# GBDT Assignment

1. **Apply GBDT on these feature sets**

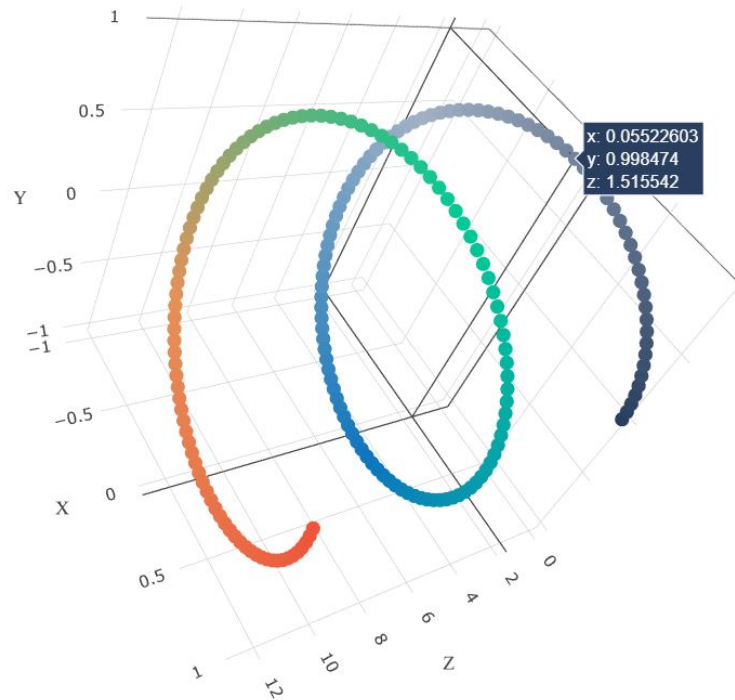   - Set 1: categorical(instead of one hot encoding, try [response coding (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(TFIDF)+ preprocessed_eassay (TFIDF)+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 (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/)](https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/): use probability values), numerical features + project_title(TFIDF W2V)+ preprocessed_eassay (TFIDF W2V)
   - Here in response encoding you need to apply the **laplase smoothing** value for test set. Laplase smoothing means, If test point is present in test but not in train then you need to apply default 0.5 as probability value for that data point (Refer the Response Encoding Image from above cell)
   - Please use atleast **35k** data points

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

   - Find the best hyper parameter which will give the maximum [AUC (https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/receiver-operating-characteristic-curve-roc-curve-and-auc-1/)](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/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*
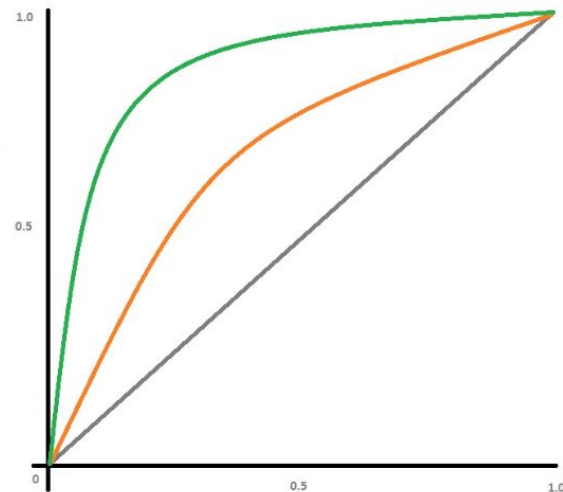
## or

- 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 (https://seaborn.pydata.org/generated/seaborn.heatmap.html)](https://seaborn.pydata.org/generated/seaborn.heatmap.html) with rows as **n_estimators**, columns as **max_depth**, and values inside the cell representing **AUC Score**

- You 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 calcualted 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 = ??        |

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

```
+----------------+-----------+------------------+----------+
|   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   |
+----------------+-----------+------------------+----------+
```

# Few Notes

1. Use atleast 35k data points
2. Use classifier.Predict_proba() method instead of predict() method while calculating roc_auc scores
3. Be sure that you are using laplase smoothing in response encoding function. Laplase smoothing means applying the default (0.5) value to test data if the test data is not present in the train set

```
In [34]:   1  import pandas as pd
           2  from sklearn.model_selection import train_test_split
           3  from sklearn.feature_extraction.text import TfidfVectorizer
           4  from sklearn.feature_extraction.text import CountVectorizer
           5  import numpy as np
           6  import seaborn as sns
           7  from sklearn.model_selection import GridSearchCV
           8  from sklearn.tree import DecisionTreeClassifier
           9  import matplotlib.pyplot as plt
          10  import pdb
          11  import lightgbm as lgb
          12
          13
          14  from sklearn.ensemble import GradientBoostingClassifier
          15
          16  import tqdm
          17  from sklearn.metrics import confusion_matrix
          18  from sklearn.metrics import accuracy_score
          19  from sklearn.metrics import roc_auc_score
          20  from sklearn.metrics import roc_curve
```

or else , you can use below code

```
In [2]:  1  import nltk
         2  from nltk.sentiment.vader import SentimentIntensityAnalyzer
         3
         4  # nltk.download('vader_lexicon')
         5
         6  sid = SentimentIntensityAnalyzer()
         7
         8  sample_sentence_1='I am happy.'
         9  ss_1 = sid.polarity_scores(sample_sentence_1)
        10  print('sentiment score for sentence 1',ss_1)
        11
        12  sample_sentence_2='I am sad.'
        13  ss_2 = sid.polarity_scores(sample_sentence_2)
        14  print('sentiment score for sentence 2',ss_2)
        15
        16  sample_sentence_3='I am going to New Delhi tommorow.'
        17  ss_3 = sid.polarity_scores(sample_sentence_3)
        18  print('sentiment score for sentence 3',ss_3)
        19
```

```
sentiment score for sentence 1 {'neg': 0.0, 'neu': 0.213, 'pos': 0.787, 'compound': 0.5719}
sentiment score for sentence 2 {'neg': 0.756, 'neu': 0.244, 'pos': 0.0, 'compound': -0.4767}
sentiment score for sentence 3 {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0}
```

