

# Gradient Boosting on Donor Choose Dataset

```
In [2]: 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
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 nltk.stem.wordnet import WordNetLemmatizer
27
28 from gensim.models import Word2Vec
29 from gensim.models import KeyedVectors
30 import pickle
31
32 from tqdm import tqdm
33 import os
34
35 from chart_studio.plotly import plotly
36 import plotly.offline as offline
37 import plotly.graph_objs as go
38 offline.init_notebook_mode()
39 from collections import Counter
```

```
In [3]: 1 project_data = pd.read_csv('train_data.csv',nrows=35000)
2 resource_data = pd.read_csv('resources.csv')
3 project_data.reset_index()
4 resource_data.reset_index()
```

```
Out[3]:
```

	index	id	description	quantity	price
0	0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95
2	2	p069063	Cory Stories: A Kid's Book About Living With Adhd	1	8.45
3	3	p069063	Dixon Ticonderoga Wood-Cased #2 HB Pencils, Bo...	2	13.59
4	4	p069063	EDUCATIONAL INSIGHTS FLUORESCENT LIGHT FILTERS...	3	24.95
5	5	p069063	Last to Finish: A Story About the Smartest Boy...	1	16.99
6	6	p069063	Mrs. Gorski, I Think I Have the Wiggle Fidgets...	1	9.95
7	7	p069063	See-N-Read 1503905CQ Reading Tool - Book Size,...	2	10.11
8	8	p096795	Brewster WPD90218 Wall Pops Flirt Dot, Set of ...	2	9.95
9	9	p096795	Brewster Wall Pops WPE99065 Peel & Stick Calyp...	2	9.02
10	10	p096795	TIME For Kids - 3-4 PRINT Bundle - 24 issues /...	40	5.01
11	11	p149007	Ahora, Spanish, Grades 6 - 12, Level 2 (min. 1...	60	7.99
12	12	p149007	Scholastic News, Grades 5/6 (min. 10 subscript...	96	5.25
13	13	p149007	Science Spin Grades 3-6 - 8 Issues / Min. 10 S...	96	0.99
14	14	p236235	PP440X - Fairy Tales Problem Solving STEM Kits	2	149.00
15	15	p052460	DD165AT - Calming Colors&#174; Easy-Clean Room...	1	129.00
16	16	p052460	DD165SB - Calming Colors&#174; Easy-Clean Room...	1	129.00
17	17	p052460	DD165SE - Calming Colors&#174; Easy-Clean Room...	1	129.00
18	18	p052460	DD165SG - Calming Colors&#174; Easy-Clean Room...	1	129.00
19	19	p233680	AA758BU - Connect & Store Book Bin - Blue	4	4.99
20	20	p233680	AA758GR - Connect & Store Book Bin - Green	4	4.99
21	21	p233680	AA758RD - Connect & Store Book Bin - Red	4	4.99
22	22	p233680	AA758RG - Connect & Store Book Bin - Orange	4	4.99
23	23	p233680	AA758VT - Connect & Store Book Bin	5	4.99
24	24	p233680	AA758YE - Connect & Store Book Bin - Yellow	5	4.99
25	25	p233680	JJ302 - Books On Wheels Mobile Library - 6 Bins	1	149.00
26	26	p233680	LX468BU - Extra Storage Bin - Blue	2	8.99
27	27	p233680	LX468GR - Extra Storage Bin - Green	2	8.99
28	28	p233680	LX468RD - Extra Storage Bin - Red	2	8.99
29	29	p233680	LX468YE - Extra Storage Bin - Yellow	2	8.99
...	...	...	...	...	...
1541242	1541242	p187432	Samsung Chromebook, 11.6" Screen, 2 GB RAM, 16...	3	202.99

	index	id	description	quantity	price
<b>1541243</b>	1541243	p187432	Sentry Folding Headphones, Black	10	7.99
<b>1541244</b>	1541244	p187432	Sentry Folding Headphones, White	4	7.99
<b>1541245</b>	1541245	p149426	Piper Computer Kit   Educational Computer that...	1	299.00
<b>1541246</b>	1541246	p238803	CARPET MY FAVORITE COLORS-7FT6INX12FT	1	314.97
<b>1541247</b>	1541247	p087783	BALL STAY N PLACE SAND FILL	2	34.07
<b>1541248</b>	1541248	p087783	BR302BU - Comfy Floor Seat - Blue	1	49.99
<b>1541249</b>	1541249	p087783	BR302RD - Comfy Floor Seat - Red	1	49.99
<b>1541250</b>	1541250	p087783	CARDINAL (PP) - CLASSROOM SELECT	3	0.00
<b>1541251</b>	1541251	p087783	CF521GR - Giant Comfy Pillow - Green	1	69.99
<b>1541252</b>	1541252	p087783	OPTION CLASS - CS NEOCLASS/NEOMOVE SHELL COLOR...	3	0.00
<b>1541253</b>	1541253	p087783	STOOL - CS NEOROK - STOOL HEIGHT 12 - RUBBER B...	3	59.47
<b>1541254</b>	1541254	p086116	Apple iPad 2 2nd generation Tablet, 1 GHz proc...	1	124.99
<b>1541255</b>	1541255	p086116	Apple iPad with Retina Display MD513LL/A (16GB...	11	367.95
<b>1541256</b>	1541256	p086116	ProCase iPad Case 9.7" 2017 - Vintage Folio St...	3	11.99
<b>1541257</b>	1541257	p086116	iPad 2 Case, iPad 3 Case, iPad 4 Case, AiSMei ...	7	10.99
<b>1541258</b>	1541258	p086116	iPad 9.7 2017 Case (New 2017 Model), EasyAcc U...	2	9.90
<b>1541259</b>	1541259	p086116	iPad Mini Case, Apple iPad Mini 2 Case, iPad M...	1	14.99
<b>1541260</b>	1541260	p228679	AA162 - First 100 Sight-Words Talking Boards	1	59.99
<b>1541261</b>	1541261	p228679	EE809 - Magnetic Fishing Poles - Set of 2	2	12.99
<b>1541262</b>	1541262	p228679	FF468 - Magnetic Sight-Word Sentence Board	1	29.99
<b>1541263</b>	1541263	p228679	TT507 - Fishing for Sight-Words - Level 1	1	21.99
<b>1541264</b>	1541264	p183340	42 PC GRADESTUFF MID SCHOOL - PACK OF 42	1	219.10
<b>1541265</b>	1541265	p183340	Rubbermaid Commercial FG9S3100GRAY Brute Tote ...	1	27.49
<b>1541266</b>	1541266	p031981	5pcs DC3V/0.1A 1.5V/0.05A 10x2.7mm Coin Mobile...	2	6.46
<b>1541267</b>	1541267	p031981	AmazonBasics 9 Volt Everyday Alkaline Batterie...	1	9.99
<b>1541268</b>	1541268	p031981	AmazonBasics AAA Performance Alkaline Batterie...	1	6.99
<b>1541269</b>	1541269	p031981	Black Electrical Tape (GIANT 3 PACK) Each Roll...	6	8.99
<b>1541270</b>	1541270	p031981	Flormoon DC Motor Mini Electric Motor 0.5-3V 1...	2	8.14
<b>1541271</b>	1541271	p031981	WAYLLSHINE 6PCS 2 x 1.5V AAA Battery Spring Cl...	2	7.39

1541272 rows × 5 columns

```
In [4]: 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)
4 project_data.project_is_approved.value_counts()
```

Number of data points in train data (35000, 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']

```
Out[4]: 1    29611
0     5389
Name: project_is_approved, dtype: int64
```

```
In [5]: 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[5]:
```

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

## 1.2 preprocessing of project\_subject\_categories

```

In [6]: 1 categories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-string
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string
7 cat_list = []
8 for i in categories:
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 spaces
13            j=j.replace('The', '') # if we have the words "The" we are going to remove it
14            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty string)
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
22 from collections import Counter
23 my_counter = Counter()
24 for word in project_data['clean_categories'].values:
25     my_counter.update(word.split())
26
27 cat_dict = dict(my_counter)
28 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))

```

```

In [7]: 1 sorted_cat_dict

```

```

Out[7]: {'Warmth': 466,
'Care_Hunger': 466,
'History_Civics': 1891,
'Music_Arts': 3280,
'AppliedLearning': 3879,
'SpecialNeeds': 4357,
'Health_Sports': 4604,
'Math_Science': 13186,
'Literacy_Language': 16763}

```

## 1.3 preprocessing of project\_subject\_subcategories

```

In [8]: 1 sub_categories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/228985
3
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string
7
8 sub_cat_list = []
9 for i in sub_categories:
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 spaces
14             j=j.replace('The', '') # if we have the words "The" we are going to remove them
15             j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty string)
16             temp +=j.strip()+" #" "abc ".strip() will return "abc", remove the trailing space
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/228985
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]))

```

```
In [9]: 1 sorted_sub_cat_dict
```

```
Out[9]: {'Economics': 100,
        'CommunityService': 157,
        'FinancialLiteracy': 177,
        'ParentInvolvement': 214,
        'Extracurricular': 250,
        'Civics_Government': 266,
        'ForeignLanguages': 271,
        'NutritionEducation': 441,
        'Warmth': 466,
        'Care_Hunger': 466,
        'SocialSciences': 592,
        'PerformingArts': 611,
        'CharacterEducation': 659,
        'TeamSports': 712,
        'Other': 792,
        'College_CareerPrep': 809,
        'Music': 969,
        'History_Geography': 1020,
        'Health_LifeScience': 1295,
        'EarlyDevelopment': 1360,
        'ESL': 1393,
        'Gym_Fitness': 1483,
        'EnvironmentalScience': 1746,
        'VisualArts': 2039,
        'Health_Wellness': 3330,
        'AppliedSciences': 3497,
        'SpecialNeeds': 4357,
        'Literature_Writing': 7064,
        'Mathematics': 8941,
        'Literacy': 10896}
```

