# **Data Science (ITE4005)**

# < Programming Assignment #2 >

Decision Tree - Gain Ratio

2016025732 이영석 2021-04-13

### 1. Summary

Gain Ratio 방식을 이용해 Decision Tree를 구축하고, Test Data의 Class Label을 Predict한다.

```
class Node():
    def __init__(self):
        self.criteria_attribute = None
        self.label_predict = None
        self.leaf = {} # key : attribute value, value: leaf node pointer
```

```
class DecisionTree():
    def __init__(self):
        self.root = Node()

    self.df_train = None
    self.df_test = None

    self.attribute = [] # (ex. age, income, etc..)
    self.label_name = None # (ex. buy_computer)
    self.label_value = [] # label_name으로 구분한 row (ex. no, yes)
```

Class는 Node, DecisionTree 두개로 구성된다.

Node의 label\_predict는 해당 노드의 예측 label 변수가 저장된다. (by Majority Voting)

criteria\_attribute는 Gain Ratio를 통해 얻어진 best\_attribute 이름이 저장된다.

DecisionTree의 df\_train/df\_test는 각각 train data와 test data가 dataframe으로 저장된다.

```
if __name__ == '__main__' :
    tree = DecisionTree()
    tree.set_data(sys.argv[1], sys.argv[2])
    tree.train_data()
    tree.test_data()
    tree.print_output(sys.argv[3])
```

주 알고리즘은 다음과 같은 순으로 구성하였다.

Setting Data -> Training Data -> Predicting Test Data's Label -> Printing Output

각 단계(모듈)는 다음과 같이 구성하였다. (자세한 내용은 2. Detailed Description 참고)

- A. set\_data(arg1, arg2)
- B. train\_data()
  - i. grow\_tree(data)
    - 1. get\_best\_attribute(data)
      - A. info(data)
      - B. expected\_info(data, attribute)
      - C. split\_info(data, attribute)
- C. test\_data()
  - i. predict(data)
- D. print\_output(arg3)

### 2. Detailed Description

- A. set\_data(arg1, arg2)
  - → train 파일과 test 파일의 dataframe화 진행.
  - → attribute, label\_name, label\_value 추출
- B. train\_data()

- i. grow\_tree(data)
  - → majority voting한 label 해당 노드에 저장
  - → Best Attribute를 받아 Data를 나누고 해당 attribute는 버리며 self-recursive call하며 확장
  - → 루트 노드부터 차근차근 확장 (Top-Down / Divide & Conquer)
  - → 해당 노드와 다음노드만을 고려 (Greedy)

```
def grow_tree(self, node, data):

# Predict label by majority voting
label, counts = np.unique(data[self.label_name], return_counts=True)
node.label_predict = label[np.argmax(counts)]

if self.is_homogenius(data):
    return

best_attribute = self.get_best_attribute(data)

# Set criteria attribute
node.criteria_attribute = best_attribute
attribute_value = np.unique(data[best_attribute])

# Recursively grow
for i in range(len(attribute_value)):
    node.add(attribute_value[i], Node())
    _data = data.get(data[best_attribute]==attribute_value[i]) #해당 attribute에서 같은 value만 가진 data select
    _data = _data.drop(best_attribute, axis=1) #해당 attribute秒 삭제
    self.grow_tree(node.leaf[attribute_value[i]], _data)
```

- ii. get best attribute(data)
  - → 모든 attribute별로 gain ratio를 구하고, 가장 큰 값(이전과 다음 노 드 엔트로피 차이가 가장 큰 값)을 가진 Attribute return

```
def get_best_attribute(self, data):
    temp = 0
    best_attribute = None
    temp_attribute = data.columns.tolist()
    del temp_attribute[-1] #Class label은 제외

for i in range(len(temp_attribute)):
    gain = self.info(data) - self.expected_info(data, temp_attribute[i])
    gain_ratio = gain / self.split_info(data, temp_attribute[i])
    if temp < gain_ratio:
        temp = gain_ratio
        best_attribute = temp_attribute[i]

return best_attribute
```

1. info(data)

- → 해당 노드에 해당하는 Entrophy값 return
- 2. expected\_info(data, attribute)
  - → test attribute로 내려간 자식 노드에 해당하는 Entrophy return
- 3. split\_info(data, attribute)
  - → gain ratio = gain / split\_info
  - → gain ratio를 구하기 위한 분모값 return

(나누는 이유: attribute value 개수 별로 가중치 주는 것을 없애기 위함)

### C. test\_data()

- → test data 열마다 predict함수 실행
- iii. predict(data)
  - → Tree의 leaf까지 이동후, 해당 노드의 label\_predict 값 새 행에 저장.
  - → ('Door == 3' 같은 Training Data에 없는 내용이 들어올 경우 try catch로 예외처리. Dictionary의 key가 없으면 해당 노드의 label\_predict를 return한다)
- D. print\_output(arg3)
  - → Class Label이 추가된 data frame을 txt 파일로 저장

### 3. Instruction for compiling

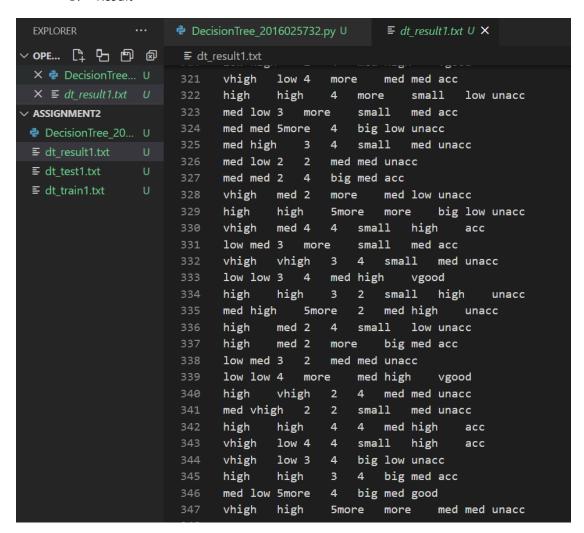
- A. python 설치
- B. train file, test file 같은 디렉토리 내 위치

C. Execute the program with three arguments: training file name, test file name, output file name

python DecisionTree\_2016025732.py dt\_train1.txt dt\_test1.txt dt\_result1.txt

#### 4. Testing

- A. 3번 과 같이 실행
- B. Result



- C. 채점 모듈 실행 결과
- → Gain Ratio (제출본)

C:#Users#이영석#Desktop#test>dt\_test.exe dt\_answer1.txt dt\_result1.txt 323 / 346 → Information Gain (대조군)

## C:₩Users₩이영석₩Desktop₩test>dt\_test.exe dt\_answer1.txt result.txt 320 / 346

근소한 차이로 Gain Ratio의 정확도가 더 높았다.

결론 : 이를통해 Information Gain의 단점(Attribute Value들의 개수가 많을수록 Best Attribute가 되기 쉬움)을 Gain Ratio가 보완한다는 것을 알 수 있었다.

다만, Gain Ratio의 단점은 unbalanced splits에 우위를 두는데, balanced splits에 우위를 두고 싶다면 Gini Index 방식을 쓰는 것이 효율적이다. Training Set의 특성에 맞춰 어떤 Decision Tree Algorithm을 선택할지가 관건이겠다.