

[illegible][illegible]

```

splits. Also keep track of the highest test accuracy observed, the tree for which the highest test accuracy
depth, the least size and the least training time observed. Repeat these steps for the next ten splits of
impurity measure. Compare the best trees (ones with highest test accuracies) obtained in the two co
In (20):
all_trees = []
avg_accuracy_ig = 0
avg_height_ig = 0
avg_size_ig = 0
avg_time_ig = 0

best_accuracy_ig = 0
best_height_ig = float('inf')
best_size_ig = float('inf')
best_time_ig = float('inf')

best_tree_test_set_ig = None
best_tree_cf_mat_ig = None
most_accurate_tree_ig = None

for split in range(10):
    train_set, test_set = RandomSplit()
    start = time()
    tree_ig = Construct_Tree(train_set, test_set, use_gini_impurity_measure = False)
    t = time() - start
    acc = tree_ig.accuracy(test_set)
    if acc > best_accuracy_ig:
        best_accuracy_ig = acc
        most_accurate_tree_ig = tree_ig
        best_tree_cf_mat_ig = tree_ig.confusion_matrix(test_set)
        best_tree_test_set_ig = test_set

    avg_accuracy_ig += acc
    h, s = tree_ig.height(), tree_ig.size()
    avg_height_ig += h
    avg_size_ig += s
    avg_time_ig += t

    best_height_ig = min(best_height_ig, h)
    best_size_ig = min(best_size_ig, s)
    best_time_ig = min(best_time_ig, t)

    all_trees.append(tree_ig)

avg_accuracy_ig /= 10
avg_height_ig /= 10
avg_size_ig /= 10
avg_time_ig /= 10

In (21):
avg_accuracy_gg = 0
avg_height_gg = 0
avg_size_gg = 0
avg_time_gg = 0

best_accuracy_gg = 0
best_height_gg = float('inf')
best_size_gg = float('inf')
best_time_gg = float('inf')

best_tree_test_set_gg = None
best_tree_cf_mat_gg = None
most_accurate_tree_gg = None

for split in range(10):
    train_set, test_set = RandomSplit()
    start = time()
    tree_gg = Construct_Tree(train_set, test_set, use_gini_impurity_measure = True)
    t = time() - start
    acc = tree_gg.accuracy(test_set)
    if acc > best_accuracy_gg:
        best_accuracy_gg = acc
        most_accurate_tree_gg = tree_gg
        best_tree_cf_mat_gg = tree_gg.confusion_matrix(test_set)
        best_tree_test_set_gg = test_set

    avg_accuracy_gg += acc
    h, s = tree_gg.height(), tree_gg.size()
    avg_height_gg += h
    avg_size_gg += s
    avg_time_gg += t

    best_height_gg = min(best_height_gg, h)
    best_size_gg = min(best_size_gg, t)
    best_time_gg = min(best_time_gg, t)

    all_trees.append(tree_gg)

avg_accuracy_gg /= 10
avg_height_gg /= 10
avg_size_gg /= 10
avg_time_gg /= 10

In (22):
print('AVG. ACCURACY (GINI IMPURITY MEASURE) :', avg_accuracy_gg)
print('AVG. HEIGHT (GINI IMPURITY MEASURE) :', avg_height_gg)
print('AVG. SIZE (GINI IMPURITY MEASURE) :', avg_size_gg)
print('AVG. TRAINING TIME (GINI IMPURITY MEASURE) :', avg_time_gg)

print('BEST ACCURACY (GINI IMPURITY MEASURE) :', best_accuracy_gg)
print('BEST HEIGHT (GINI IMPURITY MEASURE) :', best_height_gg)
print('BEST SIZE (GINI IMPURITY MEASURE) :', best_size_gg)
print('BEST TRAINING TIME (GINI IMPURITY MEASURE) :', best_time_gg)

AVG. ACCURACY (GINI IMPURITY MEASURE) : 73.75886524822695
AVG. HEIGHT (GINI IMPURITY MEASURE) : 18.9
AVG. SIZE (GINI IMPURITY MEASURE) : 211.6
AVG. TRAINING TIME (GINI IMPURITY MEASURE) : 0.5856430768966675
BEST ACCURACY (GINI IMPURITY MEASURE) : 80.141849714322
BEST HEIGHT (GINI IMPURITY MEASURE) : 15
BEST SIZE (GINI IMPURITY MEASURE) : 179
BEST TRAINING TIME (GINI IMPURITY MEASURE) : 0.5194697380065918

In (23):
print('AVG. ACCURACY (IG IMPURITY MEASURE) :', avg_accuracy_ig)
print('AVG. HEIGHT (IG IMPURITY MEASURE) :', avg_height_ig)
print('AVG. SIZE (IG IMPURITY MEASURE) :', avg_size_ig)
print('AVG. TRAINING TIME (IG IMPURITY MEASURE) :', avg_time_ig)

print('BEST ACCURACY (IG IMPURITY MEASURE) :', best_accuracy_ig)
print('BEST HEIGHT (IG IMPURITY MEASURE) :', best_height_ig)
print('BEST SIZE (IG IMPURITY MEASURE) :', best_size_ig)
print('BEST TRAINING TIME (IG IMPURITY MEASURE) :', best_time_ig)

AVG. ACCURACY (IG IMPURITY MEASURE) : 73.82978723044256
AVG. HEIGHT (IG IMPURITY MEASURE) : 20.3
AVG. SIZE (IG IMPURITY MEASURE) : 190.8
BEST TRAINING TIME (IG IMPURITY MEASURE) : 0.6182325601577958
BEST ACCURACY (IG IMPURITY MEASURE) : 79.43262411347517
BEST HEIGHT (IG IMPURITY MEASURE) : 15
BEST SIZE (IG IMPURITY MEASURE) : 179
BEST TRAINING TIME (IG IMPURITY MEASURE) : 0.516204187286377

In (25):
Plot_And_Save_Confusion_Matrix(best_tree_cf_mat_gg, 'cf3.png')