## 1.4 preprocessing of project grade categories

```
In [10]: 1 project_data['project_grade_category'].value_counts()
```

```
Out[10]: Grades PreK-2      14199
        Grades 3-5      11888
        Grades 6-8      5415
        Grades 9-12     3498
        Name: project_grade_category, dtype: int64
```

```
In [11]: 1 preprocessed_project_grade_categories= []
        2 for grade_cat in tqdm(project_data["project_grade_category"]):
        3     grade_cat = grade_cat.replace('-', '_') #Replacing(-) with(_)
        4     grade_cat = grade_cat.replace('Grades', '') #Removing grades as it is red
        5     grad_cat = ' '.join(f for f in grade_cat.split())
        6     preprocessed_project_grade_categories.append(grad_cat.strip())
```

```
100%|██████████| 35000/35000 [00:00<00:00, 584870.10it/s]
```

```
In [12]: 1 print(preprocessed_project_grade_categories[12])
2 print("="*50)
3 print(preprocessed_project_grade_categories[502])
4 print("="*50)
5 print(preprocessed_project_grade_categories[5002])
6 print("="*50)
7 # print(preprocessed_project_grade_categories[50002])
8 # print("="*50)
9 # print(preprocessed_project_grade_categories[100013])
10 # print("="*50)
```

```
6_8
=====
6_8
=====
6_8
=====
```

## 1.5 preprocessing of teacher prefix

```
In [13]: 1 project_data['teacher_prefix'].value_counts()
```

```
Out[13]: Mrs.      18352
Ms.       12530
Mr.       3364
Teacher   752
Name: teacher_prefix, dtype: int64
```

## Replacing Nan Values with maximum frequencies values i.e Mrs.

```
In [14]: 1 project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('Mrs.')
```

```
In [15]: 1 def replace_cate(lst):          # Removing (.) in Mrs.
2         return lst.replace('.', '')
3 project_data['teacher_prefix'] = project_data['teacher_prefix'].apply(replace_cate)
```

```
In [16]: 1 preprocessed_teacher_prefix = []
2 for teach_prefix in tqdm(project_data["teacher_prefix"]):
3     preprocessed_teacher_prefix.append(teach_prefix.strip())
```

```
100%|██████████| 35000/35000 [00:00<00:00, 2064764.69it/s]
```



```
In [17]: 1 print(preprocessed_teacher_prefix[1])
2 print("="*50)
3 print(preprocessed_teacher_prefix[50])
4 print("="*50)
5 project_data.teacher_prefix.value_counts()
```

```
Mr
=====
Mrs
=====
```

```
Out[17]: Mrs      18354
Ms       12530
Mr       3364
Teacher   752
Name: teacher_prefix, dtype: int64
```

## 1.6 Adding a new feature

Introducing New Features Consider these set of features for Set 5 in Assignment:

categorical data school\_state  
clean\_categories....clean\_subcategories....project\_grade\_category....teacher\_prefix

numerical data quantity....teacher\_number\_of\_previously\_posted\_projects....price

New Features

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

```
In [18]: 1 title_word_count = []
```

```
In [19]: 1 for a in project_data["project_title"] :
2         b = len(a.split())
3         title_word_count.append(b)
```

```
In [20]: 1 project_data["title_word_count"] = title_word_count
```

In [21]: 1 project\_data.head(5)

Out[21]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN	20
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr	FL	20
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms	AZ	20
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs	KY	20
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs	TX	20

## combining 4 essays into 1 essay

In [22]:

```

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)

```

## Adding a new feature Number of words in essay

In [23]: 1 essay\_word\_count=[]

In [24]:

```

1 for ess in project_data["essay"] :
2     c = len(ess.split())
3     essay_word_count.append(c)

```

In [25]: 1 project\_data["essay\_word\_count"] = essay\_word\_count

# Computing Sentiment Score

```
In [26]: 1 import nltk
2 from nltk.sentiment.vader import SentimentIntensityAnalyzer
3 import nltk
4 nltk.download('vader_lexicon')
5 # import nltk
6 #nltk.download('vader_lexicon')
7
8 sid = SentimentIntensityAnalyzer()
9
10 for_sentiment = 'a person is a person no matter how small dr seuss i teach th
11 ss = sid.polarity_scores(for_sentiment)
12
13 for k in ss:
14     print('{0}: {1}'.format(k, ss[k]), end='')
15
16 # we can use these 4 things as features/attributes (neg, neu, pos, compound)
17 # neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
```

[nltk\_data] Downloading package vader\_lexicon to  
 [nltk\_data] C:\Users\honey\AppData\Roaming\nltk\_data...  
 [nltk\_data] Package vader\_lexicon is already up-to-date!  
 neg: 0.109, neu: 0.693, pos: 0.198, compound: 0.2023,

```
In [27]: 1 SID = SentimentIntensityAnalyzer()
2 #There is NEGATIVE and POSITIVE and NEUTRAL and COMPUND SCORES
3 #http://www.nltk.org/howto/sentiment.html
4
5 negative = []
6 positive = []
7 neutral = []
8 compound = []
9 for i in tqdm(project_data['essay']):
10     j = SID.polarity_scores(i)['neg']
11     k = SID.polarity_scores(i)['neu']
12     l = SID.polarity_scores(i)['pos']
13     m = SID.polarity_scores(i)['compound']
14     negative.append(j)
15     positive.append(k)
16     neutral.append(l)
17     compound.append(m)
```

100%|██████████| 35000/35000 [05:30<00:00, 105.75it/s]

```
In [28]: 1 project_data['negative'] = negative
2 project_data['positive'] = positive
3 project_data['neutral'] = neutral
4 project_data['compound'] = compound
```

In [29]: 1 project\_data.head(2)

Out[29]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
--	------------	----	------------	----------------	--------------	-------------

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN	20
---	--------	---------	----------------------------------	-----	----	----

1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr	FL	20
---	--------	---------	----------------------------------	----	----	----

2 rows × 24 columns

In [30]: 1 project\_data.columns.values

Out[30]: array(['Unnamed: 0', 'id', 'teacher\_id', 'teacher\_prefix', 'school\_state', 'project\_submitted\_datetime', 'project\_grade\_category', 'project\_title', 'project\_essay\_1', 'project\_essay\_2', 'project\_essay\_3', 'project\_essay\_4', 'project\_resource\_summary', 'teacher\_number\_of\_previously\_posted\_projects', 'project\_is\_approved', 'clean\_categories', 'clean\_subcategories', 'title\_word\_count', 'essay', 'essay\_word\_count', 'negative', 'positive', 'neutral', 'compound'], dtype=object)

## Train Test split

In [31]:

```

1 # train test split using sklearn.model selection
2 from sklearn.model_selection import train_test_split
3
4 X_train, X_test, y_train, y_test = train_test_split(project_data, project_data['project_is_approved'],
5 #X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=

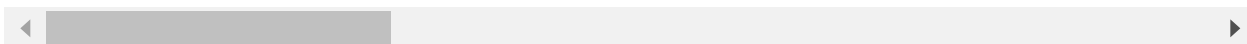
```

In [32]: 1 project\_data.head()

Out[32]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_sul
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN	20
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr	FL	20
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms	AZ	20
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs	KY	20
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs	TX	20

5 rows × 24 columns



In [33]: 1 X\_train.shape

Out[33]: (23450, 24)

In [34]:

```

1 X_train.drop(['project_is_approved'], axis=1, inplace=True)
2 X_test.drop(['project_is_approved'], axis=1, inplace=True)
3 #X_cv.drop(['project_is_approved'], axis=1, inplace=True)

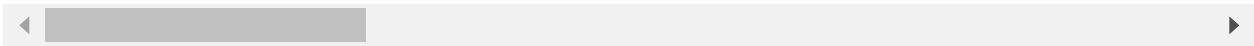
```

```
In [35]: 1 X_train.head(2)
```

Out[35]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	projec
2205	61258	p110928	da47958d7d9e9ab5550166aa5f42d500	Mrs	OH
31234	133706	p159437	d5ee4cd2ab601b71382ac1cf7094cbff	Ms	CA

2 rows × 23 columns



# 1.8 Text preprocessing

```
In [36]: 1 # printing some random reviews
2
3 print(X_train['essay'].values[0])
4 print("="*50)
5 print(X_train['essay'].values[500])
6 print("="*50)
```

This project will benefit 25, 4-5 year olds, who are first timers to the public school setting! They are excited about being in school and learning! They will be filled with lots of energy, enthusiasm, and a true love of learning! This will be a mix of both boys and girls. This group of students will include students who achieve and grow at a very rapid pace. There may be students with special needs included in this group as well. With this project, I would like to provide my students with Boogie Boards with which to practice their writing. These boards can be used to write on with a stylus or finger and then, with the touch of a button, can be erased to practice again! They will provide the students with a form of technology that is both fun and effective! They are durable and will hold the students attention. Boogie Boards will provide the students with a fun way to practice their writing for years to come! I am hoping to receive a set of 4 of these Boogie Boards so that several students can work individually or in pairs. nannan

=====

My kiddos will experience a 21st century learning environment where they get to make choices about how they learn. We are not a sit and get classroom. Traditional classroom seating is boring. I want to encourage collaboration and creative thinking, one size fits all chairs are yesterday's classroom. My students are of lower socioeconomic status and we have an extremely high percentage of free and reduced lunches. Many of them come from truly broken homes with many having a single parent to support them and many are being raised by grandparents. To give my kiddos these choices in class gives them a safe homelike feeling that they may not otherwise have. My kiddos will be using these as a seating option anytime they need to sit and complete a task, whether it be independent, partners or small groups. It will be one of the many flexible seating options that I hope to offer to my Kinders. They will have the freedoms to use these at a low table or to use them to lounge with a clipboard. They need to be free to determine what seating options work for them! Seeing a child truly engaged in their learning and making choices about the way they learn best. This helps them to learn more about who they are at a young age. Many behavior problems are virtually eliminated and they are allowed to move and learn. They are able to retain more and take an active role in their learning at a very young age. I hope to be able to get many more seating options for my kiddos to encourage and inspire them to do and be their very best in everything they do!

