

▼ Assignment 9: GBDT

▼ Response Coding: Example

Train Data			Encoded Train Data		
State	class		State_0	State_1	class
A	0		3/5	2/5	0
B	1		0/2	2/2	1
C	1		1/3	2/3	1
A	0		3/5	2/5	0
A	1		3/5	2/5	1
B	1		0/2	2/2	1
A	0		3/5	2/5	0
A	1		3/5	2/5	1
C	1		1/3	2/3	1
C	0		1/3	2/3	0

Resonse table(only from train)			
State	Class=0	Class=1	
A	3	2	
B	0	2	
C	1	2	

Test Data			Encoded Test Data	
State			State_0	State_1
A			3/5	2/5
C			1/3	2/3
D			1/2	1/2
C			1/3	2/3
B			0/2	2/2
E			1/2	1/2

The response tabel is built only on train dataset. For a category which is not there in train data and present in test data, we will encode them with default values Ex: in our test data if have State: D then we encode it as [0.5, 0.05]

1. Apply GBDT on these feature sets

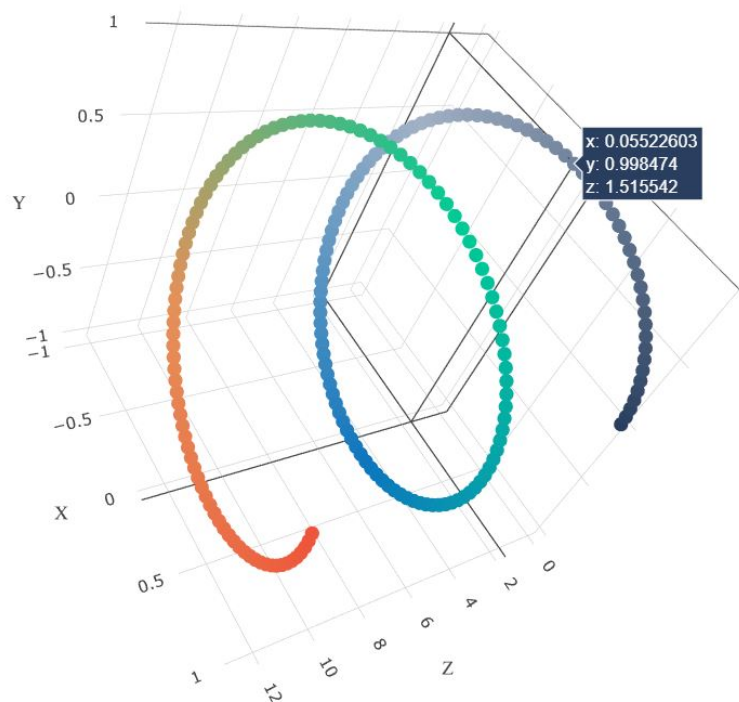
- **Set 1:** categorical(instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)+sentiment Score of eassay(check the bellow example, include all 4 values as 4 features)
- **Set 2:** categorical(instead of one hot encoding, try [response coding](#): use probability values), numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V)

2. The hyper paramter tuning (Consider any two hyper parameters)

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

3. Representation of results

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

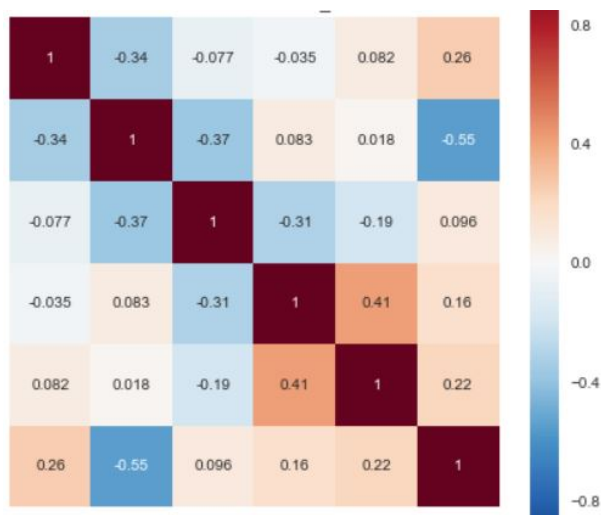


with X-axis as

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

or

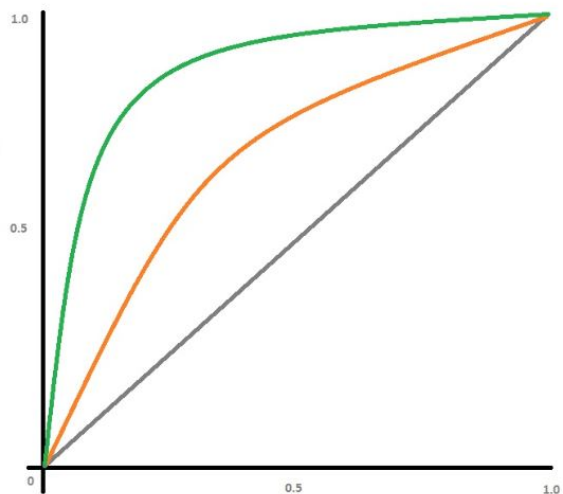
- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure



[seaborn heat maps](#) with rows as

n_estimators, 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.



- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted

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

and original labels of test data points

4. You need to summarize the results at the end of the notebook, summarize it in the table

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

format

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

↳ Drive already mounted at /content/drive; to attempt to forcibly remount, call

```
1 !pip install xgboost
```

↳ Requirement already satisfied: xgboost in /usr/local/lib/python3.6/dist-packa
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-package
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-package

```
1 !pip install plotly --upgrade
```

↳ Collecting plotly
Downloading <https://files.pythonhosted.org/packages/68/47/cec583df9ffb6142b>
|████████████████████| 13.1MB 239kB/s
Requirement already satisfied, skipping upgrade: retrying>=1.3.3 in /usr/loca
Requirement already satisfied, skipping upgrade: six in /usr/local/lib/python
Installing collected packages: plotly
Found existing installation: plotly 4.4.1
Uninstalling plotly-4.4.1:
Successfully uninstalled plotly-4.4.1
Successfully installed plotly-4.11.0

```
1 path = '/content/drive/My Drive/AAIC/ASSIGN_13/11_Donors_choose_GBDT/'
```

▼ Task 1

▼ 1.1 Loading Data

```
1 import pandas
2 data = pandas.read_csv(path+'preprocessed_data.csv',nrows=50000)
```

```
1 data.head(1)
```



school_state	teacher_prefix	project_grade_category	teacher_number_of_pro
--------------	----------------	------------------------	-----------------------

0	ca	mrs	grades_prek_2
---	----	-----	---------------

```
1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 import re
10 from tqdm import tqdm
```

1.1.0 Adding Sentiment analysis features such as NEG, POS, NEU and COMPOUND polarity scores

```
1 import nltk
2 nltk.download('vader_lexicon')
3 from nltk.sentiment.vader import SentimentIntensityAnalyzer
4
5 sid = SentimentIntensityAnalyzer()
6 NEG = []
7 NEU = []
8 POS = []
9 COMP = []
10 essays = data['essay'].values
11 for senetence in essays:
12     ss = sid.polarity_scores(senetence)
13     NEG.append(ss['neg'])
14     POS.append(ss['pos'])
15     NEU.append(ss['neu'])
16     COMP.append(ss['compound'])
17
18 df2 = pd.DataFrame({'NEG':NEG,
19                    'POS':POS,
20                    'NEU':NEU,
21                    'COMPOUND':COMP})
22
```



[nltk_data] Downloading package vader_lexicon to /root/nltk_data...

```
1 df2.head(2)
```



	NEG	POS	NEU	COMPOUND
0	0.013	0.205	0.783	0.9867
1	0.072	0.248	0.680	0.9897

```
1 data = pd.concat([data, df2], axis=1)
2 data.head(2)
```



	school_state	teacher_prefix	project_grade_category	teacher_number_of_projects
0	ca	mrs	grades_prek_2	1
1	ut	ms	grades_3_5	1

▼ 1.1.1 Processing Project Essay.

```
1 def decontracted(phrase):
2     # specific
3     phrase = re.sub(r"won't", "will not", phrase)
4     phrase = re.sub(r"can't", "can not", phrase)
5
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
```

```
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', 'it's',
6             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that',
7             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have',
8             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because',
9             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
10            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off',
```

```

10     'above', 'below', 'so', 'from', 'up', 'down', 'in', 'out', 'on', 'at',
11     'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'as',
12     'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'that',
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", 'isn't',
16     "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't",
17     'won', "won't", 'wouldn', "wouldn't"]

```

```

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

```

```
1 data.essay.iloc[:1]
```

```

0    i fortunate enough use fairy tale stem kits cl...
Name: essay, dtype: object

```

As there are some stop words present in the text of essay of already preprocessed_data. For example word `i` is present in the first essay sentence. So, I am processing the essay text again inorder to remove the stop words if there are any by chance and also because unnecessary stop words increase the dimensionality of the bow/tfidf representation.

```
1 data['essay'] = preprocess_text(data['essay'].values)
```

```

100%|██████████| 50000/50000 [00:20<00:00, 2484.92it/s]

```

```
1 data.essay.iloc[:1]
```

```

0    fortunate enough use fairy tale stem kits clas...
Name: essay, dtype: object

```

```
1 data.head(1)
```

```

0

```

school_state	teacher_prefix	project_grade_category	teacher_number_of_pro
--------------	----------------	------------------------	-----------------------

▼ 1.1.2 Loading data after processing essay

```
1 Y=data['project_is_approved']
2 X=data.drop('project_is_approved',axis=1)
```

```
1 X.head(1)
```

```
↳ 

|   | school_state | teacher_prefix | project_grade_category | teacher_number_of_pro |
|---|--------------|----------------|------------------------|-----------------------|
| 0 | ca           | mrs            | grades_prek_2          |                       |


```

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

```
1 from sklearn.model_selection import train_test_split
2
3 X_train,X_te,Y_train,Y_te = train_test_split(X,Y,test_size=0.33, stratify=Y, r
```

```
1 print(X_train.shape,Y_train.shape,X_te.shape,Y_te.shape,end=" ")
```

```
↳ (33500, 12) (33500,) (16500, 12) (16500,)
```

1.3 Make Data Model Ready: encoding project essay

▼ 1.3.1 TFIDF Representation of essay.

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,4))
4 vectorizer.fit(X_train['essay'].values)
```

```
↳
```



```
TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
                dtype=<class 'numpy.float64'>, encoding='utf-8',
```

```
1
2 essay_tr_tfidf = vectorizer.transform(X_train['essay'].values)
3 essay_te_tfidf = vectorizer.transform(X_te['essay'].values)
4
```

```
1 print(essay_tr_tfidf.shape, essay_te_tfidf.shape, end=" ")
```

```
↳ (33500, 99859) (16500, 99859)
```

▼ 1.3.2 TFIDF-W2V Representation of essay.

