데이터사이언스

# Decision Tree implement assignment

2013011800 구장회

# 순서
1. data structure - tree class , node class
2. summary of algorithm for crucial method
3. summary of rest algorithm
4. instruction for compile


# data structure

**[tree class]**

```
14  /**
15   * Tree class
16   *
17   * List<Map<List<String>, String>> data        : initial list of tuples by reading training data
18   * List<Map<List<String>, String>> tdata       : initial list of tuples by reading test data
19   * List<Map<List<String>, String>> rdata       : list of tuples for result data
20   *
21   * Map<String, List<String>> attributes        : pair of each (attribute , its possible value)
22   *
23   * List<String> label                          : list of class-label
24   * List<String> attribute                      : list of attribute name
25   *
26   * String trainingset, testset, resultset      : path to each files
27   * String classname                            : class name
28   * int numattr                                 : number of attributes
29   * Node root                                   : root node of decition tree
30   *
31   */
32  public class Tree {
33      List<Map<List<String>, String>> data, tdata, rdata;
34      Map<String, List<String>> attributes;
35      List<String> label;
36      List<String> attribute;
37      String trainingset, testset, resultset, classname = null;
38      int numattr = 0;
39      Node root;
40
41
```

**[node class]**

```
603  /**
604   * Node class
605   *
606   * String label         : node's class-label , null at default
607   * String decision      : node's split attribute , null at default
608   * String value         : value of parent node's decision attribute
609   * Node child           : its child node
610   * List<Map<List<String>, String>> list : list of tuples belong to node
611   *
612   */
613  public class Node {
614      String label, decision, value;
615      List<Node> child;
616      List<Map<List<String>, String>> list;
617
618      /* constructor */
619      public Node(List<Map<List<String>, String>> _list) {
620          this.list = _list;
621          this.child = new ArrayList<Node>();
622      }
623
624      /* add child node */
625      void addBranch(Node _child) {
626          this.child.add(_child);
627      }
628  }
629 }
```

# summary of algorithm for crucial method

[main method]

```
579    /**
580     * main method
581     *
582     * store all parameter as path to files
583     * and call work method after create decision-tree instance
584     *
585     * @param       path to training-dataset file
586     * @param       path to test-dataset file
587     * @param       path to result-dataset file
588     * @return      void
589     *
590     */
591    public static void main(String[] args) {
592        String trainingset, testset, result;
593        if (args.length != 3) {
594            System.out.println("args error!!!");
595        } else {
596            trainingset = args[0];
597            testset = args[1];
598            result = args[2];
599            new Tree().work(trainingset, testset, result);
600        }
601    }
602
```

## [work() in tree class]
- read training data
- construct decision tree
- read test data
- tree test from test data
- write result data

```
137   /**
138    * main work in DecisionTree algorithm
139    *
140    * read training-data and store all of them(all data-label pairs)
141    * construct all tree recursively and store root node
142    * read test-data and and store all of them(all data)
143    * test decision tree by test-data
144    * and write result to path of result-data
145    *
146    * @param      path to training-dataset.txt
147    * @param      path to test-dataset.txt
148    * @param      path to result-dataset.txt
149    * @return     void
150    *
151    */
152   void work(String _trainingset, String _testset, String _result) {
153       trainingset = _trainingset;
154       testset = _testset;
155       resultset = _result;
156       readTrainingData(trainingset);
157       root = id3(data, attributes, attribute);
158       readTestData(testset);
159       TreeTest();
160       writeResult(resultset);
161   }
162
```