=====

```
In [37]: 1 # https://stackoverflow.com/a/47091490/4084039
2 import re
3
4 def decontracted(phrase):
5     # specific
6     phrase = re.sub(r"won't", "will not", phrase)
7     phrase = re.sub(r"can't", "can not", phrase)
8
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase
```

```
In [38]: 1 sent = decontracted(X_train['essay'].values[20000])
2 print(sent)
3 print("="*50)
```

In my class I normally read a shared text with the students, we work on various writing pieces, and they are given time to read and write independently. When reading a shared text, we often have class discussions and Socratic Seminars. When writing, we have workshops and conferences in small groups. My students are in the 7th and 8th grade. I have students from many different races and religions. As a result, I promote and model tolerance. We often read about history, and learn about novels in context, because I believe it is very important for students gain deeper awareness of cultural context whenever we read a new text. For example, when we recently read Night by Elie Wiesel, the students also gained a deeper understanding of WWII and the Holocaust, so that they might be able to better understand and empathize with Wiesel's memoir. We also recently read To Kill a Mockingbird, by Harper Lee, and the students learned more about the Great Depression and Jim Crow South in order to better understand the plot and themes within the text. I believe that students need to be made aware of history in order to better understand not only what they read, but the world around them today, and even themselves. Most of my students do enjoy reading independently whether they are in my Integrated Co-Teaching class, or Honors class. However, I am in need of a greater variety of books for them to choose from. I am also in need of a document camera so that I might aid my students with their writing, through modeling. This will also allow students to share their work with one another, and receive feedback. I think rejuvenating my classroom library with new books will give my students a greater interest and choice in what they read. I think that acquiring the document camera will allow me to better help my students with their writing because I will be able to model writing for them, show them how to revise and edit more easily, and allow them to provide feedback to one another by sharing their work on the board.

=====



```
In [39]: 1 # \r \n \t remove from string python: http://texthandler.com/info/remove-line
2 sent = sent.replace('\r', ' ')
3 sent = sent.replace('\n', ' ')
4 sent = sent.replace('\t', ' ')
5 print(sent)
```

In my class I normally read a shared text with the students, we work on various writing pieces, and they are given time to read and write independently. When reading a shared text, we often have class discussions and Socratic Seminars. When writing, we have workshops and conferences in small groups. My students are in the 7th and 8th grade. I have students from many different races and religions. As a result, I promote and model tolerance. We often read about history, and learn about novels in context, because I believe it is very important for students gain deeper awareness of cultural context whenever we read a new text. For example, when we recently read *Night* by Elie Wiesel, the students also gained a deeper understanding of WWII and the Holocaust, so that they might be able to better understand and empathize with Wiesel's memoir. We also recently read *To Kill a Mockingbird*, by Harper Lee, and the students learned more about the Great Depression and Jim Crow South in order to better understand the plot and themes within the text. I believe that students need to be made aware of history in order to better understand not only what they read, but the world around them today, and even themselves. Most of my students do enjoy reading independently whether they are in my Integrated Co-Teaching class, or Honors class. However, I am in need of a greater variety of books for them to choose from. I am also in need of a document camera so that I might aid my students with their writing, through modeling. This will also allow students to share their work with one another, and receive feedback. I think rejuvenating my classroom library with new books will give my students a greater interest and choice in what they read. I think that acquiring the document camera will allow me to better help my students with their writing because I will be able to model writing for them, show them how to revise and edit more easily, and allow them to provide feedback to another by sharing their work on the board.

```
In [40]: 1 #remove spacial character: https://stackoverflow.com/a/5843547/4084039
2 sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
3 print(sent)
```

In my class I normally read a shared text with the students we work on various writing pieces and they are given time to read and write independently When reading a shared text we often have class discussions and Socratic Seminars When writing we have workshops and conferences in small groups My students are in the 7th and 8th grade I have students from many different races and religions As a result I promote and model tolerance We often read about history and learn about novels in context because I believe it is very important for students gain deeper awareness of cultural context whenever we read a new text For example when we recently read Night by Elie Wiesel the students also gained a deeper understanding of WWII and the Holocaust so that they might be able to better understand and empathize with Wiesel is memoir We also recently read To Kill a Mockingbird by Harper Lee and the students learned more about the Great Depression and Jim Crow South in order to better understand the plot and themes within the text I believe that students need to be made aware of history in order to better understand not only what they read but the world around them today and even themselves Most of my students do enjoy reading independently whether they are in my Integrated Co Teaching class or Honors class However I am in need of a greater variety of books for them to choose from I am also in need of a document camera so that I might aid my students with their writing through modeling This will also allow students to share their work with one another and receive feedback I think rejuvenating my classroom library with new books will give my students a greater interest and choice in what they read I think that acquiring the document camera will allow me to better help my students with their writing because I will be able to model writing for them show them how to revise and edit more easily and allow them to provide feedback to another by sharing their work on the board

```
In [41]: 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',
4             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
5             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'i',
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has',
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through',
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off',
11            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all',
12            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too',
13            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "shouldn't",
14            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn',
15            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'mustn't',
16            'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'won',
17            "won't", 'wouldn', "wouldn't"]
```

## 1.8.1 Preprocessed training data - Text

```
In [42]: 1 from tqdm import tqdm
2 def preprocess_textual(row_value):
3     preprocessed_train = []
4     for sentences in tqdm(row_value):
5         sent = decontracted(sentences)
6         sent = sent.replace('\\r', ' ')
7         sent = sent.replace('\\\"', ' ')
8         sent = sent.replace('\\n', ' ')
9         sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
10        # https://gist.github.com/sebleier/554280
11        sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
12        preprocessed_train.append(sent.lower().strip())
13    return preprocessed_train
```

```
In [43]: 1 preprocessed_essays_train=preprocess_textual(X_train['essay'].values)
```

```
100%|██████████| 23450/23450 [00:11<00:00, 2073.20it/s]
```

```
In [44]: 1 preprocessed_essays_test=preprocess_textual(X_test['essay'].values)
```

```
100%|██████████| 11550/11550 [00:05<00:00, 2078.18it/s]
```

```
In [45]: 1 preprocessed_essays_test
```

```
Out[45]: ['want introduce students world beyond know give opportunity dream big create
better future students difficult time paying attention class concerned sleepi
ng night next meal coming sometimes difficult pay attention kinds issues desp
ite high mobility extreme poverty many issues school kids active participants
learning excited classroom strive make learning environment one sanctuary kid
s forget problems focus education using donors choose get students great lite
rature students love things received lead great classroom book discussions ph
enomenal new vocabulary adventures make believe land far away places built gr
eat comprehension skills however classroom library lacking books students wan
t read say success breeds success want hep students successful possible stude
nts started book wish list classrom white board writing book titles books wan
t see classroom library kids get check books library school complain turn fin
ished reading book books want already checked project allow students read boo
k cover cover without return library chance complete books students asked nan
nan',
'students hard working bright young individuals attend school every day read
y learn meet high expectations students love challenges difficulties face eve
ry day smiles faces no less students come title school high population englis
h learners come low income household limited resources experiences many stude
...']
```

## 1.8.3 Preprocessed cross validation data

```
In [46]: 1 #preprocessed_essays_cv=preprocess_textual(X_cv['essay'].values)
```

## 1.9 preprocessing of project title

```
In [47]: 1 # printing some random project titles.
2 print(project_data['project_title'].values[1])
3 print("="*50)
4 print(project_data['project_title'].values[1501])
5 print("="*50)
6 print(project_data['project_title'].values[10001])
7 print("="*50)
8 print(project_data['project_title'].values[20001])
9 print("="*50)
```

Wanted: Projector for Hungry Learners

=====

Making Every Day at School Count

=====

Becoming 'Readerly' Readers

=====

The Beautiful Life of a Butterfly

=====

```
In [48]: 1 title = decontracted(X_train['project_title'].values[2000])
```

## 1.9.1 Preprocessing of Project Title(Train)

```
In [49]: 1 preprocessed_titles_train=preprocess_textual(X_train["project_title"].values)
100%|██████████| 23450/23450 [00:00<00:00, 49289.76it/s]
```

## 1.9.2 Preprocessing of Project Title(Test)

```
In [50]: 1 preprocessed_titles_test=preprocess_textual(X_test["project_title"])
100%|██████████| 11550/11550 [00:00<00:00, 45842.64it/s]
```

```
In [51]: 1 preprocessed_titles_test[10]
```

