

Decision Trees on Donors Choose dataset

Importing Libraries

```
In [1]: 1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 import pandas as pd
6 import numpy as np
7 import math as m
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import nltk
11 import re
12 from scipy import sparse
13 from sklearn.feature_extraction.text import TfidfVectorizer
14 from sklearn.feature_extraction.text import CountVectorizer
15 from sklearn.preprocessing import Normalizer
16 from sklearn.metrics import confusion_matrix
17 from sklearn.model_selection import train_test_split
18 from sklearn import metrics
19 from sklearn.metrics import roc_curve, auc
20 from tqdm import tqdm
21
22 import nltk
23 nltk.download('vader_lexicon')
24 from nltk.sentiment.vader import SentimentIntensityAnalyzer
25 sid = SentimentIntensityAnalyzer()
```

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

```
In [2]: 1 import tensorflow as tf
2 tf.test.gpu_device_name()
```

Out[2]: '/device:GPU:0'

Importing Dataset

```
In [9]: 1 !wget --header="Host: doc-14-3k-docs.googleusercontent.com" --header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:68.0) Gecko/20100101 Firefox/68.0" https://doc-14-3k-docs.googleusercontent.com/docs/securesc/8h8uvuh5ifib89027b7ihknt22nu5vgc/t7va9bgvmpo50k1actnt4qtpenckrn6g/1587475725000/00484516897554883881/05866892802988797180/1GU3LIJJ3zS1xLXXe-sdItSJHtI5txjV0?e=download&authuser=0&nonce=n1kem9jsfroba&user=05866892802988797180&hash=52na0kdlov88v5r8oe0k3796nblteim (https://doc-14-3k-docs.googleusercontent.com/docs/securesc/8h8uvuh5ifib89027b7ihknt22nu5vgc/t7va9bgvmpo50k1actnt4qtpenckrn6g/1587475725000/00484516897554883881/05866892802988797180/1GU3LIJJ3zS1xLXXe-sdItSJHtI5txjV0?e=download&authuser=0&nonce=n1kem9jsfroba&user=05866892802988797180&hash=52na0kdlov88v5r8oe0k3796nblteim)
Resolving doc-14-3k-docs.googleusercontent.com (doc-14-3k-docs.googleusercontent.com)... 64.233.167.132, 2a00:1450:400c:c0a::84
Connecting to doc-14-3k-docs.googleusercontent.com (doc-14-3k-docs.googleusercontent.com)|64.233.167.132|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: unspecified [text/csv]
Saving to: 'preprocessed_data.csv'

preprocessed_data.c      [          <=>          ] 118.69M  52.3MB/s   in 2.3s

2020-04-21 13:29:20 (52.3 MB/s) - 'preprocessed_data.csv' saved [124454659]
```

```
In [0]: 1 data = pd.read_csv('preprocessed_data.csv')
```

```
In [11]: 1 data.head(2)
```

Out[11]:

	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	al he2
1	ut	ms	grades_3_5	4	1	specialneeds	

```
In [0]: 1 negative = []
2 positive = []
3 neutral = []
4 compound = []
5
6 def update_sentiments(values):
7     negative.append(values["neg"])
8     positive.append(values["pos"])
9     neutral.append(values["neu"])
10    compound.append(values["compound"])
```

```
In [13]: 1 from tqdm import tqdm
2 for essay in tqdm(data["essay"]):
3     update_sentiments(sid.polarity_scores(essay))

100%|██████████| 109248/109248 [03:27<00:00, 526.52it/s]
```

```
In [0]: 1 data["neg"] = negative
2 data["pos"] = positive
3 data["neu"] = neutral      # adding new features to dataset based on Sentiment Intensity Analyzer
4 data["compound"] = compound
```

```
In [15]: 1 data.head(1)
```

```
Out[15]: school_state  teacher_prefix  project_grade_category  teacher_number_of_previously_posted_projects  project_is_approved  clean_categories  clean_s
```

	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	al hea

Splitting Data Into Train And Cross Validation(or test): Stratified Sampling

```
In [0]: 1 y = data['project_is_approved'].values
2 X = data.drop(['project_is_approved'], axis=1)
3
```

```
In [0]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y, random_state=42)
```

```
In [18]: 1 print("Total data points in Train Dataset =", len(y_train))
2 print("Total data points in Test Dataset =", len(y_test))
```

Total data points in Train Dataset = 73196
Total data points in Test Dataset = 36052

Make Data Model Ready: Encoding Eassay(text feature)

TFIDF Vectorizer

```
In [0]: 1 tfidfvectorizer = TfidfVectorizer(min_df=10, max_features=5000)
2 text_tfidf = tfidfvectorizer.fit(X_train['essay'].values) #fitting
3
4 X_train_essay_tfidf = tfidfvectorizer.transform(X_train['essay'].values)
5 X_test_essay_tfidf = tfidfvectorizer.transform(X_test['essay'].values) # transform
6
7 print("After vectorizations")
8 print(X_train_essay_tfidf.shape, y_train.shape)
9 print(X_test_essay_tfidf.shape, y_test.shape)
10 print("="*40)
```

After vectorizations
(73196, 5000) (73196,)
(36052, 5000) (36052,)
=====

TFIDF W2V

In [20]:

```
1 !wget --header="Host: doc-10-3k-docs.googleusercontent.com" --header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/80.0.3987.163 Safari/537.36" https://doc-10-3k-docs.googleusercontent.com/docs/securesc/8h8uvuh5ifib89027b7ihknt22nu5vgc/t30frjril2s6r6rt84jvdfcc5im19uck/1587476025000/00484516897554883881/05866892802988797180/1zJbDcbwvM2ueudqJHPp0b3Z9V2QrGd2r?e=download&authuser=0&nonce=uaih9emd25nk6&user=05866892802988797180&hash=939k2qd0inucqnck97ei14sqt2c7if3 (https://doc-10-3k-docs.googleusercontent.com/docs/securesc/8h8uvuh5ifib89027b7ihknt22nu5vgc/t30frjril2s6r6rt84jvdfcc5im19uck/1587476025000/00484516897554883881/05866892802988797180/1zJbDcbwvM2ueudqJHPp0b3Z9V2QrGd2r?e=download&authuser=0&nonce=uaih9emd25nk6&user=05866892802988797180&hash=939k2qd0inucqnck97ei14sqt2c7if3)
Resolving doc-10-3k-docs.googleusercontent.com (doc-10-3k-docs.googleusercontent.com)... 64.233.167.132, 2a00:1450:400c:c0a::84
Connecting to doc-10-3k-docs.googleusercontent.com (doc-10-3k-docs.googleusercontent.com)|64.233.167.132|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: unspecified [application/octet-stream]
Saving to: 'glove_vectors'

glove_vectors          [      <=>          ] 121.60M  46.6MB/s   in 2.6s

2020-04-21 13:34:29 (46.6 MB/s) - 'glove_vectors' saved [127506004]
```

