# 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 |
|---|---|
| project_id | A unique identifier for the proposed project. **Example:** `p036502` |
| project_title | Title of the project. **Examples:**<br>• `Art Will Make You Happy!`<br>• `First Grade Fun` |
| project_grade_category | Grade level of students for which the project is targeted. One of the following enumerated values:<br>• `Grades PreK-2`<br>• `Grades 3-5`<br>• `Grades 6-8`<br>• `Grades 9-12` |
| project_subject_categories | One or more (comma-separated) subject categories for the project from the following enumerated list of values:<br>• `Applied Learning`<br>• `Care & Hunger`<br>• `Health & Sports`<br>• `History & Civics`<br>• `Literacy & Language`<br>• `Math & Science`<br>• `Music & The Arts`<br>• `Special Needs`<br>• `Warmth`<br><br>**Examples:**<br>• `Music & The Arts`<br>• `Literacy & Language, Math & Science` |
| school_state | State where school is located ([Two-letter U.S. postal code](https://en.wikipedia.org/wiki/List_of_U.S._state_abbreviations#Postal_codes)). **Example:** `WY` |
| project_subject_subcategories | One or more (comma-separated) subject subcategories for the project. **Examples:**<br>• `Literacy`<br>• `Literature & Writing, Social Sciences` |
| project_resource_summary | An explanation of the resources needed for the project. **Example:**<br>• `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:<br>• `nan`<br>• `Dr.`<br>• `Mr.`<br>• `Mrs.`<br>• `Ms.`<br>• `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 `project_id` value from the `train.csv` file. **Example:** `p036502` |
| description | Desciption of the resource. **Example:** `Tenor Saxophone Reeds, Box of 25` |
| quantity | Quantity of the resource required. **Example:** `3` |
| price | Price of the resource required. **Example:** `9.95` |

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

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

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

```python
In [5]:
1  ## import all the modules
2  %matplotlib inline
3  import warnings
4  warnings.filterwarnings("ignore")
5  import scipy
6  import sqlite3
7  import pandas as pd
8  import numpy as np
9  import nltk
10 import string
11 from scipy import sparse
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14 from sklearn.feature_extraction.text import TfidfTransformer
15 from sklearn.feature_extraction.text import TfidfVectorizer
16 from sklearn.model_selection import train_test_split
17 from sklearn.feature_extraction.text import CountVectorizer
18 from sklearn.metrics import confusion_matrix
19 from sklearn.metrics import roc_auc_score
20 from sklearn import metrics
21 from sklearn.metrics import roc_curve, auc
22 from nltk.stem.porter import PorterStemmer
23 from sklearn.preprocessing import Normalizer
24 from scipy.sparse import hstack
25 from sklearn.tree import DecisionTreeClassifier
26 import re
27 from sklearn.ensemble import RandomForestClassifier
28 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
29 import string
30 from nltk.corpus import stopwords
31 from nltk.stem import PorterStemmer
32 from nltk.stem.wordnet import WordNetLemmatizer
33 import nltk
34 from nltk.sentiment import SentimentIntensityAnalyzer
35 from gensim.models import Word2Vec
36 from gensim.models import KeyedVectors
37 import pickle
38 from tqdm import tqdm
39 import os
40 from sklearn.tree import export_graphviz
41 from os import system
42 import graphviz
43 from six import StringIO
44 from chart_studio import plotly
45 import plotly.offline as offline
46 import plotly.graph_objs as go
47 offline.init_notebook_mode()
48 from collections import Counter
```

```python
In [4]:
1  !pip install six
```

Requirement already satisfied: six in c:\users\sundararaman\anaconda3\lib\site-packages (1.15.0)

## 1. Reading the Data

```python
In [6]:
1  project_data=pd.read_csv('train_data.csv')
2  resource_data=pd.read_csv('resources.csv')
```

```
In [7]:    1  ## Check the shape and attributes of the project data
           2  print("Number of data points in project train data", project_data.shape)
           3  print('-'*50)
           4  print("The attributes of data :", project_data.columns.values)
```

```
Number of data points in project 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 [8]:    1  ## Check the shape and attributes of the resource data
           2  print("Number of data points in resource train data", resource_data.shape)
           3  print(resource_data.columns.values)
           4  resource_data.head(2)
```

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

Out[8]:

|   | 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 |

## 1.1 Preprocessing Categorical Features: project_grade_category

```
In [9]:    1  print("Project grade" ,project_data['project_grade_category'].value_counts(dropna=False))
           2  ## visulaize how project grade looks like
           3  print('-'*50)
           4  print(project_data['project_grade_category'].values[1000])
           5  print(project_data['project_grade_category'].values[1500])
```

```
Project grade Grades PreK-2    44225
Grades 3-5    37137
Grades 6-8    16923
Grades 9-12    10963
Name: project_grade_category, dtype: int64
--------------------------------------------------
Grades 3-5
Grades PreK-2
```

```
In [10]:   1  # https://stackoverflow.com/questions/36383821/pandas-dataframe-apply-function-to-column-strings-based-on-other-column-value
           2  project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('Grades ','')
           3  project_data['project_grade_category'] = project_data['project_grade_category'].str.replace(' ','_')
           4  project_data['project_grade_category'] = project_data['project_grade_category'].str.replace('-','_')
           5  project_data['project_grade_category'] = project_data['project_grade_category'].str.lower()
           6  project_data['project_grade_category'].value_counts()
```

```
Out[10]: prek_2    44225
3_5    37137
6_8    16923
9_12    10963
Name: project_grade_category, dtype: int64
```

## 1.2 Preprocessing Categorical Features: project_subject_category

```
In [11]:   1  catogories = list(project_data['project_subject_categories'].values)
           2  # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
           3  # reference from course material : reference_EDA.ipynb
           4  # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
           5  # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
           6  # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
           7  cat_list = []
           8  for i in catogories:
           9      temp = ""
          10      # consider we have text like this "Math & Science, Warmth, Care & Hunger"
          11      for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
          12          if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
          13              j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
          14          j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
          15          temp+=j.strip()+" " # "#" abc ".strip() will return "abc", remove the trailing spaces
          16          temp = temp.replace('&','_') # we are replacing the & value into
          17      cat_list.append(temp.strip())
          18
          19  project_data['clean_categories'] = cat_list
          20  project_data.drop(['project_subject_categories'], axis=1, inplace=True)
          21  project_data.head(2)
          22
          23
          24  ### maintain a dict that stores count of values
          25  my_counter=Counter()
          26
          27  for word in project_data['clean_categories'].values:
          28      my_counter.update(word.split())
          29  cat_dict=dict(my_counter)
          30
          31  sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

## 1.3 Preprocessing Categorical Features: project_subject_subcategory

```
In [12]:   1  sub_catogories = list(project_data['project_subject_subcategories'].values)
           2  # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
           3
           4  # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
           5  # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
           6  # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
           7
           8  sub_cat_list = []
           9  for i in sub_catogories:
          10      temp = ""
          11      # consider we have text like this "Math & Science, Warmth, Care & Hunger"
          12      for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
          13          if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Science"
          14              j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
          15          j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
          16          temp +=j.strip()+" " # "#" abc ".strip() will return "abc", remove the trailing spaces
          17          temp = temp.replace('&','_')
          18      sub_cat_list.append(temp.strip())
          19
          20  project_data['clean_subcategories'] = sub_cat_list
          21  project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
          22
          23  # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
          24  my_counter = Counter()
          25  for word in project_data['clean_subcategories'].values:
          26      my_counter.update(word.split())
          27
          28  sub_cat_dict = dict(my_counter)
          29  sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
```

## 1.3 Preprocessing Categorical Features: school_state

```
In [13]:   1  project_data['school_state'].value_counts()
           2  ## Convert it to lower
           3  project_data['school_state'] = project_data['school_state'].str.lower()
           4  #project_data['school_state'].value_counts(dropna=False)
```

## 1.4 Preprocessing Categorical Features: Teacher_prefix

```
In [14]:   1  print(project_data['teacher_prefix'].value_counts(dropna=False))
           2  # try to remove the dots from the teacher prefix and replace nan with mrs.
           3  project_data['teacher_prefix']=project_data['teacher_prefix'].fillna('Mrs.')
           4  project_data['teacher_prefix']=project_data['teacher_prefix'].str.replace('.','')
           5  project_data['teacher_prefix']=project_data['teacher_prefix'].str.lower()
           6  project_data['teacher_prefix']=project_data['teacher_prefix'].str.strip()
```

```
Mrs.      57269
Ms.       38955
Mr.       10648
Teacher    2360
Dr.          13
NaN           3
Name: teacher_prefix, dtype: int64
```

## 1.5 Combining all the essays

```
In [15]:   1  # merge two column text dataframe:
           2  project_data["essay"] = project_data["project_essay_1"].map(str) +\
           3                          project_data["project_essay_2"].map(str) + \
           4                          project_data["project_essay_3"].map(str) + \
           5                          project_data["project_essay_4"].map(str)
```

## 1.6 Number of Words in the Essay and Title

```
In [16]:  1  source:''''https://www.geeksforgeeks.org/python-program-to-count-words-in-a-sentence/'''
          2  words_counter=[]
          3  for string in project_data['essay']:
          4      res = len(re.findall(r'\w+', string))
          5      words_counter.append(res)
          6
          7  project_data["words_in_essay"] = words_counter
          8
          9  words_counter=[]
         10
         11  for string in project_data['project_title']:
         12      res = len(re.findall(r'\w+', string))
         13      words_counter.append(res)
         14  project_data["words_in_title"] = words_counter
```

## 1.7. Preprocessing Numerical Values: price

```
In [17]:  1  ## calculate the overall count of resources and the total price for each project id
          2  price_data=resource_data.groupby('id',as_index=False).agg({'price':'sum','quantity':'sum' })
          3
```

```
In [18]:  1  project_data = pd.merge(project_data,price_data,on='id',how='left')
```

## 1.8 Preprocessing Text Features: project_title

```
In [19]:  1  # https://stackoverflow.com/a/47091490/408403
          2  def decontracted(phrase):
          3      # specific
          4      phrase = re.sub(r"won't", "will not", phrase)
          5      phrase = re.sub(r"can\'t", "can not", phrase)
          6      # general
          7      phrase = re.sub(r"n\'t", " not", phrase)
          8      phrase = re.sub(r"\'re", " are", phrase)
          9      phrase = re.sub(r"\'s", " is", phrase)
         10      phrase = re.sub(r"\'d", " would", phrase)
         11      phrase = re.sub(r"\'ll", " will", phrase)
         12      phrase = re.sub(r"\'t", " not", phrase)
         13      phrase = re.sub(r"\'ve", " have", phrase)
         14      phrase = re.sub(r"\'m", " am", phrase)
         15      return phrase
```

```
In [20]:  1  # https://gist.github.com/sebleier/554280
          2  # we are removing the words from the stop words list: 'no', 'nor', 'not'
          3  stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've",\
          4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
          5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their',\
          6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
          7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
          8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
          9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after',\
         10             'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further',\
         11             'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more',\
         12             'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
         13             's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
         14             've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn',\
         15             "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn',\
         16             "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
         17             'won', "won't", 'wouldn', "wouldn't"]
```

```
In [21]:  1  print("printing some random reviews")
          2  print(9, project_data['project_title'].values[9])
          3  print(34, project_data['project_title'].values[34])
          4  print(147, project_data['project_title'].values[147])
```

```
printing some random reviews
9 Just For the Love of Reading--\r\nPure Pleasure
34 \"Have A Ball!!!\"
147 Who needs a Chromebook?\r\nWE DO!!
```

```
In [22]:  1  # Combining all the above stundents
          2
          3  def preprocess_text(text_data):
          4      preprocessed_text = []
          5      # tqdm is for printing the status bar
          6      for sentence in tqdm(text_data):
          7          sent = decontracted(sentence)
          8          sent = sent.replace('\\r', ' ')
          9          sent = sent.replace('\\n', ' ')
         10          sent = sent.replace('\\"', ' ')
         11          sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
         12          # https://gist.github.com/sebleier/554280
         13          sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
         14          preprocessed_text.append(sent.lower().strip())
         15      return preprocessed_text
```

```
In [23]:  1  preprocessed_titles = preprocess_text(project_data['project_title'].values)
```

```
100%|████████| 109248/109248 [00:23<00:00, 4647.98it/s]
```

```
In [24]:  1  print("printing some random reviews")
          2  print(9, preprocessed_titles[9])
          3  print(34, preprocessed_titles[34])
          4  print(147, preprocessed_titles[147])
```

```
printing some random reviews
9 love reading pure pleasure
34 ball
147 needs chromebook
```

## 1.9 Preprocessing Text Features: essay

```
In [25]:  1  print("printing some random essay")
          2  print(9, project_data['essay'].values[9])
          3  print('-'*50)
          4  print(34, project_data['essay'].values[34])
          5  print('-'*50)
          6  print(147, project_data['essay'].values[147])
```

```
printing some random essay
9 Over 95% of my students are on free or reduced lunch.  I have a few who are homeless, but despite that, they come to school with an eagerness to learn.  My students are inquisitive eager learners who  embrace the challenge of not having great books and other resources  every day.  Many of them are not afforded the opportunity to engage with these big colorful pages of a book on a regular basis at home and they don't travel to the public library.  \r\nIt is my duty as a teacher to do all I can to provide each student an opportunity to succeed in every aspect of life. \r\nReading is Fundamental! My students will read these books over and over again while boosting their comprehension skills. These books will be used for read alouds, partner reading and for Independent reading. \r\nThey will engage in reading to build their \"Love for Reading\" by reading for pure enjoyment. They will be introduced to some new authors as well as some old favorites. I want my students to be ready for the 21st Century and know the pleasure of holding a good hard back book in hand. There's nothing like a good book to read!  \r\nMy students will soar in Reading, and more because of your consideration and generous funding contribution. This will help build stamina and prepare for 3rd grade. Thank you so much for reading our proposal!nannan
--------------------------------------------------
34 My students mainly come from extremely low-income families, and the majority of them come from homes where both parents work full time. Most of my students are at school from 7:30 am to 6:00 pm (2:30 to 6:00 pm in the after-school program), and they all receive free and reduced meals for breakfast and lunch. \r\n\r\n\r\nI want my students to feel  as comfortable in my classroom as they do at home. Many of my students take on multiple roles both at home as well as in their  classroom  every day.  Some of them are the caretakers of younger siblings, cooks, babysitters, academics, friends, and most of all, they are developing who they are going to become as adults.  I consider it an essential part of my job to model helping others gain knowledge in a positive manner. As a result, I have a community of students who love helping each other in and outside of the classroom. They consistently look for opportunities to support each other's learning in a kind and helpful way.I am excited to be experimenting with alternative seating in my classroom this school year. Studies have shown that giving students the option of where they sit in a classroom increases focus as well as motivation.  \r\n\r\nBy allowing students choice in the classroom, they are able to explore and create in a welcoming environment. Alternative classroom seating has been experimented with more frequently in recent years. I believe (along with many others), that every child learns differently. This does not only apply to how multiplication is memorized, or a paper is written, but applies to the space in which they are asked to work. I have had students in the past ask \"Can I work in the library? Can I work on the carpet?\" My answer was always, \"As long as you're learning, you can work wherever you want!\" \r\n\r\n\r\nWith the yoga balls and the lap-desks, I will be able to increase the options for seating in my classroom and expand its imaginable space.nannan
--------------------------------------------------
147 My students are eager to learn and make their mark on the world.\r\n\r\n\r\nThey come from a Title 1 school and need extra love.\r\n\r\n\r\nMy fourth grade students are in a high poverty area and still come to school every day to get their education. I am trying to make it fun and educational for them so they can get the most out of their schooling. I created a caring environment for the students to bloom! They deserve the best.\r\nThank you!\r\nI am requesting 1 Chromebook to access online interventions, differentiate instruction, and get extra practice. The Chromebook will be used to supplement ELA and math instruction. Students will play ELA and math games that are engaging and fun, as well as participate in assignments online. This in turn will help my students improve their skills. Having a Chromebook in the classroom would not only allow students to use the programs at their own pace, but would ensure more students are getting adequate time to use the programs. The online programs have been especially beneficial to my students with special needs. They are able to work at their level as well as be challenged with some different materials. This is making these students more confident in their abilities.\r\n\r\nThe Chromebook will allow my students to have daily access to computers and increase their computing skills.\r\nThis will change their lives for the better as they become more successful in school. Having access to technology in the classroom would help bridge the achievement gap.nannan
```

