

# Assignment 4

## Negin Baghibanzadeh

This assignment is based on content discussed in module 8 and using Decision Trees and Ensemble Models in classification and regression problems.

## Learning outcomes

- Understand how to use decision trees on a Dataset to make a prediction
- Learning hyper-parameters tuning for decision trees by using RandomGrid
- Learning the effectiveness of ensemble algorithms (Random Forest, Adaboost, Extra trees classifier, Gradient Boosted Tree)

In the first part of this assignment, you will use Classification Trees for predicting if a user has a default payment option active or not. You can find the necessary data for performing this assignment [here \(https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients\)](https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients)

This dataset is aimed at the case of customer default payments in Taiwan. From the perspective of risk management, the result of predictive accuracy of the estimated probability of default will be more valuable than the binary result of classification - credible or not credible clients. Because the real probability of default is unknown, this study presented the novel Sorting Smoothing Method to estimate the real probability of default.

Required imports for this project are given below. Make sure you have all libraries required for this project installed. You may use conda or pip based on your set up.

**NOTE:** Since data is in Excel format you need to install `xlrd` in order to read the excel file inside your pandas dataframe. You can run `pip install xlrd` to install

```
In [1]: 1 import numpy as np
        2 import pandas as pd
```

After installing the necessary libraries, proceed to download the data. Since reading the excel file won't create headers by default, we added two more operations to substitute the columns.

```
In [2]: 1 main_dataset = pd.read_excel("https://archive.ics.uci.edu/ml/machi
2 main_dataset.columns = main_dataset.iloc[0]
3 main_dataset.drop(['ID'], axis=1, inplace=True)
```

```
In [3]: 1 main_dataset.drop([0], axis=0, inplace=True)
```

In the following, you can take a look into the dataset.

```
In [4]: 1 main_dataset.head()
```

Out[4]:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	...
1	20000	2	2	1	24	2	2	-1	-1	-2	...
2	120000	2	2	2	26	-1	2	0	0	0	...
3	90000	2	2	2	34	0	0	0	0	0	...
4	50000	2	2	1	37	0	0	0	0	0	...
5	50000	1	2	1	57	-1	0	-1	0	0	...

5 rows × 24 columns

```
In [5]: 1 dataset = main_dataset.copy()
```

## Questions (15 points total)

### Question 1 (2 pts)

Build a classifier by using decision tree and calculate the confusion matrix. Try different hyper-parameters (at least two) and discuss the result.

```
In [6]: 1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.preprocessing import OneHotEncoder
```

```
In [7]: 1 import matplotlib.pyplot as plt
        2 from scipy import stats
        3 from sklearn.tree import DecisionTreeClassifier
        4 from IPython.display import SVG
        5 from graphviz import Source
        6 from sklearn import tree
        7 from sklearn.metrics import classification_report, confusion_matrix
        8 from sklearn.metrics import roc_auc_score
        9 from sklearn.metrics import f1_score
       10 from sklearn.metrics import accuracy_score
```

```
In [8]: 1 from prettytable import PrettyTable
        2 from tabulate import tabulate
```

## Cleaning data

```
In [9]: 1 columns_categorical = ['SEX', 'EDUCATION', 'MARRIAGE', 'PAY_0', 'F
```

```
In [10]: 1 pays = ["PAY_0", "PAY_2", "PAY_3", "PAY_4", "PAY_5", "PAY_6"]
        2 dataset[pays] = dataset[pays].replace(-2.0, np.nan)
```

```
In [11]: 1 dataset["MARRIAGE"].replace(0, 3, inplace=True)
```

```
In [12]: 1 dataset["EDUCATION"].replace(5, 4, inplace=True)
        2 dataset["EDUCATION"].replace(6, 4, inplace=True)
```

```
In [13]: 1 dataset.dropna(inplace=True)
```

```
In [14]: 1 target_column = "default payment next month"
        2 y = dataset[target_column]
        3 y=y.astype('int')
        4 X = dataset[[col for col in dataset.columns if col != target_column]]
```

