

# 1003 HW6

Long Chen  
lc3424@nyu.edu

April 18th, 2021

## Q1

```
1 def compute_entropy(label_array):
2     entropy = 0
3     for k in np.unique(label_array):
4         pmk = (sum(np.equal(k, label_array).astype(int)) /
5               label_array.shape[0])[0]
6         entropy += pmk * np.log(pmk)
7
8     return -entropy
9
10 def compute_gini(label_array):
11     gini = 0
12     for k in np.unique(label_array):
13         pmk = (sum(np.equal(k, label_array).astype(int)) /
14               label_array.shape[0])[0]
15         gini += pmk * (1 - pmk)
16
17     return gini
```

## Q2

```
1 class Decision_Tree(BaseEstimator):
2     def __init__(self, split_loss_function, leaf_value_estimator,
3                 depth=0, min_sample=5, max_depth=10, force_split=
4                 True):
5         self.split_loss_function = split_loss_function
6         self.leaf_value_estimator = leaf_value_estimator
7         self.depth = depth
8         self.min_sample = min_sample
9         self.max_depth = max_depth
10        self.is_leaf = False
11        self.force_split = force_split
12
13        '''
14        Author note: in my algorithm, I forces the node to split (
15        except depth or min_sample criteria is met) even if the split
16        will result in higher loss. A hyperparameter is added to the
17        constructor.
```

```

14     '''
15     def fit(self, x, y):
16         # Corner case: reaches max_depth OR does not meet
17         min_sample to split.
18         if self.depth == self.max_depth or len(y) <= self.
19         min_sample:
20             self.is_leaf = True
21             self.value = self.leaf_value_estimator(y)
22             return self
23
24         self.split_id, self.split_value, split_pos, best_loss =
25         self.find_best_feature_split(x, y)
26
27         if self.force_split or best_loss < self.split_loss_function
28         (y): # split
29             # find left and right x, y
30             lx, ly, rx, ry = None, None, None, None
31             D = np.concatenate([x, y], 1) # concatenated dataset for
32             easier sorting
33             D = np.array(sorted(D, key=lambda x: x[self.split_id]))
34             lx, ly = D[:split_pos+1, :-1], D[:split_pos+1, -1].
35             reshape(-1, 1)
36             rx, ry = D[split_pos+1:, :-1], D[split_pos+1:, -1].
37             reshape(-1, 1)
38
39             # initialize left and right node
40             self.left = Decision_Tree(self.split_loss_function,
41             self.leaf_value_estimator,
42             depth=self.depth+1, min_sample=self.min_sample
43             , max_depth=self.max_depth)
44             self.right = Decision_Tree(self.split_loss_function,
45             self.leaf_value_estimator,
46             depth=self.depth+1, min_sample=self.min_sample
47             , max_depth=self.max_depth)
48             self.left.fit(lx, ly)
49             self.right.fit(rx, ry)
50         else:
51             self.is_leaf = True
52             self.value = self.leaf_value_estimator(y)
53             return self
54
55         return self
56
57     def find_best_split(self, x_node, y_node, feature_id):
58         split_value, best_loss, best_pos = None, float('inf'), -1
59         D = np.concatenate([x_node[:, feature_id].reshape(-1, 1),
60         y_node], 1) # concatenated dataset for easier sorting
61         D = np.array(sorted(D, key=lambda x: x[0]))
62         # iterate through all datapoint intervals
63         for pos in range(len(D)-1):
64             lx, ly = D[:pos+1, 0], D[:pos+1, 1].reshape(-1, 1)
65             rx, ry = D[pos+1:, 0], D[pos+1:, 1].reshape(-1, 1)
66
67             ll = len(ly)*self.split_loss_function(ly)/len(y_node)
68             rl = len(ry)*self.split_loss_function(ry)/len(y_node)
69
70             # split condition

```

```

59         if ll + rl < best_loss:
60             split_value = (lx[-1] + rx[0]) / 2
61             best_loss = ll + rl
62             best_pos = pos
63
64         return split_value, best_loss, best_pos
65
66     def find_best_feature_split(self, x_node, y_node):
67         split_id, split_value, best_loss, best_pos = None, None,
68         float('inf'), -1
69
70         for feature_id in range(x_node.shape[1]):
71             sv, best_l, pos = self.find_best_split(x_node, y_node,
72             feature_id)
73             if best_l < best_loss:
74                 split_id, split_value, best_pos = feature_id, sv,
75                 pos
76                 best_loss = best_l
77
78         return split_id, split_value, best_pos, best_loss
79
80     def predict_instance(self, instance):
81         if self.is_leaf:
82             return self.value
83         if instance[self.split_id] <= self.split_value:
84             return self.left.predict_instance(instance)
85         else:
86             return self.right.predict_instance(instance)

