

Lab 2

Decision and Regression Trees

Description

The programs for this lab assignment are designed to implement two types of predictor models for a decision tree: one as a classification model and the other as a regression model. In this paper, I present my proposed modifications to the program for both predictor models. Parts 1-3 of this assignment pertain to the *decision_tree* Python script. So, Parts 4-5 of this assignment are implemented in the *regression_tree* Python script.

Decision Tree Classifier

Part 1

The first problem was to modify the existing code to be a recursive call to the left and right children of each tree node. In this way, we will build a decision tree classifier instead of a single-node tree, that would provide one of two classifications automatically, as was provided in the original code. Therefore, only one small change was made in the return statement (see below).

Original return statement in `_id3(3)`:

```
left = int(round(np.mean(y[less])))
right = int(round(np.mean(y[more])))
return DecisionTreeNode(best_att, thr[best_att], left, right)
```

Modified return statement in `_id3(3)`:

```
lx = x[less]; ly = y[less]
rx = x[more]; ry = y[more]
return DecisionTreeNode(best_att, thr[best_att],
                        self._id3(lx, ly, depth+1), self._id3(rx, ry, depth+1))
```

Part 2

The second task was to generate random values for each attribute to find the best threshold for them. Previously, we used the means of attributes (see `_threshold(3)`). In order to generate these random values, I wrote the following code into a `_threshold3(3)` function to generate random values between the range of the lowest and highest values of each of the attributes in the training data.

```
thr = []
VALS = 10
for i in range(x.shape[1]):
    thr.append(np.random.uniform(min(x[:, i]), max(x[:, i]), size=(VALS)))
```

In order to be used and traversed as a NumPy array, I convert this list before accessing and using it for Part 3.

Part 3

The third (and final task for the *decision tree* classification model) is to find the attribute and threshold combination that yields the highest information gain. Previously (before Part 2), we used the means of attributes. My new Python function `_threshold3(3)` is provided below, including the code for Part 2.

In the nested for loop section of this method is where the program will find the threshold for the attribute which will produce the highest information gain.

```
def _threshold3(self, x, y, orig_entropy):
    # code from Part 2 goes here
    thr = np.asarray(thr)
    entropy_attribute = np.zeros(len(thr))

    thresh = []
    for i in range(x.shape[1]):
        m = sys.maxsize; ind = 0
        for j in range(len(thr[i])):
            less = x[:, i] <= thr[i][j]
            more = ~ less
            new = self._entropy(y[less], y[more])
```

```

        if m > new:
            m = new
            ind = j
        thresh.append(thr[i][ind])
        entropy_attribute[i] = m
    gain = orig_entropy - entropy_attribute
    return np.asarray(thresh), np.argmax(gain)

```

Regression Tree

Part 4

The fourth task was to modify the existing code to be a recursive call to the left and right children of each tree node in the regression model. In this way, we will build a regression tree instead of a single-node tree, as was provided in the original code. Therefore, only one small change was made in the return statement (see below).

Original return statement in `_id3(3)`:

```

left = int(round(np.mean(y[less])))
right = int(round(np.mean(y[more])))
return RegressionTreeNode(best_att, thr[best_att], left, right)

```

Modified return statement in `_id3(3)`:

```

lx = x[less]; ly = y[less]
rx = x[more]; ry = y[more]
return RegressionTreeNode(best_att, thr[best_att],
                          self._id3(lx, ly, depth+1), self._id3(rx, ry, depth+1))

```

Much of this code is mirrored in Part 1 of this lab. The only real difference is to call the *RegressionTreeNode* class rather than the *DecisionTreeNode* class.

Part 5

For the last and final part of this assignment, we were tasked with generating several, random threshold values, and finding the one that would yield the lowest MSE. Therefore, I implemented

a similar algorithm to the one in Part 3 in the `_threshold5(3)` Python function, provided below.

```
def _threshold5(self, x, y, orig_mse):
    thr = []

    VALS = 20
    for i in range(x.shape[1]):
        thr.append(np.random.uniform(min(x[:, i]), max(x[:, i]), size=(VALS)))
    thr = np.asarray(thr)
    mse_attribute = np.zeros(len(thr))

    thresh = []
    for i in range(x.shape[1]):
        m = sys.maxsize; ind = 0
        for j in range(len(thr)):
            less = x[:, i] <= thr[i][j]
            more = ~ less
            new = self._mse(y[less], y[more])
            if new < m:
                m = new; ind = j
        mse_attribute[i] = m
        thresh.append(thr[i][ind])
    gain = orig_mse - mse_attribute
    return np.asarray(thresh), np.argmax(gain)
```

Results

For the classification decision tree, in Parts 1-3, the accuracies are measured by correctly predicted samples. For Parts 4-5, the accuracies are presented as mean squared error (MSE). All times are reported in seconds. For Parts 1-3, I ran the tests over the entire gamma ray detection dataset. For Parts 4-5, I ran the tests over the entire solar particle dataset. For each run in both models, I randomly partitioned a section of the entire dataset to be used for training or testing the model. My results are presented as the average of 100 experiments.

