

Decision Tree Algorithm

Using a Toy Dataset (Fruits)



Intuition

The decision tree 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 decision tree. Essentially, the algorithm predicts a class label for a record from the root of the 'tree'. It compares the values of the root attribute with the record's attribute. On the basis of comparison, we follow the branch corresponding to that value and jump to the next node.

Steps

- 1.) Split the data into training and testing sets
- 2.) Find the most relevant question
- 3.) Add/Build a decision tree
- 4.) Go to step 1 until arrived at the answer

Set up (Libraries Used)

```
In [308]: 1 import pandas as pd
          2 from sklearn.model_selection import train_test_split
```

Step 1.) Create the Toy Training and Testing Dataset

```
In [286]: 1  # Training Dataset
2  training = [
3      ['Green', 3, 'Apple'],
4      ['Yellow', 3, 'Apple'],
5      ['Red', 1, 'Grape'],
6      ['Red', 1, 'Grape'],
7      ['Yellow', 3, 'Lemon'],
8  ]
9
10 # Testing Dataset
11 testing = [
12     ['Green', 3, 'Apple'],
13     ['Yellow', 4, 'Apple'],
14     ['Red', 2, 'Grape'],
15     ['Red', 1, 'Grape'],
16     ['Yellow', 3, 'Lemon'],
17 ]
18
19 # Header names
20 header = ["color", "diameter", "label"]
```

Step 2.) Find the Best Question and Partition

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 [316]:

```
1  # Counts the number of each type of example in a dataset
2  def class_counts(rows):
3      counts = {}
4      for row in rows:
5          label = row[-1]
6          if label not in counts:
7              counts[label] = 0
8          counts[label] += 1
9      return counts
10
11  # Test if a value is numeric
12  def is_numeric(value):
13      return isinstance(value, int) or isinstance(value, float)
14
15  # Define the question
16  class Question:
17      def __init__(self, column, value):
18          self.column = column
19          self.value = value
20
21      def match(self, example):
22
23          val = example[self.column]
24          if is_numeric(val):
25              return val >= self.value
26          else:
27              return val == self.value
28
29      def __repr__(self):
30          condition = "=="
31          if is_numeric(self.value):
32              condition = ">="
33          return "Is %s %s %s?" % (
34              header[self.column], condition, str(self.value))
35
36  # Split between True rows and False rows based on the question
37  def partition(rows, question):
38      true_rows, false_rows = [], []
39      for row in rows:
40          if question.match(row):
41              true_rows.append(row)
42          else:
43              false_rows.append(row)
44      return true_rows, false_rows
45
46  # Calculate the Gini Impurity
47  def gini(rows):
48      counts = class_counts(rows)
49      impurity = 1
50      for lbl in counts:
51          prob_of_lbl = counts[lbl] / float(len(rows))
52          impurity -= prob_of_lbl**2
53      return impurity
54
55  # Calculate the info gained
56  def info_gain(left, right, current_uncertainty):
```

```

57     p = float(len(left)) / (len(left) + len(right))
58     return current_uncertainty - p * gini(left) - (1 - p) * gini(right)
59
60

```

Determine Current Uncertainty

```

In [287]: 1 current_uncertainty = gini(training)
          2 current_uncertainty

```

```

Out[287]: 0.6399999999999999

```

Partition the data

```

In [317]: 1 def find_best_split(rows):
          2     """Find the best question to ask by iterating over every feature /
          3     and calculating the information gain."""
          4     best_gain = 0 # keep track of the best information gain
          5     best_question = None # keep train of the feature / value that produced
          6     current_uncertainty = gini(rows)
          7     n_features = len(rows[0]) - 1 # number of columns
          8
          9     for col in range(n_features): # for each feature
         10
         11         values = set([row[col] for row in rows]) # unique values in the column
         12
         13         for val in values: # for each value
         14
         15             question = Question(col, val)
         16
         17             # try splitting the dataset
         18             true_rows, false_rows = partition(rows, question)
         19
         20             # Skip this split if it doesn't divide the
         21             # dataset.
         22             if len(true_rows) == 0 or len(false_rows) == 0:
         23                 continue
         24
         25             # Calculate the information gain from this split
         26             gain = info_gain(true_rows, false_rows, current_uncertainty)
         27
         28             # You actually can use '>' instead of '>=' here
         29             # but I wanted the tree to look a certain way for our
         30             # toy dataset.
         31             if gain >= best_gain:
         32                 best_gain, best_question = gain, question
         33
         34     return best_gain, best_question

```

Results -

```
In [321]: 1 best_gain, best_question = find_best_split(training_data)
2 print("Best Question : ",best_question)
3 print("The gain from the best question : ", best_gain )
```

Best Question : Is diameter >= 3?

The gain from the best question : 0.37333333333333324

Step 3.) Build a Decision Tree

```
In [325]: 1 # LEaf node classifies the data
2 class Leaf:
3     def __init__(self, rows):
4         self.predictions = class_counts(rows)
5
6     # Asks a question and holds references to child nodes
7 class Decision_Node:
8     def __init__(self,
9                 question,
10                true_branch,
11                false_branch):
12         self.question = question
13         self.true_branch = true_branch
14         self.false_branch = false_branch
15
16 # Builds the tree
17 def build_tree(rows):
18     gain, question = find_best_split(rows)
19
20     if gain == 0:
21         return Leaf(rows)
22
23     true_rows, false_rows = partition(rows, question)
24
25     true_branch = build_tree(true_rows)
26
27     false_branch = build_tree(false_rows)
28
29     return Decision_Node(question, true_branch, false_branch)
30
31 # To Print the tree
32 def print_tree(node, spacing=""):
33
34     if isinstance(node, Leaf):
35         print (spacing + "Predict", node.predictions)
36         return
37
38     print (spacing + str(node.question))
39
40     print (spacing + '--> True:')
41     print_tree(node.true_branch, spacing + " ")
42
43     print (spacing + '--> False:')
44     print_tree(node.false_branch, spacing + " ")
```

```
In [326]: 1 my_tree = build_tree(training)
          2 print_tree(my_tree)
```

```
Is diameter >= 3?
--> True:
    Is color == Yellow?
    --> True:
        Predict {'Apple': 1, 'Lemon': 1}
    --> False:
        Predict {'Apple': 1}
--> False:
    Predict {'Grape': 2}
```

Classification

Functions to classify and print the leaf(predictions)

```
In [327]: 1 def classify(row, node):
          2     if isinstance(node, Leaf):
          3         return node.predictions
          4
          5     if node.question.match(row):
          6         return classify(row, node.true_branch)
          7     else:
          8         return classify(row, node.false_branch)
          9
         10 def print_leaf(counts):
         11     total = sum(counts.values()) * 1.0
         12     probs = {}
         13     for lbl in counts.keys():
         14         probs[lbl] = str(int(counts[lbl] / total * 100)) + "%"
         15     return probs
```

```
In [302]: 1 print_leaf(classify(training_data[0], my_tree))
```

```
Out[302]: {'Apple': '100%'}
```

Test the Classifier on the Test data

```
In [304]: 1 for row in testing:
          2     print ("Actual: %s. Predicted: %s" %
          3           (row[-1], print_leaf(classify(row, my_tree))))
```

```
Actual: Apple. Predicted: {'Apple': '100%'}
Actual: Apple. Predicted: {'Apple': '50%', 'Lemon': '50%'}
Actual: Grape. Predicted: {'Grape': '100%'}
Actual: Grape. Predicted: {'Grape': '100%'}
Actual: Lemon. Predicted: {'Apple': '50%', 'Lemon': '50%'}
```

```
In [ ]: 1
```