## Train Test Split

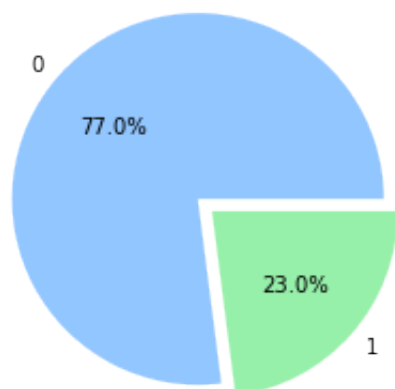
```
In [15]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25)
3 print(f"X_train.shape: {X_train.shape}")
4 print(f"X_test.shape: {X_test.shape}")
5 print(f"y_train.shape: {y_train.shape}")
6 print(f"y_test.shape: {y_test.shape}")
```

```
X_train.shape: (17579, 23)
X_test.shape: (5860, 23)
y_train.shape: (17579,)
y_test.shape: (5860,)
```

```
In [16]: 1 def check_data_balance(series, style="seaborn-pastel"):
2         with plt.style.context(style):
3             unique = series.value_counts()
4             display(unique) #show unique value counts of the target
5             plt.pie(unique, explode=[0.05]*len(unique), labels=unique)
```

```
In [17]: 1 check_data_balance(y_train)
```

```
0    13540
1     4039
Name: default payment next month, dtype: int64
```

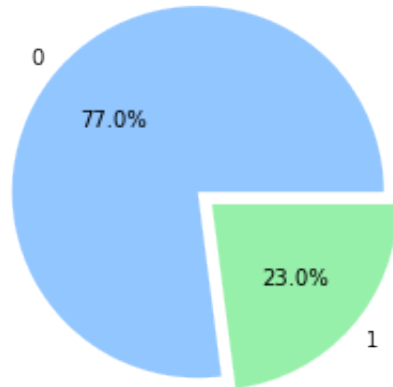


```
In [18]: 1 check_data_balance(y_test)
```

```
0    4514
```

```
1    1346
```

```
Name: default payment next month, dtype: int64
```



As you can see, our dataset is **imbalanced**. So metrics such as **ROC-AUC** or **F11-score** should be used and accuracy can't give us valid information about whether our model is having good results or not.

## Pipelines

```
In [19]: 1 pipeline_categorical = Pipeline([
2         ('onehot', OneHotEncoder(handle_unknown="ignore")),
3     ])
```

```
In [20]: 1 pipeline_full = ColumnTransformer([
2         ("categorical", pipeline_categorical, columns_categorical),
3     ])
```

Pipeline is only fitted on train data. Then both train and test data are transformed.

```
In [21]: 1 pipeline_full.fit(X_train)
2 X_train_transformed = pipeline_full.transform(X_train)
3 X_test_transformed = pipeline_full.transform(X_test)
4 print(f"X_train_transformed.shape: {X_train_transformed.shape}")
5 print(f"X_test_transformed.shape: {X_test_transformed.shape}")
```

```
X_train_transformed.shape: (17579, 68)
```

```
X_test_transformed.shape: (5860, 68)
```

## Tree Classifier

### First Decision Tree

```
In [22]: 1 tree_clf_1 = DecisionTreeClassifier(max_depth=2, criterion='entropy')
          2 tree_clf_1.fit(X_train_transformed, y_train)
```

```
Out[22]: DecisionTreeClassifier(criterion='entropy', max_depth=2)
```

```
In [23]: 1 tree_clf_1_predictions = tree_clf_1.predict(X_test_transformed)
          2 tree_clf_1_roc_auc = roc_auc_score(y_test, tree_clf_1_predictions)
          3 tree_clf_1_f1 = f1_score(y_test, tree_clf_1_predictions)
          4 tree_clf_1_accuracy = accuracy_score(y_test, tree_clf_1_predictions)
          5 tree_clf_1_confusion_matrix = confusion_matrix(y_test, tree_clf_1_predictions)
```

### Second Decision Tree

```
In [24]: 1 tree_clf_2 = DecisionTreeClassifier(max_depth=10, criterion='entropy')
          2 tree_clf_2.fit(X_train_transformed, y_train)
```

```
Out[24]: DecisionTreeClassifier(criterion='entropy', max_depth=10, splitter='random')
```

```
In [25]: 1 tree_clf_2_predictions = tree_clf_2.predict(X_test_transformed)
          2 tree_clf_2_roc_auc = roc_auc_score(y_test, tree_clf_2_predictions)
          3 tree_clf_2_f1 = f1_score(y_test, tree_clf_2_predictions)
          4 tree_clf_2_accuracy = accuracy_score(y_test, tree_clf_2_predictions)
          5 tree_clf_2_confusion_matrix = confusion_matrix(y_test, tree_clf_2_predictions)
```