Out[51]: 'leap reading new leapfrog tag reader pens'

## 1.9.2 Preprocessing of Project Title(CV)

```
In [52]: 1 #preprocessed_titles_cv=preprocess_textual(X_cv["project_title"])
```

```
In [53]: 1 #preprocessed_titles_cv[10]
```

## 1.5 Preparing data for models

```
In [54]: 1 X_train.shape
```

```
Out[54]: (23450, 23)
```

```
In [55]: 1 project_data.columns
```

```
Out[55]: Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
               'project_submitted_datetime', 'project_grade_category', 'project_title',  
               'project_essay_1', 'project_essay_2', 'project_essay_3',  
               'project_essay_4', 'project_resource_summary',  
               'teacher_number_of_previously_posted_projects', 'project_is_approved',  
               'clean_categories', 'clean_subcategories', 'title_word_count', 'essay',  
               'essay_word_count', 'negative', 'positive', 'neutral', 'compound'],  
              dtype='object')
```

we are going to consider

- school\_state : categorical data
- clean\_categories : categorical data
- clean\_subcategories : categorical data
- project\_grade\_category : categorical data
- teacher\_prefix : categorical data
  
- project\_title : text data
- text : text data
- project\_resource\_summary: text data (optinal)
  
- quantity : numerical (optinal)
- teacher\_number\_of\_previously\_posted\_projects : numerical
- price : numerical

In [56]:

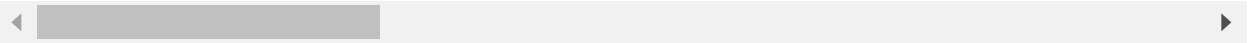
1 project\_data.head(10)

Out[56]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_su
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs	IN	20
1	140945	p258326	897464ce9ddc600bcd1151f324dd63a	Mr	FL	20
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms	AZ	20
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs	KY	20
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs	TX	20
5	141660	p154343	a50a390e8327a95b77b9e495b58b9a6e	Mrs	FL	20
6	21147	p099819	9b40170bfa65e399981717ee8731efc3	Mrs	CT	20
7	94142	p092424	5bfd3d12fae3d2fe88684bbac570c9d2	Ms	GA	20
8	112489	p045029	487448f5226005d08d36bdd75f095b31	Mrs	SC	20

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	project_su
9	158561 p001713	140eeac1885c820ad5592a409a3a8994	Ms	NC	21

10 rows × 24 columns



## 1.5.1 Vectorizing Categorical data

```
In [57]: 1 def Responsetable(table, col) :
2         cat = table[col].unique()
3         freq_Pos = []
4         for i in cat :
5             freq_Pos.append(len(table.loc[(table[col] == i) & (table['project_is_
6
7
8         freq_Neg = []
9         for i in cat :
10            freq_Neg.append(len(table.loc[(table[col] == i) & (table['project_is_
11
12            encoded_Pos = []
13            for i in range(len(cat)) :
14                encoded_Pos.append(freq_Pos[i]/(freq_Pos[i] + freq_Neg[i]))
15
16            encoded_Neg = []
17            encoded_Neg[:] = [1 - x for x in encoded_Pos]
18
19            encoded_Pos_val = dict(zip(cat, encoded_Pos))
20            encoded_Neg_val = dict(zip(cat, encoded_Neg))
21
22            return encoded_Pos_val, encoded_Neg_val
```

```

In [58]: 1 def Responsecode(train,test) :
2         pos_cleancat, neg_cleancat = Responsetable(train,'clean_categories')
3         pos_cleansubcat, neg_cleansubcat = Responsetable(train,'clean_subcategories')
4         pos_schoolstate, neg_schoolstate = Responsetable(train, 'school_state')
5         pos_teacherprefix, neg_teacherprefix = Responsetable(train, 'teacher_prefix')
6         pos_projgradecat, neg_projgradecat = Responsetable(train, 'project_grade_category')
7         dftrain = pd.DataFrame()
8         dftrain['clean_cat_pos'] = train['clean_categories'].map(pos_cleancat)
9         dftrain['clean_cat_neg'] = train['clean_categories'].map(neg_cleancat)
10        dftrain['clean_subcat_pos'] = train['clean_subcategories'].map(pos_cleansubcat)
11        dftrain['clean_subcat_neg'] = train['clean_subcategories'].map(neg_cleansubcat)
12        dftrain['school_state_pos'] = train['school_state'].map(pos_schoolstate)
13        dftrain['school_state_neg'] = train['school_state'].map(neg_schoolstate)
14        dftrain['teacher_prefix_pos'] = train['teacher_prefix'].map(pos_teacherprefix)
15        dftrain['teacher_prefix_neg'] = train['teacher_prefix'].map(neg_teacherprefix)
16        dftrain['proj_grade_cat_pos'] = train['project_grade_category'].map(pos_projgradecat)
17        dftrain['proj_grade_cat_neg'] = train['project_grade_category'].map(neg_projgradecat)
18        dftrain = pd.DataFrame()
19        dftrain['clean_cat_pos'] = train['clean_categories'].map(pos_cleancat).fillna(0)
20        dftrain['clean_cat_neg'] = train['clean_categories'].map(neg_cleancat).fillna(0)
21        dftrain['clean_subcat_pos'] = train['clean_subcategories'].map(pos_cleansubcat).fillna(0)
22        dftrain['clean_subcat_neg'] = train['clean_subcategories'].map(neg_cleansubcat).fillna(0)
23        dftrain['school_state_pos'] = train['school_state'].map(pos_schoolstate).fillna(0)
24        dftrain['school_state_neg'] = train['school_state'].map(neg_schoolstate).fillna(0)
25        dftrain['teacher_prefix_pos'] = train['teacher_prefix'].map(pos_teacherprefix).fillna(0)
26        dftrain['teacher_prefix_neg'] = train['teacher_prefix'].map(neg_teacherprefix).fillna(0)
27        dftrain['proj_grade_cat_pos'] = train['project_grade_category'].map(pos_projgradecat).fillna(0)
28        dftrain['proj_grade_cat_neg'] = train['project_grade_category'].map(neg_projgradecat).fillna(0)
29        return dftrain

```

```

In [59]: 1 X_new_train=X_train
2         X_new_train['project_is_approved']=y_train.to_frame()['project_is_approved']
3         X_new_test=X_test
4         X_new_test['project_is_approved']=y_test.to_frame()['project_is_approved']
5         newTrain,newTest = Responsecode(X_new_train,X_new_test)

```

```

In [60]: 1 print(newTrain.shape,newTest.shape)

```

```

(23450, 10) (11550, 10)

```

```

In [61]: 1 def mergeEncoding(table, p, n) :
2         lstPos = table[p].values.tolist()
3         lstNeg = table[n].values.tolist()
4         frame = pd.DataFrame(list(zip(lstNeg, lstPos)))
5         return frame

```



```
In [62]: 1 #Clean Categories
2 X_train_clean_cat_ohe = mergeEncoding(newTrain, 'clean_cat_pos', 'clean_cat_neg')
3 X_test_clean_cat_ohe = mergeEncoding(newTest, 'clean_cat_pos', 'clean_cat_neg')
4 print(X_train_clean_cat_ohe.shape)
5 X_test_clean_cat_ohe.shape
```

(23450, 2)

Out[62]: (11550, 2)

```
In [63]: 1 #Clean SUB Categories
2 X_train_clean_subcat_ohe = mergeEncoding(newTrain, 'clean_subcat_pos', 'clean_subcat_neg')
3 X_test_clean_subcat_ohe = mergeEncoding(newTest, 'clean_subcat_pos', 'clean_subcat_neg')
4 print(X_train_clean_subcat_ohe.shape)
5 print(X_test_clean_subcat_ohe.shape)
```

(23450, 2)

(11550, 2)

```
In [64]: 1 #Project Grade Category
2 X_train_grade_ohe = mergeEncoding(newTrain, 'proj_grade_cat_pos', 'proj_grade_cat_neg')
3 X_test_grade_ohe = mergeEncoding(newTest, 'proj_grade_cat_pos', 'proj_grade_cat_neg')
4 print(X_train_grade_ohe.shape)
5 print(X_test_grade_ohe.shape)
```

(23450, 2)

(11550, 2)

```
In [65]: 1 #School State
2 X_train_state_ohe = mergeEncoding(newTrain, 'school_state_pos', 'school_state_neg')
3 X_test_state_ohe = mergeEncoding(newTest, 'school_state_pos', 'school_state_neg')
4 print(X_train_state_ohe.shape)
5 print(X_test_state_ohe.shape)
```

(23450, 2)

(11550, 2)

```
In [66]: 1 #Teacher Prefix
2 X_train_teacher_ohe = mergeEncoding(newTrain, 'teacher_prefix_pos', 'teacher_prefix_neg')
3 X_test_teacher_ohe = mergeEncoding(newTest, 'teacher_prefix_pos', 'teacher_prefix_neg')
4 print(X_train_teacher_ohe.shape)
5 print(X_test_teacher_ohe.shape)
```

(23450, 2)

(11550, 2)

## TFIDF

**tfidf(train essays)**

```
In [67]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vectorizer_tfidf_essay = TfidfVectorizer(min_df=10,max_features=5000) #Consid
4 vectorizer_tfidf_essay.fit(preprocessed_essays_train)
5
6 text_tfidf_train = vectorizer_tfidf_essay.transform(preprocessed_essays_train)
7 print("Shape of matrix after TfidfVectorizer ",text_tfidf_train.shape)
```

