

DonorsChoose

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

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

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

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

About the DonorsChoose Data Set

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

Feature		Description
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: Art Will Make You Happy! First Grade Fun
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: Grades PreK-2 Grades 3-5 Grades 6-8 Grades 9-12
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: Applied Learning Care & Hunger Health & Sports History & Civics Literacy & Language Math & Science Music & The Arts Special Needs Warmth
<code>school_state</code>		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
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: Music & The Arts Literacy & Language, Math & Science
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: My students need hands on literacy materials to manage sensory needs!</code

Feature	Description
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values:
	nan
	Dr.
	Mr.
	Mrs.
	Ms.
	Teacher.
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. Example: 2

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

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

Feature	Description
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
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 <code>0</code> indicates the project was not approved, and a value of <code>1</code> indicates the project was approved.

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

1.1 Reading Data

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

In [3]: ▶

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

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

The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
'project_submitted_datetime' 'project_grade_category'
'project_subject_categories' 'project_subject_subcategories'
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']

In [4]: ▶

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

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

Out[4]:

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

1.1 Preprocessing Categorical Features: project_grade_category

In [5]: ▶

```
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 [6]: 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[6]: 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 [7]: 1 categories = 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
8 def _process_cat_subcat(categories):
9     cat_list = []
10    for i in categories:
11        temp = ""
12        # consider we have text like this "Math & Science, Warmth, Care & Hunger"
13        for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
14            if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math", "&", "Science"
15                j=j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
16                j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science"=>"Math&Science"
17                temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
18                temp = temp.replace('&', '_') # we are replacing the & value into
19        cat_list.append(temp.strip())
20    return cat_list
21
22
23
24 project_data['clean_categories'] = _process_cat_subcat(categories)
25 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
26 project_data.head(2)
27
28 ### maintain a dict that
29 my_counter=Counter()
30 for word in project_data['clean_categories'].values:
31     my_counter.update(word.split())
32 cat_dict=dict(my_counter)
33 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
34
35
```

1.3 Preprocessing Categorical Features: project_subject_category

```
In [8]: ▶ 1 sub_categories = list(project_data['project_subject_subcategories'].values)
2 project_data['clean_subcategories'] = _process_cat_subcat(sub_categories)
3 project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)
4
5 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
6 my_counter = Counter()
7 for word in project_data['clean_subcategories'].values:
8     my_counter.update(word.split())
9
10 sub_cat_dict = dict(my_counter)
11 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
```

1.4 Preprocessing Categorical Features: school_state

In [9]:

```
1 project_data['school_state'].value_counts()
2 ## Convert it to Lower
3 project_data['school_state'] = project_data['school_state'].str.lower()
4 print(project_data['school_state'].value_counts(dropna=False))
5 print('No nan values in this feature')
```

```
ca    15388
tx     7396
ny     7318
fl     6185
nc     5091
il     4350
ga     3963
sc     3936
mi     3161
pa     3109
in     2620
mo     2576
oh     2467
la     2394
ma     2389
wa     2334
ok     2276
nj     2237
az     2147
va     2045
wi     1827
al     1762
ut     1731
tn     1688
ct     1663
md     1514
nv     1367
ms     1323
ky     1304
or     1242
mn     1208
co     1111
ar     1049
id       693
ia       666
ks       634
nm       557
dc       516
hi       507
me       505
wv       503
nh       348
ak       345
de       343
ne       309
sd       300
ri       285
mt       245
nd       143
wy        98
vt        80
```

Name: school_state, dtype: int64

No nan values in this feature

1.5 Preprocessing Categorical Features: Teacher_prefix

```
In [10]: 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
```

1.6 Combining all the essays

```
In [11]: 1 print('Number of nan values in essay1 is ', len(project_data[project_data["project_essay_1"].isna()==True]))
2 print('Number of nan values in essay2 is ', len(project_data[project_data["project_essay_2"].isna()==True]))
3 print('Number of nan values in essay3 is ', len(project_data[project_data["project_essay_3"].isna()==True]))
4 print('Number of nan values in essay4 is ', len(project_data[project_data["project_essay_4"].isna()==True]))
```

```
Number of nan values in essay1 is  0
Number of nan values in essay2 is  0
Number of nan values in essay3 is  105490
Number of nan values in essay4 is  105490
```

```
In [12]: 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)
```

1.7 Number of Words in the Essay and Title

```
In [13]: 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
```

1.8. Preprocessing Numerical Values: price


```
In [14]: 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 ##merge into the project_Data
4 project_data = pd.merge(project_data,price_data,on='id',how='left')
```

1.9 Preprocessing Text Features: project_title , essay

```
In [15]: 1 # https://stackoverflow.com/a/47091490/4084039
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
16 # https://gist.github.com/sebleier/554280
17 # we are removing the words from the stop words List: 'no', 'nor', 'not'
18 stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
19             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
20             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
21             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
22             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
23             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
24             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
25             'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
26             'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \
27             'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
28             's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
29             've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
30             "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
31             "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
32             'won', "won't", 'wouldn', "wouldn't"]
33
34 print("printing some random reviews")
35 print(9, project_data['project_title'].values[9])
36 print(34, project_data['project_title'].values[34])
37 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!!
```

```

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

```

```
1 preprocessed_titles = preprocess_text(project_data['project_title'].values)
2 #merge the column in the project_data
3 project_data['processed_title']=preprocessed_titles
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109248 [00:08<00:00, 13383.18it/s]
```

```
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 read pure pleasur
34 ball
147 need chromebook
```

```
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])
5 preprocessed_essays = preprocess_text(project_data['essay'].values)
```

```
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!!
```

