# Data Science [ITE4005]
# **Programming Assignment #2**
## : Build a decision tree, and classify the test set

2015004120

Software Engineering

Park Seonha

## 0. Environment

- OS : Windows
- Language : Python3 (version : 3.4.4)
- Executable file : dt.py (../Programming_Assignment_2/project_dt/dt.py)

## 1. Summary of Algorithm

This programming assignment's goal if "build a decision tree, and classify the test set using it." And, it can be seperated "generate a decision tree" and "get test set and classify with decision tree"

I used gain ratio to select proper attribute to maximize the decrement of entropy, and generated tree recursively with classified in each node.

Generating tree, check the node's attribute name list or attribute list is empty, and if it is, return majority class in that node. And, if all values are in one same class, return that class. If node is not in that case, calculate gain ratio and select proper attribute and make child tree. Child tree has all classified tuples in parent tree.

To calculate Gain ratio, first, calculate expected information and residual information in dataset.

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2(p_i) \quad Info_a(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$

and get information gain using $Info(D)$ and $Info_a(D)$. Information gain $Gain(A)$ is
$$Gain(A) = Info(D) - Info_a(D)$$

Since this gain is biased towards attributes with large amount of values, to normalize this gain to gain ratio, calculate splitinfo.

$$SplitInfo_a(D) = -\sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|} \quad , \quad GainRatio(A) = Gain(A)/SplitInfo(A)$$

Choose the attribute which has the maximum gain ratio, and classify the tuples with that attribute.

After generate whole decision tree, get test tuples and trace the tree during the answer is found or tree is end and can't get answer. If tracing finished and can't get answer, find the label which has largest value.

## 2. Description of codes

This assignment program has 3 source code files. "DecisionTree.py", "TreeNode.py", "dt.py". "dt.py" has the main routine of this program.

### Treenode.py

```
3    class TreeNode:
4        value = ""
5        childs = []
6
7        def __init__(self,val, child):
8            self.value = val
9            if (isinstance(child,dict)):
10               self.childs = child.keys()
11
12       def __repr__(self):
13           return self.value
14
```

Treenode.py has a class named TreeNode. As its name shows, this contains each node's attribute name(value) and list of child node(childs). __init__(self, val, child) is initializer of this class, __repr__(self) returns its attribute name.

### DecisionTree.py

This file is implementing the method which generates tree recursively and functions need for implementing that method.

```
153    # Generate Tree recursively
154    def GenerateTree(attribute_list, attribute_name):
155        target = attribute_name[len(attribute_name)-1]
156        attribute_list = attribute_list[:]
157        values = []
158        idx = attribute_name.index(target)
159        for attr in attribute_list:
160            values.append(attr[idx])
161
162        # if attribute_list is empty or no choice of attribute name, return majority class
163        if not attribute_list or (len(attribute_name) - 1) <= 0:
164            return majorityClass(attribute_name, attribute_list, target)
165
166        #if all valaues are in same class
167        elif values.count(values[0]) == len(values):
168            return values[0]
169
170        else :
171            # select attribute using gain ratio and generate tree
172            selected = GainRatio(makeDataPartition(attribute_name, attribute_list), attribute_name, len(attribute_list))
173            tree = {selected:{}}
174
175
176            for value in getValues(attribute_list, attribute_name, selected):
177                all_attribute = getEntries(attribute_list, attribute_name, selected, value)
178                new_attribute = attribute_name[:]
179                new_attribute.remove(selected)
180
181                #generate child tree
182                child = GenerateTree(all_attribute, new_attribute)
183                tree[selected][value] = child
184                tree[selected]["num"] = getNumOfAnswer(attribute_list)
185        return tree
```

GenerateTree(attribute_list, attribute_name) is the function which generate tree recursively. In this Function, tree is generated in dictionary(key, value) form. Its terminal condition is ○ attribute list (list of all tuples) is empty or attribute name list is empty so can't choice next node ○ all tuples are in same class. In other case, calculate gain ratio and generate sub-tree with its attribute list and attribute name list except selected

just before.

After generate sub-tree, append the key named 'num' which has the number of each class label in that node's attribute list.

