

Homework - 3

Problem 1: Decision Trees based on Entropy

(1) The Entropy for a given node t is given by -

$$\text{Entropy}(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

→ The overall entropy before splitting is,

$$\begin{aligned} \text{Entropy}(t) &= -\sum_j p(j|t) \cdot \log_2 p(j|t) \\ &= -\left[\frac{4}{9} \log_2 \left(\frac{4}{9} \right) + \frac{5}{9} \log_2 \left(\frac{5}{9} \right) \right] \\ &= -\left[(-0.5199) + (-0.4711) \right] \\ &= \underline{\underline{0.9910}} \end{aligned}$$

→ The contingency tables after splitting on attribute a_1 ,

	$a_1 = T$	$a_1 = F$
$+$	3	1
$-$	1	4

→ Entropy for attribute a_1 ,

$$\begin{aligned} E_{a_1} &= -\left[\frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right] \\ &= -\left[(-0.3112) + (-0.5) \right] \\ &= \underline{\underline{0.8112}} \end{aligned}$$

$$\begin{aligned}
 E_{a_1=F} &= - \left[\frac{1}{5} \log_2 \left(\frac{1}{5} \right) + \frac{4}{5} \log_2 \left(\frac{4}{5} \right) \right] \\
 &= - \left[(-0.4643) + (-0.2575) \right] \\
 &= \underline{\underline{0.7218}}
 \end{aligned}$$

⇒ Calculating the Information gain,

$$\Delta_{a_1} = \text{Entropy}(+) - \left[\frac{4}{9} \times \text{Entropy}(a_1=T) + \frac{5}{9} \times \text{Entropy}(a_1=F) \right]$$

$$= 0.9910 - \left[\frac{4}{9} \times (0.8112) + \frac{5}{9} (0.7218) \right]$$

$$= 0.9910 - [\cancel{0.3005} \ 0.76153]$$

$$\boxed{\Delta_{a_1} = 0.22946}$$

⇒ Since a_2 is a continuous attribute, the contingency table after splitting on a_2 ,

Sorted values 1 3 4 5 6 7 8

Split position	0.5		2		3.5		4.5		5.5		6.5		7.5		8.5	
	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>	<=	>
+	0	4	1	3	1	3	2	2	2	2	3	1	4	0	4	0
-	0	5	0	5	1	4	1	4	3	2	3	2	4	1	5	0
Gain	0		0.1427		0.00248		0.0728		0.00722		0.0183		0.1021		0	

Now,

⇒ split 1:

$$H(1) = - \sum_{i=1}^5 p(j|t) \log_2 p(j|t)$$

$$= - \left[\frac{4}{9} \log_2 \left(\frac{4}{9} \right) + \frac{5}{9} \log_2 \left(\frac{5}{9} \right) \right]$$

$$= 0.99107$$

Information Gain = $E(t) - H(1)$

$$= 0.9910 - 0.9910$$

$$= \underline{0}$$

⇒ split 2:

$$H(2: <=) = - \left(\frac{1}{1} \log_2 \left(\frac{1}{1} \right) + 0 \log_2 0 \right)$$

$$= 0$$

$$H(2: >) = - \left(\frac{3}{8} \log_2 \left(\frac{3}{8} \right) + \frac{5}{8} \log_2 \left(\frac{5}{8} \right) \right)$$

$$= - \left[(-0.5306) + (-0.4237) \right]$$

$$= 0.9543$$

Weighted Average = $\left[\frac{1}{9} \times 0 + \frac{8}{9} \times 0.9543 \right]$

$$= 0.8482$$

Information Gain = $0.9910 - 0.8482$

$$= \underline{0.1427}$$

⇒ split 3:

$$H(3.5: <=) = - \left(\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right)$$

$$= - \left[(-0.5) + (-0.5) \right] = 1$$

$$\begin{aligned}
 H(3.5: \geq) &= - \left[\frac{3}{7} \log_2 \left(\frac{3}{7} \right) + \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right] \\
 &= - \left[(-0.5238) + (-0.4613) \right] \\
 &= \underline{\underline{0.9851}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Weighted Average} &= \left(\frac{2}{9} \times 1 + \frac{7}{9} \times 0.98523 \right) \\
 &= 0.98851
 \end{aligned}$$

$$\begin{aligned}
 \text{Information Gain} &= 0.9910 - 0.98851 \\
 &= \underline{\underline{0.00248}}
 \end{aligned}$$