[illegible]

printing some random essay

9 95 student free reduc lunch homeless despit come school eager learn student inquisit eager learner embrac challeng not great book resourc everi day mani not afford opportun e ngag big color page book regular basi home not travel public librari duti teacher provid student opportun succeed everi aspect life read fundament student read book boost compr ehens skill book use read aloud partner read independ read engag read build love read read pure enjoy introduc new author well old favorit want student readi 21st centuri know pleasur hold good hard back book hand noth like good book read student soar read consider generous fund contribut help build stamina prepar 3rd grade thank much read propos nan nan

34 student main come extrem low incom famili major come home parent work full time student school 7 30 6 00 pm 2 30 6 00 pm school program receiv free reduc meal breakfast lunc h want student feel comfort classroom home mani student take multipl role home well school sometim caretak younger sibl cook babysitt academ friend develop go becom adult consi d essenti part job model help other gain knowledg posit manner result communiti student love help outsid classroom consist look opportun support learn kind help way excit exper i altern seat classroom school year studi shown give student option sit classroom increas focus well motiv allow student choic classroom abl explor creat welcom environ altern classroom seat experi frequent recent year believ along mani other everi child learn differ not appli multipl memor paper written appli space ask work student past ask work lib rari work carpet answer alway long learn work wherev want yoga ball lap desk abl increas option seat classroom expand imagin space nannan

147 student eager learn make mark world come titl 1 school need extra love fourth grade student high poverti area still come school everi day get educ tri make fun educ get sch ool creat care environ student bloom deserv best thank request 1 chromebook access onlin intervent differenti instruct get extra practic chromebook use supplement ela math inst ruct student play ela math game engag fun well particip assign onlin turn help student improv skill chromebook classroom would not allow student use program pace would ensur st udent get adequ time use program onlin program espec benefici student special need abl work level well challeng differ materi make student confid abil chromebook would allow s tudent daili access comput increas comput skill chang live better becom success school access technolog classroom would help bridg achiev gap nannan

1.9 Creating Sentiment Columns

[illegible]

2 Train, Test, CV Split

```
In [22]: 1 # train test split using sklearn.model selection
2 sampled_data = project_data.sample(50000)
3 X_train, X_test, y_train, y_test = train_test_split(sampled_data, sampled_data['project_is_approved'], test_size=0.20, stratify = sampled_data['project_is_approved'], random_state=0)
4 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 [23]: 1 ## drop the y labels from splits
2 X_train.drop(['project_is_approved'], axis=1, inplace=True)
```

```
In [24]: 1 X_train.head(2)
```

Out[24]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_submitted_datetime	project_grade_category	project_title	project_essay_1	project_essay_2	...	words_in_essay	words
69191	102275	p200146	db88eb4e11230df8462b304830e627e2	mrs	ca	2016-08-25 19:19:04	6_8	Classroom Chromebook Set	I teach middle school and they are already tec...	Our school shares a set of Chromebooks between...	...	232	
90301	23896	p218117	4d283fd47f0fe5396e141369003fd47d	mr	va	2016-06-13 22:13:59	3_5	SHOW IT MOVE IT	My students learn by best when they are presen...	Having an interactive whiteboard in our classr...	...	211	