Shape of matrix after TfidfVectorizer (23450, 5000)

### **tfidf(test essays)**

```
In [68]: 1 text_tfidf_test = vectorizer_tfidf_essay.transform(preprocessed_essays_test)
2 print("Shape of matrix after TfidfVectorizer ",text_tfidf_test.shape)
```

Shape of matrix after TfidfVectorizer (11550, 5000)

### **tfidf(cv essays)**

```
In [69]: 1 #text_tfidf_cv = vectorizer_tfidf_essay.transform(preprocessed_essays_cv)
2 #print("Shape of matrix after TfidfVectorizer ",text_tfidf_cv.shape)
```

### **tfidf(train Titles)**

```
In [70]: 1 vectorizer_tfidf_titles = TfidfVectorizer(min_df=10)
2
3 vectorizer_tfidf_titles.fit(preprocessed_titles_train)
4 title_tfidf_train = vectorizer_tfidf_titles.transform(preprocessed_titles_train)
5 print("Shape of matrix after TfidfVectorizer ",title_tfidf_train.shape)
```

Shape of matrix after TfidfVectorizer (23450, 1217)

### **tfidf(test titles)**

```
In [71]: 1 title_tfidf_test = vectorizer_tfidf_titles.transform(preprocessed_titles_test)
2 print("Shape of matrix after TfidfVectorizer ",title_tfidf_test.shape)
```

Shape of matrix after TfidfVectorizer (11550, 1217)

### **tfidf(cv titles)**

```
In [72]: 1 #title_tfidf_cv = vectorizer_tfidf_titles.transform(preprocessed_titles_cv)
2 #print("Shape of matrix after TfidfVectorizer ",title_tfidf_cv.shape)
```

## 1.5.2.3 Using Pretrained Models: TFIDF Weighted W2V

In [73]:

```

1
2 # stronging variables into pickle files python: http://www.jessicayung.com/ho
3 # make sure you have the glove_vectors file
4 with open('glove_vectors', 'rb') as f:
5     model = pickle.load(f)
6     glove_words = set(model.keys())

```

### train essays

In [74]:

```

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

```

In [75]:

```

1 def tfidfWV(preprocessed_data):
2     # average Word2Vec
3     # compute average word2vec for each review.
4     tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored
5     for sentence in tqdm(preprocessed_data): # for each review/sentence
6         vector = np.zeros(300) # as word vectors are of zero length
7         tf_idf_weight = 0; # num of words with a valid vector in the sentence/
8         for word in sentence.split(): # for each word in a review/sentence
9             if (word in glove_words) and (word in tfidf_words):
10                 vec = model[word] # getting the vector for each word
11                 # here we are multiplying idf value(dictionary[word]) and the
12                 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.
13                 vector += (vec * tf_idf) # calculating tfidf weighted w2v
14                 tf_idf_weight += tf_idf
15             if tf_idf_weight != 0:
16                 vector /= tf_idf_weight
17             tfidf_w2v_vectors.append(vector)
18     print("The length of TFIDF word to vec is ", len(tfidf_w2v_vectors))
19     print("The length of TFIDF word to vec of index 0 is ", len(tfidf_w2v_vect
20     return tfidf_w2v_vectors

```

In [76]:

```
1 tfidf_w2v_vectors_train=tfidfWV(preprocessed_essays_train)
```

100%|██████████| 23450/23450 [00:35<00:00, 656.44it/s]

The length of TFIDF word to vec is 23450

The length of TFIDF word to vec of index 0 is 300

### Test essays

```
In [77]: 1 tfidf_w2v_vectors_test=tfidfWV(preprocessed_essays_test)
```

100%|██████████| 11550/11550 [00:26<00:00, 431.10it/s]

The length of TFIDF word to vec is 11550

The length of TFIDF word to vec of index 0 is 300

## cv essays

```
In [78]: 1 #tfidf_w2v_vectors_cv = tfidfWV(preprocessed_essays_cv)
```

## train titles

```
In [79]: 1 tfidf_model = TfidfVectorizer()
2 tfidf_model.fit(preprocessed_titles_train)
3 # we are converting a dictionary with word as a key, and the idf as a value
4 dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_))
5 tfidf_words = set(tfidf_model.get_feature_names())
```

```
In [80]: 1 tfidf_w2v_vectors_titles_train = tfidfWV(preprocessed_titles_train)
```

100%|██████████| 23450/23450 [00:00<00:00, 28506.97it/s]

The length of TFIDF word to vec is 23450

The length of TFIDF word to vec of index 0 is 300

## test titles

```
In [81]: 1 tfidf_w2v_vectors_titles_test = tfidfWV(preprocessed_titles_test)
```

100%|██████████| 11550/11550 [00:00<00:00, 27768.46it/s]

The length of TFIDF word to vec is 11550

The length of TFIDF word to vec of index 0 is 300

## cv titles

```
In [82]: 1 #tfidf_w2v_vectors_titles_cv = tfidfWV(preprocessed_titles_cv)
```

# 1.12 Vectorizing Numerical features

Various numerical features are :

- 1.Price
- 2.Quantity
- 3.Number of Projects previously proposed by Teacher
- 4.Title word Count ( introduced by us)
- 5.Essay word Count ( introduced by us)
- 6.Negative Intensity
- 7.Positive Intensity
- 8.Neutral Intensity

```
In [83]: 1 # https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-index  
2 price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'})  
3 price_data.head(4)
```

```
Out[83]:
```

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21
2	p000003	298.97	4
3	p000004	1113.69	98

```
In [84]: 1 price_data.head(4)
```

```
Out[84]:
```

	id	price	quantity
0	p000001	459.56	7
1	p000002	515.89	21
2	p000003	298.97	4
3	p000004	1113.69	98

In [85]: 1 X\_train

Out[85]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	price
2205	61258	p110928	da47958d7d9e9ab5550166aa5f42d500	Mrs	OH	16500
31234	133706	p159437	d5ee4cd2ab601b71382ac1cf7094cbff	Ms	CA	15000
7957	89555	p090465	5c0c521aa9aa0cf31e32e4fa1fbb3deb	Ms	CA	15000
9891	49379	p038756	17561728716402d8cf0c380c3c358813	Mrs	WI	16500

```
In [86]: 1 # join two dataframes in python:
2 X_train = pd.merge(X_train, price_data, on='id')
3 X_test = pd.merge(X_test, price_data, on='id')
4 #X_cv = pd.merge(X_cv, price_data, on='id', how='left')
```

In [87]: 1 X\_train

Out[87]:

	Unnamed: 0	id	teacher_id	teacher_prefix	school_state	price
0	61258	p110928	da47958d7d9e9ab5550166aa5f42d500	Mrs	OH	16500
1	133706	p159437	d5ee4cd2ab601b71382ac1cf7094cbff	Ms	CA	15000
2	89555	p090465	5c0c521aa9aa0cf31e32e4fa1fbb3deb	Ms	CA	15000
3	49379	p038756	17561728716402d8cf0c380c3c358813	Mrs	WI	16500

```

In [88]: 1 from sklearn.preprocessing import Normalizer
2
3 normalizer = Normalizer()
4
5 # normalizer.fit(X_train['price'].values)
6 # this will rise an error Expected 2D array, got 1D array instead:
7 # array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
8 # Reshape your data either using
9 # array.reshape(-1, 1) if your data has a single feature
10 # array.reshape(1, -1) if it contains a single sample.
11
12 price_for_test=X_test['price'] #Box Plot Purpose
13
14 normalizer.fit(X_train['price'].values.reshape(1,-1))
15
16 price_train = normalizer.transform(X_train['price'].values.reshape(-1,1))
17 #price_cv = normalizer.transform(X_cv['price'].values.reshape(-1,1))
18 price_test = normalizer.transform(X_test['price'].values.reshape(-1,1))
19
20 print("After vectorizations")
21 print(price_train.shape, y_train.shape)
22 #print(price_cv.shape, y_cv.shape)
23 print(price_test.shape, y_test.shape)
24 print("=="*100)

```

After vectorizations

(23450, 1) (23450,)

(11550, 1) (11550,)

=====  
=====

## 2) Quantity

```
In [89]: 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 = normalizer.transform(X_train['quantity'].values.reshape(-1,1))
13 #quantity_cv = normalizer.transform(X_cv['quantity'].values.reshape(-1,1))
14 quantity_test = normalizer.transform(X_test['quantity'].values.reshape(-1,1))
15
16 print("After vectorizations")
17 print(quantity_train.shape, y_train.shape)
18 #print(quantity_cv.shape, y_cv.shape)
19 print(quantity_test.shape, y_test.shape)
20 print("="*100)
```

After vectorizations

(23450, 1) (23450,)

(11550, 1) (11550,)

=====

### 3) Number of Projects previously proposed by Teacher



```

In [90]: 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)
11
12 prev_projects_train = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values)
13 #prev_projects_cv = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values)
14 number_of_previously_posted_projects_dtrees=X_test['teacher_number_of_previously_posted_projects'].values
15 prev_projects_test = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values)
16
17 print("After vectorizations")
18 print(prev_projects_train.shape, y_train.shape)
19 #print(prev_projects_cv.shape, y_cv.shape)
20 print(prev_projects_test.shape, y_test.shape)
21 print("=="*100)

```

After vectorizations

(23450, 1) (23450,)

(11550, 1) (11550,)

