

9_DonorsChoose_RF_GBDT

March 21, 2020

1 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

How to scale current manual processes and resources to screen 500,000 projects so that they can
How to increase the consistency of project vetting across different volunteers to improve t
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.

1.1 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

`project_title` | Title of the project. **Examples:**

Art Will Make You Happy!

First Grade Fun

`project_grade_category` | Grade level of students for which the project is targeted. One of the following enumerated values:

Grades PreK-2

Grades 3-5

Grades 6-8

Grades 9-12

`project_subject_categories` | One or more (comma-separated) subject categories for the project from the following enumerated list of values:

Applied Learning

Care & Hunger
 Health & Sports
 History & Civics
 Literacy & Language
 Math & Science
 Music & The Arts
 Special Needs
 Warmth

Examples:

Music & The Arts
 Literacy & Language, Math & Science

school_state | State where school is located ([Two-letter U.S. postal code](#)). **Example:** WY
project_subject_subcategories | One or more (comma-separated) subject subcategories for the project. **Examples:**

Literacy
 Literature & Writing, Social Sciences

project_resource_summary | An explanation of the resources needed for the project. **Example:**

My students need hands on literacy materials to manage sensory needs!

project_essay_1 | First application essay

project_essay_2 | *Second application essay* **project_essay_3** | Third application essay

project_essay_4 | *Fourth application essay* **project_submitted_datetime** | Datetime when project application was submitted. **Example:** 2016-04-28 12:43:56.245

teacher_id | A unique identifier for the teacher of the proposed project. **Example:** bdf8baa8fedef6bfeec7ae4ff1c15c56

teacher_prefix | Teacher's title. One of the following enumerated values:

nan
 Dr.
 Mr.
 Mrs.
 Ms.
 Teacher.

teacher_number_of_previously_posted_projects | Number of project applications previously submitted by the same teacher. **Example:** 2

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

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

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

Feature	Description
quantity	Quantity of the resource required.
price	Price of the resource required. Example: 9.95

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

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

Label	Description
project_is_approved	Advisory 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.

1.1.1 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]: %matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
```

```

import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

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

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

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

from tqdm import tqdm
import os

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

```

1.2 1.1 Reading Data

```

In [0]: from google.colab import drive
        drive.mount('/content/drive')

```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-

Enter your authorization code:

ûûûûûûûûûûûû

Mounted at /content/drive

```

In [0]: project_data = pd.read_csv('/content/drive/My Drive/Assignments_DonorsChoose_2018/train
        resource_data = pd.read_csv('/content/drive/My Drive/Assignments_DonorsChoose_2018/res

```

```
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

1.3 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/4117234

        # 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
                    j=j.replace('The','') # if we have the words "The" we are going to replace
                j = j.replace(' ','') # we are replacing all the ' '(space) with ''(empty) ex: "Math & Science"
                temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing space
            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]))

```

1.4 1.3 preprocessing of project_subject_subcategories

```

In [0]: sub_categories = list(project_data['project_subject_subcategories'].values)
        # remove special characters from list of strings python: https://stackoverflow.com/a/4

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

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 "The"
            j=j.replace('The','') # if we have the words "The" we are going to replace
        j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex: "Math & Science"
        temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing space
        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/22898595/408403
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]))

```

Preprocessing of project_grade_category

```

In [0]: project_grade = list(project_data['project_grade_category'].values)
        # remove special characters from list of strings python: https://stackoverflow.com/a/4

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

```

```

grade_cat_list = []
for i in project_grade:
    # consider we have text like this:
    for j in i.split(' '): # # split by spae
        j=j.replace('Grades','')# clean grades from the row
    grade_cat_list.append(j.strip())

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

```

1.5 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  ...                                     essay
0      160221  ...  My students are English learners that are work...
1      140945  ...  Our students arrive to our school eager to lea...

[2 rows x 18 columns]

```

```

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)

```

My students are English learners that are working on English as their second or third languages

=====

The 51 fifth grade students that will cycle through my classroom this year all love learning, a

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, row

=====

My kindergarten students have varied disabilities ranging from speech and language delays, cog

=====

```
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, cogn
=====

```
In [0]: # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-p
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, cogn

```
In [0]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

My kindergarten students have varied disabilities ranging from speech and language delays cogn

```
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'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 't',
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "tl
```



```
'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'h
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as
'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through
'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'o
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any
'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too
's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'n
've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't"
'hadn't', 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'migh
'mustn't', 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn'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 [01:03<00:00, 1730.28it/s]

1.4 Preprocessing of project_title

```
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['project_title'].values[2000])
        print(sent)
        print("="*50)
```

Steady Stools for Active Learning

=====

```
In [0]: # \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-p
        sent = sent.replace('\r', ' ')
        sent = sent.replace('\n', ' ')
        sent = sent.replace('\t', ' ')
        print(sent)
```

Steady Stools for Active Learning

```
In [0]: #remove spacial character: https://stackoverflow.com/a/5843547/4084039
        sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
        print(sent)
```

Steady Stools for Active Learning

```
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'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him',
                    'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'theirs',
                    'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that's", 'they',
                    'them', 'their', 'theirs', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had',
                    'do', 'does', 'isn't', 'are', 'aren't', 'wasn't', 'were', 'weren't', 'couldn't', 'couldn't',
                    'couldn't', 'couldn't', 'didn't', "didn't", "hadn't", 'hasn't', "hasn't", 'haven't',
                    "haven't", 'isn't', "isn't", 'ma', 'mightn't', "mightn't", 'mustn't', 'needn't',
                    "needn't", 'shan't', "shan't", 'shouldn't', "shouldn't", 'won't', "won't", 'wouldn't',
                    "wouldn't"]
```

```
In [0]: # Combining all the above stundents
        from tqdm import tqdm
        preprocessed_titles = []
        # tqdm is for printing the status bar
        for sentence in tqdm(project_data['project_title'].values):
            sent = decontracted(sentence)
            sent = sent.replace('\r', ' ')
```

```

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

```

100%| 109248/109248 [00:02<00:00, 41614.09it/s]

In [0]: resource_data.head(2)

```

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

```

In [0]: price_data = resource_data.groupby('id').agg({'price': 'sum', 'quantity': 'sum'}).reset_index()
price_data.head(2)

```

Out[0]:
   id  price  quantity
0  p000001  459.56      7
1  p000002  515.89     21

```

In [0]: project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('null')

project_data.head(2)

```

Out[0]:
   Unnamed: 0  ...  essay
0      160221  ...  My students are English learners that are work...
1      140945  ...  Our students arrive to our school eager to lea...

```

[2 rows x 18 columns]

In [0]: *#Join train & Resource dataset*
join two dataframes in python:

data = pd.merge(project_data, price_data, on='id', how='left')

In [0]: approved_price = data[data['project_is_approved']==1]['price'].values

rejected_price = data[data['project_is_approved']==0]['price'].values

In [0]: data.head(2)

```

Out[0]:
   Unnamed: 0  id  ...  price  quantity
0      160221  p253737  ...  154.6      23
1      140945  p258326  ...  299.0      1

```

[2 rows x 20 columns]

Train Test split

```
In [0]: project_data = data.sample(n=50000)
        project_data.head(3)
```

```
Out [0]:      Unnamed: 0      id  ...  price quantity
100368      84919  p015166  ...  580.96          4
90247      72722  p162752  ...  422.83         33
22663      161971  p049465  ...  199.00          1

[3 rows x 20 columns]
```

```
In [0]: y = project_data['project_is_approved'].values
        X = project_data
        X.head(1)
```

```
Out [0]:      Unnamed: 0      id  ...  price quantity
100368      84919  p015166  ...  580.96          4

[1 rows x 20 columns]
```

```
In [0]: y=project_data['project_is_approved'].values
```

```
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=0.2, stratify=y)
```

```
In [0]: # Combining all the above students
        from tqdm import tqdm
        train_preprocessed_essays = []
        # tqdm is for printing the status bar
        for sentence in tqdm(X_train['essay'].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.lower() not in stopwords)
            train_preprocessed_essays.append(sent.lower().strip())
```

```
100%| 40000/40000 [00:22<00:00, 1753.73it/s]
```

```
In [0]: # Combining all the above students
        from tqdm import tqdm
        test_preprocessed_essays = []
        # tqdm is for printing the status bar
        for sentence in tqdm(X_test['essay'].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.lower() not in stopwords)
test_preprocessed_essays.append(sent.lower().strip())

```

100%| 10000/10000 [00:05<00:00, 1742.62it/s]

```

In [0]: # after preprocessing
preprocessed_essays[20000]

```

Out[0]: 'my kindergarten students varied disabilities ranging speech language delays cognitive

```

In [0]: # Combining all the above students
from tqdm import tqdm
train_preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(X_train['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.lower() not in stopwords)
    train_preprocessed_titles.append(sent.lower().strip())

```

100%| 40000/40000 [00:00<00:00, 41573.60it/s]

```

In [0]: # Combining all the above students
from tqdm import tqdm
test_preprocessed_titles = []
# tqdm is for printing the status bar
for sentence in tqdm(X_test['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.lower() not in stopwords)
    test_preprocessed_titles.append(sent.lower().strip())

```

100%| 10000/10000 [00:00<00:00, 40372.28it/s]

```
In [0]: preprocessed_titles[2000]
```

```
Out[0]: 'steady stools active learning'
```

1.6 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_title', 'project_essay_1',  
              'project_essay_2', 'project_essay_3', 'project_essay_4',  
              'project_resource_summary',  
              'teacher_number_of_previously_posted_projects', 'project_is_approved',  
              'clean_categories', 'clean_subcategories', 'grade_cat_list', 'essay',  
              'price', 'quantity'],  
             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

2 Response coding for Categorical Data

```
In [0]: #http://www.saedsayad.com/encoding.htm  
        #https://gist.github.com/marnixkoops/e68815d30474786e2b293682ed7cdb01  
        #https://www.slideshare.net/0xdata/feature-engineering-83511751
```