2 rows × 27 columns

3. VECTORIZING DATA

3.1 One hot encoding on Categorical : (categories,subcategories,schoolstate,teacher_prefix,projectgrade,sentiment_columns)

```
In [25]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean categories
3
4 cat_vectorize = CountVectorizer(lowercase=False, binary=True)
5 cat_vectorize.fit(X_train['clean_categories'].values)
6
7 train_categories = cat_vectorize.transform(X_train['clean_categories'].values)
8 test_categories = cat_vectorize.transform(X_test['clean_categories'].values)
9 cv_categories = cat_vectorize.transform(X_cv['clean_categories'].values)
10
11 print(cat_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ", train_categories.shape)
13 print("Shape of matrix of Test data after one hot encoding ", test_categories.shape)
14 print("Shape of matrix of CV data after one hot encoding ", cv_categories.shape)
```

```
['AppliedLearning', 'Care_Hunger', 'Health_Sports', 'History_Civics', 'Literacy_Language', 'Math_Science', 'Music_Arts', 'SpecialNeeds', 'Warmth']
Shape of matrix of Train data after one hot encoding (26800, 9)
Shape of matrix of Test data after one hot encoding (10000, 9)
Shape of matrix of CV data after one hot encoding (13200, 9)
```

```
In [26]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
2 ## clean subcategories
3
4 subcat_vectorize = CountVectorizer(lowercase=False, binary=True)
5 subcat_vectorize.fit(X_train['clean_subcategories'].values)
6
7 train_subcategories = subcat_vectorize.transform(X_train['clean_subcategories'].values)
8 test_subcategories = subcat_vectorize.transform(X_test['clean_subcategories'].values)
9 cv_subcategories = subcat_vectorize.transform(X_cv['clean_subcategories'].values)
10
11 print(subcat_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",train_subcategories.shape)
13 print("Shape of matrix of Test data after one hot encoding ",test_subcategories.shape)
14 print("Shape of matrix of CV data after one hot encoding ",cv_subcategories.shape)
15
```

```
['AppliedSciences', 'Care_Hunger', 'CharacterEducation', 'Civics_Government', 'College_CareerPrep', 'CommunityService', 'ESL', 'EarlyDevelopment', 'Economics', 'EnvironmentalSc
ience', 'Extracurricular', 'FinancialLiteracy', 'ForeignLanguages', 'Gym_Fitness', 'Health_LifeScience', 'Health_Wellness', 'History_Geography', 'Literacy', 'Literature_Writin
g', 'Mathematics', 'Music', 'NutritionEducation', 'Other', 'ParentInvolvement', 'PerformingArts', 'SocialSciences', 'SpecialNeeds', 'TeamSports', 'VisualArts', 'Warmth']
Shape of matrix of Train data after one hot encoding (26800, 30)
Shape of matrix of Test data after one hot encoding (10000, 30)
Shape of matrix of CV data after one hot encoding (13200, 30)
```

```
In [27]: ▶ 1 # we use count vectorizer to convert the values into one hot vectors
2 ## school state
3
4
5 sklstate_vectorize = CountVectorizer(lowercase=False, binary=True)
6 sklstate_vectorize.fit(X_train['school_state'].values)
7
8 sklstate_train = sklstate_vectorize.transform(X_train['school_state'].values)
9 sklstate_test = sklstate_vectorize.transform(X_test['school_state'].values)
10 sklstate_cv = sklstate_vectorize.transform(X_cv['school_state'].values)
11
12 print(sklstate_vectorize.get_feature_names())
13 print("Shape of matrix of Train data after one hot encoding ",sklstate_train.shape)
14 print("Shape of matrix of Test data after one hot encoding ",sklstate_test.shape)
15 print("Shape of matrix of CV data after one hot encoding ",sklstate_cv.shape)
```

```
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'ks', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd',
'ne', 'nh', 'nj', 'nm', 'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv', 'wy']
Shape of matrix of Train data after one hot encoding (26800, 51)
Shape of matrix of Test data after one hot encoding (10000, 51)
Shape of matrix of CV data after one hot encoding (13200, 51)
```

```
In [28]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## teacher_prefix
3
4 teacher_prefix_vectorize = CountVectorizer(lowercase=False, binary=True)
5 teacher_prefix_vectorize.fit(X_train['teacher_prefix'].values)
6
7 teacher_prefix_train = teacher_prefix_vectorize.transform(X_train['teacher_prefix'].values)
8 teacher_prefix_test = teacher_prefix_vectorize.transform(X_test['teacher_prefix'].values)
9 teacher_prefix_cv = teacher_prefix_vectorize.transform(X_cv['teacher_prefix'].values)
10
11 print(teacher_prefix_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",teacher_prefix_train.shape)
13 print("Shape of matrix of Test data after one hot encoding ",teacher_prefix_test.shape)
14 print("Shape of matrix of CV data after one hot encoding ",teacher_prefix_cv.shape)
```

```
['dr', 'mr', 'mrs', 'ms', 'teacher']
```

```
Shape of matrix of Train data after one hot encoding (26800, 5)
```

```
Shape of matrix of Test data after one hot encoding (10000, 5)
```

```
Shape of matrix of CV data after one hot encoding (13200, 5)
```

```
In [29]: 1 # we use count vectorizer to convert the values into one hot vectors
2 ## project_grade
3
4 proj_grade_vectorize = CountVectorizer(lowercase=False, binary=True)
5 proj_grade_vectorize.fit(X_train['project_grade_category'].values)
6
7 proj_grade_train = proj_grade_vectorize.transform(X_train['project_grade_category'].values)
8 proj_grade_test = proj_grade_vectorize.transform(X_test['project_grade_category'].values)
9 proj_grade_cv = proj_grade_vectorize.transform(X_cv['project_grade_category'].values)
10
11 print(proj_grade_vectorize.get_feature_names())
12 print("Shape of matrix of Train data after one hot encoding ",proj_grade_train.shape)
13 print("Shape of matrix of Test data after one hot encoding ",proj_grade_test.shape)
14 print("Shape of matrix of CV data after one hot encoding ",proj_grade_cv.shape)
```

```
['3_5', '6_8', '9_12', 'prek_2']
```

```
Shape of matrix of Train data after one hot encoding (26800, 4)
```

```
Shape of matrix of Test data after one hot encoding (10000, 4)
```

```
Shape of matrix of CV data after one hot encoding (13200, 4)
```

3.2. Vectorizing Numerical Features :

3.2.1 Price , Quantity, teacher_number_of_previously_posted_projects,words_in_essay,words_in_title

```
In [30]: ▶ 1 def normalize_numerical(tr,te,cv):
2           normalizer = Normalizer()
3           # normalizer.fit(tr.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           normalizer.fit(tr.values.reshape(1,-1))
10
11          tr_norm = normalizer.transform(tr.values.reshape(1,-1))
12          cv_norm = normalizer.transform(cv.values.reshape(1,-1))
13          te_norm = normalizer.transform(te.values.reshape(1,-1))
14
15          ## reshaping
16          tr_norm=tr_norm.reshape(-1,1)
17          cv_norm=cv_norm.reshape(-1,1)
18          te_norm=te_norm.reshape(-1,1)
19
20          print("After vectorizations")
21          print(tr_norm.shape, y_train.shape)
22          print(cv_norm.shape, y_cv.shape)
23          print(te_norm.shape, y_test.shape)
24          print("=*100)
25          return tr_norm,cv_norm,te_norm
```



```
In [31]: 1 X_train_price_norm,X_cv_price_norm,X_test_price_norm=normalize_numerical(X_train['price'],X_test['price'],X_cv['price'])
2 quantity_train_norm,quantity_cv_norm,quantity_test_norm=normalize_numerical(X_train['quantity'],X_test['quantity'],X_cv['quantity'])
3 prev_projects_train_norm,prev_projects_cv_norm,prev_projects_test_norm=normalize_numerical(X_train['teacher_number_of_previously_posted_projects'],
4 X_test['teacher_number_of_previously_posted_projects'],
5 X_cv['teacher_number_of_previously_posted_projects'])
6 essay_word_count_train_norm,essay_word_count_cv_norm,essay_word_count_test_norm=normalize_numerical(X_train['words_in_essay'],
7 X_test['words_in_essay'],
8 X_cv['words_in_essay'])
9 title_word_count_train_norm,title_word_count_cv_norm,title_word_count_test_norm=normalize_numerical(X_train['words_in_title'],
10 X_test['words_in_title'],
11 X_cv['words_in_title'])
12
13
```

