

Data Mining (CS4038)

Quiz No.4

Solution

Roll No: _____

Section: B

Date: 01-04-2024

Question No.1

(10 marks)

Consider the following dataset of Passed and Failed exams with three attributes: studied(Y/N), slept(Y/N) and cheated(Y/N) and answer the questions given below. Compute the Gini Index for all of these attributes and identify which is best attribute to serve as root node.

Student ID	Studied	Slept	Cheated	Result
1	Yes	No	No	Passed
2	Yes	No	Yes	Failed
3	No	Yes	No	Failed
4	Yes	Yes	Yes	Failed
5	Yes	Yes	No	Passed
6	No	No	Yes	Failed
7	No	No	No	Failed
8	No	Yes	Yes	Failed

$$Gini(t) = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

$$Gini(Result) = 1 - \left(\frac{2}{8}\right)^2 - \left(\frac{6}{8}\right)^2 = \boxed{0.375}$$

Result
Passed 2
Filed $\frac{6}{8}$

~~Gini(Studied/Yes)~~

x3

$$Gini(Yes/studied) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

studied
Passed Yes No
2 0
Failed $\frac{2}{4}$ $\frac{4}{4}$

$$Gini(No/studied) = 1 - \left(\frac{0}{4}\right)^2 - \left(\frac{4}{4}\right)^2 = 0$$

$$Gini(Studied) = \frac{4}{8}(0.5) + \frac{4}{8}(0) = \boxed{0.25}$$

$$Gini(Yes/slept) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 \quad \underline{\underline{+3}}$$

$$= 0.375$$

	<u>Slept</u>	
	yes	No
passed	1	1
Failed	$\frac{3}{4}$	$\frac{3}{4}$

$$Gini(No/slept) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2$$

$$= 0.375$$

$$Gini(slept) = \frac{4}{8}(0.375) + \frac{4}{8}(0.375)$$

$$= 0.1875 + 0.1875 = \boxed{0.375}$$

$$Gini(Yes/cheated) = 1 - \left(\frac{0}{4}\right)^2 - \left(\frac{4}{4}\right)^2 \quad \underline{\underline{+3}}$$

$$= 0$$

	<u>cheated</u>	
	yes	No
passed	0	2
Failed	$\frac{4}{4}$	$\frac{2}{4}$

$$Gini(No/cheated) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2$$

$$= 0.5$$

$$Gini(cheated) = \frac{4}{8}(0) + \frac{4}{8}(0.5)$$

$$= \boxed{0.25}$$

As the Gini Index value of studied and cheated attributes is the lowest and same so, ~~both~~ anyone of these 2 attributes can serve as Root node in Decision tree. +1

Data Mining (CS4038)

Quiz No.4

Roll No: _____

Section: A

Date: 01-04-2024

Question No.1

(10 marks)

Consider the following dataset of Passed and Failed exams with three attributes: studied(Y/N), slept(Y/N) and cheated(Y/N) and answer the questions given below. Compute gain ratio of all these attributes and identify which is the best attribute to serve as root node.

Student ID	Studied	Slept	Cheated	Result
1	Yes	No	No	Passed
2	Yes	No	Yes	Failed
3	No	Yes	No	Failed
4	Yes	Yes	Yes	Failed
5	Yes	Yes	No	Passed
6	No	No	Yes	Failed
7	No	No	No	Failed
8	No	Yes	Yes	Failed

$$\text{Entropy}(t) = - \sum_{i=0}^{c-1} P_i(t) \log_2 P_i(t)$$

$$\begin{aligned} \text{Entropy}(\text{Result}) &= - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) \quad \begin{array}{l} \text{Passed} \quad 2 \\ \text{Failed} \quad 6 \end{array} \\ &= 0.5 + 0.3112 \\ &= \boxed{0.8112} \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Studied}) &= - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \quad \begin{array}{l} \text{studied} \\ \text{yes} \quad \text{no} \\ \text{passed} \quad 2 \quad 0 \\ \text{failed} \quad \frac{2}{4} \quad \frac{4}{4} \end{array} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{no}|\text{studied}) &= - \frac{0}{4} \log_2 \left(\frac{0}{4} \right) - \frac{4}{4} \log_2 \left(\frac{4}{4} \right) \\ &= 0 \end{aligned}$$

$$\text{Entropy}(\text{studied}) = \frac{4}{8} (1) + \frac{4}{8} (0)$$

$$\begin{aligned} \text{Gain info} &= \text{Ent}(\text{Result}) - \text{Ent}(\text{studied}) \\ &= 0.8112 - 0.5 = 0.3112 \end{aligned}$$

$$\text{Split info} = - \sum_{i=1}^K \frac{n_i}{n} \log_2 \frac{n_i}{n} = - \frac{4}{8} \log_2 \frac{4}{8} + \frac{4}{8} \log_2 \frac{4}{8} = 1$$

$$\text{Gain Ratio} = \frac{\text{Gain info}}{\text{Split info}} = \frac{0.3112}{1} = \boxed{0.3112}$$

$$Ent(Yes|slept) = -\frac{1}{4} \log_2(\frac{1}{4}) - \frac{3}{4} \log_2(\frac{3}{4})$$

$$= 0.8112$$

$$Ent(no|slept) = -\frac{1}{4} \log_2(\frac{1}{4}) - \frac{3}{4} \log_2(\frac{3}{4})$$

$$= 0.8112$$

$$Ent(slept) = \frac{4}{8} (0.8112) + \frac{4}{8} (0.8112)$$

$$= 0.8112$$

$$Gain_{info} = 0.8112 - 0.8112 = 0$$

$$Split_{info} = -\frac{4}{8} \log_2 \frac{4}{8} - \frac{4}{8} \log_2 \frac{4}{8} = 1$$

$$Gain_{Ratio} = \frac{0}{1} = \boxed{0}$$

$$Ent(Yes|cheated) = -\frac{0}{4} \log_2(\frac{0}{4}) - \frac{4}{4} \log_2(\frac{4}{4})$$

$$= 0$$

	cheated	
	yes	no
Passed	0	2
Failed	$\frac{4}{4}$	$\frac{2}{4}$

$$Ent(No|cheated) = -\frac{2}{4} \log_2(\frac{2}{4}) - (\frac{2}{4}) \log_2(\frac{2}{4})$$

$$= 0.5 + 0.5 = 1$$

$$Ent(cheated) = \frac{4}{8} (0) + \frac{4}{8} (1) = 0.5$$

$$Gain_{info} = 0.8112 - 0.5 = 0.3112$$

$$Split_{info} = -\frac{4}{8} \log_2 \frac{4}{8} - \frac{4}{8} \log_2 \frac{4}{8} = 1$$

$$Gain_{Ratio} = \frac{0.3112}{1} = \boxed{0.3112}$$

So, As the Gain Ratio of Studied and ⁺¹cheated attributes is highest and same so anyone of these attributes can serve as root node for Decision tree.