For Part 1, the accuracy of the tree has improved in training and testing, see Table 1; time, however, had stayed about the same.

	Training Time	Training Accuracy	Testing Time	Testing Accuracy
Original	0.038	0.735	0.005	0.74
Modified	0.654	0.881	0.0243	0.842

Table 1: Part 1 Results

For Part 2, there were no experiments run; only runs were to debug the implemented code. Results will be tested for 10 and 20 randomly generated threshold values. This same code is used for Part 5, also, and results are presented for 10, 20, and 50 randomly generated threshold values, also. See Table 2 for results in Part 3 and Table 4 for results in Part 5.

For Part 3, which includes the code from Part 2 (shown here for 10, 20, and 50 values), I did not find a decline in accuracy (as expected). However, I did not find a significant increase in accuracy either (about 0.7% increase from Part 1). Time, however, did increase greatly, which makes sense, since my implemented algorithm runs in $O(n^2)$.

	Training		Testing	
	Time	Accuracy	Time	Accuracy
10 Threshold Values	3.819	0.888	0.09	0.849
20 Threshold Values	7.165	0.892	0.89	0.849
50 Threshold Values	16.098	0.895	0.91	0.849

Table 2: Part 3 Results

For Part 4, as expected, the time increased because we went from a one-node regression tree to a multi-node tree. The accuracy, of course, also increased (lower MSE means higher accuracy). See Table 3 below for results.

	Training Time	Training Accuracy	Testing Time	Testing Accuracy
Original	0.029	0.0869	0.0065	0.0364
Modified	4.387	0.0021	0.0504	0.0033

Table 3: Part 4 Results

For Part 5, some peculiar results are shown in the average time. It is expected for the time to increase when calculating MSE for more threshold values for each attribute. This might be explained by my computer's running extra processes at the time of these experiments. The accuracy, though already low, did increase (as shown in Table 4).

	Training		Testing	
	Time	Accuracy	Time	Accuracy
10 Threshold Values	3.234	0.00272	0.149	0.00259
20 Threshold Values	2.73	0.00253	0.648	0.0029
50 Threshold Values	2.86	0.00251	0.067	0.00285

Table 4: Part 5 Results

Conclusions and Discussion

For Part 1, as expected, the time to train and test increased (since we were building a tree which contained more than the one root node). Further, the accuracy improved by about 10%. Since I ran the tests over the entire gamma ray detection dataset. Each time, I randomly partitioned a section of the entire dataset to be used for training or testing the model to simulate more randomness in my testing and training data.

For Part 3, I expected my accuracy to improve, though perhaps not significantly. As shown in my results, the accuracy of the classification model improved by under 1%.

For Parts 4-5, I was happy with the results of the significant accuracy increase for the solar particle dataset, since it was expected. However, I'd have expected the time to increase with the number of randomly generated threshold values we used for the attributes. My results show that something happened to make the tests with 10 take longer than those with 20 threshold values to check for each attribute. I attribute this abnormality to unstable computer resource management.

Compared to my knn program, the accuracies in my decision tree models are not as good for classification but are better in terms of regression. Regarding time, the use of the decision trees are definitely less time-consuming. This is expected because trees allow us to use branches to come to predictions rather than calculating multiple distances and sorting each of these calculations for each new testing sample. Of course, my implementation of knn was not very optimized (it could have been improved). However, still comparing the run times for both algorithms (knn vs. decision trees), still decision trees win by a lot. This time would also add up when we use larger datasets.

In terms of accuracy, it seems that using either 10, 20 or 50 threshold values did not change accuracy much. There was no change for testing accuracy in the classification model. The training accuracy, however, was increased, but not by much. It might, however, be preferable to use more threshold values, but it would depend on the resources available to the program's environment. For instance, if using the classifier, we can reduce the number of threshold values

to get approximately the same amount of testing accuracy, but we would be gaining time (cutting it significantly). For the regression model, the time stayed relatively the same despite the number of threshold values generated. Similarly, the accuracy (as MSE), did improve, but only in small increments, which decreased as I increased the number of threshold values.