=====

## 4) title word count

```

In [91]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['title_word_count'].values.reshape(1,-1))
4
5 title_word_count_train = normalizer.transform(X_train['title_word_count'].values)
6 #title_word_count_cv = normalizer.transform(X_cv['title_word_count'].values)
7 title_word_count_test = normalizer.transform(X_test['title_word_count'].values)
8
9 print("After vectorizations")
10 print(title_word_count_train.shape, y_train.shape)
11 #print(title_word_count_cv.shape, y_cv.shape)
12 print(title_word_count_test.shape, y_test.shape)
13 print("=="*100)

```

After vectorizations

(23450, 1) (23450,)

(11550, 1) (11550,)

=====

## 5) essay word count

```
In [92]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['essay_word_count'].values.reshape(1,-1))
4
5 essay_word_count_train = normalizer.transform(X_train['essay_word_count'].values.reshape(1,-1))
6 #essay_word_count_cv = normalizer.transform(X_cv['essay_word_count'].values.reshape(1,-1))
7 essay_word_count_test = normalizer.transform(X_test['essay_word_count'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(essay_word_count_train.shape, y_train.shape)
11 #print(essay_word_count_cv.shape, y_cv.shape)
12 print(essay_word_count_test.shape, y_test.shape)
```

After vectorizations  
(23450, 1) (23450,)  
(11550, 1) (11550,)

## 6) Positive Intensity

```
In [93]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['positive'].values.reshape(1,-1))
4
5 positive_intensity_train = normalizer.transform(X_train['positive'].values.reshape(1,-1))
6 #positive_intensity_cv = normalizer.transform(X_cv['positive'].values.reshape(1,-1))
7 positive_intensity_test = normalizer.transform(X_test['positive'].values.reshape(1,-1))
8
9 print("After vectorizations")
10 print(positive_intensity_train.shape, y_train.shape)
11 #print(positive_intensity_cv.shape, y_cv.shape)
12 print(positive_intensity_test.shape, y_test.shape)
```

After vectorizations  
(23450, 1) (23450,)  
(11550, 1) (11550,)

## 7) Negative Intensity

```
In [94]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['negative'].values.reshape(1,-1))
4
5 negative_intensity_train = normalizer.transform(X_train['negative'].values.re
6 #negative_intensity_cv = normalizer.transform(X_cv['negative'].values.reshape
7 negative_intensity_test = normalizer.transform(X_test['negative'].values.res
8
9 print("After vectorizations")
10 print(negative_intensity_train.shape, y_train.shape)
11 #print(negative_intensity_cv.shape, y_cv.shape)
12 print(negative_intensity_test.shape, y_test.shape)
```

After vectorizations  
(23450, 1) (23450,)  
(11550, 1) (11550,)

## 8) Neutral Intensity

```
In [95]: 1 normalizer = Normalizer()
2
3 normalizer.fit(X_train['negative'].values.reshape(1,-1))
4
5 neutral_intensity_train = normalizer.transform(X_train['neutral'].values.res
6 #neutral_intensity_cv = normalizer.transform(X_cv['neutral'].values.reshape(-
7 neutral_intensity_test = normalizer.transform(X_test['neutral'].values.res
8
9 print("After vectorizations")
10 print(neutral_intensity_train.shape, y_train.shape)
11 #print(neutral_intensity_cv.shape, y_cv.shape)
12 print(neutral_intensity_test.shape, y_test.shape)
```

After vectorizations  
(23450, 1) (23450,)  
(11550, 1) (11550,)

### 1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

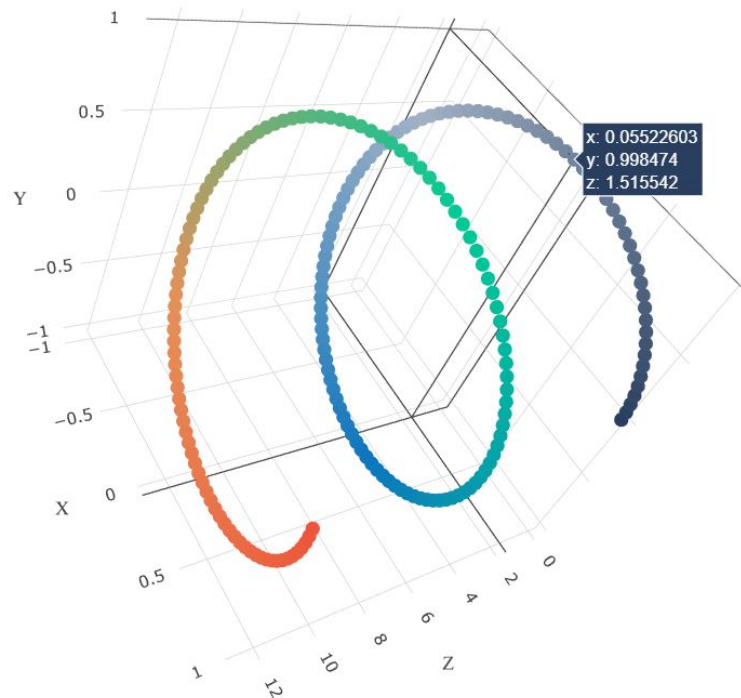
- **Set 1:** categorical, numerical features + preprocessed\_essay (TFIDF) + Sentiment scores(preprocessed\_essay)
- **Set 2:** categorical, numerical features + preprocessed\_essay (TFIDF W2V) + Sentiment scores(preprocessed\_essay)

### 2. The hyper paramter tuning (best depth in range [1, 3, 10, 30], and the best min\_samples\_split in range [5, 10, 100, 500])

- 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/) (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(use gridsearch cv or randomsearch cv)/simple cross validation data(you can write your own for loops refer sample solution)

### 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 **min\_sample\_split**, 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

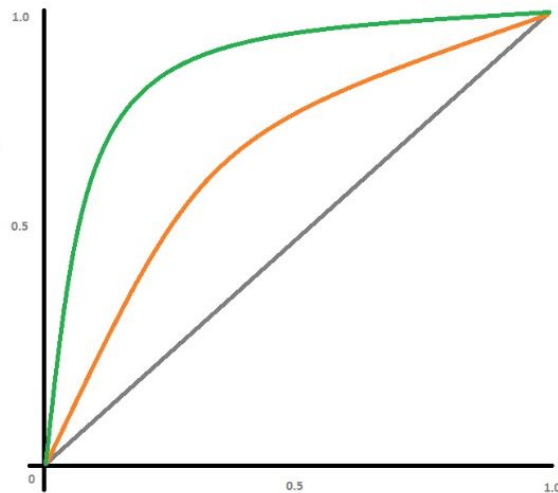
- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure



seaborn heat maps (<https://seaborn.pydata.org/generated/seaborn.heatmap.html>) with rows as **min\_sample\_split**, columns as **max\_depth**, and values inside the cell representing **AUC Score**

- You choose either of the plotting techniques out of 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. Make sure that you are using predict\_proba method to calculate AUC curves, because AUC is calculated on class probabilities and not on class labels.



- 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/) (<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 = ??

- Once after you plot the confusion matrix with the test data, get all the false positive data points
  - Plot the WordCloud(<https://www.geeksforgeeks.org/generating-word-cloud-python/> (<https://www.geeksforgeeks.org/generating-word-cloud-python/>)) with the words of essay text of these false positive data points
  - Plot the box plot with the price of these false positive data points
  - Plot the pdf with the teacher\_number\_of\_previously\_posted\_projects of these false positive data points

```
In [96]: 1 from scipy.sparse import hstack
2 X_train = hstack((X_train_clean_cat_ohe , X_train_clean_subcat_ohe, X_train_s
3                 X_train_teacher_ohe ,text_tfidf_train,title_tfidf_train, qu
4                 prev_projects_train, price_train,title_word_count_train,
5                 essay_word_count_train, positive_intensity_train, negative_
6 print(X_train.shape, y_train.shape)
7 print(type(X_train))
```

```
(23450, 6235) (23450,)
<class 'scipy.sparse.csr.csr_matrix'>
```

```
In [97]: 1 print(X_train_clean_cat_ohe.shape)
2 print(neutral_intensity_train.shape)
```

```
(23450, 2)
(23450, 1)
```

```
In [98]: 1 X_test = hstack((X_test_clean_cat_ohe ,X_test_clean_subcat_ohe ,X_test_state
2                 X_test_teacher_ohe ,text_tfidf_test, title_tfidf_test, quant
3                 prev_projects_test, price_test,title_word_count_test,
4                 essay_word_count_test, title_word_count_test, negative_inte
5                 neutral_intensity_test)).tocsr()
6 print(X_test.shape, y_test.shape)
7 print(type(X_test))
```

```
(11550, 6235) (11550,)
<class 'scipy.sparse.csr.csr_matrix'>
```

## We don't need CV data since it is automatically done by the optimization function of python

```
In [99]: 1 # X_cv = hstack((cv_categories_one_hot ,sub_categories_one_hot_cv,school_stat
2 #               teacher_prefix_categories_one_hot_cv,text_tfidf_cv, title_
3 #               prev_projects_cv, price_cv, title_word_count_cv,
4 #               essay_word_count_cv ,positive_intensity_cv, negative_inte
5 #               neutral_intensity_cv)).tocsr()
6 # print(X_cv.shape, y_cv.shape)
7 # print(type(X_cv))
```

## Hyperparameter Tuning

```
In [100]: 1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.model_selection import GridSearchCV
3 GBDT = GradientBoostingClassifier()
4
5 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100]}
6
7 classifier_5 = GridSearchCV(GBDT, parameters, cv=3, scoring='roc_auc', verbose=
8 classifier_5.fit(X_train, y_train)
```