# Decision Tree

## Task - 1

### 1.1 Loading Data

```
In [3]:  1  #make sure you are loading atleast 50k datapoints
         2  #you can work with features of preprocessed_data.csv for the assignment.
         3  data = pd.read_csv('preprocessed_data.csv', nrows = 35000)
```

```
In [4]:   1  data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35000 entries, 0 to 34999
Data columns (total 9 columns):
 #   Column                                          Non-Null Count  Dtype
---  ------                                          --------------  -----
 0   school_state                                    35000 non-null  object
 1   teacher_prefix                                  35000 non-null  object
 2   project_grade_category                          35000 non-null  object
 3   teacher_number_of_previously_posted_projects    35000 non-null  int64
 4   project_is_approved                             35000 non-null  int64
 5   clean_categories                                35000 non-null  object
 6   clean_subcategories                             35000 non-null  object
 7   essay                                           35000 non-null  object
 8   price                                           35000 non-null  float64
dtypes: float64(1), int64(2), object(6)
memory usage: 2.4+ MB
```

```
In [5]:   1  data.head(2)
```

Out[5]:

| | school_state | teacher_prefix | project_grade_category | teacher_number_of_previously_posted_projects | project_is_approved | clean_categories | clean_s |
|---|---|---|---|---|---|---|---|
| **0** | ca | mrs | grades_prek_2 | 53 | 1 | math_science | a hea |
| **1** | ut | ms | grades_3_5 | 4 | 1 | specialneeds | |

```
In [6]:  1  # saperating yi's
         2  data_target = data['project_is_approved'].values
         3  data1 = data
         4  data = data.drop(['project_is_approved'],axis=1)
         5
         6
```

```
In [7]:  1  # data distribution
         2
         3  data1['project_is_approved'].value_counts()
```

```
Out[7]:  1    29629
         0     5371
         Name: project_is_approved, dtype: int64
```

```
In [8]:  1  data_target.shape
```

```
Out[8]:  (35000,)
```

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

```
In [9]:  1  # 1. Split your data. train test split
         2  from sklearn.model_selection import train_test_split
         3
         4
         5  X_train, X_test, y_train, y_test = train_test_split(data,data_target, test_size=0.20, stratify=data_target)
         6  X_train1, X_cv, y_train1, y_cv = train_test_split(X_train, y_train, test_size=0.20, stratify=y_train)
```

```
In [10]:  1  X_train.shape, X_test.shape, data.shape,y_train.shape
          2
```

```
Out[10]:  ((28000, 8), (7000, 8), (35000, 8), (28000,))
```

### Text Feature - " essay "

**3. perform tfidf vectorization.**

```
In [11]:   1   # tfidf for feature "essay" used for set-1
           2
           3   vectorizer = TfidfVectorizer(ngram_range= (1,4), max_features=5000, min_df=10)
           4   vectorizer_ess = vectorizer.fit(X_train['essay'].values)
           5
           6   X_tr_essay_tfidf = vectorizer_ess.transform(X_train['essay'].values)
           7   X_te_essay_tfidf = vectorizer_ess.transform(X_test['essay'].values)
           8   X_cv_essay_tfidf = vectorizer_ess.transform(X_cv['essay'].values)
           9
          10   print("After vectorizations")
          11   print(X_tr_essay_tfidf.shape, y_train.shape)
```

After vectorizations
(28000, 5000) (28000,)

```
In [12]:   1   # X_tr_essay_tfidf = X_tr_essay_tfidf.toarray()
           2   # X_tr_essay_tfidf = X_te_essay_tfidf.toarray()
           3   # X_cv_essay_tfidf = X_cv_essay_tfidf.toarray()
           4   type(X_tr_essay_tfidf)
```

Out[12]:  scipy.sparse.csr.csr_matrix


**4. perform tfidf w2v vectorization of text data.**


**loading glove words**


```
In [13]:   1   #please use below code to load glove vectors
           2   import pickle
           3   with open('glove_vectors', 'rb') as f:
           4       model = pickle.load(f)
           5       glove_words =  set(model.keys())
```

```python
# part2 - average Word2Vec
# compute tfidf-word2vec    for each review.
def cal_tfidf_word2vec(essay,vectorizer):  # essay is series

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


    tfidf_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list

    for sentence in essay.values: # for each review/sentence
        vector = np.zeros(300) # as word vectors are of zero length
        tf_idf_weight =0; # num of words with a valid vector in the sentence/review

        for word in sentence.split(): # for each word in a review/sentence
            if (word in glove_words) and (word in tfidf_words):
                vec = model[word] # getting the vector for each word

                # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sen
                tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for
                vector += (vec * tf_idf) # calculating tfidf weighted w2v
                tf_idf_weight += tf_idf

        if tf_idf_weight != 0:
            vector /= tf_idf_weight
        tfidf_w2v_vectors.append(vector)


    return np.array(tfidf_w2v_vectors)
```

```
In [15]:  1  %%time
          2
          3  X_tr_essay_tfidf_w2v = cal_tfidf_word2vec(X_train['essay'], vectorizer_ess)
          4  X_te_essay_tfidf_w2v = cal_tfidf_word2vec(X_test['essay'], vectorizer_ess)
          5  X_cv_essay_tfidf_w2v = cal_tfidf_word2vec(X_cv['essay'], vectorizer_ess)
          6
          7  print("essay feature as tfidf_w2v :",X_tr_essay_tfidf_w2v.shape, X_te_essay_tfidf_w2v.shape, X_cv_essay_tfidf_w2v.sh
          8
```

```
essay feature as tfidf_w2v : (28000, 300) (7000, 300) (5600, 300)
Wall time: 1min 54s
```

```
In [16]:  1  print("essay feature as tfidf_w2v :",X_tr_essay_tfidf_w2v.shape, X_te_essay_tfidf_w2v.shape, X_cv_essay_tfidf_w2v.sh
          2
```

```
essay feature as tfidf_w2v : (28000, 300) (7000, 300) (5600, 300)
```

# BOW for 5 features

5. perform encoding of categorical features.