### Third Decision Tree

```
In [26]: 1 tree_clf_3 = DecisionTreeClassifier(max_depth=20, criterion='gini')
          2 tree_clf_3.fit(X_train_transformed, y_train)
```

```
Out[26]: DecisionTreeClassifier(max_depth=20, splitter='random')
```

```
In [27]: 1 tree_clf_3_predictions = tree_clf_3.predict(X_test_transformed)
          2 tree_clf_3_roc_auc = roc_auc_score(y_test, tree_clf_3_predictions)
          3 tree_clf_3_f1 = f1_score(y_test, tree_clf_3_predictions)
          4 tree_clf_3_accuracy = accuracy_score(y_test, tree_clf_3_predictions)
          5 tree_clf_3_confusion_matrix = confusion_matrix(y_test, tree_clf_3_predictions)
```

### Fourth Decision Tree

```
In [28]: 1 tree_clf_4 = DecisionTreeClassifier(max_depth=10, criterion='gini')
          2 tree_clf_4.fit(X_train_transformed, y_train)
```

Out[28]: DecisionTreeClassifier(max\_depth=10)

```
In [29]: 1 tree_clf_4_predictions = tree_clf_4.predict(X_test_transformed)
          2 tree_clf_4_roc_auc = roc_auc_score(y_test, tree_clf_4_predictions)
          3 tree_clf_4_f1 = f1_score(y_test, tree_clf_4_predictions)
          4 tree_clf_4_accuracy = accuracy_score(y_test, tree_clf_4_predictions)
          5 tree_clf_4_confusion_matrix = confusion_matrix(y_test, tree_clf_4_predictions)
```

### Decision Trees Testing Results

```
In [30]: 1 tree_data = [
          2     ["Decision Tree 1", tree_clf_1_roc_auc, tree_clf_1_f1, tree_clf_1_accuracy],
          3     ["Decision Tree 2", tree_clf_2_roc_auc, tree_clf_2_f1, tree_clf_2_accuracy],
          4     ["Decision Tree 3", tree_clf_3_roc_auc, tree_clf_3_f1, tree_clf_3_accuracy],
          5     ["Decision Tree 4", tree_clf_4_roc_auc, tree_clf_4_f1, tree_clf_4_accuracy],
          6 ]
          7
          8 head = ["Model", "ROC-AUC Score", "F1 Score", "Accuracy"]
          9
         10 print(tabulate(tree_data, headers=head, tablefmt="grid"))
```

Model	ROC-AUC Score	F1 Score	Accuracy
Decision Tree 1	0.637826	0.435451	0.811945
Decision Tree 2	0.672473	0.504155	0.816724
Decision Tree 3	0.66191	0.483738	0.807679
Decision Tree 4	0.668973	0.497898	0.816553

```
In [31]: 1 confusion_matrix_data = [
2         ["Decision Tree 1", tree_clf_1_confusion_matrix],
3         ["Decision Tree 2", tree_clf_2_confusion_matrix],
4         ["Decision Tree 3", tree_clf_3_confusion_matrix],
5         ["Decision Tree 4", tree_clf_4_confusion_matrix]
6     ]
7
8     head = ["Model", "Confusion Matrix"]
9
10    print(tabulate(confusion_matrix_data, headers=head, tablefmt="grid"))
```

Model	Confusion Matrix
Decision Tree 1	[[4333 181] [ 921 425]]
Decision Tree 2	[[4240 274] [ 800 546]]
Decision Tree 3	[[4205 309] [ 818 528]]
Decision Tree 4	[[4252 262] [ 813 533]]

The accuracy in all models, is way more than ROC-AUC and F1-score which is because the data is unbalanced so Accuracy isn't a good method for evaluating the models.

As you can see in the confusion matrix, all the models, predict the '0' cases better than they predict '1' cases. The ratio of right-predicted-ones to wrong-predicted-ones is way less than the ratio of right-predicted-zeros to wrong-predicted-zeros.