3 Response Coding - School state

```
In [0]: # code for response coding with Laplace smoothing.  
        # alpha : used for laplace smoothing
```

```
def get_cat_fea_dict(alpha, feature, df):  
    value_count = X_train[feature].value_counts()  
  
    cat_dict = dict()
```

```

        # denominator will contain the number of time that particular feature occurred in w
        for i, denominator in value_count.items():
            vec = []
            for k in range(0,2):

                cls_cnt = X_train.loc[(X_train['project_is_approved']==k) & (X_train[feature

                # cls_cnt.shape[0](numerator) will contain the number of time that particu
                vec.append((cls_cnt.shape[0] + alpha*1)/ (denominator + 2*alpha))

            cat_dict[i]=vec
        return cat_dict

def get_cat_feature(alpha, feature, df):
    cat_dict = get_cat_fea_dict(alpha, feature, df)

    value_count = X_train[feature].value_counts()

    cat_fea = []
    # for every feature values in the given data frame we will check if it is there
    # if not we will add [1/2,1/2] to cat_fea
    for index, row in df.iterrows():
        if row[feature] in dict(value_count).keys():
            cat_fea.append(cat_dict[row[feature]])
        else:
            cat_fea.append([1/2,1/2])

    return cat_fea

```

```

In [0]: #response-coding of the clean_categories
        # alpha is used for laplace smoothing
        alpha = 1

        # train gene feature
        tr_clean_categories = np.array(get_cat_feature(alpha, "clean_categories", X_train))

        # test gene feature
        te_clean_categories = np.array(get_cat_feature(alpha, "clean_categories", X_test))

In [0]: print("The shape of train_clean_categories feature:", tr_clean_categories.shape)
        print("The shape of test_clean_categories feature:", te_clean_categories.shape)

The shape of train_clean_categories feature: (40000, 2)
The shape of test_clean_categories feature: (10000, 2)

```

```

In [0]: #response-coding of the clean_subcategories
        # alpha is used for laplace smoothing

```

```

alpha = 1

# train gene feature
tr_clean_subcategories = np.array(get_cat_feature(alpha, "clean_subcategories", X_train))

# test gene feature
te_clean_subcategories = np.array(get_cat_feature(alpha, "clean_subcategories", X_test))

print("The shape of tr_clean_subcategories feature:", tr_clean_subcategories.shape)
print("The shape of te_clean_subcategories feature:", te_clean_subcategories.shape)

```

The shape of tr_clean_subcategories feature: (40000, 2)
The shape of te_clean_subcategories feature: (10000, 2)

```

In [0]: project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('null')

#response-coding of the clean_subcategories
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
tr_prefix = np.array(get_cat_feature(alpha, "teacher_prefix", X_train))

# test gene feature
te_prefix = np.array(get_cat_feature(alpha, "teacher_prefix", X_test))

print("The shape of train_teacher_prefix feature:", tr_prefix.shape)
print("The shape of test_teacher_prefix feature:", te_prefix.shape)

```

The shape of train_teacher_prefix feature: (40000, 2)
The shape of test_teacher_prefix feature: (10000, 2)

```

In [0]: #response-coding of the clean_subcategories
# alpha is used for laplace smoothing
alpha = 1

# train gene feature
tr_school_state = np.array(get_cat_feature(alpha, "school_state", X_train))

# test gene feature
te_school_state = np.array(get_cat_feature(alpha, "school_state", X_test))

print("The shape of train_school_state feature:", tr_school_state.shape)
print("The shape of test_school_state feature:", te_school_state.shape)

```


The shape of train_school_state feature: (40000, 2)
The shape of test_school_state feature: (10000, 2)

```
In [0]: #response-coding of the clean_subcategories
        # alpha is used for laplace smoothing
        alpha = 1

        # train gene feature
        tr_project_grade = np.array(get_cat_feature(alpha, "grade_cat_list", X_train))

        # test gene feature
        te_project_grade = np.array(get_cat_feature(alpha, "grade_cat_list", X_test))

        print("The shape of tr_clean_project_grade_category feature:", tr_project_grade.shape)
        print("The shape of te_clean_project_grade_category feature:", te_project_grade.shape)
```

The shape of tr_clean_project_grade_category feature: (40000, 2)
The shape of te_clean_project_grade_category feature: (10000, 2)

3.0.1 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 pr
        vectorizer = CountVectorizer(min_df=10)
        text_bow = vectorizer.fit_transform(preprocessed_essays)
        print("Shape of matrix after one hot encodig ",text_bow.shape)
```

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

```
In [0]: # We are considering only the words which appeared in at least 10 documents(rows or pr
        vectorizer = CountVectorizer(min_df=10)
        text_bow = vectorizer.fit_transform(preprocessed_titles)
        print("Shape of matrix after one hot encodig ",text_bow.shape)
```

Shape of matrix after one hot encodig (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 encodig ",text_tfidf.shape)
```

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

```
In [0]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=10)
text_tfidf = vectorizer.fit_transform(preprocessed_titles)
print("Shape of matrix after one hot encodig ",text_tfidf.shape)
```

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

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 preproced_texts:
    words.extend(i.split(' '))

for i in preproced_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 corpus", \
      len(inter_words), "(", np.round(len(inter_words)/len(words)*100,3), "%)")

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

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

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

'''

```

Out[0]: '\n# Reading glove vectors in python: <https://stackoverflow.com/a/38230349/4084039>\nde

```

In [0]: # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-
# make sure you have the glove_vectors file
with open('/content/drive/My Drive/Assignments_DonorsChoose_2018/glove_vectors', 'rb')
    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 [00:34<00:00, 3187.45it/s]

109248
300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
train_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(train_preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    train_avg_w2v_vectors.append(vector)

print(len(train_avg_w2v_vectors))
print(len(train_avg_w2v_vectors[0]))
```

100%|| 40000/40000 [00:11<00:00, 3517.78it/s]

40000
300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
test_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(test_preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    test_avg_w2v_vectors.append(vector)

print(len(test_avg_w2v_vectors))
print(len(test_avg_w2v_vectors[0]))
```

100%|| 10000/10000 [00:02<00:00, 3482.86it/s]

10000
300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles): # 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_vectors1.append(vector)

print(len(avg_w2v_vectors1))
print(len(avg_w2v_vectors1[0]))
```

100%|| 109248/109248 [00:01<00:00, 64013.01it/s]

109248
300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
train_avg_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(train_preprocessed_titles): # 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
    train_avg_w2v_vectors1.append(vector)

print(len(train_avg_w2v_vectors1))
print(len(train_avg_w2v_vectors1[0]))
```

100%|| 40000/40000 [00:00<00:00, 58548.53it/s]

40000

300

```
In [0]: # average Word2Vec
        # compute average word2vec for each review.
        test_avg_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this l
        for sentence in tqdm(test_preprocessed_titles): # 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
            test_avg_w2v_vectors1.append(vector)

        print(len(test_avg_w2v_vectors1))
        print(len(test_avg_w2v_vectors1[0]))
```

100%|| 10000/10000 [00:00<00:00, 59421.44it/s]

10000

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((se
        tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
        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 [03:35<00:00, 506.76it/s]

109248

300

```

In [0]: # average Word2Vec
# compute average word2vec for each review.
train_tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(train_preprocessed_essays): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((se
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    train_tfidf_w2v_vectors.append(vector)

print(len(train_tfidf_w2v_vectors))
print(len(train_tfidf_w2v_vectors[0]))

```

100%|| 40000/40000 [01:10<00:00, 564.69it/s]

40000

300

```

In [0]: # average Word2Vec
# compute average word2vec for each review.
test_tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(test_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((se
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    test_tfidf_w2v_vectors.append(vector)

print(len(test_tfidf_w2v_vectors))
print(len(test_tfidf_w2v_vectors[0]))

100%|| 10000/10000 [00:17<00:00, 565.64it/s]

10000
300

```

```

In [0]: # average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(preprocessed_titles): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((se
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors1.append(vector)

print(len(tfidf_w2v_vectors1))
print(len(tfidf_w2v_vectors1[0]))

```


100%|| 109248/109248 [00:03<00:00, 31200.33it/s]

109248

300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
train_tfidf_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(train_preprocessed_titles): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))) # g
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    train_tfidf_w2v_vectors1.append(vector)

print(len(train_tfidf_w2v_vectors1))
print(len(train_tfidf_w2v_vectors1[0]))
```

100%|| 40000/40000 [00:01<00:00, 31964.05it/s]

40000

300

```
In [0]: # average Word2Vec
# compute average word2vec for each review.
test_tfidf_w2v_vectors1 = []; # the avg-w2v for each sentence/review is stored in this
for sentence in tqdm(test_preprocessed_titles): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split()))) # g
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # g
```

```

        vector += (vec * tf_idf) # calculating tfidf weighted w2v
        tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    test_tfidf_w2v_vectors1.append(vector)

print(len(test_tfidf_w2v_vectors1))
print(len(test_tfidf_w2v_vectors1[0]))

100%|| 10000/10000 [00:00<00:00, 29814.97it/s]

10000
300

```

3.0.2 1.5.3 Vectorizing Numerical features

```

In [0]: # check this one: https://www.youtube.com/watch?v=0HQq0cln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.p
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. ...
# 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 st
print(f"Mean : {price_scaler.mean_[0]}, Standard deviation : {np.sqrt(price_scaler.var_

# Now standardize the data with above maen and variance.
price_standardized = price_scaler.transform(project_data['price'].values.reshape(-1, 1)
tr_price_standardized = price_scaler.transform(X_train['price'].values.reshape(-1, 1))
#cv_price_standardized = price_scaler.transform(X_cv['price'].values.reshape(-1, 1))
te_price_standardized = price_scaler.transform(X_test['price'].values.reshape(-1, 1))

Mean : 298.96676780000007, Standard deviation : 368.6852344377611

In [0]: price_standardized

Out[0]: array([[ 0.76486175],
               [ 0.3359593 ],
               [-0.27114394],
               ...,
               [-0.40407576],
               [-0.51368146],
               [-0.78960788]])

```

```
In [0]: # check this one: https://www.youtube.com/watch?v=0H0q0cln3Z4&t=530s
# standardization sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.p
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. ...]
# Reshape your data either using array.reshape(-1, 1)

quantity_scalar = StandardScaler()
quantity_scalar.fit(project_data['quantity'].values.reshape(-1,1)) # finding the mean
print(f"Mean : {quantity_scalar.mean_[0]}, Standard deviation : {np.sqrt(quantity_scalar

# Now standardize the data with above mean and variance.
quantity_standardized = quantity_scalar.transform(project_data['quantity'].values.reshape(-1,1))
tr_quantity_standardized = quantity_scalar.transform(X_train['quantity'].values.reshape(-1,1))
#cv_quantity_standardized = quantity_scalar.transform(X_cv['quantity'].values.reshape(-1,1))
te_quantity_standardized = quantity_scalar.transform(X_test['quantity'].values.reshape(-1,1))