To implement this function, I implemented functions majorityClass(attributes, data, target), GainRatio(D,attribute_name, total_attr), getValues(data, attributes, attribute), getEntries(data, attributes, selected, val), getNumOfAnswer(attribute_list).

```python
# follow majority rule to finish tree
def majorityClass(attributes, data, target):
    freq_value = {}
    index = attributes.index(target)
    for entry in data:
        if entry[index] in freq_value:
            freq_value[entry[index]] += 1
        else:
            freq_value[entry[index]] = 1
    max = 0
    major_label = ""
    for key in freq_value.keys():
        if freq_value[key] > max:
            max = freq_value[key]
            major_label = key
    return major_label
```

In majorityClass(attribute, data, target) function, get all class label attribute[index]'s number in parameter data, find most frequent class label and return it.

```python
# get proper attributes's value in attribute
def getValues(data, attributes, attribute):
    index = attributes.index(attribute)
    values = []
    for entry in data:
        if entry[index] not in values:
            values.append(entry[index])
    return values
```

getValues(data, attributes, attribute) is the function to get class label in specific attribute attributes[index]. It returns like ['high', 'med', 'low'], ['unacc','acc','good', vgood]

```python
# get proper entries in all attribute
def getEntries(data, attributes, selected, val):
    all_entries = [[]]
    index = attributes.index(selected)
    for entry in data:
        if entry[index] == val:
            new_entry = []
            for i in range(0,len(entry)):
                if i != index :
                    new_entry.append(entry[i])
            all_entries.append(new_entry)
    all_entries.remove([])
    return all_entries
```

getEntries(data, attributes, selected, val) is the function to get all proper entries in data. It returns like ['high','no','excellent','yes']

```python
143     # check all answers in the node
144     def getNumOfAnswer(attribute_list):
145         result_dict = {}
146         for attr in attribute_list:
147             if attr[len(attr)-1] not in result_dict:
148                 result_dict[attr[len(attr)-1]] = 1
149             else:
150                 result_dict[attr[len(attr) - 1]] += 1
151         return result_dict
```

getNumOfAnswer(attribute_list) makes dictionary whose key is last class label like ['yes','no'], ['unacc','acc','good', vgood], and value is the number of each class label. The result of this function will be append in tree as {'num': {'yes':3, 'no':2}}

```python
23    # calculate gain and splitinfo to get gain ratio
24    def GainRatio(D, attribute_name, total_attr):
25        #get Info(D)
26        info_dic = D[attribute_name[len(attribute_name)-1]]
27        info_val_list = []
28        for k, v in info_dic.items():
29            if k in v:
30                info_val_list.append(v[k])
31
32        info_d = 0.0
33        for i in info_val_list:
34            if i > 0:
35                info_d -= (float(i)/float(total_attr)) * math.log2((float(i)/float(total_attr)))
36            else:
37                info_d -= 0
38        del D[attribute_name[len(attribute_name)-1]]
```

GainRatio(D, attribute_name, total_attr) is the function which calculate gain ratio with data partition D, attribute name list, and the number of total attribute. This function returns attribute name which has the maximum gain ratio. So it has several calculation part. First, calculate $Info(D)$.

```python
40        #calculate gain
41        gain_list = {}
42        for k,v in D.items():
43            calc_info = 0.0
44            for k1, v1 in v.items():
45                total_in_class = 0
46                list_in_class = []
47                for k2, v2 in v1.items():
48                    list_in_class.append(v2)
49                    total_in_class += v2
50                for num in list_in_class:
51                    if num > 0:
52                        info_dj = (float(num)/float(total_in_class)) * math.log2(float(num)/float(total_in_class))
53                    else:
54                        info_dj = 0
55                    calc_info -= float(total_in_class)/float(total_attr) * info_dj
56            gain_list[k] = info_d - calc_info

58        # calculate splitInfo
59        split_info = {}
60        for k, v in D.items():
61            calc_info = 0.0
62            list_in_class = []
63            for k1, v1 in v.items():
64                total_in_class = 0
65                for k2, v2 in v1.items():
66                    total_in_class += v2
67                list_in_class.append(total_in_class)
68            for num in list_in_class:
69                if num > 0:
70                    calc_info -= (float(num)/float(total_attr)) * math.log2(float(num)/float(total_attr))
71                else:
72                    calc_info -= 0
73            split_info[k] = calc_info
```

Make the dictionary gain_list whose key is attribute name and value is the gain of each attribute. Same way, calculate splitInfo too.

```python
75        # do gain/splitInfo to get gain ratio
76        gain_ratio = {}
77        for k, v in gain_list.items():
78            for k1, v1 in split_info.items():
79                if k == k1 :
80                    gain_ratio[k] = v / v1
81
82        # return max value in the list of gain ratio
83        attr_name = max(gain_ratio, key=gain_ratio.get)
84
85        return attr_name
```

Bind two dictionary as gain_ratio. find same key in gain_list and split_info, calculate gain ratio and insert the dictionary gain_ratio. Find max value of gain ratio, get its attribute name, and return it.

