

# E10 Decision Tree

---

18340215 张天祎

2020 年 11 月 25 日

## 目录

<b>1</b>	<b>Datasets</b>	<b>2</b>
<b>2</b>	<b>Decision Tree</b>	<b>3</b>
2.1	ID3 . . . . .	3
2.2	C4.5 and CART . . . . .	4
<b>3</b>	<b>Tasks</b>	<b>5</b>
<b>4</b>	<b>Codes and Results</b>	<b>5</b>

# 1 Datasets

The UCI dataset (<http://archive.ics.uci.edu/ml/index.php>) is the most widely used dataset for machine learning. If you are interested in other datasets in other areas, you can refer to <https://www.zhihu.com/question/63383992/answer/222718972>.

Today's experiment is conducted with the **Adult Data Set** which can be found in <http://archive.ics.uci.edu/ml/datasets/Adult>.

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	48842	<b>Area:</b>	Social
<b>Attribute Characteristics:</b>	Categorical, Integer	<b>Number of Attributes:</b>	14	<b>Date Donated</b>	1996-05-01
<b>Associated Tasks:</b>	Classification	<b>Missing Values?</b>	Yes	<b>Number of Web Hits:</b>	1305515

You can also find 3 related files in the current folder, `adult.name` is the description of **Adult Data Set**, `adult.data` is the training set, and `adult.test` is the testing set. There are 14 attributes in this dataset:

>50K, <=50K.

1. age: continuous.
2. workclass: Private, Self-emp-not-inc, Self-emp-inc, Federal-gov, Local-gov, State-gov, Without-pay, Never-worked.
3. fnlwgt: continuous.
4. education: Bachelors, Some-college, 11th, HS-grad, Prof-school, Assoc-acdm, Assoc-voc, 9th, 7th-8th, 12th, Masters, 5. 1st-4th, 10th, Doctorate, 5th-6th, Preschool.
5. education-num: continuous.
6. marital-status: Married-civ-spouse, Divorced, Never-married, Separated, Widowed, Married-spouse-absent, Married-AF-spouse.
7. occupation: Tech-support, Craft-repair, Other-service, Sales, Exec-managerial, Prof-specialty, Handlers-cleaners, Machine-op-inspct, Adm-clerical, Farming-fishing, Transport-moving, Priv-house-serv, Protective-serv, Armed-Forces.
8. relationship: Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried.
9. race: White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black.
10. sex: Female, Male.

11. capital-gain: continuous.
12. capital-loss: continuous.
13. hours-per-week: continuous.
14. native-country: United-States , Cambodia , England , Puerto-Rico , Canada , Germany , Outlying-US(Guam-USVI-etc) , India , Japan , Greece , South , China , Cuba , Iran , Honduras , Philippines , Italy , Poland , Jamaica , Vietnam , Mexico , Portugal , Ireland , France , Dominican-Republic , Laos , Ecuador , Taiwan , Haiti , Columbia , Hungary , Guatemala , Nicaragua , Scotland , Thailand , Yugoslavia , El-Salvador , Trinidad&Tobago , Peru , Hong , Holand-Netherlands .

**Prediction task is to determine whether a person makes over 50K a year.**

## 2 Decision Tree

### 2.1 ID3

ID3 (Iterative Dichotomiser 3) was developed in 1986 by Ross Quinlan. The algorithm creates a multiway tree, finding for each node (i.e. in a greedy manner) the categorical feature that will yield the largest information gain for categorical targets. Trees are grown to their maximum size and then a pruning step is usually applied to improve the ability of the tree to generalise to unseen data.

#### **ID3 Algorithm:**

1. Begins with the original set  $S$  as the root node.
2. Calculate the entropy of every attribute  $a$  of the data set  $S$ .
3. Partition the set  $S$  into subsets using the attribute for which the resulting entropy after splitting is minimized; or, equivalently, information gain is maximum.
4. Make a decision tree node containing that attribute.
5. Recur on subsets using remaining attributes.

#### **Recursion on a subset may stop in one of these cases:**

- every element in the subset belongs to the same class; in which case the node is turned into a leaf node and labelled with the class of the examples.
- there are no more attributes to be selected, but the examples still do not belong to the same class. In this case, the node is made a leaf node and labelled with the most common class of the examples in the subset.
- there are no examples in the subset, which happens when no example in the parent set was found to match a specific value of the selected attribute.

