

Assignment 2: Decision Tree

2019094511 김준표

1. Summary of algorithm

해당 과제에서는 gain ratio 방식을 이용하여 decision tree를 생성했다.

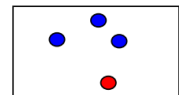
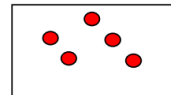
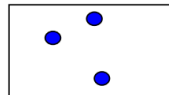
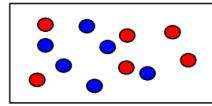
Gain ratio 방식:

이는 특정 attribute에 대해 information gain을 구한다. Information gain은 분류 전 집단에서 entropy를 구하고, 기준 attribute에 대해 분류 후 entropy를 구한 뒤 둘의 차이를 반환한 값이다.

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

$$Gain(A) = Info(D) - Info_A(D)$$



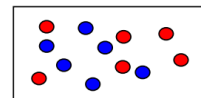
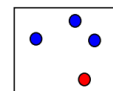
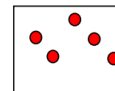
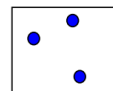
Gain ratio은 위의 방식으로 구한 information gain에서 attribute마다 가지는 value set의 요소 수가 달라 발생하는 bias를 보정하기 위한 방식이다. 이는 이전에 구한 gain 값을 split info값으로 나눠 반환한다.

□ **Gain Ratio** takes the number and size of branches into account when choosing a feature

$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$Gain(A) = Info(D) - Info_A(D)$$

$$GainRatio(A) = Gain(A) / SplitInfo_A(D)$$



2. Code description

main 함수에서는 input file로부터 받은 train file과 test file을 data frame에 넣어준다. 이때, csv 형식으로 읽고 쓰도록 한다. 이후 decision tree를 생성하고 분류를 실시한 뒤 output으로 result file을 작성한다.

```

126 ∨ def main():
127     train_file = sys.argv[1]
128     test_file = sys.argv[2]
129     result_file = sys.argv[3]
130
131     train_df = pd.read_csv(train_file, sep='\t')
132     test_df = pd.read_csv(test_file, sep='\t')
133
134     dt = DecisionTree(train_df)
135
136     df = dt.classify(test_df)
137     df.to_csv(result_file, index=False, sep='\t')
138
139
140 ∨ if __name__ == "__main__":
141     main()
142

```

Node class는 decision tree의 노드를 저장하는 구조체로 해당 노드의 attribute, leaf 여부, label, child node를 저장한다.

```

5 class Node:
6     def __init__(self, attr, isLeaf = False, label = None):
7         self.attr = attr          # DB의 카테고리 (ex: attr 1, attr 2, ... , attr n)
8         self.isleaf = isLeaf      # 노드 속성 (ex: leaf node/internal node)
9         self.label = label        # 데이터의 최종 분류 기준 (ex: yes/no)
10        self.child = dict()

```

DecisionTree class는 tree를 초기화하고, gain ratio를 계산하고, 이를 바탕으로 decision tree를 생성하고, data frame을 받아와 label을 분류해주는 메소드로 이루어져있다.

```

13     """
14     Decision tree을 초기화한다.
15     """
16     def __init__(self, df):
17         # Data frame의 마지막 column에 label 정보가 담겨있음
18         self.label = df.columns[-1]
19         # Data frame의 attribute와 attribute의 value set들 추출하기
20         self.attrs = {attr: df[attr].unique() for attr in df.columns if attr!=self.label}
21         # Decision tree의 root 노드
22         self.root = self.create_dt(df)

```

Decision tree는 label (최종 결정), attribute와 value set 집합의 딕셔너리, root 노드를 갖는다.

```

24     """
25     Data frame의 label 값들을 분류하여 각 값이 나올 확률을 구한다.
26     """
27     def probability(self, df):
28         values = df[self.label]
29         total_cnt = len(values)
30         return [float(value)/float(total_cnt) for value in values.value_counts()]
31
32     """
33     확률을 받아와 entropy를 구한다.
34     """
35     def entropy(self, prob):
36         return sum(-1.0 * (p * math.log(p + 1e-15, 2)) for p in prob if p)