### Question 2 (4 pts)

Try to build the decision tree which you built for the previous question, but this time by RandomizedSearchCV over hyper-parameters. Compare the results.

```
In [32]: 1 from sklearn.model_selection import RandomizedSearchCV
```



```
In [33]: 1 param = {'max_depth': [5, 10, 20, 50, None],
2             'max_features': [1, 10, 30, X_train_transformed.shape[1]]
3             'splitter' : ['best', 'random'],
4             'criterion' : ['gini', 'entropy'],
5             'min_samples_leaf' : [1, 2, 5],
6             'min_impurity_decrease' : [0.0, 0.1, 0.5]
7             }
```

```
In [34]: 1 rnd_search_tree = RandomizedSearchCV(DecisionTreeClassifier(), param
2 rnd_search_tree.fit(X_train_transformed, y_train)
3 rnd_search_tree.best_params_
4 rnd_search_tree.best_score_
```

```
In [35]: 1 rnd_tree_clf_predictions = rnd_search_tree.predict(X_test_transfor
2 rnd_tree_clf_roc_auc = roc_auc_score(y_test, rnd_tree_clf_predicti
3 rnd_tree_clf_f1 = f1_score(y_test, rnd_tree_clf_predictions)
4 rnd_tree_clf_accuracy = accuracy_score(y_test, rnd_tree_clf_predic
5 rnd_tree_clf_confusion_matrix = confusion_matrix(y_test, rnd_tree_
```

```
In [36]: 1 tabel_data = [
2     ["Decision Tree ", rnd_tree_clf_roc_auc, rnd_tree_clf_f1, rnd_
3 ]
4 head = ["RandomizedSearchCV Model", "ROC-AUC Score", "F1 Score", "
5 print(tabulate(tabel_data, headers=head, tablefmt="grid"))
```

```
+-----+-----+-----+-----+
----+-----+-----+-----+
| RandomizedSearchCV Model | ROC-AUC Score | F1 Score | Accur
acy | Best Score | Confusion Matrix |
+=====+=====+=====+=====+
===+=====+=====+=====+
| Decision Tree          | 0.660085 | 0.480406 | 0.807
679 | 0.480793 | [[4212 302]
|
|
|          | [ 825 521]]
+-----+-----+-----+-----+
----+-----+-----+-----+
```

```
In [37]: 1 rnd_search_tree.best_params
```

```
Out[37]: {'splitter': 'best',
'min_samples_leaf': 1,
'min_impurity_decrease': 0.0,
'max_features': 30,
'max_depth': None,
'criterion': 'gini'}
```

The ROC-AUC score and F1-score of the decision tree model built using RandomizedSearchCV, is almost the same as the best decision tree model(second model) that we made changing the hyperparameters in the previous section.

### Question 3 (6 pts)

Try to build the same classifier by using following ensemble models. For each of these models calculate accuracy and at least for two in the list below, plot the learning curves.

- Random Forest
- AdaBoost
- Extra Trees Classifier
- Gradient Boosted Trees

```
In [38]: 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.ensemble import AdaBoostClassifier
3 from sklearn.ensemble import ExtraTreesClassifier
4 from sklearn.ensemble import GradientBoostingClassifier
5 from sklearn.model_selection import learning_curve
```

Since using all the possible parameters of the classifier in the RandomizedSearchCV could make a lot of possible situations and checking even most of those situations could take a lot of time, only some of the parameters are selected and used in RandomizedSearchCV.

### Random Forest

```
In [39]: 1 rf_params = {'n_estimators' : [2, 10, 50],
2               'criterion' : ['gini', 'entropy'],
3               'max_features': [1, 10, 30, X_train_transformed.shape
4               'max_depth': [5, 10, 20, None],
5               'bootstrap': [True, False]
6               }
```

```
In [40]: 1 rf = RandomForestClassifier()
2         rf_random = RandomizedSearchCV(estimator=rf, param_distributions=r
3         rf_random.fit(X_train_transformed, y_train)
4         rf_random_best_params = rf_random.best_params_
5         rf_random_best_score = rf_random.best_score_
6         rf_random_predictions = rf_random.predict(X_test_transformed)
7         rf_random_roc_auc = roc_auc_score(y_test, rf_random_predictions)
8         rf_random_f1 = f1_score(y_test, rf_random_predictions)
9         rf_random_accuracy = accuracy_score(y_test, rf_random_predictions)
10        rf_random_confusion_matrix = confusion_matrix(y_test, rf_random_pr
```

```
In [41]: 1 tabel_data.append(["Random Forest", rf_random_roc_auc, rf_random_f
2         rf_random_accuracy, rf_random_best_score, rf_ra
```

```
In [42]: 1 rf_random_best_params
```

```
Out[42]: {'n_estimators': 50,
          'max_features': 68,
          'max_depth': 10,
          'criterion': 'entropy',
          'bootstrap': False}
```