Fitting 3 folds for each of 12 candidates, totalling 36 fits

```
[Parallel(n_jobs=-1)]: Done 2 tasks      | elapsed: 27.4s
[Parallel(n_jobs=-1)]: Done 9 tasks      | elapsed: 50.5s
[Parallel(n_jobs=-1)]: Done 16 tasks     | elapsed: 2.5min
[Parallel(n_jobs=-1)]: Done 25 out of 36 | elapsed: 6.5min remaining: 2.9mi
n
[Parallel(n_jobs=-1)]: Done 29 out of 36 | elapsed: 29.0min remaining: 7.0mi
n
[Parallel(n_jobs=-1)]: Done 33 out of 36 | elapsed: 30.0min remaining: 2.7mi
n
[Parallel(n_jobs=-1)]: Done 36 out of 36 | elapsed: 42.1min finished
```

```
Out[100]: GridSearchCV(cv=3, error_score='raise',
                    estimator=GradientBoostingClassifier(criterion='friedman_mse', init=None,
                    learning_rate=0.1, loss='deviance', max_depth=3,
                    max_features=None, max_leaf_nodes=None,
                    min_impurity_decrease=0.0, min_impurity_split=None,
                    min_samples_leaf=1, min_samples_split=2,
                    min_weight_fraction_leaf=0.0, n_estimators=100,
                    presort='auto', random_state=None, subsample=1.0, verbose=0,
                    warm_start=False),
                    fit_params=None, iid=True, n_jobs=-1,
                    param_grid={'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 10
0]},
                    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
                    scoring='roc_auc', verbose=10)
```

```
In [101]: 1 max_depth_best_param=classifier_5.best_params_['max_depth']
2 min_samples_split_best_param=classifier_5.best_params_['min_samples_split']
3 print(classifier_5.best_params_)

{'max_depth': 5, 'min_samples_split': 100}
```

```
In [102]: 1 import plotly.offline as offline
2 import plotly.graph_objs as go
3 offline.init_notebook_mode()
4 import numpy as np
```

```
In [103]: 1 x1 = [5, 10, 100, 500] #min_sample_split
2 y1 = [1, 3, 10, 30] #max_depth
3 z1 = classifier_5.cv_results_['mean_test_score']
4
5 x2 = [5, 10, 100, 500] #min_sample_split
6 y2 = [1, 3, 10, 30] #max_depth
7 z2 = classifier_5.cv_results_['mean_train_score']
8
```

```
In [104]: 1 # https://plot.ly/python/3d-axes/
2 trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'Cross Validation')
3 trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Train')
4 data = [trace1, trace2]
5
6 layout = go.Layout(scene = dict(
7     xaxis = dict(title='min_sample_split'),
8     yaxis = dict(title='max_depth'),
9     zaxis = dict(title='AUC'),))
10
11 fig = go.Figure(data=data, layout=layout)
12 offline.iplot(fig, filename='3d-scatter-colorscale')
```



```

In [105]: 1 def batch_predict(clf, data):
2           # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
3           # not the predicted outputs
4
5           y_data_pred = []
6           tr_loop = data.shape[0] - data.shape[0]%1000
7           # consider you X_tr shape is 49041, then your cr_loop will be 49041 - 490
8           # in this for loop we will iterate until the last 1000 multiplier
9           for i in range(0, tr_loop, 1000):
10              y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
11              # we will be predicting for the last data points
12              y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
13
14           return y_data_pred

```

```

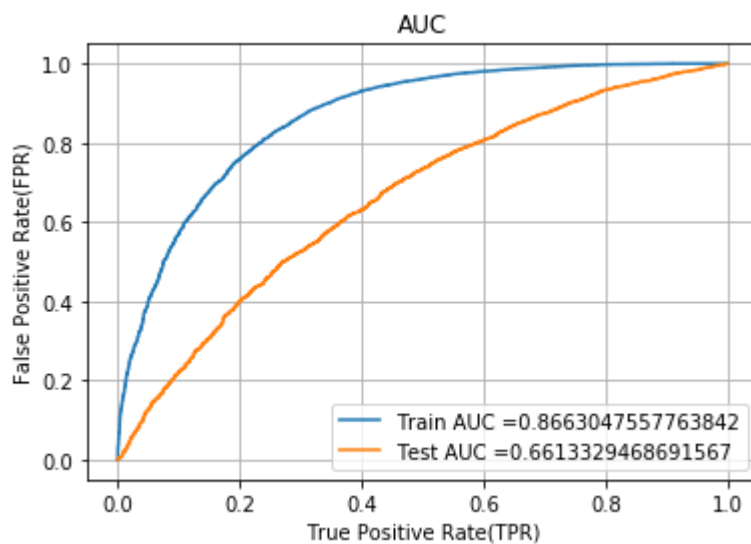
In [106]: 1 def predict(proba, threshold, fpr, tpr):
2           t = threshold[np.argmax(fpr*(1-tpr))]
3           print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold")
4           predictions = []
5           global predictions1
6           for i in proba:
7               if i>=t:
8                   predictions.append(1)
9               else:
10                  predictions.append(0)
11           predictions1 = predictions
12           return predictions

```

```

In [107]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3 from sklearn.ensemble import GradientBoostingClassifier
4
5 classifier = GradientBoostingClassifier(max_depth = max_depth_best_param, min
6 classifier.fit(X_train, y_train)
7
8 y_train_pred = batch_predict(classifier,X_train)
9 y_test_pred = batch_predict(classifier,X_test)
10
11 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
12 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
13
14 plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_t
15 plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
16 plt.legend()
17 plt.xlabel("True Positive Rate(TPR)")
18 plt.ylabel("False Positive Rate(FPR)")
19 plt.title("AUC")
20 plt.grid()
21 plt.show()

```



## Confusion Matrix For Train & Test Data

```

In [108]: 1 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
2 import seaborn as sns; sns.set()
3 print("For the Train Data")
4 con_m_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
5 print("For the Test Data")
6 con_m_test = confusion_matrix(y_test, predict(y_test_pred, te_thresholds, tes
7
8 key = (np.asarray([[ 'TN', 'FP'], [ 'FN', 'TP'] ])))
9 fig, ax = plt.subplots(1,2, figsize=(12,5))
10
11 labels_train = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value
12 labels_test = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value
13
14 sns.heatmap(con_m_train, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDI
15 sns.heatmap(con_m_test, linewidths=.5, xticklabels=['PREDICTED : NO', 'PREDI
16
17 ax[0].set_title('Train Set')
18 ax[1].set_title('Test Set')
19
20 plt.show()
21

```

For the Train Data

the maximum value of  $tpr \cdot (1 - fpr)$  0.615496647765272 for threshold 0.851

For the Test Data

the maximum value of  $tpr \cdot (1 - fpr)$  0.3831900515096276 for threshold 0.851



## Set2 TFIDF W2V

```
In [109]: 1 from scipy.sparse import hstack
2 X_train = np.hstack((X_train_clean_cat_ohe,X_train_clean_subcat_ohe,X_train_s
3                     X_train_teacher_ohe,tfidf_w2v_vectors_train,tfidf_w2v_vecto
4                     prev_projects_train, price_train,title_word_count_train,
5                     essay_word_count_train,
6                     positive_intensity_train,
7                     negative_intensity_train,
8                     neutral_intensity_train))
9 print(X_train.shape, y_train.shape)
10 print(type(X_train))
11 # X_train=X_train.tocsr()
```

```
(23450, 618) (23450,)
<class 'numpy.ndarray'>
```

```
In [110]: 1 X_test = np.hstack((X_test_clean_cat_ohe ,X_test_clean_subcat_ohe ,X_test_st
2                     X_test_teacher_ohe ,tfidf_w2v_vectors_test, tfidf_w2v_vector
3                     prev_projects_test, price_test,title_word_count_test,
4                     essay_word_count_test, title_word_count_test, negative_inte
5                     neutral_intensity_test))
6 print(X_test.shape, y_test.shape)
7 print(type(X_test))
```

```
(11550, 618) (11550,)
<class 'numpy.ndarray'>
```

```
In [111]: 1 # X_cv = hstack((cv_categories_one_hot ,sub_categories_one_hot_cv,school_stat
2 #                     teacher_prefix_categories_one_hot_cv,tfidf_w2v_vectors_cv,
3 #                     prev_projects_cv, price_cv, title_word_count_cv,
4 #                     essay_word_count_cv ,positive_intensity_cv, negative_inte
5 #                     neutral_intensity_cv)).tocsr()
6 # print(X_cv.shape, y_cv.shape)
7 # print(type(X_cv))
```

```
In [112]: 1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.model_selection import GridSearchCV
3 GBDT = GradientBoostingClassifier()
4
5 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100]}
6
7 classifier_5 = GridSearchCV(GBDT, parameters, cv=3, scoring='roc_auc', verbose=10)
8 classifier_5.fit(X_train, y_train)
```

