

index

March 20, 2022

1 Hyperparameter Tuning and Pruning in Decision Trees - Lab

1.1 Introduction

In this lab, you will use the titanic dataset to see the impact of tree pruning and hyperparameter tuning on the predictive performance of a decision tree classifier. Pruning reduces the size of decision trees by removing nodes of the tree that do not provide much predictive power to classify instances. Decision trees are the most susceptible out of all the machine learning algorithms to overfitting and effective pruning can reduce this likelihood.

1.2 Objectives

In this lab you will:

- Determine the optimal hyperparameters for a decision tree model and evaluate the model performance

1.3 Import necessary libraries

Let's first import the libraries you'll need for this lab.

```
[1]: import numpy as np
import pandas as pd
import matplotlib as mpl
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import roc_curve, auc
plt.style.use('seaborn')
```

1.4 Import the data

The titanic dataset, available in 'titanic.csv', is all cleaned up and preprocessed for you so that you can focus on pruning and optimization. Import the dataset and print the first five rows of the data:

```
[2]: # Import the data
df = pd.read_csv("titanic.csv")
```

1.5 Create training and test sets

- Assign the 'Survived' column to y
- Drop the 'Survived' and 'PassengerId' columns from df, and assign the resulting DataFrame to X
- Split X and y into training and test sets. Assign 30% to the test set and set the random_state to SEED

```
[4]: # Create X and y
y = df["Survived"]
X = df.drop(columns = ["Survived", "PassengerId"], axis = 1)

# Split into training and test sets
SEED = 1
X_train, X_test, y_train, y_test = train_test_split(X,y,
                                                    test_size = 0.3,
                                                    random_state = SEED)
```

1.6 Train a vanilla classifier

Note: The term “vanilla” is used for a machine learning algorithm with its default settings (no tweaking/tuning).

- Instantiate a decision tree
 - Use the 'entropy' criterion and set the random_state to SEED
- Fit this classifier to the training data

```
[7]: # Train the classifier using training data
dt = DecisionTreeClassifier(criterion = "entropy", random_state = SEED)
dt.fit(X_train, y_train)
```

```
[7]: DecisionTreeClassifier(criterion='entropy', random_state=1)
```

1.7 Make predictions

- Create a set of predictions using the test set
- Using y_test and y_pred, calculate the AUC (Area under the curve) to check the predictive performance

```
[9]: # Make predictions using test set
y_pred = dt.predict(X_test)

# Check the AUC of predictions
false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test, y_pred)
roc_auc = auc(false_positive_rate, true_positive_rate)
roc_auc
```

```
[9]: 0.7367718101733446
```

1.8 Maximum Tree Depth

Let's first check for the best depth parameter for our decision tree:

- Create an array for `max_depth` values ranging from 1 - 32
- In a loop, train the classifier for each depth value (32 runs)
- Calculate the training and test AUC for each run
- Plot a graph to show under/overfitting and the optimal value
- Interpret the results

```
[18]: # Identify the optimal tree depth for given data
import seaborn as sns
max_depth = range(1, 33)

auc_values_test = []
auc_values_train = []

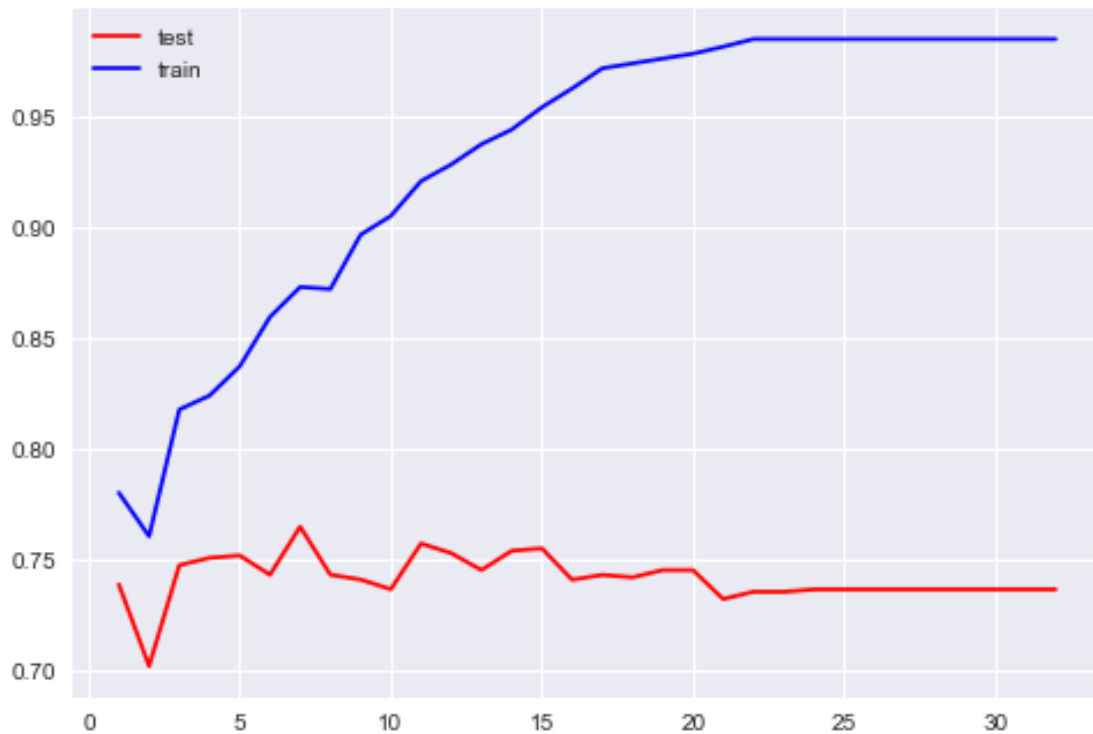
for i in max_depth:
    dt = DecisionTreeClassifier(criterion = "entropy",
                               random_state = SEED,
                               max_depth = i)

    dt.fit(X_train, y_train)

    ## Test Prediction
    y_pred_test = dt.predict(X_test)
    fpr_test, tpr_test, thresholds_test = roc_curve(y_test, y_pred_test)
    roc_auc_test = auc(fpr_test, tpr_test)
    auc_values_test.append(roc_auc_test)

    ## Train Prediction
    y_pred_train = dt.predict(X_train)
    fpr_train, tpr_train, thresholds_train = roc_curve(y_train, y_pred_train)
    roc_auc_train = auc(fpr_train, tpr_train)
    auc_values_train.append(roc_auc_train)
#     sns.lineplot(x = i, y = roc_auc)

sns.lineplot(x = range(1, 33), y = auc_values_test, color = "r", label = "test")
sns.lineplot(x = range(1, 33), y = auc_values_train, color = "blue", label = "train");
```



```
[28]: # Your observations here

## We see overfitting for values more than 3 because the AUCs converges fast
## Around depth 2 we have minimum AUCs and since for test set for depths
## between 3 and 15 we have almost a flat diagram with some fluctuations
## we should find the maximum depth within this range. However, for max_depth
## more than 3 the diagrams converge rapidly and we can conclude that the
## optimal max_depth is 3

# Training error decreases with increasing tree depth - clear sign of
↳ overfitting
# Test error increases after depth=3 - nothing more to learn from deeper trees
# (some fluctuations, but not stable)
# Training and test errors rise rapidly between the depths of 2 and 3
# Optimal value seen here is 3
```

1.9 Minimum Sample Split

Now check for the best `min_samples_splits` parameter for our decision tree

- Create an array for `min_sample_splits` values ranging from 0.1 - 1 with an increment of 0.1
- In a loop, train the classifier for each `min_samples_splits` value (10 runs)