```

entropy 계산은 알고리즘 소개의 공식을 따른다. 계산을 위해 data frame의 label (최종 결정)값을 분류하여 각 label이 나올 확률을 계산하는 함수를 추가해준다.

```

38     '''
39     특정 attribute를 인자로 받아 해당 attribute의 information gain을 구한다.
40     '''
41     def information_gain(self, df, attr):
42         info_D = self.entropy(self.probability(df))
43         info_A = 0.0
44         for a in self.attrs[attr]:
45             # 속성이 a인 row들을 뺌 (ex: attr(income)이 a (high, medium, low)인 애들을 각각 뺌)
46             tmp_df = df[df[attr] == a]
47             # 각각에 대해 entropy를 계산하고 가중치를 곱해줌
48             info_A += (len(tmp_df) / len(df)) * self.entropy(self.probability(tmp_df))
49
50         gain = info_D - info_A
51         return gain
52
53     '''
54     Split information 값을 계산한다.
55     '''
56     def split_info(self, df, attr):
57         splitInfo = 0.0
58         for a in self.attrs[attr]:
59             # 속성이 a인 row들을 뺌
60             tmp_df = df[df[attr] == a]
61             p = len(tmp_df) / len(df)
62             splitInfo += (-1.0) * p * (math.log(p + 1e-15, 2))
63
64         return splitInfo
65
66     '''
67     Information gain과 split information을 통해 gain ratio를 구한다.
68     '''
69     def gain_ratio(self, df, attr):
70         gain = self.information_gain(df, attr)
71         splitInfo = self.split_info(df, attr)
72
73         return gain / splitInfo

```

Gain ratio 역시 알고리즘을 따라 작성했다. Data frame을 받아와 특정 속성에 대해, 해당 속성의 value set을 따라 data frame을 나눠 entropy를 계산하고 이를 통해 information gain을 구하고 split info 값을 구하여 gain ratio를 구한다.

```

75     '''
76     Decision tree를 생성한다.
77     '''
78     def create_dt(self, df):
79         # Data frame이 분류가 완벽하여 label이 한가지만 존재하는 경우
80         # Label을 정한다.
81         if df[self.label].nunique() == 1:
82             label = df[self.label].unique()[0]
83             return Node(None, True, label)
84
85         # 분류 방법으로 gain ratio를 이용한다.
86         result_dict = {attr: self.gain_ratio(df, attr) for attr in df.columns if attr != self.label}
87         selected_attr = sorted(result_dict.items(), key=lambda x: x[1], reverse = True)[0][0]
88
89         # 노드를 생성하고 선택된 속성을 기준으로 split한다.
90         node = Node(selected_attr, False, None)
91
92         a = self.attrs[node.attr]
93         for i in range(len(a)):
94             tmp_df = df[df[node.attr].values == a[i]]
95             tmp_df = tmp_df.drop(columns=node.attr, axis=1)
96             if (len(tmp_df) > 0):
97                 node.child[a[i]] = self.create_dt(tmp_df)
98
99         return node

```

Decision tree 생성은 data frame의 모든 attribute에 대해서 gain ratio를 적용해보고 이를 정렬하

여 gain ratio의 반환 값이 제일 커지는 attribute를 골라 해당 attribute를 기준으로 노드를 split하고 child를 생성한다.

```
101     '''
102     Data가 주어진 경우 해당 data의 leaf node를 찾는다.
103     '''
104     def find_leaf(self, data):
105         node = self.root
106         a = data[node.attr]
107         while not node.isleaf:
108             for (child, next_node) in node.child.items():
109                 if a in child:
110                     break
111             node = next_node
112             if node.attr:
113                 a = data[node.attr]
114
115         return node.label
116
117     '''
118     Data frame을 받아서 label을 분류해준다.
119     '''
120     def classify(self, df):
121         label = [self.find_leaf(df.loc[i]) for i in range(len(df))]
122         df[self.label] = label
123         return df
```

분류 작업은 생성된 decision tree의 노드를 따라 attribute를 기준으로 값을 비교하고 leaf node까지 순회하여 최종 label을 반환한다.

3. Compiling method

실행환경:

OS: Windows 10

Python version: Python 3.8.3

실행방법:

```
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2>
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> python ./dt.py dt_train.txt dt_test.txt dt_result.txt
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> |
```

python dt.py train_file_name.txt test_file_name.txt result_file_name.txt

위의 실행 명령어를 입력하면 아래와 같이 test.txt에 대해 result.txt 파일이 생성되는 것을 확인할 수 있다.

dt_test.txt	dt_result.txt
1 age income student credit_rating	1 age income student credit_rating Class:buys_computer
2 <=30 low no fair	2 <=30 low no fair no
3 <=30 medium yes fair	3 <=30 medium yes fair yes
4 31...40 low no fair	4 31...40 low no fair yes
5 >40 high no fair	5 >40 high no fair yes
6 >40 low yes excellent	6 >40 low yes excellent no
7	7

4. Other specification

dt_test.exe 결과

```
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2>  
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> ./dt_test.exe dt_answer.txt dt_result.txt  
5 / 5  
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> █
```

첫번째 예제인 buy computer에 대해서 5 / 5의 결과를 확인할 수 있었다.

```
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2>  
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> ./dt_test.exe dt_answer1.txt dt_result1.txt  
317 / 346  
PS C:\Users\pyo99\CSE\DS\2022_ite4005_2019094511\Assignment#2> █
```

두번째 예제인 car evaluation에 대해서 317 / 346의 결과를 확인할 수 있었다.