To implement this function, we need to make data partition D.

```
111    # make data partition D with all attribute and attribute names to calculate gain ratio
112    def makeDataPartition(attribute_name, attribute_list):
113        possible_name = {}
114        for i in range(len(attribute_name)):
115            ith_class = []
116            for attr in attribute_list:
117                if attr[i] not in ith_class:
118                    ith_class.append(attr[i])
119            possible_name[attribute_name[i]] = ith_class
120        class_label = possible_name[attribute_name[len(attribute_name)-1]]
121
122        data_partition = {}
123        for i in range(len(attribute_name)):
124            ith_attr = {}
125
126            for attr in attribute_list:
127                if attr[i] not in ith_attr:
128                    ith_attr[attr[i]] = {attr[len(attribute_name)-1] : 1}
129                else :
130                    if attr[len(attribute_name)-1] not in ith_attr[attr[i]]:
131                        ith_attr[attr[i]][attr[len(attribute_name)-1]] = 1
132                    else :
133                        ith_attr[attr[i]][attr[len(attribute_name) - 1]] += 1
134
135            for k,v in ith_attr.items():
136                for label in class_label:
137                    if label not in v:
138                        v[label] = 0
139
140            data_partition[attribute_name[i]] = ith_attr
141        return data_partition
```

makeDataPartition(attribute_name, attribute_list) do that role in this program. Find the number of all class label in each attribute. It returns dictionary like {'credit_rating': {'fair': {'no': 0, 'yes': 3}, 'excellent': {'no': 2, 'yes': 0}}, 'income': {'medium': {'no': 1, 'yes': 2}, 'low': {'no': 1, 'yes': 1}}, 'student': {'no': {'no': 1, 'yes': 1}, 'yes': {'no': 1, 'yes': 2}}, 'Class:buys_computer': {'no': {'no': 2, 'yes': 0}, 'yes': {'no': 0, 'yes': 3}}}

**dt.py**

This file has the main routine of program like file I/O, test program.

```
3    # written in Python3
4    import DecisionTree
5    import TreeNode
6    import sys
7
     if __name__ == '__main__':
         main()
```

It includes other two file and sys module, and main() function is in this file.

```
50   def main():
51       # get training file, test file and output file's name from command line argument
52       train_file = sys.argv[1]
53       test_file = sys.argv[2]
54       output_file = sys.argv[3]
```

In main function, get command line argument using sys module.

```
56       # get training file and make list of columns
57       with open(train_file) as f:
58           train_data = f.readlines()
59       train_data = [d.strip() for d in train_data]
```

open training dataset file.

```python
60
61         # make train_data string to attribute list and get each attribute name and total number of columns
62         attribute_list = []
63         total_attribute = 0
64         for line in train_data:
65             each_line = line.split("\t")
66             attribute_list.append(each_line)
67             total_attribute += 1
68         attribute_name = attribute_list[0]
69         attribute_list.pop(0)
70         total_attribute -= 1
```

attribute_list is the list of entry tuples. Tuples are splited in list. while generating attribute_list, count total number of attribute too.

```python
72         possible_name = {}
73
74         # get possible classes
75         for i in range(len(attribute_name)):
76             ith_class = []
77             for attr in attribute_list:
78                 if attr[i] not in ith_class:
79                     ith_class.append(attr[i])
80             possible_name[attribute_name[i]] = ith_class
```

possible name has each attribute and class label. It contains like {'age': ['<=30', '31...40', '>40'], 'credit_rating': ['fair', 'excellent'], 'Class:buys_computer': ['no', 'yes'], 'student': ['no', 'yes'], 'income': ['high', 'medium', 'low']}

```python
84         # get decision class's label and generate tree
85         decision_label = attribute_name[len(attribute_name) - 1]
86         tree = DecisionTree.GenerateTree(attribute_list, attribute_name)
```

