

Exercises:

$$2\text{-a: } 1 - (1/2)^2 - (1/2)^2 = 0.5$$

$$2\text{-b: } 0$$

2-c:

$$1 - (0.6)^2 - (0.4)^2 = 0.48$$

$$1 - (0.4)^2 - (0.6)^2 = 0.48$$

$$0.5 \times 0.48 + 0.5 \times 0.48 = 0.48$$

2-d:

$$1 - (1/4)^2 - (3/4)^2 = 0.375$$

$$1 - (0/8)^2 - (8/8)^2 = 0$$

$$1 - (1/8)^2 - (7/8)^2 = 0.218$$

$$4/20 \times 0.375 + 8/20 \times 0.218 = 0.16252$$

2-e:

$$1 - (3/5)^2 - (2/5)^2 = 0.48$$

$$1 - (3/7)^2 - (4/7)^2 = 0.4898$$

$$1 - (2/4)^2 - (2/4)^2 = 0.5$$

$$1 - (2/4)^2 - (2/4)^2 = 0.5$$

$$5/20 \times 0.48 + 7/20 \times 0.4898 + 4/20 \times 0.5 + 4/20 \times 0.5 = 0.4914$$

2-f: Car Type

2-g: Everyone is different and has no commonality with others

3-a: Entropy =  $-\frac{4}{9} \times \log_2\left(\frac{4}{9}\right) - \frac{5}{9} \times \log_2\left(\frac{5}{9}\right) = 0.9911$

3-b:

$$\text{Entropy} = \frac{4}{9} \times \left[ -\frac{1}{4} \times \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \times \log_2\left(\frac{3}{4}\right) \right] + \frac{5}{9} \times \left[ -\frac{1}{5} \times \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \times \log_2\left(\frac{4}{5}\right) \right] = 0.7616$$

a1:  $0.9911 - 0.7616 = 0.2294$

$$\text{Entropy} = \frac{5}{9} \times \left[ -\frac{2}{5} \times \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \times \log_2\left(\frac{3}{5}\right) \right] + \frac{4}{9} \times \left[ -\frac{2}{4} \times \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \times \log_2\left(\frac{2}{4}\right) \right] = 0.9839$$

a2:  $0.9911 - 0.9839 = 0.0072$

3-c:

a3	Class label	Split point	Entropy	Info Gain
1.0	+	2.0	0.848	0.143
3.0	-	3.5	0.989	0.003
4.0	+	4.5	0.918	0.073
5.0	-	5.5	0.984	0.007
5.0	-			
6.0	+	6.5	0.973	0.018
7.0	+	7.5	0.889	0.102
7.0	-			

3-d: a1

3-e:

a1's classification error rate =  $\frac{2}{9}$

a2's classification error rate =  $\frac{4}{9}$

So a1 is the best division

3-f:

a1=0.167+0.178=0.345

$$a_2 = 0.267 + 0.222 = 0.489$$

So  $a_1$  is the best division

5-a:

$$E = -0.4 \log 0.4 - 0.6 \log 0.6 = 0.9710$$

$$\Delta A = E - 107EA = T - 103EA = F = 0.2813$$

$$\Delta B = E - 104EB = T - 106EB = F = 0.2565$$

Therefore, the decision tree induction algorithm selects the A attribute

5-b:

$$\text{GINI} : G = 1 - (0.4)^2 - (0.6)^2 = 0.48$$

$$\text{GINI}_{A=T} = 1 - (0.74)^2 - (0.26)^2 = 0.4898$$

$$\text{GINI}_{A=F} = 1 - (0.30)^2 - (0.70)^2 = 0.48$$

$$EA = \text{GINI} - 107\text{GINI}_{A=T} - 103\text{GINI}_{A=F} = 0.1371$$

$$\text{GINI}_{B=T} = 1 - (0.43)^2 - (0.57)^2 = 0.3750$$

$$\text{GINI}_{B=F} = 1 - (0.61)^2 - (0.39)^2 = 0.2778$$

$$EB = \text{GINI} - 104\text{GINI}_{B=T} - 106\text{GINI}_{B=F} = 0.1633$$

Therefore, the decision tree induction algorithm selects the B attribute

5-c:

Information gain examines the contribution of features to the entire data, not to specific categories, so generally it can only be used for global feature selection

The Gini coefficient is a feature selection method similar to the information entropy, which is used for the impurity of the data. When making feature selection, we can choose the one with the largest  $\Delta \text{Gini}(X)$ .

18-a:  $0.5 \times 0 + 0.5 \times 1 = 50\%$  (The number of positive examples and negative examples are equal)

18-b:  $0.5 \times 0.8 + 0.5 \times 0.2 = 50\%$

18-c:  $2/3 \times 0 + 1/3 \times 1 = 1/3 = 33.3\%$

18-d:  $2/3 \times 1/3 + 1/3 \times 2/3 = 4/9 = 44.4\%$