

1.5.1.4 Confusion Matrix for Train and Test data

```
In [44]: def find_best_threshold(proba, threshold, fpr, tpr):
t = threshold[np.argmax(tpr*(1-fpr))]
# (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
print("the maximum value of fpr*(1-tpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
return t

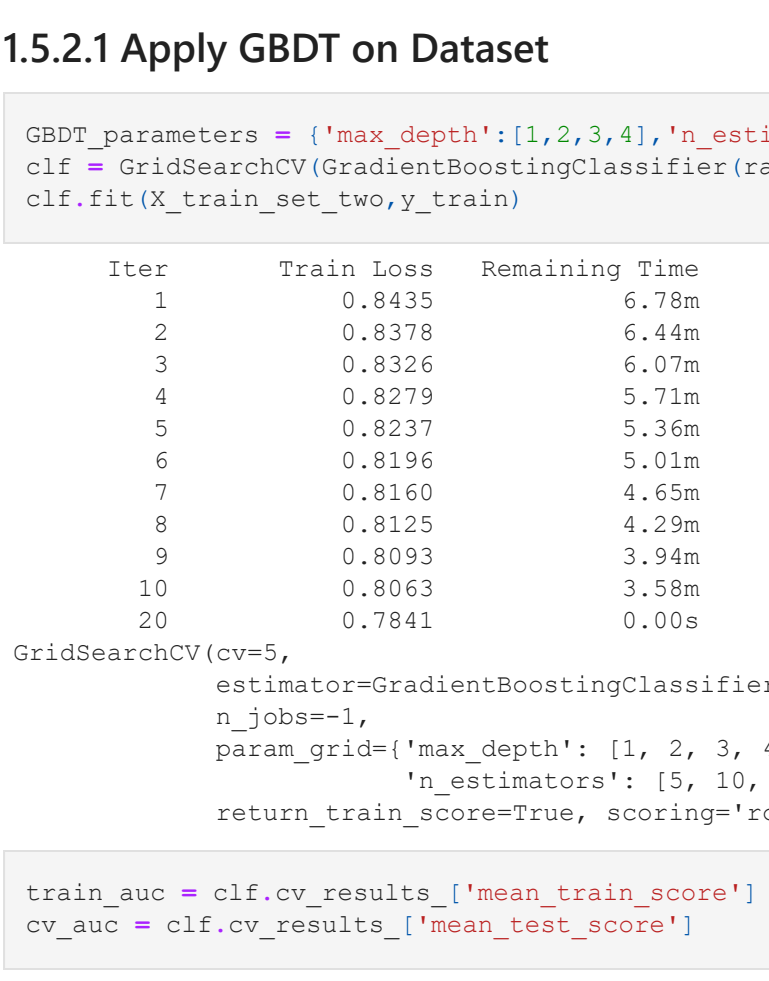
In [45]: def predict_with_best_threshold(proba, t):
predictions = []
for i in proba:
    if i>=t:
        predictions.append(1)
    else:
        predictions.append(0)
return predictions

In [46]: print("=="*100)
best_threshold = find_best_threshold(y_train_proba, tr_thresholds, train_fpr, train_tpr)
conf_matr_train_set_one = confusion_matrix(y_train, predict_with_best_threshold(y_train_proba, best_threshold))
print("Train confusion matrix")
print(conf_matr_train_set_one)

=====
Train confusion matrix
[[ 7039 4044]
 [20529 41584]]

In [47]: sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_train_set_one, annot=True,annot_kws={"size": 16}, fmt='g')
plt.title("Train data confusion matrix")
plt.xlabel("Predicted Class")
plt.ylabel("True class")

Out[47]: Text(26.5, 0.5, 'True class')
```