In [0]:

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

In [0]:

```
1 tfidf_model = TfidfVectorizer()
2 tfidf_model.fit(X_train['essay'].values)
3 dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
4 tfidf_words = set(tfidf_model.get_feature_names())
```

In [23]:

```
1 #TFIDF W2V for train dataset
2 train_tfidf_w2v_essays = [] # the tfidf-w2v for each essay is stored in this list
3 for sentence in tqdm(X_train['essay']):
4     vector = np.zeros(300)
5     tf_idf_weight = 0;
6     for word in sentence.split():
7         if (word in glove_words) and (word in tfidf_words):
8             vec = model[word]
9             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
10            vector += (vec * tf_idf)
11            tf_idf_weight += tf_idf
12        if tf_idf_weight != 0:
13            vector /= tf_idf_weight
14        train_tfidf_w2v_essays.append(vector)
15 X_train_essay_tfidf_w2v= sparse.csr_matrix(train_tfidf_w2v_essays)
16 print("After vectorizations")
17 print(X_train_essay_tfidf_w2v.shape, y_train.shape)
18 print("="*40)
```

100%|██████████| 73196/73196 [02:31<00:00, 483.50it/s]

After vectorizations
(73196, 300) (73196,)
=====

In [24]:

```
1 #TFIDF W2V for test dataset
2 test_tfidf_w2v_essays = [] # the tfidf-w2v for each essay is stored in this list
3 for sentence in tqdm(X_test['essay']):
4     vector = np.zeros(300)
5     tf_idf_weight = 0;
6     for word in sentence.split():
7         if (word in glove_words) and (word in tfidf_words):
8             vec = model[word]
9             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split()))
10            vector += (vec * tf_idf)
11            tf_idf_weight += tf_idf
12        if tf_idf_weight != 0:
13            vector /= tf_idf_weight
14        test_tfidf_w2v_essays.append(vector)
15 X_test_essay_tfidf_w2v= sparse.csr_matrix(test_tfidf_w2v_essays)
16 print("After vectorizations")
17 print(X_test_essay_tfidf_w2v.shape, y_test.shape)
18 print("="*40)
```

100%|██████████| 36052/36052 [01:15<00:00, 479.44it/s]

After vectorizations
(36052, 300) (36052,)
=====

Make Data Model Ready: Categorical Features

Response Encoding

```
In [0]: 1 def fit(feature):
2         X_train['class_label']=y_train # adding 'project_is_approved' column to x_train
3         count = X_train[ feature ].value_counts() # getting value counts(denominator) of each category
4         feature_dictionary = dict()
5         for i, denominator in count.items():
6             vector = []
7             for j in range(2):
8                 compare =X_train.loc[ ( X_train['class_label'] == j ) & (X_train[feature] == i ) ]
9                 vector.append( len( compare) / denominator )
10            feature_dictionary[i] = vector # adding probability of each class label for a pariticular category of feat
11        return feature_dictionary
12    def transform(feature, df ):
13        feature_dictionary = fit(feature)
14        count = X_train[feature].value_counts()
15        f=[]
16        for cat in df[feature]:
17            if cat in dict( count ).keys():# transform test data with trainning probabilities
18                f.append( feature_dictionary[cat] )
19            else:
20                f.append([0.5, 0.05])
21        return f
```

Encoding Categorical Features: School State

```
In [0]: 1 X_train_state_rc =np.array(transform('school_state',X_train))
2         X_test_state_rc =np.array(transform('school_state',X_test))
3         print("After vectorizations")
4         print(X_train_state_rc.shape, y_train.shape)
5         print(X_test_state_rc.shape, y_test.shape)
6         print("="*40)
```

After vectorizations
(73196, 2) (73196,)
(36052, 2) (36052,)
=====

Encoding Categorical Features: teacher_prefix

```
In [0]: 1 X_train_teacher_rc =np.array(transform('teacher_prefix',X_train))
2         X_test_teacher_rc = np.array(transform('teacher_prefix',X_test))
3         print("After vectorizations")
4         print(X_train_teacher_rc.shape, y_train.shape)
5         print(X_test_teacher_rc.shape, y_test.shape)
6         print("="*40)
```

After vectorizations
(73196, 2) (73196,)
(36052, 2) (36052,)
=====

Encoding Categorical Features: project_grade_category

```
In [0]: 1 X_train_grade_rc = np.array(transform('project_grade_category',X_train))
2         X_test_grade_rc = np.array(transform('project_grade_category',X_test))
3
4         print("After vectorizations")
5         print(X_train_grade_rc.shape, y_train.shape)
6         print(X_test_grade_rc.shape, y_test.shape)
7         print("="*40)
```

After vectorizations
(73196, 2) (73196,)
(36052, 2) (36052,)
=====

Encoding Categorical Features: clean_categories

```
In [0]: 1 X_train_category_rc =np.array(transform('clean_categories',X_train))
2         X_test_category_rc = np.array(transform('clean_categories',X_test))
3
4         print("After vectorizations")
5         print(X_train_category_rc.shape, y_train.shape)
6         print(X_test_category_rc.shape, y_test.shape)
7         print("="*40)
```

After vectorizations
(73196, 2) (73196,)
(36052, 2) (36052,)
=====

Encoding Categorical Features: clean_subcategories

```
In [0]: 1 X_train_subcategory_rc = np.array(transform('clean_subcategories',X_train))
2 X_test_subcategory_rc = np.array(transform('clean_subcategories',X_test))
3 print("After vectorizations")
4 print(X_train_subcategory_rc.shape, y_train.shape)
5 print(X_test_subcategory_rc.shape, y_test.shape)
6 print("="*40)
```

```
After vectorizations
(73196, 2) (73196,)
(36052, 2) (36052,)
=====
```

