

In [47]:

```
1 from __future__ import print_function
```

In [48]:

```
1 from sklearn.datasets import load_breast_cancer
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import cross_val_score
```

In [126]:

```
1 cancer = load_breast_cancer()# 导入数据
2 X_train,X_test,y_train,y_test = train_test_split(cancer.data,cancer.target,stratify=cancer.target)
3 header = cancer.feature_names
4 print(cancer.feature_names)
5 cancer.target_names
6 y_train = y_train.reshape((426,1))
7 y_test = y_test.reshape((143,1))
8 traindata = np.hstack((X_train,y_train)) #所给程序中标签和特征在同一数组中, 在这里把他们弄到同一数组中
9 testdata = np.hstack((X_test,y_test)) #所给程序中标签和特征在同一数组中, 在这里把他们弄到同一数组中
```

```
['mean radius' 'mean texture' 'mean perimeter' 'mean area'
 'mean smoothness' 'mean compactness' 'mean concavity'
 'mean concave points' 'mean symmetry' 'mean fractal dimension'
 'radius error' 'texture error' 'perimeter error' 'area error'
 'smoothness error' 'compactness error' 'concavity error'
 'concave points error' 'symmetry error' 'fractal dimension error'
 'worst radius' 'worst texture' 'worst perimeter' 'worst area'
 'worst smoothness' 'worst compactness' 'worst concavity'
 'worst concave points' 'worst symmetry' 'worst fractal dimension']
```

In [127]:

```

1 class Question:
2     """A Question is used to partition a dataset.
3
4     This class just records a 'column number' (e.g., 0 for Color) and a
5     'column value' (e.g., Green). The 'match' method is used to compare
6     the feature value in an example to the feature value stored in the
7     question. See the demo below.
8     """
9
10    def __init__(self, column, value):
11        self.column = column
12        self.value = value
13
14    def match(self, example):
15        # Compare the feature value in an example to the
16        # feature value in this question.
17        val = example[self.column]
18        if is_numeric(val):
19            return val >= self.value
20        else:
21            return val == self.value
22
23    def __repr__(self):
24        # This is just a helper method to print
25        # the question in a readable format.
26        condition = "=="
27        if is_numeric(self.value):
28            condition = ">="
29        return "Is %s %s %s?" % (
30            header[self.column], condition, str(self.value))

```

In [128]:

```

1 def class_counts(rows):
2     """Counts the number of each type of example in a dataset."""
3     counts = {} # a dictionary of label -> count.
4     for row in rows:
5         # in our dataset format, the label is always the last column
6         label = row[-1]
7         if label not in counts:
8             counts[label] = 0
9         counts[label] += 1
10    return counts

```

In [129]:

```

1 def is_numeric(value):
2     """Test if a value is numeric."""
3     return isinstance(value, int) or isinstance(value, float)

```

In [130]:

```
1 Question(1, 3)
```

Out[130]:

Is mean texture >= 3?

In [131]:

```
1 def partition(rows, question):
2     """Partitions a dataset.
3
4     For each row in the dataset, check if it matches the question. If
5     so, add it to 'true rows', otherwise, add it to 'false rows'.
6     """
7     true_rows, false_rows = [], []
8     for row in rows:
9         if question.match(row):
10             true_rows.append(row)
11         else:
12             false_rows.append(row)
13     return true_rows, false_rows
```

In [132]:

```
1 def gini(rows):
2     """Calculate the Gini Impurity for a list of rows.
3
4     There are a few different ways to do this, I thought this one was
5     the most concise. See:
6     https://en.wikipedia.org/wiki/Decision_tree_learning#Gini_impurity
7     """
8     counts = class_counts(rows)
9     impurity = 1
10    for lbl in counts:
11        prob_of_lbl = counts[lbl] / float(len(rows))
12        impurity -= prob_of_lbl**2
13    return impurity
```

In [133]:

```
1 lebal = [['malignant'],
2           ['benign']]
```

In [134]:

```
1 gini(lebal)
```

Out[134]:

0.5

In [135]:

```
1 def info_gain(left, right, current_uncertainty):
2     """Information Gain.
3
4     The uncertainty of the starting node, minus the weighted impurity of
5     two child nodes.
6     """
7     p = float(len(left)) / (len(left) + len(right))
8     return current_uncertainty - p * gini(left) - (1 - p) * gini(right)
```

In [136]:

```
1 current_uncertainty = gini(traindata)
2 current_uncertainty
```

Out[136]:

0.46786351914302715

In [137]:

```
1 def find_best_split(rows):
2     """Find the best question to ask by iterating over every feature / value
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 it
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
```

In [138]:

```
1 best_gain, best_question = find_best_split(traindata)
2 best_question
```

Out[138]:

Is mean concave points >= 0.04938?

In [139]:

```
1 class Leaf:
2     """A Leaf node classifies data.
3
4     This holds a dictionary of class (e.g., "Apple") -> number of times
5     it appears in the rows from the training data that reach this leaf.
6     """
7
8     def __init__(self, rows):
9         self.predictions = class_counts(rows)
```

In [140]:

```
1 class Decision_Node:
2     """A Decision Node asks a question.
3
4     This holds a reference to the question, and to the two child nodes.
5     """
6
7     def __init__(self,
8                 question,
9                 true_branch,
10                false_branch):
11         self.question = question
12         self.true_branch = true_branch
13         self.false_branch = false_branch
```

In [141]:

```
1 def build_tree(rows):
2     """Builds the tree.
3
4     Rules of recursion: 1) Believe that it works. 2) Start by checking
5     for the base case (no further information gain). 3) Prepare for
6     giant stack traces.
7     """
8
9     # Try partitioning the dataset on each of the unique attribute,
10    # calculate the information gain,
11    # and return the question that produces the highest gain.
12    gain, question = find_best_split(rows)
13
14    # Base case: no further info gain
15    # Since we can ask no further questions,
16    # we'll return a leaf.
17    if gain == 0:
18        return Leaf(rows)
19
20    # If we reach here, we have found a useful feature / value
21    # to partition on.
22    true_rows, false_rows = partition(rows, question)
23
24    # Recursively build the true branch.
25    true_branch = build_tree(true_rows)
26
27    # Recursively build the false branch.
28    false_branch = build_tree(false_rows)
29
30    # Return a Question node.
31    # This records the best feature / value to ask at this point,
32    # as well as the branches to follow
33    # depending on the answer.
34    return Decision_Node(question, true_branch, false_branch)
```

In [142]:

```
1 def print_tree(node, spacing=""):
2     """World's most elegant tree printing function."""
3
4     # Base case: we've reached a leaf
5     if isinstance(node, Leaf):
6         print (spacing + "Predict", node.predictions)
7         return
8
9     # Print the question at this node
10    print (spacing + str(node.question))
11
12    # Call this function recursively on the true branch
13    print (spacing + '--> True:')
14    print_tree(node.true_branch, spacing + " ")
15
16    # Call this function recursively on the false branch
17    print (spacing + '--> False:')
18    print_tree(node.false_branch, spacing + " ")
```

In [143]:

```
1 def print_tree(node, spacing=""):
2     """World's most elegant tree printing function."""
3
4     # Base case: we've reached a leaf
5     if isinstance(node, Leaf):
6         print (spacing + "Predict", node.predictions)
7         return
8
9     # Print the question at this node
10    print (spacing + str(node.question))
11
12    # Call this function recursively on the true branch
13    print (spacing + '--> True:')
14    print_tree(node.true_branch, spacing + " ")
15
16    # Call this function recursively on the false branch
17    print (spacing + '--> False:')
18    print_tree(node.false_branch, spacing + " ")
```