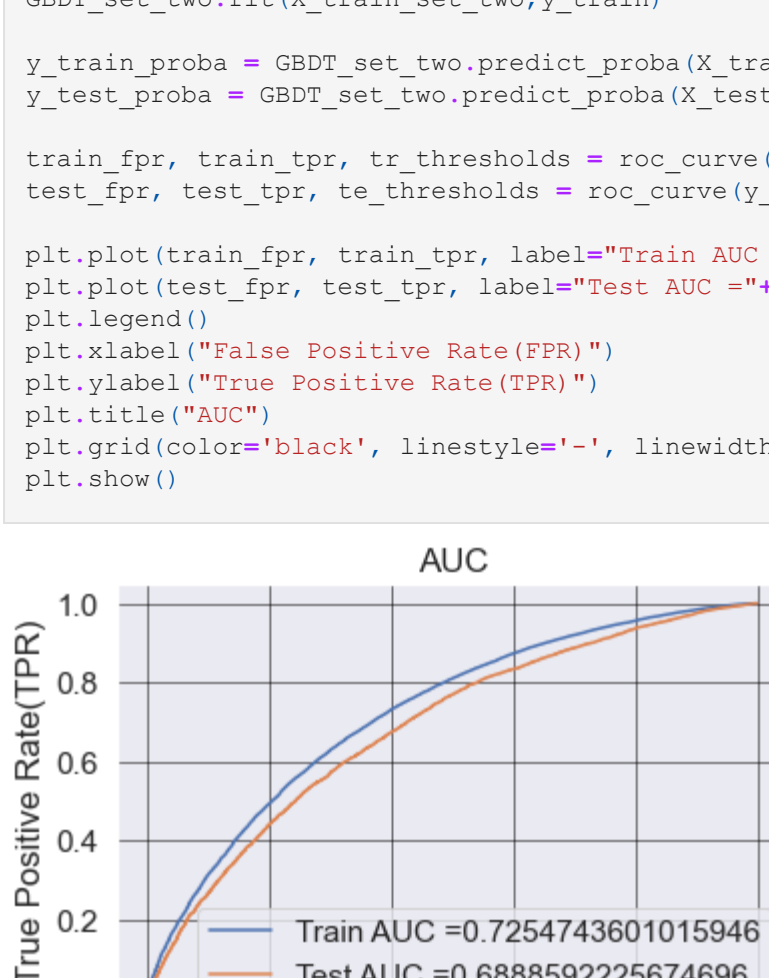


```
In [48]: print("=="*100)
print("Test confusion matrix")
conf_matr_test_set_one = confusion_matrix(y_test, predict_with_best_threshold(y_test_proba, best_threshold))
print(conf_matr_test_set_one)

=====
Test confusion matrix
[[ 3235 2224]
 [10308 20285]]

In [49]: sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_test_set_one, annot=True,annot_kws={"size": 16}, fmt='g')
plt.title("Test data confusion matrix")
plt.xlabel("Predicted Class")
plt.ylabel("True class")

Out[49]: Text(26.5, 0.5, 'True class')
```



1.5.2 Set 2:

1.5.2.1 Apply GBDT on Dataset

```
In [59]: GBDT_parameters = {'max_depth':[1,2,3,4], 'n_estimators':[5,10,15,20]}
clf = GridSearchCV(GradientBoostingClassifier(random_state=0, verbose = 1), GBDT_parameters, scoring = 'roc_auc',
clf.fit(X_train_set_two,y_train)

Iter      Train Loss      Remaining Time
1          0.8435          6.78m
2          0.8378          6.44m
3          0.8326          6.07m
4          0.8279          5.71m
5          0.8237          5.36m
6          0.8186          5.01m
7          0.8160          4.65m
8          0.8125          4.29m
9          0.8093          3.94m
10         0.8063          3.58m
20         0.7841          0.09s

Out[59]: GridSearchCV(cv=estimator=GradientBoostingClassifier(random_state=0, verbose=1),
n_jobs=-1,
param_grid={'max_depth': [1, 2, 3, 4],
param_name='n_estimators': [5, 10, 15, 20]},
return_train_score=True, scoring='roc_auc')

In [60]: train_auc = clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']

In [61]: print("Best score = ",clf.best_score_)
print("Best Hyper parameters = ",clf.best_params_)

Best score = 0.6927180923216303
Best Hyper parameters = {'max_depth': 4, 'n_estimators': 20}
```

1.5.1.2 Hyperparameter vs AUC Plot

```
In [62]: #Below code is obtained from: https://towardsdatascience.com/using-3d-visualizations-to-tune-hyperparameters-6
max_scores = pd.DataFrame(clf.cv_results_.groupby(['param_n_estimators', 'param_max_depth']).max())
max_scores = max_scores.unstack()[['mean_test_score', 'mean_train_score']]
plt.figure(figsize=(15, 5))
plt.subplot(1,2,1)
plt.title("Train Dataset")
sns.heatmap(max_scores.mean_train_score, annot=True, fmt='.4g')
plt.subplot(1,2,2)
plt.title("CV Dataset")
sns.heatmap(max_scores.mean_test_score, annot=True, fmt='.4g')

Out[62]: <AxesSubplot:title='center': 'CV Dataset', xlabel='param_max_depth', ylabel='param_n_estimators'>
```