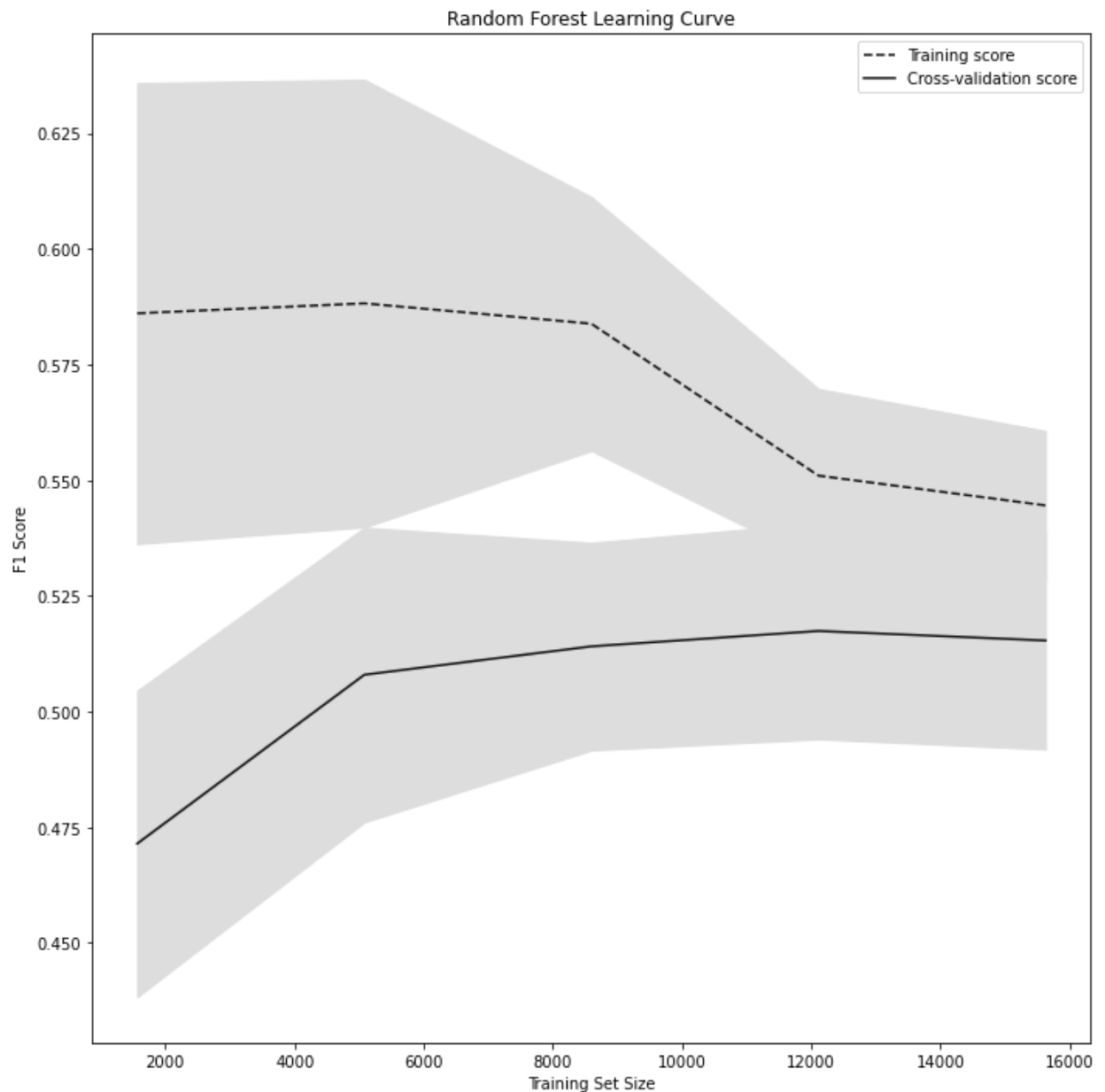
```
In [43]: 1 rf_train_sizes, rf_train_scores, rf_test_scores = learning_curve(r
2         cv=9, scor
```

```
In [44]: 1 rf_train_mean = np.mean(rf_train_scores, axis=1)
2         rf_train_std = np.std(rf_train_scores, axis=1)
3         rf_test_mean = np.mean(rf_test_scores, axis=1)
4         rf_test_std = np.std(rf_test_scores, axis=1)
```

```

In [45]: 1 plt.subplots(1, figsize=(10,10))
          2 plt.plot(rf_train_sizes, rf_train_mean, '--', color="#111111", label="Training score")
          3 plt.plot(rf_train_sizes, rf_test_mean, color="#111111", label="Cross-validation score")
          4
          5 plt.fill_between(rf_train_sizes, rf_train_mean - rf_train_std, rf_train_mean + rf_train_std, color="#cccccc")
          6 plt.fill_between(rf_train_sizes, rf_test_mean - rf_test_std, rf_test_mean + rf_test_std, color="#cccccc")
          7
          8 plt.title("Random Forest Learning Curve")
          9 plt.xlabel("Training Set Size"), plt.ylabel("F1 Score"), plt.legend()
         10 plt.tight_layout()
         11 plt.show()

