics import roc_curve, auc
neigh = LogisticRegression(penalty='l1', C=best_c)
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

y_train_pred3 = batch_predict(neigh, X_tr4)
y_test_pred3 = batch_predict(neigh, X_te4)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred3)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred3)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves")
plt.grid()
plt.show()
```



In [0]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred3, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred3, best_t)))
```

```
=====

the maximum value of tpr*(1-fpr) 0.4646923436105996 for threshold 0.835
Train confusion matrix
[[ 2349  1114]
 [ 5978 13004]]
Test confusion matrix
[[1576  970]
 [4489 9465]]
```

In [0]:

```
#Confusion matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred3, best_t))
print(conf_mat)

#For heatmap
```

```

import seaborn as sns
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred3, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt='.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

print('='*100)

print("Test Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred3, best_t))
print(conf_mat)

#For heatmap

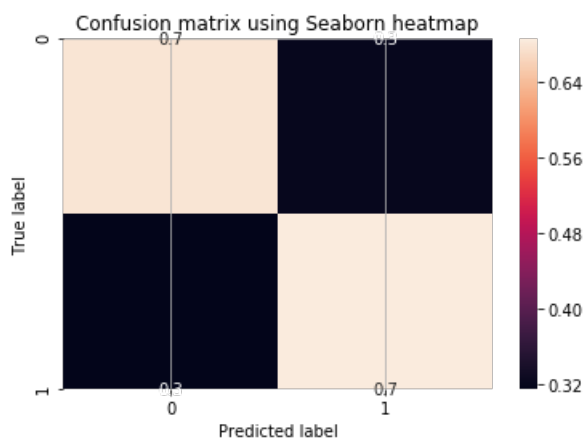
import seaborn as sns
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred3, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt= '.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4646923436105996 for threshold 0.835

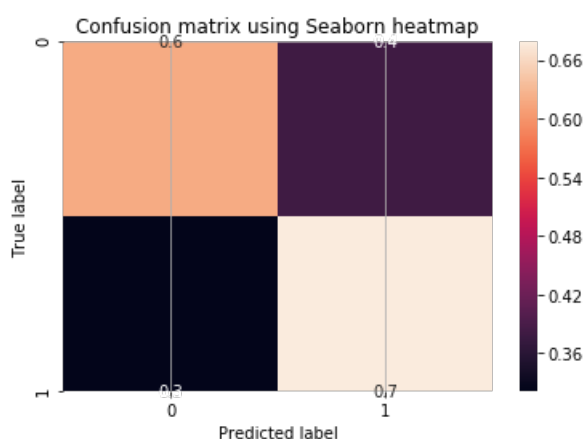
Train Confusion matrix

```
[[ 2349  1114]
 [ 5978 13004]]
```



Test Confusion matrix

```
[[1576  970]
 [4489 9465]]
```



2.5 Logistic Regression with added Features `Set 5`

Set 5:

In [0]:

```
from sklearn.linear_model import LogisticRegression
import matplotlib.pyplot as plt
from sklearn.metrics import roc_auc_score

train_auc = []
cv_auc = []
K = [10**-5, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1, 10**2]
for i in tqdm(K):
    neigh = LogisticRegression(penalty='l1', C=i, class_weight='balanced')
    neigh.fit(X_tr5, y_train)

    y_train_pred = batch_predict(neigh, X_tr5)
    y_cv_pred = batch_predict(neigh, X_cv5)

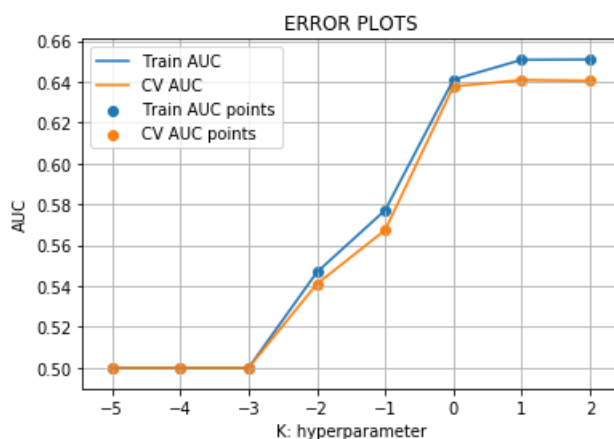
    # roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
    # not the predicted outputs
    train_auc.append(roc_auc_score(y_train, y_train_pred))
    cv_auc.append(roc_auc_score(y_cv, y_cv_pred))

plt.plot([np.log10(i) for i in K], train_auc, label='Train AUC')
plt.plot([np.log10(i) for i in K], cv_auc, label='CV AUC')

plt.scatter([np.log10(i) for i in K], train_auc, label='Train AUC points')
plt.scatter([np.log10(i) for i in K], cv_auc, label='CV AUC points')

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

100% |██████████| 8/8 [00:12<00:00, 2.83s/it]



Testing the performance of the model on test data, plotting ROC Curves

In [0]:

```
# from the error plot we choose C such that, we will have maximum AUC on cv data

#here we are choosing the best_c based on forloop results

best_c = 1
```