# [id3() in tree class]

```java
/**
 * crucial method for recursively constructing decision-tree(specially each node)
 *
 * create new node from list of tuples
 *
 * if # of attribute is zero, then return node after majority-vote and store its label.
 * if list of tuples are uniformly homogeneous, then return node after store its label.
 *
 * choose best attribute for split from (information gain or gain ratio or gini index).
 * it is done simply by removing comment!!!
 *
 * after choosing best attribute, store it as decision attribute for node.
 *
 * for each value of best attribute,
 * create sub list of tuples which has some value at best attribute
 * create set pair of (attribute-its possible value) by removing best attribute at all tuples in sub list
 * create list of attribute by removing best attribute
 *
 * finally create child node by calling same method recursively
 * and store child node's value same as each value of best attribute
 * and make pair parent-child node relationship
 *
 * @param      list of tuples
 * @param      set pair of (attribute-its possible value)
 * @param      list of attribute
 *
 * @return     node which is constructed
 *
 */
Node id3(List<Map<List<String>, String>> tuples, Map<String, List<String>> _attributes, List<String> _attribute) {
    List<Map<List<String>, String>> _tuples = new ArrayList<Map<List<String>, String>>(tuples);
    Node node = new Node(_tuples);

    if (_attribute.size() == 0) {
        node.label = majorityVote(_tuples);
        return node;
    }

    Map<String, Integer> dictionary = summarizeExamples(_tuples);
    for (String key : dictionary.keySet()) {
        if (dictionary.get(key) == _tuples.size()) {
            node.label = key;
            return node;
        }
    }

    //String bestattr = getBestAttr(getInfoGain(_tuples, _attributes, _attribute));
    //String bestattr = getBestAttr(getGainRatio(_tuples, _attributes, _attribute));
    String bestattr = getBestAttr(getGiniIndex(_tuples, _attributes, _attribute));

    int idx = _attribute.indexOf(bestattr);
    node.decision = bestattr;

    for (String value : _attributes.get(bestattr)) {
        int size = 0;

        for (Map<List<String>, String> tuple : _tuples) {
            for (List<String> key : tuple.keySet()) {
                size = Math.max(size, key.size());
            }
        }
        if (size == 0) {
            node.label = majorityVote(_tuples);
            return node;
        }
```

```java
            List<Map<List<String>, String>> subexamples = new ArrayList<>();
            for (Map<List<String>, String> tuple : _tuples) {
                for (List<String> line : tuple.keySet()) {
                    if(line.size()>idx) {
                        if (line.size() == size && line.get(idx).equalsIgnoreCase(value)) {
                            subexamples.add(tuple);
                        }
                    }
                }
            }
            if (subexamples.size() != 0) {
                Map<String, List<String>> subattributes = new HashMap<String, List<String>>(_attributes);
                subattributes.remove(bestattr);

                List<String> subattribute = new ArrayList<String>(_attribute);
                subattribute.remove(idx);

                for (Map<List<String>, String> tuple : subexamples) {
                    for (List<String> line : tuple.keySet()) {
                        line.remove(idx);
                    }
                }
                Node child = id3(subexamples, subattributes, subattribute);
                child.value = value;
                node.addBranch(child);
            }
        }
        return node;
    }
```

# summary of rest algorithm

**[readTrainingData & readTestData]**

```
42    /**
43     * read training-dataset file
44     * store all tuples(all data-label pair)
45     * store classname & all possible class-label
46     * store all attibutes & all possible value of each attribute
47     *
48     * @param    path to training-dataset.txt
49     * @return   void
50     * @exception FileNotFoundException e
51     * @exception IOException e
52     */
53    void readTrainingData(String _trainingset) {
99
100    /**
101     * read test-dataset file
102     * store all tuples(all data)
103     *
104     * @param    path to test-dataset.txt
105     * @return   void
106     * @exception FileNotFoundException e
107     * @exception IOException e
108     */
109    void readTestData(String _testset) {
136
137    /**
```

**[writeResult & TreeTest]**

```java
163    /**
164     * write result to path of result-data
165     *
166     * add induced class-label at each tuples of test-data
167     *
168     * @param     path to result-dataset.txt
169     * @return    void
170     * @exception FileNotFoundException e
171     */
172    void writeResult(String resultset) {
193
194    /**
195     * construct result-data by adding class-label at each tuples
196     * in test-data.
197     * each class-label is induced by calling findLabel method
198     *
199     * @param     void
200     * @return    void
201     */
202    void TreeTest() {
203        rdata = new ArrayList<>();
204        for (Map<List<String>, String> line : tdata) {
205            for (List<String> tuple : line.keySet()) {
206                String tlabel = findLabel(tuple, root);
207                Map<List<String>, String> temp = new HashMap<>();
208                temp.put(tuple, tlabel);
209                rdata.add(temp);
210            }
211        }
212    }
213
```