**response coding**

```python
def rc_table(data,y):
    word_dict_of_prob = {}
    table = []
    c0,c1 = [],[]   # count of x for both class
    uni_x = list(set(data))

    data_y_0 = [x for i,x in enumerate(data) if y[i] == 0 ]
    data_y_1 = [x  for i,x in enumerate(data) if y[i] == 1 ]

    for x in uni_x:

        c0.append(data_y_0.count(x))
        c1.append(data_y_1.count(x))

    for i,x in enumerate(uni_x):
        try:
            prob_x_0 = c0[i]/(c0[i]+c1[i])
            prob_x_1 = c1[i]/(c0[i]+c1[i])
            table.append([prob_x_0, prob_x_1])
        except:
            pdb.set_trace()

    for i, x in enumerate(data):

        for i2, j in enumerate(uni_x):
            if x == j :
                word_dict_of_prob[x] = [table[i2][0],table[i2][1]]


    return word_dict_of_prob

def response_coding(word_dict_of_prob, data):  # taking dictionary type and series type value

    rc_data = []
    uni_words = word_dict_of_prob.keys()

    for i in data:
        if i in uni_words:
            rc_data.append(word_dict_of_prob[i])

        else :
```

```
42            rc_data.append([0.5,0.5])
43
44     return np.array(rc_data)
45
46
47
```

In [18]:
```
1 # word_dict_of_prob = rc_table(X_train['school_state'].values[:100], y_train[:100])
2 # # word_dict_of_prob
3 # responce_codeing(word_dict_of_prob, X_train['school_state'].values[100:1000])
```

In [19]:
```
1 # data.head()
```

```python
In [20]:   1  #responce encoding of 5 feature in data
           2
           3  # responce encoding of School state
           4
           5  word_dict_of_prob = rc_table(X_train['school_state'].values, y_train) # return word_dict_of_prob for both class
           6
           7  X_tr_state_rc = response_coding(word_dict_of_prob, X_train['school_state'].values)
           8  X_te_state_rc = response_coding(word_dict_of_prob, X_test['school_state'].values)
           9  X_cv_state_rc = response_coding(word_dict_of_prob, X_cv['school_state'].values)
          10
          11
          12  # responce encoding of teacher prefix
          13  word_dict_of_prob = rc_table(X_train['teacher_prefix'].values, y_train)
          14
          15  X_tr_teacher_rc = response_coding(word_dict_of_prob, X_train['teacher_prefix'].values)
          16  X_te_teacher_rc = response_coding(word_dict_of_prob, X_test['teacher_prefix'].values)
          17  X_cv_teacher_rc = response_coding(word_dict_of_prob, X_cv['teacher_prefix'].values)
          18
          19  # responce encoding of project grade category
          20  word_dict_of_prob = rc_table(X_train['project_grade_category'].values, y_train)
          21
          22  X_tr_grade_rc = response_coding(word_dict_of_prob, X_train['project_grade_category'].values)
          23  X_te_grade_rc = response_coding(word_dict_of_prob, X_test['project_grade_category'].values)
          24  X_cv_grade_rc = response_coding(word_dict_of_prob, X_cv['project_grade_category'].values)
          25
          26  # responce encoding of clean_subcategories
          27
          28  word_dict_of_prob = rc_table(X_train['clean_categories'].values, y_train)
          29
          30  X_tr_cate_rc = response_coding(word_dict_of_prob, X_train['clean_categories'].values)
          31  X_te_cate_rc = response_coding(word_dict_of_prob, X_test['clean_categories'].values)
          32  X_cv_cate_rc = response_coding(word_dict_of_prob, X_cv['clean_categories'].values)
          33
          34
          35
          36  # responce encoding of clean_subcategories
          37
          38  word_dict_of_prob = rc_table(X_train['clean_subcategories'].values, y_train)
          39
          40  X_tr_sub_cate_rc = response_coding(word_dict_of_prob, X_train['clean_subcategories'].values)
          41  X_te_sub_cate_rc = response_coding(word_dict_of_prob, X_test['clean_subcategories'].values)
```

```
42  X_cv_sub_cate_rc = response_coding(word_dict_of_prob, X_cv['clean_subcategories'].values)
43
44
45
46  # 5-responce encoding features
47  X_tr_state_rc.shape, X_tr_teacher_rc.shape, X_tr_grade_rc.shape, X_tr_cate_rc.shape, X_tr_sub_cate_rc.shape
```

Out[20]:  ((28000, 2), (28000, 2), (28000, 2), (28000, 2), (28000, 2))

**Normalization**

6. perform encoding of numerical features

```
In [21]:  1  # normalization of 2 numerical features - price and project_posted_by teacher
          2  from sklearn.preprocessing import Normalizer
          3
          4
          5  # normalization of numerical features - price
          6  normalizer = Normalizer()
          7  normalizer.fit(data['price'].values.reshape(1,-1))     # normalizer.fit(X_train['price'].values) # showing error as
          8
          9  X_tr_price_norm = normalizer.transform(X_train['price'].values.reshape(-1,1))
         10  X_te_price_norm = normalizer.transform(X_test['price'].values.reshape(-1,1))
         11  X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(-1,1))
         12
         13
         14
         15  # normalization of numerical features - teacher_number_of_previously_posted_projects
         16  normalizer = Normalizer()
         17  normalizer.fit(data['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
         18
         19  X_tr_projects_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,
         20  X_te_projects_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1
         21  X_cv_projects_norm = normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
         22
         23
         24  X_tr_price_norm.shape,X_tr_projects_norm.shape
```

Out[21]:  ((28000, 1), (28000, 1))

**sentiment analyzer for text data**

   1. calculate sentiment scores for the essay feature

```
In [22]:   1  # sentiment analyzer
           2
           3  sid = SentimentIntensityAnalyzer()
           4
           5  def sentiment_score(analyzer_obj, data):
           6      sentiment_score = []
           7
           8      for sentence in data.values:
           9
          10          senti = sid.polarity_scores(sentence)
          11          #sentiment_score = np.append(sentiment_score,list(senti.values()))
          12          sentiment_score.append(list(senti.values()))
          13
          14
          15      return np.array(sentiment_score)
          16
          17  X_tr_sentiment_score = sentiment_score(sid, X_train['essay'])
          18  X_te_sentiment_score = sentiment_score(sid, X_test['essay'])
          19  X_cv_sentiment_score = sentiment_score(sid, X_cv['essay'])
          20
          21  X_tr_sentiment_score.shape ,X_te_sentiment_score.shape , X_cv_sentiment_score.shape
```

Out[22]:  ((28000, 4), (7000, 4), (5600, 4))

**making 2 set with diff-2 feature**

Set 1: categorical, numerical features + preprocessed_essay (TFIDF W2V) + Sentiment scores(preprocessed_essay) + responce coding