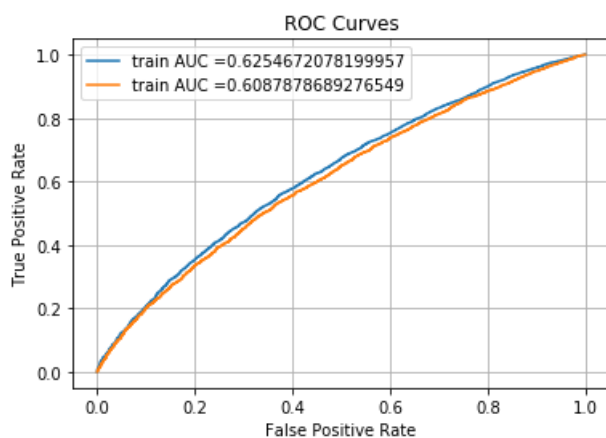
In [0]:

```
neigh = LogisticRegression(penalty='l1',C=best_c)
neigh.fit(X_tr5, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs

y_train_pred4 = batch_predict(neigh, X_tr5)
y_test_pred4 = batch_predict(neigh, X_te5)

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred4)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred4)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curves")
plt.grid()
plt.show()
```



In [0]:

```
print("="*100)
from sklearn.metrics import confusion_matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train confusion matrix")
print(confusion_matrix(y_train, predict_with_best_t(y_train_pred4, best_t)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict_with_best_t(y_test_pred4, best_t)))
```

=====

the maximum value of $tpr \cdot (1 - fpr)$ 0.34964199863737044 for threshold 0.847

Train confusion matrix

```
[[ 2148  1315]
 [ 8282 10700]]
```

Test confusion matrix

```
[[1482 1064]
 [5958 7996]]
```

In [0]:

```
#Confusion matrix
best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
print("Train Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred4, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred4, best_t))
```



```

conf_mat=confusion_matrix(y_train, predict_with_best_t(y_train_pred4, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt='.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

print('='*100)

print("Test Confusion matrix")
from sklearn.metrics import confusion_matrix
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred4, best_t))
print(conf_mat)

#For heatmap

import seaborn as sns
conf_mat=confusion_matrix(y_test, predict_with_best_t(y_test_pred4, best_t))
conf_mat_normalized=conf_mat.astype('float')/conf_mat.sum(axis=1)[:,np.newaxis]
sns.heatmap(conf_mat_normalized,annot=True, fmt= '.1f')
plt.ylabel('True label')
plt.xlabel('Predicted label')
plt.title("Confusion matrix using Seaborn heatmap")
plt.grid()
plt.show()

```

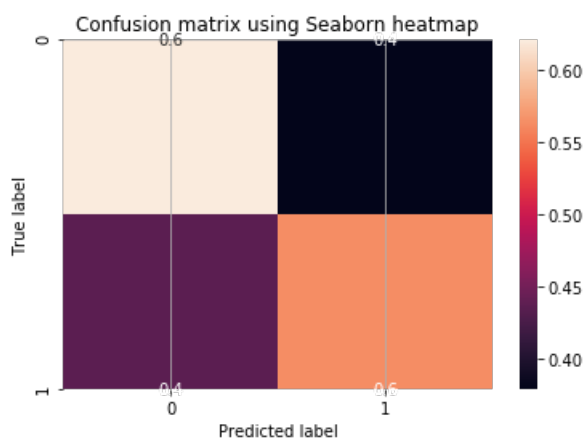
the maximum value of $tpr \cdot (1 - fpr)$ 0.34964199863737044 for threshold 0.847

Train Confusion matrix

```

[[ 2148  1315]
 [ 8282 10700]]

```

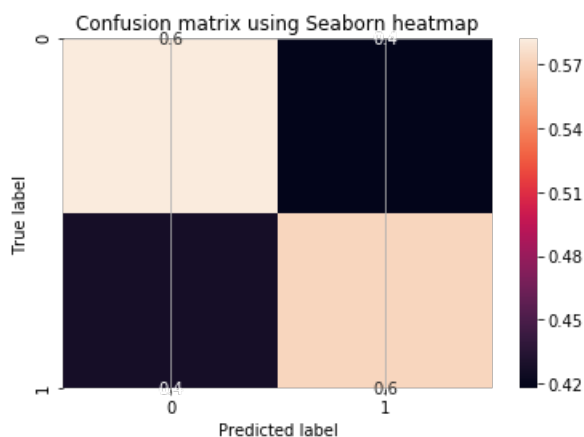


Test Confusion matrix

```

[[1482 1064]
 [5958 7996]]

```



3. Conclusion

In [0]:

```
#http://zetcode.com/python/prettytable/
```

```
from prettytable import PrettyTable
```

```
x = PrettyTable()
```

```
x.field_names = ["Vectorizer", "Regularizer", "Hyper Parameter", "AUC"]
```

```
x.add_row(["BOW", "L1", 0.1, 0.827])
```

```
x.add_row(["TFIDF", "L1", 1, 0.831])
```

```
x.add_row(["W2V", "L1", 1, 0.827])
```

```
x.add_row(["TFIDFW2V", "L1", 1, 0.835])
```

```
x.add_row(["Custom", "L1", 1, 0.847])
```

```
print(x)
```

Vectorizer	Regularizer	Hyper Parameter	AUC
BOW	L1	0.1	0.827
TFIDF	L1	1	0.831
W2V	L1	1	0.827
TFIDFW2V	L1	1	0.835
Custom	L1	1	0.847