```
In [26]:  1  preprocessed_essays = preprocess_text(project_data['essay'].values)
```

```
100%|████████| 109248/109248 [08:15<00:00, 220.50it/s]
```

```
In [27]:  1  print("printing some random essay")
          2  print(9, preprocessed_essays[9])
          3  print('-'*50)
          4  print(34, preprocessed_essays[34])
          5  print('-'*50)
          6  print(147, preprocessed_essays[147])
          7
          8  #merge the column in the project_data
          9  project_data['processed_essay']=preprocessed_essays
```

printing some random essay
9 95 students free reduced lunch homeless despite come school eagerness learn students inquisitive eager learners embrace challenge not great books resources every day many not afforded opportunity engage big colorful pages book regular basis home no
t travel public library duty teacher provide student opportunity succeed every aspect life reading fundamental students read books boosting comprehension skills books used read alouds partner reading independent reading engage reading build love read
ing reading pure enjoyment introduced new authors well old favorites want students ready 21st century know pleasure holding good hard back book hand nothing like good book read students soar reading consideration generous funding contribution help bu
ild stamina prepare 3rd grade thank much reading proposal nannan
--------------------------------------------------
34 students mainly come extremely low income families majority come homes parents work full time students school 7 30 6 00 pm 2 30 6 00 pm school program receive free reduced meals breakfast lunch want students feel comfortable classroom home many st
udents take multiple roles home well school sometimes caretakers younger siblings cooks babysitters academics friends developing going become adults consider essential part job model helping others gain knowledge positive manner result community stud
ents love helping outside classroom consistently look opportunities support learning kind helpful way excited experimenting alternative seating classroom school year studies shown giving students option sit classroom increases focus well motivation a
llowing students choice classroom able explore create welcoming environment alternative classroom seating experimented frequently recent years believe along many others every child learns differently not apply multiplication memorized paper written a
pplies space asked work students past ask work library work carpet answer always long learning work wherever want yoga balls lap desks able increase options seating classroom expand imaginable space nannan
--------------------------------------------------
147 students eager learn make mark world come title 1 school need extra love fourth grade students high poverty area still come school every day get education trying make fun educational get schooling created caring environment students bloom deserve
best thank requesting 1 chromebook access online interventions differentiate instruction get extra practice chromebook used supplement ela math instruction students play ela math games engaging fun well participate assignments online turn help studen
ts improve skills chromebook classroom would not allow students use programs pace would ensure students getting adequate time use programs online programs especially beneficial students special needs able work level well challenged different material
s making students confident abilities chromebook would allow students daily access computers increase computing skills change lives better become successful school access technology classroom would help bridge achievement gap nannan

### 1.10 Preprocessing Text Features: Project Title

```
In [28]:  1  print("printing some random project_titles")
          2  print(9, project_data['project_title'].values[9])
          3  print('-'*50)
          4  print(34, project_data['project_title'].values[34])
          5  print('-'*50)
          6  print(147, project_data['project_title'].values[147])
```

printing some random project_titles
9 Just For the Love of Reading--\r\nPure Pleasure
--------------------------------------------------
34 \"Have A Ball!!!\"
--------------------------------------------------
147 Who needs a Chromebook?\r\nWE DO!!

```
In [29]:  1  processed_title = preprocess_text(project_data['project_title'].values)
```

100%|████████| 109248/109248 [00:16<00:00, 6788.75it/s]

```
In [30]:  1  print("printing some random title")
          2  print(9, preprocessed_titles[9])
          3  print('-'*50)
          4  print(34, preprocessed_titles[34])
          5  print('-'*50)
          6  print(147, preprocessed_titles[147])
          7
          8  #merge the column in the project_data
          9  project_data['processed_title']=processed_title
```

printing some random title
9 love reading pure pleasure
--------------------------------------------------
34 ball
--------------------------------------------------
147 needs chromebook

### 1.11 Creating sentiment columns

```
In [32]:   1  ## craete the sentiment columns using  essay
           2  neg=[]
           3  pos=[]
           4  neu=[]
           5  compound=[]
           6  sentiment_model=SentimentIntensityAnalyzer()
           7  for text in project_data['processed_essay']:
           8      pol_scores = sentiment_model.polarity_scores(text)
           9      neg.append(pol_scores['neg'])
          10      pos.append(pol_scores['pos'])
          11      neu.append(pol_scores['neu'])
          12      compound.append(pol_scores['compound'])
          13
          14  project_data['pos']=pos
          15  project_data['neg']=neg
          16  project_data['neu']=neu
          17  project_data['compound']=compound
```

### 1.12 Dropping redundant columns before splitting

```
In [33]:  1  project_data.drop(columns=['Unnamed: 0','project_essay_1', 'project_essay_2', 'project_essay_3','project_essay_4','project_title','essay'],inplace=True)
```

## 2 Train,Test,CV Split

```
In [52]:  1  # train test split using sklearn.model selection
          2  ## Considering only 50k points
          3  project_data=project_data.sample(50000,random_state=1)
          4  X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'], test_size=0.33, stratify = project_data['project_is_approved'],random_state=0)
          5  X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33, stratify=y_train,random_state=0)
```

```
In [53]:  1  X_train.head(2)
```

Out[53]:

| | id | teacher_id | teacher_prefix | school_state | project_submitted_datetime | project_grade_category | project_resource_summary | teacher_number_of_previously_posted_projects | project_is_approved | clean_categories | ... | words_in_essay | words_in_title | price | quantity | processed_e |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 38452 | p126775 | 8544a1fcc7f73eddd7de295c66b47738 | mrs | ma | 2016-08-25 08:24:00 | 9_12 | My students need items to create a coffee cart... | 0 | 0 | SpecialNeeds | ... | 210 | 2 | 315.01 | 9 | learning give community se other |
| 38286 | p248274 | 2ea294df350d76c01bc4c8aa019fd2bc | mrs | mt | 2016-07-13 13:15:16 | 3_5 | My students need a tinker / maker station!\r\n... | 0 | 0 | Math_Science Literacy_Language | ... | 233 | 6 | 609.19 | 6 | teach awe school prou staff thr |

2 rows × 21 columns

## 3. Vectorization on categorical and text features

### 3.1 Response Encoding on categorical features :

1. Categorical features are  : School state,categories,sub_categories,project_grade,teacher_prefix

```
In [58]:   1  ## craete a response table for the categoical data :
           2  def response_dict(col_name):
           3      col_cat=X_train[col_name].unique()
           4      res_no,res_yes={},{}
           5      for cat in col_cat:
           6          # total number of non approvals (y=0)
           7          no=len(X_train[(X_train[col_name] == cat) & (X_train['project_is_approved'] == 0)])
           8          # total number of non approvals (y=1)
           9          yes=len(X_train[(X_train[col_name] == cat) & (X_train['project_is_approved'] == 1)])
          10          #total
          11          total=len(X_train[X_train[col_name] == cat])
          12          #store the probability value for the category for both class labels
          13          res_no.update({cat :no/total })
          14          res_yes.update({cat :yes/total })
          15      return res_no,res_yes
          16
          17  def response_encoding(data,col_name):
          18      #store class 0 and class 1 probability value for each
          19      resp_no,resp_yes=response_dict(col_name)
          20      data['class_0']=0.5
          21      data['class_1']=0.5
          22      data['class_0']=data[col_name].map(resp_no)
          23      data['class_1']=data[col_name].map(resp_yes)
          24      return data['class_0'].values,data['class_1'].values
```

```python
In [59]:  1  ## vectorized categorical variables for our train data
          2  teacher_prefix_tr_0,teacher_prefix_tr_1=response_encoding(X_train,'teacher_prefix')
          3  proj_grade_tr_0,proj_grade_tr_1=response_encoding(X_train,'project_grade_category')
          4  cat_tr_0,cat_tr_1=response_encoding(X_train,'clean_categories')
          5  subcat_tr_0,subcat_tr_1=response_encoding(X_train,'clean_subcategories')
          6  sklst_tr_0,sklst_tr_1=response_encoding(X_train,'school_state')
          7
          8  ## merge into 2d
          9  #https://stackoverflow.com/questions/8372399/zip-with-list-output-instead-of-tuple
         10  vectorize_tr_trpr=np.array(list(zip(teacher_prefix_tr_0,teacher_prefix_tr_1)))
         11  vectorize_tr_proj=np.array(list(zip(proj_grade_tr_0,proj_grade_tr_1)))
         12  vectorize_tr_cat=np.array(list(zip(cat_tr_0,cat_tr_1)))
         13  vectorize_tr_subcat=np.array(list(zip(subcat_tr_0,subcat_tr_1)))
         14  vectorize_tr_sklst=np.array(list(zip(sklst_tr_0,sklst_tr_1)))
         15
         16  ## print shapes
         17  print(vectorize_tr_trpr.shape)
         18  print(vectorize_tr_proj.shape)
         19  print(vectorize_tr_cat.shape)
         20  print(vectorize_tr_subcat.shape)
         21  print(vectorize_tr_sklst.shape)
```

```
(22445, 2)
(22445, 2)
(22445, 2)
(22445, 2)
(22445, 2)
```

```python
In [60]:  1  ## vectorized categorical variables for our test data
          2  teacher_prefix_te_0,teacher_prefix_te_1=response_encoding(X_test,'teacher_prefix')
          3  proj_grade_te_0,proj_grade_te_1=response_encoding(X_test,'project_grade_category')
          4  cat_te_0,cat_te_1=response_encoding(X_test,'clean_categories')
          5  subcat_te_0,subcat_te_1=response_encoding(X_test,'clean_subcategories')
          6  sklst_te_0,sklst_te_1=response_encoding(X_test,'school_state')
          7
          8  ## merge into 2d
          9  #https://stackoverflow.com/questions/8372399/zip-with-list-output-instead-of-tuple
         10  vectorize_te_trpr=np.array(list(zip(teacher_prefix_te_0,teacher_prefix_te_1)))
         11  vectorize_te_proj=np.array(list(zip(proj_grade_te_0,proj_grade_te_1)))
         12  vectorize_te_cat=np.array(list(zip(cat_te_0,cat_te_1)))
         13  vectorize_te_subcat=np.array(list(zip(subcat_te_0,subcat_te_1)))
         14  vectorize_te_sklst=np.array(list(zip(sklst_te_0,sklst_te_1)))
         15
         16  ## print shapes
         17  print(vectorize_te_trpr.shape)
         18  print(vectorize_te_proj.shape)
         19  print(vectorize_te_cat.shape)
         20  print(vectorize_te_subcat.shape)
         21  print(vectorize_te_sklst.shape)
```

```
(16500, 2)
(16500, 2)
(16500, 2)
(16500, 2)
(16500, 2)
```

```python
In [61]:  1  ## vectorized categorical variables for our cv data
          2  teacher_prefix_cv_0,teacher_prefix_cv_1=response_encoding(X_cv,'teacher_prefix')
          3  proj_grade_cv_0,proj_grade_cv_1=response_encoding(X_cv,'project_grade_category')
          4  cat_cv_0,cat_cv_1=response_encoding(X_cv,'clean_categories')
          5  subcat_cv_0,subcat_cv_1=response_encoding(X_cv,'clean_subcategories')
          6  sklst_cv_0,sklst_cv_1=response_encoding(X_cv,'school_state')
          7
          8  ## merge into 2d
          9  #https://stackoverflow.com/questions/8372399/zip-with-list-output-instead-of-tuple
         10  vectorize_cv_trpr=np.array(list(zip(teacher_prefix_cv_0,teacher_prefix_cv_1)))
         11  vectorize_cv_proj=np.array(list(zip(proj_grade_cv_0,proj_grade_cv_1)))
         12  vectorize_cv_cat=np.array(list(zip(cat_cv_0,cat_cv_1)))
         13  vectorize_cv_subcat=np.array(list(zip(subcat_cv_0,subcat_cv_1)))
         14  vectorize_cv_sklst=np.array(list(zip(sklst_cv_0,sklst_cv_1)))
         15
         16  ## print shapes
         17  print(vectorize_cv_trpr.shape)
         18  print(vectorize_cv_proj.shape)
         19  print(vectorize_cv_cat.shape)
         20  print(vectorize_cv_subcat.shape)
         21  print(vectorize_cv_sklst.shape)
```

```
(11055, 2)
(11055, 2)
(11055, 2)
(11055, 2)
(11055, 2)
```

```python
In [62]:  1  ## drop the y labels from splits
          2  X_train.drop(['project_is_approved'], axis=1, inplace=True)
          3  X_test.drop(['project_is_approved'], axis=1, inplace=True)
          4  X_cv.drop(['project_is_approved'], axis=1, inplace=True)
```

### 3.2 Vectorizing Text data

#### 3.2.1 BOW on Essay data

```python
In [63]:  1  ##Considering the words that appeared in atleast 10 documents
          2
          3  bow_essay = CountVectorizer(min_df=10,max_features=5000)
          4  bow_essay.fit(X_train['processed_essay'])
          5
          6  bow_essay_train = bow_essay.transform(X_train['processed_essay'])
          7
          8  print("Shape of matrix after one hot encoding ",bow_essay_train.shape)
          9
         10  ## tranform Test data
         11
         12  bow_essay_test = bow_essay.transform(X_test['processed_essay'])
         13
         14  print("Shape of matrix after one hot encoding ",bow_essay_test.shape)
         15
         16
         17  ## Teansform cv data
         18
         19  bow_essay_cv = bow_essay.transform(X_cv['processed_essay'])
         20  print("Shape of matrix after one hot encoding ",bow_essay_cv.shape)
```

```
Shape of matrix after one hot encoding  (22445, 5000)
Shape of matrix after one hot encoding  (16500, 5000)
Shape of matrix after one hot encoding  (11055, 5000)
```

#### 3.2.2 BOW on Title data

```python
In [64]:  1  ##Considering the words that appeared in atleast 10 documents
          2
          3  bow_title = CountVectorizer(min_df=10,max_features=5000)
          4  bow_title.fit(X_train['processed_title'])
          5
          6  bow_title_train = bow_title.transform(X_train['processed_title'])
          7
          8  print("Shape of matrix after one hot encoding ",bow_title_train.shape)
          9
         10  ## tranform Test data
         11
         12  bow_title_test = bow_title.transform(X_test['processed_title'])
         13
         14  print("Shape of matrix after one hot encoding ",bow_title_test.shape)
         15
         16
         17  ## Teansform cv data
         18
         19  bow_title_cv = bow_title.transform(X_cv['processed_title'])
         20  print("Shape of matrix after one hot encoding ",bow_title_cv.shape)
```

```
Shape of matrix after one hot encoding  (22445, 1146)
Shape of matrix after one hot encoding  (16500, 1146)
Shape of matrix after one hot encoding  (11055, 1146)
```

#### 3.2.3 TFIDF on Essay data

In [65]:
```python
##Considering the words that appeared in atleast 10 documents

tfidf_essay = TfidfVectorizer(min_df=10,max_features=5000)
tfidf_essay.fit(X_train['processed_essay'])

tfidf_essay_train = tfidf_essay.transform(X_train['processed_essay'])