\Rightarrow Split 4 :

$$\begin{aligned}
 H(4.5: <=) &= - \left[\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] \\
 &= - \left[(-0.3899) + (-0.5283) \right] \\
 &= 0.9182
 \end{aligned}$$

$$\begin{aligned}
 H(4.5: \geq) &= - \left[\frac{2}{6} \log_2 \left(\frac{2}{6} \right) + \frac{4}{6} \log_2 \left(\frac{4}{6} \right) \right] \\
 &= - \left[(-0.5283) + (-0.3899) \right] \\
 &= 0.9182
 \end{aligned}$$

$$\begin{aligned}
 \text{Weighted Average} &= \left(\frac{3}{9} \times 0.9182 + \frac{6}{9} \times 0.9182 \right) \\
 &= \underline{\underline{0.9182}}
 \end{aligned}$$

$$\text{Information Gain} = A = 0.9910 - 0.9182 = \underline{\underline{0.0728}}$$

→ split 5:

$$\begin{aligned}H(5.5 : \leq) &= - \left[\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] \\&= - \left[(-0.5287) + (-0.4421) \right] \\&= 0.9708\end{aligned}$$

$$\begin{aligned}H(5.5 : >) &= - \left[\frac{2}{4} \log_2 \left(\frac{2}{4} \right) + \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] \\&= 1\end{aligned}$$

$$\begin{aligned}\text{Weighted Average} &= \left[\frac{5}{9} \times 0.9708 + 1 \times \frac{4}{9} \right] \\&= \underline{\underline{0.9837}}\end{aligned}$$

$$\begin{aligned}\text{Information Gain} &= 0.9910 - 0.9837 \\&= \underline{\underline{0.00722}}\end{aligned}$$

→ split 6:

$$\begin{aligned}H(6.5 : \leq) &= - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right] \\&= 1\end{aligned}$$

$$\begin{aligned}H(6.5 : >) &= - \left[\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] \\&= - \left[(-0.5283) + (-0.38997) \right] \\&= \underline{\underline{0.9182}}\end{aligned}$$

$$\text{Weighted average} = \left[\frac{6}{9} \times 1 + \frac{3}{9} \times 0.9182 \right] = \underline{\underline{0.9727}}$$

$$\begin{aligned}\text{Information Gain} &= 0.9910 - 0.9727 \\ &= \underline{\underline{0.0183}}\end{aligned}$$

→ Split 7:

$$H(7.5 : <=) = - \left[\frac{4}{8} \log_2 \left(\frac{4}{8} \right) + \frac{4}{8} \log_2 \left(\frac{4}{8} \right) \right] = 1$$

$$H(7.5 : >) = - [0 \log_2 0 + 1 \log_2 1] = 0$$

$$\text{weighted average} = \frac{8}{9} \times 1 = \underline{\underline{0.8889}}$$

$$\begin{aligned}\text{Information gain} &= 0.9910 - 0.8889 \\ &= 0.1021\end{aligned}$$

→ Split 8:

$$H(8.5 : <=) = - \left[\frac{4}{9} \log_2 \left(\frac{4}{9} \right) + \frac{5}{9} \log_2 \left(\frac{5}{9} \right) \right]$$

$$= - [(-0.51996) + (-0.4711)]$$

$$= 0.9910$$

$$H(8.5 : >) = - [0 \log_2(0) + 0 \log_2(0)] = 0$$

$$\text{Information gain} = 0.9910 - 0.9910 = 0$$

⇒ For attribute a_2 the best split is 2.0 with information gain 0.1427

⇒ But attribute a_1 has higher information gain

⇒ Hence, best attribute for the first splitting is attribute a_1 .

(2) If we use "Instance" as another attribute and split point = 2.5, we can get best split point for instance

count matrix < 2.5

	<=	>
+	2	2
-	0	5

$$\text{Entropy}_{(2.5: \leq)} = - \left[\frac{2}{2} \log_2 \left(\frac{2}{2} \right) + 0 \log_2 0 \right] = 0$$

$$\text{Entropy}_{(2.5: >)} = - \left[\frac{2}{7} \log_2 \left(\frac{2}{7} \right) + \frac{5}{7} \log_2 \left(\frac{5}{7} \right) \right] = \underline{\underline{0.86312}}$$

$$\text{Information Gain} = H(t) - \frac{7}{9} \times \text{Entropy}$$

$$= 0.99107 - 0.6713$$

$$= \underline{\underline{0.3198}}$$

By Comparing both attribute, instance has high ~~attribute~~ information gain. Therefore instance as an attribute is a root node and it can be used as an attribute for a decision in a tree.