In [144]:

```
1 my_tree = build_tree(traindata)
```

In [145]:

```
1 print_tree(my_tree)
```

```
Is mean concave points >= 0.04938?
--> True:
  Is worst concavity >= 0.2249?
  --> True:
    Is mean texture >= 14.26?
    --> True:
      Is worst perimeter >= 97.65?
      --> True:
        Predict {0.0: 140}
      --> False:
        Is worst texture >= 26.38?
        --> True:
          Predict {0.0: 3}
        --> False:
          Predict {1.0: 3}
      --> False:
        Predict {1.0: 3}
    --> False:
      Is worst fractal dimension >= 0.06599?
      --> True:
        Predict {1.0: 11}
      --> False:
        Predict {0.0: 2}
    --> False:
      Is worst radius >= 16.89?
      --> True:
        Is worst texture >= 20.24?
        --> True:
          Is concave points error >= 0.01033?
          --> True:
            Predict {1.0: 2}
          --> False:
            Predict {0.0: 12}
        --> False:
          Predict {1.0: 4}
      --> False:
        Is area error >= 49.11?
        --> True:
          Is worst symmetry >= 0.2179?
          --> True:
            Predict {1.0: 1}
          --> False:
            Predict {0.0: 2}
        --> False:
          Predict {1.0: 243}
```


In [146]:

```
1 def classify(row, node):
2     """See the 'rules of recursion' above."""
3
4     # Base case: we've reached a leaf
5     if isinstance(node, Leaf):
6         return node.predictions
7
8     # Decide whether to follow the true-branch or the false-branch.
9     # Compare the feature / value stored in the node,
10    # to the example we're considering.
11    if node.question.match(row):
12        return classify(row, node.true_branch)
13    else:
14        return classify(row, node.false_branch)
```

In [147]:

```
1 classify(traindata[1], my_tree)
```

Out[147]:

```
{1.0: 243}
```

In [148]:

```
1 def print_leaf(counts):
2     """A nicer way to print the predictions at a leaf."""
3     total = sum(counts.values()) * 1.0
4     probs = {}
5     for lbl in counts.keys():
6         probs[lbl] = str(int(counts[lbl] / total * 100)) + "%"
7     return probs
```

In [149]:

```
1 print_leaf(classify(traindata[0], my_tree))
```

Out[149]:

```
{0.0: '100%'}
```

In [150]:

```
1 print_leaf(classify(traindata[0], my_tree))
```

Out[150]:

```
{0.0: '100%'}
```

In [151]:

```
1 testdata
```

Out[151]:

```
array([[8.219e+00, 2.070e+01, 5.327e+01, ..., 3.322e-01, 1.486e-01,
        1.000e+00],
       [1.225e+01, 1.794e+01, 7.827e+01, ..., 3.113e-01, 8.132e-02,
        1.000e+00],
       [9.295e+00, 1.390e+01, 5.996e+01, ..., 3.681e-01, 8.982e-02,
        1.000e+00],
       ...,
       [1.831e+01, 2.058e+01, 1.208e+02, ..., 3.074e-01, 7.863e-02,
        0.000e+00],
       [1.705e+01, 1.908e+01, 1.134e+02, ..., 3.109e-01, 9.061e-02,
        0.000e+00],
       [1.170e+01, 1.911e+01, 7.433e+01, ..., 3.487e-01, 6.958e-02,
        1.000e+00]])
```

In [152]:

```
1 for row in testdata:
2     print ("Actual: %s. Predicted: %s" %
3           (row[-1], print_leaf(classify(row, my_tree))))#此处1和0分别代表是否有癌症
```

```
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 0.0. Predicted: {0.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {0.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
Actual: 1.0. Predicted: {1.0: '100%'}
```

用现有模型来对数据集进行训练和测试

In [153]:

```
1 DTC = DecisionTreeClassifier().fit(X_train,y_train)
2 accuracy = cross_val_score(DTC, X_test, y_test, scoring='accuracy', cv=5)
3 print("准确率:",accuracy.mean())
```

准确率: 0.9233990147783251