```
1 import pickle
2 with open(path+'glove_vectors', 'rb') as f:
3     model = pickle.load(f)
4     glove_words = set(model.keys())
```

```
1 # average Word2Vec
2 # compute average word2vec for each review.
3 def tfidf_w2v_transform(preprocessed_essays):
4     tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in
5     for sentence in tqdm(preprocessed_essays): # 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/review
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 tf value
12                 tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
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
19     return tfidf_w2v_vectors
20
```

▼ Train data

```
1 # S = ["abc def pqr", "def def def abc", "pqr pqr def"]
2 tfidf_model = TfidfVectorizer(min_df=10, ngram_range=(1,4))
3 tfidf_model.fit(X_train['essay'].values)
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())
```

```
1 essay_tr_tfidf_w2v = np.array(tfidf_w2v_transform(X_train['essay']))
```

```

1 essay_tr_tfidf_w2v = np.array(tfidf_w2v_transform(X_train['essay']))
2 essay_te_tfidf_w2v = np.array(tfidf_w2v_transform(X_te['essay']))
3 #print(preprocessed_essays[:3])

```

```

↳ 100%|██████████| 33500/33500 [01:05<00:00, 511.80it/s]
100%|██████████| 16500/16500 [00:32<00:00, 508.49it/s]

```

```

1 print(essay_tr_tfidf_w2v.shape, essay_te_tfidf_w2v.shape, end=" ")

```

```

↳ (33500, 300) (16500, 300)

```

1.4 Make Data Model Ready: encoding numerical, categorical features

▼ 1.3.3 Encoding Categorical features

```

1 data.columns

```

```

↳ Index(['school_state', 'teacher_prefix', 'project_grade_category',
        'teacher_number_of_previously_posted_projects', 'project_is_approved',
        'clean_categories', 'clean_subcategories', 'essay', 'price', 'NEG',
        'POS', 'NEU', 'COMPOUND'],
        dtype='object')

```

▼ School State

```

1 df = pd.DataFrame({'ftr' : X_train['school_state'], 'label' : Y_train})
2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'], index = set(df['ftr']))
3 for i in set(df['ftr']):
4     sam_df.loc[i].prob_0 = len(df[(df['ftr'] == (i)) & (df['label'] == 0)]) / len(df)
5     sam_df.loc[i].prob_1 = len(df[(df['ftr'] == (i)) & (df['label'] == 1)]) / len(df)

```

```

1 cat_data_tr = pd.DataFrame(columns= ['prob_0', 'prob_1'], index = df.index)
2 for i in list(df.index.values):
3     j = df.loc[i]['ftr'].values
4     #print(j[0])
5     cat_data_tr.loc[i] = {'prob_0':sam_df.loc[j[0]].prob_0 , 'prob_1':sam_df.loc[j[0]].prob_1}
6 cat_data_tr= cat_data_tr.values
7 print(cat_data_tr.shape)

```

```

↳ (33500, 2)

```

```

1 df3 = X_te['school_state']
2 cat_data_te = pd.DataFrame(columns= ['prob_0', 'prob_1'], index = df3.index)
3 for i in (df3.index):
4     if df3.loc[i] not in sam_df.index:
5         cat_data_te.loc[i] = {'prob_0':0.5 , 'prob_1':0.5}
6     else:

```

```

7     cat_data_te.loc[i] = {'prob_0':sam_df.loc[df3.loc[i]].prob_0 , 'prob_1':sa
8     cat_data_te = cat_data_te.values

```

```

1 state_tr = cat_data_tr
2 state_te = cat_data_te

```

```

1 print(state_tr.shape,state_te.shape,end=" ")

```

➞ (33500, 2) (16500, 2)

▼ Teacher prefix

```

1 df = pd.DataFrame({'ftr' : X_train['teacher_prefix'], 'label' : Y_train})
2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'],index = set(df['ftr']) )
3 for i in set(df['ftr']):
4     sam_df.loc[i].prob_0 = len(df[(df['ftr'] == (i)) & (df['label'] == 0)])/ len
5     sam_df.loc[i].prob_1 = len(df[(df['ftr'] == (i)) & (df['label'] == 1)])/ len

```

```

1 cat_data_tr = pd.DataFrame(columns= ['prob_0','prob_1'],index = df.index)
2 for i in list(df.index.values):
3     j = df.loc[i]['ftr'].values
4     #print(j[0])
5     cat_data_tr.loc[i] = {'prob_0':sam_df.loc[j[0]].prob_0 , 'prob_1':sam_df.loc
6 cat_data_tr= cat_data_tr.values
7 print(cat_data_tr.shape)

```

➞ (33500, 2)

```

1 df3 = X_te['teacher_prefix']
2 cat_data_te = pd.DataFrame(columns= ['prob_0','prob_1'],index = df3.index)
3 for i in (df3.index):
4     if df3.loc[i] not in sam_df.index:
5         cat_data_te.loc[i] = {'prob_0':0.5 , 'prob_1':0.5}
6     else:
7         cat_data_te.loc[i] = {'prob_0':sam_df.loc[df3.loc[i]].prob_0 , 'prob_1':sa
8 cat_data_te = cat_data_te.values

```