Mean : 16.94668, Standard deviation : 26.19335558834721
```

```
In [0]: quantity_standardized
```

```
Out[0]: array([[ -0.49427344],
               [  0.61287757],
               [-0.6088063 ],
               ...,
               [-0.4179182 ],
               [-0.22703009],
               [  1.26189712]])
```

__ Computing Sentiment Scores __

```
In [0]: import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

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

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest
for learning my students learn in many different ways using all of our senses and mult
of techniques to help all my students succeed students in my class come from a variety
for wonderful sharing of experiences and cultures including native americans our school
learners which can be seen through collaborative student project based learning in and
in my class love to work with hands on materials and have many different opportunities
mastered having the social skills to work cooperatively with friends is a crucial aspe
```

montana is the perfect place to learn about agriculture and nutrition my students love in the early childhood classroom i have had several kids ask me can we try cooking with and create common core cooking lessons where we learn important math and writing concepts food for snack time my students will have a grounded appreciation for the work that went into of where the ingredients came from as well as how it is healthy for their bodies this project nutrition and agricultural cooking recipes by having us peel our own apples to make honey and mix up healthy plants from our classroom garden in the spring we will also create a shared with families students will gain math and literature skills as well as a life lesson nannan'

```
ss = sid.polarity_scores(for_sentiment)
```

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

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

D:\installed\Anaconda3\lib\site-packages\nlTK\twitter_init_.py:20: UserWarning:

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

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

4 Assignment 9: RF and GBDT

Response Coding: Example

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

Apply both Random Forrest and GBDT on these feature sets

Set 1: categorical (instead of one hot encoding, try one hot encoding

Set 2: categorical (instead of one hot encoding, try one hot encoding

Set 3: categorical (instead of one hot encoding, try one hot encoding

Set 4: categorical (instead of one hot encoding, try one hot encoding

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

Consider the following range for hyperparameters n_estimators = [10, 50, 100]

max_depth = [2, 3, 4, 5, 6, 7, 8, 9, 10]

Find the best hyper parameter which will give the maximum cross validation score

Find the best hyper parameter using simple cross validation data

You can write your own for loops to do this task

```

    </ul>
</li>
<br>
<li>
<strong>Representation of results</strong>
    <ul>
<li>You need to plot the performance of model both on train data and cross validation data for
<img src='3d_plot.JPG' width=500px> with X-axis as <strong>n_estimators</strong>, Y-axis as <strong>or</strong></li>
    <p style="text-align:center;font-size:30px;color:red;"><strong>or</strong></p> <br>
<li>You need to plot the performance of model both on train data and cross validation data for
<img src='heat_map.JPG' width=300px> <a href='https://seaborn.pydata.org/generated/seaborn.heat'>
<li>You can choose either of the plotting techniques: 3d plot or heat map</li>
<li>Once after you found the best hyper parameter, you need to train your model with it, and f
<img src='train_test_auc.JPG' width=300px></li>
<li>Along with plotting ROC curve, you need to print the <a href='https://www.appliedaicourse.'>
<img src='confusion_matrix.png' width=300px></li>
    </ul>
<br>
<li><strong>Conclusion</strong>
    <ul>
<li>You need to summarize the results at the end of the notebook, summarize it in the table for
    <img src='summary.JPG' width=400px>
</li>
    </ul>

```

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-leakag, make sure to split your data first and then vectorize it.
 3. While vectorizing your data, apply the method fit_transform() on you train data, and apply the method transform() on cv/test data.
 4. For more details please go through this link.
2. Random Forest and GBDT

2.2 Make Data Model Ready: encoding numerical, categorical features

5 Set1 - BoW

```

In [0]: print(X_train.shape, y_train.shape)
        print(X_test.shape, y_test.shape)

        print("="*100)

vectorizer = CountVectorizer(min_df=10, ngram_range=(1,4), max_features=10000)
vectorizer.fit(train_preprocessed_essays) # 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(train_preprocessed_essays)
X_test_essay_bow = vectorizer.transform(test_preprocessed_essays)

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

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)
# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")

```

```

(40000, 20) (40000,)
(10000, 20) (10000,)

```

```

=====
After vectorizations

```

```

(40000, 10000) (40000,)
(10000, 10000) (10000,)

```

```

=====
NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME

```

```

In [0]: print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)

```

```

print("="*100)

```

```

vectorizer = CountVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer.fit(train_preprocessed_titles) # 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(train_preprocessed_titles)
X_test_title_bow = vectorizer.transform(test_preprocessed_titles)

```

```

print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)

```

```

print(X_test_title_bow.shape, y_test.shape)
print("="*100)

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)
# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")

(40000, 20) (40000,)
(10000, 20) (10000,)
=====
After vectorizations
(40000, 2879) (40000,)
(10000, 2879) (10000,)
=====
NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME

```

6 Set2 - TfIdf

```

In [0]: print(X_train.shape, y_train.shape)
print(X_test.shape, y_test.shape)

print("="*100)

vectorizer = TfidfVectorizer(min_df=10, max_features=10000)
vectorizer.fit(train_preprocessed_essays) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tf = vectorizer.transform(train_preprocessed_essays)
X_test_essay_tf = vectorizer.transform(test_preprocessed_essays)

print("After vectorizations")
print(X_train_essay_tf.shape, y_train.shape)
print(X_test_essay_tf.shape, y_test.shape)
print("="*100)

```

```

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)
# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")

```

```

(40000, 20) (40000,)
(10000, 20) (10000,)

```

=====

After vectorizations

```

(40000, 10000) (40000,)
(10000, 10000) (10000,)

```

=====

NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME

```

In [0]: print(X_train.shape, y_train.shape)
        print(X_test.shape, y_test.shape)

```

```

print("="*100)

```

```

vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,4), max_features=10000)
vectorizer.fit(train_preprocessed_titles) # fit has to happen only on train data

```

```

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tf = vectorizer.transform(train_preprocessed_titles)
X_test_title_tf = vectorizer.transform(test_preprocessed_titles)

```

```

print("After vectorizations")
print(X_train_title_tf.shape, y_train.shape)
print(X_test_title_tf.shape, y_test.shape)
print("="*100)

```

```

# print("YOU SHOULD NOT DO SOMETHING LIKE THIS")
# vectorizer = CountVectorizer()
# x_train_bow = vectorizer.fit_transform(X_train['essay'].values)
# x_cv_bow = vectorizer.fit_transform(X_cv['essay'].values)
# x_test_bow = vectorizer.fit_transform(X_test['essay'].values)

# print(x_train_bow.shape, y_train.shape)

```



```

# print(x_cv_bow.shape, y_cv.shape)
# print(x_test_bow.shape, y_test.shape)

print("NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME")

(40000, 20) (40000,)
(10000, 20) (10000,)
=====
After vectorizations
(40000, 2879) (40000,)
(10000, 2879) (10000,)
=====
NOTE: THE NUMBER OF COLUMNS IN EACH OF THE VECTOR WONT BE SAME

```

7 Set3 - Avg Word 2 Vec

```

In [0]: #https://stackoverflow.com/questions/21015674/list-object-has-no-attribute-shape
        #List to Numpy array
        #for Essays

X_train_essay_avgw2v = np.array(train_avg_w2v_vectors)
#X_cv_essay_avgw2v = np.array(cv_avg_w2v_vectors)
X_test_essay_avgw2v = np.array(test_avg_w2v_vectors)

#similarly, we are doing it for titles

X_train_title_avgw2v = np.array(train_avg_w2v_vectors1)
#X_cv_title_avgw2v = np.array(cv_avg_w2v_vectors1)
X_test_title_avgw2v = np.array(test_avg_w2v_vectors1)

In [0]: #For Essays - Avgw2v
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

print("After vectorizations")
print(X_train_essay_avgw2v.shape, y_train.shape)
#print(X_cv_essay_avgw2v.shape, y_cv.shape)
print(X_test_essay_avgw2v.shape, y_test.shape)
print("="*100)

(40000, 20) (40000,)
(10000, 20) (10000,)
=====
After vectorizations

```

```
(40000, 300) (40000,)  
(10000, 300) (10000,)  
=====
```

```
In [0]: #For Titles - Avgw2v  
        print(X_train.shape, y_train.shape)  
        #print(X_cv.shape, y_cv.shape)  
        print(X_test.shape, y_test.shape)  
  
        print("="*100)  
  
        print("After vectorizations")  
        print(X_train_title_avgw2v.shape, y_train.shape)  
        #print(X_cv_title_avgw2v.shape, y_cv.shape)  
        print(X_test_title_avgw2v.shape, y_test.shape)  
        print("="*100)
```

```
(40000, 20) (40000,)  
(10000, 20) (10000,)  
=====
```

```
After vectorizations  
(40000, 300) (40000,)  
(10000, 300) (10000,)  
=====
```

8 Set4 -TfIdf Weighted W2vec

```
In [0]: #https://stackoverflow.com/questions/21015674/list-object-has-no-attribute-shape  
        #List to Numpy array  
        #for Essays
```

```
X_train_es_tfidf_w2v = np.array(train_tfidf_w2v_vectors)  
#X_cv_es_tfidf_w2v = np.array(cv_tfidf_w2v_vectors)  
X_test_es_tfidf_w2v = np.array(test_tfidf_w2v_vectors)  
  
#similarly, we are doing it for titles  
  
X_train_title_tfidf_w2v = np.array(train_tfidf_w2v_vectors1)  
#X_cv_title_tfidf_w2v = np.array(cv_tfidf_w2v_vectors1)  
X_test_title_tfidf_w2v = np.array(test_tfidf_w2v_vectors1)
```

```
In [0]: #For Essays - TfIdf weighted W2vec  
        print(X_train.shape, y_train.shape)  
        #print(X_cv.shape, y_cv.shape)  
        print(X_test.shape, y_test.shape)
```

```

print("="*100)

print("After vectorizations")
print(X_train_es_tfidf_w2v.shape, y_train.shape)
#print(X_cv_es_tfidf_w2v.shape, y_cv.shape)
print(X_test_es_tfidf_w2v.shape, y_test.shape)
print("="*100)

```

```

(40000, 20) (40000,)
(10000, 20) (10000,)

```

```

=====
After vectorizations
(40000, 300) (40000,)
(10000, 300) (10000,)
=====

```

