

HW4-Programming

● Graded

Student

Brian Bertness

Total Points

30 / 30 pts

Autograder Score

16.0 / 16.0

Passed Tests

test_combined_metric (test_DecisionTree.TestDecisionTree) (4/4)

test_entropy (test_DecisionTree.TestDecisionTree) (4/4)

test_generate_tree (test_DecisionTree.TestDecisionTree) (4/4)

test_gini_index (test_DecisionTree.TestDecisionTree) (4/4)

Question 2

Sanity Check

14 / 14 pts

✓ - 0 pts Correct

- 14 pts Incorrect

Autograder Results

test_combined_metric (test_DecisionTree.TestDecisionTree) (4/4)

The ground truth combined_entropy is: 2.959717

Your combined_entropy is: 2.959717

test_entropy (test_DecisionTree.TestDecisionTree) (4/4)

The ground truth entropy is: 3.311140

Your entropy is: 3.311140

test_generate_tree (test_DecisionTree.TestDecisionTree) (4/4)

For testing, we assume min entropy is 0.1

The testing for generate tree is conducted on the test set rather than training set here

The groundtruth tree node list is:

[42, [30, [44, [33, [None], [19, [2, [None], [4, [None], [5, [None], [None]]]], [None]]], [28, [None], [None]]], [

Your tree node list is:

[42, [30, [44, [33, [None], [19, [2, [None], [4, [None], [5, [None], [None]]]], [None]]], [28, [None], [None]]], [

The tree node list follows the following pattern: [root, [left], [right]]

Please run with Python3 rather than Python2 code.

test_gini_index (test_DecisionTree.TestDecisionTree) (4/4)

The ground truth gini_index is: 0.898545

Your gini_index is: 0.898545

Submitted Files

```
1  import numpy as np
2
3
4  # You are going to implement functions for this file.
5
6  class Tree_node:
7      """
8      Data structure for nodes in the decision-tree
9      """
10     def __init__(self,):
11         self.feature = None # index of the selected feature (for non-leaf node)
12         self.label = None # class label (for leaf node), if not leaf node, label will be None
13         self.left_child = None # left child node
14         self.right_child = None # right child node
15
16
17     class Decision_tree:
18         """
19         Decision tree with binary features
20         """
21         def __init__(self,min_entropy, metric='entropy'):
22             self.metricname = metric
23             self.min_entropy = min_entropy
24             self.root = None
25
26         def fit(self,train_x,train_y):
27             # construct the decision-tree with recursion
28             self.root = self.generate_tree(train_x,train_y)
29
30         def predict(self,test_x):
31             # iterate through all samples
32             prediction = []
33             for i in range(len(test_x)):
34                 cur_data = test_x[i]
35                 # traverse the decision-tree based on the features of the current sample
36                 cur_node = self.root
37                 while True:
38                     if cur_node.label != None:
39                         break
40                     elif cur_node.feature == None:
41                         print("You haven't selected the feature yet")
42                         exit()
43                     else:
44                         if cur_data[cur_node.feature] == 0:
45                             cur_node = cur_node.left_child
46                         else:
47                             cur_node = cur_node.right_child
48                 prediction.append(cur_node.label)
49
```

```

50     prediction = np.array(prediction)
51
52     return prediction
53
54     # ----- You are going to implement this function -----
55     # use recursion to build up the tree. Starting from the root node, you can call itself to determine
56     # what is the left_child and what is the right_child
57     # return: cur_node: the current tree node you create (Type: Tree_Node)
58     def generate_tree(self, data, label):
59         # initialize the current tree node
60         cur_node = Tree_node()
61
62         # compute the node entropy or gini index and determine if the current node is a leaf node
63         # Specifically, if entropy/gini_index (ie, self.metric(label)) < min_entropy,
64         # determine what will be the label (by choosing the label with the largest count) for this leaf
65         node
66         # and directly return the leaf node
67         node_entropy = self.metric(label)
68         if node_entropy < self.min_entropy:
69             # ----- Add your lines here -----
70             cur_node.label = np.argmax( np.bincount(label) )
71
72             return cur_node
73
74         # select the feature that will best split the current non-leaf node
75         # assign the feature index to cur_node.feature
76         # ----- Add your line here -----
77         selected_feature = self.select_feature( data, label )
78         cur_node.feature = selected_feature
79
80         # split the data based on the selected feature
81         # if the selected feature of the data equals to 0, assign the data and corresponding point to
82         left_x, left_y
83         # otherwise assigned to right_x, right_y
84         select_x = data[:, selected_feature]
85         left_x = data[select_x==0,:]
86         left_y = label[select_x==0,]
87         right_x = data[select_x==1,:]
88         right_y = label[select_x==1,]
89
90         # determine cur_node.left_child and cur_node.right_child by call itself, with left_x, left_y and
91         right_x, right_y
92         # ----- Add your line here -----
93         cur_node.left_child = self.generate_tree( left_x, left_y)
94         cur_node.right_child = self.generate_tree( right_x, right_y)
95
96         return cur_node
97
98         # select the feature that maximize the information gains
99         # return: best_feat, which is the index of the feature
100        def select_feature(self, data, label):

```

```

97     best_feat = 0
98     min_entropy = float('inf')
99
100    # iterate through all features and compute their corresponding entropy
101    for i in range(len(data[0])):
102        # split data based on i-th feature
103        split_x = data[:,i]
104        left_y = label[split_x==0,]
105        right_y = label[split_x==1,]
106
107        # compute the combined entropy which weightedly combine the entropy/gini of left_y and
right_y
108        cur_entropy = self.combined_metric(left_y,right_y)
109
110        # select the feature with minimum entropy (set best_feat)
111        if cur_entropy < min_entropy:
112            min_entropy = cur_entropy
113            best_feat = i
114
115    return best_feat
116
117    # ----- You are going to implement this function -----
-----
118    # weightedly combine the entropy/gini of left_y and right_y
119    # the weights are [len(left_y)/(len(left_y)+len(right_y)), len(right_y)/(len(left_y)+len(right_y))]
120    # return: result
121    def combined_metric(self,left_y,right_y):
122        # compute the entropy of a potential split
123        result = 0
124        left = len(left_y)/(len(left_y)+len(right_y))
125        right = len(right_y)/(len(left_y)+len(right_y))
126        result = left * self.metric(left_y) + right * self.metric( right_y )
127
128    return result
129
130    # ----- You are going to implement this function -----
-----
131    # compute entropy/gini_index based on the labels
132    # entropy = -sum_i p_i*log2(p_i+1e-15) (add 1e-15 inside the log when computing the entropy to
prevent numerical issue)
133    # gini_index = 1 - sum_i p_i^2
134    def metric(self, label):
135        result = 0
136
137        if self.metricname == 'entropy':
138            class_label, count = np.unique(label,return_counts=True)
139            count = count/len(label)
140            result = -np.sum( count * np.log2( count + 1e-15 ))
141
142
143        elif self.metricname == 'gini_index':
144            class_label, count = np.unique(label,return_counts=True)

```

```
145     count = count/len(label)
146     result = 1 - np.sum( count**2 )
147
148     return result
149
```