```

1 tchr_prefix_tr = cat_data_tr
2 tchr_prefix_te = cat_data_te

```

```

1 print(tchr_prefix_tr.shape,tchr_prefix_te.shape,end=" ")

```

➞ (33500, 2) (16500, 2)

▼ Project grade_category

```

1 df = pd.DataFrame({'ftr' : X_train['project_grade_category'], 'label' : Y_train)
2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'],index = set(df['ftr']) )

```

```

2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'],index = set(df['ftr']) ,
3 for i in set(df['ftr']):
4     sam_df.loc[i].prob_0 = len(df[(df['ftr'] == (i)) & (df['label'] == 0)])/ len
5     sam_df.loc[i].prob_1 = len(df[(df['ftr'] == (i)) & (df['label'] == 1)])/ len

```

```

1 cat_data_tr = pd.DataFrame(columns= ['prob_0','prob_1'],index = df.index)
2 for i in list(df.index.values):
3     j = df.loc[i][['ftr']].values
4     #print(j[0])
5     cat_data_tr.loc[i] = {'prob_0':sam_df.loc[j[0]].prob_0 , 'prob_1':sam_df.loc
6 cat_data_tr= cat_data_tr.values
7 print(cat_data_tr.shape)

```

➞ (33500, 2)

```

1 df3 = X_te['project_grade_category']
2 cat_data_te = pd.DataFrame(columns= ['prob_0','prob_1'],index = df3.index)
3 for i in (df3.index):
4     if df3.loc[i] not in sam_df.index:
5         cat_data_te.loc[i] = {'prob_0':0.5 , 'prob_1':0.5}
6     else:
7         cat_data_te.loc[i] = {'prob_0':sam_df.loc[df3.loc[i]].prob_0 , 'prob_1':sa
8 cat_data_te = cat_data_te.values

```

```

1 grade_tr = cat_data_tr
2 grade_te = cat_data_te

```

```

1 print(grade_tr.shape,grade_te.shape,end=" ")

```

➞ (33500, 2) (16500, 2)

▼ Subject categories

```

1 df = pd.DataFrame({'ftr' : X_train['clean_categories'], 'label' : Y_train})
2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'],index = set(df['ftr']) ,
3 for i in set(df['ftr']):
4     sam_df.loc[i].prob_0 = len(df[(df['ftr'] == (i)) & (df['label'] == 0)])/ len
5     sam_df.loc[i].prob_1 = len(df[(df['ftr'] == (i)) & (df['label'] == 1)])/ len

```

```

1 cat_data_tr = pd.DataFrame(columns= ['prob_0','prob_1'],index = df.index)
2 for i in list(df.index.values):
3     j = df.loc[i][['ftr']].values
4     #print(j[0])
5     cat_data_tr.loc[i] = {'prob_0':sam_df.loc[j[0]].prob_0 , 'prob_1':sam_df.loc
6 cat_data_tr= cat_data_tr.values
7 print(cat_data_tr.shape)

```

➞ (33500, 2)

```

1 df3 = X_te['clean_categories']

```

```

2 cat_data_te = pd.DataFrame(columns= ['prob_0','prob_1'],index = df3.index)
3 for i in (df3.index):
4     if df3.loc[i] not in sam_df.index:
5         cat_data_te.loc[i] = {'prob_0':0.5 , 'prob_1':0.5}
6     else:
7         cat_data_te.loc[i] = {'prob_0':sam_df.loc[df3.loc[i]].prob_0 , 'prob_1':sam_df.loc[df3.loc[i]].prob_1}
8 cat_data_te = cat_data_te.values

```

```

1 subject_tr = cat_data_tr
2 subject_te = cat_data_te

```

```

1 print(subject_tr.shape,subject_te.shape,end=" ")

```

```

↳ (33500, 2) (16500, 2)

```

▼ Subject sub categories

```

1 df = pd.DataFrame({'ftr' : X_train['clean_subcategories'], 'label' : Y_train})
2 sam_df = pd.DataFrame(columns = ['prob_0', 'prob_1'],index = set(df['ftr']))
3 for i in set(df['ftr']):
4     sam_df.loc[i].prob_0 = len(df[(df['ftr'] == (i)) & (df['label'] == 0)])/ len(df[(df['ftr'] == (i)) & (df['label'] == 0) | (df['label'] == 1)])
5     sam_df.loc[i].prob_1 = len(df[(df['ftr'] == (i)) & (df['label'] == 1)])/ len(df[(df['ftr'] == (i)) & (df['label'] == 0) | (df['label'] == 1)])

```

```

1 print(df.shape,sam_df.shape)

```

```

↳ (33500, 2) (342, 2)

```

```

1 cat_data_tr = pd.DataFrame(columns= ['prob_0','prob_1'],index = df.index)
2 for i in list(df.index.values):
3     j = df.loc[i]['ftr'].values
4     #print(j[0])
5     cat_data_tr.loc[i] = {'prob_0':sam_df.loc[j[0]].prob_0 , 'prob_1':sam_df.loc[j[0]].prob_1}
6 cat_data_tr= cat_data_tr.values
7 print(cat_data_tr.shape)

```

```

↳ (33500, 2)

```

```

1 df3 = X_te['clean_subcategories']
2 cat_data_te = pd.DataFrame(columns= ['prob_0','prob_1'],index = df3.index)
3 for i in (df3.index):
4     if df3.loc[i] not in sam_df.index:
5         cat_data_te.loc[i] = {'prob_0':0.5 , 'prob_1':0.5}
6     else:
7         cat_data_te.loc[i] = {'prob_0':sam_df.loc[df3.loc[i]].prob_0 , 'prob_1':sam_df.loc[df3.loc[i]].prob_1}
8 cat_data_te = cat_data_te.values

```