```

In [0]: #For Titles - TfIdf Weighted W2Vec
print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)

print("After vectorizations")
print(X_train_title_tfidf_w2v.shape, y_train.shape)
#print(X_cv_title_tfidf_w2v.shape, y_cv.shape)
print(X_test_title_tfidf_w2v.shape, y_test.shape)
print("="*100)

```

```

(40000, 20) (40000,)
(10000, 20) (10000,)

```

```

=====
After vectorizations
(40000, 300) (40000,)
(10000, 300) (10000,)
=====

```

```

In [0]: # Concatinating all the features

```

```

In [0]: #BoW

```

```

In [0]: from scipy.sparse import hstack
        # with the same hstack function we are concatenating a sparse matrix and a dense matrix
X_Bow_train = hstack((X_train_essay_bow, X_train_title_bow, tr_price_standardized, tr_prefix,
                      tr_school_state, tr_project_grade, tr_clean_subcategories))

print(X_Bow_train.shape, y_train.shape)

```

(40000, 12891) (40000,)

```
In [0]: X_Bow_test = hstack((X_test_essay_bow, X_test_title_bow, te_price_standardized, te_quantity_ratio,
                             te_prefix, te_school_state, te_project_grade, te_clean_subcategories))

print(X_Bow_test.shape, y_test.shape)
```

(10000, 12891) (10000,)

```
In [0]: #tfIdf
```

```
In [0]: from scipy.sparse import hstack
        # with the same hstack function we are concatenating a sparse matrix and a dense matrix
        X_tf_train = hstack((X_train_essay_tf, X_train_title_tf, tr_price_standardized, tr_quantity_ratio,
                             tr_prefix, tr_school_state, tr_project_grade, tr_clean_subcategories))

print(X_tf_train.shape, y_train.shape)
```

(40000, 12891) (40000,)

```
In [0]: X_tf_test = hstack((X_test_essay_tf, X_test_title_tf, te_price_standardized, te_quantity_ratio,
                             te_prefix, te_school_state, te_project_grade, te_clean_subcategories))

print(X_tf_test.shape, y_test.shape)
```

(10000, 12891) (10000,)

```
In [0]: #Avg W2V
```

```
In [0]: from scipy.sparse import hstack
        # with the same hstack function we are concatenating a sparse matrix and a dense matrix
        X_avg_w2v_train = np.hstack((X_train_essay_avgw2v, X_train_title_avgw2v, tr_price_standardized, tr_quantity_ratio,
                                      tr_prefix, tr_school_state, tr_project_grade, tr_clean_subcategories))

print(X_avg_w2v_train.shape, y_train.shape)
```

(40000, 612) (40000,)

```
In [0]: X_avg_w2v_test = np.hstack((X_test_essay_avgw2v, X_test_title_avgw2v, te_price_standardized, te_quantity_ratio,
                                      te_prefix, te_school_state, te_project_grade, te_clean_subcategories))

print(X_avg_w2v_test.shape, y_test.shape)
```

```
(10000, 612) (10000,)
```

```
In [0]: #Tfidf W2V
```

```
In [0]: from scipy.sparse import hstack
        # with the same hstack function we are concatenating a sparse matrix and a dense matrix
        X_tf_w2v_train = np.hstack((X_train_es_tfidf_w2v, X_train_title_tfidf_w2v, tr_price_standar
                                     tr_prefix, tr_school_state, tr_project_grade, tr_clean_subcate

        print(X_tf_w2v_train.shape, y_train.shape)
```

```
(40000, 612) (40000,)
```

```
In [0]: X_tf_w2v_test = np.hstack((X_test_es_tfidf_w2v, X_test_title_tfidf_w2v, te_price_standar
                                     te_prefix, te_school_state, te_project_grade, te_clean_subcate

        print(X_tf_w2v_test.shape, y_test.shape)
```

```
(10000, 612) (10000,)
```

2.4 Applying Random Forest

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

8.0.1 2.4.1 Applying Random Forests on BOW, SET 1

```
In [0]: from sklearn.ensemble import RandomForestClassifier
        from sklearn.model_selection import RandomizedSearchCV

        parameters = {'n_estimators': [500,1000,1200,1500] ,
                       'max_depth' : [5,6,7,8,9,10,15,20]}

        clf_rf = RandomForestClassifier(random_state=11,class_weight='balanced')

        clf=RandomizedSearchCV(clf_rf ,param_distributions = parameters, scoring="roc_auc", cv=
        clf.fit(X_Bow_train,y_train)
```

```
Out[0]: RandomizedSearchCV(cv=10, error_score=nan,
                           estimator=RandomForestClassifier(bootstrap=True,
                                                                ccp_alpha=0.0,
                                                                class_weight='balanced',
                                                                criterion='gini',
                                                                max_depth=None,
```

```

max_features='auto',
max_leaf_nodes=None,
max_samples=None,
min_impurity_decrease=0.0,
min_impurity_split=None,
min_samples_leaf=1,
min_samples_split=2,
min_weight_fraction_leaf=0.0,
n_estimators=100,
n_jobs=None,
oob_score=False,
random_state=11, verbose=0,
warm_start=False),
iid='deprecated', n_iter=10, n_jobs=None,
param_distributions={'max_depth': [5, 6, 7, 8, 9, 10, 15,
                                   20],
                    'n_estimators': [500, 1000, 1200,
                                       1500]},
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=True, scoring='roc_auc', verbose=0)

```

In [0]: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchC

```

a=clf.best_params_['n_estimators']
p = clf.best_params_['max_depth']
print(clf.best_score_)
print(a)
print(p)

```

0.7145762122407694

1200

20

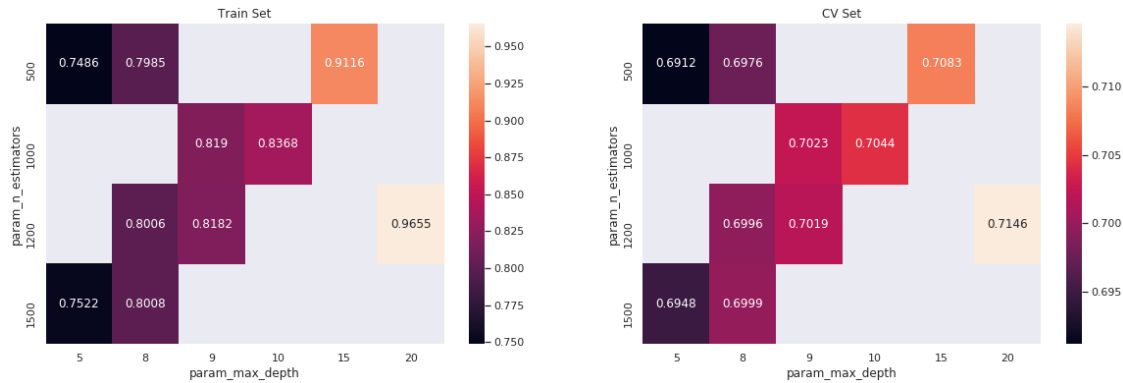
In [0]: <https://towardsdatascience.com/using-3d-visualizations-to-tune-hyperparameters-of-ml-1>

<https://github.com/woelop/Medium-posts/blob/master/3d%20cross%20validation/ML%206%20->

```

import seaborn as sns; sns.set()
max_scores1 = pd.DataFrame(clf.cv_results_).groupby(['param_n_estimators', 'param_max_c
fig, ax = plt.subplots(1,2, figsize=(20,6))
sns.heatmap(max_scores1.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores1.mean_test_score, annot = True, fmt='.4g', ax=ax[1])
ax[0].set_title('Train Set')
ax[1].set_title('CV Set')
plt.show()

```



In [0]: *# Train new model*

```
DT_tfidf = RandomForestClassifier(max_depth=p, n_estimators=a, class_weight='balanced')
DT_tfidf.fit(X_tf_train, y_train)
```

#https://scikitlearn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.h

```
y_train_pred1 = DT_tfidf.predict_proba(X_tf_train)[:,1]
```

```
y_test_pred1 = DT_tfidf.predict_proba(X_tf_test)[:,1]
```

```
train_fpr1, train_tpr1, tr_thresholds1 = roc_curve(y_train, y_train_pred1)
```

```
test_fpr1, test_tpr1, te_thresholds1 = roc_curve(y_test, y_test_pred1)
```

```
plt.plot(train_fpr1, train_tpr1, label="train AUC =" + str(auc(train_fpr1, train_tpr1)))
```

```
plt.plot(test_fpr1, test_tpr1, label="test AUC =" + str(auc(test_fpr1, test_tpr1)))
```