**[findLabel]**

```java
214    /**
215     * inducing class-label at each tuples
216     *
217     * if the node has no label
218     *
219     * find childnode which has value as its parent node's decision attribute.
220     * and call same method recursively simply change node to childnode.
221     * if there is no childnode which has same value then, do majority vote.
222     *
223     * @param      a tuple which has no class-label
224     * @param      a node at decision tree
225     * @return     induced class-label
226     *
227    */
228    String findLabel(List<String> tuple, Node node) {
229        boolean flag = false;
230        if (node.label == null) {
231            for (Node childnode : node.child) {
232                if (childnode.value.equalsIgnoreCase(tuple.get(attribute.indexOf(node.decision)))) {
233                    flag = true;
234                    return findLabel(tuple, childnode);
235                }
236            }
237            if (!flag) {
238                return majorityVote(node.list);
239            }
240        }
241        return node.label;
242    }
243
```

**[summarizeExamples]**

```java
244    /**
245     * construct pair of (class-label , its count)
246     * by searching all tuples and count numbers of class-label
247     * for checking whether tuples is homogenious or not
248     *
249     * @param     list of tuples
250     * @return    pair of (class-label , count)
251     *
252     */
253    Map<String, Integer> summarizeExamples(List<Map<List<String>, String>> examples) {
254        List<Map<List<String>, String>> temp = new ArrayList<Map<List<String>, String>>(examples);
255        Map<String, Integer> ret = new HashMap<>();
256        for (String key : label) {
257            ret.put(key, 0);
258        }
259        for (int i = 0; i < temp.size(); i++) {
260            String fkey = (String) temp.get(i).values().toArray()[0];
261            Integer count = ret.get(fkey);
262            ret.put(fkey, (count + 1));
263        }
264        return ret;
265    }
266
267
```

**[majority-vote]**

```java
/**
 * search all tuples and count numbers of class-label
 * and return major class-label
 *
 * @param     list of tuples
 * @return    class-label which is major in the tuples
 *
 */
String majorityVote(List<Map<List<String>, String>> _tuples) {
    List<Map<List<String>, String>> tmp = new ArrayList<Map<List<String>, String>>(_tuples);

    String ret = null;
    int[] temp = new int[label.size()];
    for (int i = 0; i < tmp.size(); i++) {
        String fkey = (String) tmp.get(i).values().toArray()[0];
        temp[label.indexOf(fkey)]++;
    }
    int max = -1;
    for (int i = 0; i < label.size(); i++) {
        if (temp[i] > max) {
            max = temp[i];
            ret = label.get(i);
        }
    }
    return ret;
}
```

**[getBestAttr]**

```java
/**
 *
 * @param     pair of (attribute - its gain for any metric)
 * @return    best attribute for split node
 *
 */
String getBestAttr(Map<String, Double> infogain) {
    String ret = "";
    double min = 100.0;
    for (String key : infogain.keySet()) {
        if (infogain.get(key) < min) {
            min = infogain.get(key);
            ret = key;
        }
    }
    return ret;
}
```

**[getGiniIndex]**

```java
    /**
     * calculating gini-index
     *
     * @param      list of tuples
     * @param      pair of (attribute and its possible value)
     * @param      list of attributes
     * @return     pair of (attribute - its gain) for gini-index
     *
     */
    Map<String,Double> getGiniIndex(List<Map<List<String>, String>> _tuples, Map<String, List<String>> _attributes,
            List<String> _attribute) {
        Map<String, Double> gainstore = new HashMap<>();
        int arr[] = new int[label.size()];
        for (int i = 0; i < _tuples.size(); i++) {
            String fkey = (String) _tuples.get(i).values().toArray()[0];
            arr[label.indexOf(fkey)]++;
        }
        for (String attr : _attributes.keySet()) {
            int idx = _attribute.indexOf(attr);
            double global = _tuples.size();
            double gain = 0;
            double info_a = 0;
            for (String value : _attributes.get(attr)) {
                double local = 0;
                int[] Arr = new int[label.size()];
                for (int i = 0; i < _tuples.size(); i++) {
                    for (List<String> key : _tuples.get(i).keySet()) {
                        if (key.get(idx).equalsIgnoreCase(value)) {
                            local++;
                            String fkey = (String) _tuples.get(i).values().toArray()[0];
                            Arr[label.indexOf(fkey)]++;
                        }
                    }
                }
                info_a += local / global * gini(Arr);
            }
            gain = gini(arr) - info_a;
            gain = Math.round(gain * 1000) / 1000.0;
            gainstore.put(attr, gain);
        }
        return gainstore;
    }
```

