

MIS 637 ASSIGNMENT-3

POORVI RAUT - 20009560

chapter 6, page 127, Problems 5-10

Occupation	Gender	Age	Salary
Service	Female	45	\$48,000
	Male	25	\$25,000
Management	Male	33	\$85,000
	Female	25	\$45,000
Management	Male	35	\$65,000
	Male	26	\$45,000
Sales	Female	45	\$70,000
	Female	40	\$50,000
Sales	Male	30	\$40,000
	Male	25	\$25,000
Staff	Female	50	\$40,000
Staff	Male	25	\$25,000

Consider the data in above Table. The target variable is salary. start by discretizing salary as follows :

less than \$ 35,000 level 1

\$35,000 to less than \$ 45,000 Level 2

\$45,000 to less than \$ 55,000 level 3

Above \$ 55,000 level 4.

Construct a classification & regression tree to classify salary based on the other variables. Do as much as you can by hand, before turning to the software.

Occupation	Gender	Age	Salary(\$)	level	
Service	female	45	\$48,000	level 3	
Service	male	25	\$25,000	level 1	
Service	male	33	\$85,000	level 2	
Management	male	25	\$45,000	level 3	
Management	Female	35	\$65,000	level 4	
Management	male	26	\$45,000	level 3	
Management	Female	45	\$70,000	level 4	

CONTD page 2.

Occupation	Gender	Age	Salary (\$)	level
Sales	female	40	\$50,000	level 3
Sales	male	30	\$40,000	level 2
staff	female	50	\$40,000	level 2
staff	male	25	\$25,000	level 1

Candidate split for $t = \text{Root Node}$.

Candidate Split	left child Node, t_L	Right child node, t_R
1	Occupation : Service	Occupation : management, sales, staff
2.	Occupation : management	Occupation : service, sales, staff
3	Occupation : sales	Occupation : service, management, staff
4	Occupation: staff	Occupation: service, management, sales
5	gender : male	Gender : female
6	Age < 30	Age ≥ 30
7	Age < 40	Age ≥ 40

$\phi(S|t)$ ← optimality measure of split (maximizes split over all possible splits)
Let $\phi(S|t)$ be measure of "goodness" of candidate split S at node t

where formula is

$$\phi(S|t) = 2P_L P_R \sum_{j=1}^{\# \text{classes}} |P(j|t_L) - P(j|t_R)|$$

t_L = Left child node of node t

t_R = Right child node of node t

P_L = number of records at t_L

number of records in training set

P_R = number of records at t_R

number of records at training set

$P(j|t_L)$ = number of class j records at t_L / number of records at t

$P(j|t_R)$ = number of class j records at t_R / number of records at t

$$D(t) = \sum_{j=1}^{\# \text{ classes}} |P(j|t_L) - P(j|t_R)|$$

$D(t)$ is large when distance between $P(j|t_L)$ & $P(j|t_R)$ is maximized across each value of target variable.

∴ The values for our candidate split for $t = \text{Root}$ will be determined as for example t_L

1st candidate value for occupation = service

$$P_L = \frac{3}{11} = 0.2727 \quad \text{since service as occupation in dataset} = 3$$

Total value = 11

$$P_R = \frac{8}{11} = 0.7272 \quad // \underline{\text{other occupation excluding service}}$$

Total

$$P(j|t_L) = \frac{\text{level 1 (service)}}{\text{Total Service in left child node}} = \frac{1}{3} = 0.333$$

$$\frac{\text{level 2 (service)}}{\text{service in } t_L} = \frac{1}{3} = 0.33$$

$$\frac{\text{level 3 (service)}}{\text{service in } t_L} = \frac{1}{3} = 0.33$$

$$\frac{\text{level 4 of salary in service}}{\text{service in } t_L} = \frac{0}{3} = 0$$

$$P(j|t_R) = \frac{\text{level 1 (other in } t_R)}{\text{other occupation in } t_R} = \frac{1}{8} = 0.125$$

$$\frac{\text{level 2 (other in } t_R)}{\text{other occupation in } t_R} = \frac{2}{8} = 0.25$$

$$\frac{\text{level 3 (other in } t_R)}{\text{other occupations in } t_R} = \frac{3}{8} = 0.375$$

$$\frac{\text{level 4 (other in } t_R)}{\text{other occupations in } t_R} = \frac{2}{8} = 0.25$$

$$\text{Given } 2PL \times PR = 2 \times 0.2727 \times 0.7273 = 0.3966$$

$$\begin{aligned}\Phi(S|t) &= \sum_{j=1}^4 |P(j|t_L) - P(j|t_R)| \\ &= |\sum P(j|t_L) - \sum P(j|t_R)|\end{aligned}$$

$$\begin{aligned}&= |P(\text{level 1}|t_L) - P(\text{level 1}|t_R)| + |P(\text{level 2}|t_L) - P(\text{level 2}|t_R)| \\ &\quad + |P(\text{level 3}|t_L) - P(\text{level 3}|t_R)| + |P(\text{level 4}|t_L) - P(\text{level 4}|t_R)| \\ &= |0.333 - 0.125| + |0.33 - 0.25| + |0.33 - 0.375| + |0 - 0.25| \\ &= 0.208 + 0.08 + 0.045 + 0.25\end{aligned}$$

$$\boxed{\Phi(S|t) = 0.583}$$

$$\begin{aligned}\phi(S|t) &= \Phi(S|t) \times 2PL \times PR \\ &= 0.583 \times 0.3966 \\ &= \boxed{0.2312}\end{aligned}$$