decision_label is the attribute name to predict in testset. and tree is the dictionary which has decision tree.

```python
88         # get test input from test file
89         with open(test_file) as f:
90             test_input = f.readlines()
91         test_input = [d.strip() for d in test_input]
92         test_input.pop(0)
93         test_data = [[]]
94         for line in test_input:
95             test_data.append(line.split("\t"))
96         test_data.remove([])
```

open test_file and get test input in list.

```python
98          # predict decision class and write in output file
99          f = open(output_file,"w")
.00          del attribute_name[attribute_name.index(decision_label)]
.01          for name in attribute_name:
.02              f.write(name + "\t")
.03          f.write(decision_label+"\n")
```

open output file and write test input's attribute name.

```python
104        for entry in test_data:
105            # call decision tree
106            tempDict = tree.copy()
107            parentDict = tempDict
108            rootDict = tree.copy()
109            rootDict = rootDict[list(rootDict.keys())[0]]
110            result = ""
```

for each testset tuple, copy decision tree dictionary and make it as TreeNode,

```
111            # trace tree while the answer is found or tree is end
112            while isinstance(tempDict, dict):
113                root = TreeNode.TreeNode(list(tempDict.keys())[0], tempDict[list(tempDict.keys())[0]])
114
115                tempDict = tempDict[list(tempDict.keys())[0]]
116                index = attribute_name.index(root.value)
117                value = entry[index]
118
119                if (value in list(tempDict.keys())):
120                    child = TreeNode.TreeNode(value, tempDict[value])
121                    result = tempDict[value]
122                    parentDict = tempDict
123                    tempDict = tempDict[value]
124                # can't find the entry in tree, follow majority vote
125                else:
126                    result = getMaxLabel(parentDict, tempDict,rootDict)
127                    break
```

trace tree while the value is in tree and find the result. Find the attribute value in tree. If it is not exist call the function getMaxLabel(parent_attr, this_attr, root_attr) to choose result.

```
128            for i in entry:
129                f.write(i+"\t")
130            f.write(result+"\n")
131        f.close()
```

After choice result, write entry in output file and repeat with next entry. After predict all entry, close output file. The program is end in this part.

```
8    #get label which has largest value. to estimate non-leaf ended entry
9    def getMaxLabel(parent_attr, this_attr, root_attr):
10       max_array = []
11       max_val = 0
12       #find largest number label
13       for k , v in this_attr["num"].items():
14           if v > max_val:
15               max_array = []
16               max_array.append(k)
17               max_val = v
18           elif v == max_val:
19               max_array.append(k)
21       #if two or more label has same value, find parent node's largest label
22       if len(max_array) > 1:
23           parent_max = 0
24           parent_array = []
25           for k, v in parent_attr["num"].items():
26               if v > parent_max:
27                   parent_array = []
28                   parent_array.append(k)
29                   parent_max = v
30               elif v == parent_max:
31                   parent_array.append(k)
32           #if two or more label has same value in parent node, find root node's largest label
33           if len(parent_array) > 1:
34               root_array = []
35               root_val = 0
36               for k, v in root_attr["num"].items():
37                   if v > root_val:
38                       root_array = []
39                       root_array.append(k)
40                       root_val = v
41                   elif v == root_val:
42                       root_array.append(k)
43               return root_array[0]
44           else :
45               return parent_array[0]
46       else:
47           return max_array[0]
```