**[gini]**

```java
/**
 * calculating gini-index value at each integer array
 *
 * @param   integer array
 * @return  its value of gini-index
 *
 */
double gini(int[] arr) {
    double ret = 1;
    double sum = 0;
    for (int i = 0; i < arr.length; i++) {
        sum += arr[i];
    }

    for (int i = 0; i < arr.length; i++) {
        double p = arr[i] / sum;
        ret -= Math.pow(p, 2);
    }
    return Math.round(ret * 1000) / 1000.0;
}
```

**[getGainRatio]**

```java
/**
 * calculating gain-ratio
 *
 * @param   list of tuples
 * @param   pair of (attribute and its possible value)
 * @param   list of attributes
 * @return  pair of (attribute - its gain) for gain-ratio
 *
 */
Map<String, Double> getGainRatio(List<Map<List<String>, String>> _tuples, Map<String, List<String>> _attributes,
        List<String> _attribute) {
    Map<String, Double> gainstore = new HashMap<>();
    gainstore = getInfoGain(_tuples, _attributes, _attribute);

    for (String attr : _attributes.keySet()) {
        int idx = _attribute.indexOf(attr);
        int[] arr = new int[_attributes.get(attr).size()];
        int index = 0;
        for (String value : _attributes.get(attr)) {
            int local = 0;
            for (int i = 0; i < _tuples.size(); i++) {
                for (List<String> key : _tuples.get(i).keySet()) {
                    if (key.get(idx).equalsIgnoreCase(value)) {
                        local++;
                    }
                }
            }
            arr[index]=local;
            index++;
        }
        double gain = gainstore.get(attr);
        gainstore.put(attr, gain/entropy(arr));
    }
    return gainstore;
}
```