After vectorizations

(26800, 1) (26800,)

(13200, 1) (13200,)

(10000, 1) (10000,)

After vectorizations

(26800, 1) (26800,)

(13200, 1) (13200,)

(10000, 1) (10000,)

After vectorizations

(26800, 1) (26800,)

(13200, 1) (13200,)

(10000, 1) (10000,)

After vectorizations

(26800, 1) (26800,)

(13200, 1) (13200,)

(10000, 1) (10000,)

After vectorizations

(26800, 1) (26800,)

(13200, 1) (13200,)

(10000, 1) (10000,)

3.2.2 Vectorize Sentiment Columns


```
In [32]: 1  ### vectorize pos
2  Normalize=Normalizer()
3  Normalize.fit(X_train['pos'].values.reshape(1,-1))
4  sentiment_pos_train_norm=Normalize.transform(X_train['pos'].values.reshape(1,-1))
5  sentiment_pos_test_norm=Normalize.transform(X_test['pos'].values.reshape(1,-1))
6  sentiment_pos_cv_norm=Normalize.transform(X_cv['pos'].values.reshape(1,-1))
7
8  sentiment_pos_train_norm=sentiment_pos_train_norm.reshape(-1,1)
9  sentiment_pos_test_norm=sentiment_pos_test_norm.reshape(-1,1)
10 sentiment_pos_cv_norm=sentiment_pos_cv_norm.reshape(-1,1)
11 print('After vectorizations : ')
12 print(sentiment_pos_train_norm.shape)
13 print(sentiment_pos_test_norm.shape)
14 print(sentiment_pos_cv_norm.shape)
```

After vectorizations :
(26800, 1)
(10000, 1)
(13200, 1)

```
In [33]: 1  ### vectorize neu
2  Normalize=Normalizer()
3  Normalize.fit(X_train['neu'].values.reshape(1,-1))
4  sentiment_neu_train_norm=Normalize.transform(X_train['neu'].values.reshape(1,-1))
5  sentiment_neu_test_norm=Normalize.transform(X_test['neu'].values.reshape(1,-1))
6  sentiment_neu_cv_norm=Normalize.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)
10 print('After vectorizations : ')
11 print(sentiment_neu_train_norm.shape)
12 print(sentiment_neu_test_norm.shape)
13 print(sentiment_neu_cv_norm.shape)
```

After vectorizations :
(26800, 1)
(10000, 1)
(13200, 1)

```
In [34]: 1  ### vectorize compound
2  Normalize=Normalizer()
3  Normalize.fit(X_train['compound'].values.reshape(1,-1))
4  sentiment_compound_train_norm=Normalize.transform(X_train['compound'].values.reshape(1,-1))
5  sentiment_compound_test_norm=Normalize.transform(X_test['compound'].values.reshape(1,-1))
6  sentiment_compound_cv_norm=Normalize.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)
10 print('After vectorizations : ')
11 print(sentiment_compound_train_norm.shape)
12 print(sentiment_compound_test_norm.shape)
13 print(sentiment_compound_cv_norm.shape)
```

After vectorizations :
(26800, 1)
(10000, 1)
(13200, 1)

```
In [35]: 1  ### vectorize neg
2  Normalize=Normalizer()
3  Normalize.fit(X_train['neg'].values.reshape(1,-1))
4  sentiment_neg_train_norm=Normalize.transform(X_train['neg'].values.reshape(1,-1))
5  sentiment_neg_test_norm=Normalize.transform(X_test['neg'].values.reshape(1,-1))
6  sentiment_neg_cv_norm=Normalize.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)
10 print('After vectorizations : ')
11 print(sentiment_neg_train_norm.shape)
12 print(sentiment_neg_test_norm.shape)
13 print(sentiment_neg_cv_norm.shape)
```

After vectorizations :

```
(26800, 1)
(10000, 1)
(13200, 1)
```

Assignment_11: DonorsChoose_TruncatedSVD

- **step 1** Select the top 2k words from essay text and project_title (concatenate essay text with project title and then find the top 2k words) based on their `idf_` (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html) values
- **step 2** Compute the co-occurrence matrix with these 2k words, with window size=5 (ref (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2vec/>))
- **step 3** Use [TruncatedSVD](http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html) (<http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>) on calculated co-occurrence matrix and reduce its dimensions, choose the number of components (`n_components`) using [elbow method](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/) (<https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/pca-code-example-using-non-visualization/>).