```
plt.plot([0, 1], [0, 1], 'r--')
```

```
plt.xlim([0, 1])
```

```
plt.ylim([0, 1])
```

```
plt.legend()
```

```
plt.xlabel("False Positive Rate")
```

```
plt.ylabel("True Positive Rate")
```

```
plt.title("ERROR PLOTS")
```

```
plt.grid(True)
```

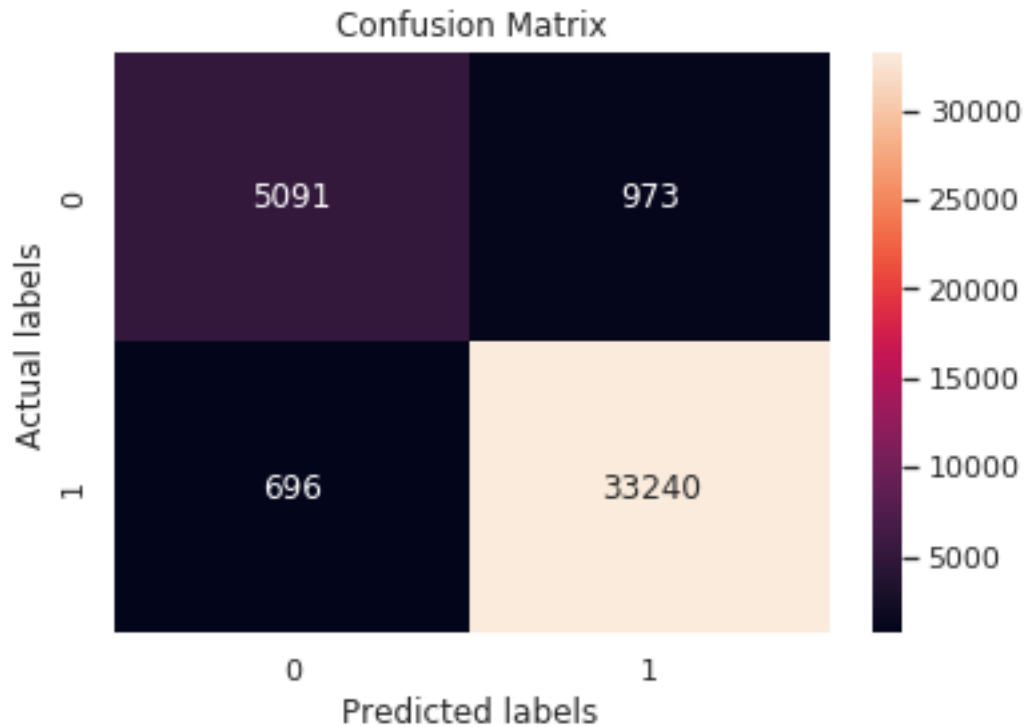
```
plt.show()
```



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On train")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, DT_tfidf.predict(X_tf_train)), annot=True, ax = ax)

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

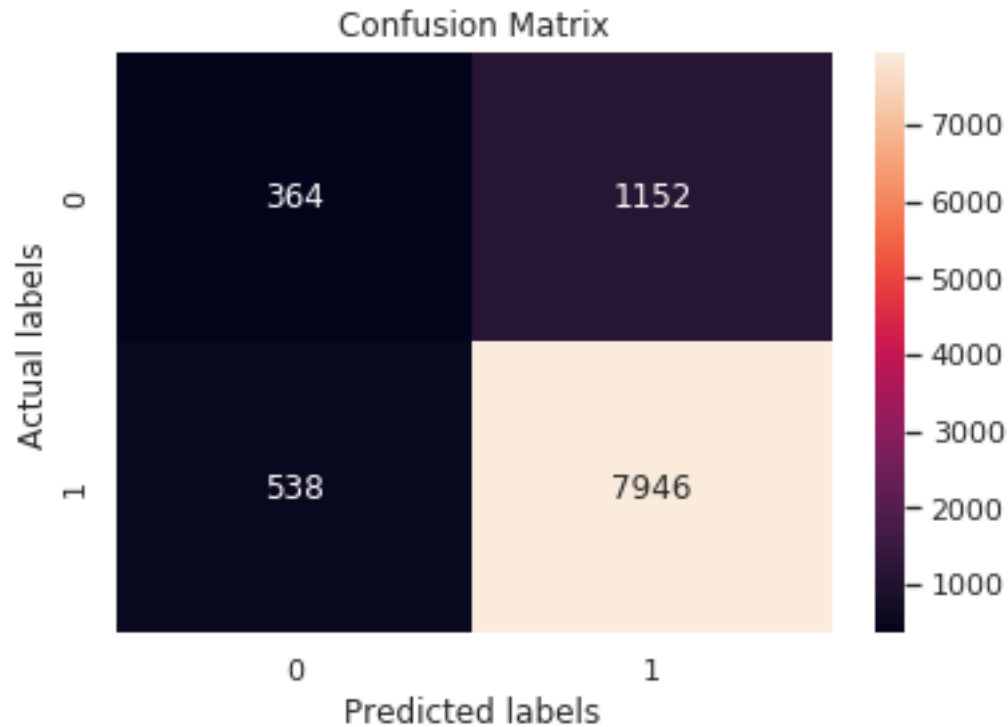
Confusin Matrix On train



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On test")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, DT_tfidf.predict(X_tf_test)), annot=True, ax = ax

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

Confusin Matrix On test



8.0.2 2.4.2 Applying Random Forests on TFIDE, SET 2

In [0]: *# Please write all the code with proper documentation*

```
In [0]: from sklearn.ensemble import RandomForestClassifier
        from sklearn.model_selection import RandomizedSearchCV
```

```
parameters1 = {'n_estimators': [5,10,50,100,200,500,1000] ,
               'max_depth' : [2,3,4,5,7,9,10]}
```

```
clf_rf1 = RandomForestClassifier(random_state=11,class_weight='balanced')
```

```
clf1=RandomizedSearchCV(clf_rf1 ,param_distributions = parameters1, scoring="roc_auc",
                        clf1.fit(X_tf_train,y_train)
```

```
Out[0]: RandomizedSearchCV(cv=10, error_score=nan,
                          estimator=RandomForestClassifier(bootstrap=True,
                                                              ccp_alpha=0.0,
                                                              class_weight='balanced',
                                                              criterion='gini',
                                                              max_depth=None,
```

```

max_features='auto',
max_leaf_nodes=None,
max_samples=None,
min_impurity_decrease=0.0,
min_impurity_split=None,
min_samples_leaf=1,
min_samples_split=2,
min_weight_fraction_leaf=0.0,
n_estimators=100,
n_jobs=None,
oob_score=False,
random_state=11, verbose=0,
warm_start=False),
iid='deprecated', n_iter=10, n_jobs=None,
param_distributions={'max_depth': [2, 3, 4, 5, 7, 9, 10],
                    'n_estimators': [5, 10, 50, 100, 200,
                                     500, 1000]},
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=True, scoring='roc_auc', verbose=0)

```

```

In [0]: #https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchC
a1=clf1.best_params_['n_estimators']
p1 = clf1.best_params_['max_depth']
print(clf1.best_score_)
print(a1)
print(p1)

```

0.7096677969578249

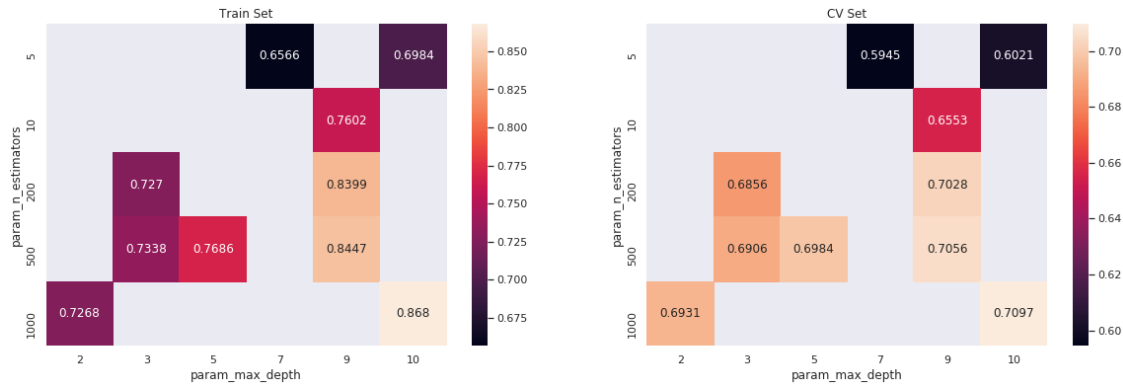
1000

10

```

In [0]: #https://towardsdatascience.com/using-3d-visualizations-to-tune-hyperparameters-of-ml-
#https://github.com/xoelop/Medium-posts/blob/master/3d%20cross%20validation/ML%206%20-
import seaborn as sns; sns.set()
max_scores1 = pd.DataFrame(clf1.cv_results_).groupby(['param_n_estimators', 'param_max_
fig, ax = plt.subplots(1,2, figsize=(20,6))
sns.heatmap(max_scores1.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores1.mean_test_score, annot = True, fmt='.4g', ax=ax[1])
ax[0].set_title('Train Set')
ax[1].set_title('CV Set')
plt.show()

```



In [0]: *# Train new model*

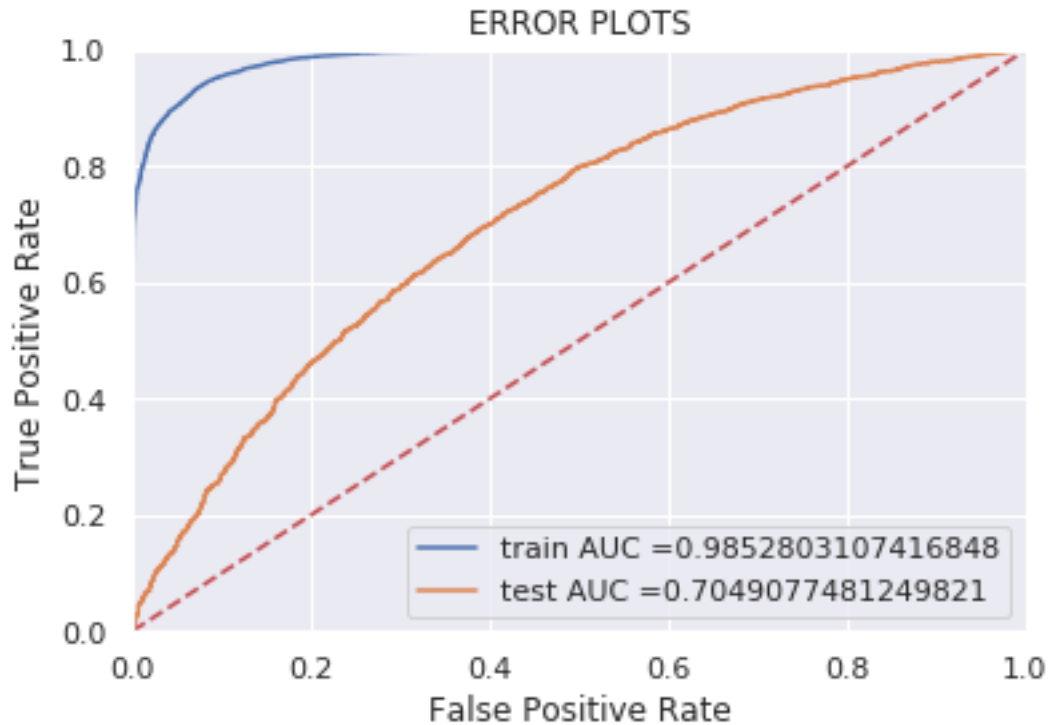
```
DT_tfidf1 = RandomForestClassifier(max_depth=p1, n_estimators=a1, class_weight='balance')
DT_tfidf1.fit(X_tf_train, y_train)
```

```
#https://scikitlearn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
y_train_pred1 = DT_tfidf.predict_proba(X_tf_train)[:,1]
y_test_pred1 = DT_tfidf.predict_proba(X_tf_test)[:,1]
```

```
train_fpr1, train_tpr1, tr_thresholds1 = roc_curve(y_train, y_train_pred1)
test_fpr1, test_tpr1, te_thresholds1 = roc_curve(y_test, y_test_pred1)
plt.plot(train_fpr1, train_tpr1, label="train AUC =" + str(auc(train_fpr1, train_tpr1)))
plt.plot(test_fpr1, test_tpr1, label="test AUC =" + str(auc(test_fpr1, test_tpr1)))
plt.legend()
```