## [getInfoGain]

```java
 * calculating info-gain
 *
 * @param      list of tuples
 * @param      pair of (attribute and its possible value)
 * @param      list of attributes
 * @return     pair of (attribute - its gain) for info-gain
 *
 */
Map<String, Double> getInfoGain(List<Map<List<String>, String>> _tuples, Map<String, List<String>> _attributes,
        List<String> _attribute) {

    Map<String, Double> gainstore = new HashMap<>();
    int infoDArray[] = new int[label.size()];
    for (int i = 0; i < _tuples.size(); i++) {
        String fkey = (String) _tuples.get(i).values().toArray()[0];
        infoDArray[label.indexOf(fkey)]++;
    }
    for (String attr : _attributes.keySet()) {
        int idx = _attribute.indexOf(attr);
        double global = _tuples.size();
        double gain = 0;
        double info_a = 0;
        for (String value : _attributes.get(attr)) {
            double local = 0;
            int[] infoADArray = new int[label.size()];
            for (int i = 0; i < _tuples.size(); i++) {
                for (List<String> key : _tuples.get(i).keySet()) {
                    if (key.get(idx).equalsIgnoreCase(value)) {
                        local++;
                        String fkey = (String) _tuples.get(i).values().toArray()[0];
                        infoADArray[label.indexOf(fkey)]++;
                    }
                }
            }
            info_a += local / global * entropy(infoADArray);
        }
        gain = entropy(infoDArray) - info_a;
        gain = Math.round(gain * 1000) / 1000.0;
        gainstore.put(attr, gain);
    }
    return gainstore;
}
```
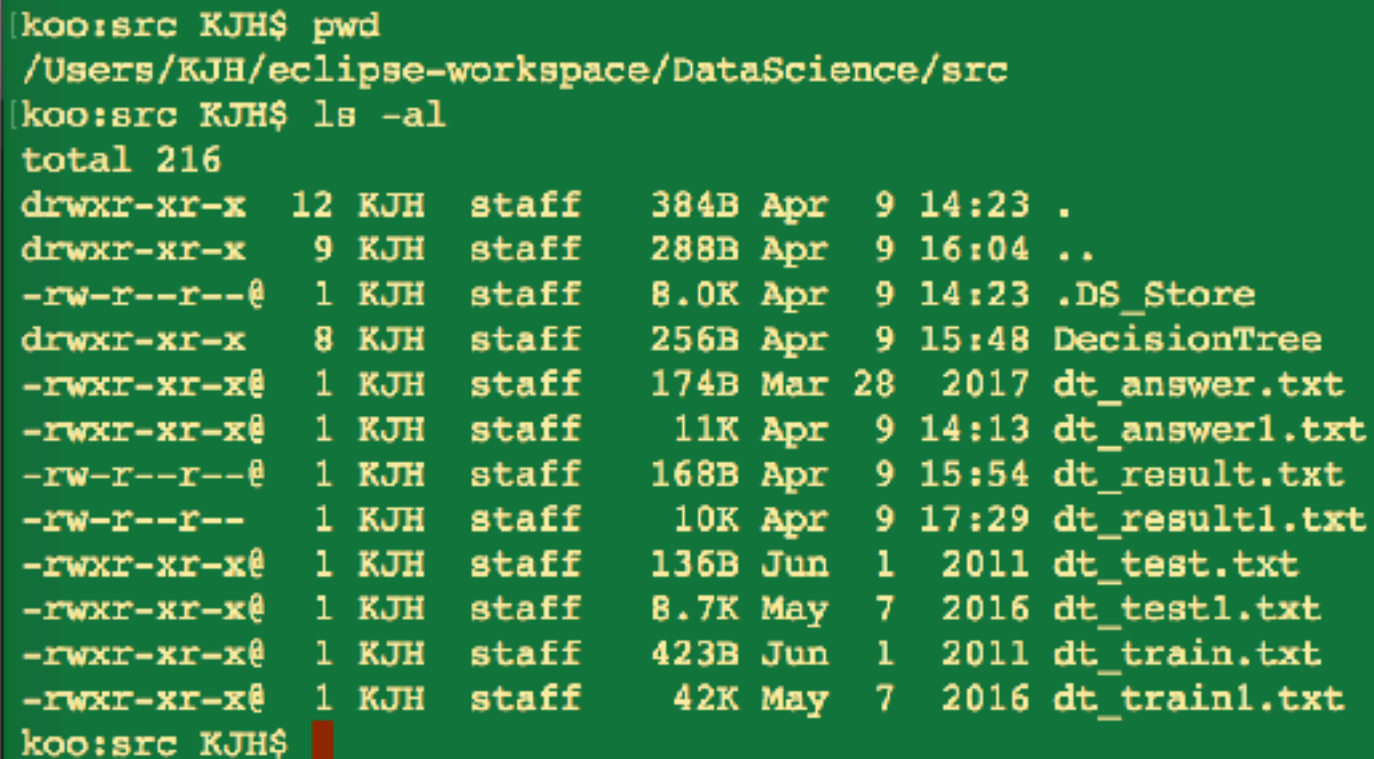
**[entropy]**

```java
/**
 * calculating entropy value at each integer array
 *
 * @param      integer array
 * @return     its value of info gain
 *
 */
double entropy(int[] arr) {
    double ret = 0;
    double sum = 0;
    for (int i = 0; i < arr.length; i++) {
        sum += arr[i];
    }
    for (int i = 0; i < arr.length; i++) {
        double p = arr[i] / sum;
        if (p != 0)
            ret -= p * Math.log(p) / Math.log(2);
    }
    return Math.round(ret * 1000) / 1000.0;
}
```