- Calculate the training and test AUC for each run
- Plot a graph to show under/overfitting and the optimal value
- Interpret the results

```
[29]: # Identify the optimal min-samples-split for given data
import seaborn as sns
min_sample_splits = np.linspace(0.1, 1, num = 10)

auc_values_test = []
auc_values_train = []

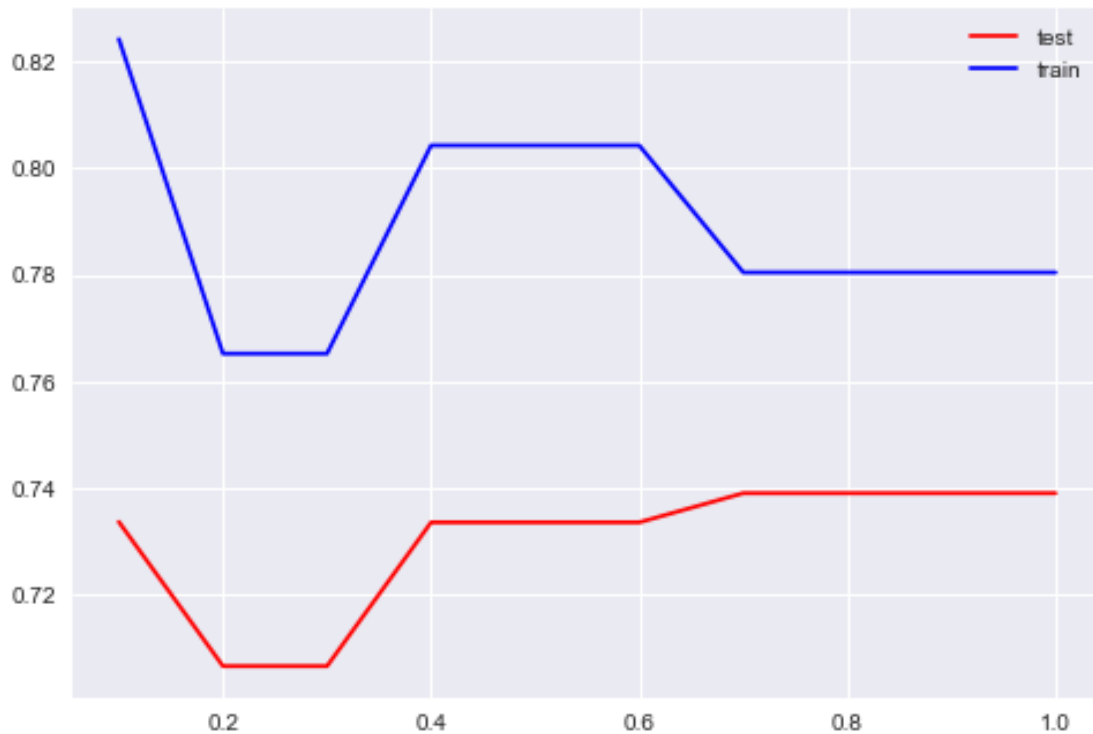
for i in min_sample_splits:
    dt = DecisionTreeClassifier(criterion = "entropy",
                               random_state = SEED,
                               min_samples_split = i)

    dt.fit(X_train, y_train)

    ## Test Prediction
    y_pred_test = dt.predict(X_test)
    fpr_test, tpr_test, thresholds_test = roc_curve(y_test, y_pred_test)
    roc_auc_test = auc(fpr_test, tpr_test)
    auc_values_test.append(roc_auc_test)

    ## Train Prediction
    y_pred_train = dt.predict(X_train)
    fpr_train, tpr_train, thresholds_train = roc_curve(y_train, y_pred_train)
    roc_auc_train = auc(fpr_train, tpr_train)
    auc_values_train.append(roc_auc_train)
#     sns.lineplot(x = i, y = roc_auc)

sns.lineplot(x = min_sample_splits, y = auc_values_test,
             color = "r", label = "test")
sns.lineplot(x = min_sample_splits, y = auc_values_train,
             color = "blue", label = "train");
```



```
[30]: # Your observations here

# AUC for both test and train data stabilizes at 0.7
# Further increase in minimum sample split does not improve learning
```

