

# Random Forest Algorithm

Using a Toy Dataset (Fruits)



## Intuition

The random forest algorithm is a type of supervised learning which can be used for solving both classification and regression problems. In this example we will solve classification problems. Therefore, our algorithm is called a categorical variable random forest. The 'forest' is an ensemble of decision trees, usually trained with the 'bagging' method. The general idea of the bagging method is that a combination of learning models increases the overall result.

## Steps

- 1.) Split the data into training and testing sets
- 2.) Build a decision Tree
- 3.) Build forest by repeating steps 1 and 2 for 'n' number of times to create 'n' number of trees.
- 4.) Evaluate

## Set up (Libraries Used)

```
In [33]: 1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from math import sqrt
4 from random import randrange
```

## Step 1.) Create the Toy Dataset

```
In [34]: 1 # Training Dataset
2 dataset = [
3     ['Green', 3, 'Apple'],
4     ['Yellow', 3, 'Apple'],
5     ['Red', 1, 'Grape'],
6     ['Red', 1, 'Grape'],
7     ['Yellow', 3, 'Lemon'],
8     ['Green', 3, 'Apple'],
9     ['Yellow', 4, 'Apple'],
10    ['Red', 2, 'Grape'],
11    ['Red', 1, 'Grape'],
12    ['Yellow', 3, 'Lemon']
13 ]
14
15 # Header names
16 header = ["color", "diameter", "label"]
```

## Step 2.) Decision Tree

Define functions that will help find the best partition. To find the best partition means to find the best question to split the data.

In [35]:

```
1  # Convert string column to float
2  def str_column_to_float(dataset, column):
3      for row in dataset:
4          row[column] = float(row[column].strip())
5
6  # Convert string column to integer
7  def str_column_to_int(dataset, column):
8      class_values = [row[column] for row in dataset]
9      unique = set(class_values)
10     lookup = dict()
11     for i, value in enumerate(unique):
12         lookup[value] = i
13     for row in dataset:
14         row[column] = lookup[row[column]]
15     return lookup
16
17 # Split a dataset into k folds
18 def cross_validation_split(dataset, n_folds):
19     dataset_split = list()
20     dataset_copy = list(dataset)
21     fold_size = int(len(dataset) / n_folds)
22     for i in range(n_folds):
23         fold = list()
24         while len(fold) < fold_size:
25             index = randrange(len(dataset_copy))
26             fold.append(dataset_copy.pop(index))
27         dataset_split.append(fold)
28     return dataset_split
29
30 # Calculate accuracy percentage
31 def accuracy_metric(actual, predicted):
32     correct = 0
33     for i in range(len(actual)):
34         if actual[i] == predicted[i]:
35             correct += 1
36     return correct / float(len(actual)) * 100.0
37
38 # Evaluate an algorithm using a cross validation split
39 def evaluate_algorithm(dataset, algorithm, n_folds, *args):
40     folds = cross_validation_split(dataset, n_folds)
41     scores = list()
42     for fold in folds:
43         train_set = list(folds)
44         train_set.remove(fold)
45         train_set = sum(train_set, [])
46         test_set = list()
47         for row in fold:
48             row_copy = list(row)
49             test_set.append(row_copy)
50             row_copy[-1] = None
51         predicted = algorithm(train_set, test_set, *args)
52         actual = [row[-1] for row in fold]
53         accuracy = accuracy_metric(actual, predicted)
54         scores.append(accuracy)
55     return scores
56
```

```

57 # Split a dataset based on an attribute and an attribute value
58 def test_split(index, value, dataset):
59     left, right = list(), list()
60     for row in dataset:
61         if row[index] < value:
62             left.append(row)
63         else:
64             right.append(row)
65     return left, right
66
67 # Calculate the Gini index for a split dataset
68 def gini_index(groups, classes):
69     # count all samples at split point
70     n_instances = float(sum([len(group) for group in groups]))
71     # sum weighted Gini index for each group
72     gini = 0.0
73     for group in groups:
74         size = float(len(group))
75         # avoid divide by zero
76         if size == 0:
77             continue
78         score = 0.0
79         # score the group based on the score for each class
80         for class_val in classes:
81             p = [row[-1] for row in group].count(class_val) / size
82             score += p * p
83         # weight the group score by its relative size
84         gini += (1.0 - score) * (size / n_instances)
85     return gini
86
87 # Select the best split point for a dataset
88 def get_split(dataset, n_features):
89     class_values = list(set(row[-1] for row in dataset))
90     b_index, b_value, b_score, b_groups = 999, 999, 999, None
91     features = list()
92     while len(features) < n_features:
93         index = randrange(len(dataset[0])-1)
94         if index not in features:
95             features.append(index)
96     for index in features:
97         for row in dataset:
98             groups = test_split(index, row[index], dataset)
99             gini = gini_index(groups, class_values)
100             if gini < b_score:
101                 b_index, b_value, b_score, b_groups = index, row[index], gini, groups
102     return {'index':b_index, 'value':b_value, 'groups':b_groups}
103
104 # Create a terminal node value
105 def to_terminal(group):
106     outcomes = [row[-1] for row in group]
107     return max(set(outcomes), key=outcomes.count)
108
109 # Create child splits for a node or make terminal
110 def split(node, max_depth, min_size, n_features, depth):
111     left, right = node['groups']
112     del(node['groups'])
113     # check for a no split