```



## AdaBoost

```
In [46]: 1 ada_boost_params = {'n_estimators' : [2, 10, 50],
2                               'learning_rate' : [0.1, 0.5, 1],
3                               'algorithm' : ['SAMME.R', 'SAMME']}
4
```

```
In [47]: 1 ada_boost = AdaBoostClassifier()
2   ada_random = RandomizedSearchCV(estimator = ada_boost, param_distr
3   ada_random.fit(X_train_transformed, y_train)
4   ada_random_best_params = ada_random.best_params_
5   ada_random_best_score = ada_random.best_score_
6   ada_random_predictions = ada_random.predict(X_test_transformed)
7   ada_random_roc_auc = roc_auc_score(y_test, ada_random_predictions)
8   ada_random_f1 = f1_score(y_test, ada_random_predictions)
9   ada_random_accuracy = accuracy_score(y_test, ada_random_prediction
10  ada_random_confusion_matrix = confusion_matrix(y_test, ada_random_
```

```
In [48]: 1 tabel_data.append(["Ada Boost", ada_random_roc_auc, ada_random_f1,
2                               ada_random_accuracy, ada_random_best_score, ada
```

```
In [49]: 1 ada_random_best_params
```

```
Out[49]: {'n_estimators': 50, 'learning_rate': 1, 'algorithm': 'SAMME.R'}
```