Set 2: categorical, numerical features + preprocessed_essay (TFIDF ) + Sentiment scores(preprocessed_essay) + responce coding

```
In [23]:    1  from scipy.sparse import hstack
            2  # for set 1 - categorical, numerical features + preprocessed_essay (TFIDF) + Sentiment scores
            3
            4  S1_X_tr = hstack((X_tr_state_rc, X_tr_teacher_rc, X_tr_grade_rc, X_tr_cate_rc, X_tr_sub_cate_rc, X_tr_sentiment_scor
            5  S1_X_te = hstack((X_te_state_rc, X_te_teacher_rc, X_te_grade_rc, X_te_cate_rc, X_te_sub_cate_rc, X_te_sentiment_scor
            6  S1_X_cv = hstack((X_cv_state_rc, X_cv_teacher_rc, X_cv_grade_rc, X_cv_cate_rc, X_cv_sub_cate_rc, X_cv_sentiment_scor
            7
            8
            9
           10  print("Final Data matrix for set 1")
           11  print(S1_X_tr.shape, S1_X_te.shape, S1_X_cv.shape)
           12  print("="*100)
           13
           14
           15
```

```
Final Data matrix for set 1
(28000, 5016) (7000, 5016) (5600, 5016)
====================================================================================================
```

```
In [24]:   1
           2
           3   # for set 2 -  categorical, numerical features + preprocessed_essay (TFIDF W2v) + Sentiment scores
           4
           5
           6   # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
           7   S2_X_tr = np.hstack((X_tr_state_rc, X_tr_teacher_rc, X_tr_grade_rc, X_tr_cate_rc, X_tr_sub_cate_rc, X_tr_sentiment_s
           8   S2_X_te = np.hstack((X_te_state_rc, X_te_teacher_rc, X_te_grade_rc, X_te_cate_rc, X_te_sub_cate_rc, X_te_sentiment_s
           9   S2_X_cv = np.hstack((X_cv_state_rc, X_cv_teacher_rc, X_cv_grade_rc, X_cv_cate_rc, X_cv_sub_cate_rc, X_cv_sentiment_s
          10
          11   print("Final Data matrix for set 2")
          12   print(S1_X_tr.shape, S1_X_te.shape, S1_X_cv.shape)
          13   print("="*100)
          14
          15
```

```
Final Data matrix for set 2
(28000, 5016) (7000, 5016) (5600, 5016)
====================================================================================================
```

# modeling on set1 feature

**9. Perform hyperparameter tuning and plot either heatmap or 3d plot.**

```
In [26]:   1
```

```
In [27]:   1   y_train[y_train==1].shape, y_train[y_train==0].shape, S1_X_tr.shape, y_train.shape
```

Out[27]: ((23703,), (4297,), (28000, 5016), (28000,))

```
In [ ]:   1
```

```
In [28]:   1   %%time
           2   # hyperparameter with sklearn
           3   # took long time to execute
           4
           5   parameters = {'max_depth': [1, 3, 10, 30], 'min_child_samples':[10, 20, 40, 80] }
           6
           7   model = lgb.LGBMClassifier( class_weight= 'balanced' )
           8   clf_hyper = GridSearchCV(model, parameters, scoring ="roc_auc", n_jobs = -1).fit(S1_X_tr,y_train)
           9
          10   print('best parameter',clf_hyper.best_params_)
          11
          12
```

```
best parameter {'max_depth': 30, 'min_child_samples': 80}
Wall time: 16min 9s
```

```
In [29]:   1   # clf_hyper.predict(X_train)
           2   S1_X_tr.shape
```

Out[29]:  (28000, 5016)

```
In [29]:   1   # without reg_lambda model is overfitting , hence we manually hyper-tune reg_lambda parameter
           2
           3
           4   print('best parameter',clf_hyper.best_params_)
           5   train_auc = roc_auc_score(y_train, clf_hyper.predict(S1_X_tr))
           6   test_auc = roc_auc_score(y_test, clf_hyper.predict(S1_X_te))
           7
           8   print(train_auc, test_auc)
```

```
best parameter {'max_depth': 30, 'min_child_samples': 80}
0.8607903222830907 0.6569168409137902
```

## best model

{'max_depth': 30, 'min_data_in_leaf': 40}

```
In [33]:   1  #shape of dataset
           2
           3  S1_X_tr.shape
```

Out[33]:  (28000, 5016)

```
In [42]:   1
           2
           3  odel with lightGbm  ref - https://www.kaggle.com/code/prashant111/lightgbm-classifier-in-python/notebook ; https://lightg
           4  =lgb.LGBMClassifier(reg_lambda = 10000, max_depth =30, min_child_samples = 80 ,class_weight= 'balanced',  n_estimators=
           5  .fit(S1_X_tr,y_train)
           6
           7  ing class
           8  d = best_clf.predict(S1_X_tr)
           9  d = best_clf.predict(S1_X_te)
          10
          11  ting probability
          12  b = best_clf.predict_proba(S1_X_tr)
          13  b = best_clf.predict_proba(S1_X_te)
          14
          15  c = roc_auc_score(y_train, tr_y_pred)
          16  = roc_auc_score(y_test, te_y_pred)
          17
          18  rain auc :-",train_auc,"test auc :-",test_auc)
```

train auc :- 0.6831163724946182 test auc :- 0.626181156674089
Wall time: 27.3 s

```
In [43]:   1  print("train auc :-",train_auc,"test auc :-",test_auc)
```

train auc :- 0.6831163724946182 test auc :- 0.626181156674089

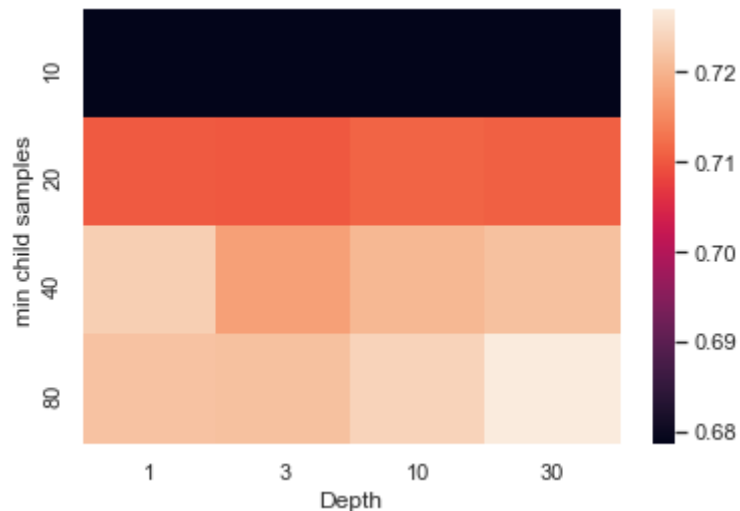## 1.3 Representation of results

**Heatmap on mean_test_score**

    corrensponding to the each hyperparameter

```
In [36]:   1  # reshaping 'mean_test_score' in (4,5) for heatmap
           2  cv_accu = clf_hyper.cv_results_['mean_test_score']
           3  cv_accu = np.reshape(cv_accu, (4,4))
           4  cv_accu.shape
```