**ID3 shortcomings:**

- ID3 does not guarantee an optimal solution.
- ID3 can overfit the training data.
- ID3 is harder to use on continuous data.

**Entropy:**

Entropy  $H(S)$  is a measure of the amount of uncertainty in the set  $S$ .

$$H(S) = \sum_{x \in X} -p(x) \log_2 p(x)$$

where

- $S$  is the current dataset for which entropy is being calculated
- $X$  is the set of classes in  $S$
- $p(x)$  is the proportion of the number of elements in class  $x$  to the number of elements in set  $S$ .

**Information gain:**

Information gain  $IG(A)$  is the measure of the difference in entropy from before to after the set  $S$  is split on an attribute  $A$ . In other words, how much uncertainty in  $S$  was reduced after splitting set  $S$  on attribute  $A$ .

$$IG(S, A) = H(S) - \sum_{t \in T} p(t)H(t) = H(S) - H(S | A)$$

where

- $H(S)$  is the entropy of set  $S$
- $T$  is the subsets created from splitting set  $S$  by attribute  $A$  such that  $S = \cup_{t \in T} t$
- $p(t)$  is the proportion of the number of elements in  $t$  to the number of elements in set  $S$
- $H(t)$  is the entropy of subset  $t$ .

**2.2 C4.5 and CART**

C4.5 is the successor to ID3 and removed the restriction that features must be categorical by dynamically defining a discrete attribute (based on numerical variables) that partitions the continuous attribute value into a discrete set of intervals. C4.5 converts the trained trees (i.e. the output of the ID3 algorithm) into sets of if-then rules. The accuracy of each rule is then evaluated to determine the order in which they should be applied. Pruning is done by removing a rule's precondition if the accuracy of the rule improves without it.

C5.0 is Quinlan's latest version release under a proprietary license. It uses less memory and builds smaller rulesets than C4.5 while being more accurate.

CART (Classification and Regression Trees) is very similar to C4.5, but it differs in that it supports numerical target variables (regression) and does not compute rule sets. CART constructs binary trees using the feature and threshold that yield the largest information gain at each node.

### 3 Tasks

- Given the training dataset `adult.data` and the testing dataset `adult.test`, please accomplish the prediction task to determine whether a person makes over 50K a year in `adult.test` by using ID3 (or C4.5, CART) algorithm (C++ or Python), and compute the accuracy.
  1. You can process the continuous data with **bi-partition** method.
  2. You can use prepruning or postpruning to avoid the overfitting problem.
  3. You can assign probability weights to solve the missing attributes (data) problem.
- Please finish the experimental report named `E10_YourNumber.pdf`, and send it to `ai_2020@foxmail.com`

### 4 Codes and Results

```
1 import math
2 import copy
3 import random
4 from collections import Counter
5
6 # 超参数
7 MAX_LEFT_DEPTH = 6
8
9 # 读入并格式化数据
10 def load(file_name):
11     f = open(file_name, 'rt')
12     adult_list = f.readlines()
13     adult_msg = []
14     for i in adult_list:
15         curr = i.split(',')
16         for j in range(1, len(curr)):
17             curr[j] = curr[j][1:] # 去空格
18         adult_msg.append(curr)
19     # 检查是否有不合规的数据, 数据只有最后一行多了一个回车, 测试集还有第一行不是数据
20     for i in adult_msg:
```

```

21         if len(i) != 15:
22             # print(i)
23             adult_msg.remove(i)
24     return [i[:-1] for i in adult_msg], [i[-1][-1] for i in adult_msg] # 这里标
        签删去了回车符
25
26 class node:
27     def __init__(self, leaf=False, label=None, attribute=None, son=[]):
28         self.leaf = leaf
29         self.attribute = attribute
30         self.label = label
31         self.son = son
32
33     def add_son(self, son):
34         if son[1] != None:
35             self.son.append(son)
36
37     # for test
38     def printinf(self):
39         if self.leaf:
40             # print("leaf label: "+self.label, end=" ")
41             print("leaf", end="_")
42         else:
43             print(self.attribute, end="_")
44
45 class DecisionTree:
46     def __init__(self, train_data, train_label):
47         # 预处理连续值数据
48         self.con_attlist = [0, 2, 4, 10, 11, 12]
49         # 预处理离散值数据
50         tem_dis_list = [1, 3, 5, 6, 7, 8, 9, 13]
51         self.dis_attlist = {}
52         for j in tem_dis_list:
53             self.dis_attlist[j] = list(set([i[j] for i in train_data]))
54         # 预处理标签
55         self.labelist = list(set(train_label))
56
57     def treeLearn(self, data, label, attributes, entropy=1, option='TD3'):