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

## tranform Test data

tfidf_essay_test = tfidf_essay.transform(X_test['processed_essay'])

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


## Teansform cv data

tfidf_essay_cv = tfidf_essay.transform(X_cv['processed_essay'])
print("Shape of matrix after one hot encoding ",tfidf_essay_cv.shape)
```

```
Shape of matrix after one hot encoding  (22445, 5000)
Shape of matrix after one hot encoding  (16500, 5000)
Shape of matrix after one hot encoding  (11055, 5000)
```

### 3.2.4 TFIDF on Title data

In [66]:
```python
##Considering the words that appeared in atleast 10 documents

tfidf_title = TfidfVectorizer(min_df=10,max_features=5000)
tfidf_title.fit(X_train['processed_title'])

tfidf_title_train = tfidf_title.transform(X_train['processed_title'])

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

## tranform Test data

tfidf_title_test = tfidf_title.transform(X_test['processed_title'])

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


## Teansform cv data

tfidf_title_cv = tfidf_title.transform(X_cv['processed_title'])
print("Shape of matrix after one hot encoding ",tfidf_title_cv.shape)
```

```
Shape of matrix after one hot encoding  (22445, 1146)
Shape of matrix after one hot encoding  (16500, 1146)
Shape of matrix after one hot encoding  (11055, 1146)
```

### 3.2.5 Weighted tfidf on Essay data using Pretrained Models

In [67]:
```python
# 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
## Glove vectors are global vectors for words which has vector every word in 300d .
## for read more :https://nlp.stanford.edu/projects/glove/
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words =  set(model.keys())
```

In [68]:
```python
# average Word2Vec on train
# compute average word2vec for each review.
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]

tfidf_model_essay = TfidfVectorizer()
tfidf_model_essay.fit(X_train["processed_essay"])

# we are converting a essay_dictionary with word as a key, and the idf as a value
essay_dictionary = dict(zip(tfidf_model_essay.get_feature_names(), list(tfidf_model_essay.idf_)))
tfidf_words_essay = set(tfidf_model_essay.get_feature_names())



def tfidf_w2v_vectors(data,glove_words,essay_dictionary,tfidf_words_essay):
    tf_vector=[]
    for sentence in data: # 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_essay):
                vec = model[word] # getting the vector for each word
                # here we are multiplying idf value(essay_dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
                tf_idf = essay_dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf
        if tf_idf_weight != 0:
            vector /= tf_idf_weight
        tf_vector.append(vector)
    print(len(tf_vector))
    print(len(tf_vector[0]))
    return tf_vector

tfidf_w2v_vectors_train=tfidf_w2v_vectors(X_train['processed_essay'],glove_words,essay_dictionary,tfidf_words_essay)
tfidf_w2v_vectors_cv=tfidf_w2v_vectors(X_cv['processed_essay'],glove_words,essay_dictionary,tfidf_words_essay)
tfidf_w2v_vectors_test=tfidf_w2v_vectors(X_test['processed_essay'],glove_words,essay_dictionary,tfidf_words_essay)
```

```
22445
300
11055
300
16500
300
```

### 3.2.6 Weighted tfidf on Title data using Pretrained Models

In [69]:
```python
# average Word2Vec on train
# compute average word2vec for each review.
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]

tfidf_model_title = TfidfVectorizer()
tfidf_model_title.fit(X_train["processed_title"])
# we are converting a title_dictionary with word as a key, and the idf as a value
title_dictionary = dict(zip(tfidf_model_title.get_feature_names(), list(tfidf_model_title.idf_)))
tfidf_words_title = set(tfidf_model_title.get_feature_names())

tfidf_w2v_vectors_title_train=tfidf_w2v_vectors(X_train['processed_title'],glove_words,title_dictionary,tfidf_words_title)
tfidf_w2v_vectors_title_cv=tfidf_w2v_vectors(X_cv['processed_title'],glove_words,title_dictionary,tfidf_words_title)
tfidf_w2v_vectors_title_test=tfidf_w2v_vectors(X_test['processed_title'],glove_words,title_dictionary,tfidf_words_title)
```

```
22445
300
11055
300
16500
300
```

### 3.2.7 Avg W2V on Essay data using Pretrained Models

```
In [70]:  ▶  1  # average Word2Vec on train
             2  # compute average word2vec for each review.
             3
             4
             5  def avg_w2v_vector(data,glove_words):
             6      tf_vector=[]
             7      for sentence in tqdm(data): # for each review/sentence
             8          vector = np.zeros(300) # as word vectors are of zero length
             9          cnt_words =0; # num of words with a valid vector in the sentence/review
            10          for word in sentence.split(): # for each word in a review/sentence
            11              if word in glove_words:
            12                  vector += model[word]
            13                  cnt_words += 1
            14          if cnt_words != 0:
            15              vector /= cnt_words
            16          tf_vector.append(vector)
            17
            18      print(len(tf_vector))
            19      print(len(tf_vector[0]))
            20      return tf_vector
            21
            22  # average Word2Vec on CV
            23  # compute average word2vec for each review.
            24  avg_w2v_vectors_train = avg_w2v_vector(X_train['processed_essay'],glove_words)
            25  avg_w2v_vectors_cv = avg_w2v_vector(X_cv['processed_essay'],glove_words)
            26  avg_w2v_vectors_test = avg_w2v_vector(X_test['processed_essay'],glove_words)
            27
```

```
100%|██████████| 22445/22445 [00:11<00:00, 1922.51it/s]
  2%|▎         | 258/11055 [00:00<00:04, 2561.30it/s]

22445
300

100%|██████████| 11055/11055 [00:07<00:00, 1432.30it/s]
  2%|▎         | 364/16500 [00:00<00:09, 1692.43it/s]

11055
300

100%|██████████| 16500/16500 [00:07<00:00, 2249.37it/s]

16500
300
```

***3.2.8 Avg W2V on Title data using Pretrained Models***

```
In [71]:  ▶  1  # average Word2Vec on CV
             2  # compute average word2vec for each review.
             3  avg_w2v_vectors_title_train = avg_w2v_vector(X_train['processed_title'],glove_words)
             4  avg_w2v_vectors_title_cv = avg_w2v_vector(X_cv['processed_title'],glove_words)
             5  avg_w2v_vectors_title_test = avg_w2v_vector(X_test['processed_title'],glove_words)
             6
```

```
100%|██████████| 22445/22445 [00:01<00:00, 17903.75it/s]
 20%|██        | 2157/11055 [00:00<00:00, 21413.49it/s]

22445
300

100%|██████████| 11055/11055 [00:00<00:00, 16922.96it/s]
  0%|          | 0/16500 [00:00<?, ?it/s]

11055
300

100%|██████████| 16500/16500 [00:00<00:00, 19486.67it/s]

16500
300
```

## 4. Vectorizing Numerical Features

### 4.1 Price

```
In [72]:  ▶  1  normalizer = Normalizer()
             2  # normalizer.fit(X_train['price'].values)
             3  # this will rise an error Expected 2D array, got 1D array instead:
             4  # array=[105.22 215.96  96.01 ... 368.98   80.53 709.67].
             5  # Reshape your data either using
             6  # array.reshape(-1, 1) if your data has a single feature
             7  # array.reshape(1, -1)  if it contains a single sample.
             8  normalizer.fit(X_train['price'].values.reshape(1,-1))
             9
            10  X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1))
            11  X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(1,-1))
            12  X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))
            13
            14  print("After vectorizations")
            15  print(X_train_price_norm.shape, y_train.shape)
            16  print(X_cv_price_norm.shape, y_cv.shape)
            17  print(X_test_price_norm.shape, y_test.shape)
            18  print("="*100)
            19
            20  ## reshaping
            21  X_train_price_norm=X_train_price_norm.reshape(-1,1)
            22  X_cv_price_norm=X_cv_price_norm.reshape(-1,1)
            23  X_test_price_norm=X_test_price_norm.reshape(-1,1)
```

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

### 4.2 Quantity

```
In [73]:  ▶  1  normalizer = Normalizer()
             2
             3  # normalizer.fit(X_train['price'].values)
             4  # this will rise an error Expected 2D array, got 1D array instead:
             5  # array=[105.22 215.96  96.01 ... 368.98   80.53 709.67].
             6  # Reshape your data either using
             7  # array.reshape(-1, 1) if your data has a single feature
             8  # array.reshape(1, -1)  if it contains a single sample.
             9
            10  normalizer.fit(X_train['quantity'].values.reshape(1,-1))
            11
            12  quantity_train_norm = normalizer.transform(X_train['quantity'].values.reshape(1,-1))
            13  quantity_cv_norm = normalizer.transform(X_cv['quantity'].values.reshape(1,-1))
            14  quantity_test_norm = normalizer.transform(X_test['quantity'].values.reshape(1,-1))
            15
            16  print("After vectorizations")
            17  print(quantity_train_norm.shape, y_train.shape)
            18  print(quantity_cv_norm.shape, y_cv.shape)
            19  print(quantity_test_norm.shape, y_test.shape)
            20  print("="*100)
            21
            22  ## reshaping
            23  quantity_train_norm=quantity_train_norm.reshape(-1,1)
            24  quantity_cv_norm=quantity_cv_norm.reshape(-1,1)
            25  quantity_test_norm=quantity_test_norm.reshape(-1,1)
```

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

### 4.3 Number of Previously posted projects

```
In [74]:    1  normalizer = Normalizer()
            2
            3  # normalizer.fit(X_train['price'].values)
            4  # this will rise an error Expected 2D array, got 1D array instead:
            5  # array=[105.22 215.96  96.01 ... 368.98  80.53 709.67].
            6  # Reshape your data either using
            7  # array.reshape(-1, 1) if your data has a single feature
            8  # array.reshape(1, -1)  if it contains a single sample.
            9
           10  normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
           11
           12  prev_projects_train_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
           13  prev_projects_cv_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
           14  prev_projects_test_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
           15
           16  print("After vectorizations")
           17  print(prev_projects_train_norm.shape, y_train.shape)
           18  print(prev_projects_cv_norm.shape, y_cv.shape)
           19  print(prev_projects_test_norm.shape, y_test.shape)
           20  print("="*100)
           21
           22  ## reshaping
           23  prev_projects_train_norm=prev_projects_train_norm.reshape(-1,1)
           24  prev_projects_cv_norm=prev_projects_cv_norm.reshape(-1,1)
           25  prev_projects_test_norm=prev_projects_test_norm.reshape(-1,1)
```

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

## 4.4 Title Word counts

```
In [75]:    1  normalizer = Normalizer()
            2
            3  normalizer.fit(X_train['words_in_title'].values.reshape(1,-1))
            4
            5  title_word_count_train_norm = normalizer.transform(X_train['words_in_title'].values.reshape(1,-1))
            6  title_word_count_cv_norm = normalizer.transform(X_cv['words_in_title'].values.reshape(1,-1))
            7  title_word_count_test_norm = normalizer.transform(X_test['words_in_title'].values.reshape(1,-1))
            8
            9  print("After vectorizations")
           10  print(title_word_count_train_norm.shape, y_train.shape)
           11  print(title_word_count_cv_norm.shape, y_cv.shape)
           12  print(title_word_count_test_norm.shape, y_test.shape)
           13  print("="*100)
           14
           15  ## reshaping
           16  title_word_count_train_norm=title_word_count_train_norm.reshape(-1,1)
           17  title_word_count_cv_norm=title_word_count_cv_norm.reshape(-1,1)
           18  title_word_count_test_norm=title_word_count_test_norm.reshape(-1,1)
```

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

## 4.5 Essay Words Counts

```
In [76]:    1  normalizer = Normalizer()
            2
            3  normalizer.fit(X_train['words_in_essay'].values.reshape(1,-1))
            4
            5  essay_word_count_train_norm = normalizer.transform(X_train['words_in_essay'].values.reshape(1,-1))
            6  essay_word_count_cv_norm = normalizer.transform(X_cv['words_in_essay'].values.reshape(1,-1))
            7  essay_word_count_test_norm = normalizer.transform(X_test['words_in_essay'].values.reshape(1,-1))
            8
            9  print("After vectorizations")
           10  print(essay_word_count_train_norm.shape, y_train.shape)
           11  print(essay_word_count_cv_norm.shape, y_cv.shape)
           12  print(essay_word_count_test_norm.shape, y_test.shape)
           13
           14  ## reshaping
           15  essay_word_count_train_norm=essay_word_count_train_norm.reshape(-1,1)
           16  essay_word_count_cv_norm=essay_word_count_cv_norm.reshape(-1,1)
           17  essay_word_count_test_norm=essay_word_count_test_norm.reshape(-1,1)
```

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

## 4.6 Vectorizing sentiment Columns

```
In [81]:    1  ### vectorize pos
            2  Normalize_pos=Normalizer()
            3  Normalize_pos.fit(X_train['pos'].values.reshape(1,-1))
            4  sentiment_pos_train_norm=Normalize_pos.transform(X_train['pos'].values.reshape(1,-1))
            5  sentiment_pos_test_norm=Normalize_pos.transform(X_test['pos'].values.reshape(1,-1))
            6  sentiment_pos_cv_norm=Normalize_pos.transform(X_cv['pos'].values.reshape(1,-1))
            7  sentiment_pos_train_norm=sentiment_pos_train_norm.reshape(-1,1)
            8  sentiment_pos_test_norm=sentiment_pos_test_norm.reshape(-1,1)
            9  sentiment_pos_cv_norm=sentiment_pos_cv_norm.reshape(-1,1)
```

```
In [82]:    1  ### vectorize neg
            2  Normalize_neg=Normalizer()
            3  Normalize_neg.fit(X_train['neg'].values.reshape(1,-1))
            4  sentiment_neg_train_norm=Normalize_neg.transform(X_train['neg'].values.reshape(1,-1))
            5  sentiment_neg_test_norm=Normalize_neg.transform(X_test['neg'].values.reshape(1,-1))
            6  sentiment_neg_cv_norm=Normalize_neg.transform(X_cv['neg'].values.reshape(1,-1))
            7  sentiment_neg_train_norm=sentiment_neg_train_norm.reshape(-1,1)
            8  sentiment_neg_test_norm=sentiment_neg_test_norm.reshape(-1,1)
            9  sentiment_neg_cv_norm=sentiment_neg_cv_norm.reshape(-1,1)
```

```
In [83]:    1  ### vectorize compound
            2  Normalize_co=Normalizer()
            3  Normalize_co.fit(X_train['compound'].values.reshape(1,-1))
            4  sentiment_compound_train_norm=Normalize_co.transform(X_train['compound'].values.reshape(1,-1))
            5  sentiment_compound_test_norm=Normalize_co.transform(X_test['compound'].values.reshape(1,-1))
            6  sentiment_compound_cv_norm=Normalize_co.transform(X_cv['compound'].values.reshape(1,-1))
            7  sentiment_compound_train_norm=sentiment_compound_train_norm.reshape(-1,1)
            8  sentiment_compound_test_norm=sentiment_compound_test_norm.reshape(-1,1)
            9  sentiment_compound_cv_norm=sentiment_compound_cv_norm.reshape(-1,1)
```