I would like to see how *pruning* would affect the results, especially timewise. If we could prune the tree so that branches can make guesses quicker, then we can not only reduce time but also reduce the workload of the decision tree. In this way, the model won't have learned extra detail and will therefore be more generalizable but still as accurate.

Appendix

The appendix includes the source code of the described program implemented in Python 3.7. Methods/Files authored by me are denoted with an author tag (*@author mmafr*). The code is also available at https://github.com/mahdafr/19w_cs5361-labs on GitHub.

```

1  # Program to build a one-node decision tree
2  # Programmed by Olac Fuentes
3  # Last modified September 16, 2019
4  import sys
5
6  import numpy as np
7  import time
8
9
10 class DecisionTreeNode(object):
11     # Constructor
12     def __init__(self, att, thr, left, right):
13         self.attribute = att
14         self.threshold = thr
15         # left and right are either binary classifications or references to
16         # decision tree nodes
17         self.left = left
18         self.right = right
19
20     def print_tree(self, indent=''):
21         # If prints the right subtree, corresponding to the condition x[attribute] >
22         # threshold
23         # above the condition stored in the node
24         if self.right in [0, 1]:
25             print(indent + ' ', 'class=', self.right)
26         else:
27             self.right.print_tree(indent + ' ')
28
29         print(indent, 'if x[' + str(self.attribute) + '] <=', self.threshold)
30
31         if self.left in [0, 1]:
32             print(indent + ' ', 'class=', self.left)
33         else:
34             self.left.print_tree(indent + ' ')
35
36 class DecisionTreeClassifier(object):
37     # Constructor
38     def __init__(self, max_depth=10, min_samples_split=10, min_accuracy=1):
39         self.max_depth = max_depth
40         self.min_samples_split = min_samples_split
41         self.min_accuracy = min_accuracy
42
43     def fit(self, x, y):
44         self.root = self._id3(x, y, depth=0)
45
46     def predict(self, x_test):
47         pred = np.zeros(len(x_test), dtype=int)
48         for i in range(len(x_test)):
49             pred[i] = self._classify(self.root, x_test[i])
50         return pred
51
52     def _id3(self, x, y, depth):
53         orig_entropy = self._entropy(y, [])
54         mean_val = np.mean(y)
55
56         # if accuracy not attained and cannot go further in tree
57         if depth >= self.max_depth or len(y) <= self.min_samples_split or max(
58             [mean_val, 1 - mean_val]) >= self.min_accuracy:
59             return int(round(mean_val))
60
61         # @author mahdafr for part2-part3
62         thr, best_att = self._threshold3(x, y, orig_entropy)
63
64         less = x[:, best_att] <= thr[best_att]
65         more = ~ less

```



```

66     # @author mahdafr for part1
67     lx = x[less]; ly = y[less] # int(round(np.mean(y[less])))
68     rx = x[more]; ry = y[more] # int(round(np.mean(y[more])))
69     return DecisionTreeNode(best_att, thr[best_att], self._id3(lx,ly,depth+1),
70                             self._id3(rx,ry,depth+1))
71
72 # original code
73 def _threshold(self, x, y, orig_entropy):
74     thr = np.mean(x, axis=0)
75     entropy_attribute = np.zeros(len(thr))
76
77     # foreach training example, find entropy for each attribute
78     for i in range(x.shape[1]):
79         less = x[:, i] <= thr[i]
80         more = ~ less
81         entropy_attribute[i] = self._entropy(y[less], y[more])
82     gain = orig_entropy - entropy_attribute
83     # print('Gain:',gain)
84     return thr, np.argmax(gain)
85
86 # @author mahdafr for part3
87 def _threshold3(self, x, y, orig_entropy):
88     thr = [] # np.mean(x, axis=1)
89
90     # @author mahdafr for part2
91     # generate random values for each attribute
92     VALS = 20
93     for i in range(x.shape[1]):
94         thr.append(np.random.uniform(min(x[:,i]),max(x[:,i]),size=(VALS)))
95     thr = np.asarray(thr)
96     entropy_attribute = np.zeros(len(thr))
97
98     thresh = []
99     # find entropy
100     for i in range(x.shape[1]):
101         m = sys.maxsize; ind = 0
102         for j in range(len(thr[i])):
103             less = x[:, i] <= thr[i][j]
104             more = ~ less
105             new = self._entropy(y[less], y[more])
106             if new < m:
107                 m = new
108                 ind = j
109             thresh.append(thr[i][ind])
110         entropy_attribute[i] = m
111     gain = orig_entropy - entropy_attribute
112     # print('Gain:',gain)
113     return np.asarray(thresh), np.argmax(gain)
114
115 def _entropy(self, l, m):
116     ent = 0
117     for p in [l, m]:
118         if len(p) > 0:
119             pp = sum(p) / len(p)
120             pn = 1 - pp
121             if pp < 1 and pp > 0:
122                 ent -= len(p) * (pp * np.log2(pp) + pn * np.log2(pn))
123     ent = ent / (len(l) + len(m))
124     return ent
125
126 def _classify(self, dt_node, x):
127     if dt_node in [0, 1]:
128         return dt_node
129     if x[dt_node.attribute] <= dt_node.threshold:
130         return self._classify(dt_node.left, x)
131     else:

```

```

131         return self._classify(dt_node.right, x)
132
133     def display(self):
134         print('Model:')
135         self.root.print_tree()
136
137
138 x = []; y = []
139 infile = open("magic04.txt", "r")
140 for line in infile:
141     y.append(int(line[-2:-1] == 'g'))
142     x.append(np.fromstring(line[:-2], dtype=float, sep=','))
143 infile.close()
144
145 xa = np.zeros((len(x), len(x[0])))
146 for i in range(len(xa)):
147     xa[i] = x[i]
148 x = xa
149
150 x = np.array(x).astype(np.float32)
151 y = np.array(y)
152
153 TESTS = 100
154 test_acc = 0; train_acc = 0
155 test_time = 0; train_time = 0
156 print('Tests ' + str(TESTS))
157 for i in range(TESTS):
158     # Split data into training and testing
159     ind = np.random.permutation(len(y))
160     split_ind = int(len(y) * 0.8)
161     x_train = x[ind[:split_ind]]
162     x_test = x[ind[split_ind:]]
163     y_train = y[ind[:split_ind]]
164     y_test = y[ind[split_ind:]]
165
166     model = DecisionTreeClassifier()
167     start = time.time()
168     model.fit(x_train, y_train)
169     train_time += time.time() - start
170
171     train_pred = model.predict(x_train)
172     start = time.time()
173     test_pred = model.predict(x_test)
174     test_time += time.time() - start
175
176     train_acc += np.sum(train_pred == y_train) / len(train_pred)
177     test_acc += np.sum(test_pred == y_test) / len(test_pred)
178     # model.display()
179
180 print('Elapsed_time training {0:.6f} '.format(train_time/TESTS))
181 print('Elapsed_time testing {0:.6f} '.format(test_time/TESTS))
182 print('train accuracy:', train_acc/TESTS)
183 print('test accuracy:', test_acc/TESTS)
184

```

```

1  # Program to build a one-node regression tree
2  # Programmed by Olac Fuentes
3  # Last modified September 16, 2019
4  import sys
5
6  import numpy as np
7  import time
8
9
10 class RegressionTreeNode(object):
11     # Constructor
12     def __init__(self, att, thr, left, right):
13         self.attribute = att
14         self.threshold = thr
15         # left and right are either binary classifications or references to
16         # decision tree nodes
17         self.left = left
18         self.right = right
19
20     def print_tree(self, indent=''):
21         # If prints the right subtree, corresponding to the condition x[attribute] >
22         # threshold
23         # above the condition stored in the node
24         if isinstance(self.right, np.float32):
25             print(indent + ' ', 'pred=', self.right)
26         else:
27             self.right.print_tree(indent + ' ')
28
29         print(indent, 'if x[' + str(self.attribute) + '] <=', self.threshold)
30
31         if isinstance(self.left, np.float32):
32             print(indent + ' ', 'pred=', self.left)
33         else:
34             self.left.print_tree(indent + ' ')
35
36 class DecisionTreeRegressor(object):
37     # Constructor
38     def __init__(self, max_depth=10, min_samples_split=5, max_mse=0.001):
39         self.max_depth = max_depth
40         self.min_samples_split = min_samples_split
41         self.max_mse = max_mse
42
43     def fit(self, x, y):
44         self.root = self._id3(x, y, depth=0)
45
46     def predict(self, x_test):
47         pred = np.zeros(len(x_test), dtype=np.float32)
48         for i in range(len(x_test)):
49             pred[i] = self._predict(self.root, x_test[i])
50         return pred
51
52     def _id3(self, x, y, depth):
53         orig_mse = np.var(y)
54         # print('original mse:', orig_mse)
55         mean_val = np.mean(y)
56         if depth >= self.max_depth or len(y) <= self.min_samples_split or orig_mse <=
57             self.max_mse:
58             return mean_val
59
60         # @author mahdafr part5
61         thr, best_att = self._threshold(x, y, orig_mse)
62
63         # print('mse best attribute:', mse_attribute[best_att])
64         less = x[:, best_att] <= thr[best_att]
65         more = ~ less