```

```

58     # 数据集是否空
59     if data == []:
60         # 数据缺失时，需要分配叶节点防止出现决策时找不到的情况
61         # 策略一，随机，TD3:0.80
62         return node(True, self.labelist[random.randint(0, len(self.labelist)
63             - 1)])
64         # 策略二，硬分配一个，TD3:0.79~0.81
65         # return node(True, self.labelist[0])
66     # 是否全为同一标签
67     same = True
68     for i in label[1:]:
69         if i != label[0]:
70             same = False
71             break
72     if same:
73         return node(True, label[0])
74     # 是否没有属性
75     if option == 'C4.5':
76         # if len(attributes) < 7: # 同ID3
77         if len(attributes) < MAX_LEFT_DEPTH:
78             return node(True, Counter(label).most_common(1)[0][0])
79     else: # TD3变为是否无离散值属性:
80         flag = False
81         for i in attributes:
82             if i in self.dis_attlist:
83                 flag = True
84             if not flag:
85                 return node(True, Counter(label).most_common(1)[0][0])
86     # 一般情况
87     Entropys = [(i, self.calEntropy(data, label, i, entropy, option)) for i
88         in attributes]
89     if option == 'C4.5': # C4.5 比较增益率
90         pickatt, entropy = max(Entropys, key=lambda x: x[1][2])
91     else: # 比较信息熵
92         pickatt, entropy = min(Entropys, key=lambda x: x[1][0])
93     mid_val = entropy[1]
94     curr_node = copy.deepcopy(node(False, attribute=pickatt)) # 这里有个，需要
95     深复制bug

```

```

93     par_map = self.partition(data, label, pickatt, mid_val)
94     tem_attributes = copy.deepcopy(attributes)
95     tem_attributes.remove(pickatt)
96     for i in par_map:
97         curr_node.add_son((i, self.treeLearn(par_map[i][0], par_map[i][1],
98             tem_attributes, entropy[0], option)))
99
100     return curr_node
101
102 def calEntropy(self, data, label, attribute, fa_entropy, option='TD3'):
103     att_data = [i[attribute] for i in data]
104     length = len(data)
105     true_label = self.labelist[0]
106     sum = 0
107     IV = 0
108     if attribute in self.con_attlist:
109         # print(option)
110         if option=='C4.5':
111             # 二分法
112             att_data = [int(i) for i in att_data]
113             max_val = max(att_data)
114             min_val = min_val
115             for i in att_data:
116                 # 考虑到该数据集中有大量无用的，做了一点优化，效果不大0
117                 if i != 0 and i < min_val:
118                     min_val = i
119             mid_val = (max_val+min_val)//2
120             true_count = false_count = [0, 0]
121             for i in zip(att_data, label):
122                 ind = i[0] > mid_val
123                 if i[1] == true_label:
124                     true_count[ind] += 1
125                 else:
126                     false_count[ind] += 1
127             for i in range(2):
128                 ci = true_count[i] + false_count[i]
129                 pi = ci / length
130                 if true_count[i]: # 防止出现数学错，这里不会影响熵的结果

