

Assignment On
Data Mining (MCSE 642)
Assignment – 01

Submitted to
Dr. Md Tauhid Bin Iqbal

Submitted by:
Mohammed Morad Hossen
Student ID: MCE 079 055 32
Dept of Computer Science and Engineering
Stamford University Bangladesh

Income

Name: Mohammed Morka
ID: MCE 079 055 32

Income is a nominal feature. It can be High, medium, Low.
I will summarize the final decision for Income feature

Income	Yes	No	Number of Instances
High	2	2	4
Medium	4	2	6
Low	3	1	4

$$\text{Gini}(\text{Income} = \text{High}) = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 1 - 0.25 - 0.25 \\ = 0.5$$

$$\text{Gini}(\text{Income} = \text{Medium}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 1 - 0.44 - 0.11 = 0.45$$

$$\text{Gini}(\text{Income} = \text{Low}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 1 - 0.563 - 0.063 \\ = 0.374$$

Then we will calculate weighted sum of Gini index for Income feature

$$= \frac{4}{14}(0.5) + \frac{6}{14}(0.45) + \frac{4}{14}(0.374) \\ = 0.143 + 0.193 + 0.107 \\ = 0.443$$

STUDENT

Student is a nominal feature. It can be Yes, No.
I will summarize the final decision for student feature

Student	Yes	No	Number of Instances
Yes	6	1	7
No	3	4	7

$$\text{Gini}(\text{Student} = \text{Yes}) = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 1 - 0.73 - 0.02 \\ = 0.25$$

$$\text{Gini}(\text{Student} = \text{No}) = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 1 - 0.18 - 0.33 \\ = 0.49$$

Then we will calculate weighted sum of Gini Index for ~~Income~~ ^{Student} feature

$$= \cancel{\left(\frac{7}{7}\right)}(0.25) + \cancel{\left(\frac{7}{7}\right)} \\ = \left(\frac{7}{14}\right)(0.25) + \left(\frac{7}{14}\right)0.49 \\ = 0.125 + 0.245 = 0.37$$

AG6

Age is a nominal feature. It can be ≤ 30 , $31-40$, > 40 . I will summarize the final decision for Age feature.

Age	Yes	No	Number of Instances
≤ 30	2	3	5
$31-40$	4	0	4
> 40	3	2	5

$$\text{Gini}(\text{Age} = \leq 30) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 1 - 0.16 - 0.36 \\ = 0.48$$

$$\text{Gini}(\text{Age} = 31-40) = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 1 - 1 - 0 = 0$$

$$\text{Gini}(\text{Age} > 40) = 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 1 - 0.36 - 0.16 \\ = 0.48$$

~~Credit Rating: Credit Rating is a nominal feature. It can be~~

Then we will calculate sum of GINI Index for Age

$$\text{feature} = \left(\frac{5}{14}\right) \times 0.48 + \left(\frac{4}{14}\right) \times 0 + \left(\frac{5}{14}\right) \times 0.48 \\ = 0.171 + 0.171 = 0.34$$

Credit-Rating

(3)

Credit Rating is a nominal feature. It can be Fair, Excellent. It will summarize the final decision for Credit Rating.

Credit-Rating	Yes	No	Number of Instances.
Fair	6	2	8
Excellent	3	3	6

$$\text{Gini}(C.R = \text{Fair}) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 1 - 0.56 - 0.063 = 0.38$$

$$\text{Gini}(C.R = \text{Excellent}) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 1 - 0.25 - 0.25 = 0.50$$

Then we will calculate weighted sum of GINI index for Credit Rating

$$= \left(\frac{8}{14}\right) \times 0.38 + \left(\frac{6}{14}\right) (0.5) \\ = 0.22 + 0.21 = 0.43$$

Feature	GINI Index
Age	0.34
Income	0.443
Student	0.37
Credit-Rating	0.43

Hence Age is the winner feature because its cost is lowest (0.34). The Age attribute will be considered as the root node.

(4)

Calculating Entropy

$$E(S) = \sum_{i=1}^c -P_i (\log_2 P_i)$$

Buy Computer	
Yes	No
9	5

Entropy (Buy-Computer)

$$= -\left(\frac{9}{14}\right) \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$= -0.64 \times (-0.64) - 0.36 \times (-1.47)$$

$$= 0.41 + 0.53 = 0.94$$

Income	Buyer Buy Computer		Number of Instances
	Yes	No	
High	2	2	4
Medium	4	2	6
Low	3	1	4

$$\begin{aligned} \text{Entropy (Income, High)} &= -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \left(\frac{2}{4}\right) \log_2 \left(\frac{2}{4}\right) \\ &= -0.5(-1) - (0.5)(-1) \\ &= 0.5 + 0.5 = 1 \end{aligned}$$

$$\begin{aligned} \text{Entropy (Income, Medium)} &= -4/6 \log_2 (4/6) - \frac{2}{6} \log_2 (2/6) \\ &= -0.67(-0.58) - 0.33(-1.6) \\ &= 0.39 + 0.53 \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{Entropy (Income, Low)} &= -\frac{3}{4} \log_2 (3/4) - \frac{1}{4} \log_2 (1/4) \\ &= -0.75(-0.42) - 0.25(-2) \\ &= 0.32 + 0.5 = 0.82 \end{aligned}$$

(5)