getMaxLabel(parent_attr, this_attr, root_attr) is the function follows the majority rule. use the dictionary whose key is 'num'. and find max value in num. If max value is duplicate, find same thing in parent node's 'num', and if parent node's max value is duplicate, find in root node.
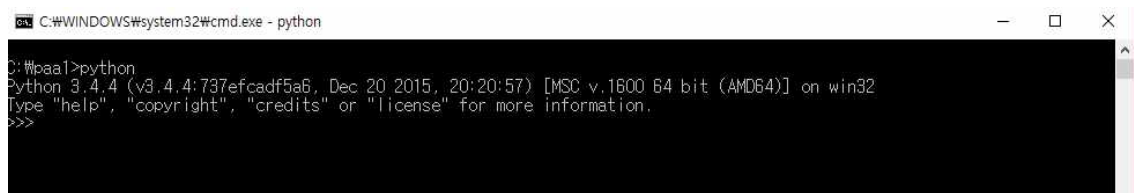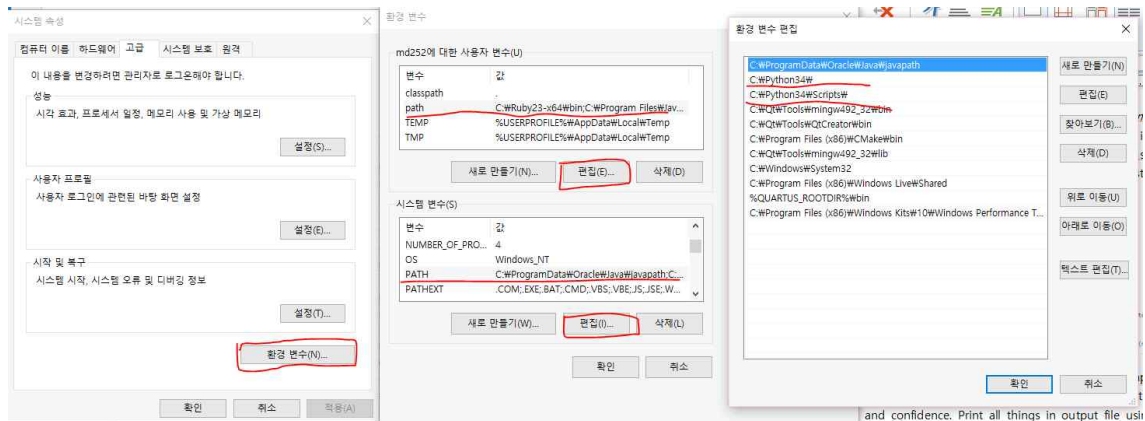
# 3. Instructions for compiling this code

This code is written in Python3 and tested in Python 3.4.4. Since this code has **log2**(to calculate gain ratio), and it only exist in Python3 math module, to run this code, we need python3 interpreter.

We can get python 3.4.4 in (https://www.python.org/downloads/release/python-344/)

If python2 is already installed in PC, change path Python2x to Python34 in advanced system settings – environment variables.



After environment variables setting, we can get python3 default in cmd.

In code directory(python file is in /project_dt), we can run this program on cmd "python dt.py [training data file name] [test set file name] [output file name]".
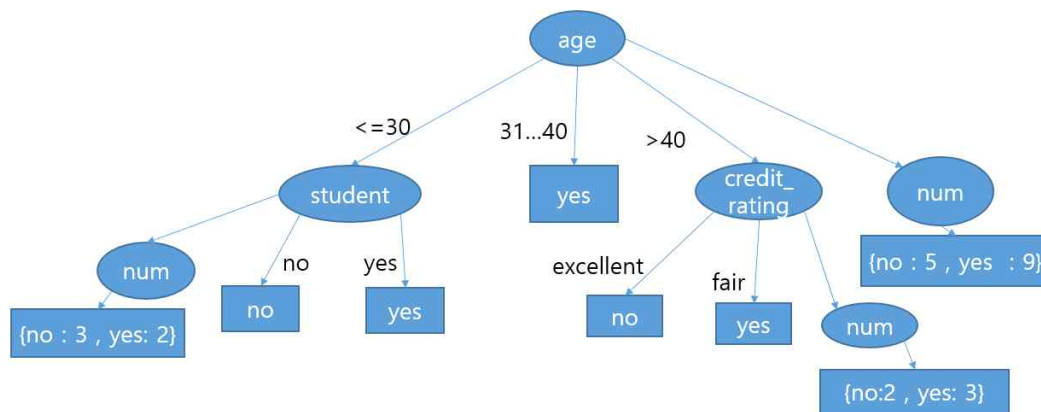


continue in next picture

## 4. Other Specifications

I implemented this tree with all attribute node should have the number of its class labels as child of themselves. It was used to classify by majority vote which can't reach leaf node. So tree looks like below.

In dictionary shaped tree, they all have key named 'num' and it contains number of class labels the node has.



tree generated by dt_train.txt



the accuracy of this program