Out[36]: (4, 4)

```
In [132]:  1
           2  np.random.seed(0)
           3  sns.set_theme()
           4
           5  pd_cv_accu = pd.DataFrame(cv_accu , columns = [1, 3, 10, 30], index =  [10, 20, 40, 80])
           6
           7  # with ->  vmin= min(cv_accu), vmax= max(cv_accu)
           8  ax = sns.heatmap(pd_cv_accu, vmin=np.min(cv_accu), vmax=np.max(cv_accu))
           9  plt.xlabel("Depth")
          10  plt.ylabel("min child samples ")
          11  plt.show()
```
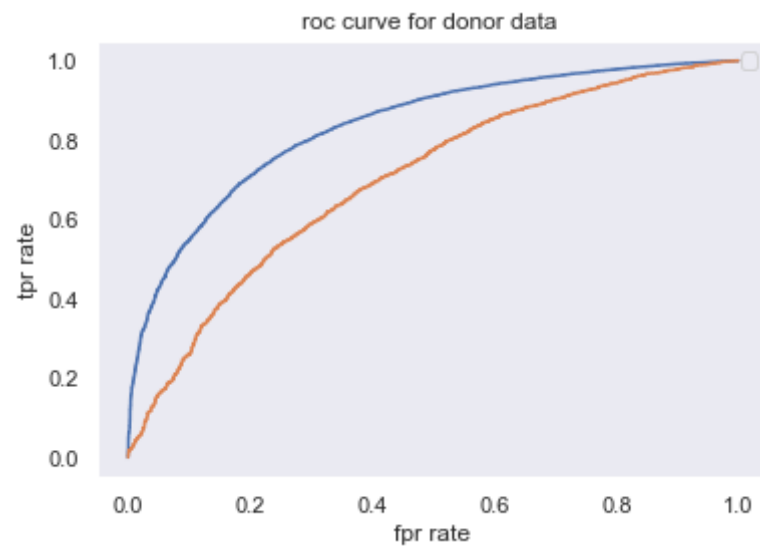


**ROC plot for train and test data**

```
In [39]:   1
           2  tr_fpr,tr_tpr,_ = roc_curve(y_train,tr_y_prob[:,1])
           3  te_fpr,te_tpr,_ = roc_curve(y_test,te_y_prob[:,1])
           4
```

```
In [40]:   1
           2
           3  plt.plot(tr_fpr,tr_tpr)
           4  plt.plot(te_fpr,te_tpr)
           5  plt.xlabel("fpr rate")
           6  plt.ylabel("tpr rate")
           7
           8  plt.title("roc curve for donor data ")
           9  plt.legend()
          10  plt.grid()
          11  plt.show()
```

No handles with labels found to put in legend.



**confusion matrix**

```
In [41]:  1  ar = confusion_matrix(y_train,tr_y_pred)
          2  ar2 = confusion_matrix(y_test,te_y_pred)
          3  tn, fp, fn, tp = confusion_matrix(y_train,tr_y_pred).ravel()
          4  ar,ar2
```

Out[41]: (array([[ 3347,   950],
                [ 6343, 17360]], dtype=int64),
          array([[ 633,   441],
                [1792, 4134]], dtype=int64))

**12. Find all the false positive data points and plot wordcloud of essay text and pdf of teacher_number_of_previously_posted_projects.**

```
In [42]:  1  false_positive = []
          2
          3  # false_positive.append    if x= 0 and y=1
          4
          5  for idx,x in enumerate(zip(y_train,tr_y_pred)):
          6
          7      y_true,y_pred = x
          8
          9      if y_true== 0 and y_pred == 1 :
         10          false_positive.append(idx)
         11
         12  false_positive = np.array(false_positive)
         13  false_positive.shape,type(false_positive)
```
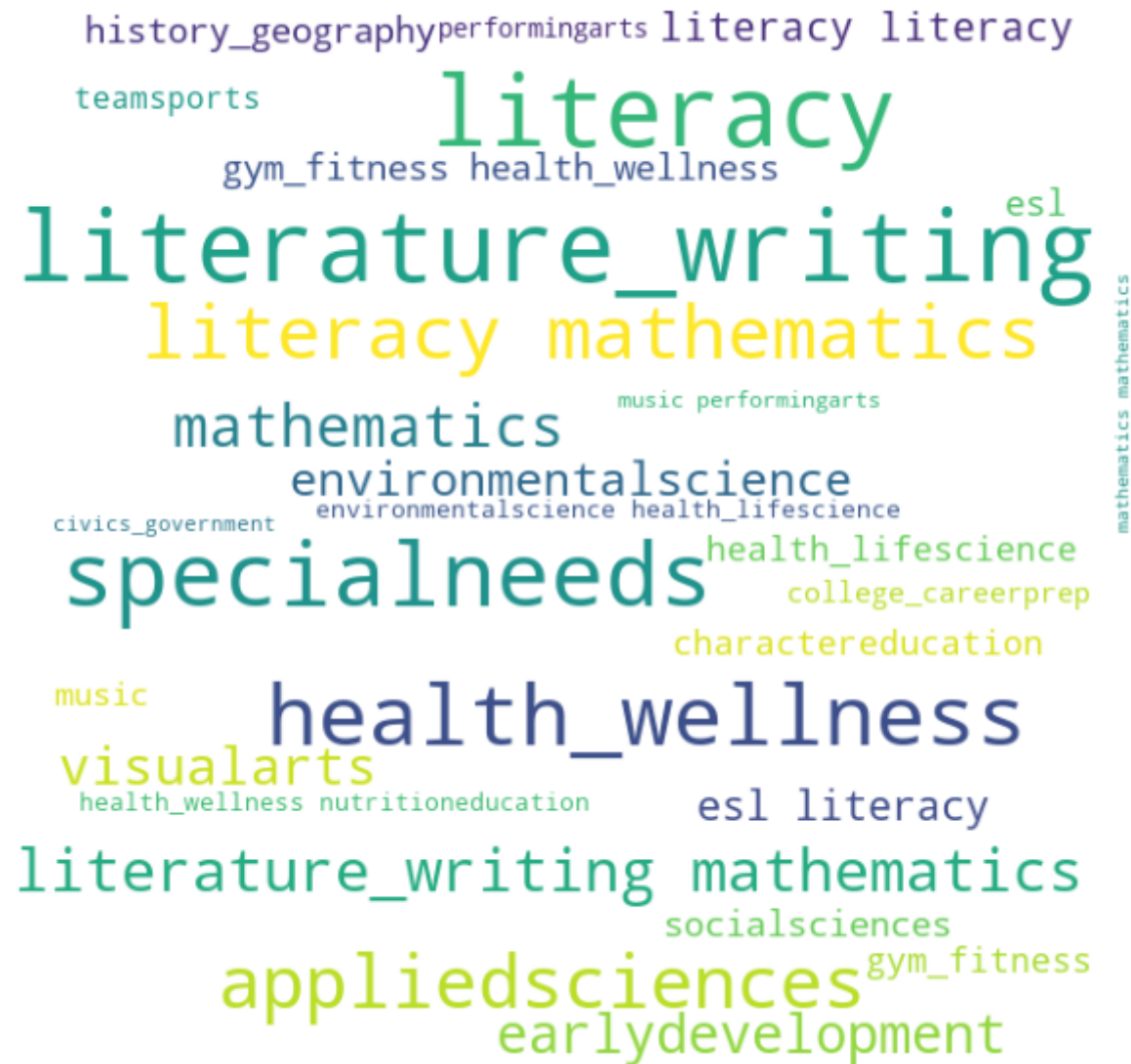
Out[42]: ((950,), numpy.ndarray)