```

1 print(df3.shape,cat_data_te.shape)

```

```

↳ (16500,) (16500, 2)

```

```

1  subj_subct_tr = cat_data_tr
2  subj_subct_te = cat_data_te

```

```

1  print(subj_subct_tr.shape,subj_subct_te.shape,end=" ")

```

```

↳ (33500, 2) (16500, 2)

```

▼ 1.3.4 Encoding Numerical features

▼ teacher_number_of_previously_posted_projects

```

1  #For z-score normalization the values will be both negative and positive. Multi
2  #So we use Min-Max normalization.
3  #from sklearn.preprocessing import StandardScaler
4  from sklearn.preprocessing import MinMaxScaler
5
6  scaler = MinMaxScaler()
7  scaler.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
8
9  tchr_prj_tr = scaler.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
10 tchr_prj_te = scaler.transform(X_te['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))

```

```

1  print(list(tchr_prj_tr[:5]),'\n',list(tchr_prj_te[:5]))

```

```

↳ [array([0.00561798]), array([0.09269663]), array([0.02808989]), array([0.00561798]), array([0.00842697]), array([0.00561798]), array([0.01966292]), array([0.0361798])]

```

```

1  #from sklearn.preprocessing import MinMaxScaler
2
3  #scaler = MinMaxScaler()
4  #scaler.fit(data['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
5  #data['teacher_number_of_previously_posted_projects']=scaler.transform(data['teacher_number_of_previously_posted_projects'].values.reshape(-1, 1))
6

```

▼ price

```

1  # Min-Max Normalization
2
3  scaler = MinMaxScaler()
4  scaler.fit(X_train['price'].values.reshape(-1, 1))
5
6  price_tr = scaler.transform(X_train['price'].values.reshape(-1, 1))
7  price_te = scaler.transform(X_te['price'].values.reshape(-1, 1))

```

```

1  print(list(price_tr[:5]),'\n',list(price_te[:5]))

```

```
↳ [array([0.00185336]), array([0.07942025]), array([0.03560584]), array([0.0789
      [array([0.07126968]), array([0.00790852]), array([0.02513283]), array([0.021
```

```
1 # Min-Max Normalization Use only one Normalization but not both.
2
3 #scaler = MinMaxScaler()
4 #scaler.fit(data['price'].values.reshape(-1, 1))
5 #data['price']=scaler.transform(data['price'].values.reshape(-1, 1))
```

▼ 1.4 Stacking all features to form set1 and set2 features.

```
1 data.columns
```

```
↳ Index(['school_state', 'teacher_prefix', 'project_grade_category',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'essay', 'price', 'NEG',
      'POS', 'NEU', 'COMPOUND'],
      dtype='object')
```

▼ 1.4.1. Set1(TFIDF) features

```
1 from scipy.sparse import hstack
2
3 X_train_tfidf = hstack((essay_tr_tfidf, state_tr.astype('float64'), tchr_pref:
4 X_te_tfidf = hstack((essay_te_tfidf, state_te.astype('float64'), tchr_prefix_t
5
6 print("Final Data matrix")
7 print(X_train_tfidf.shape, Y_train.shape)
8 print(X_te_tfidf.shape, Y_te.shape)
9 print("="*100)
```

```
↳ Final Data matrix
(33500, 99875) (33500,)
(16500, 99875) (16500,)
=====
```

▼ 1.4.2 Set2(TFIDF-W2V) features

```
1 from scipy.sparse import hstack
2 from scipy.sparse import csr_matrix
3
4 X_train_tfidf_w2v = hstack((csr_matrix(essay_tr_tfidf_w2v), state_tr.astype('f
5 X_te_tfidf_w2v = hstack((csr_matrix(essay_te_tfidf_w2v), state_te.astype('floa
6
7 print("Final Data matrix")
8 print(X_train_tfidf_w2v.shape, Y_train.shape)
9 print(X_te_tfidf_w2v.shape, Y_te.shape)
10 print("="*100)
```

```
↳ Final Data matrix
(33500, 316) (33500,)
(16500, 316) (16500,)
```

GBDT (xgboost/lightgbm)

Applying Models on different kind of featurization as mentioned in the instructions

Apply GBDT on different kind of featurization as mentioned in the instructions

For Every model that you work on make sure you do the step 2 and step 3 of instructions

▼ 2. Hyper parameter Tuning

▼ 2.0 functions for plotly plotting

```
1 import plotly.offline as offline
2 import plotly.graph_objs as go
3 offline.init_notebook_mode()
4 import numpy as np
```

↳

```
1 #https://stackoverflow.com/questions/47230817/plotly-notebook-mode-with-google
2 def enable_plotly_in_cell():
3     import IPython
4     from plotly.offline import init_notebook_mode
5     display(IPython.core.display.HTML(''<script src="/static/components/require
6     init_notebook_mode(connected=False)
```

▼ 2.1 Set1 (TF-IDF) Features