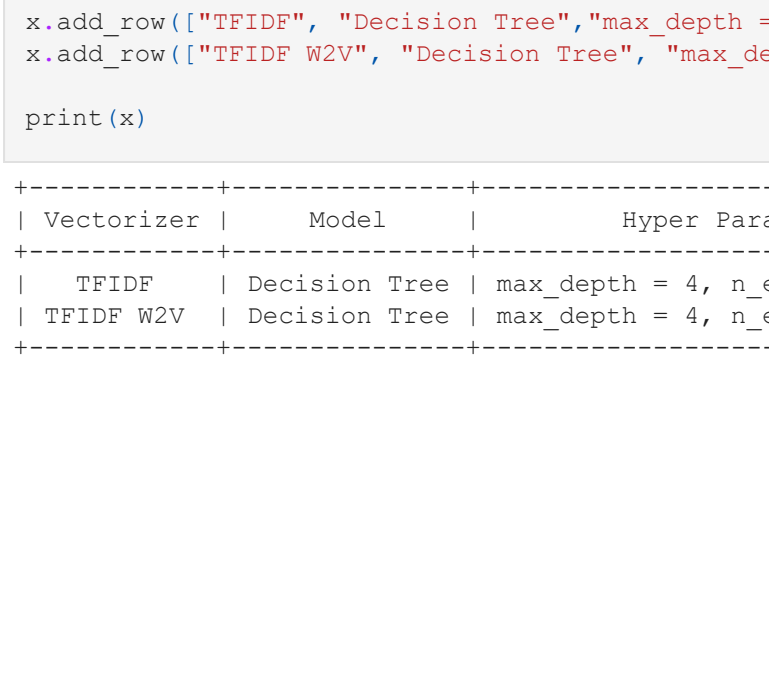
1.5.1.3 Plot ROC-AUC curve of model with best hyperparameters

```
In [79]: GBDT_set_two=GradientBoostingClassifier(max_depth = clf.best_params_["max_depth"],n_estimators = clf.best_params_
GBDT_set_two.fit(X_train_set_two,y_train)

y_train_proba = GBDT_set_two.predict_proba(X_train_set_two)[:,1]
y_test_proba = GBDT_set_two.predict_proba(X_test_set_two)[:,1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_proba)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_proba)

plt.plot(train_fpr, train_tpr, label="Train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="Test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(color='black', linestyle='-', linewidth=0.5)
plt.show()
```



1.5.1.4 Confusion Matrix for Train and Test data

```
In [80]: def find_best_threshold(proba, threshold, fpr, tpr):
t = threshold[np.argmax(tpr*(1-fpr))]
# (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high
print("the maximum value of fpr*(1-tpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
return t

In [81]: def predict_with_best_threshold(proba, t):
predictions = []
for i in proba:
    if i>=t:
        predictions.append(1)
    else:
        predictions.append(0)
return predictions

In [82]: print("=="*100)
best_threshold = find_best_threshold(y_train_proba, tr_thresholds, train_fpr, train_tpr)
conf_matr_train_set_two = confusion_matrix(y_train, predict_with_best_threshold(y_train_proba, best_threshold))
print("Train confusion matrix")
print(conf_matr_train_set_two)

=====
the maximum value of fpr*(1-tpr) 0.44411710039776253 for threshold 0.847
Train confusion matrix
[[ 7459 3624]
 [21125 40988]]

In [83]: sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_train_set_two, annot=True,annot_kws={"size": 16}, fmt='g')
plt.title("Train data confusion matrix")
plt.xlabel("Predicted Class")
plt.ylabel("True class")

Out[83]: Text(26.5, 0.5, 'True class')
```



```
In [84]: print("=="*100)
print("Test confusion matrix")
conf_matr_test_set_two = confusion_matrix(y_test, predict_with_best_threshold(y_test_proba, best_threshold))
print(conf_matr_test_set_two)

=====
Test confusion matrix
[[ 3411 2048]
 [10712 19881]]

In [85]: sns.set(font_scale=1.4)#for label size
sns.heatmap(conf_matr_test_set_two, annot=True,annot_kws={"size": 16}, fmt='g')
plt.title("Test data confusion matrix")
plt.xlabel("Predicted Class")
plt.ylabel("True class")

Out[85]: Text(26.5, 0.5, 'True class')
```



3. Summary

as mentioned in the step 4 of instructions

```
In [88]: # http://setcode.com/python/prettytable/
from PrettyTable import PrettyTable

x = PrettyTable()
x.field_names = ["Vectorizer", "Model", "Hyper Parameter", "AUC"]

x.add_row(["TFIDF", "Decision Tree", "max_depth = 4, n_estimators = 20", "0.6778"])
x.add_row(["TFIDF W2V", "Decision Tree", "max_depth = 4, n_estimators = 20", "0.6888"])

print(x)

+-----+
| Vectorizer | Model | Hyper Parameter | AUC |
+-----+-----+
| TFIDF | Decision Tree | max_depth = 4, n_estimators = 20 | 0.6778 |
| TFIDF W2V | Decision Tree | max_depth = 4, n_estimators = 20 | 0.6888 |
+-----+-----+

=====
the maximum value of fpr*(1-tpr) 0.44411710039776253 for threshold 0.847
Train confusion matrix
[[ 7459 3624]
 [21125 40988]]

In [89]:
```