```
In [84]:    1  ### vectorize neu
            2  Normalize_neu=Normalizer()
            3  Normalize_neu.fit(X_train['neu'].values.reshape(1,-1))
            4  sentiment_neu_train_norm=Normalize_neu.transform(X_train['neu'].values.reshape(1,-1))
            5  sentiment_neu_test_norm=Normalize_neu.transform(X_test['neu'].values.reshape(1,-1))
            6  sentiment_neu_cv_norm=Normalize_neu.transform(X_cv['neu'].values.reshape(1,-1))
            7  sentiment_neu_train_norm=sentiment_neu_train_norm.reshape(-1,1)
            8  sentiment_neu_test_norm=sentiment_neu_test_norm.reshape(-1,1)
            9  sentiment_neu_cv_norm=sentiment_neu_cv_norm.reshape(-1,1)
```

# Assignment 9: RF

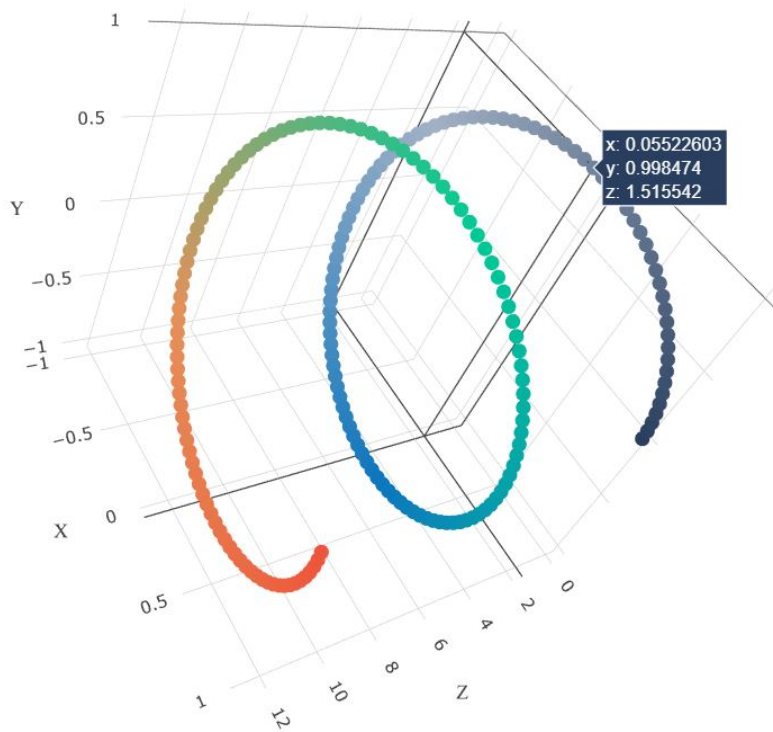1. Apply both Random Forrest and GBDT on these feature sets

- Set 1: categorical(instead of one hot encoding, try response coding (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(BOW) + preprocessed_eassay (BOW)
- Set 2: categorical(instead of one hot encoding, try response coding (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)
- Set 3: categorical(instead of one hot encoding, try response coding (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(AVG W2V)+ preprocessed_eassay (AVG W2V). Here for this set take **20K** datapoints only.
- Set 4: categorical(instead of one hot encoding, try response coding (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V). Here for this set take **20K** datapoints only.

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

- Consider the following range for hyperparameters **n_estimators** = [10, 50, 100, 150, 200, 300, 500, 1000], **max_depth** = [2, 3, 4, 5, 6, 7, 8, 9, 10]
- Find the best hyper parameter which will give the maximum AUC (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/) value
- Find the best hyper parameter using simple cross validation data
- You can write your own for loops to do this task
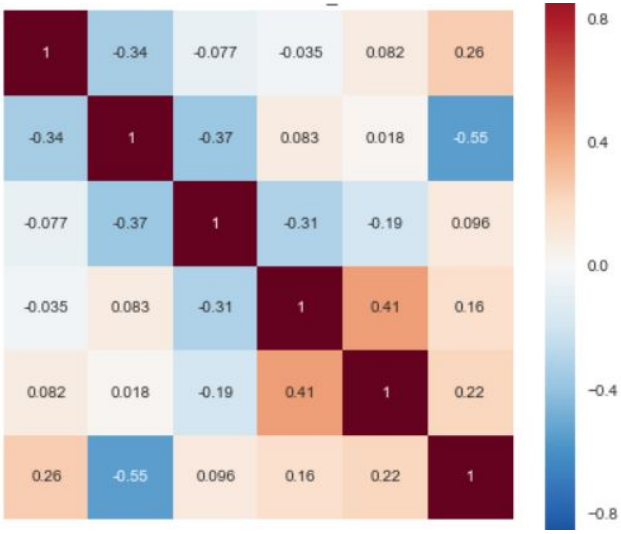
3. **Representation of results**

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

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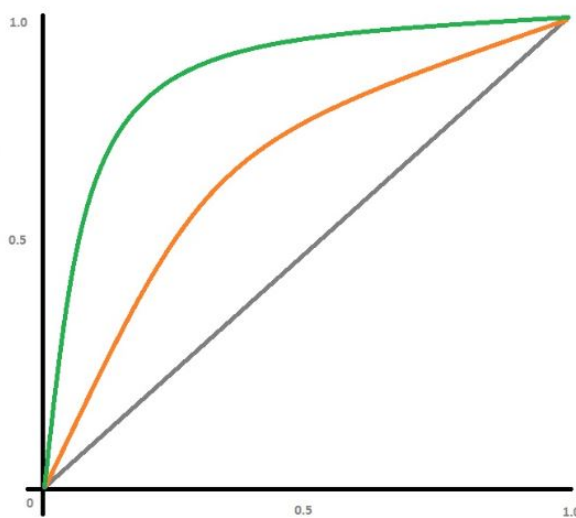
## or

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



seaborn heat maps (https://seaborn.pydata.org/generated/seaborn.heatmap.html) with rows as **n_estimators**, columns as **max_depth**, and values inside the cell representing **AUC Score**
- You can choose either of the plotting techniques: 3d plot or heat map
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.



- Along with plotting ROC curve, you need to print the confusion matrix (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/confusion-matrix-tpr-fpr-fnr-tnr-1/) with predicted and original labels of test data points

|  | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | TN = ?? | FP = ?? |
| Actual: YES | FN = ?? | TP = ?? |

4. **Conclusion**

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link (http://zetcode.com/python/prettytable/)

```
+-----------+--------+-----------------+------+
| Vectorizer | Model | Hyper parameter | AUC |
+-----------+--------+-----------------+------+
|    BOW    | Brute |        7        | 0.78 |
+-----------+--------+-----------------+------+
|   TFIDF   | Brute |        12       | 0.79 |
+-----------+--------+-----------------+------+
|    W2V    | Brute |        10       | 0.78 |
+-----------+--------+-----------------+------+
|  TFIDFW2V | Brute |        6        | 0.78 |
+-----------+--------+-----------------+------+
```

## 5. 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 instrucations

### 5.1 Applying Random Forests on BOW, SET 1

```
In [85]:
1  # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2
3  X_tr = hstack((vectorize_tr_trpr, vectorize_tr_proj,vectorize_tr_cat,vectorize_tr_subcat,
4                 vectorize_tr_sklst,bow_essay_train,bow_title_train,sentiment_pos_train_norm,
5                 sentiment_neg_train_norm,sentiment_compound_train_norm,sentiment_neu_train_norm,
6                 X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
7                 essay_word_count_train_norm)).tocsr()
8
9  X_te = hstack((vectorize_te_trpr, vectorize_te_proj,vectorize_te_cat,vectorize_te_subcat,
10                vectorize_te_sklst,bow_essay_test,bow_title_test,sentiment_pos_test_norm,
11                sentiment_neg_test_norm,sentiment_compound_test_norm,sentiment_neu_test_norm,
12                X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
13                essay_word_count_test_norm)).tocsr()
14
15 X_cr = hstack((vectorize_cv_trpr, vectorize_cv_proj,vectorize_cv_cat,vectorize_cv_subcat,
16                vectorize_cv_sklst,bow_essay_cv,bow_title_cv,sentiment_pos_cv_norm,
17                sentiment_neg_cv_norm,sentiment_compound_cv_norm,sentiment_neu_cv_norm,
18                X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
19                essay_word_count_cv_norm)).tocsr()
20
21
22 print(X_tr.shape)
23 print(X_te.shape)
24 print(X_cr.shape)
```

```
(22445, 6165)
(16500, 6165)
(11055, 6165)
```

```
In [86]:
1  print("Final Data matrix")
2  print(X_tr.shape, y_train.shape)
3  print(X_cr.shape, y_cv.shape)
4  print(X_te.shape, y_test.shape)
5  print("="*100)
```

```
Final Data matrix
(22445, 6165) (22445,)
(11055, 6165) (11055,)
(16500, 6165) (16500,)
====================================================================================================
```

### 5.1.1 Write loop to find best hyperparameters and do simple cross validation
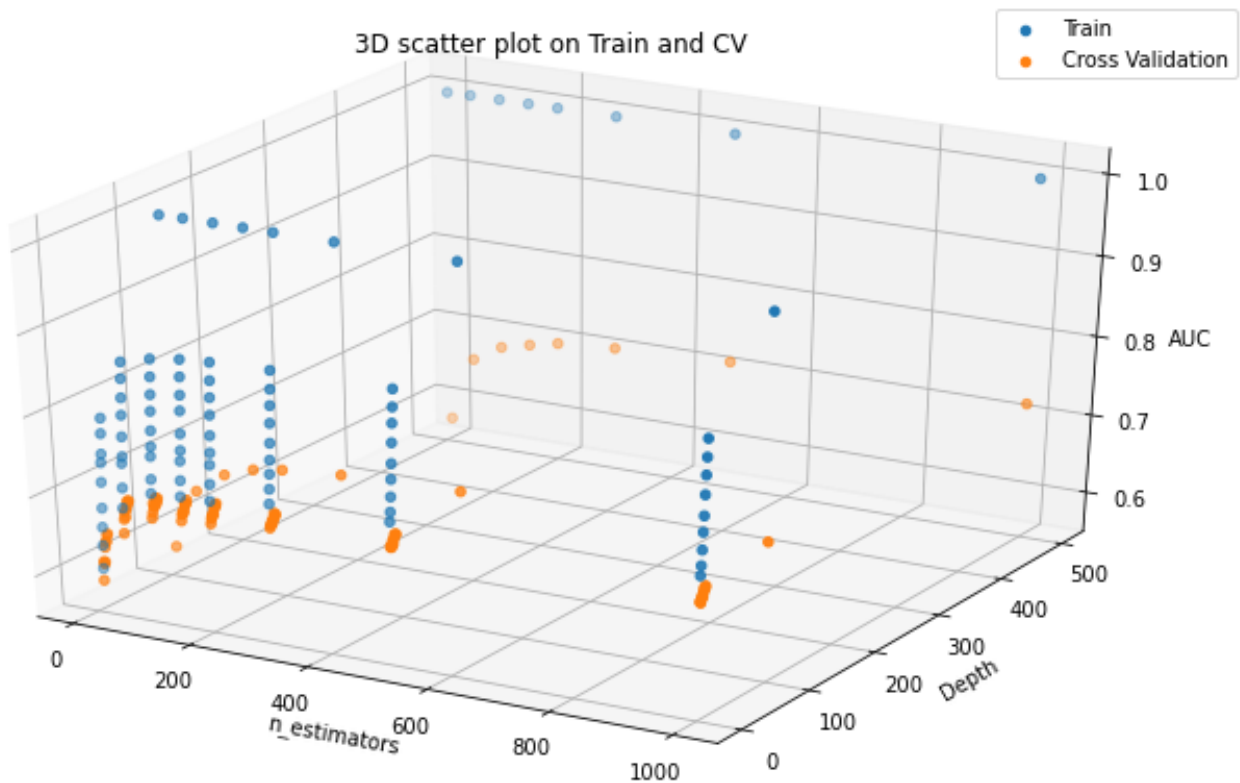
```
In [87]:  ▶   1  ## Iterative looop to find best parameter
              2  n_estimators = [10, 50, 100, 150, 200, 300, 500, 1000]
              3  max_depth = [2, 3, 4, 5, 6, 7, 8, 9, 10,100,500]
              4
              5  def hyperparam_auc(max_depth,n_estimators,X_tr,X_cr,y_train,y_cv):
              6      train_auc=[]
              7      cv_auc=[]
              8      es,de=[],[]
              9      for estimators in tqdm(n_estimators):
             10          for depth in max_depth:
             11              RF=RandomForestClassifier(n_estimators=estimators,max_depth=depth,
             12                                        class_weight='balanced',
             13                                        random_state=4)
             14              RF.fit(X_tr,y_train)
             15              #predict train and cv
             16              train_predicts=RF.predict_proba(X_tr)[:,1]
             17              cv_predicts=RF.predict_proba(X_cr)[:,1]
             18              #Store train and cv auc score in dict
             19              train_auc.append(roc_auc_score(y_train,train_predicts))
             20              cv_auc.append(roc_auc_score(y_cv,cv_predicts))
             21              es.append(estimators)
             22              de.append(depth)
             23      return train_auc,cv_auc,es,de
             24
             25  train_auc,cv_auc,es,de=hyperparam_auc(max_depth,n_estimators,X_tr,X_cr,y_train,y_cv)
```

```
100%|████████| 8/8 [1:47:17<00:00, 804.64s/it]
```

## 5.1.2 Representation using 3D scatter plot

```
In [88]:  ▶   1  from mpl_toolkits.mplot3d import Axes3D
              2  plt.figure(figsize=(12,7))
              3  ax = plt.axes(projection='3d')
              4  zipped_=list(map(list,zip(es,de)))
              5  x1=[i[0] for i in zipped_]
              6  y1=[i[1] for i in zipped_]
              7  # Data for three-dimensional scattered points
              8  # reference : https://jakevdp.github.io/PythonDataScienceHandbook/04.12-three-dimensional-plotting.html
              9  ax.scatter3D(x1, y1, train_auc,label='Train')
             10  ax.scatter3D(x1, y1, cv_auc,label='Cross Validation')
             11  ax.set_xlabel('n_estimators')
             12  ax.set_ylabel('Depth')
             13  ax.set_zlabel('AUC')
             14  ax.set_title('3D scatter plot on Train and CV')
             15  ax.legend()
```

Out[88]: <matplotlib.legend.Legend at 0x227d7cbb0b8>



## 5.1.3 Finding best parameter for training the model

```
In [89]:  ▶   1  ## finding the best pair of hyperparameter with max AUC using a function
              2
              3  def find_best_hyperparam(x1,y1,cv_auc):
              4      zipped_=list(zip(x1,y1,cv_auc))
              5      max_auc=max(cv_auc)
              6      print('Max_auc is',max_auc)
              7      for i in zipped_:
              8          if i[2] == max_auc:
              9              best_params={'n_estimators':i[0],'depth':i[1]}
             10          else:
             11              pass
             12      return best_params
```
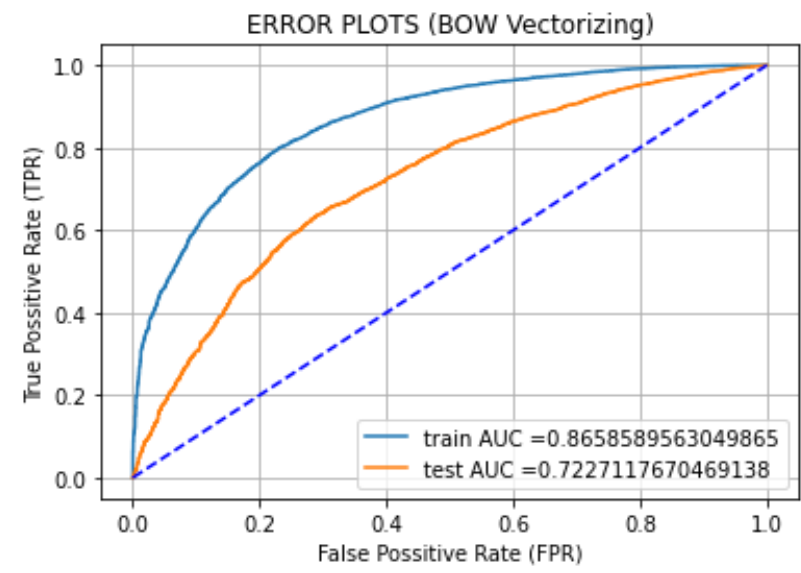
```
In [93]:  ▶   1  best_params=find_best_hyperparam(x1,y1,cv_auc)
              2  best_params
```