1.10 Minimum Sample Leafs

Now check for the best `min_samples_leafs` parameter value for our decision tree

- Create an array for `min_samples_leafs` values ranging from 0.1 - 0.5 with an increment of 0.1
- In a loop, train the classifier for each `min_samples_leafs` value (5 runs)
- Calculate the training and test AUC for each run
- Plot a graph to show under/overfitting and the optimal value
- Interpret the results

```
[31]: # Calculate the optimal value for minimum sample leafs
import seaborn as sns
min_samples_leafs = np.linspace(0.1, 0.5, num = 5)

auc_values_test = []
```

```

auc_values_train = []

for i in min_samples_leafs:
    dt = DecisionTreeClassifier(criterion = "entropy",
                                random_state = SEED,
                                min_samples_leaf = i)

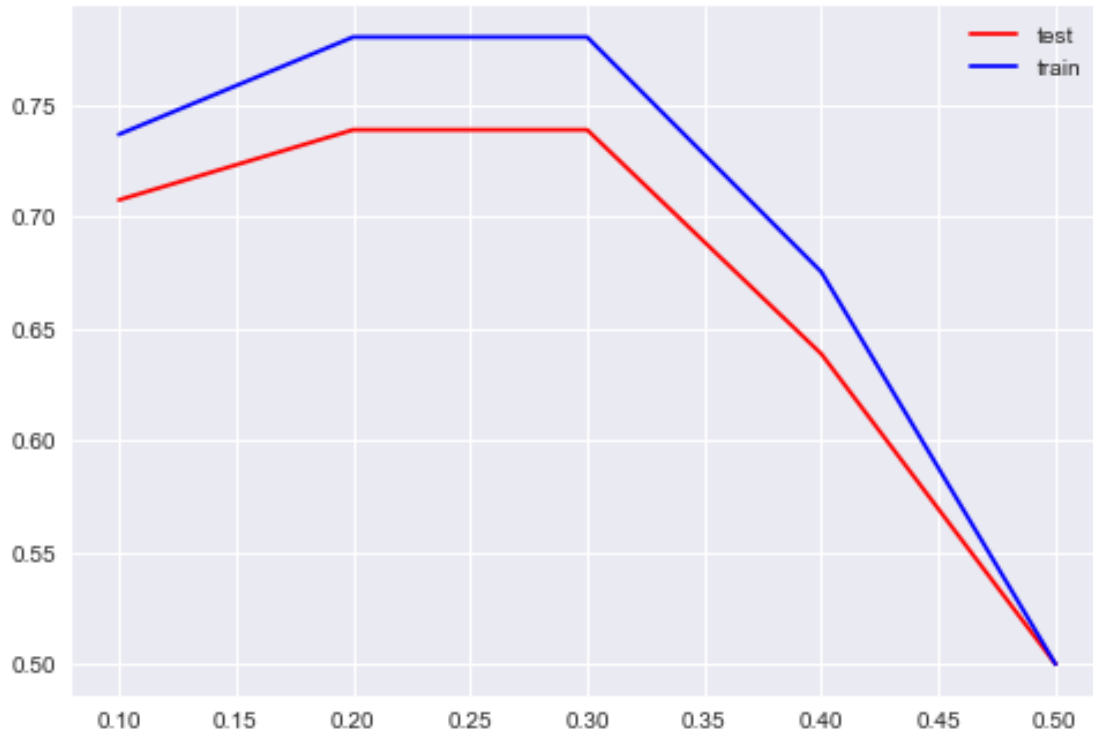
    dt.fit(X_train, y_train)

    ## Test Prediction
    y_pred_test = dt.predict(X_test)
    fpr_test, tpr_test, thresholds_test = roc_curve(y_test, y_pred_test)
    roc_auc_test = auc(fpr_test, tpr_test)
    auc_values_test.append(roc_auc_test)

    ## Train Prediction
    y_pred_train = dt.predict(X_train)
    fpr_train, tpr_train, thresholds_train = roc_curve(y_train, y_pred_train)
    roc_auc_train = auc(fpr_train, tpr_train)
    auc_values_train.append(roc_auc_train)
#     sns.lineplot(x = i, y = roc_auc)

sns.lineplot(x = min_samples_leafs, y = auc_values_test,
              color = "r", label = "test")
sns.lineplot(x = min_samples_leafs, y = auc_values_train,
              color = "blue", label = "train");

```



```
[ ]: # Your observations here

# AUC gives best value between 0.2 and 0.3 for both test and training sets
# The accuracy drops down if we continue to increase the parameter value
```