Entropy (Buy Computer, Income)

$$= P(\text{High}) E(2, 2) + P(\text{medium}) E(4, 2) + P(\text{Low}) E(3, 1)$$

$$= \frac{4}{14} \times 1 + \frac{6}{14} \times 0.92 + \frac{4}{14} \times 0.82$$

$$= 0.29 + 0.4 + 0.23 = 0.92$$

Information Gain (Buy Computer, Income)

$$= E(\text{Buy Computer}) - E(\text{Buy Computer, Income})$$

$$= 0.94 - 0.92$$

$$= 0.02$$

Student

Student	Buy Computer		Number of Instances
	Yes	No	
Yes	6	1	7
No	3	4	7

$$\text{Entropy}(\text{Student, Yes}) = -\frac{6}{7} \log_2(6/7) - \frac{1}{7} (\log_2 1/7)$$

$$= -0.86(-0.22) - 0.14(-2.84)$$

$$= 0.19 + 0.4 = 0.59$$

$$\text{Entropy}(\text{Student, No}) = -\frac{3}{7} \log_2(3/7) - \frac{4}{7} \log_2(4/7)$$

$$= -0.43(-1.22) - 0.57(-0.81)$$

$$= 0.52 + 0.46 = 0.98$$

$$\text{Entropy}(\text{Student, Buy computer}) = P(\text{Yes}) E(6, 1) + P(\text{No}) E(3, 4)$$

$$= \frac{7}{14} \times 0.59 + \frac{7}{14} (0.98)$$

$$= 0.34 + 0.49 = 0.79$$

Information Gain (Buy Computer, Student) ^⑥ - ^⑥

$$= E(\text{Buy Computer}) - E(\text{Buy Computer, Student})$$

$$= 0.94 - 0.79 = 0.15$$

Age

Age	Buy Computer		Number of Instances
	Yes	No	
<=30	2	3	5
31-40	4	0	4
>40	3	2	5

$$\text{Entropy}(\text{Age}, \leq 30) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= -0.4(-1.32) - 0.6(-0.74)$$

$$= 0.53 + 0.44 = 0.97$$

$$\text{Entropy}(\text{Age}, 31-40) = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4}$$

$$= -1(0) - 0 = 0$$

$$\text{Entropy}(\text{Age}, > 40) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= -0.6(-0.74) - 0.4(-1.32)$$

$$= 0.44 + 0.53 = 0.97$$

$$\text{Entropy}(\text{Age, Buy Computer}) = P(\leq 30)E(2,3) + P(31-40)E(4,0) + P(>40)E(3,2)$$

$$= \frac{5}{14}(0.97) + \frac{4}{14}(0) + \frac{5}{14}(0.97)$$

$$= 0.35 + 0.35 = 0.7$$

Information Gain (Buy Computer, Age)

$$= E(\text{Buy Computer}) - E(\text{Buy Computer, Age})$$

$$= 0.94 - 0.7 = 0.24$$

Credit-Rating

(4)

Credit-Rating	Buy Computer		Number of Instances
	Yes	No	
Fair	6	2	8
Excellent	3	3	6

$$\begin{aligned}\text{Entropy}(\text{C.R., Fair}) &= -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \\ &= -0.75(-0.42) - 0.25(-2) \\ &= 0.33 + 0.5 = 0.83\end{aligned}$$

$$\begin{aligned}\text{Entropy}(\text{C.R., Excellent}) &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= -0.5(-1) - (0.5)(-1) \\ &= 0.5 + 0.5 = 1\end{aligned}$$

$$\begin{aligned}\text{Entropy}(\text{C.R., Buy Computer}) &= \frac{8}{14} (0.83) + \frac{6}{14} (1) \\ &= \cancel{0.57} = 0.47 + 0.43 = 0.9\end{aligned}$$

$$\begin{aligned}\text{Informain Gain}(\text{Buy Computer, Credit-Rating}) &= E(\text{Buy Computer}) - E(\text{Buy Computer, C.R.}) \\ &= 0.94 - 0.9 = 0.04\end{aligned}$$

feature	Information Gain
Age	0.24
Income	0.02
Student	0.15
Credit-Rating	0.04

Age feature produce highest score. So Age will be considered as the root node.

$$\begin{aligned}
 \text{Gain Ratio (Decision, Age)} &= \frac{\text{I.Gain (Decision, Age)}}{\text{Splitinfo (Decision, Age)}} \\
 &= \frac{0.24}{0.7} = 0.34
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain Ratio (Decision, Income)} &= \frac{\text{I.G (Decision, Income)}}{\text{Splitinfo (Decision, Income)}} \\
 &= \frac{0.02}{0.92} = 0.02
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain Ratio (Decision, Student)} &= \frac{\text{I.G (Decision, Student)}}{\text{Splitinfo (Decision, Student)}} \\
 &= \frac{0.15}{0.79} = 0.19
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain Ratio} &= (\text{Decision}, \text{Credit-Ration}) \\
 &= \frac{I.G. (\text{Decision}, C.R.)}{\text{Splitinf} (\text{Decision}, C.R.)} \\
 &= \frac{0.04}{0.9} = 0.04
 \end{aligned}$$

feature	Gain Ratio
Age	0.34
Income	0.02
Student	0.19
Credit Rating	0.04

hence Age feature has highest gain ratio and has been selected as root node.