

Machine Learning

Quiz - 1 - Decision Trees

- Kinshuk Vasisht

- Roll Number : 19

- M.Sc. Computer Science

Q3. b) Instances :

f_1	f_2	L
10	3	1
6	3	1
2	7	0
9	6	1
4	6	0

Given splits : $f_1 > 2$

$f_1 > 4$

$f_2 > 3$

$f_2 > 6$

Now, Information Gain (Δ_{info}) = $\phi_{parent} - \phi_{split}$

(Node entropy for parent

- Split entropy based on decision)

$$\Rightarrow \Delta_{info} = - \sum_{i=1}^2 p_m^i \log_2(p_m^i) + \sum_{i=1}^2 \frac{N_{m_i}}{N_m} \sum$$

$$\Delta_{info} = - \underbrace{\sum_{i=1}^2 p_m^i \log_2(p_m^i)}_{\text{Node}} + \underbrace{\sum_{j=1}^2 \frac{N_m^j}{N_m} \sum_{i=1}^2 p_{mj}^i \log p_{mj}^i}_{\text{Split}}$$

where p_m^i = probability of class i at node m

p_{mj}^i = probability of class i at node m , branch j

$$\text{Now, } \phi_{parent} = - \left[\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right]$$

1) Now, split entropy for $f_1 > 2$

	$L=0$	$L=1$
$f_1 > 2$	1	2
$f_1 \leq 2$	1	1

$$\therefore \phi_{f_1 > 2} = - \frac{3}{5} \sum_i p_{m1}^i \log p_{m1}^i$$

$$- \frac{2}{5} \sum_i p_{m2}^i \log p_{m2}^i$$

$$\Rightarrow -\frac{3}{5} \left[\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right] - \frac{2}{5} \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right]$$

$$= +\frac{3}{5} \left[\frac{\log_2 3}{3} + \frac{2 \log_2 1.5}{3} \right] + \frac{2}{5} \times 1 =$$

2) Now, split entropy for $f_1 > 4$

	$L=0$	$L=1$	
$f_1 > 4$	0	3	$\therefore \phi_{f_1 > 4} = -\frac{3}{5} \sum_i p_{m1}^i \log_2(p_{m1}^i) - \frac{2}{5} \sum_i p_{m2}^i \log_2(p_{m2}^i)$
$f_1 \leq 4$	2	0	

$$\Rightarrow \phi_{f_1 > 4} = -\frac{3}{5} \left[\frac{0}{3} \log_2 \frac{0}{3} + \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right] - \frac{2}{5} \left[\frac{2}{2} \log_2 \frac{2}{2} + \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$= (3/5) \cdot 0 + (2/5) \cdot 0 = \underline{0}$$

3) Now, split entropy for $f_2 > 3$

	$L=0$	$L=1$	
$f_2 > 3$	1	1	$\therefore \phi_{f_2 > 3} = -\frac{2}{5} \sum_i p_{m1}^i \log_2(p_{m1}^i) - \frac{3}{5} \sum_i p_{m2}^i \log_2(p_{m2}^i)$
$f_2 \leq 3$	1	2	

$$\Rightarrow \phi_{f_2 > 3} = -\frac{2}{5} \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] - \frac{3}{5} \left[\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right]$$

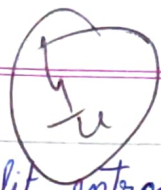
$$= +\frac{2}{5} \times 1 + \frac{3}{5} \left[\frac{\log_2(3)}{3} + \frac{2 \log_2(1.5)}{3} \right] =$$

4) Now, split entropy for $f_2 > 6$

	$L=0$	$L=1$	
$f_2 > 6$	1	0	$\phi_{f_2 > 6} = -\frac{1}{5} \sum_i p_{m1}^i \log(p_{m1}^i) - \frac{4}{5} \sum_i p_{m2}^i \log(p_{m2}^i)$
$f_2 \leq 6$	1	3	

$$\Rightarrow \phi_{f_2 > 6} = -\frac{1}{5} \left[\frac{1}{1} \log \left(\frac{1}{1} \right) + \frac{0}{1} \log \left(\frac{0}{1} \right) \right] - \frac{4}{5} \left[\frac{1}{4} \log \left(\frac{1}{4} \right) + \frac{3}{4} \log \left(\frac{3}{4} \right) \right]$$

$$= \frac{1}{5} \times 0 + \frac{4}{5} \left[\frac{2}{4} + \left(-\frac{3}{4} \log \left(\frac{4}{3} \right) \right) \right] = \underline{1.14}$$



Now, split entropy for $f_i > 4$ is the lowest (0)
 Corresponding to the same, the gain obtained by subtracting
 split entropy from the node entropy of the parent is the
 highest.

∴ Best Split : $f_i > 4$ ($f_i > 4$)

Q5 Every node of the ~~the~~ decision tree makes a split over an
 attribute different from that of its parent, based on the
 best split criteria (obtaining lowest heterogeneity).
 Following the given splits, the primary split is made
 over $x > x_1$, then on $y < y_1$ & $y > y_3$ on left
 & right subtrees. So we have — :