- The shape of the matrix after TruncatedSVD will be 2000*n, i.e. each row represents a vector form of the corresponding word.
- Vectorize the essay text and project titles using these word vectors. (while vectorizing, do ignore all the words which are not in top 2k words)

- **step 4** Concatenate these truncatedSVD matrix, with the matrix with features
 - **school_state** : categorical data
 - **clean_categories** : categorical data
 - **clean_subcategories** : categorical data
 - **project_grade_category** :categorical data
 - **teacher_prefix** : categorical data
 - **quantity** : numerical data
 - **teacher_number_of_previously_posted_projects** : numerical data
 - **price** : numerical data
 - **sentiment score's of each of the essay** : numerical data
 - **number of words in the title** : numerical data
 - **number of words in the combine essays** : numerical data
 - **word vectors calculated in step 3** : numerical data
- **step 5:** Apply GBDT on matrix that was formed in **step 4** of this assignment, **DO REFER THIS BLOG: XGBOOST DMATRIX** (<https://www.kdnuggets.com/2017/03/simple-xgboost-tutorial-iris-dataset.html>)
- **step 6:**Hyper parameter tuning (Consider any two hyper parameters)
 - 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 paramter using k-fold cross validation or simple cross validation data
 - Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3.3 Vectorizing Essay and Title data : Creating custom Word Vectors

```
In [36]: ▶ 1  ## To create Cooccurrence matrix first select top 2k words from essay and title
2  def top_k_words(data,k):
3      top_k=[]
4      # concatenate words from essay and title
5      corpus=pd.DataFrame(columns=['title_essay'])
6      corpus['title_essay']=data['processed_essay'].map(str) + \
7                          data['processed_title'].map(str)
8      # using tfidf model we are calculating idf for each word
9      tfidf_model = TfidfVectorizer()
10     tfidf_model.fit(corpus['title_essay'].values)
11     # we are converting a dictionary with word as a key, and the idf as a value
12     dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
13     tfidf_words = set(tfidf_model.get_feature_names())
14     # sort the dictionary based on idf values in descending order (i.e) selecting rare words
15     sort_d=sorted(dictionary.items(),key=lambda x:(x[1],x[0]),reverse=True)
16     #append top k words
17     top_2k=[sort_d[i][0] for i in range(k)]
18     return top_2k
```

```
In [37]: ▶ 1  top_2000_words=top_k_words(X_train,2000)
```

Build Co-occurrence matrix from scratch

In [38]:

```

1  # build co-occurrence matrix
2  def get_indices(var,w):
3      idx=[]
4      for i in range(len(var)):
5          if var[i] == w :
6              idx.append(i)
7      return idx
8  def get_co_matrix(win_size,entire_corpus,wrd):
9      df=pd.DataFrame(columns=wrd,index=wrd)
10     df=df.fillna(0)
11     for i in tqdm(range(len(entire_corpus))):
12         split_=entire_corpus[i].split()
13         imp_words=[]
14         for i in split_:
15             ## store the words that are present in top 2000 words
16             if i in wrd and i not in imp_words:imp_words.append(i)
17         for wrd1 in imp_words :
18             for wrd2 in imp_words :
19                 # find all the indices where word exists in the corpus
20                 indices=get_indices(split_,wrd1)
21                 for id_ in indices:
22                     ## if the word exists in the start
23                     if id_ == 0 :
24                         if wrd2 in split_[id_+1:id_+win_size+1]:
25                             df.at[wrd1,wrd2]+=1
26                     ## if the word exists in the end
27                     elif id_ == len(split_)-1 :
28                         if wrd2 in split_[id_-1:id_-win_size-1:-1] :
29                             df.at[wrd1,wrd2]+=1
30                     ## if the word exists in between
31                     else:
32                         ## if the word exists in the starting range less than the window size i.e
33                         ## if window=5 and index of word is at 2 ,we cant slice 5 words before as it will show error
34                         if id_ <=win_size :
35                             if wrd2 in split_[id_-1::-1]:
36                                 df.at[wrd1,wrd2]+=1
37                             if wrd2 in split_[id_+1:id_+win_size+1]:
38                                 df.at[wrd1,wrd2]+=1
39                         ## if the word exists in the ending range less than the window size i.e
40                         ## if window=5 and index of word is at 2 ,we cant slice 5 words after as it will show error
41                         elif id_ >= len(split_)-win_size-1:
42                             if wrd2 in split_[id_+1:]:
43                                 df.at[wrd1,wrd2]+=1
44                             if wrd2 in split_[id_-1:id_-win_size-1:-1]:
45                                 df.at[wrd1,wrd2]+=1
46                     # if words are anywhere in between
47                     else :
48                         if wrd2 in split_[id_+1:id_+win_size+1] :
49                             df.at[wrd1,wrd2]+=1
50                         if wrd2 in split_[id_-1:id_-win_size-1:-1] :
51                             df.at[wrd1,wrd2]+=1
52     return df

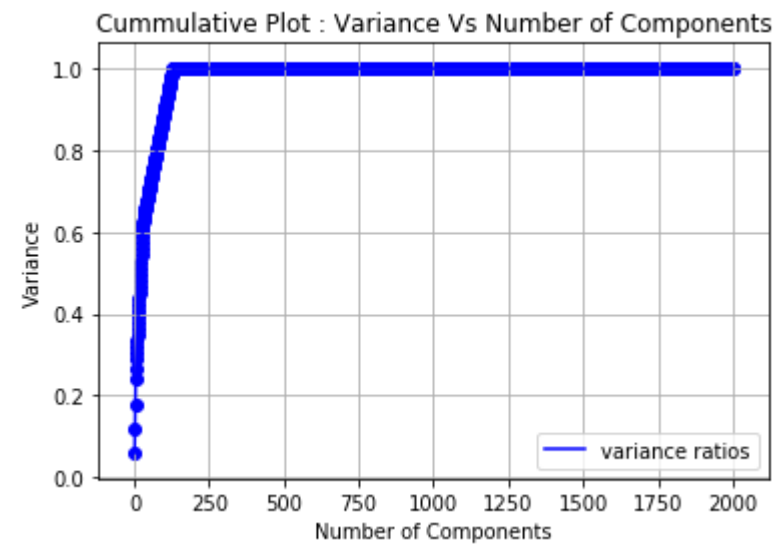
```