```
plt.plot([0, 1], [0, 1], 'r--')
plt.xlim([0, 1])
plt.ylim([0, 1])
```

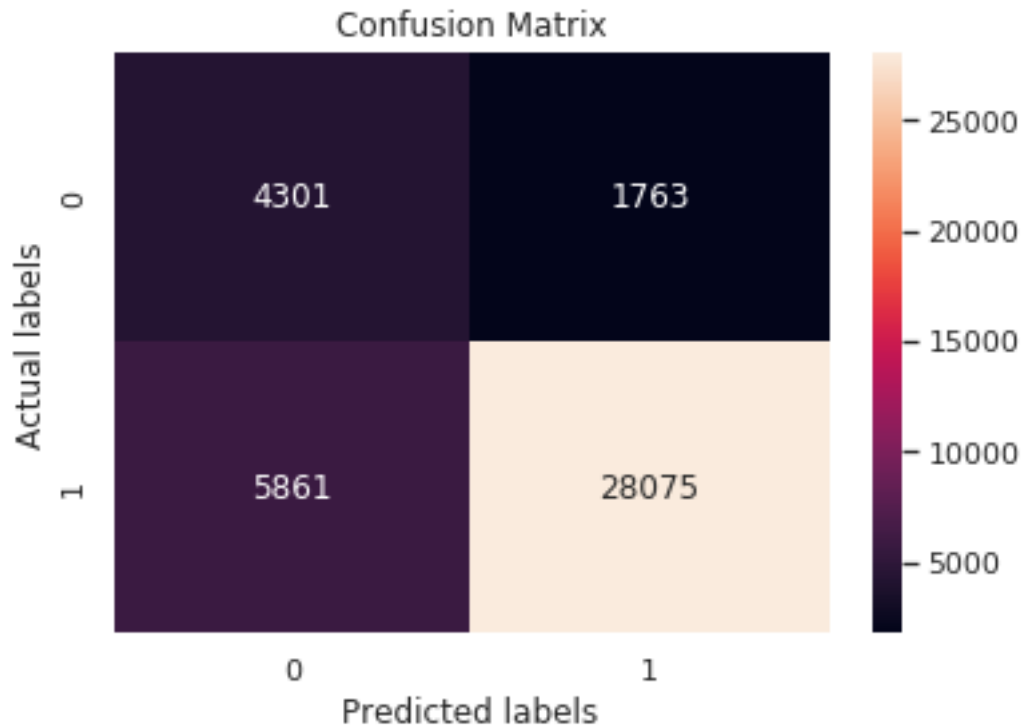
```
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On train")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, DT_tfidf1.predict(X_tf_train)), annot=True, ax =

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

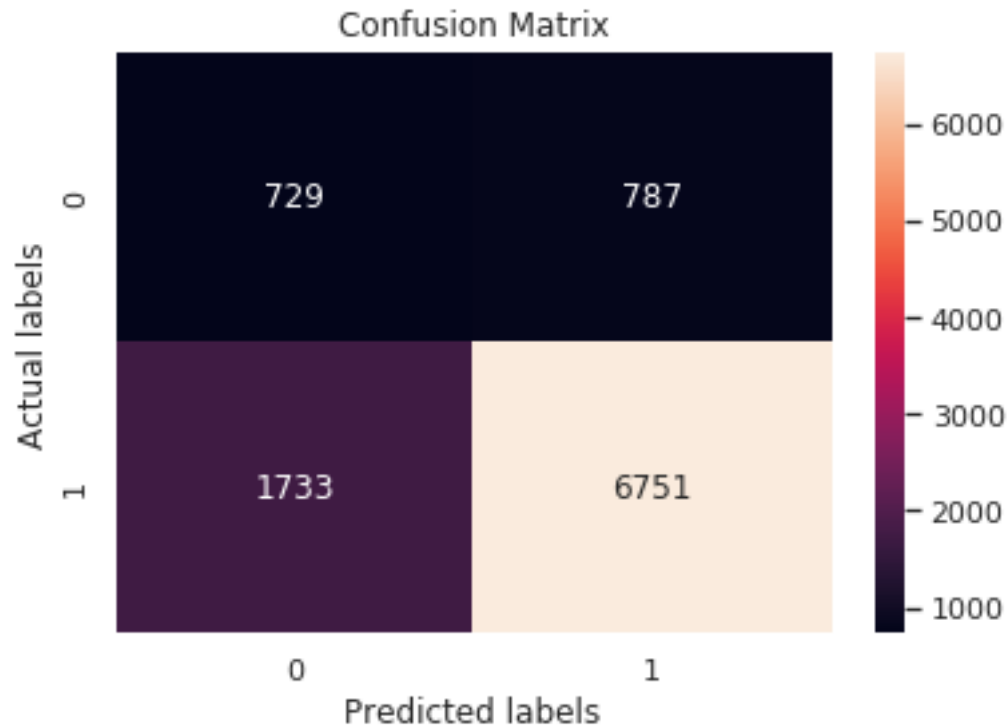
Confusin Matrix On train



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On test")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, DT_tfidf1.predict(X_tf_test)), annot=True, ax = ax

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

Confusin Matrix On test



8.0.3 2.4.3 Applying Random Forests on AVG W2V, SET 3

In [0]: *# Please write all the code with proper documentation*

```
In [0]: from sklearn.ensemble import RandomForestClassifier
        from sklearn.model_selection import RandomizedSearchCV
```

```
parameters2 = {'n_estimators': [5,10,50,100,200,500,1000] ,
               'max_depth' : [2,3,4,5,7,9,10]}
```

```
clf_rf2 = RandomForestClassifier(random_state=11,class_weight='balanced')
```

```
clf2=RandomizedSearchCV(clf_rf2 ,param_distributions = parameters2, scoring="roc_auc",
                        clf2.fit(X_avg_w2v_train,y_train)
```

```
Out[0]: RandomizedSearchCV(cv=10, error_score=nan,
                          estimator=RandomForestClassifier(bootstrap=True,
                                                              ccp_alpha=0.0,
                                                              class_weight='balanced',
                                                              criterion='gini',
                                                              max_depth=None,
```

```

max_features='auto',
max_leaf_nodes=None,
max_samples=None,
min_impurity_decrease=0.0,
min_impurity_split=None,
min_samples_leaf=1,
min_samples_split=2,
min_weight_fraction_leaf=0.0,
n_estimators=100,
n_jobs=None,
oob_score=False,
random_state=11, verbose=0,
warm_start=False),
iid='deprecated', n_iter=10, n_jobs=None,
param_distributions={'max_depth': [2, 3, 4, 5, 7, 9, 10],
                    'n_estimators': [5, 10, 50, 100, 200,
                                     500, 1000]},
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=True, scoring='roc_auc', verbose=0)

```

```

In [0]: #https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchC
a2=clf2.best_params_['n_estimators']
p2 = clf2.best_params_['max_depth']
print(clf2.best_score_)
print(a2)
print(p2)

```

0.704711971629938

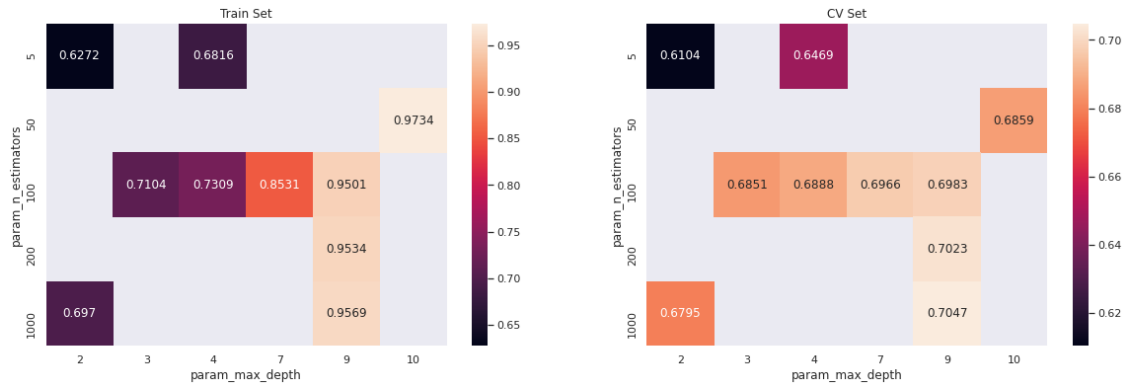
1000

9

```

In [0]: #https://towardsdatascience.com/using-3d-visualizations-to-tune-hyperparameters-of-ml-
#https://github.com/xoelop/Medium-posts/blob/master/3d%20cross%20validation/ML%206%20-
import seaborn as sns; sns.set()
max_scores1 = pd.DataFrame(clf2.cv_results_).groupby(['param_n_estimators', 'param_max_
fig, ax = plt.subplots(1,2, figsize=(20,6))
sns.heatmap(max_scores1.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores1.mean_test_score, annot = True, fmt='.4g', ax=ax[1])
ax[0].set_title('Train Set')
ax[1].set_title('CV Set')
plt.show()

