

Tree data structure: Decision Tree

Hoàng-Nguyên Vũ

1 Lý thuyết về cây quyết định

1.1 Cây quyết định là gì?

Cây quyết định (Decision Tree) là một mô hình học máy dựa trên cấu trúc cây để đưa ra quyết định hoặc dự đoán. Mỗi nút trong cây biểu diễn một điều kiện trên một thuộc tính, nhánh là kết quả của một điều kiện, và lá chứa giá trị dự đoán cuối cùng.

1.2 Ứng dụng của cây quyết định

- **Phân loại (Classification):** Xác định nhóm của một đối tượng dựa trên dữ liệu đầu vào.
- **Hồi quy (Regression):** Dự đoán giá trị số lượng dựa trên dữ liệu đầu vào.
- **Chẩn đoán y tế:** Xác định bệnh dựa trên triệu chứng.
- **Tài chính:** Dự đoán khả năng vỡ nợ của khách hàng.

2 Thước đo Gini trong cây quyết định

2.1 Thước đo Gini là gì?

Gini Index là một chỉ số đo độ tinh khiết của một tập dữ liệu. Nó giúp xác định mức độ hỗn loạn của tập dữ liệu và được sử dụng để chọn thuộc tính tốt nhất để chia nhánh trong cây quyết định.

2.2 Công thức tính Gini

Gini Index được tính bằng công thức:

$$Gini = 1 - \sum_{i=1}^n p_i^2$$

Trong đó:

- p_i là tỷ lệ mẫu thuộc lớp i trong tập dữ liệu.
- n là số lượng lớp.

Nếu một tập dữ liệu chỉ chứa một lớp duy nhất, Gini Index bằng 0 (hoàn toàn tinh khiết). Nếu tỷ lệ của các lớp là đều nhau, Gini Index đạt giá trị cao nhất.

2.3 Ví dụ tính Gini Index

Giả sử có tập dữ liệu phân loại như sau:

Tuổi	Mua sản phẩm?
22	<i>Yes</i>
35	<i>No</i>
26	<i>Yes</i>
45	<i>No</i>
30	<i>Yes</i>
40	<i>No</i>

Tỷ lệ của hai lớp:

- $p_{Yes} = \frac{3}{6} = 0.5$
- $p_{No} = \frac{3}{6} = 0.5$

$$Gini = 1 - (0.5^2 + 0.5^2) = 1 - (0.25 + 0.25) = 0.5$$

3 Cách sử dụng Gini trong cây quyết định

3.1 Các bước xây dựng cây quyết định dựa trên Gini

1. Tính Gini của tập dữ liệu ban đầu.
2. Chia tập dữ liệu thành các nhóm nhỏ dựa trên một thuộc tính.
3. Tính Gini trung bình có trọng số của các nhóm con.
4. Chọn thuộc tính có **Gini thấp nhất** để làm nút phân nhánh.
5. Lặp lại cho đến khi đạt điều kiện dừng.

4 Cài đặt cây quyết định sử dụng Gini

Dưới đây là một triển khai đơn giản của cây quyết định sử dụng Gini Index trong Python:

```

1 import numpy as np
2
3 def gini_index(groups, classes):
4     """Tính toán chỉ số Gini của một tập hợp"""
5     total_samples = float(sum([len(group) for group in groups]))
6     gini = 0.0
7     for group in groups:
8         size = float(len(group))
9         if size == 0:
10             continue

```

```

11     score = sum([(row[-1] == c) for row in group for c in classes]) /
        size
12     gini += (1.0 - sum([score ** 2 for c in classes])) * (size /
        total_samples)
13     return gini
14
15 # Ví dụ dữ liệu
16 dataset = [
17     [2.8, 'Yes'],
18     [1.2, 'No'],
19     [3.6, 'Yes'],
20     [4.5, 'No'],
21     [5.1, 'Yes']
22 ]
23
24 # Chia tập dữ liệu theo một giá trị ngưỡng
25 def split_data(dataset, feature_index, threshold):
26     left = [row for row in dataset if row[feature_index] < threshold]
27     right = [row for row in dataset if row[feature_index] >= threshold]
28     return left, right
29
30 # Ví dụ tính Gini cho một cách chia dữ liệu
31 groups = split_data(dataset, 0, 3.0)
32 classes = ['Yes', 'No']
33 gini = gini_index(groups, classes)
34 print(f'Gini Index: {gini:.4f}')

```

5 Ví dụ: Xây dựng cây quyết định sử dụng OOP

5.1 Lớp TreeNode - Biểu diễn một nút trong cây

Chúng ta định nghĩa một lớp `TreeNode` để lưu trữ thông tin về các nút trong cây.

```

1 class TreeNode:
2     def __init__(self, feature_index=None, threshold=None, left=None,
        right=None, label=None):
3         """
4         Khởi tạo một nút trong cây quyết định.
5         - feature_index: Chỉ số thuộc tính được chọn để chia.
6         - threshold: Ngưỡng giá trị để phân chia dữ liệu.
7         - left: Nhánh trái của cây.
8         - right: Nhánh phải của cây.
9         - label: Nhãn dự đoán nếu là nút lá.
10        """
11        self.feature_index = feature_index
12        self.threshold = threshold
13        self.left = left
14        self.right = right
15        self.label = label