```

### Q3

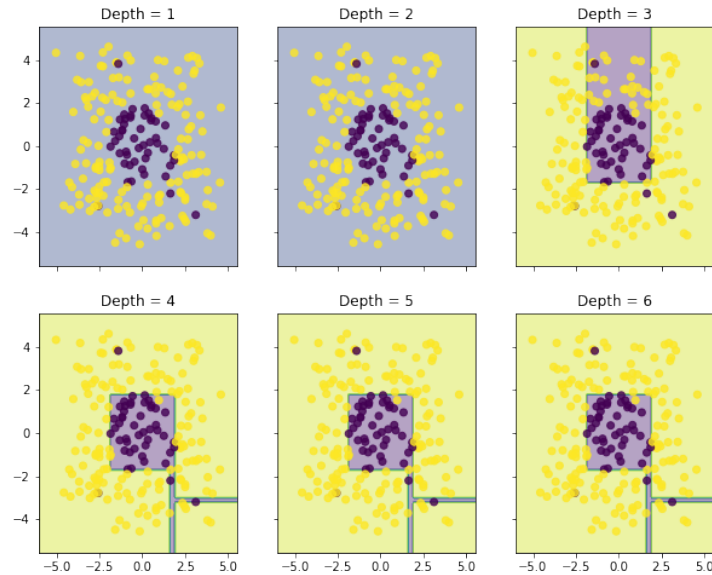


Figure 1: Q3 result.

## Q4

```

1 def mean_absolute_deviation_around_median(y):
2     mae, m = 0, np.mean(y)
3     for yi in y:
4         mae += np.abs(yi - m)[0]
5
6     return mae / len(y)

```

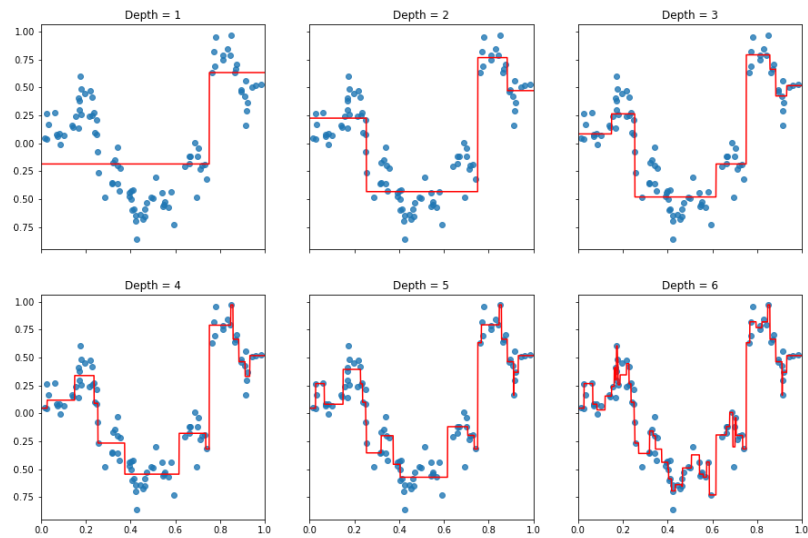


Figure 2: Q4 result.

## Q5

```

1 class gradient_boosting():
2     def __init__(self, n_estimator, pseudo_residual_func,
3                   learning_rate=0.01,
4                   min_sample=5, max_depth=5):
5
6         self.n_estimator = n_estimator
7         self.pseudo_residual_func = pseudo_residual_func
8         self.learning_rate = learning_rate
9         self.min_sample = min_sample
10        self.max_depth = max_depth
11        self.estimators = [] #will collect the n_estimator models
12
13    def fit(self, train_data, train_target):
14        # do f0
15        res = self.pseudo_residual_func(train_target.reshape(-1),
16                                         np.zeros(train_target.shape).reshape(-1))
17        h0 = DecisionTreeRegressor(max_depth=self.max_depth,
18                                   min_samples_leaf=self.min_sample, criterion='mse')

```

```

16     h0.fit(train_data, res)
17     self.estimatedors.append(h0)
18
19     # do f[1, M]
20     for m in range(self.n_estimator):
21         step = 0
22         for i in range(len(self.estimatedors)):
23             step += self.learning_rate * self.estimatedors[i].
24             predict(train_data)
25
26         res = self.pseudo_residual_func(train_target.reshape
27         (-1), step)
28         hm = DecisionTreeRegressor(max_depth=self.max_depth,
29         min_samples_leaf=self.min_sample, criterion='mse')
30         hm.fit(train_data, res)
31         self.estimatedors.append(hm)
32
33     def predict(self, test_data):
34         test_predict = self.learning_rate * sum([estimator.predict(
35         test_data) for estimator in self.estimatedors])
36
37         return test_predict