	ABC	PQR	DEF
ABC	0	3	3
PQR	3	0	2
DEF	3	2	0

10 rows × 2000 columns

```
Out[42]: TruncatedSVD(algorithm='randomized', n_components=1999, n_iter=5,
    random_state=None, tol=0.0)
```

```
In [43]: 1 no_components=list(range(1,df.shape[1]))
2 #percentage of variance contributed by each feature
3 percent_variance=trunc_svd.explained_variance_ / np.sum(trunc_svd.explained_variance_)
4 #cumulative plot to determine number of features required to preserve maximum variance
5 plt.plot(no_components,np.cumsum(percent_variance),color='b',label='variance ratios')
6 plt.scatter(no_components,np.cumsum(percent_variance),color='b')
7 plt.grid()
8 plt.legend()
9 plt.xlabel('Number of Components')
10 plt.ylabel('Variance')
11 plt.title('Cumulative Plot : Variance Vs Number of Components')
12 plt.show()
```



Observations :

1. We see that 250 features are preserving 99% of variance

```
In [44]: 1 ## Training with 250 features
2 trunc_svd=TruncatedSVD(n_components=250)
3 trunc_svd.fit(df)
4 co_occur_mat=trunc_svd.transform(df)
```

```
In [45]: 1 print('*'*20)
2 print('Shape of Co-occurrence matrix after applying truncated SVD: ',co_occur_mat.shape)
```

Shape of Co-occurrence matrix after applying truncated SVD: (2000, 250)

```
In [46]: 1 #vectorize(essay ,title) using final co ocuurance matrix
2 # 1. initiaize ur vector A 1 x 250 as zeros
3 # 2. word_Count=0
4 # for every word in the essay
5 # 3. if word in co_ocur_matrix :
6 # get the word vector from co_occur and add to A
7 # word ocunt ++
8 # 4.divide by word count
9
10 idx_dict=dict(zip(top_2000_words,range(0,2001)))
11 def vectorize_text(data,co_mat,col_name):
12     vectorized_text=[]
13     for sentence in tqdm(data[col_name]):
14         vector=np.zeros(250)
15         word_count=0
16         for word in sentence.split():
17             if word in top_2000_words :
18                 vector+=co_mat[idx_dict[word]]
19                 word_count+=1
20         if word_count!=0:
21             vector = vector/word_count
22         vectorized_text.append(vector)
23     return vectorized_text
```

```
In [47]: 1 vectorize_tr_essay=vectorize_text(X_train,co_occur_mat,'processed_essay')
2 vectorize_te_essay=vectorize_text(X_test,co_occur_mat,'processed_essay')
3 vectorize_cv_essay=vectorize_text(X_cv,co_occur_mat,'processed_essay')
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 26800/26800 [01:51<00:00, 241.22it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 10000/10000 [00:41<00:00, 241.18it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 13200/13200 [00:54<00:00, 240.49it/s]
```

```
In [48]: 1 vectorize_tr_title=vectorize_text(X_train,co_occur_mat,'processed_title')
2 vectorize_te_title=vectorize_text(X_test,co_occur_mat,'processed_title')
3 vectorize_cv_title=vectorize_text(X_cv,co_occur_mat,'processed_title')
```

```
100%|████████████████████████████████████████████████████████████████████████████████| 26800/26800 [00:03<00:00, 8722.78it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 10000/10000 [00:01<00:00, 8821.56it/s]
100%|████████████████████████████████████████████████████████████████████████████████| 13200/13200 [00:01<00:00, 8058.04it/s]
```

```
In [49]: 1 ### print shapes
2 print('*'*50)
3 print('Shape of vectorized essay train ',len(vectorize_tr_essay),len(vectorize_tr_essay[0]))
4 print('Shape of vectorized essay test ',len(vectorize_te_essay),len(vectorize_te_essay[0]))
5 print('Shape of vectorized essay cv ',len(vectorize_cv_essay),len(vectorize_cv_essay[0]))
```

```
*****
Shape of vectorized essay train  26800 250
Shape of vectorized essay test  10000 250
Shape of vectorized essay cv  13200 250
```

4. Apply XGBoost after stacking all the features

```
In [50]: ▶ 1 X_tr = hstack((train_categories, train_subcategories,sklstate_train,teacher_prefix_train,
2               proj_grade_train,vectorize_tr_title,vectorize_tr_essay,
3               sentiment_neg_train_norm,sentiment_pos_train_norm,
4               sentiment_neu_train_norm,sentiment_compound_train_norm,
5               X_train_price_norm,quantity_train_norm,
6               prev_projects_train_norm,title_word_count_train_norm,
7               essay_word_count_train_norm)).tocsr()
8
9 X_te = hstack((test_categories, test_subcategories,sklstate_test,teacher_prefix_test,
10              proj_grade_test,vectorize_te_title,vectorize_te_essay,
11              sentiment_neg_test_norm,sentiment_pos_test_norm,
12              sentiment_neu_test_norm,sentiment_compound_test_norm,
13              X_test_price_norm,quantity_test_norm,
14              prev_projects_test_norm,title_word_count_test_norm,
15              essay_word_count_test_norm)).tocsr()
16
17 X_cr = hstack((cv_categories, cv_subcategories,sklstate_cv,teacher_prefix_cv,
18              proj_grade_cv,vectorize_cv_title,vectorize_cv_essay
19              ,sentiment_neg_cv_norm,sentiment_pos_cv_norm,
20              sentiment_neu_cv_norm,sentiment_compound_cv_norm,
21              X_cv_price_norm,quantity_cv_norm,
22              prev_projects_cv_norm,title_word_count_cv_norm,
23              essay_word_count_cv_norm)).tocsr()
24
25
26 print(X_tr.shape)
27 print(X_te.shape)
28 print(X_cr.shape)
```

```
(26800, 608)
```

```
(10000, 608)
```

```
(13200, 608)
```