```

```

65     # print('subtree mse:', np.var(y[less]), np.var(y[more]))
66     # @author mahdafr for part4
67     lx = x[less]; ly = y[less]
68     rx = x[more]; ry = y[more]
69     return RegressionTreeNode(best_att, thr[best_att], self._id3(lx, ly, depth +
70     1), self._id3(rx, ry, depth + 1))
71     # return RegressionTreeNode(best_att, thr[best_att], np.mean(y[less]),
72     np.mean(y[more]))
73
74 # original code
75 def _threshold(self, x, y, orig_mse):
76     thr = np.mean(x, axis=0)
77     mse_attribute = np.zeros(len(thr))
78     for i in range(x.shape[1]):
79         less = x[:, i] <= thr[i]
80         more = ~ less
81         mse_attribute[i] = self._mse(y[less], y[more])
82     gain = orig_mse - mse_attribute
83     # print('Gain:', gain)
84     return thr, np.argmax(gain)
85
86 # @author mahdafr for part5
87 def _threshold5(self, x, y, orig_mse):
88     thr = [] # np.mean(x, axis=1)
89
90     # @author mahdafr for part2
91     # generate random values for each attribute
92     VALS = 20
93     for i in range(x.shape[1]):
94         thr.append(np.random.uniform(min(x[:, i]), max(x[:, i]), size=(VALS)))
95     thr = np.asarray(thr)
96     mse_attribute = np.zeros(len(thr))
97
98     thresh = []
99     for i in range(x.shape[1]):
100         m = sys.maxsize; ind = 0
101         for j in range(len(thr)):
102             less = x[:, i] <= thr[i][j]
103             more = ~ less
104             new = self._mse(y[less], y[more])
105             if new < m:
106                 m = new; ind = j
107         mse_attribute[i] = m
108         thresh.append(thr[i][ind])
109     gain = orig_mse - mse_attribute
110     # print('Gain:', gain)
111     return np.asarray(thresh), np.argmax(gain)
112
113 def _mse(self, l, m):
114     err = np.append(l - np.mean(l), m - np.mean(m)) # It will issue a warning if
115     either l or m is empty
116     return np.mean(err * err)
117
118 def _predict(self, dt_node, x):
119     if isinstance(dt_node, np.float32):
120         return dt_node
121     if x[dt_node.attribute] <= dt_node.threshold:
122         return self._predict(dt_node.left, x)
123     else:
124         return self._predict(dt_node.right, x)
125
126 def display(self):
127     print('Model:')
128     self.root.print_tree()

```

```

128 TESTS = 100
129 test_acc = 0; train_acc = 0
130 test_time = 0; train_time = 0
131 print('Tests ' + str(TESTS))
132 dir = 'D:\Google Drive\skool\CS 5361\datasets\lab1\'
133 for i in range(TESTS):
134     skip = np.random.randint(40,50)
135     x_train = np.load(dir + 'x_ray_data_train.npy')[::skip]
136     y_train = np.load(dir + 'x_ray_target_train.npy')[::skip]
137     x_test = np.load(dir + 'x_ray_data_test.npy')[::skip]
138     y_test = np.load(dir + 'x_ray_target_test.npy')[::skip]
139
140     model = DecisionTreeRegressor()
141     start = time.time()
142     model.fit(x_train, y_train)
143     train_time += time.time() - start
144     pred = model.predict(x_train)
145     train_acc += np.mean(np.square(pred - y_train))
146
147     start = time.time()
148     pred = model.predict(x_test)
149     test_time += time.time() - start
150     test_acc += np.mean(np.square(pred - y_test))
151     # model.display()
152
153 print('Elapsed_time training {0:.6f} '.format(train_time/TESTS))
154 print('Elapsed_time testing {0:.6f} '.format(test_time/TESTS))
155 print('Mean square error training set:', train_acc/TESTS)
156 print('Mean square error test set:', test_acc/TESTS)
157
158

```