1.11 Maximum Features

Now check for the best `max_features` parameter value for our decision tree

- Create an array for `max_features` values ranging from 1 - 12 (1 feature vs all)
- In a loop, train the classifier for each `max_features` value (12 runs)
- Calculate the training and test AUC for each run
- Plot a graph to show under/overfitting and the optimal value
- Interpret the results

```
[32]: # Find the best value for optimal maximum feature size
import seaborn as sns
max_features = list(range(1, X_train.shape[1]))

auc_values_test = []
auc_values_train = []
```



```

for i in max_features:
    dt = DecisionTreeClassifier(criterion = "entropy",
                               random_state = SEED,
                               max_features = i)

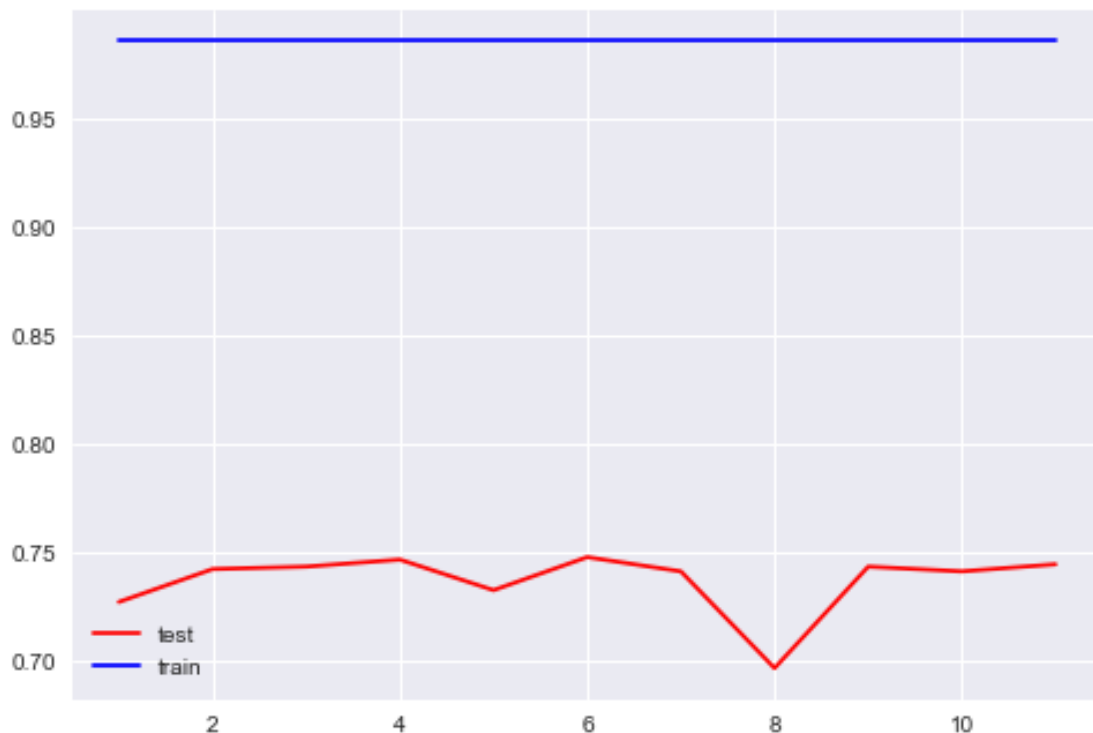
    dt.fit(X_train, y_train)

    ## Test Prediction
    y_pred_test = dt.predict(X_test)
    fpr_test, tpr_test, thresholds_test = roc_curve(y_test, y_pred_test)
    roc_auc_test = auc(fpr_test, tpr_test)
    auc_values_test.append(roc_auc_test)

    ## Train Prediction
    y_pred_train = dt.predict(X_train)
    fpr_train, tpr_train, thresholds_train = roc_curve(y_train, y_pred_train)
    roc_auc_train = auc(fpr_train, tpr_train)
    auc_values_train.append(roc_auc_train)

sns.lineplot(x = max_features, y = auc_values_test,
             color = "r", label = "test")
sns.lineplot(x = max_features, y = auc_values_train,
             color = "blue", label = "train");

```



```
[ ]: # Your observations here

# No clear effect on the training dataset - flat AUC
# Some fluctuations in test AUC but not definitive enough to make a judgement
# Highest AUC value seen at 6
```

1.12 Re-train the classifier with chosen values

Now we will use the best values from each training phase above and feed it back to our classifier. Then we can see if there is any improvement in predictive performance.

- Train the classifier with the optimal values identified
- Compare the AUC of the new model with the earlier vanilla decision tree AUC
- Interpret the results of the comparison

```
[34]: # Train a classifier with optimal values identified above
dt = DecisionTreeClassifier(
    criterion = "entropy",
    random_state = SEED,
    max_depth = 3,
    min_samples_leaf = 0.25,
    min_samples_split = 0.7,
    max_features = 6
)

dt.fit(X_train, y_train)

y_pred_test = dt.predict(X_test)

false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test,
    ↪ y_pred_test)
roc_auc = auc(fpr_test, tpr_test)
roc_auc
```

```
[34]: 0.7443876101165104
```

```
[87]: X_test_copy = X_test.copy()
y_test_copy = list(y_test)
X_test_copy.reset_index(inplace = True, drop = True)
y_test_copy
i=6
type(X_test_copy)
```

```
y_pred_test = dt.predict(X_test_copy.iloc[i].to_frame().T)
# X_test_copy.iloc[0].to_frame().T
print(y_pred_test)
print(y_test_copy[i])
```

```
[1]
```

```
0
```

```
[35]: # Your observations here
```

```
# We improved the AUC from 0.73 in the vanilla classifier to 0.74 with some ↵
↵tuning.
# Due to randomness, results may slightly differ, there is some improvement in ↵
↵most cases.
# With more complicated (and bigger) datasets,
# we might see an even bigger improvement in AUC/accuracy of the classifier.
```

In the next section, we shall talk about hyperparameter tuning using a technique called “grid-search” to make this process even more granular and decisive.

1.13 Summary

In this lesson, we looked at tuning a decision tree classifier in order to avoid overfitting and increasing the generalization capabilities of the classifier. For the titanic dataset, we see that identifying optimal parameter values can result in some improvements towards predictions. This idea will be exploited further in upcoming lessons and labs.