Encoding Numerical Features**Encoding Numerical Feature :price**

```
In [0]: 1 normalizer = Normalizer()
2 normalizer.fit(X_train['price'].values.reshape(1,-1)) #fitting
3
4 X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(1,-1)) #transform
5 X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(1,-1))
6
7 X_train_price_norm =X_train_price_norm.reshape(-1,1)
8 X_test_price_norm=X_test_price_norm.reshape(-1,1)
9
10 print("After vectorizations")
11 print(X_train_price_norm.shape, y_train.shape)
12 print(X_test_price_norm.shape, y_test.shape)
13 print("="*40)
```

```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

Encoding Numerical Features:teacher_number_of_previously_posted_projects

```
In [0]: 1 normalizer = Normalizer()
2 normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1)) #fitting
3 X_train_submission_norm = normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
4 X_test_submission_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
5
6 X_train_submission_norm =X_train_submission_norm.reshape(-1,1)
7 X_test_submission_norm=X_test_submission_norm.reshape(-1,1)
8
9 print("After vectorizations")
10 print(X_train_submission_norm.shape, y_train.shape)
11 print(X_test_submission_norm.shape, y_test.shape)
12 print("="*40)
```

```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

Encoding Numerical Features:neg

```
In [0]: 1 normalizer = Normalizer()
2 normalizer.fit(X_train['neg'].values.reshape(1,-1)) #fitting
3
4 X_train_neg_norm = normalizer.transform(X_train['neg'].values.reshape(1,-1)) #transform
5 X_test_neg_norm = normalizer.transform(X_test['neg'].values.reshape(1,-1))
6
7 X_train_neg_norm =X_train_neg_norm.reshape(-1,1)
8 X_test_neg_norm=X_test_neg_norm.reshape(-1,1)
9
10 print("After vectorizations")
11 print(X_train_neg_norm.shape, y_train.shape)
12 print(X_test_neg_norm.shape, y_test.shape)
13 print("="*40)
```

```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

Encoding Numerical Features:pos


```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

```
After vectorizations
(73196, 1) (73196,)
(36052, 1) (36052,)
=====
```

```
SHAPE OF TRAIN AND TEST AFTER STACKING
(73196, 5016)
(36052, 5016)
```

SET-2

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

```
In [0]: 1 from scipy.sparse import hstack
2 X_tr_set_two = hstack((X_train_essay_tfidf_w2v, X_train_state_rc, X_train_teacher_rc, X_train_grade_rc, X_train_price_rc))
3 X_te_set_two = hstack((X_test_essay_tfidf_w2v, X_test_state_rc, X_test_teacher_rc, X_test_grade_rc, X_test_price_norm_rc))
```

```
In [0]: 1 print("SHAPE OF TRAIN AND TEST AFTER STACKING")
2 print(X_tr_set_two.shape)
3 print(X_te_set_two.shape)
```

```
SHAPE OF TRAIN AND TEST AFTER STACKING
(73196, 312)
(36052, 312)
```

GRADIENT BOOSTING CLASSIFIER USING GRID SEARCH CROSS VALIDATION (SET - 1)

```
In [0]: 1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.model_selection import GridSearchCV
3 parameters = {"max_depth": [1, 2, 3, 4], "n_estimators": [5, 10, 15, 20]}
4 clf = GridSearchCV(GradientBoostingClassifier(), parameters, cv=5, scoring='roc_auc', return_train_score=True, n_jobs=-1)
5 clf.fit(X_tr_set_one, y_train)
```

```
Out[76]: GridSearchCV(cv=5, error_score=nan,
                    estimator=GradientBoostingClassifier(ccp_alpha=0.0,
                                                         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,
                                                         n_iter_no_change=None,
                                                         presort='deprecated',
                                                         random_state=None,
                                                         subsample=1.0, tol=0.0001,
                                                         validation_fraction=0.1,
                                                         verbose=0, warm_start=False),
                    iid='deprecated', n_jobs=-1,
                    param_grid={'max_depth': [1, 2, 3, 4],
                                'n_estimators': [5, 10, 15, 20]},
                    pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                    scoring='roc_auc', verbose=0)
```

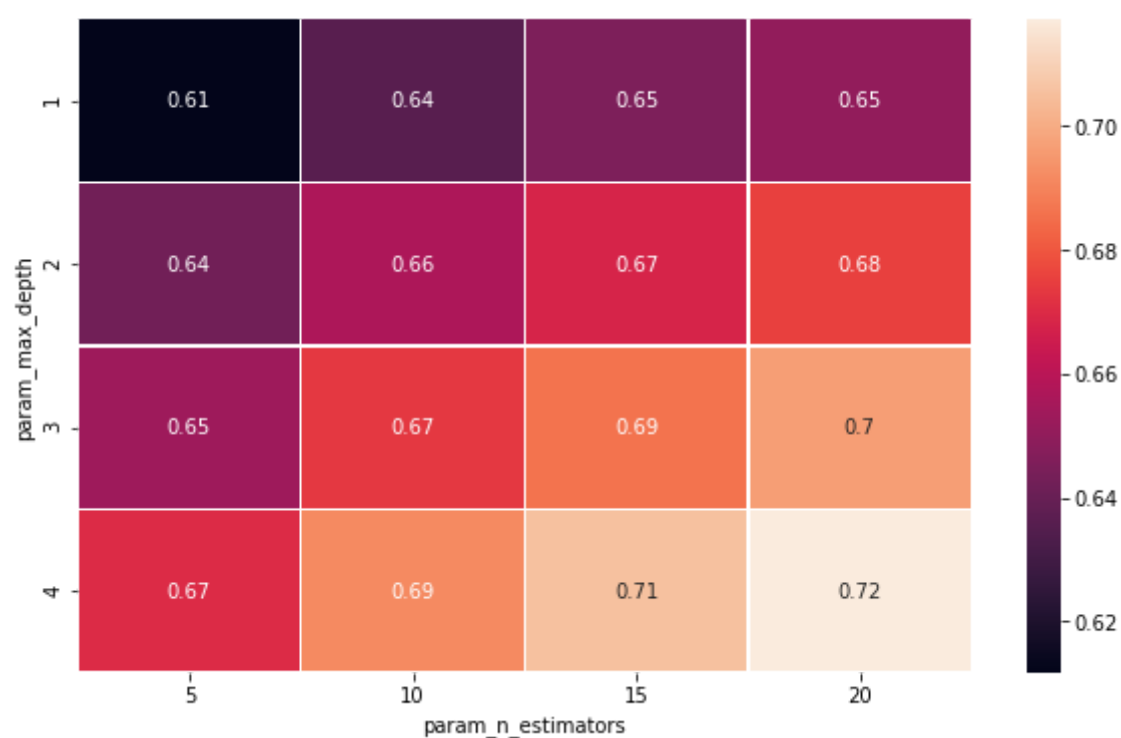
```
In [0]: 1 train_auc= clf.cv_results_['mean_train_score']
2 train_auc_std= clf.cv_results_['std_train_score']
3 cv_auc = clf.cv_results_['mean_test_score']
4 cv_auc_std= clf.cv_results_['std_test_score']
```

```
In [0]: 1 print('Best score: ',clf.best_score_)
2 print('Best Hyper parameters: ',clf.best_params_)
```

```
Best score: 0.684567328741456
Best Hyper parameters: {'max_depth': 4, 'n_estimators': 20}
```