```

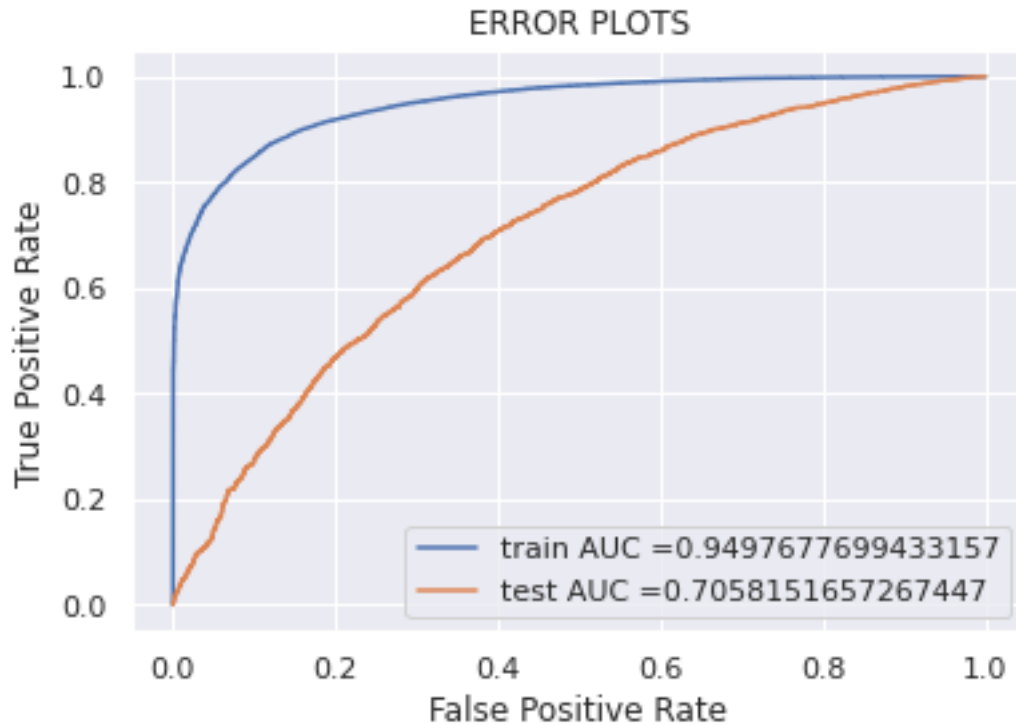



In [0]: # Train new model

```
DT_tfidf2 = RandomForestClassifier(max_depth=p2, n_estimators=a2, class_weight='balance')
DT_tfidf2.fit(X_avg_w2v_train, y_train)
```

```
#https://scikitlearn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
y_train_pred1 = DT_tfidf2.predict_proba(X_avg_w2v_train)[:,1]
y_test_pred1 = DT_tfidf2.predict_proba(X_avg_w2v_test)[:,1]
```

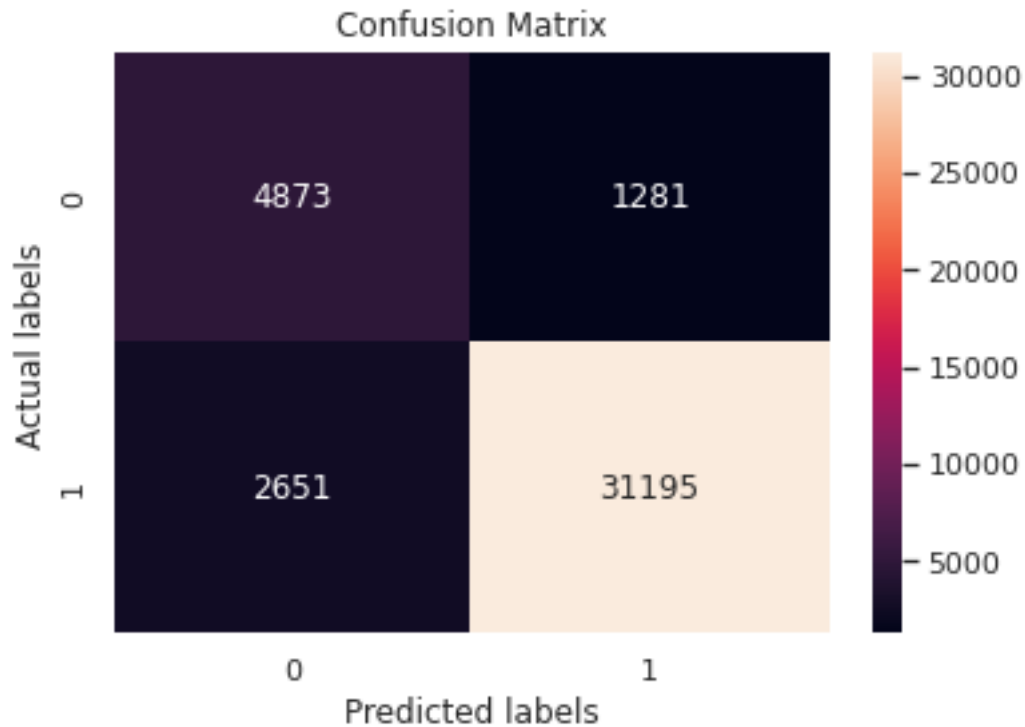
```
train_fpr1, train_tpr1, tr_thresholds1 = roc_curve(y_train, y_train_pred1)
test_fpr1, test_tpr1, te_thresholds1 = roc_curve(y_test, y_test_pred1)
plt.plot(train_fpr1, train_tpr1, label="train AUC =" + str(auc(train_fpr1, train_tpr1)))
plt.plot(test_fpr1, test_tpr1, label="test AUC =" + str(auc(test_fpr1, test_tpr1)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On train")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, DT_tfidf2.predict(X_avg_w2v_train)), annot=True,

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

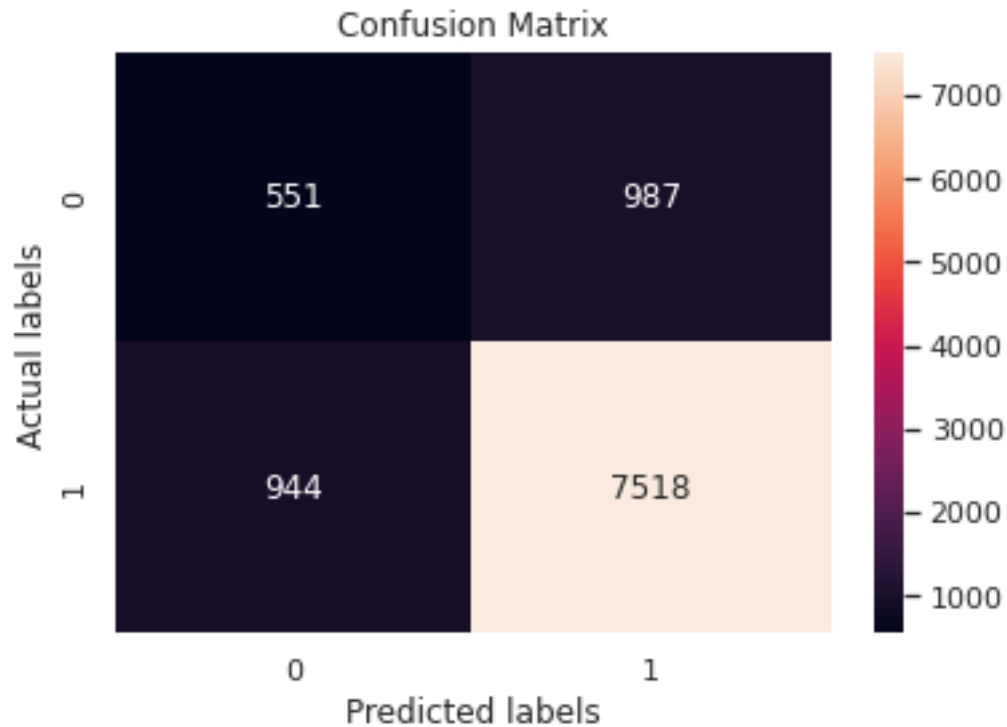
Confusin Matrix On train



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On test")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, DT_tfidf2.predict(X_avg_w2v_test)), annot=True, ax=

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

Confusin Matrix On test



8.0.4 2.4.4 Applying Random Forests on TFIDF W2V, SET 4

In [0]: *# Please write all the code with proper documentation*

```
In [0]: from sklearn.ensemble import RandomForestClassifier
        from sklearn.model_selection import RandomizedSearchCV
```

```
parameters3 = {'n_estimators': [5,10,50,100,200,500,1000] ,
               'max_depth' : [2,3,5,7,9,10]}
```

```
clf_rf3 = RandomForestClassifier(random_state=11,class_weight='balanced')
```

```
clf3=RandomizedSearchCV(clf_rf3 ,param_distributions = parameters3, scoring="roc_auc",
                        clf3.fit(X_tf_w2v_train,y_train)
```

```
Out[0]: RandomizedSearchCV(cv=10, error_score=nan,
                          estimator=RandomForestClassifier(bootstrap=True,
                                                              ccp_alpha=0.0,
                                                              class_weight='balanced',
                                                              criterion='gini',
                                                              max_depth=None,
```

```

max_features='auto',
max_leaf_nodes=None,
max_samples=None,
min_impurity_decrease=0.0,
min_impurity_split=None,
min_samples_leaf=1,
min_samples_split=2,
min_weight_fraction_leaf=0.0,
n_estimators=100,
n_jobs=None,
oob_score=False,
random_state=11, verbose=0,
warm_start=False),
iid='deprecated', n_iter=10, n_jobs=None,
param_distributions={'max_depth': [2, 3, 5, 7, 9, 10],
                    'n_estimators': [5, 10, 50, 100, 200,
                                     500, 1000]},
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=True, scoring='roc_auc', verbose=0)

```

```

In [0]: #https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchC
a3=clf3.best_params_['n_estimators']
p3 = clf3.best_params_['max_depth']
print(clf3.best_score_)
print(a3)
print(p3)

```

0.6986299347424445

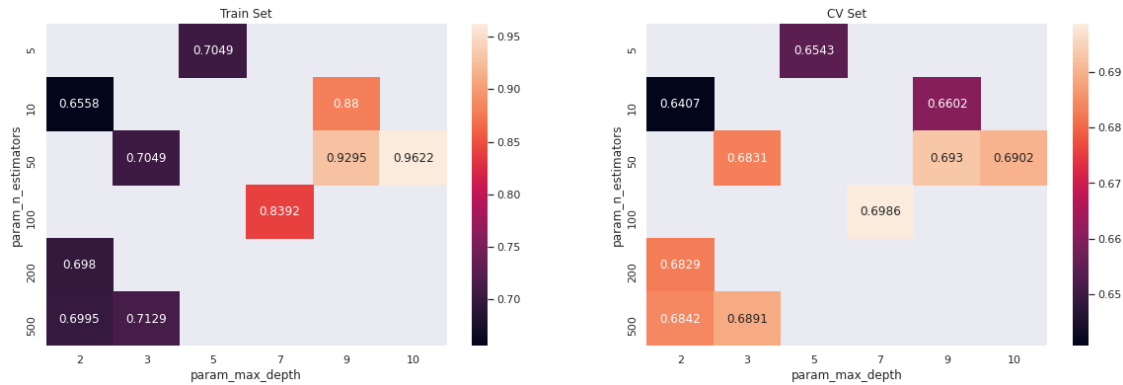
100

7