```
Max_auc is 0.7247584111051009
```

Out[93]: {'n_estimators': 1000, 'depth': 100}

```
In [99]:  ▶   1  RF=RandomForestClassifier(class_weight='balanced',n_estimators=best_params['n_estimators'],
              2                            max_depth=best_params['depth'],min_samples_split=1000)
              3  RF.fit(X_tr,y_train)
              4
              5  ## Predict the test
              6  train_predicts=RF.predict_proba(X_tr)[:,1]
              7  test_predicts=RF.predict_proba(X_te)[:,1]
              8
              9  ## Store fpr and tpr rates
             10
             11  train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, train_predicts)
             12  test_fpr, test_tpr, te_thresholds = roc_curve(y_test, test_predicts)
```

```
In [100]:  ▶   1  #plot
               2  plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
               3  plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
               4  plt.legend()
               5  plt.xlabel("False Positive Rate (FPR)")
               6  plt.ylabel("True Positive Rate (TPR)")
               7  plt.title("ERROR PLOTS (BOW Vectorizing)")
               8  plt.plot([0, 1], [0, 1],'b--')
               9  plt.grid()
              10  plt.show()
              11
              12
              13  ## Store auc results in a dictionary
              14  results_dict={'BOW' : {'trainauc':str(auc(train_fpr, train_tpr)),
              15                         'testauc': str(auc(test_fpr, test_tpr)),
              16                         'max_depth' : best_params['depth'],
              17                         'n_estimators' :best_params['n_estimators']} }
```



## 5.1.4 Confusion Matrix on train and test data

```
In [102]:  ▶   1  ## Finding best threshold for predictions
               2  def best_threshold(thresholds,fpr,tpr):
               3      t=thresholds[np.argmax(tpr*(1-fpr))]
               4      # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
               5      print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
               6      return t
               7
               8  def predict_with_best_t(proba, threshold):
               9      predictions = []
              10      for i in proba:
              11          if i>=threshold:
              12              predictions.append(1)
              13          else:
              14              predictions.append(0)
              15      return predictions
```

```
In [104]: ▶    1  print("="*100)
                2  from sklearn.metrics import confusion_matrix
                3  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
                4  print("Train confusion matrix")
                5  print(confusion_matrix(y_train, predict_with_best_t(train_predicts, best_t)))
                6  print("Test confusion matrix")
                7  print(confusion_matrix(y_test, predict_with_best_t(test_predicts, best_t)))
```

```
==================================================================================================
the maximum value of tpr*(1-fpr) 0.6152378147120736 for threshold 0.513
Train confusion matrix
[[ 2676   742]
 [ 4075 14952]]
Test confusion matrix
[[ 1321  1191]
 [ 3010 10978]]
```
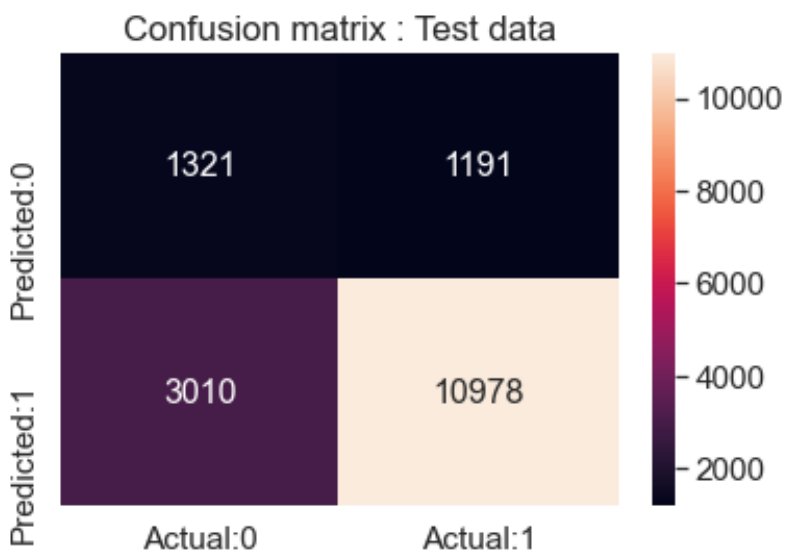
```
In [105]: ▶    1  ### PLOT the matrix for Train
                2  #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
                3  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
                4  df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(train_predicts, best_t))
                5                       , range(2), range(2))
                6  # plt.figure(figsize=(10,7))
                7  sns.set(font_scale=1.4) # for label size
                8  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',
                9              xticklabels=['Actual:0','Actual:1']
               10              ,yticklabels=['Predicted:0','Predicted:1']) # font size
               11  plt.title('Confusion matrix : Train data')
               12  plt.show()
```



Confusion matrix : Train data

```
In [106]: ▶    1  ### PLOT the matrix for Train
                2  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
                3  df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(test_predicts, best_t)), range(2), range(2))
                4  # plt.figure(figsize=(10,7))
                5  sns.set(font_scale=1.4) # for label size
                6  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',xticklabels=['Actual:0','Actual:1']
                7              ,yticklabels=['Predicted:0','Predicted:1']) # font size
                8  plt.title('Confusion matrix : Test data')
                9  plt.show()
```



Confusion matrix : Test data

### 5.2 Applying Random Forests on TFIDF, SET 2

```
In [107]: ▶    1  # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
                2
                3  X_tr = hstack((vectorize_tr_trpr, vectorize_tr_proj,vectorize_tr_cat,vectorize_tr_subcat,
                4                vectorize_tr_sklst,tfidf_essay_train,tfidf_title_train,sentiment_pos_train_norm,
                5                sentiment_neg_train_norm,sentiment_compound_train_norm,sentiment_neu_train_norm,
                6                X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
                7                essay_word_count_train_norm)).tocsr()
                8
                9  X_te = hstack((vectorize_te_trpr, vectorize_te_proj,vectorize_te_cat,vectorize_te_subcat,
               10                vectorize_te_sklst,tfidf_essay_test,tfidf_title_test,sentiment_pos_test_norm,
               11                sentiment_neg_test_norm,sentiment_compound_test_norm,sentiment_neu_test_norm,
               12                X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
               13                essay_word_count_test_norm)).tocsr()
               14
               15  X_cr = hstack((vectorize_cv_trpr, vectorize_cv_proj,vectorize_cv_cat,vectorize_cv_subcat,
               16                vectorize_cv_sklst,tfidf_essay_cv,tfidf_title_cv,sentiment_pos_cv_norm,
               17                sentiment_neg_cv_norm,sentiment_compound_cv_norm,sentiment_neu_cv_norm,
               18                X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
               19                essay_word_count_cv_norm)).tocsr()
               20
               21
               22  print(X_tr.shape)
               23  print(X_te.shape)
               24  print(X_cr.shape)
```

```
(22445, 6165)
(16500, 6165)
(11055, 6165)
```

```
In [108]: ▶    1  print("Final Data matrix")
                2  print(X_tr.shape, y_train.shape)
                3  print(X_cr.shape, y_cv.shape)
                4  print(X_te.shape, y_test.shape)
                5  print("="*100)
```

```
Final Data matrix
(22445, 6165) (22445,)
(11055, 6165) (11055,)
(16500, 6165) (16500,)
==================================================================================================
```

### 5.2.1 Write loop to find best hyperparameters and do simple cross validation

```
In [109]: ▶    1  ## Iterative looop to find best parameter
                2  n_estimators = [10, 50, 100, 150, 200, 300, 500, 1000]
                3  max_depth = [2, 3, 4, 5, 6, 7, 8, 9, 10,100,500]
                4
                5
                6  train_auc,cv_auc,es,de=hyperparam_auc(max_depth,n_estimators,X_tr,X_cr,y_train,y_cv)
```
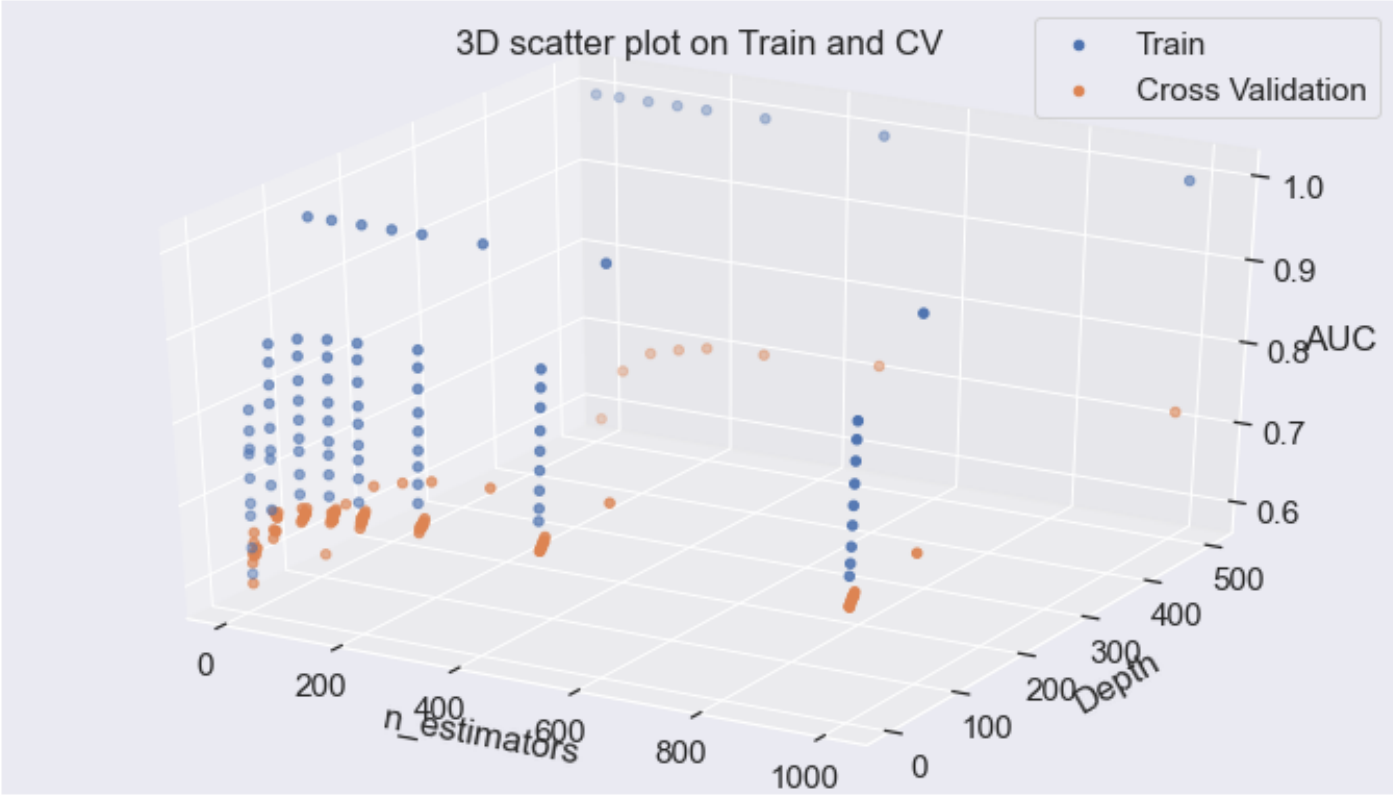
```
100%|██████████| 8/8 [2:15:51<00:00, 1018.92s/it]
```

### 5.2.2 Representation using 3D scatter plot

```python
In [115]:  1  #from mpl_toolkits.mplot3d import Axes3D
           2  plt.figure(figsize=(12,7))
           3  ax = plt.axes(projection='3d')
           4  zipped_=list(map(list,zip(es,de)))
           5  x1=[i[0] for i in zipped_]
           6  y1=[i[1] for i in zipped_]
           7  # Data for three-dimensional scattered points
           8  # reference : https://jakevdp.github.io/PythonDataScienceHandbook/04.12-three-dimensional-plotting.html
           9  ax.scatter3D(x1, y1, train_auc,label='Train')
          10  ax.scatter3D(x1, y1, cv_auc,label='Cross Validation')
          11  ax.set_xlabel('n_estimators')
          12  ax.set_ylabel('Depth')
          13  ax.set_zlabel('AUC')
          14  ax.set_title('3D scatter plot on Train and CV')
          15  ax.legend()
```

Out[115]: <matplotlib.legend.Legend at 0x227d7c3a3c8>
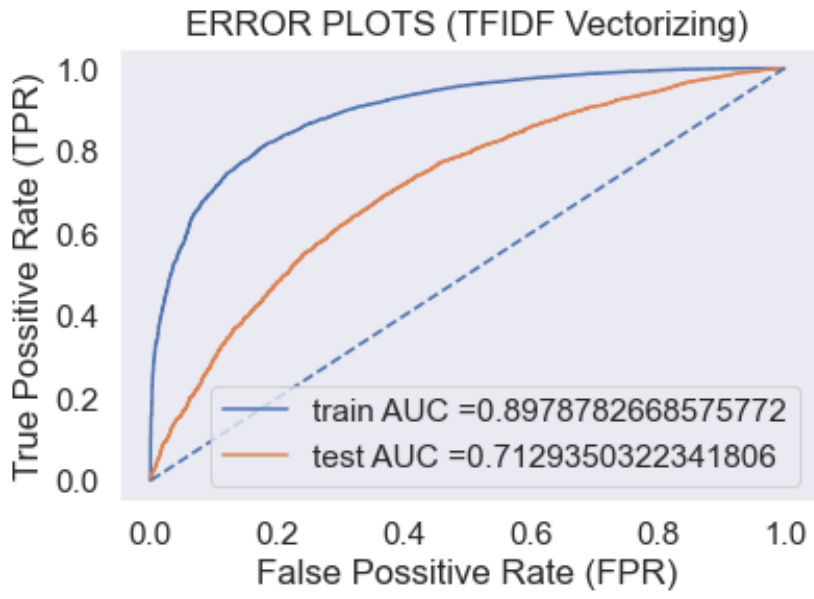


### 5.2.3 Finding best parameter for training the model

```python
In [116]:  1  best_params=find_best_hyperparam(x1,y1,cv_auc)
           2  best_params
```

Max_auc is 0.7205517807686973

Out[116]: {'n_estimators': 1000, 'depth': 100}

```python
In [117]:  1  RF=RandomForestClassifier(class_weight='balanced',n_estimators=best_params['n_estimators'],
           2                            max_depth=best_params['depth'],random_state=4,min_samples_split=1000)
           3  RF.fit(X_tr,y_train)
           4
           5  ## Predict the test
           6  train_predicts=RF.predict_proba(X_tr)[:,1]
           7  test_predicts=RF.predict_proba(X_te)[:,1]
           8
           9  ## Store fpr and tpr rates
          10
          11  train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, train_predicts)
          12  test_fpr, test_tpr, te_thresholds = roc_curve(y_test, test_predicts)
```