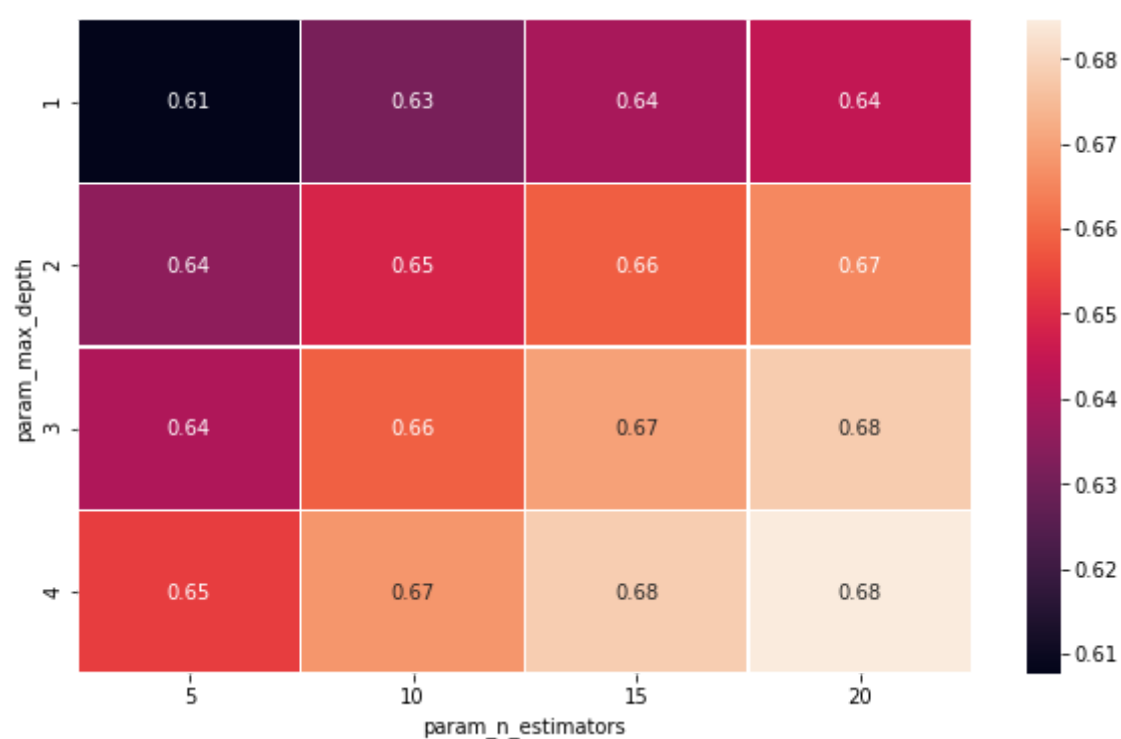
Plotting Hyperparameter v/s Auc**Roc Plot Of Train And Test Data****Train data**

```
In [0]: 1 import pandas as pd
2 pvt = pd.pivot_table(pd.DataFrame(clf.cv_results_),
3     values='mean_train_score', index='param_max_depth', columns='param_n_estimators') #https://stackoverflow.com/ques
4 plt.figure(figsize=(10,6))
5 ax=sns.heatmap(pvt,annot=True,linewidths=.5)
```



Test data

```
In [0]: 1 import pandas as pd
2 pvt = pd.pivot_table(pd.DataFrame(clf.cv_results_),
3     values='mean_test_score', index='param_max_depth', columns='param_n_estimators') #https://stackoverflow.com/quest
4 plt.figure(figsize=(10,6))
5 ax=sns.heatmap(pvt,annot=True,linewidths=.5)
```