```

In [0]: #https://towardsdatascience.com/using-3d-visualizations-to-tune-hyperparameters-of-ml-
#https://github.com/xoelop/Medium-posts/blob/master/3d%20cross%20validation/ML%206%20-
import seaborn as sns; sns.set()
max_scores1 = pd.DataFrame(clf3.cv_results_).groupby(['param_n_estimators', 'param_max_
fig, ax = plt.subplots(1,2, figsize=(20,6))
sns.heatmap(max_scores1.mean_train_score, annot = True, fmt='.4g', ax=ax[0])
sns.heatmap(max_scores1.mean_test_score, annot = True, fmt='.4g', ax=ax[1])
ax[0].set_title('Train Set')
ax[1].set_title('CV Set')
plt.show()

```

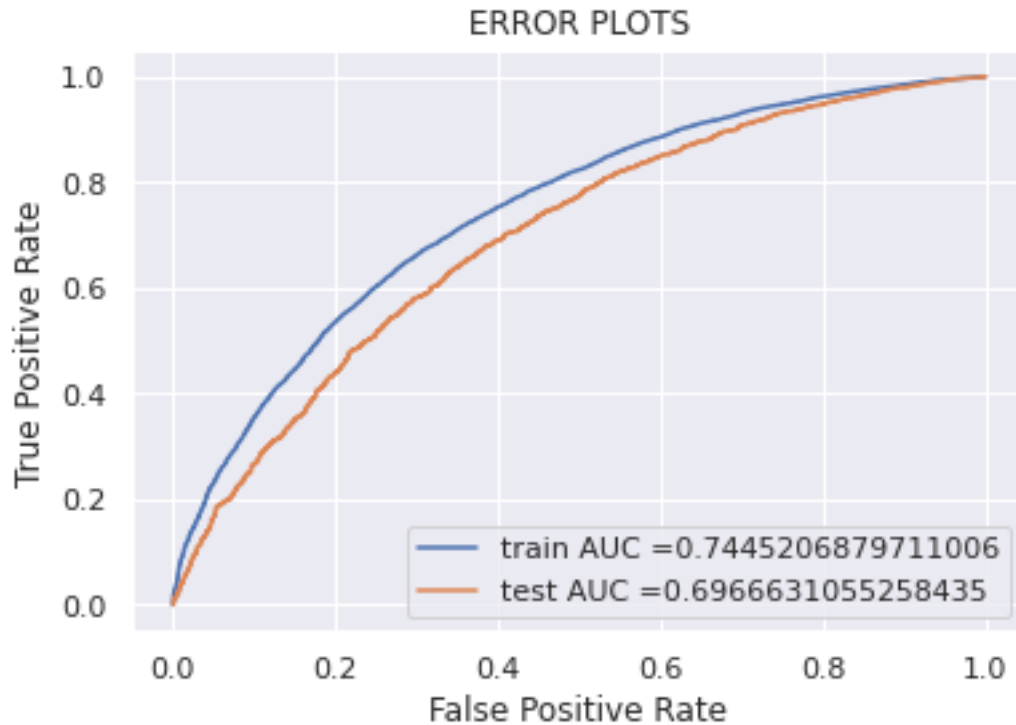


In [0]: *# Train new model*

```
DT_tfidf3 = RandomForestClassifier(max_depth=p3, n_estimators=a3, class_weight='balance')
DT_tfidf3.fit(X_avg_w2v_train, y_train)
```

```
#https://scikitlearn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
y_train_pred1 = DT_tfidf3.predict_proba(X_tf_w2v_train)[:,1]
y_test_pred1 = DT_tfidf3.predict_proba(X_tf_w2v_test)[:,1]
```

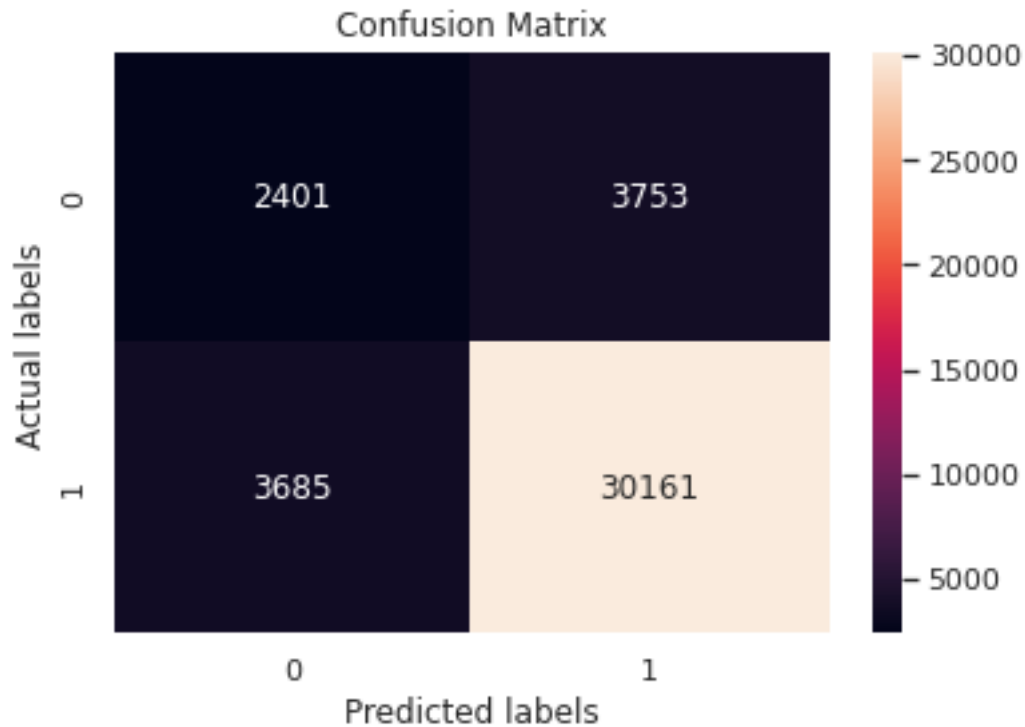
```
train_fpr1, train_tpr1, tr_thresholds1 = roc_curve(y_train, y_train_pred1)
test_fpr1, test_tpr1, te_thresholds1 = roc_curve(y_test, y_test_pred1)
plt.plot(train_fpr1, train_tpr1, label="train AUC =" + str(auc(train_fpr1, train_tpr1)))
plt.plot(test_fpr1, test_tpr1, label="test AUC =" + str(auc(test_fpr1, test_tpr1)))
plt.legend()
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On train")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_train, DT_tfidf3.predict(X_tf_w2v_train)), annot=True, a

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

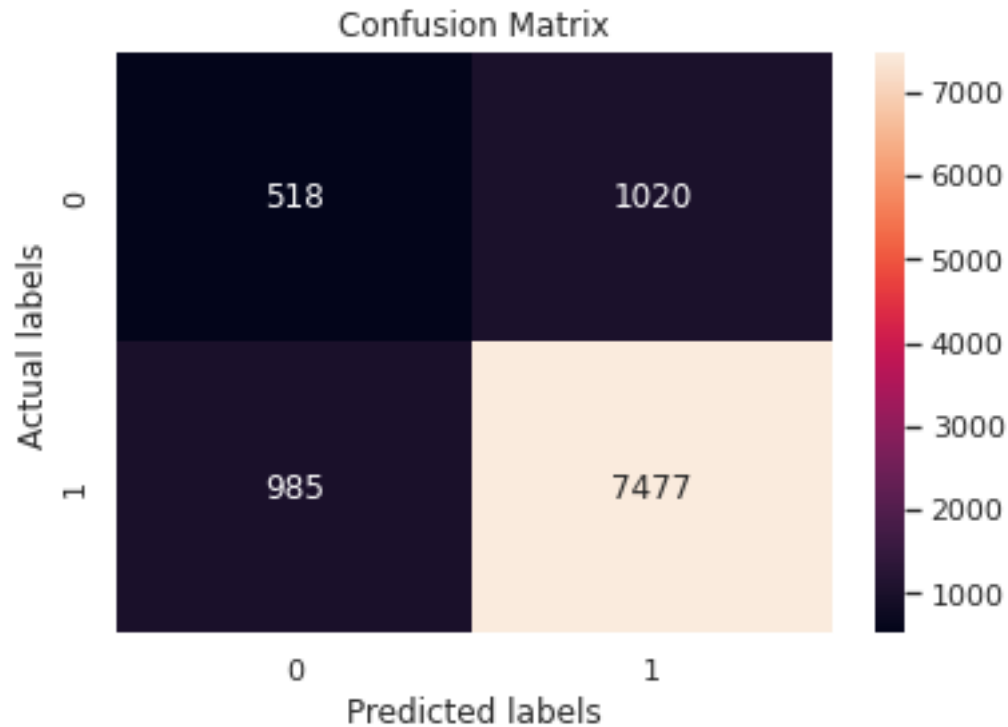
Confusin Matrix On train



```
In [0]: #https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
#Train dataset - COnfusion Matrix
import seaborn as sns
import matplotlib.pyplot as plt
print("Confusin Matrix On test")
ax= plt.subplot()
sns.heatmap(confusion_matrix(y_test, DT_tfidf3.predict(X_tf_w2v_test)), annot=True, ax=

# labels, title and ticks
ax.set_xlabel('Predicted labels');
ax.set_ylabel('Actual labels');
ax.set_title('Confusion Matrix');
```

Confusin Matrix On test



3. Conclusion

In [0]: *# Please compare all your models using Prettytable library*

In [2]: *# Please compare all your models using Prettytable library*

```
from prettytable import PrettyTable
```

```
p = PrettyTable()
```

```
p.field_names = ["Vetorizer", "Model", "n_estimators", "max_depth", "Test AUC"]
```

```
p.add_row(["SET 1- Bow", "Random Forest", "1200", "20", "0.704" ])
```

```
p.add_row(["SET 2- TF_IDF", "Random Forest", "1000", "10", "0.704"])
```

```
p.add_row(["SET 3- AVG w2v", "Random Forest", "1000", "9", "0.705" ])
```

```
p.add_row(["SET 4- TF_IDF w2v", "Random Forest", "100", "7", "0.696" ])
```

```
print(p)
```

Vetorizer	Model	n_estimators	max_depth	Test AUC
SET 1- Bow	Random Forest	1200	20	0.704
SET 2- TF_IDF	Random Forest	1000	10	0.704
SET 3- AVG w2v	Random Forest	1000	9	0.705

SET 4- TF_IDF w2v Random Forest	100		7		0.696	
+-----+-----+-----+-----+-----+						