```

```

114     if not left or not right:
115         node['left'] = node['right'] = to_terminal(left + right)
116         return
117     # check for max depth
118     if depth >= max_depth:
119         node['left'], node['right'] = to_terminal(left), to_terminal(r
120         return
121     # process left child
122     if len(left) <= min_size:
123         node['left'] = to_terminal(left)
124     else:
125         node['left'] = get_split(left, n_features)
126         split(node['left'], max_depth, min_size, n_features, depth+1)
127     # process right child
128     if len(right) <= min_size:
129         node['right'] = to_terminal(right)
130     else:
131         node['right'] = get_split(right, n_features)
132         split(node['right'], max_depth, min_size, n_features, depth+1)
133
134     # Build a decision tree
135     def build_tree(train, max_depth, min_size, n_features):
136         root = get_split(train, n_features)
137         split(root, max_depth, min_size, n_features, 1)
138         return root
139
140
141     # Make a prediction with a decision tree
142     def predict(node, row):
143         if row[node['index']] < node['value']:
144             if isinstance(node['left'], dict):
145                 return predict(node['left'], row)
146             else:
147                 return node['left']
148         else:
149             if isinstance(node['right'], dict):
150                 return predict(node['right'], row)
151             else:
152                 return node['right']

```

### Step 3.) Random Forest

```

In [36]: 1 # Create a random subsample from the dataset with replacement
2 def subsample(dataset, ratio):
3     sample = list()
4     n_sample = round(len(dataset) * ratio)
5     while len(sample) < n_sample:
6         index = randrange(len(dataset))
7         sample.append(dataset[index])
8     return sample
9
10 # Make a prediction with a list of bagged trees
11 def bagging_predict(trees, row):
12     predictions = [predict(tree, row) for tree in trees]
13     return max(set(predictions), key=predictions.count)
14
15 # Random Forest Algorithm
16 def random_forest(train, test, max_depth, min_size, sample_size, n_trees):
17     trees = list()
18     for i in range(n_trees):
19         sample = subsample(train, sample_size)
20         tree = build_tree(sample, max_depth, min_size, n_features)
21         trees.append(tree)
22     predictions = [bagging_predict(trees, row) for row in test]
23     return predictions

```

## Evaluate

```

In [39]: 1 # evaluate algorithm
2 n_folds = 5
3 max_depth = 10
4 min_size = 1
5 sample_size = 1.0
6 n_features = int(sqrt(len(dataset[0])-1))
7 for n_trees in [1, 5, 10]:
8     scores = evaluate_algorithm(dataset, random_forest, n_folds, max_depth, min_size, sample_size, n_trees)
9     print('Trees: %d' % n_trees)
10    print('Scores: %s' % scores)
11    print('Mean Accuracy: %.3f%%' % (sum(scores)/float(len(scores))))

```

```

Trees: 1
Scores: [50.0, 100.0, 50.0, 100.0, 50.0]
Mean Accuracy: 70.000%
Trees: 5
Scores: [100.0, 50.0, 100.0, 100.0, 50.0]
Mean Accuracy: 80.000%
Trees: 10
Scores: [50.0, 100.0, 50.0, 100.0, 0.0]
Mean Accuracy: 60.000%

```

In [ ]:

1