```

Q6

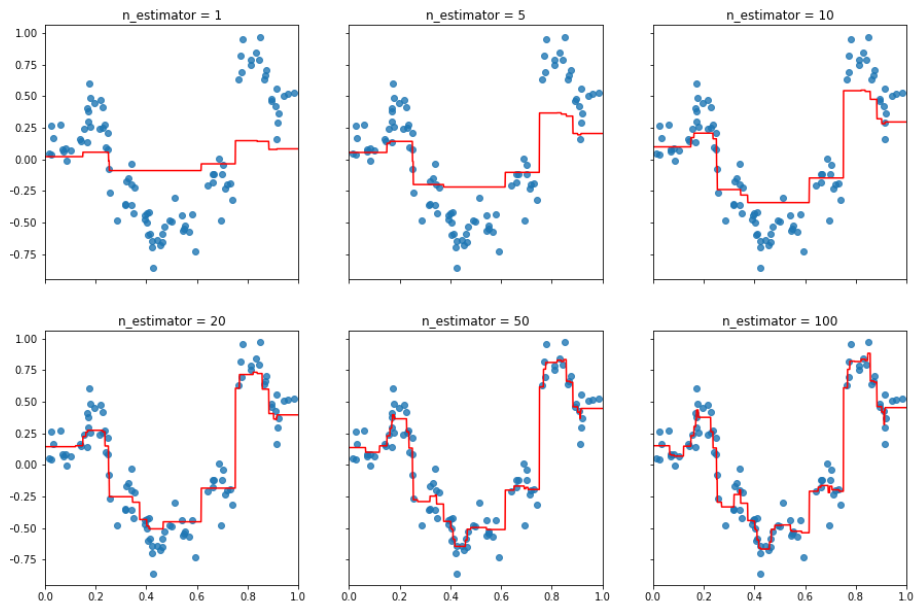


Figure 3: Q6 result.

**Q7**

$$\begin{aligned} -g_m &= - \left( \frac{\delta}{\delta f_{m-1}(x_j)} \sum_{i=1}^m \ell(y_i, f_{m-1}(x_i)) \right)_{j=1}^n \\ &= - \left( \frac{\delta}{\delta f_{m-1}(x_j)} \sum_{i=1}^m \ln(1 + \exp(-y_i f_{m-1}(x_i))) \right)_{j=1}^n \end{aligned}$$

Dimension of  $g_m$  is  $n$ .

**Q8**

$$h_m = \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n ((-g_m)_i - h(x_i))^2$$

**Q9**

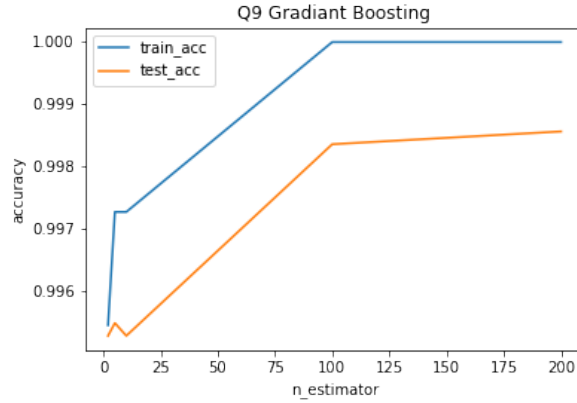


Figure 4: Q9 result.

**Q10 (Optional)**

Random forest is an ensemble method which uses bagged decision trees with modified tree-growing procedure to reduce the dependence between trees (a randomly chosen subset of features for splitting criteria). This will avoid dominance by strong features across the trees and will thus reduces variance (which is often high for a single decision tree).

## Q11 (Optional)

```
1 for n_estimator in [2, 5, 10, 100, 200]:
2     clf = RandomForestClassifier(criterion='entropy', max_depth=3,
3     n_estimators=n_estimator)
4     clf.fit(X_train, y_train)
5     print('n_estimator: {:3}, Train accuracy: {:.15}, test accuracy:
6     {:.15}'.format(n_estimator, accuracy_score(clf.predict(X_train),
7     y_train), accuracy_score(clf.predict(X_test), y_test)))
```

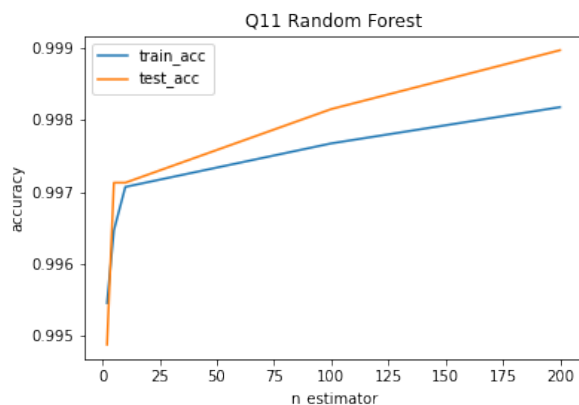


Figure 5: Q11 result.

## Q12 (Optional)

We do observe less overfitting on the random forest results in comparison with gradient boosted trees. The best training accuracy is achieved by gradient boosted tree (acc=1.0). This is as expected, since random forest tends to reduce variance and thus leading to less overfitting.

If I am understanding the question correctly, the “best performing method” refers to the model with best train accuracy. A practical disadvantage is that if we have a new data point that is very different to any training samples, the prediction may be extremely bad.

In terms of test accuracy, random forest is a slightly ahead (0.9990 vs. 0.9986).