4.1 Writing own loop to find best hyper parameter using simple cross validation data

[illegible]

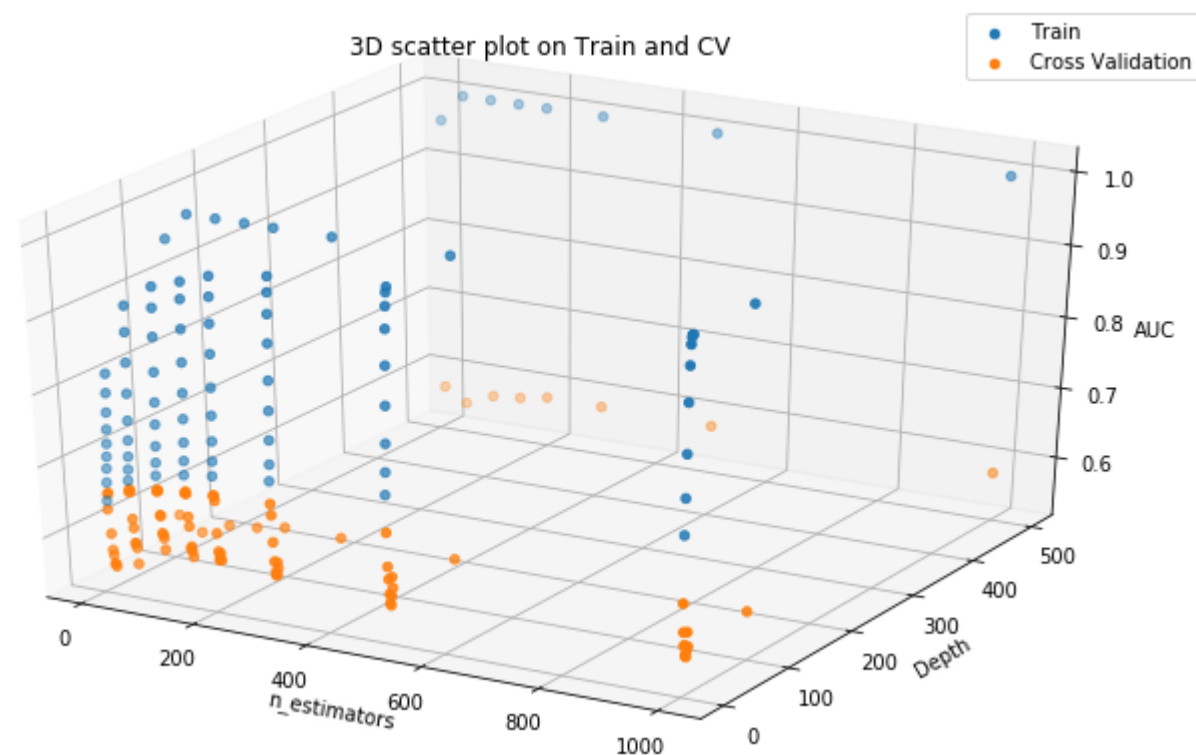
4.2 Representation using 3D scatter plot

```

In [52]: 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[52]: <matplotlib.legend.Legend at 0x257189d6518>



4.3 Finding best parameter for training the model

```

In [53]: 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

```

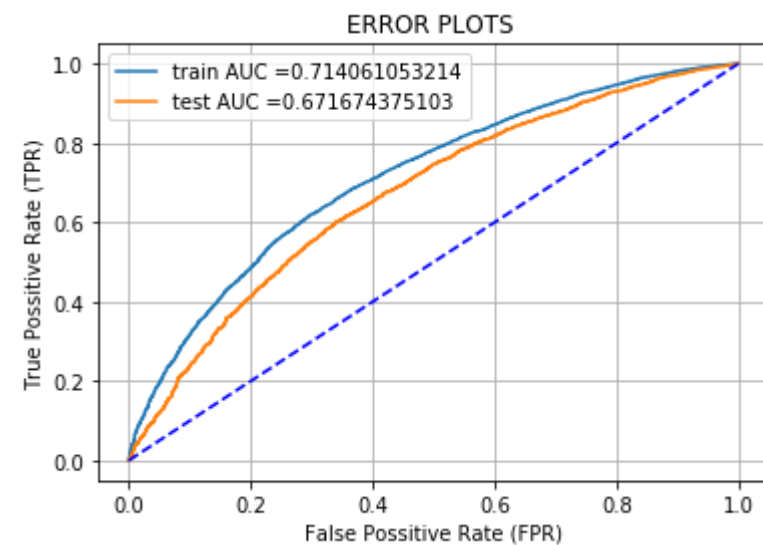
```
In [54]: 1 best_params=find_best_hyperparam(x1,y1,cv_auc)
        2 best_params
```

Max_auc is 0.682584717305

Out[54]: {'n_estimators': 300, 'depth': 2}