**word cloud**

```python
In [43]:  1  # Python program to generate WordCloud for all false positive datapoints'essay feacture
          2
          3  # importing all necessary modules
          4  from wordcloud import WordCloud, STOPWORDS
          5
          6  # Reads 'Youtube04-Eminem.csv' file
          7  df = false_positive
          8
          9  comment_words = ''
         10  stopwords = set(STOPWORDS)
         11  print("------------ WordCloud for all false pasitive(- negetive words actually classified as positive)words occur --
         12  # iterate through the csv file
         13  for v in df:
         14      val = data1.iloc[v,6]
         15      # typecaste each val to string
         16      val = str(val)
         17
         18      # split the value
         19      tokens = val.split()
         20
         21      # Converts each token into Lowercase
         22      for i in range(len(tokens)):
         23          tokens[i] = tokens[i].lower()
         24
         25      comment_words += " ".join(tokens)+" "
         26
         27  wordcloud = WordCloud(width = 600, height = 600,
         28                  background_color ='white',
         29                  stopwords = stopwords,
         30                  min_font_size = 10).generate(comment_words)
         31
         32  # plot the WordCloud image
         33  plt.figure(figsize = (8, 8), facecolor = None)
         34  plt.imshow(wordcloud)
         35  plt.axis("off")
         36  plt.tight_layout(pad = 0)
         37
         38  plt.show()
         39
```

here - these are the negetive points by which we consider poject is_approved(=1)

**observation**

```
 - projects related to helth_wellness, applidsciences, early_development, literature_writing, literacy, heath are
 pridicted positive(accepted of funders)  but it is actually negative.
```
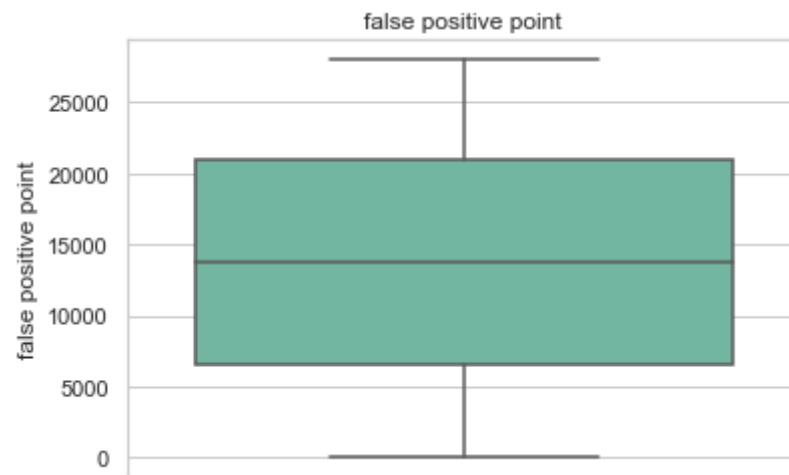
**Box plot**

```
-  with the `price` of these `false positive data points`
```