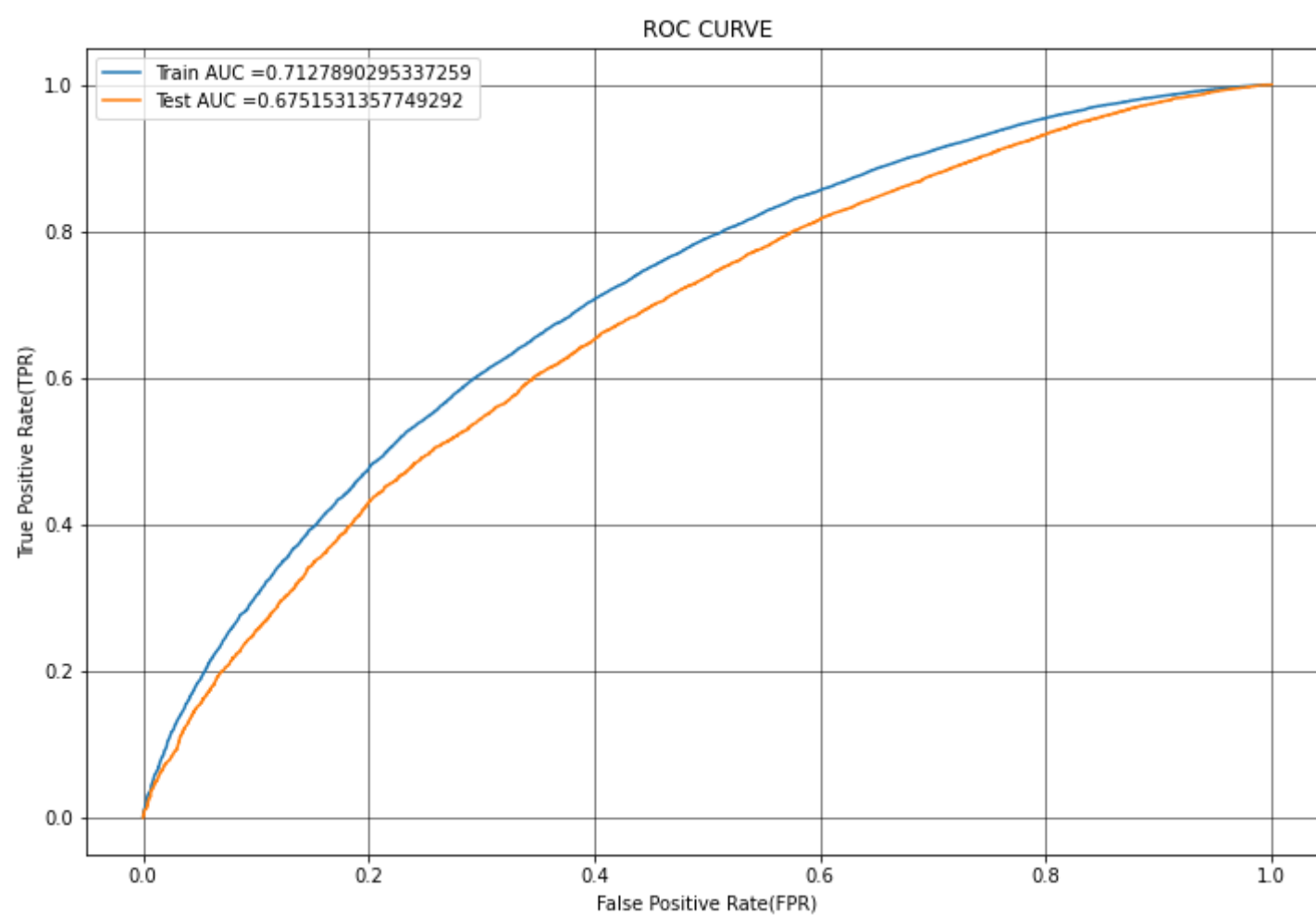


Roc Plot Of Train And Test Data


```

In [0]: 1 model_set1=GradientBoostingClassifier(max_depth = clf.best_params_["max_depth"], n_estimators= clf.best_params_["n_es
2 model_set1.fit(X_tr_set_one,y_train)
3 y_train_probs = clf.predict_proba(X_tr_set_one)[: ,1] # converting train and test output into probability
4 y_test_probs= clf.predict_proba(X_te_set_one )[: ,1]
5
6 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_probs) # storing values of fpr and tpr
7 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_probs)
8
9 plt.figure(figsize=(12,8))
10 plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("False Positive Rate(FPR)")
14 plt.ylabel("True Positive Rate(TPR)")
15 plt.title("ROC CURVE")
16 plt.grid(color='black',lw=0.5)

```



Confusion Matrix

```

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

```

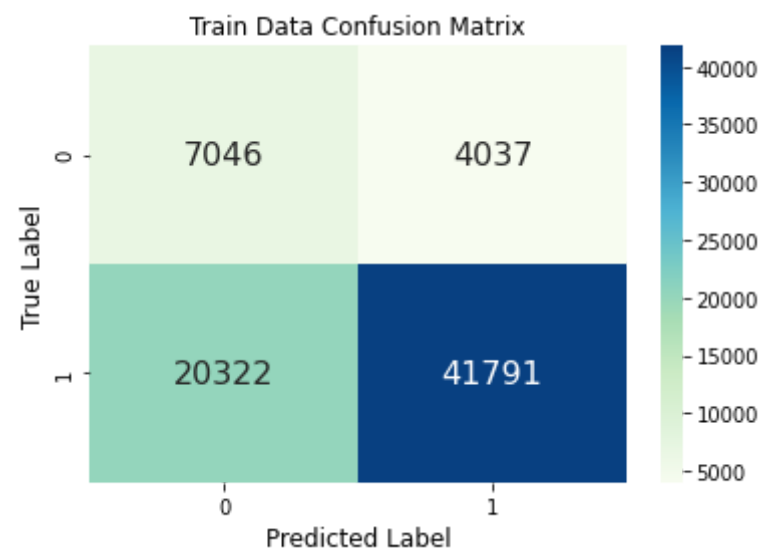
Train Data

```
In [0]: 1 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
2 cm=metrics.confusion_matrix(y_train,predict_with_best_t(y_train_probs, best_t)) # https://stackoverflow.com/question
3
4 print("CONFUSION MATRIX OF TRAIN DATA")
5 print("\n")
6 print(cm)
7 sns.heatmap(cm, annot=True, fmt='d',cmap='GnBu',annot_kws = {"size":16})
8 plt.ylabel('True Label',size=12)
9 plt.xlabel('Predicted Label',size=12)
10 plt.title('Train Data Confusion Matrix',size=12)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.4277456121087118 for threshold 0.846
CONFUSION MATRIX OF TRAIN DATA

```
[[ 7046  4037]
 [20322 41791]]
```

Out[112]: Text(0.5, 1.0, 'Train Data Confusion Matrix')



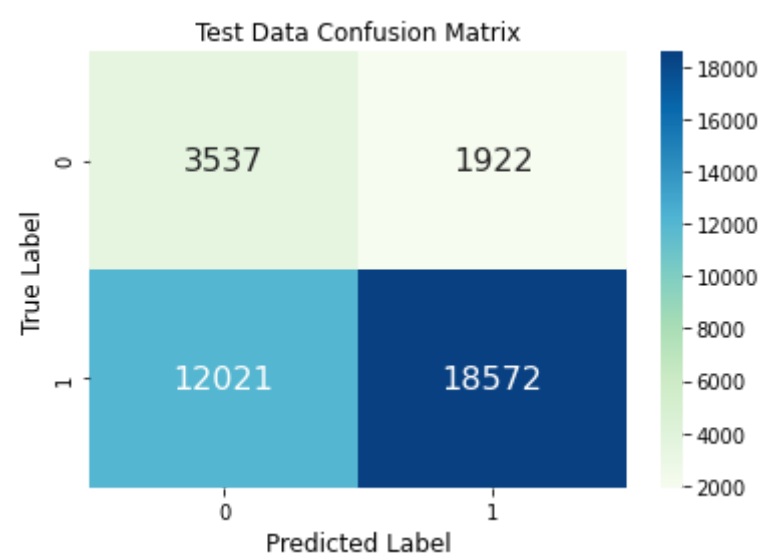
Test Data