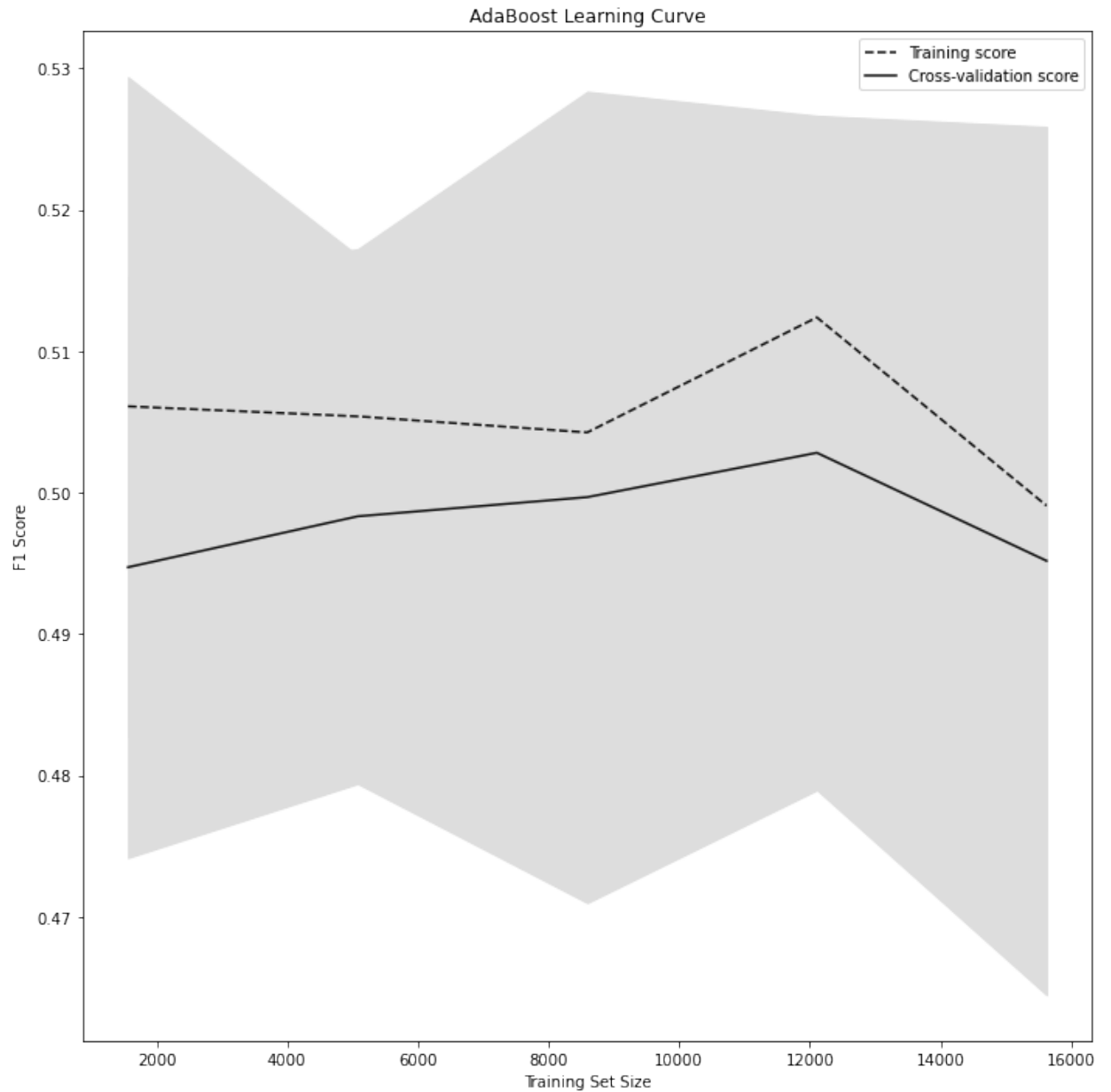
```
In [50]: 1 ada_train_sizes, ada_train_scores, ada_test_scores = learning_curve
2                                               cv=9, scor
```

```
In [51]: 1 ada_train_mean = np.mean(ada_train_scores, axis=1)
2   ada_train_std = np.std(ada_train_scores, axis=1)
3   ada_test_mean = np.mean(ada_test_scores, axis=1)
4   ada_test_std = np.std(ada_test_scores, axis=1)
```

```

In [52]: 1 plt.subplots(1, figsize=(10,10))
          2 plt.plot(ada_train_sizes, ada_train_mean, '--', color="#111111",
          3 plt.plot(ada_train_sizes, ada_test_mean, color="#111111", label="C
          4
          5 plt.fill_between(ada_train_sizes, ada_train_mean - ada_train_std,
          6 plt.fill_between(ada_train_sizes, ada_test_mean - ada_test_std, ac
          7
          8 plt.title("AdaBoost Learning Curve")
          9 plt.xlabel("Training Set Size"), plt.ylabel("F1 Score"), plt.leg
         10 plt.tight_layout()
         11 plt.show()

```



## Extra Trees Classifier

```
In [53]: 1 extra_tree_clf_params = {'n_estimators' : [2, 10, 50],
2                                     'criterion' : ['gini', 'entropy'],
3                                     'max_features': [1, 10, 30, X_train_trans
4                                     'max_depth': [5, 10, 20, None],
5                                     'bootstrap': [True, False]
6                                     }
```

```
In [54]: 1 extra_tree_clf = ExtraTreesClassifier()
2 extra_tree_clf_random = RandomizedSearchCV(estimator = extra_tree_clf,
3                                             param_distributions=extra_tree_clf_params,
4                                             cv=5, n_jobs=-1, verbose=1)
5 extra_tree_clf_random.fit(X_train_transformed, y_train)
6 extra_tree_clf_random.best_params_ = ada_random.best_params_
7 extra_tree_clf_random.best_score_ = ada_random.best_score_
8 extra_tree_clf_random.predictions = ada_random.predict(X_test_transformed)
9 extra_tree_clf_random_roc_auc = roc_auc_score(y_test, extra_tree_clf_random.predict(X_test_transformed))
10 extra_tree_clf_random_f1 = f1_score(y_test, extra_tree_clf_random.predict(X_test_transformed))
11 extra_tree_clf_random_accuracy = accuracy_score(y_test, extra_tree_clf_random.predict(X_test_transformed))
12 extra_tree_clf_random_confusion_matrix = confusion_matrix(y_test, extra_tree_clf_random.predict(X_test_transformed))
```

```
In [55]: 1 tabel_data.append(["Extra Tree Classifier", extra_tree_clf_random.best_params_,
2                                     extra_tree_clf_random_f1, extra_tree_clf_random_roc_auc,
3                                     extra_tree_clf_random_accuracy, extra_tree_clf_random_confusion_matrix])
```

```
In [56]: 1 extra_tree_clf_random.best_params_
```

```
Out[56]: {'n_estimators': 50, 'learning_rate': 1, 'algorithm': 'SAMME.R'}
```