```python
In [118]:  1  #plot
           2  plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
           3  plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
           4  plt.legend()
           5  plt.xlabel("False Positive Rate (FPR)")
           6  plt.ylabel("True Positive Rate (TPR)")
           7  plt.title("ERROR PLOTS (TFIDF Vectorizing)")
           8  plt.plot([0, 1], [0, 1],'b--')
           9  plt.grid()
          10  plt.show()
          11
          12
          13  ## Store auc results in a dictionary
          14  results_dict.update({'TFIDF' : {'trainauc':str(auc(train_fpr, train_tpr)),
          15                                  'testauc': str(auc(test_fpr, test_tpr)),
          16                                  'max_depth' : best_params['depth'],
          17                                  'n_estimators' :best_params['n_estimators']} })
```
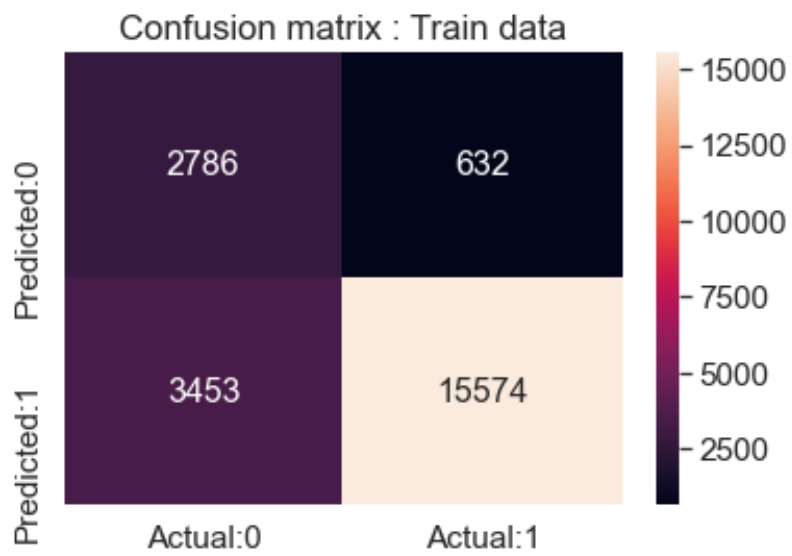


### 5.2.4 Confusion Matrix on train and test data
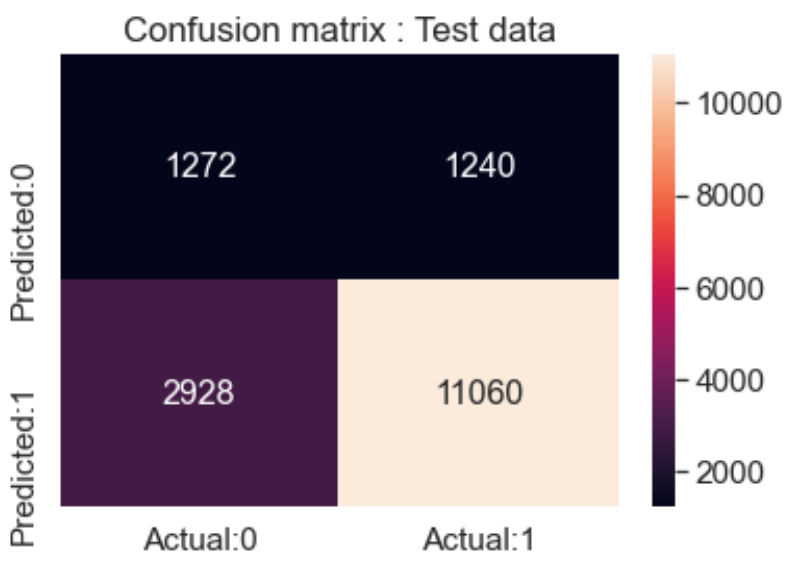
```python
In [119]:  1  print("="*100)
           2  from sklearn.metrics import confusion_matrix
           3  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
           4  print("Train confusion matrix")
           5  print(confusion_matrix(y_train, predict_with_best_t(train_predicts, best_t)))
           6  print("Test confusion matrix")
           7  print(confusion_matrix(y_test, predict_with_best_t(test_predicts, best_t)))
```

```
====================================================================================================
the maximum value of tpr*(1-fpr) 0.6671736812794408 for threshold 0.513
Train confusion matrix
[[ 2786    632]
 [ 3453 15574]]
Test confusion matrix
[[ 1272   1240]
 [ 2928 11060]]
```

```python
In [120]:  1  ### PLOT the matrix for Train
           2  #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
           3  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
           4  df_cm = pd.DataFrame(confusion_matrix(y_train, predict_with_best_t(train_predicts, best_t))
           5                       , range(2), range(2))
           6  # plt.figure(figsize=(10,7))
           7  sns.set(font_scale=1.4) # for label size
           8  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',
           9              xticklabels=['Actual:0','Actual:1']
          10              ,yticklabels=['Predicted:0','Predicted:1']) # font size
          11  plt.title('Confusion matrix : Train data')
          12  plt.show()
```

```
In [121]:    1  ### PLOT the matrix for Train
             2  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
             3  df_cm = pd.DataFrame(confusion_matrix(y_test, predict_with_best_t(test_predicts, best_t)), range(2), range(2))
             4  # plt.figure(figsize=(10,7))
             5  sns.set(font_scale=1.4) # for label size
             6  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',xticklabels=['Actual:0','Actual:1']
             7             ,yticklabels=['Predicted:0','Predicted:1']) # font size
             8  plt.title('Confusion matrix : Test data')
             9  plt.show()
```



## 5.3 Applying Random Forests on AVG W2V, SET 3

```
In [122]:    1  ## converting to csr_matrix to avoid error : "could not broadcast input array from shape
             2  ## reference : https://stackoverflow.com/questions/24924940/convert-list-and-list-of-lists-to-scipy-sparse-arrays
             3  avg_w2v_vectors_train=scipy.sparse.csr_matrix(avg_w2v_vectors_train)
             4  avg_w2v_vectors_title_train=scipy.sparse.csr_matrix(avg_w2v_vectors_title_train)
             5  avg_w2v_vectors_test=scipy.sparse.csr_matrix(avg_w2v_vectors_test)
             6  avg_w2v_vectors_title_test=scipy.sparse.csr_matrix(avg_w2v_vectors_title_test)
             7  avg_w2v_vectors_cv=scipy.sparse.csr_matrix(avg_w2v_vectors_cv)
             8  avg_w2v_vectors_title_cv=scipy.sparse.csr_matrix(avg_w2v_vectors_title_cv)
             9
```

```
In [123]:    1  # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039b
             2  X_tr = hstack((vectorize_tr_trpr, vectorize_tr_proj,vectorize_tr_cat,vectorize_tr_subcat,
             3          vectorize_tr_sklst,avg_w2v_vectors_train,avg_w2v_vectors_title_train,sentiment_pos_train_norm,
             4          sentiment_neg_train_norm,sentiment_compound_train_norm,sentiment_neu_train_norm,
             5          X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
             6          essay_word_count_train_norm)).tocsr()
             7
             8
             9  X_te = hstack((vectorize_te_trpr, vectorize_te_proj,vectorize_te_cat,vectorize_te_subcat,
            10          vectorize_te_sklst,avg_w2v_vectors_test,avg_w2v_vectors_title_test,sentiment_pos_test_norm,
            11          sentiment_neg_test_norm,sentiment_compound_test_norm,sentiment_neu_test_norm,
            12          X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
            13          essay_word_count_test_norm)).tocsr()
            14
            15  X_cr = hstack((vectorize_cv_trpr, vectorize_cv_proj,vectorize_cv_cat,vectorize_cv_subcat,
            16          vectorize_cv_sklst,avg_w2v_vectors_cv,avg_w2v_vectors_title_cv,sentiment_pos_cv_norm,
            17          sentiment_neg_cv_norm,sentiment_compound_cv_norm,sentiment_neu_cv_norm,
            18          X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
            19          essay_word_count_cv_norm)).tocsr()
            20
            21
            22  print(X_tr.shape)
            23  print(X_te.shape)
            24  print(X_cr.shape)
```

```
(22445, 619)
(16500, 619)
(11055, 619)
```

```
In [382]:    1  print("Final Data matrix")
             2  print(X_tr.shape, y_train.shape)
             3  print(X_cr.shape, y_cv.shape)
             4  print(X_te.shape, y_test.shape)
             5  print("="*100)
```

```
Final Data matrix
(22445, 619) (22445,)
(11055, 619) (11055,)
(16500, 619) (16500,)
====================================================================================================
```

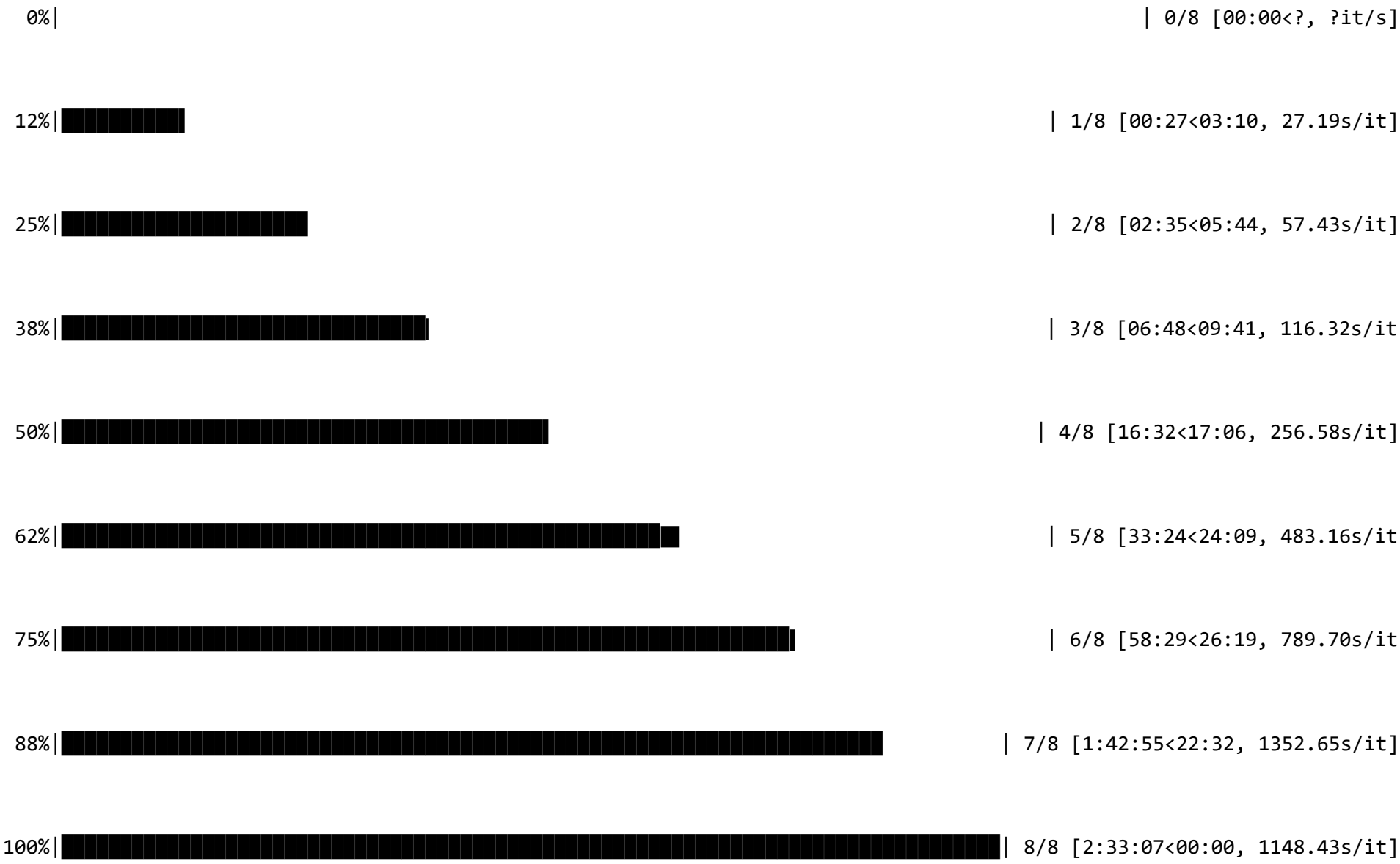```
In [384]:    1  ### taking only 20k points for this set from the randomly sampled 50k dataset.
             2  X_tr=X_tr[:10000]
             3  X_cr=X_cr[:5000]
             4  X_te=X_te[:5000]
             5  y_tr=y_train[:10000]
             6  y_cr=y_cv[:5000]
             7  y_te=y_test[:5000]
```

```
In [385]:    1  print("Final Data matrix")
             2  print(X_tr.shape, y_tr.shape)
             3  print(X_cr.shape, y_cr.shape)
             4  print(X_te.shape, y_te.shape)
             5  print("="*100)
```

```
Final Data matrix
(10000, 619) (10000,)
(5000, 619) (5000,)
(5000, 619) (5000,)
====================================================================================================
```

### 5.3.1 Write loop to find best hyperparameters and do simple cross validation
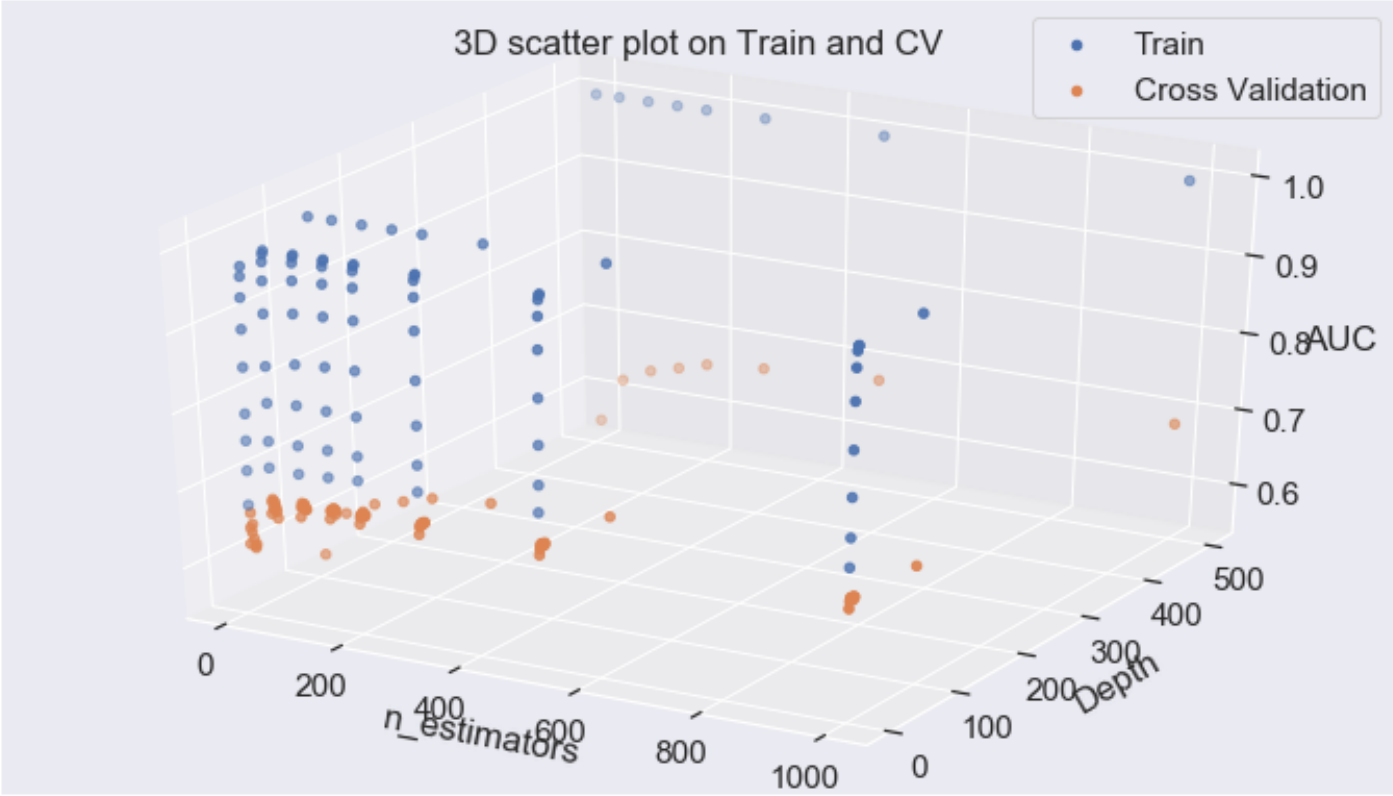
```
In [386]:    1  ## Iterative looop to find best parameter
             2  n_estimators = [10, 50, 100, 150, 200, 300, 500, 1000]
             3  max_depth = [2, 3, 4, 5, 6, 7, 8, 9, 10,100,500]
             4
             5  train_auc,cv_auc,es,de=hyperparam_auc(max_depth,n_estimators,X_tr,X_cr,y_tr,y_cr)
```

```
  0%|                                                        | 0/8 [00:00<?, ?it/s]

 12%|██████                                                  | 1/8 [00:27<03:10, 27.19s/it]

 25%|████████████                                            | 2/8 [02:35<05:44, 57.43s/it]

 38%|██████████████████                                      | 3/8 [06:48<09:41, 116.32s/it]

 50%|████████████████████████                                | 4/8 [16:32<17:06, 256.58s/it]

 62%|██████████████████████████████                          | 5/8 [33:24<24:09, 483.16s/it]

 75%|████████████████████████████████████                    | 6/8 [58:29<26:19, 789.70s/it]

 88%|██████████████████████████████████████████              | 7/8 [1:42:55<22:32, 1352.65s/it]

100%|████████████████████████████████████████████████████████| 8/8 [2:33:07<00:00, 1148.43s/it]
```