# instruction for complile

**[Enviorment]**

- OS : Mac OS
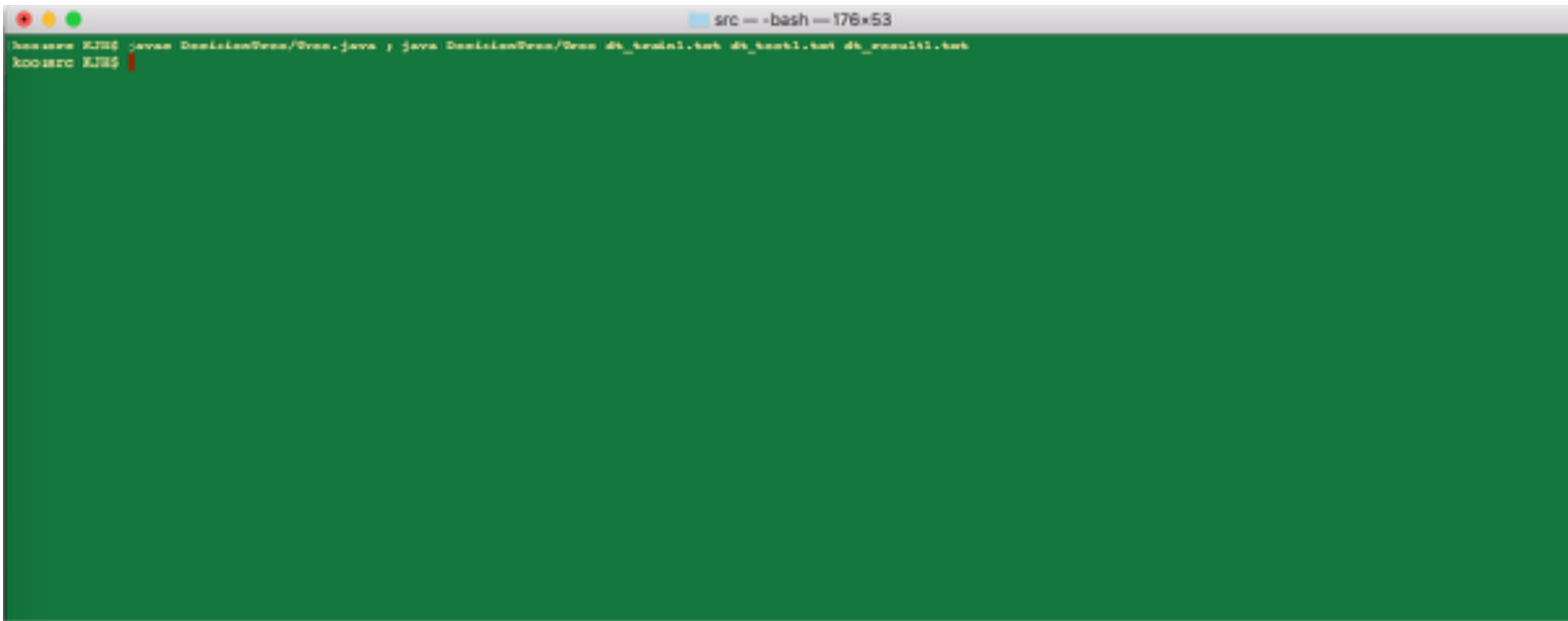- Language : Java

**[Screenshot - 실행 전]**

```
[koo:src KJH$ pwd
/Users/KJH/eclipse-workspace/DataScience/src
[koo:src KJH$ ls -al
total 216
drwxr-xr-x  12 KJH  staff   384B Apr  9 14:23 .
drwxr-xr-x   9 KJH  staff   288B Apr  9 16:04 ..
-rw-r--r--@  1 KJH  staff   8.0K Apr  9 14:23 .DS_Store
drwxr-xr-x   8 KJH  staff   256B Apr  9 15:48 DecisionTree
-rwxr-xr-x@  1 KJH  staff   174B Mar 28  2017 dt_answer.txt
-rwxr-xr-x@  1 KJH  staff    11K Apr  9 14:13 dt_answer1.txt
-rw-r--r--@  1 KJH  staff   168B Apr  9 15:54 dt_result.txt
-rw-r--r--   1 KJH  staff    10K Apr  9 17:29 dt_result1.txt
-rwxr-xr-x@  1 KJH  staff   136B Jun  1  2011 dt_test.txt
-rwxr-xr-x@  1 KJH  staff   8.7K May  7  2016 dt_test1.txt
-rwxr-xr-x@  1 KJH  staff   423B Jun  1  2011 dt_train.txt
-rwxr-xr-x@  1 KJH  staff    42K May  7  2016 dt_train1.txt
koo:src KJH$
```

**[Screenshot - 실행 방법]**

compile: $javac DecisionTree/Tree.java
execute : $java DecisionTree/Tree dt_train1.txt dt_test1.tst dt_result1.txt