## Gradient Boosted Trees

```
In [57]: 1 gradient_boosting_clf_params = {'loss' : ['deviance', 'exponential'],
2                                     'learning_rate' : [0.1, 0.5, 1],
3                                     'n_estimators' : [2, 10, 50],
4                                     'criterion' : ['squared_error', 'entropy'],
5                                     'max_features': [1, 10, 30, X_train_transformed],
6                                     'max_depth': [5, 10, 20, None]
7                                     }
```

```
In [58]: 1 gradient_boosting_clf = GradientBoostingClassifier()
2 gradient_boosting_clf_random = RandomizedSearchCV(estimator = grad
3                                                    param_distributi
4                                                    scoring='f1')
5 gradient_boosting_clf_random.fit(X_train_transformed, y_train)
6 gradient_boosting_clf_random_best_params = ada_random.best_params_
7 gradient_boosting_clf_best_score = ada_random.best_score_
8 gradient_boosting_clf_random_predictions = ada_random.predict(X_te
9 gradient_boosting_clf_random_roc_auc = roc_auc_score(y_test, gradi
10 gradient_boosting_clf_random_f1 = f1_score(y_test, gradient_boosti
11 gradient_boosting_clf_random_accuracy = accuracy_score(y_test, gra
12 gradient_boosting_clf_random_confusion_matrix=confusion_matrix(y_t
```

The score on this train-test partition for these parameters will be set to nan. Details:

Traceback (most recent call last):

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/model\_selection/\_validation.py", line 593, in \_fit\_and\_score

estimator.fit(X\_train, y\_train, \*\*fit\_params)

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/ensemble/\_gb.py", line 504, in fit

n\_stages = self.\_fit\_stages()

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/ensemble/\_gb.py", line 561, in \_fit\_stages

raw\_predictions = self.\_fit\_stage()

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/ensemble/\_gb.py", line 214, in \_fit\_stage

tree.fit(X, residual, sample\_weight=sample\_weight,

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/tree/\_classes.py", line 1247, in fit

super().fit()

File "/Users/Negin/opt/anaconda3/lib/python3.8/site-packages/sklearn/tree/\_classes.py", line 350, in fit

```
In [59]: 1 tabel_data.append(["Gradient Boosted Trees", gradient_boosting_clf
2                       gradient_boosting_clf_random_f1, gradient_boos
3                       gradient_boosting_clf_best_score, gradient_boos
```

```
In [60]: 1 gradient_boosting_clf_random_best_params
```

```
Out[60]: {'n_estimators': 50, 'learning_rate': 1, 'algorithm': 'SAMME.R'}
```



```
In [61]: 1 print(tabulate(tabel_data, headers=head, tablefmt="grid"))
```

RandomizedSearchCV Model				ROC-AUC Score	F1 Score	Accuracy
Best Score   Confusion Matrix						
Decision Tree				0.660085	0.480406	0.807
679	0.480793	[[4212	302]			
		[ 825	521]]			
Random Forest				0.672805	0.504854	0.817
235	0.510989	[[4243	271]			
		[ 800	546]]			
Ada Boost				0.676246	0.511888	0.821
331	0.513985	[[4264	250]			
		[ 797	549]]			
Extra Tree Classifier				0.676246	0.511888	0.821
331	0.513985	[[4264	250]			
		[ 797	549]]			
Gradient Boosted Trees				0.676246	0.511888	0.821
331	0.513985	[[4264	250]			
		[ 797	549]]			

#### Question 4 (3 pts)

Discuss and compare the results for the all past three questions.

- How does changing hyperparms effect model performance?
- Why do you think certain models performed better/worse?
- How does this performance line up with known strengths/weakness of these models?

Changing the hyperparameters did change the performance of the models but changes were not that huge.

Depending on the way each decision tree is formed, the shape and result of the decision trees could differ a lot. The models that split the input by more important features first, usually have better results. Things like this could change the output of a decision tree a lot.

We know that decision trees can overfit easily, In all of our results, the models predicted zeros better than ones. Since our data is imbalanced and has more zeros, I think all of our models are overfit and the ratio of right-predicted-ones to wrong-predicted-ones is way less than the ratio of right-predicted-zeros to wrong-predicted-zeros.

In [ ]:

1