

机器学习课程 第4次作业

黄昊 20204205

4.1 显然成立：构造这样一颗决策树：第一层判断特征向量的第一个分量，第二层判断第二个... 以此类推。由于数据各不相同，故这样构造出来1决策树，必然能分到一个叶节点，且只有一个数据符合。根据这个构造方法，每个数据到达叶节点的路径各不相同，且一定完全符合（因为各不冲突），故训练误差为0.

4.2 把训练误差作为训练准则容易出现泛化能力差的问题。

4.3

4.4

4.8 算法见下页。如果属性取值较多但属性少，BFS比DFS空间消耗更大；若属性多但属性值少，则DFS比BFS空间消耗更大，DFS有爆栈的风险。

Algorithm 1: 决策树生成算法——基于广度优先搜索**Data:** 训练集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ 属性集 $A = \{a_1, a_2, \dots, a_d\}$ 最大高度 MaxDepth **Result:** 决策树 T

```

1 生成节点  $N$ ，节点信息包括数据集  $D$ ，属性集  $A$ ，高度信息  $h$ ；
2 记录决策树  $T$  的根为  $N$ ；
3 生成节点队列  $Q$ ；
4 将  $N$  压入队列  $Q$  的队尾；
5 while 节点队列  $Q$  非空 do
6   从节点队列  $Q$  中取出队首节点  $N$ ；
7   if 节点  $N.D$  中样本全属于同一类别  $C$  then
8     | 将  $N$  标记为  $C$  类叶节点； continue；
9   end
10  if 节点  $N.h$  已达到  $\text{MaxDepth}$  OR  $N.A = \emptyset$  OR  $N.D$  中样本在  $N.A$  上的取值相同 then
11    | 将  $N$  标记为叶节点，其类别标记为  $N.D$  中样本最多的类； continue；
12  end
13  从  $N.A$  中选择最优划分属性  $a_*$ ；
14  for  $a_*$  的每一个值  $a_*^v$  do
15    为  $N$  生成一个分支；令  $D_v$  为  $D$  中在  $a_*$  上取值为  $a_*^v$  的样本子集；
16    if  $D_v$  为空 then
17      | 将分支节点标记为叶节点，其类别表及为  $D$  中样本最多的类； continue；
18    else
19      | 生成节点  $N_s$ ，节点信息包括数据集  $D_v$ ，属性集  $A \setminus \{a_*\}$ ，高度信息  $N.h + 1$ ；
20      | 将  $N_s$  压入节点队列  $Q$ 
21    end
22  end
23 end
24 return 决策树  $T$ 

```