```
1 import xgboost as xgb
2 from sklearn.model_selection import GridSearchCV
3
4 depth = [3, 4, 5] # 6, 7]
5 no_of_est = [100, 110, 120] #, 130, 140, 150]
6 hyper_param = {'max_depth' : depth, 'n_estimators' : no_of_est }
7
8 clf = xgb.XGBClassifier(booster = 'gbtree', objective = 'binary:logistic', lea
9 classifier = GridSearchCV(clf, hyper_param, cv=5, return_train_score=True, scoi
```



```

9 classifier = GradientBoostingClassifier(hyper_param, cv=5, random_state=100)
10 classifier.fit(X_train_tfidf, Y_train)
11
12 results = pd.DataFrame.from_dict(classifier.cv_results_)

```

```

1 results = results.sort_values(['rank_test_score'])
2 results.head()

```

```

↳

```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_dep
5	230.624216	0.996491	1.094910	0.011630	
4	214.412764	1.369309	1.105326	0.016504	
8	268.062441	38.728753	1.016270	0.231020	
7	263.698027	1.218970	1.109874	0.034966	
3	195.380259	1.210561	1.115658	0.023528	

```

1 train_auc= results['mean_train_score']
2 cv_auc = results['mean_test_score']
3 depth = results['param_max_depth']
4 no_of_est = results['param_n_estimators']
5 # we need not generate mesh grid explicitly as depth and min_samples combined
6 #print(depth)
7 #print(no_of_est)

```

```

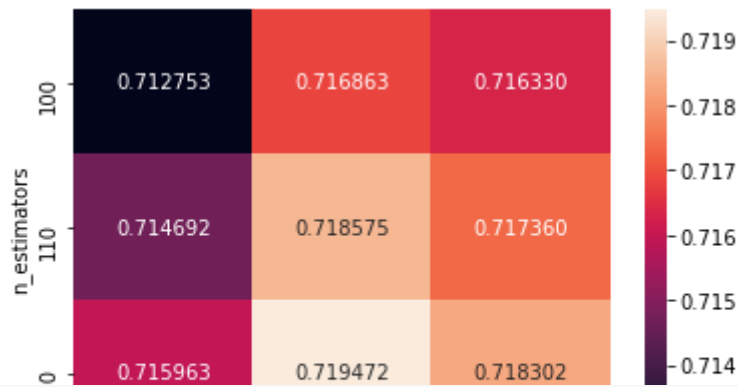
1 import seaborn as sns
2 cv_auc_values = pd.DataFrame({'max_depth':depth,
3                               'n_estimators': no_of_est,
4                               'auc': cv_auc})
5 cv_auc_values = cv_auc_values.pivot("n_estimators","max_depth", "auc")
6 sns.heatmap(cv_auc_values, annot=True, fmt="f")
7 plt.show()

```

```

↳

```



```

1 results = results.sort_values(['rank_test_score'])
2 best_depth = int(results[results['rank_test_score'] == 1].param_max_depth)
3 best_n_estimators = int(results[results['rank_test_score'] == 1].param_n_estimators)
4 results.head(2)

```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_depth
5	230.624216	0.996491	1.094910	0.011630	10
4	214.412764	1.369309	1.105326	0.016504	12

```

1 print('GBDT having depth {} and no of estimators(boosting stages) {} gives the best score')

```

```

➞ GBDT having depth 4 and no of estimators(boosting stages) 120 gives the best score

```

From the above plot and table, we can observe that best hyper parameters are max_depth = 10 and n_estimators = 500 and we take these values for testing.

▼ 2.1.1 Testing

```

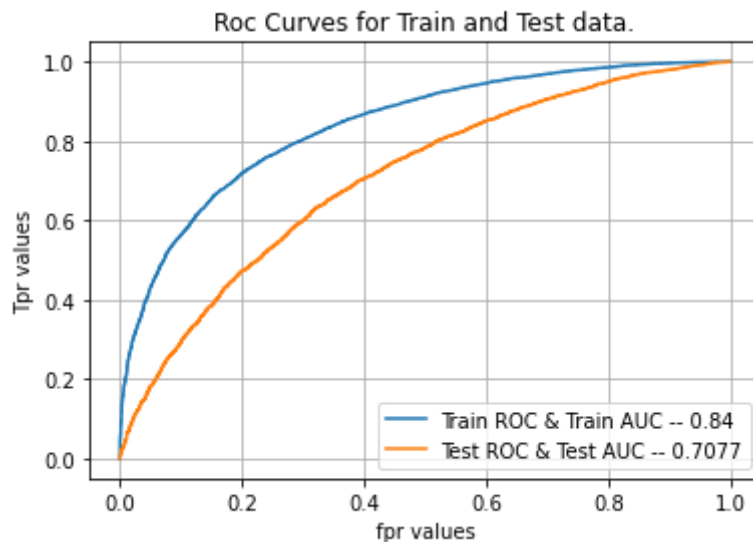
1 from sklearn.metrics import roc_curve,roc_auc_score,confusion_matrix
2
3 clf = xgb.XGBClassifier(max_depth = best_depth, n_estimators = best_n_estimators)
4 clf.fit(X_train_tfidf,Y_train)
5
6 y_prob_tr = clf.predict_proba(X_train_tfidf)[:,-1]
7 fpr_tr, tpr_tr, thresholds_tr = roc_curve(Y_train,y_prob_tr)
8 plt.plot(fpr_tr, tpr_tr)
9
10 y_prob_te = clf.predict_proba(X_te_tfidf)[:,-1]
11 fpr_te, tpr_te, thresholds_te = roc_curve(Y_te,y_prob_te)
12 plt.plot(fpr_te, tpr_te)
13
14 AUC_tr_tfidf = roc_auc_score(Y_train, y_prob_tr)

```