```

PREDICTED LABELS	TRUE LABELS		
	Patient	Non Patient	
Patient	True Pos 66 46.81%	False Pos 11 7.80%	-0.4
	False Neg 17 12.06%	True Neg 47 33.33%	-0.3
Non Patient			-0.2
			-0.1

In [26]:

```
Plot_And_Save_Confusion_Matrix(best_tree_of_mat_ig, 'cf4.png')
```

PREDICTED LABELS	TRUE LABELS	
	Patient	Non Patient
Patient	True Pos 63 44.68%	False Pos 16 11.35%
Non Patient	False Neg 13 9.22%	True Neg 49 34.75%

In [27]:

```
best_tree = None
best_tree_test_acc = None
if best_accuracy_ig > best_accuracy_gg :
    best_tree = most_accurate_tree_ig
```

```
else
    best_tree = most_accurate_tree.py
    best_test_set = best_test_test_set.py
```

```

for idx, h in enumerate(test_training_history['height_progress']):
    if test_h in test_accuracy_vs_height:
        best_accuracy_vs_height[h] = test_training_history['test_accuracy']
    else:
        best_accuracy_vs_height[h] = max(test_accuracy_vs_height[h],
                                          test_training_history['test_acc'])
    if test_training_history['test_accuracy_progress'][idx] > highest_overall_
highest_overall_test_acc = test_training_history['test_accuracy_progress']
tree_with_optimal_depth = test_training_history['tree_progress'][idx]

print('OPTIMAL DEPTH LIMIT FOR THIS DATASET : ', tree_with_optimal_depth.height())

OPTIMAL DEPTH LIMIT FOR THIS DATASET : 12

In (30):
best_accuracy_vs_height

Out[30]:
0: 65.95744680851064,
1: 63.829787234042556,
2: 73.049646338007092,
3: 74.46808510638297,
4: 75.1773049646339,
5: 75.1773049646339,
6: 76.59574468085106,
7: 78.01418439716312,
8: 79.43026411347517,
9: 80.85106382978724,

```

```

11: 80.85106382978724,
12: 82.26950354609929,
13: 83.80.85106382978724,
14: 79.43262411347517,
15: 81.56028386794327,
16: 81.56028386794327,
17: 80.1418439716312,
18: 79.43262411347517,
19: 79.43262411347517,
20: 79.43262411347517,
21: 79.43262411347517,
22: 80.1418439716312,
23: 73.7886524822695,
24: 72.34042553191489,
25: 75.177304964539,
26: 73.0496453907092,
27: 73.7886524822695,
28: 73.0496453907092,
29: 74.46809310638297

```

```

In [36]: fontsize=12
plt.style.use('seaborn-whitegrid')
current_kris = plt.gca()
graph = current_kris.plot(best_accuracy_vs_height.keys(), best_accuracy_vs_height,
plt.xlabel('DEPTH OF DECISION TREE', fontsize=fontsize)
plt.ylabel('AUC TEST ACCURACY', fontsize=fontsize)
plt.title('HIGHEST OBSERVED ACCURACY V/S DEPTH OF DT', fontsize=13, fontweight='b'

```

HIGHEST OBSERVED ACCURACY VS DEPTH OF DT

DEPTH OF DECISION TREE	%AGE TEST ACCURACY
0	65.0
1	64.0
2	75.0
3	75.0
4	76.0
5	76.0
6	77.0
7	77.0
8	78.0
9	78.0
10	78.0
11	78.0
12	79.0
13	78.0
14	76.0
15	77.0
16	78.0
17	78.0
18	78.0
19	78.0
20	78.0
21	73.0
22	72.0
23	74.0
24	73.0
25	74.0
26	74.0
27	74.0
28	74.0
29	74.0
30	74.0

```
In [34]: best_accuracy_vs_size = dict()
for tree in all_trees :
    for idx, s in enumerate(tree.training_history['size_progress']) :
        if not s in best_accuracy_vs_size :
            best_accuracy_vs_size[s] = tree.training_history['test_accuracy_progress']
```

```
best_accuracy_vs_size[s] = max(best_accuracy_vs_size[s],
                                tree.training_history['test_accou

In [35]: best_accuracy_vs_size

[1: 65.957446808051064,
 3: 63.829787230404556,
 5: 73.04964539007092,
 7: 73.04964539007092,
 9: 73.04964539007092,
11: 74.46808510638297,
13: 73.04964539007092,
15: 73.04964539007092,
17: 75.177304964539,
19: 75.177304964539,
21: 74.46808510638297,
23: 74.46808510638297,
25: 74.46808510638297,
27: 73.75886524822695,
29: 73.75886524822695,
31: 74.46808510638297,
33: 74.46808510638297,
35: 74.46808510638297,
37: 75.177304964539,
39: 75.177304964539,
```