```
In [0]: 1 best_t = find_best_threshold(te_thresholds, test_fpr, test_tpr)
2 cm=metrics.confusion_matrix(y_test,predict_with_best_t(y_test_probs, best_t))
3
4 print("CONFUSION MATRIX OF TEST DATA")
5 print('\n')
6 print(cm)
7 sns.heatmap(cm, annot=True, fmt='d',cmap='GnBu',annot_kws = {"size":16})
8 plt.ylabel('True Label',size=12)
9 plt.xlabel('Predicted Label',size=12)
10 plt.title('Test Data Confusion Matrix',size=12)
```

the maximum value of $tpr \cdot (1 - fpr)$ 0.39333136004500213 for threshold 0.85
CONFUSION MATRIX OF TEST DATA

```
[[ 3537  1922]
 [12021 18572]]
```

Out[113]: Text(0.5, 1.0, 'Test Data Confusion Matrix')



GRADIENT BOOSTING CLASSIFIER USING GRID SEARCH CROSS VALIDATION (SET - 2)

```
In [0]: 1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.model_selection import GridSearchCV
3 parameters = {"max_depth":[1,2,3,4],"n_estimators":[5,10,15,20] }
4 clf = GridSearchCV(GradientBoostingClassifier(), parameters, cv=5, scoring='roc_auc',return_train_score=True,n_jobs=
5 clf.fit(X_tr_set_two,y_train)
```

```
Out[114]: GridSearchCV(cv=5, error_score=nan,
                      estimator=GradientBoostingClassifier(ccp_alpha=0.0,
                                                            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,
                                                            n_iter_no_change=None,
                                                            presort='deprecated',
                                                            random_state=None,
                                                            subsample=1.0, tol=0.0001,
                                                            validation_fraction=0.1,
                                                            verbose=0, warm_start=False),
                      iid='deprecated', n_jobs=-1,
                      param_grid={'max_depth': [1, 2, 3, 4],
                                   'n_estimators': [5, 10, 15, 20]},
                      pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                      scoring='roc_auc', verbose=0)
```

```
In [0]: 1 train_auc= clf.cv_results_['mean_train_score']
2 train_auc_std= clf.cv_results_['std_train_score']
3 cv_auc = clf.cv_results_['mean_test_score']
4 cv_auc_std= clf.cv_results_['std_test_score']
```

```
In [0]: 1 print('Best score: ',clf.best_score_)
2 print('Best Hyper parameters: ',clf.best_params_)
```

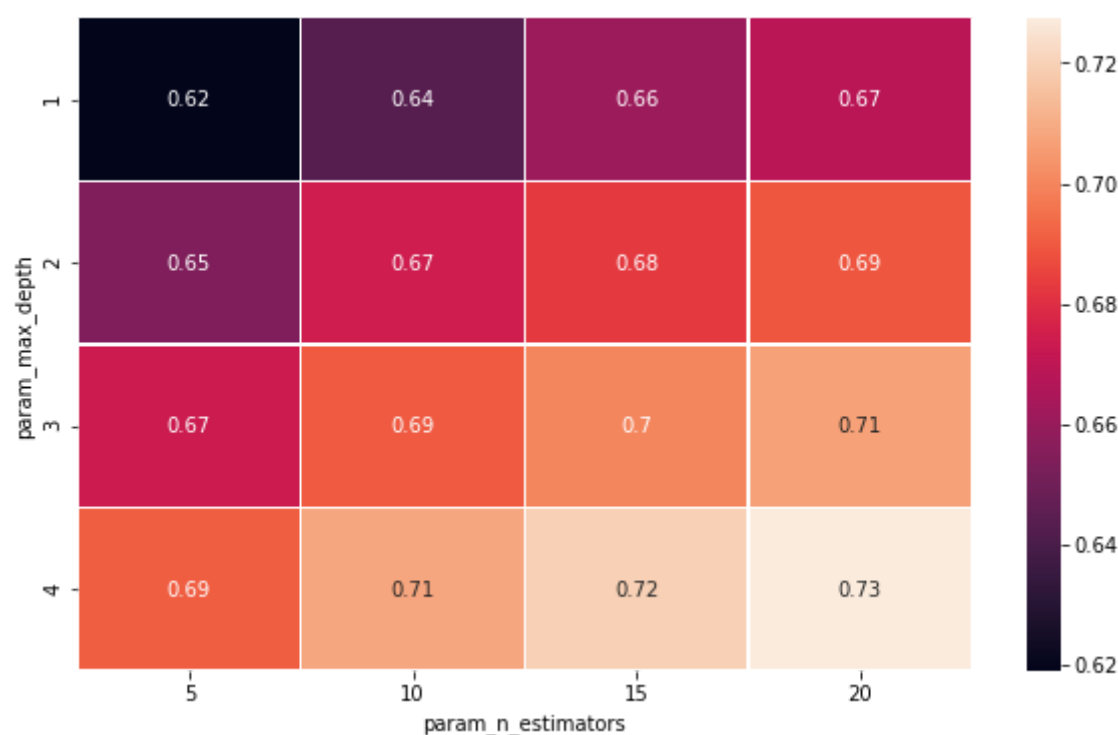
Best score: 0.6902930550494676

Best Hyper parameters: {'max_depth': 4, 'n_estimators': 20}

Plotting Hyperparameter v/s Auc

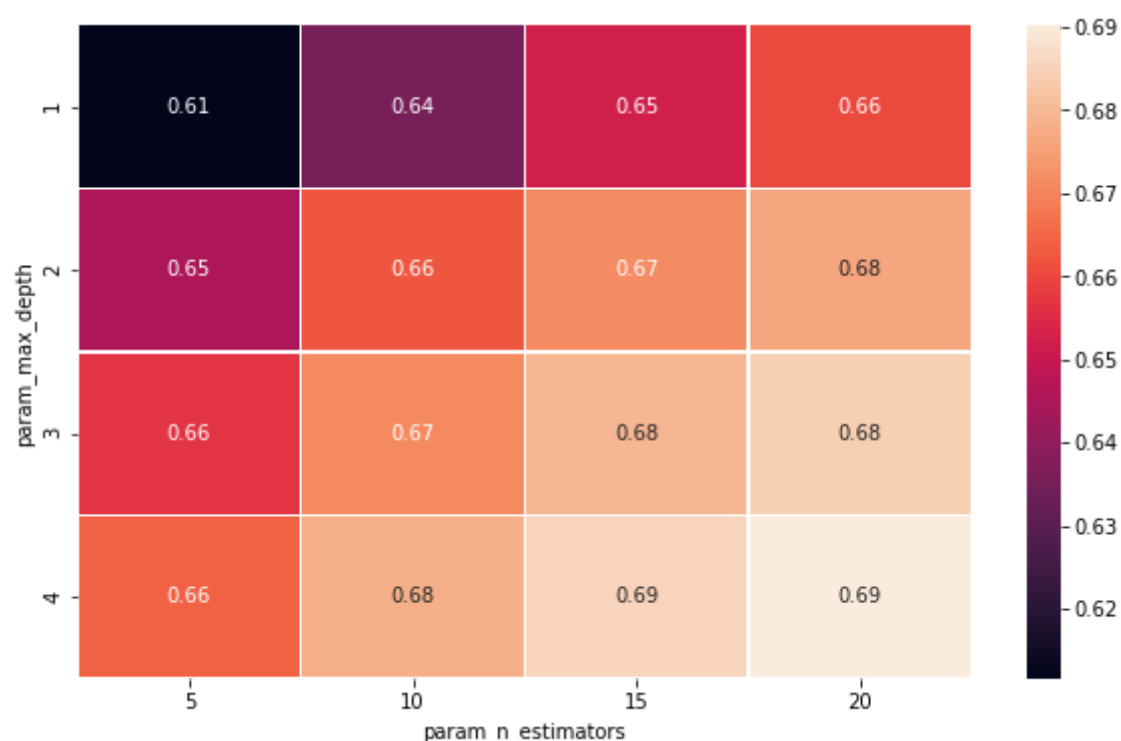
Train Data