```

15 AUC_te_tfidf = roc_auc_score(Y_te, y_prob_te)
16
17 plt.grid()
18 plt.legend( ('Train ROC & Train AUC -- '+str(np.round(AUC_tr_tfidf, decimals=
19 plt.xlabel('fpr values')
20 plt.ylabel('Tpr values')
21 plt.title('Roc Curves for Train and Test data.')
22 plt.show()
23

```



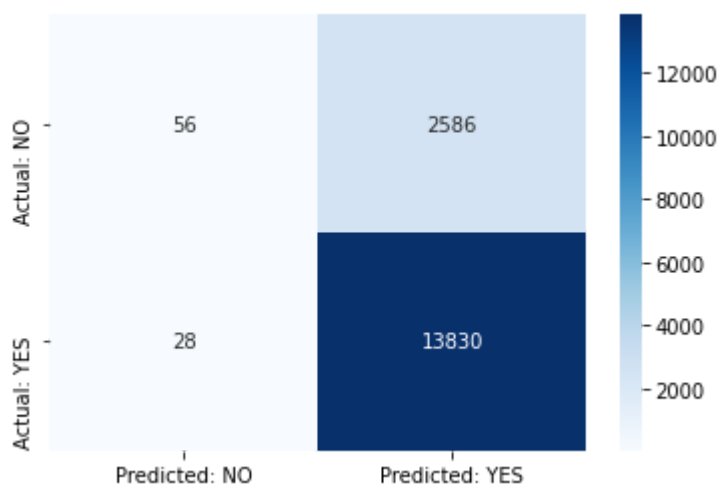
```

1 y_pred_tfidf=clf.predict(X_te_tfidf)
2 cnf_mtrx_tfidf = confusion_matrix(Y_te, y_pred_tfidf)
3
4 print('Confusion matrix for the Test data.')
5 #print(cnf_mtrx_tfidf)
6 x=pd.DataFrame(cnf_mtrx_tfidf, columns = ["Predicted: NO","Predicted: YES"], :
7 sns.heatmap(x, annot=True,fmt="d", cmap='Blues')
8 plt.show()

```



Confusion matrix for the Test data.



▼ 2.2 Set2(TF-IDF Weighted W2V) Features dup

```

1 import xgboost as xgb

```

```

1 import xgboost as xgb
2 from sklearn.model_selection import GridSearchCV
3
4 depth = [3, 4, 5] #, 6, 7]
5 no_of_est = [100, 110, 120] #, 130, 140, 150]
6 hyper_param = {'max_depth' : depth, 'n_estimators' : no_of_est }
7
8 clf = xgb.XGBClassifier(booster = 'gbtree',objective = 'binary:logistic', lea
9
10 classifier = GridSearchCV(clf, hyper_param, cv=5,return_train_score=True, sco
11 classifier.fit(X_train_tfidf_w2v, Y_train)
12
13 results = pd.DataFrame.from_dict(classifier.cv_results_)

```

```

1 results = results.sort_values(['rank_test_score'])
2 results.head(2)

```

```

↳

```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_dep
8	208.193783	27.187026	1.374231	0.332773	
4	163.315378	0.197586	1.507982	0.022203	

```

1 train_auc= results['mean_train_score']
2 cv_auc = results['mean_test_score']
3 depth = results['param_max_depth']
4 no_of_est = results['param_n_estimators']
5 # we need not generate mesh grid explicitly as depth and min_samples combined
6 #print(depth)
7 #print(no_of_est)

```

```

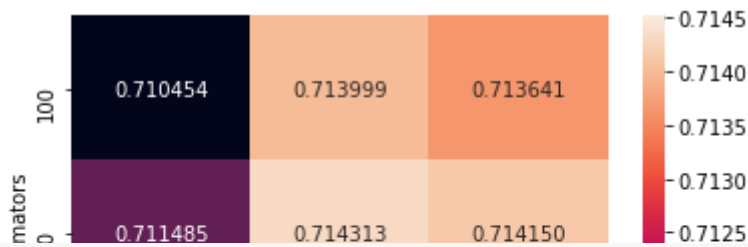
1 import seaborn as sns
2 cv_auc_values = pd.DataFrame({'max_depth':depth,
3                               'n_estimators': no_of_est,
4                               'auc': cv_auc})
5 cv_auc_values = cv_auc_values.pivot("n_estimators","max_depth", "auc")
6 sns.heatmap(cv_auc_values, annot=True, fmt="f")
7 plt.show()

```

```

↳

```



```

1 results = results.sort_values(['rank_test_score'])
2 best_depth_tfidf_w2v = int(results[results['rank_test_score'] == 1].param_max_depth)
3 best_n_estimators_tfidf_w2v = int(results[results['rank_test_score'] == 1].param_n_estimators)
4 results.head(2)

```

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max_depth
8	208.193783	27.187026	1.374231	0.332773	100
4	163.315378	0.197586	1.507982	0.022203	0