Fitting 3 folds for each of 12 candidates, totalling 36 fits

```
[Parallel(n_jobs=-1)]: Done 2 tasks      | elapsed: 1.7min
[Parallel(n_jobs=-1)]: Done 9 tasks      | elapsed: 3.0min
[Parallel(n_jobs=-1)]: Done 16 tasks     | elapsed: 9.7min
[Parallel(n_jobs=-1)]: Done 25 out of 36 | elapsed: 31.6min remaining: 13.9min
[Parallel(n_jobs=-1)]: Done 29 out of 36 | elapsed: 96.8min remaining: 23.4min
[Parallel(n_jobs=-1)]: Done 33 out of 36 | elapsed: 101.4min remaining: 9.2min
[Parallel(n_jobs=-1)]: Done 36 out of 36 | elapsed: 120.2min finished
```

```
Out[112]: GridSearchCV(cv=3, error_score='raise',
                    estimator=GradientBoostingClassifier(criterion='friedman_mse', init=None,
                    learning_rate=0.1, loss='deviance', max_depth=3,
                    max_features=None, max_leaf_nodes=None,
                    min_impurity_decrease=0.0, min_impurity_split=None,
                    min_samples_leaf=1, min_samples_split=2,
                    min_weight_fraction_leaf=0.0, n_estimators=100,
                    presort='auto', random_state=None, subsample=1.0, verbose=0,
                    warm_start=False),
                    fit_params=None, iid=True, n_jobs=-1,
                    param_grid={'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100]},
                    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
                    scoring='roc_auc', verbose=10)
```

```
In [113]: 1 max_depth_best_param=classifier_5.best_params_['max_depth']
2 min_samples_split_best_param=classifier_5.best_params_['min_samples_split']
3 print(classifier_5.best_params_)

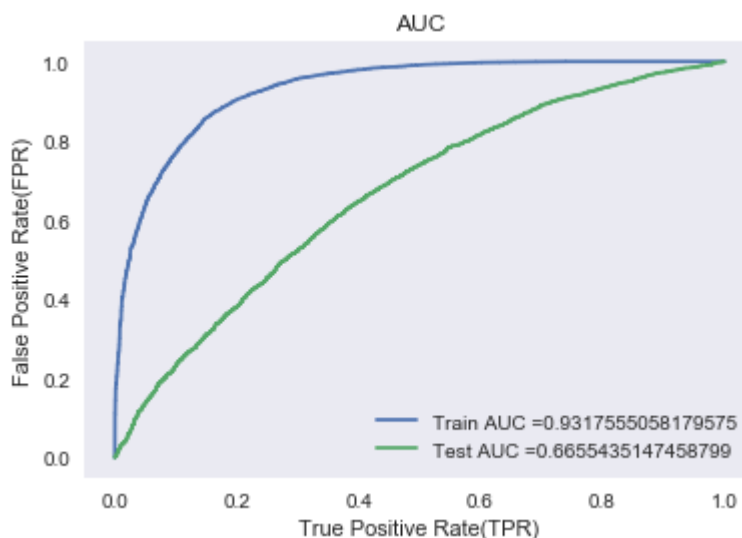
{'max_depth': 5, 'min_samples_split': 10}
```

```
In [114]: 1 import plotly.offline as offline
2 import plotly.graph_objs as go
3 offline.init_notebook_mode()
4 import numpy as np
```

```
In [115]: 1 x1 = [1, 3, 10, 30]#min_sample_split
2 y1 = [5, 10, 100, 500]#max_depth
3 z1 = classifier_5.cv_results_['mean_test_score']
4
5 x2 = [1, 3, 10, 30] #min_sample_split
6 y2 = [5, 10, 100, 500] #max_depth
7 z2 = classifier_5.cv_results_['mean_train_score']
8
```

```
In [116]: 1 # https://plot.ly/python/3d-axes/
2 trace1 = go.Scatter3d(x=x1,y=y1,z=z1, name = 'Cross Validation')
3 trace2 = go.Scatter3d(x=x2,y=y2,z=z2, name = 'Train')
4 data = [trace1, trace2]
5
6 layout = go.Layout(scene = dict(
7     xaxis = dict(title='min_sample_split'),
8     yaxis = dict(title='max_depth'),
9     zaxis = dict(title='AUC'),))
10
11 fig = go.Figure(data=data, layout=layout)
12 offline.iplot(fig, filename='3d-scatter-colorscale')
```

```
In [117]: 1 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve
2 from sklearn.metrics import roc_curve, auc
3 from sklearn.ensemble import GradientBoostingClassifier
4
5 classifier = GradientBoostingClassifier(max_depth = max_depth_best_param, min
6 classifier.fit(X_train, y_train)
7
8 y_train_pred = batch_predict(classifier,X_train)
9 y_test_pred = batch_predict(classifier,X_test)
10
11 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_pred)
12 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
13
14 plt.plot(train_fpr, train_tpr, label="Train AUC =" + str(auc(train_fpr, train_t
15 plt.plot(test_fpr, test_tpr, label="Test AUC =" + str(auc(test_fpr, test_tpr)))
16 plt.legend()
17 plt.xlabel("True Positive Rate(TPR)")
18 plt.ylabel("False Positive Rate(FPR)")
19 plt.title("AUC")
20 plt.grid()
21 plt.show()
```



## Confusion Matrix For Train & Test Data

```

In [118]: 1 #https://www.quantinsti.com/blog/creating-heatmap-using-python-seaborn
2 import seaborn as sns; sns.set()
3 print("For the Train Data")
4 con_m_train = confusion_matrix(y_train, predict(y_train_pred, tr_thresholds,
5 print("For the Test Data")
6 con_m_test = confusion_matrix(y_test, predict(y_test_pred, te_thresholds, tes
7
8 key = (np.asarray([[ 'TN', 'FP'], [ 'FN', 'TP'] ])))
9 fig, ax = plt.subplots(1,2, figsize=(12,5))
10
11 labels_train = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value
12 labels_test = (np.asarray(["{0} = {1:.2f}" .format(key, value) for key, value
13
14 sns.heatmap(con_m_train, linewidths=.5, xticklabels=[ 'PREDICTED : NO', 'PREDI
15 sns.heatmap(con_m_test, linewidths=.5, xticklabels=[ 'PREDICTED : NO', 'PREDIC
16
17 ax[0].set_title('Train Set')
18 ax[1].set_title('Test Set')
19
20 plt.show()

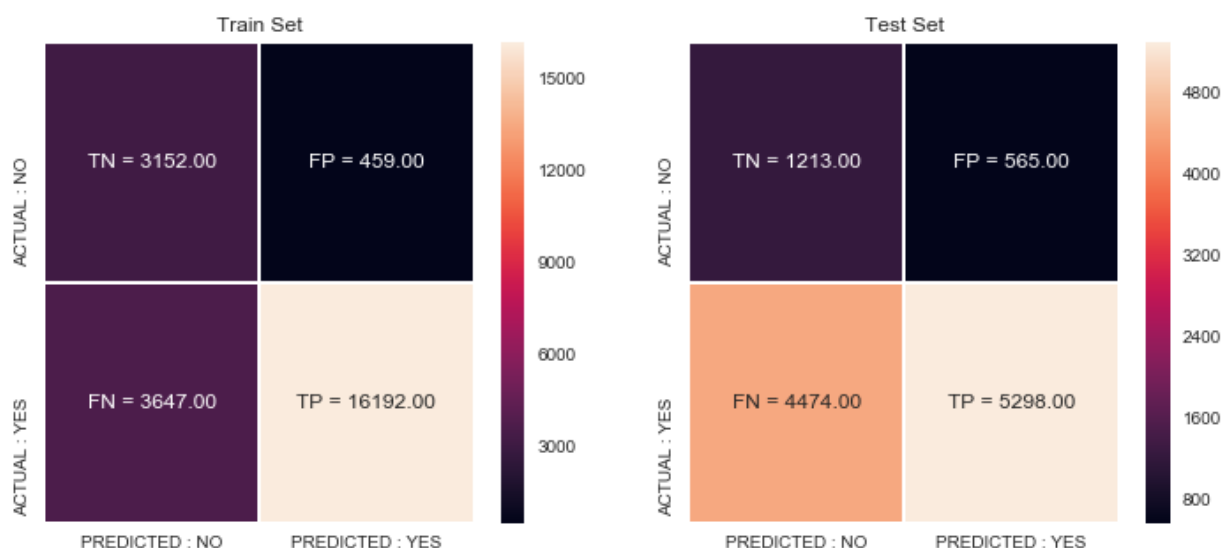
```

For the Train Data

the maximum value of  $tpr \cdot (1 - fpr)$  0.7302713875219472 for threshold 0.836

For the Test Data

the maximum value of  $tpr \cdot (1 - fpr)$  0.3893065607896025 for threshold 0.876





# Conclusion

```
In [120]: 1 # Please compare all your models using Prettytable Library
2 # http://zetcode.com/python/prettytable/
3 from prettytable import PrettyTable
4 TB = PrettyTable()
5 TB.field_names = ["Model", "Hyperparameter", "Train_AUC", "Test_Auc"]
6 TB.title = "Decision Tree"
7 TB.add_row(["TFIDF", "Depth:5 | Samp_Split:100", 0.8663, 0.6613])
8 TB.add_row(["TFIDF-W2V", "Depth:5 | Samp_Split:10", 0.93175, 0.665543])
9 print(TB)
```

```
+-----+
|                                     |
|               Decision Tree        |
|-----+-----+-----+-----+
|  Model  | Hyperparameter | Train_AUC | Test_Auc |
|-----+-----+-----+-----+
|  TFIDF  | Depth:5 | Samp_Split:100 | 0.8663 | 0.6613 |
| TFIDF-W2V | Depth:5 | Samp_Split:10 | 0.93175 | 0.665543 |
+-----+-----+-----+-----+
```

```
In [ ]: 1
```