### 5.3.2 Representation using 3D scatter plot

```
In [387]:    1  #from mpl_toolkits.mplot3d import Axes3D
             2  plt.figure(figsize=(12,7))
             3  ax = plt.axes(projection='3d')
             4  zipped_=list(map(list,zip(es,de)))
             5  x1=[i[0] for i in zipped_]
             6  y1=[i[1] for i in zipped_]
             7  # Data for three-dimensional scattered pointsf
             8  # reference : https://jakevdp.github.io/PythonDataScienceHandbook/04.12-three-dimensional-plotting.html
             9  ax.scatter3D(x1, y1, train_auc,label='Train')
            10  ax.scatter3D(x1, y1, cv_auc,label='Cross Validation')
            11  ax.set_xlabel('n_estimators')
            12  ax.set_ylabel('Depth')
            13  ax.set_zlabel('AUC')
            14  ax.set_title('3D scatter plot on Train and CV')
            15  ax.legend()
```

Out[387]: <matplotlib.legend.Legend at 0x1922aa88b70>



### 5.3.3 Finding best parameter for training the model
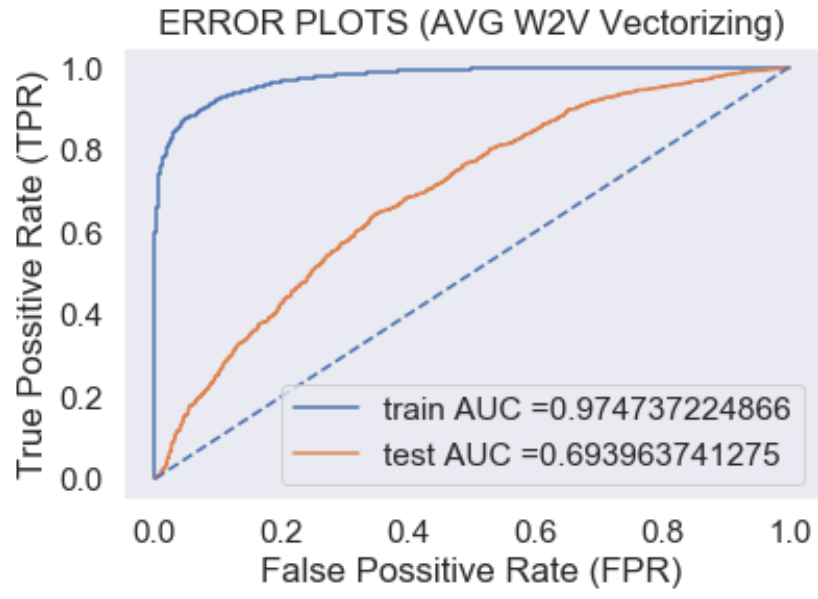
```
In [388]:    1  best_params=find_best_hyperparam(x1,y1,cv_auc)
             2  best_params
```

Max_auc is 0.692991403363

Out[388]: {'n_estimators': 1000, 'depth': 7}

```
In [389]:    1  RF=RandomForestClassifier(class_weight='balanced',n_estimators=best_params['n_estimators'],
             2                            max_depth=best_params['depth'],random_state=4)
             3  RF.fit(X_tr,y_tr)
             4
             5  ## Predict the test
             6  train_predicts=RF.predict_proba(X_tr)[:,1]
             7  test_predicts=RF.predict_proba(X_te)[:,1]
             8
             9  ## Store fpr and tpr rates
            10
            11  train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, train_predicts)
            12  test_fpr, test_tpr, te_thresholds = roc_curve(y_te, test_predicts)
```

```
In [390]:    1  #plot
             2  plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
             3  plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
             4  plt.legend()
             5  plt.xlabel("False Positive Rate (FPR)")
             6  plt.ylabel("True Positive Rate (TPR)")
             7  plt.title("ERROR PLOTS (AVG W2V Vectorizing)")
             8  plt.plot([0, 1], [0, 1],'b--')
             9  plt.grid()
            10  plt.show()
            11
            12
            13  ## Store auc results in a dictionary
            14  results_dict.update({'AVGW2v' : {'trainauc':str(auc(train_fpr, train_tpr)),
            15                       'testauc': str(auc(test_fpr, test_tpr)),
            16                       'max_depth' : best_params['depth'],
            17                       'n_estimators' :best_params['n_estimators']} })
```
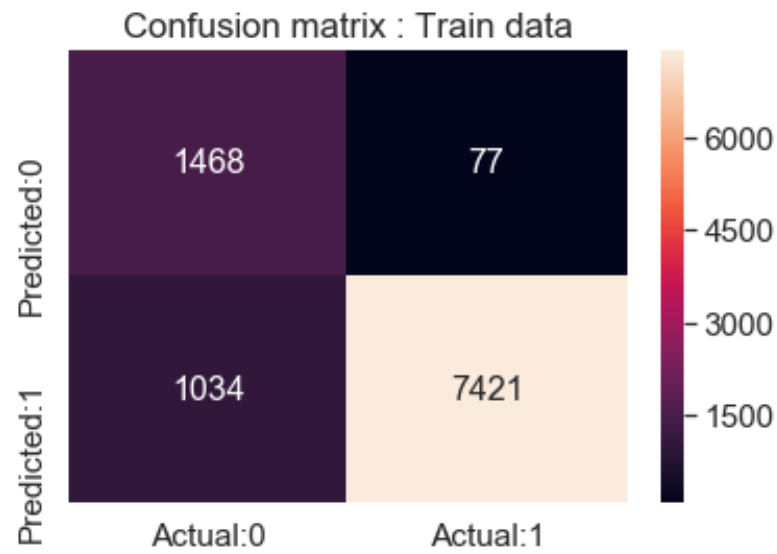


### 5.3.4 Confusion Matrix on train and test data

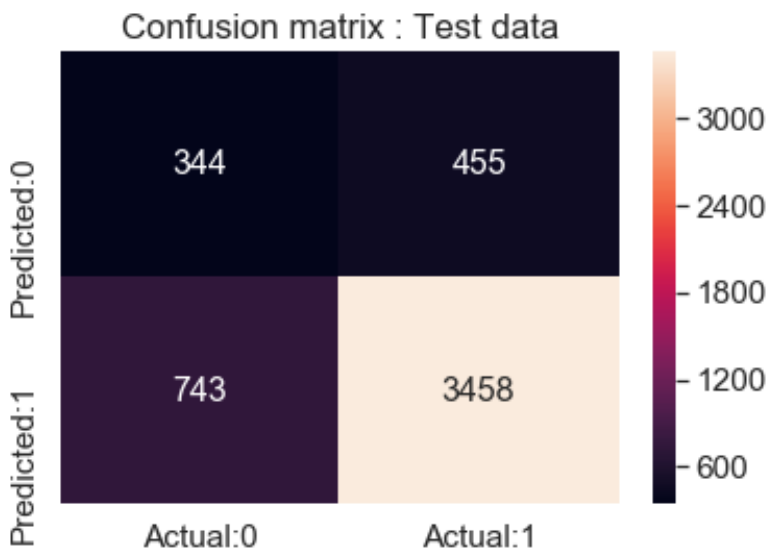```
In [391]:    1  print("="*100)
             2  from sklearn.metrics import confusion_matrix
             3  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
             4  print("Train confusion matrix")
             5  print(confusion_matrix(y_tr, predict_with_best_t(train_predicts, best_t)))
             6  print("Test confusion matrix")
             7  print(confusion_matrix(y_te, predict_with_best_t(test_predicts, best_t)))
```

```
====================================================================================================
the maximum value of tpr*(1-fpr) 0.833962248263 for threshold 0.537
Train confusion matrix
[[1468   77]
 [1034 7421]]
Test confusion matrix
[[ 344  455]
 [ 743 3458]]
```

```
In [392]:    1  ### PLOT the matrix for Train
             2  #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
             3  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
             4  df_cm = pd.DataFrame(confusion_matrix(y_tr, predict_with_best_t(train_predicts, best_t))
             5                       , range(2), range(2))
             6  # plt.figure(figsize=(10,7))
             7  sns.set(font_scale=1.4) # for label size
             8  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',
             9              xticklabels=['Actual:0','Actual:1']
            10              ,yticklabels=['Predicted:0','Predicted:1']) # font size
            11  plt.title('Confusion matrix : Train data')
            12  plt.show()
```

```
In [393]:    1  ### PLOT the matrix for Train
             2  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
             3  df_cm = pd.DataFrame(confusion_matrix(y_te, predict_with_best_t(test_predicts, best_t)), range(2), range(2))
             4  # plt.figure(figsize=(10,7))
             5  sns.set(font_scale=1.4) # for label size
             6  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',xticklabels=['Actual:0','Actual:1']
             7               ,yticklabels=['Predicted:0','Predicted:1']) # font size
             8  plt.title('Confusion matrix : Test data')
             9  plt.show()
```



## 5.4 Applying Random Forests on TFIDF W2V SET 4

```
In [394]:    1  ## converting to csr_matrix to avoid error : "could not broadcast input array from shape
             2  ## reference : https://stackoverflow.com/questions/24924940/convert-list-and-list-of-lists-to-scipy-sparse-arrays
             3  tfidf_w2v_vectors_train=scipy.sparse.csr_matrix(tfidf_w2v_vectors_train)
             4  tfidf_w2v_vectors_title_train=scipy.sparse.csr_matrix(tfidf_w2v_vectors_title_train)
             5  tfidf_w2v_vectors_test=scipy.sparse.csr_matrix(tfidf_w2v_vectors_test)
             6  tfidf_w2v_vectors_title_test=scipy.sparse.csr_matrix(tfidf_w2v_vectors_title_test)
             7  tfidf_w2v_vectors_cv=scipy.sparse.csr_matrix(tfidf_w2v_vectors_cv)
             8  tfidf_w2v_vectors_title_cv=scipy.sparse.csr_matrix(tfidf_w2v_vectors_title_cv)
             9
```

```
In [395]:    1
             2  X_tr = hstack((vectorize_tr_trpr, vectorize_tr_proj,vectorize_tr_cat,vectorize_tr_subcat,
             3                 vectorize_tr_sklst,tfidf_w2v_vectors_train,tfidf_w2v_vectors_title_train,
             4                 sentiment_pos_train_norm,sentiment_neg_train_norm,sentiment_compound_train_norm,sentiment_neu_train_norm,
             5                 X_train_price_norm,quantity_train_norm,prev_projects_train_norm,title_word_count_train_norm,
             6                 essay_word_count_train_norm)).tocsr()
             7
             8  X_te = hstack((vectorize_te_trpr, vectorize_te_proj,vectorize_te_cat,vectorize_te_subcat,
             9                 vectorize_te_sklst,tfidf_w2v_vectors_test,tfidf_w2v_vectors_title_test,
            10  sentiment_pos_test_norm,sentiment_neg_test_norm,sentiment_compound_test_norm,sentiment_neu_test_norm,
            11                 X_test_price_norm,quantity_test_norm,prev_projects_test_norm,title_word_count_test_norm,
            12                 essay_word_count_test_norm)).tocsr()
            13
            14  X_cr = hstack((vectorize_cv_trpr, vectorize_cv_proj,vectorize_cv_cat,vectorize_cv_subcat,
            15                 vectorize_cv_sklst,tfidf_w2v_vectors_cv,tfidf_w2v_vectors_title_cv,sentiment_pos_cv_norm,
            16                 sentiment_neg_cv_norm,sentiment_compound_cv_norm,sentiment_neu_cv_norm,
            17                 X_cv_price_norm,quantity_cv_norm,prev_projects_cv_norm,title_word_count_cv_norm,
            18                 essay_word_count_cv_norm)).tocsr()
            19
            20
            21  print(X_tr.shape)
            22  print(X_te.shape)
            23  print(X_cr.shape)
```

```
(22445, 619)
(16500, 619)
(11055, 619)
```

```
In [396]:    1  print("Final Data matrix")
             2  print(X_tr.shape, y_train.shape)
             3  print(X_cr.shape, y_cv.shape)
             4  print(X_te.shape, y_test.shape)
             5  print("="*100)
```

```
Final Data matrix
(22445, 619) (22445,)
(11055, 619) (11055,)
(16500, 619) (16500,)
====================================================================================================
```
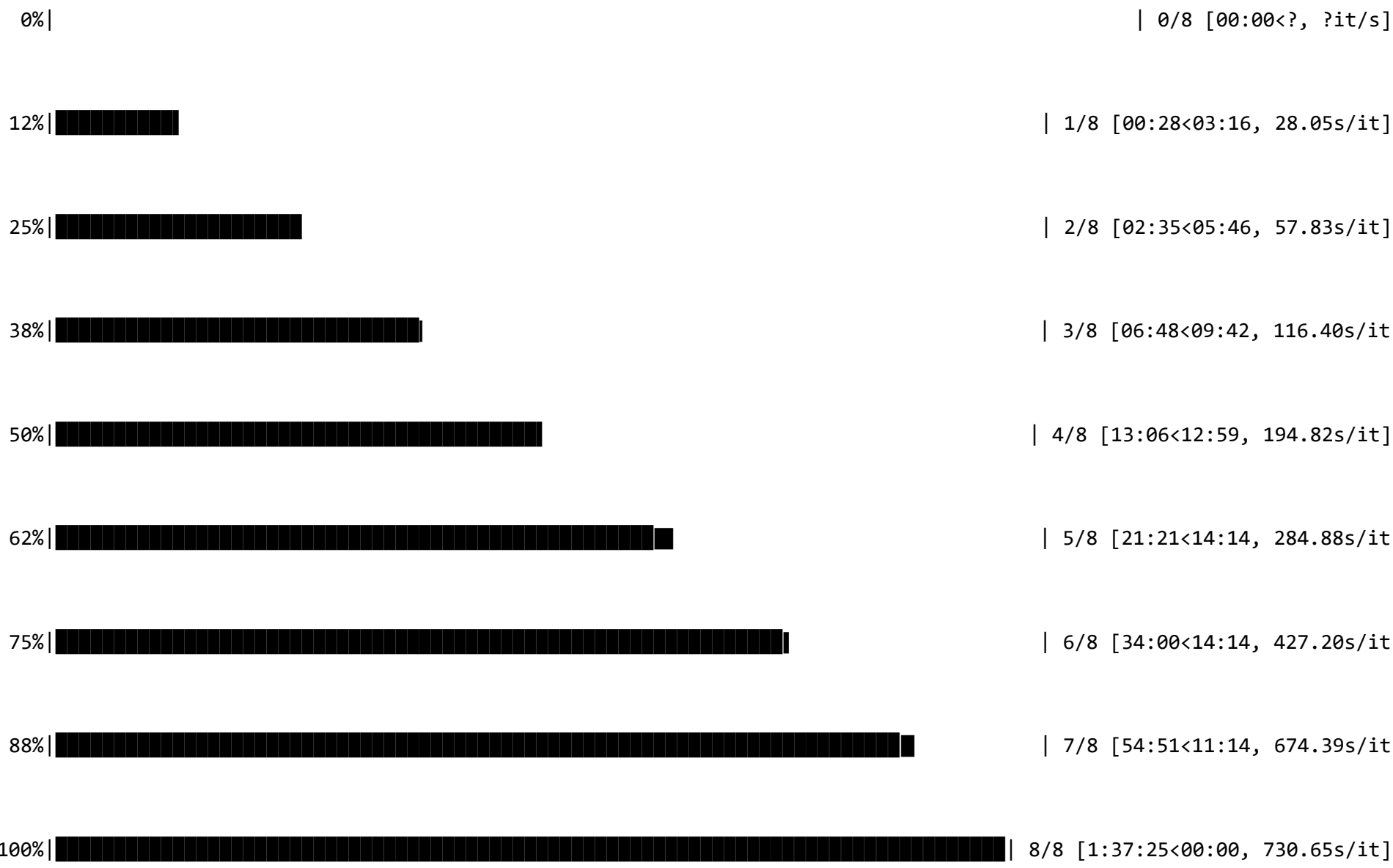
```
In [397]:    1  ### taking only 20k points for this set from the randomly sampled 50k dataset.
             2  X_tr=X_tr[:10000]
             3  X_cr=X_cr[:5000]
             4  X_te=X_te[:5000]
             5  y_tr=y_train[:10000]
             6  y_cr=y_cv[:5000]
             7  y_te=y_test[:5000]
```

