

数据挖掘

Hw3

11周周评

2021010619



清华大学

Tsinghua University

一、①

采用cart法找根结点

职业类型

	是	否
职员	1	3
工	3	0
管	2	1

$$gini(职业类型) = \frac{3}{10} \times (1 - (\frac{1}{5})^2 - (\frac{3}{5})^2) + \frac{3}{10} \times (1 - (\frac{3}{4})^2 - (\frac{1}{4})^2) + \frac{4}{10} \times (1 - (\frac{2}{5})^2 - (\frac{3}{5})^2) = \frac{17}{60}$$

婚姻

	是	否
已婚	3	1
单身	2	2
离	1	1

$$gini(婚姻) = \frac{4}{10} (1 - (\frac{1}{4})^2 - (\frac{3}{4})^2) + \frac{4}{10} (1 - (\frac{2}{4})^2 - (\frac{2}{4})^2) + \frac{2}{10} (1 - (\frac{1}{2})^2 - (\frac{1}{2})^2) = \frac{9}{20}$$

教育水平

	是	否
非	1	4
高	5	0

$$gini(教育水平) = \frac{5}{10} \times (1 - 1) + \frac{5}{10} \times (1 - (\frac{1}{5})^2 - (\frac{4}{5})^2) = \frac{4}{5}$$

$gini(教育水平) < gini(职业) < gini(婚姻)$ 选教育水平为根

教育水平为高，定是，教育水平为低，进一步计算下一个结点

职业类型

	是	否
职员	0	3
工	1	0
管	0	1

$$gini(职业类型) = \frac{1}{5} \times (1 - 1) + \frac{1}{5} \times (1 - 1) + \frac{1}{5} \times (1 - 1) = 0$$

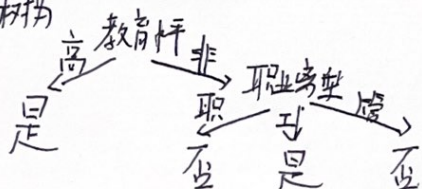
婚姻

	是	否
已婚	0	1
单身	1	2
离	0	1

$$gini(婚姻) = \frac{1}{5} \times (1 - 1) + \frac{1}{5} \times (1 - (\frac{1}{3})^2 - (\frac{2}{3})^2) + \frac{1}{5} \times (1 - 1) = \frac{4}{15}$$

$gini(婚姻) > gini(职业类型)$ 选职业类型

决策树为



给定非高等 → 职员 → 单身

对金钱不感兴趣的人



清华大学

Tsinghua University

② 职业类型

$$P(\text{职}|\text{是}) = \frac{1}{6} \quad P(\text{职}|\text{否}) = \frac{0.3}{4}$$

$$P(\text{工}|\text{是}) = \frac{3}{6} \quad P(\text{工}|\text{否}) = 0$$

$$P(\text{管}|\text{是}) = \frac{2}{6} \quad P(\text{管}|\text{否}) = \frac{1}{4}$$

婚姻

$$P(\text{已}|\text{是}) = \frac{3}{6} \quad P(\text{已}|\text{否}) = \frac{1}{4}$$

$$P(\text{单}|\text{是}) = \frac{2}{6} \quad P(\text{单}|\text{否}) = \frac{2}{4}$$

$$P(\text{空}|\text{是}) = \frac{1}{6} \quad P(\text{空}|\text{否}) = \frac{1}{4}$$

教育

$$P(\text{非}|\text{是}) = \frac{1}{6} \quad P(\text{非}|\text{否}) = 1$$

$$P(\text{高}|\text{是}) = \frac{5}{6} \quad P(\text{高}|\text{否}) = 0$$

$$P(\text{是}, \text{职}, \text{单}, \text{非}) = P(\text{职}|\text{是}) P(\text{单}|\text{是}) P(\text{非}|\text{是}) P(\text{是}) = \frac{1}{6} \times \frac{2}{6} \times \frac{1}{6} \times \frac{6}{10} = \frac{1}{180}$$

$$P(\text{否}, \text{职}, \text{单}, \text{非}) = P(\text{职}|\text{否}) P(\text{单}|\text{否}) P(\text{非}|\text{否}) P(\text{否}) = \frac{0.3}{4} \times \frac{2}{4} \times 1 \times \frac{4}{10} = \frac{3}{20}$$

$$\therefore P(\text{否}, \text{职}, \text{单}, \text{非}) > P(\text{是}, \text{职}, \text{单}, \text{非}) \quad P(\text{否}|\text{职}, \text{单}, \text{非}) > P(\text{是}|\text{单}, \text{职}, \text{非})$$

\therefore 不感兴趣



清华大学

Tsinghua University

	(1,3)	(7,5)
(2,3)	✓	
(1,2)	✓	
(0,1)	✓	
(3,3)	✓	
(2,2)	✓	
(4,4)	✓	
(5,4)	✓	
(6,1)	✓	
(5,6)	✓	
(6,5)	✓	
(7,4)	✓	
(7,6)	✓	

更新
→
中心

$$C_1 = (2.5, \frac{11}{6})$$
$$C_2 = (6, \frac{14}{5})$$

再次迭代, 发现分组仍是

$$C_1: \{ (2,3), (1,2), (3,1), (3,3), (2,2), (4,2) \}$$

$$C_2: \{ (5,4), (6,5), (5,6), (6,5), (7,4), (7,6) \} \text{ 不变}$$

∴ 中心为 $C_1(2.5, \frac{11}{6})$, $C_2(6, \frac{14}{5})$, 分组为