```



```

129         px = true_count[i] / ci
130         sum += pi * (-px * math.log(px))
131     IV -= pi * math.log2(pi) if pi != 0 else 0
132     gain_ratio = (fa_entropy - sum) / IV if IV != 0 else 2 # IV 小是
        较好的分类
133     return sum, mid_val, gain_ratio # 2 这里
        只是一个相对较大的数
134 # td3 下连续值属性不考虑, 直接返回一个较大的熵 2
135     return 2, None, None
136 else:
137     true_count = dict([(i, 0) for i in self.dis_attlist[attribute]])
138     false_count = dict([(i, 0) for i in self.dis_attlist[attribute]])
139     for i in zip(att_data, label):
140         if i[1] == true_label:
141             true_count[i[0]] += 1
142         else:
143             false_count[i[0]] += 1
144     for i in self.dis_attlist[attribute]:
145         ci = true_count[i] + false_count[i]
146         pi = ci / length
147         if true_count[i]: # 防止出现数学错, 这里不会影响熵的结果
148             px = true_count[i] / ci
149             sum += pi * (-px * math.log2(px))
150         IV -= pi * math.log2(pi) if pi != 0 else 0
151         gain_ratio = (fa_entropy - sum) / IV if IV != 0 else 2
152     return sum, None, gain_ratio
153
154 def partition(self, data, label, attribute, mid_val):
155     par_map = {}
156     if attribute in self.dis_attlist:
157         for i in self.dis_attlist[attribute]:
158             par_map[i] = [[], []]
159         for i in range(len(data)):
160             par_map[data[i][attribute]][0].append(data[i])
161             par_map[data[i][attribute]][1].append(label[i])
162     else:
163         ind1 = ">" + str(mid_val)
164         ind2 = "<=" + str(mid_val)

```

```

165         par_map[ind1] = [[], []]
166         par_map[ind2] = [[], []]
167         for i in range(len(data)):
168             ind = ind1 if int(data[i][attribute])>mid_val else ind2
169             par_map[ind][0].append(data[i])
170             par_map[ind][1].append(label[i])
171         return par_map
172
173 # 没有考虑测试集出现训练集中未出现值的情况
174 def test(self, root, test_data):
175     rst = []
176     for i in test_data:
177         curr = root
178         while(not curr.leaf):
179             att = curr.attribute
180             val = i[att]
181             if att in self.dis_attlist:
182                 for j in curr.son:
183                     if j[0] == val:
184                         curr = j[1]
185                     break
186             else:
187                 num = int(curr.son[0][0][1:])
188                 if int(val)>num:
189                     curr = curr.son[0][1]
190                 else:
191                     curr = curr.son[1][1]
192             rst.append(curr.label)
193     return rst
194
195 # for test
196 def printinf(self, root):
197     print("root:", end=" ")
198     que = []
199     que.append(root)
200     while(len(que)):
201         tem = []

```

```

202         while(len(que)):
203             curr = que[0]
204             que.pop(0)
205             curr.printinf()
206             if not curr.leaf:
207                 # print(curr.attribute , end=" ")
208                 for i in curr.son:
209                     tem.append(i[1])
210             # else:
211                 # print("leaf",end=" ")
212         que = tem
213         print("")
214
215 def main():
216     # pre_train
217     train_data, train_label = load('adult.data')
218     test_data, test_label = load('adult.test')
219     test_label = [i[:-1] for i in test_label] # test_label 后面多了一个 '.'
220     tree = DecisionTree(train_data, train_label)
221
222     # train
223     root_1 = tree.treeLearn(train_data, train_label, list(range(14)))
224     root_2 = tree.treeLearn(train_data, train_label, list(range(14)),option='C4
        .5')
225
226     # test
227     # DecisionTree.printinf(root)
228     rst_label1 = tree.test(root_1, test_data)
229     rst_label2 = tree.test(root_2, test_data)
230     testlen = len(test_label)
231     count1 = count2 = 0
232     for i in range(testlen):
233         if rst_label1[i] == test_label[i]:
234             count1 += 1
235         if rst_label2[i] == test_label[i]:
236             count2 += 1
237     print("\u00ID3\u00accuracy:\u00"+str(count1/testlen*100)+"%")

```

```
238     print("\nC4.5 accuracy: "+str(count2/testlen*100)+"%")
239     print("with MAX_DEPTH:", 14 - MAX_LEFT_DEPTH)
240
241 if __name__ == "__main__":
242     main()
```

结果如下：

```
ID3 accuracy: 80.38204041520791%
C4.5 accuracy: 81.46919722375775%
with MAX_DEPTH: 8
```

- ID3 只使用了离散型的属性，效果也还不错。
- C4.5 相比 ID3 有很多改进，但总体上是类似的，改进难度不大。(1) 使用了连续型的属性，在不考虑大量无意义的 0 的前提下使用了二分的方法进行分割。(2) 因为连续属性二分的熵增益在数学上显然比离散属性的熵增益要低，为了合理的评估属性选取，引入（熵）增益率的概念，以此为基础进行属性选取。(3) 树太深很容易产生过拟合，需要人为设定超参数树深。
- CART 是一棵二叉树，代码修改量比较大，这里没有再做尝试。