```
1 print('GBDT having depth {} and no of estimators(boosting stages) {} gives the best performance')
```

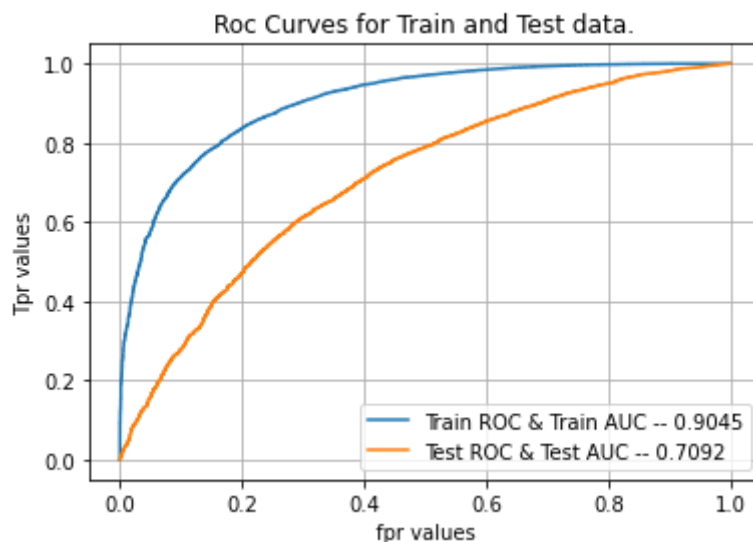
GBDT having depth 5 and no of estimators(boosting stages) 120 gives the best performance

2.2.1 Testing

```

1 from sklearn.metrics import roc_curve,roc_auc_score,confusion_matrix
2
3 clf = xgb.XGBClassifier(max_depth = best_depth_tfidf_w2v, n_estimators = best_n_estimators_tfidf_w2v)
4 clf.fit(X_train_tfidf_w2v,Y_train)
5
6 y_prob_tr = clf.predict_proba(X_train_tfidf_w2v)[:,-1]
7 fpr_tr, tpr_tr, thresholds_tr = roc_curve(Y_train,y_prob_tr)
8 plt.plot(fpr_tr, tpr_tr)
9
10 y_prob_te = clf.predict_proba(X_te_tfidf_w2v)[:,-1]
11 fpr_te, tpr_te, thresholds_te = roc_curve(Y_te,y_prob_te)
12 plt.plot(fpr_te, tpr_te)
13
14 AUC_tr_tfidf_w2v = roc_auc_score(Y_train, y_prob_tr)
15 AUC_te_tfidf_w2v = roc_auc_score(Y_te, y_prob_te)
16
17 plt.grid()
18 plt.legend( ('Train ROC & Train AUC -- '+str(np.round(AUC_tr_tfidf_w2v, decimal_places=4))+' , '
19            'Test ROC & Test AUC -- '+str(np.round(AUC_te_tfidf_w2v, decimal_places=4))+' , '
20            'Train AUC -- '+str(np.round(AUC_tr_tfidf_w2v, decimal_places=4))+' , '
21            'Test AUC -- '+str(np.round(AUC_te_tfidf_w2v, decimal_places=4))+' ')
22 plt.xlabel('fpr values')
23 plt.ylabel('Tpr values')
24 plt.title('Roc Curves for Train and Test data.')
25 plt.show()

```



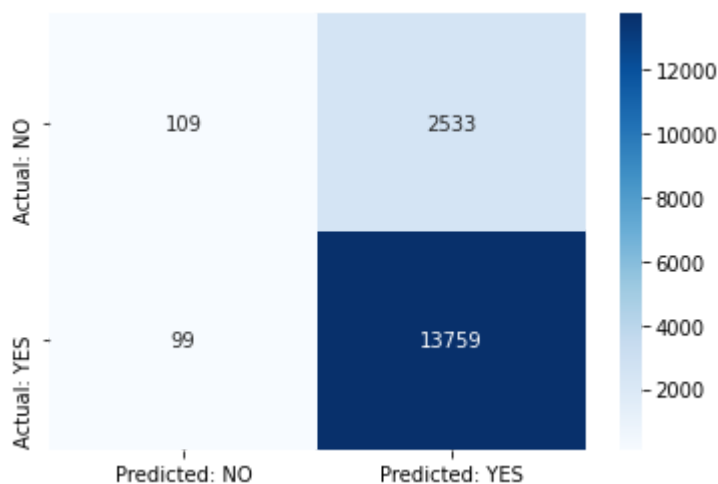
```

1 y_pred_tfidf_w2v=clf.predict(X_te_tfidf_w2v)
2 cnf_mtrx_tfidf_w2v = confusion_matrix(Y_te, y_pred_tfidf_w2v)
3
4 print('Confusion matrix for the Test data.')
5 #print(cnf_mtrx_tfidf)
6 x=pd.DataFrame(cnf_mtrx_tfidf_w2v, columns = ["Predicted: NO","Predicted: YES"]
7 sns.heatmap(x, annot=True,fmt="d", cmap='Blues')
8 plt.show()

```



Confusion matrix for the Test data.



▼ 3. Summary

```

1 # http://zetcode.com/python/prettytable/
2 from prettytable import PrettyTable
3
4 x = PrettyTable()
5 x.field_names = ["Vectorizer","Model", "Hyper parameters (depth and no of est:
6
7 for i in range(1):
8     x.add_row(['TFIDF', "GBDT", (best depth,best n estimators) , np.round(AUC

```

```
9     x.add_row(["TFIDF-W2V", "GBDT" , (best_depth_tfidf_w2v,best_n_estimators_tfidf_w2v)])
10 print(x)
```

Vectorizer	Model	Hyper parameters (depth and no of estimators)	Train
TFIDF	GBDT	(4, 120)	0.8
TFIDF-W2V	GBDT	(5, 120)	0.90