```
In [55]: 1 xg=xgb.XGBClassifier(learning_rate=0.05,n_estimators=best_params['n_estimators'],
        2                               max_depth=best_params['depth'],random_state=4)
        3 xg.fit(X_tr,y_train)
        4
        5 ## Predict the test
        6 train_predicts=xg.predict_proba(X_tr)[:,-1]
        7 test_predicts=xg.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 [56]: 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 Possitive Rate (FPR)")
        6 plt.ylabel("True Possitive Rate (TPR)")
        7 plt.title("ERROR PLOTS ")
        8 plt.plot([0, 1], [0, 1], 'b--')
        9 plt.grid()
       10 plt.show()
       11
```



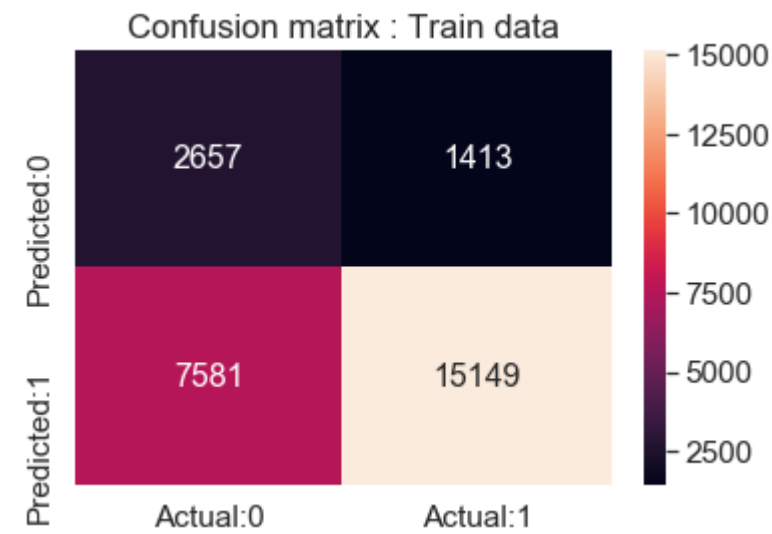
4.4 Confusion Matrix on train and test data

```
In [57]: 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
```

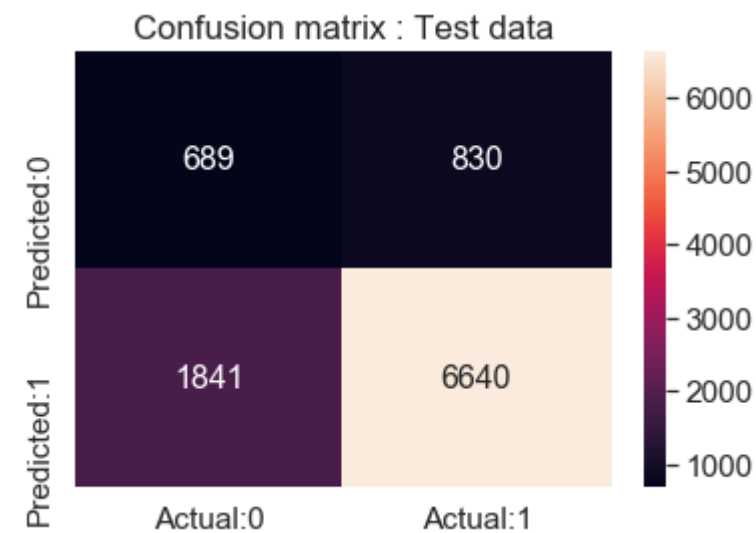
```
In [58]: 1  print("="*100)
2  best_t=best_threshold(tr_thresholds,train_fpr, train_tpr)
3  print("Train confusion matrix")
4  print(confusion_matrix(y_train, predict_with_best_t(train_predicts, best_t)))
5  print("Test confusion matrix")
6  print(confusion_matrix(y_test, predict_with_best_t(test_predicts, best_t)))
```

```
=====
the maximum value of tpr*(1-fpr) 0.43509257808 for threshold 0.84
Train confusion matrix
[[ 2657  1413]
 [ 7581 15149]]
Test confusion matrix
[[ 689  830]
 [1841 6640]]
```

```
In [59]: 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 [60]: 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()
```



Observations :

- 1.We can observe from train and test we are getting majority True positives
- 2.Least number of data falls in False negative,which refers as least number of projects were incorrectly predicted as not approved in both Test and Train.
- 3.For a model to perform well we need High True Positive Rate and Low False Positive Rate.From the above our train data has True Positive Rate as 91.4% and False Positive Rate as 74%
- 4.In our test data : True Positive Rate as 88.8% and False Positive Rate as 72.7%.

Conclusions :

- 1.Due to memory issue we sampled 50000 data and split into train,test and cv
- 2.Found top 2000 words with highest idf_ values and constructed co-occurrence matrix.
- 3.Using TruncatedSVD we found optimal number of components that preserves maximum variance of co-occurrence matrix to be 250
- 4.After applying XGBoost we found optimal number of depth to be 2 and number of estimators to be 300.
- 5.After vectorizing text data using co-occurrence matrix adds significant value to our model.
- 6.On the cv data we found max AUC to be 0.68.
- 7.After hyperparameter tuning and training model .Found train AUC to be 0.71 and test AUC to be 0.67

In [61]:

▶

1

https://xgboost.readthedocs.io/en/latest/python/python_intro.html

2

reference : <https://jakevdp.github.io/PythonDataScienceHandbook/04.12-three-dimensional-plotting.html>

In []:

▶

1