```
In [0]: 1 import pandas as pd
2 pvt = pd.pivot_table(pd.DataFrame(clf.cv_results_),
3                       values='mean_train_score', index='param_max_depth', columns='param_n_estimators') #https://stackoverflow.com/ques
4 plt.figure(figsize=(10,6))
5 ax=sns.heatmap(pvt,annot=True,linewidths=.5)
```



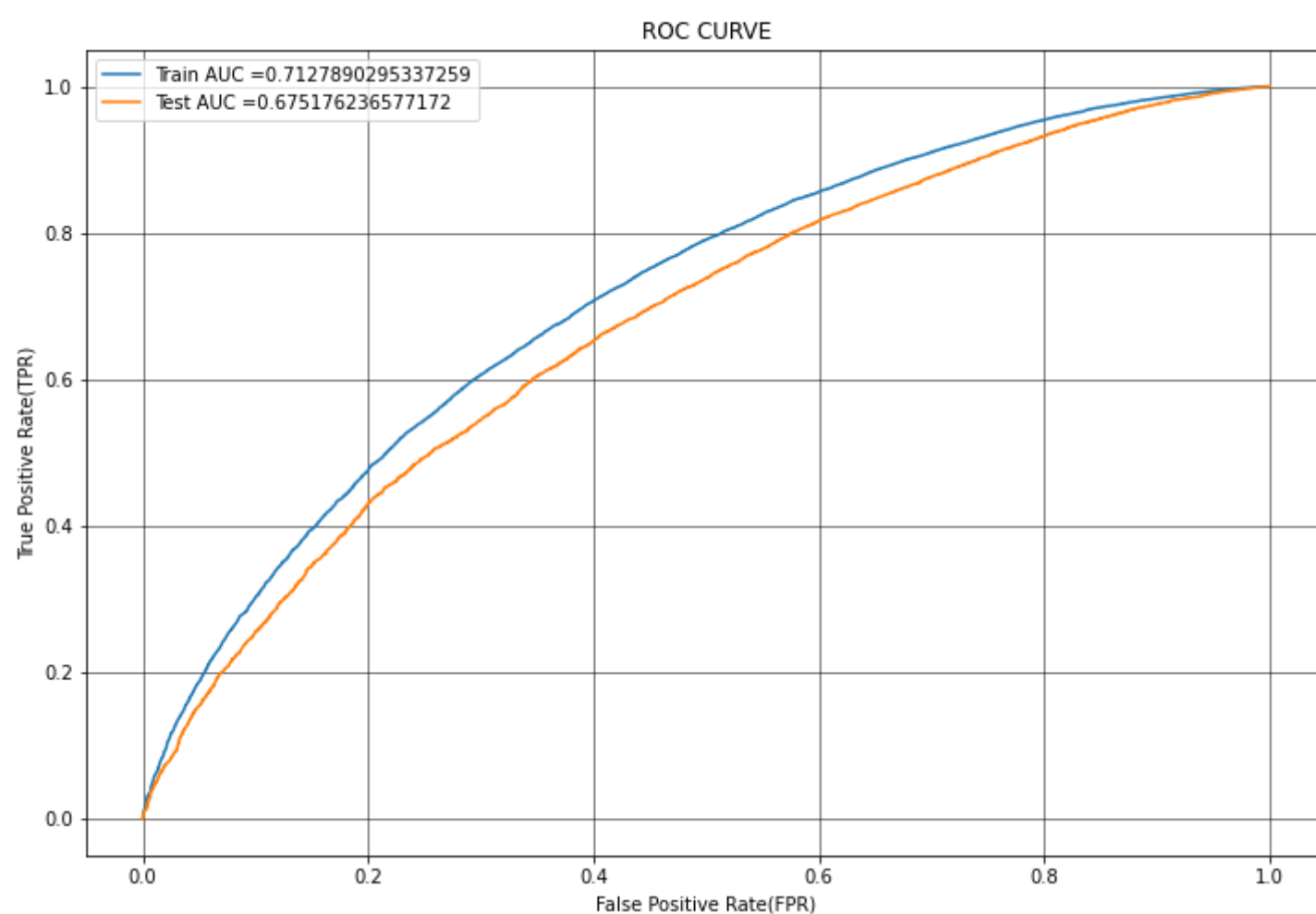
Cross Validation data

```
In [0]: 1 import pandas as pd
2 pvt = pd.pivot_table(pd.DataFrame(clf.cv_results_),
3     values='mean_test_score', index='param_max_depth', columns='param_n_estimators') #https://stackoverflow.com/quest
4 plt.figure(figsize=(10,6))
5 ax=sns.heatmap(pvt,annot=True,linewidths=.5)
```



Roc Plot Of Train And Test Data