```

5.2 Lớp DecisionTree - Xây dựng cây quyết định

Chúng ta xây dựng lớp DecisionTree với các phương thức chính để chia tập dữ liệu và xây dựng cây.

```

1 import numpy as np
2
3 class DecisionTree:
4     def __init__(self, max_depth=3):
5         """
6         Khởi tạo cây quyết định với độ sâu tối đa.
7         """
8         self.max_depth = max_depth
9         self.root = None
10
11     def gini_index(self, groups, classes):
12         """
13         Tính chỉ số Gini cho một tập hợp.
14         """
15         total_samples = sum([len(group) for group in groups])
16         gini = 0.0
17         for group in groups:
18             size = len(group)
19             if size == 0:
20                 continue
21             score = 0.0
22             for class_val in classes:
23                 proportion = [row[-1] for row in group].count(class_val) /
24                 size
25                 score += proportion ** 2
26             gini += (1.0 - score) * (size / total_samples)
27         return gini
28
29     def split_data(self, dataset, feature_index, threshold):
30         """
31         Chia tập dữ liệu dựa trên một thuộc tính và ngưỡng giá trị.
32         """
33         left = [row for row in dataset if row[feature_index] < threshold]
34         right = [row for row in dataset if row[feature_index] >= threshold]
35         return left, right
36
37     def best_split(self, dataset):
38         """
39         Tìm thuộc tính tốt nhất để chia tập dữ liệu.
40         """
41         class_values = list(set(row[-1] for row in dataset))
42         best_index, best_threshold, best_score, best_groups = None, None,
43         float('inf'), None
44         for index in range(len(dataset[0]) - 1):
45             for row in dataset:
46                 groups = self.split_data(dataset, index, row[index])
47                 gini = self.gini_index(groups, class_values)

```

```

47         if gini < best_score:
48             best_index, best_threshold, best_score, best_groups =
index, row[index], gini, groups
49         return best_index, best_threshold, best_groups
50
51     def build_tree(self, dataset, depth=0):
52         """
53         Xây dựng cây quyết định đệ quy.
54         """
55         class_values = [row[-1] for row in dataset]
56
57         # Điều kiện dừng: Nếu chỉ có một lớp hoặc đạt đến độ sâu tối đa
58         if len(set(class_values)) == 1 or depth >= self.max_depth:
59             return TreeNode(label=max(set(class_values), key=class_values.
count))
60
61         # Tìm thuộc tính và giá trị ngưỡng tốt nhất để chia dữ liệu
62         feature_index, threshold, (left, right) = self.best_split(dataset)
63
64         # Nếu không thể chia tiếp, tạo nút lá
65         if not left or not right:
66             return TreeNode(label=max(set(class_values), key=class_values.
count))
67
68         # Xây dựng nhánh trái và nhánh phải
69         left_node = self.build_tree(left, depth + 1)
70         right_node = self.build_tree(right, depth + 1)
71
72         return TreeNode(feature_index, threshold, left_node, right_node)
73
74     def fit(self, dataset):
75         """
76         Huấn luyện cây quyết định bằng cách xây dựng cây từ dữ liệu đầu và
o.
77         """
78         self.root = self.build_tree(dataset)
79
80     def print_tree(self, node=None, depth=0):
81         """
82         In ra cây quyết định theo dạng phân cấp.
83         """
84         if node is None:
85             node = self.root
86
87         if node.label is not None:
88             print(f"{' ' * depth} [Leaf] Label: {node.label}")
89         else:
90             print(f"{' ' * depth} [Node] Feature {node.feature_index} <=
{node.threshold}")
91             self.print_tree(node.left, depth + 1)
92             self.print_tree(node.right, depth + 1)

```

5.3 Ví dụ sử dụng - Xây dựng cây từ tập dữ liệu

Chúng ta thử nghiệm xây dựng một cây quyết định với một tập dữ liệu đơn giản.

```

1 # Tập dữ liệu đơn giản: [Thuộc tính, Nhân]
2 dataset = [
3     [2.8, 'Yes'],
4     [1.2, 'No'],
5     [3.6, 'Yes'],
6     [4.5, 'No'],
7     [5.1, 'Yes']
8 ]
9
10 # Khởi tạo và huấn luyện cây quyết định
11 tree = DecisionTree(max_depth=3)
12 tree.fit(dataset)
13
14 # In ra cây quyết định
15 print("Cây quyết định được xây dựng:")
16 tree.print_tree()
```

5.4 Kết quả mong đợi

Sau khi chạy đoạn mã trên, cây quyết định sẽ được hiển thị theo dạng phân cấp:

Cây quyết định được xây dựng:

```

[Node] Feature 0 <= 3.6
  [Node] Feature 0 <= 2.8
    [Leaf] Label: No
    [Leaf] Label: Yes
  [Node] Feature 0 <= 4.5
    [Leaf] Label: Yes
    [Leaf] Label: No
```

6 Bài tập thực hành

Bài 1: Tính chỉ số Gini cho một tập dữ liệu

Tính toán chỉ số Gini cho tập dữ liệu sau:

Lương	Đủ điều kiện vay vốn?
50	<i>Yes</i>
20	<i>No</i>
30	<i>No</i>
70	<i>Yes</i>
40	<i>No</i>
60	<i>Yes</i>

Bài 2: Mở rộng cây quyết định

- Thêm nhiều thuộc tính hơn vào tập dữ liệu (ví dụ: Tuổi, Điểm tín dụng).
- Cải tiến thuật toán để xây dựng cây với nhiều thuộc tính.
- In ra cây theo cách dễ đọc hơn.