43: 75.177304964539,
44: 75.177304964539,
47: 74.46808510638297,
49: 74.46808510638297,
50: 74.46808510638297,
53: 74.46808510638297,
55: 75.177304964539,
57: 75.177304964539,
59: 76.59574468085107,
61: 76.59574468085107,
63: 76.59574468085107,
65: 76.59574468085107,
67: 75.177304964539,
69: 75.177304964539,
71: 75.88652482269504,
73: 75.88652482269504,
75: 76.59574468085107,
77: 76.59574468085107,
79: 77.30496453900709,
81: 77.30496453900709,
83: 78.01418439716312,
85: 77.30496453900709,
87: 77.30496453900709,
89: 77.30496453900709,
91: 77.30496453900709,
93: 77.30496453900709,
95: 77.30496453900709,
97: 78.01418439716312,

101: 78.0148439716312,
103: 78.0148439716312,
105: 77.3049643900709,
107: 77.3049643900709,
109: 78.0148439716312,
111: 78.7234025531913,
113: 79.43262411347517,
115: 80.1458439716312,
117: 79.43262411347517,
119: 80.1458439716312,
121: 80.1458439716312,
123: 80.1458439716312,
125: 80.1458439716312,
127: 80.85106382978724,
129: 80.85106382978724,
131: 80.85106382978724,
133: 80.85106382978724,
135: 80.85106382978724,
137: 80.85106382978724,
139: 80.85106382978724,
141: 80.85106382978724,
143: 80.85106382978724,
145: 80.85106382978724,
147: 79.43262411347517,
149: 80.1458439716312,
151: 80.85106382978724,
153: 80.85106382978724,
155: 80.85106382978724,

159: 81.56028368794327,
161: 81.56028368794327,
163: 81.56028368794327,
165: 81.56028368794327,
167: 81.56028368794327,
169: 81.56028368794327,
171: 80.1418439716312,
173: 80.1418439716312,
175: 80.1418439716312,
177: 79.43262413347517,
179: 79.43262413347517,
181: 80.851036382978724,
183: 80.851036382978724,
185: 80.1418439716312,
187: 79.43262413347517,
189: 79.43262413347517,
191: 79.43262413347517,
193: 79.43262413347517,
195: 80.1418439716312,
197: 79.72340425531915,
199: 78.72340425531915,
201: 78.72340425531915,
203: 78.72340425531915,
205: 78.72340425531915,
207: 78.72340425531915,
209: 78.72340425531915,
211: 80.1418439716312,
213: 77.3046643800708

```
In [38]: fontsize=12
plt.style.use('seaborn-whitegrid')
currentAxis = plt.gca()
graph = currentAxis.plot(BEST_ACCURACY_V/S_SIZE.keys(), best_accuracy_v/s_size.values,
    plt.xlabel('SIZE OF EXCITON TRIP', fontsize=fontsize)
    plt.ylabel('AGEE TEST ACCURACY', fontsize=fontsize)
    plt.title('HIGHEST OBSERVED ACCURACY V/S SIZE OF DT', fontsize=13, fontweight='bold')
    plt.savefig('plot_size.png')
plt.savefig('plot_size.png')
```

PART 04

The tree with the highest test accuracy obtained in the second part was pruned using a valid statistical test on the post-order fashion, i.e. the child nodes must be checked for pruning before the parent nodes. Find it (within some predefined bounds) by pruning copies of the original tree with variable thresholds, with a small value at every iteration. Choose the threshold that produces the tree with the maximum test accuracy. Compare the structural difference between the tree produced by pruning the original tree with the most optimal threshold value. Compare the structural difference between the trees before and after pruning.

```
In [39]: test_set = best_tree_test_set
```

```
In [40]:
threshold = 0.1
optimal_threshold = None
highest_pruned_tree_accuracy = -1
while threshold <= 10:
    tree = best_tree.copy()
    tree.prune(threshold = threshold)
    acc = tree.accuracy(test_get)
    if highest_pruned_tree_accuracy < acc:
        highest_pruned_tree_accuracy = acc
        optimal_threshold = threshold
    print('THRESHOLD: ', round(threshold, 1), ' | ACCURACY: ', round(acc, 3))
    threshold += 0.1
```