```
In [0]: 1 model_set2=GradientBoostingClassifier(max_depth = clf.best_params_["max_depth"], n_estimators= clf.best_params_["n_es
2 model_set2.fit(X_tr_set_one,y_train)
3 y_train_probs = model_set2.predict_proba(X_tr_set_one)[: ,1] # converting train and test output into probability
4 y_test_probs= model_set2.predict_proba(X_te_set_one )[: ,1]
5
6 train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_probs) # storing values of fpr and tpr
7 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_probs)
8
9 plt.figure(figsize=(12,8))
10 plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("False Positive Rate(FPR)")
14 plt.ylabel("True Positive Rate(TPR)")
15 plt.title("ROC CURVE")
16 plt.grid(color='black',lw=0.5)
```



Confusion Matrix

```

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

```

Train Data

```

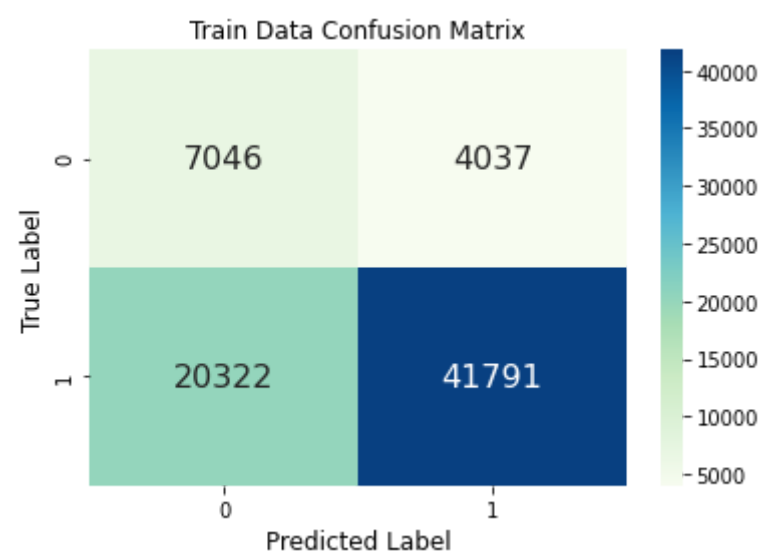
In [0]: 1 best_t = find_best_threshold(tr_thresholds, train_fpr, train_tpr)
2       cm=metrics.confusion_matrix(y_train,predict_with_best_t(y_train_probs, best_t)) # https://stackoverflow.com/question
3
4       print("CONFUSION MATRIX OF TRAIN DATA")
5       print("\n")
6       print(cm)
7       sns.heatmap(cm, annot=True, fmt='d',cmap='GnBu',annot_kws = {"size":16})
8       plt.ylabel('True Label',size=12)
9       plt.xlabel('Predicted Label',size=12)
10      plt.title('Train Data Confusion Matrix',size=12)

```

the maximum value of tpr*(1-fpr) 0.4277456121087118 for threshold 0.846
CONFUSION MATRIX OF TRAIN DATA

```
[[ 7046  4037]
 [20322 41791]]
```

Out[121]: Text(0.5, 1.0, 'Train Data Confusion Matrix')



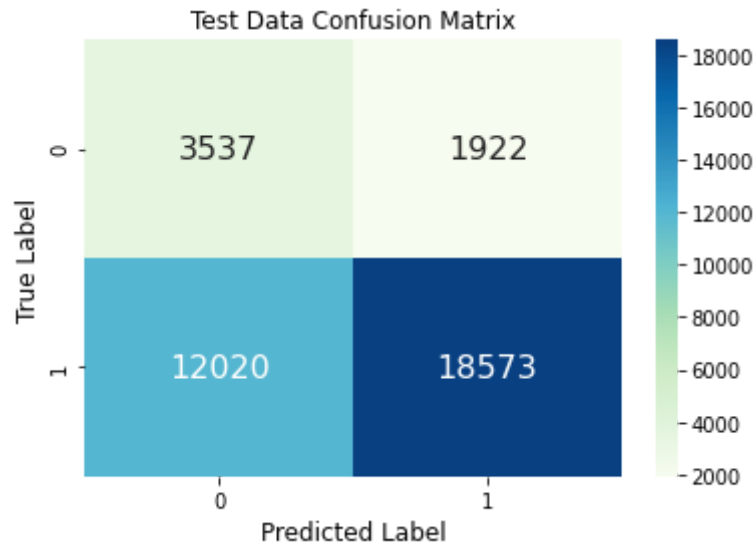
Test Data

```
In [0]: 1 best_t = find_best_threshold(te_thresholds, test_fpr, test_tpr)
2 cm=metrics.confusion_matrix(y_test,predict_with_best_t(y_test_probs, best_t))
3
4 print("CONFUSION MATRIX OF TEST DATA")
5 print('\n')
6 print(cm)
7 sns.heatmap(cm, annot=True, fmt='d',cmap='GnBu',annot_kws = {"size":16})
8 plt.ylabel('True Label',size=12)
9 plt.xlabel('Predicted Label',size=12)
10 plt.title('Test Data Confusion Matrix',size=12)
```

the maximum value of tpr*(1-fpr) 0.3933525387742744 for threshold 0.85
CONFUSION MATRIX OF TEST DATA

```
[[ 3537  1922]
 [12020 18573]]
```

Out[122]: Text(0.5, 1.0, 'Test Data Confusion Matrix')



Summary

```
In [26]: 1 from prettytable import PrettyTable
2 from prettytable import ALL as ALL
3 table=PrettyTable(hrules=ALL)
4 table.field_names = [ "Sl.NO", "Vectorizer", "Model", "Hyper Parameter", "Test-AUC"] # # http://zetcode.com/python/prettytable/
5 table.add_row([1,"TFIDF", "GRADIENT BOOSTING CLASSIFIER", "max_depth =4 , n_estimators=20", 0.67515])
6 table.add_row([2,"TFIDF W2V", "GRADIENT BOOSTING CLASSIFIER", " max_depth =4 , n_estimators=20", 0.67517])
7 print(table)
```

Sl.NO	Vectorizer	Model	Hyper Parameter	Test-AUC
1	TFIDF	GRADIENT BOOSTING CLASSIFIER	max_depth =4 , n_estimators=20	0.67515
2	TFIDF W2V	GRADIENT BOOSTING CLASSIFIER	max_depth =4 , n_estimators=20	0.67517

```
In [0]: 1
```