```
In [398]:    1  print("Final Data matrix")
             2  print(X_tr.shape, y_tr.shape)
             3  print(X_cr.shape, y_cr.shape)
             4  print(X_te.shape, y_te.shape)
             5  print("="*100)
```

```
Final Data matrix
(10000, 619) (10000,)
(5000, 619) (5000,)
(5000, 619) (5000,)
====================================================================================================
```

### 5.4.1 Write loop to find best hyperparameters and do simple cross validation
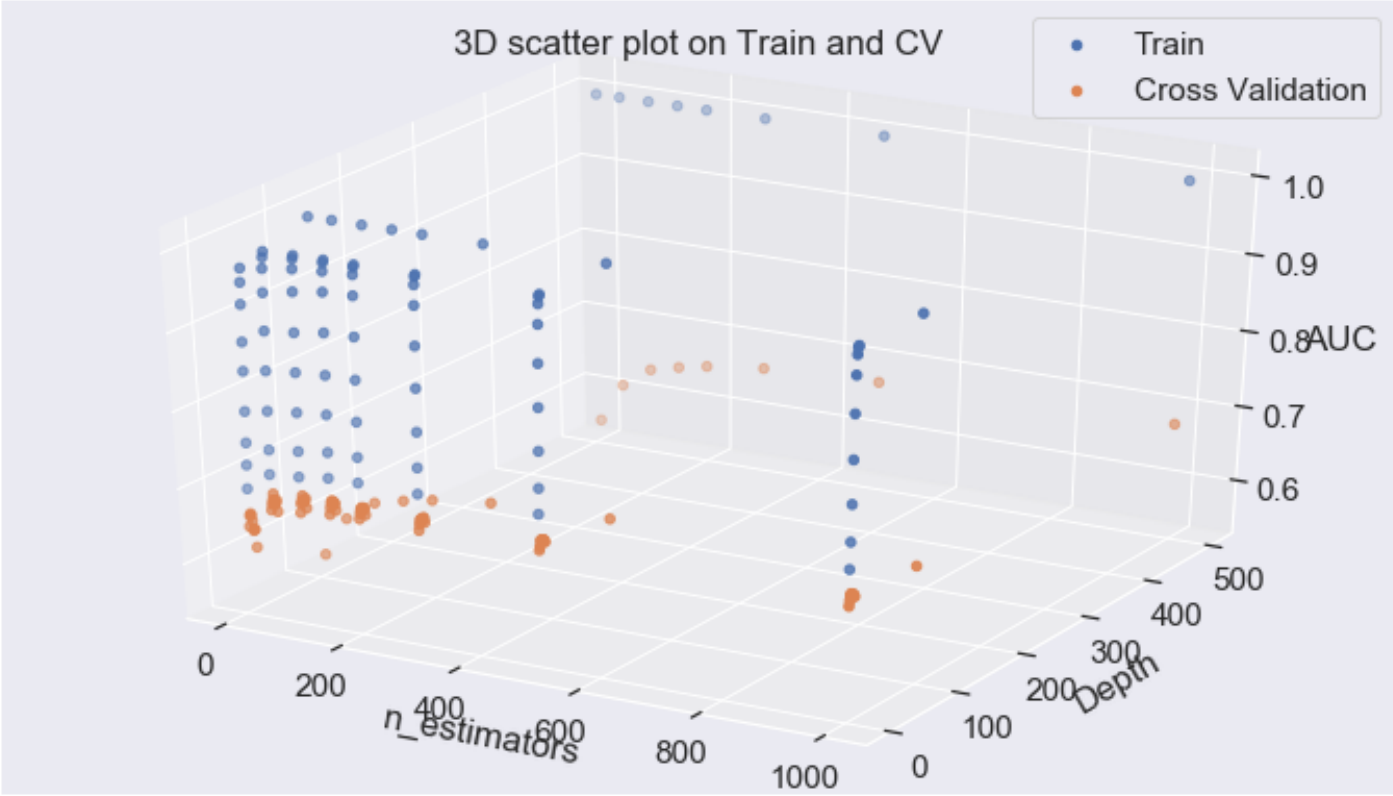
```
In [399]:    1  ## Iterative looop to find best parameter
             2  n_estimators = [10, 50, 100, 150, 200, 300, 500, 1000]
             3  max_depth = [2, 3, 4, 5, 6, 7, 8, 9, 10,100,500]
             4  train_auc,cv_auc,es,de=hyperparam_auc(max_depth,n_estimators,X_tr,X_cr,y_tr,y_cr)
```

```
  0%|                                                      | 0/8 [00:00<?, ?it/s]
 12%|██████                                                | 1/8 [00:28<03:16, 28.05s/it]
 25%|████████████                                          | 2/8 [02:35<05:46, 57.83s/it]
 38%|██████████████████                                    | 3/8 [06:48<09:42, 116.40s/it]
 50%|████████████████████████                              | 4/8 [13:06<12:59, 194.82s/it]
 62%|██████████████████████████████                        | 5/8 [21:21<14:14, 284.88s/it]
 75%|████████████████████████████████████                  | 6/8 [34:00<14:14, 427.20s/it]
 88%|██████████████████████████████████████████            | 7/8 [54:51<11:14, 674.39s/it]
100%|██████████████████████████████████████████████████████| 8/8 [1:37:25<00:00, 730.65s/it]
```

### 5.4.2 Representation using 3D scatter plot

```
In [400]: ▶    1  #from mpl_toolkits.mplot3d import Axes3D
               2  plt.figure(figsize=(12,7))
               3  ax = plt.axes(projection='3d')
               4  zipped_=list(map(list,zip(es,de)))
               5  x1=[i[0] for i in zipped_]
               6  y1=[i[1] for i in zipped_]
               7  # Data for three-dimensional scattered points
               8  # reference : https://jakevdp.github.io/PythonDataScienceHandbook/04.12-three-dimensional-plotting.html
               9  ax.scatter3D(x1, y1, train_auc,label='Train')
              10  ax.scatter3D(x1, y1, cv_auc,label='Cross Validation')
              11  ax.set_xlabel('n_estimators')
              12  ax.set_ylabel('Depth')
              13  ax.set_zlabel('AUC')
              14  ax.set_title('3D scatter plot on Train and CV')
              15  ax.legend()
```

Out[400]: <matplotlib.legend.Legend at 0x19228ee2da0>



### 5.4.3 Finding best parameter for training the model
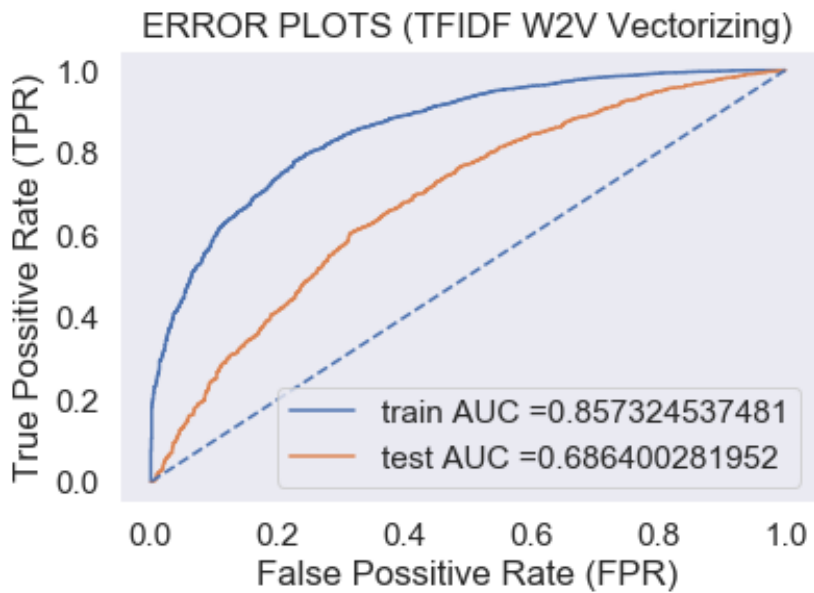
```
In [401]: ▶    1  best_params=find_best_hyperparam(x1,y1,cv_auc)
               2  best_params
```

Max_auc is 0.695786664206

Out[401]: {'n_estimators': 150, 'depth': 5}

```
In [402]: ▶    1  RF=RandomForestClassifier(class_weight='balanced',n_estimators=best_params['n_estimators'],
               2                            max_depth=best_params['depth'],random_state=4)
               3  RF.fit(X_tr,y_tr)
               4
               5  ## Predict the test
               6  train_predicts=RF.predict_proba(X_tr)[:,1]
               7  test_predicts=RF.predict_proba(X_te)[:,1]
               8
               9  ## Store fpr and tpr rates
              10
              11  train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, train_predicts)
              12  test_fpr, test_tpr, te_thresholds = roc_curve(y_te, test_predicts)
```

```
In [408]: ▶    1  #plot
               2  plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, train_tpr)))
               3  plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_tpr)))
               4  plt.legend()
               5  plt.xlabel("False Positive Rate (FPR)")
               6  plt.ylabel("True Possitive Rate (TPR)")
               7  plt.title("ERROR PLOTS (TFIDF W2V Vectorizing)")
               8  plt.plot([0, 1], [0, 1],'b--')
               9  plt.grid()
              10  plt.show()
              11
              12
              13  ## Store auc results in a dictionary
              14  results_dict.update({'TFIDFW2V' : {'trainauc':str(auc(train_fpr, train_tpr)),
              15                      'testauc': str(auc(test_fpr, test_tpr)),
              16                      'max_depth' : best_params['depth'],
              17                      'n_estimators' :best_params['n_estimators']} })
```



### 5.4.4 Confusion Matrix on train and test data

```
In [405]: ▶    1  print("="*100)
               2  #from sklearn.metrics import confusion_matrix
               3  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
               4  print("Train confusion matrix")
               5  print(confusion_matrix(y_tr, predict_with_best_t(train_predicts, best_t)))
               6  print("Test confusion matrix")
               7  print(confusion_matrix(y_te, predict_with_best_t(test_predicts, best_t)))
```

```
====================================================================================================
the maximum value of tpr*(1-fpr) 0.60257827945 for threshold 0.501
Train confusion matrix
[[1195  350]
 [1868 6587]]
Test confusion matrix
[[ 414  385]
 [1017 3184]]
```
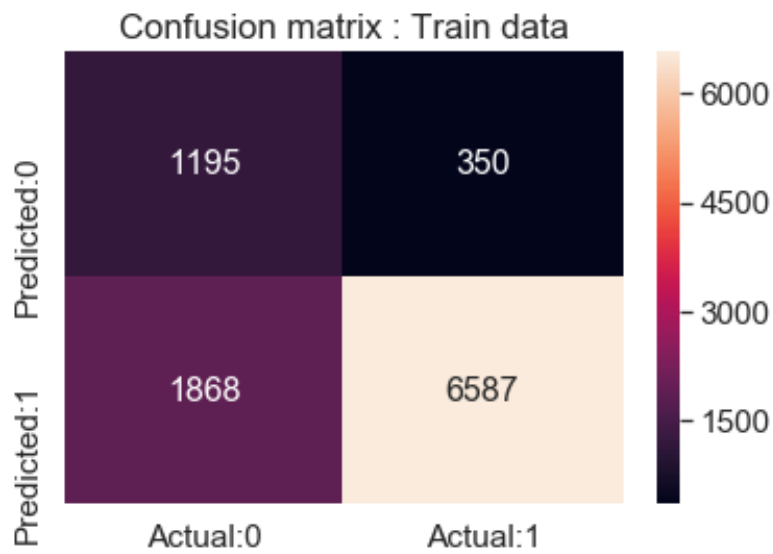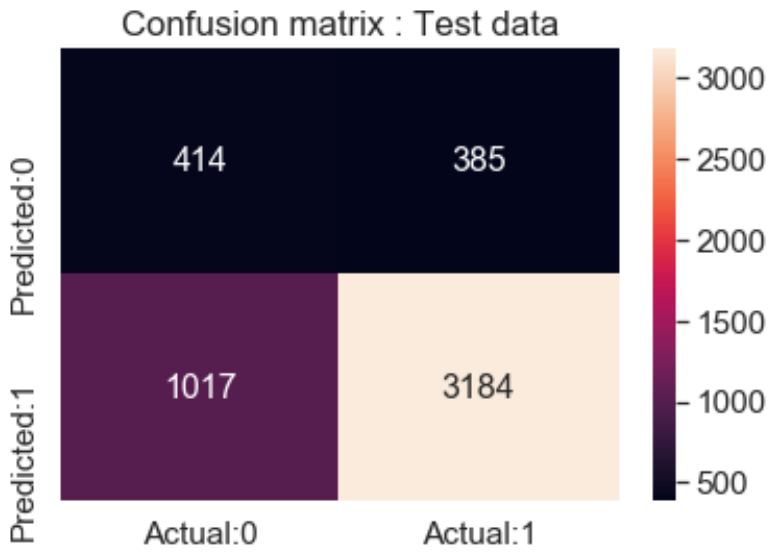
```
In [406]: ▶    1  ### PLOT the matrix for Train
               2  #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
               3  # source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
               4  df_cm = pd.DataFrame(confusion_matrix(y_tr, predict_with_best_t(train_predicts, best_t))
               5                       , range(2), range(2))
               6  # plt.figure(figsize=(10,7))
               7  sns.set(font_scale=1.4) # for Label size
               8  sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',
               9              xticklabels=['Actual:0','Actual:1']
              10              ,yticklabels=['Predicted:0','Predicted:1']) # font size
              11  plt.title('Confusion matrix : Train data')
              12  plt.show()
```

In [407]: ▶
```
### PLOT the matrix for Train
# source : https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
df_cm = pd.DataFrame(confusion_matrix(y_te, predict_with_best_t(test_predicts, best_t)), range(2), range(2))
# plt.figure(figsize=(10,7))
sns.set(font_scale=1.4) # for label size
sns.heatmap(df_cm, annot=True, annot_kws={"size": 16},fmt='g',xticklabels=['Actual:0','Actual:1']
            ,yticklabels=['Predicted:0','Predicted:1']) # font size
plt.title('Confusion matrix : Test data')
plt.show()
```