In [45]:
```
1
2  sns.set_theme(style="whitegrid")
3  tips = sns.load_dataset("tips")
4  ax = sns.boxplot(y=false_positive, orient = "h", palette="Set2")
5  plt.ylabel('false positive point')
6  plt.title("false positive point")
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\_core.py:1296: UserWarning: Horizontal orientation ignored with only `y` specified.
  warnings.warn(single_var_warning.format("Horizontal", "y"))
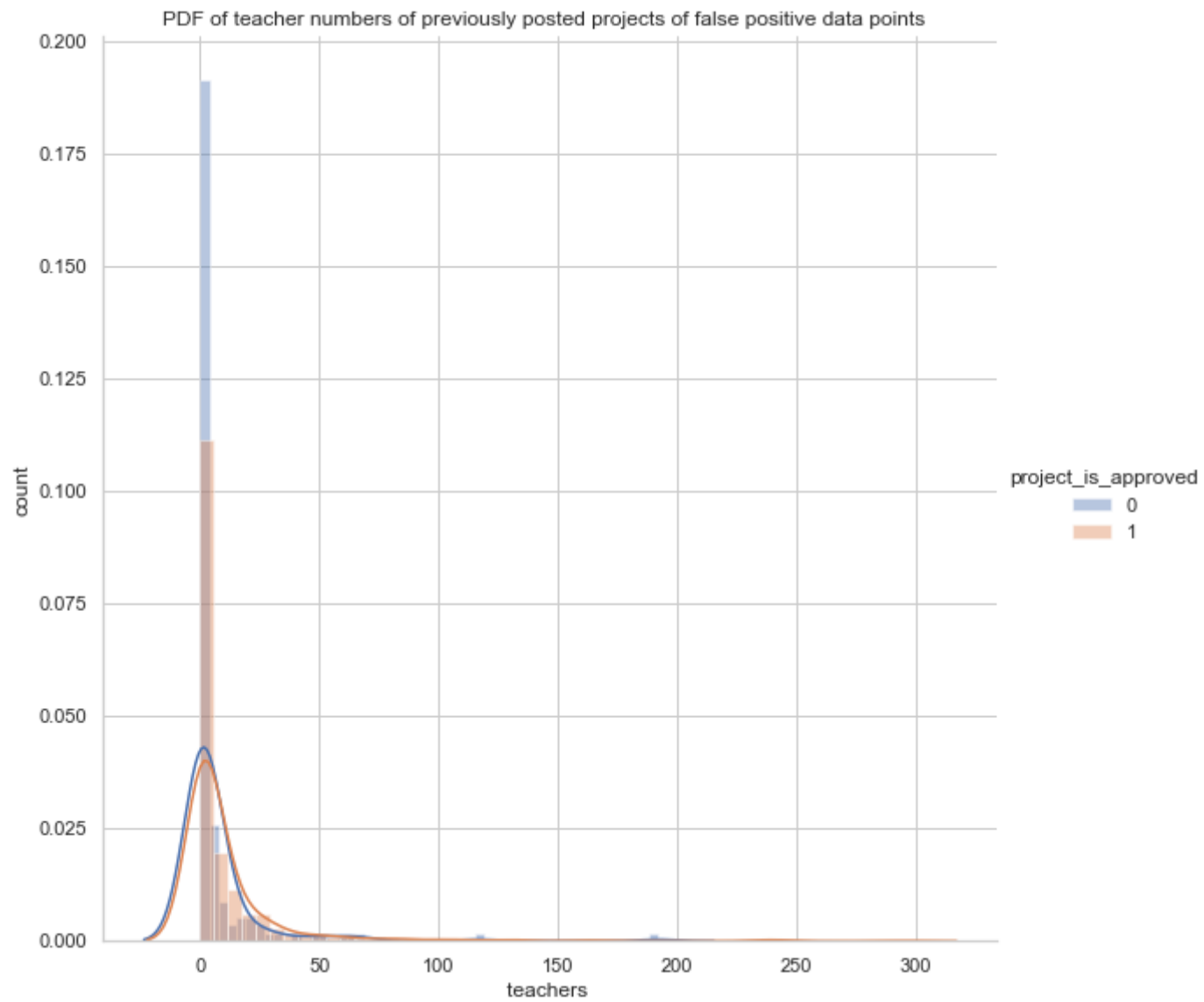
Out[45]: Text(0.5, 1.0, 'false positive point')



**pdf**

In [46]:
```
1  false_pt_df = data1.loc[false_positive]
2
```

```python
import warnings
# adding column


# Op_Year

warnings.filterwarnings('ignore')

sns.FacetGrid(false_pt_df, hue="project_is_approved", size=8) \
    .map(sns.distplot, "teacher_number_of_previously_posted_projects") \
    .add_legend()\
    .set(xlabel='teachers ', ylabel='count');
plt.title("PDF of teacher numbers of previously posted projects of false positive data points  ")
plt.show();
```

PDF of teacher numbers of previously posted projects of false positive data points

**observation:**

> most of teacher did not previously posted poject classied as positive pts by the model

## set_2 features

feature - Set 2: categorical, numerical features + preprocessed_essay (TFIDF W2V) + Sentiment scores(preprocessed_essay)

```
In [52]:   1  import time
           2
```

```
In [53]:   1  %%time
           2  # 4. perform hyperparameter tuning and plot either heatmap or 3d plot.
           3
           4
           5
           6  parameters = {'max_depth': [10, 30, 50, 70], 'min_child_samples':[10, 20, 40, 60], 'reg_lambda':[.001,.01,1,100] }
           7
           8  model2 = lgb.LGBMClassifier( class_weight= 'balanced', n_estimators= 150)
           9
          10
          11  clf_hyper = GridSearchCV( model2, parameters, scoring = 'roc_auc', n_jobs = -1 , cv = 10 ).fit(S2_X_tr,y_train)
          12
          13
          14
```

Wall time: 1h 20min 19s

```
In [54]:   1  print('best parameter',clf_hyper.best_params_)
           2
```

best parameter {'max_depth': 10, 'min_child_samples': 60, 'reg_lambda': 100}

**best model**

```
In [63]:    1  %%time
            2
            3  best_m2 = lgb.LGBMClassifier(reg_lambda = 1000, max_depth =clf_hyper.best_params_['max_depth'],min_child_samples = c
            4
            5  best_m2.fit(S2_X_tr,y_train)
            6
            7  tr_y_pred = best_m2.predict(S2_X_tr)
            8  te_y_pred = best_m2.predict(S2_X_te)
            9
           10  tr_accu_auc = roc_auc_score(y_train,tr_y_pred)
           11  te_accu_auc = roc_auc_score(y_test,te_y_pred)
           12
           13  tr_accu_auc,te_accu_auc
```

Wall time: 7.18 s

Out[63]:  (0.7289158224031622, 0.6147507653361037)


## Task - 2

```
In [27]:    1  %%time
            2  # 1. write your code in following steps for task 2
            3
            4  print("shape of dataset set1",S1_X_tr.shape)
            5
            6  S1_X_tr_cpy = S1_X_tr.toarray()
            7  imp_features = best_clf.feature_importances_
            8
            9  # 2. select all non zero features
           10  list_idx = [i for i,x in enumerate(imp_features) if x>0  ]
           11
           12  # 3. Update your dataset i.e. X_train,X_test and X_cv so that it contains all rows and only non zero features
           13  data_set_compress = S1_X_tr_cpy[:,list_idx]
           14
           15
           16
           17
           18
           19  # split data
           20  X, X_cv, y, y_cv = train_test_split(data_set_compress, y_train, test_size=.20, stratify= y_train)
           21
           22  print("after compressing",data_set_compress.shape)
```

```
shape of dataset set1 (28000, 5016)
after compressing (28000, 869)
Wall time: 1.12 s
```

```
In [28]:   1  %%time
           2  # 4. perform hyperparameter tuning and plot either heatmap or 3d plot.
           3
           4  parameters = {'max_depth': [1, 3, 10, 30], 'min_child_samples':[10, 20, 40, 80] }
           5
           6  parameters = {'max_depth': [10, 30, 50, 70], 'min_child_samples':[10, 20, 40, 60],  }
           7
           8  model = lgb.LGBMClassifier( class_weight= 'balanced', n_estimators= 150)
           9
          10
          11  clf_hyper3 = GridSearchCV( model, parameters, scoring = 'roc_auc', n_jobs = 1 , cv = 10 ).fit(X,y)
          12
          13
          14
          15
          16
          17  print('best parameter',clf_hyper3.best_params_)
          18
```

```
best parameter {'max_depth': 30, 'min_child_samples': 40}
Wall time: 23min 11s
```

```
best parameter {'max_depth': 10, 'min_child_samples': 20}
Wall time: 29min 50s
```

## 5. Fit the best model. Plot ROC AUC curve and confusion matrix similar to model 1.

```python
# best_m3 = GradientBoostingClassifier(max_depth = clf_hyper2.best_params_['max_depth'], min_samples_split =clf_hype
# best_m3 = GradientBoostingClassifier(max_depth = 30, min_samples_split =40, reg_lambda = 1000)
best_m3 = lgb.LGBMClassifier(reg_lambda = 10000, max_depth =30,min_child_samples = 40 ,class_weight= 'balanced', n_e

best_m3.fit(X,y)

tr_y_pred = best_m3.predict(X)
te_y_pred = best_m3.predict(X_cv)

tr_accu_auc = roc_auc_score(y,tr_y_pred)
te_accu_auc = roc_auc_score(y_cv,te_y_pred)

tr_accu_auc,te_accu_auc
```
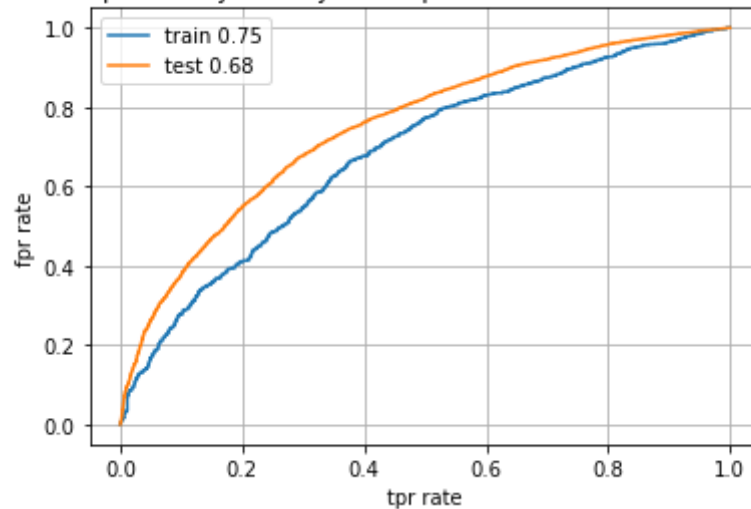
Out[35]: (0.6876605696006691, 0.6399886654918001)

**ROC Curve**

```python
# predicting probability
tr_y_prob = best_m3.predict_proba(X)
te_y_prob = best_m3.predict_proba(X_cv)

# calculating roc curve
fpr1, tpr1,_ = roc_curve(y, tr_y_prob[:,1])
fpr2, tpr2,_ = roc_curve(y_cv, te_y_prob[:,1])
```

```
In [38]:    1  from sklearn.metrics import auc
            2
            3  plt.plot(fpr2,tpr2, label= "train %0.2f"% auc(fpr1,tpr1))
            4  plt.plot(fpr1,tpr1, label= "test %0.2f"% auc(fpr2,tpr2))
            5
            6  plt.xlabel("tpr rate")
            7  plt.ylabel("fpr rate")
            8  plt.legend()
            9  plt.title("ROC Curve on pridicted y data by model performed on set1 with feature selection ")
           10
           11  plt.grid()
           12  plt.show()
```

ROC Curve on pridicted y data by model performed on set1 with feature selection

**confusion matrix**

```
In [ ]:    1  tr_cm = confusion_matrix(y,tr_y_pred)
           2  te_cm = confusion_matrix(y_cv,te_y_pred)
           3
           4  tr_cm,te_cm
```

## Tabulate your results

```
In [ ]:    1  # !pip install prettytable
```

```
In [44]:    1  from prettytable import PrettyTable
            2
            3  # Specify the Column Names while initializing the Table
            4  myTable = PrettyTable(["Vectorizer", "Model", "Hyper Parameter","Train AUC"," Test AUC" ])
            5
            6  # Add rows
            7  myTable.add_rows(
            8      [
            9
           10          ["TFIDF", "Decision Tree", "{max_depth :10 ,  min_samples_split :500}"," 0.6571","0.5796" ],
           11
           12          ["TFIDF W2V", "Decision Tree", " {'max_depth': 3, 'min_samples_split': 5}"," 0.5926","0.5689" ],
           13
           14          ["TFIDF -imp feature", "Decision Tree", " {'max_depth': 10, 'min_samples_split': 500}"," 0.6678","0.5762 \n"
           15
           16          ["TFIDF + rc", "LGBMClassifier", "{'max_depth': 30, 'min_child_samples': 80, 'reg_lambda': 10000}","0.6831",
           17          ["TFIDF-W2V + rc", "LGBMClassifier", "{'max_depth': 10, 'min_child_samples': 60, 'reg_lambda': 1000}"," 0.72
           18          ["TFIDF-imp feature", "LGBMClassifier", "{'max_depth': 30, 'min_child_samples': 40, 'reg_lambda': 10000}","0
           19
           20
           21
           22      ]
           23  )
           24
           25  print(myTable)
           26
           27
           28
```

```
+--------------------+---------------+-----------------------------------------------------------+-----------+
-----------+
|      Vectorizer    |     Model     |                     Hyper Parameter                       | Train AUC |
Test AUC |
+--------------------+---------------+-----------------------------------------------------------+-----------+
-----------+
|        TFIDF       | Decision Tree |           {max_depth :10 ,  min_samples_split :500}       |   0.6571  |
0.5796  |
|      TFIDF W2V      | Decision Tree |            {'max_depth': 3, 'min_samples_split': 5}        |   0.5926  |
0.5689  |
| TFIDF -imp feature | Decision Tree |           {'max_depth': 10, 'min_samples_split': 500}      |   0.6678  |
```

```
0.5762   |
|                   |               |                                                          |          |
|
|     TFIDF + rc    | LGBMClassifier |  {'max_depth': 30, 'min_child_samples': 80, 'reg_lambda': 1000} |   0.6831 |
0.6261  |
|   TFIDF-W2V + rc   | LGBMClassifier |  {'max_depth': 10, 'min_child_samples': 60, 'reg_lambda': 1000} |   0.7289 |
0.6147  |
| TFIDF-imp feature | LGBMClassifier | {'max_depth': 30, 'min_child_samples': 40, 'reg_lambda': 10000} |   0.6876 |
0.6399  |
+-------------------+---------------+----------------------------------------------------------+----------+
----------+
```

**note :**