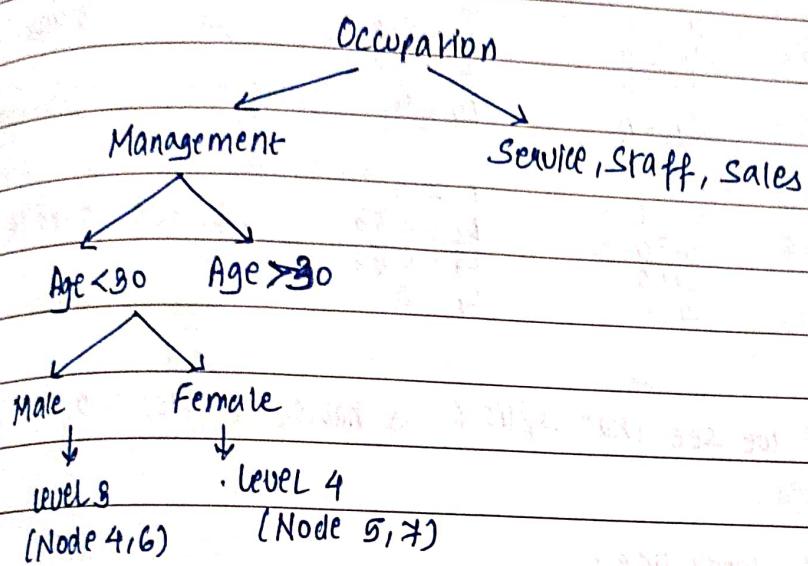
Similarly we calculate values for other splits and determine value of $\phi(S|t)$ which needs to be maximized.

	P_L	P_R	$P(j t_L)$	$P(j t_R)$	$\Phi(S t)$	$2PL \times PR$	$\phi(S t)$
Service 1. t_L	0.2727	0.7273	$L_1 = 0.333$ $L_2 = 0.33$ $L_3 = 0.33$ $L_4 = 0$	$L_1 = 0.125$ $L_2 = 0.25$ $L_3 = 0.375$ $L_4 = 0.25$	0.583	0.3966	0.2312
Mgmt 2. t_L	0.3636	0.6363	$L_1 = 0$ $L_2 = 0$ $L_3 = 0.5$ $L_4 = 0.5$	$L_1 = 0.285$ $L_2 = 0.428$ $L_3 = 0.285$ $L_4 = 0$	1.428	0.4626	0.66
Sales 3. t_L	0.1818	0.8181	$L_1 = 0$ $L_2 = 0.35$ $L_3 = 0.5$ $L_4 = 0$	$L_1 = 0.22$ $L_2 = 0.22$ $L_3 = 0.33$ $L_4 = 0$	0.889	0.2974	0.2643
Staff 4. t_L	0.1818	0.8181	$L_1 = 0.5$ $L_2 = 0.5$ $L_3 = 0$ $L_4 = 0$	$L_1 = 0.11$ $L_2 = 0.22$ $L_3 = 0.44$ $L_4 = 0.22$	1.833	0.2974	0.3964
Male 5. t_L	0.5454	0.4546	$L_1 = 0.33$ $L_2 = 0.33$ $L_3 = 0.33$ $L_4 = 0$	$L_1 = 0$ $L_2 = 0.20$ $L_3 = 0.40$ $L_4 = 0.40$	0.933	0.4958	0.4626
Age < 30 6. t_L	0.3636	0.6363	$L_1 = 0.5$ $L_2 = 0$ $L_3 = 0.5$ $L_4 = 0$	$L_1 = 0$ $L_2 = 0.428$ $L_3 = 0.285$ $L_4 = 0.28$	1.4285	0.4627	0.66
≥ 30							

P_L	P_R	$P(j t_L)$	$P(j t_R)$	$\Phi(S t)$	$2PLPK$	$\Phi(S t')$
0.6363	0.3636	$L_1 = 0.285$ $L_2 = 0.285$ $L_3 = 0.285$ $L_4 = 0.142$	$L_1 = 0$ $L_2 = 0.25$ $L_3 = 0.5$ $L_4 = 0.25$	0.642	0.4624	0.297

From above table, now $S & G$ has highest value of $\Phi(S|t)$
∴ Random Key : 2

∴ The decision tree looks like :-



Split	Left child node t_L	Right child node t_R
1	Occupation : Service	Occupation : Sales, Staff
2	Occupation : Sales	Occupation : Service, Staff
3	Occupation : Staff	Occupation : Service, Sales
4	Gender : Male	Gender : Female
5	Age < 30	Age >= 30
6	Age < 40	Age >= 40

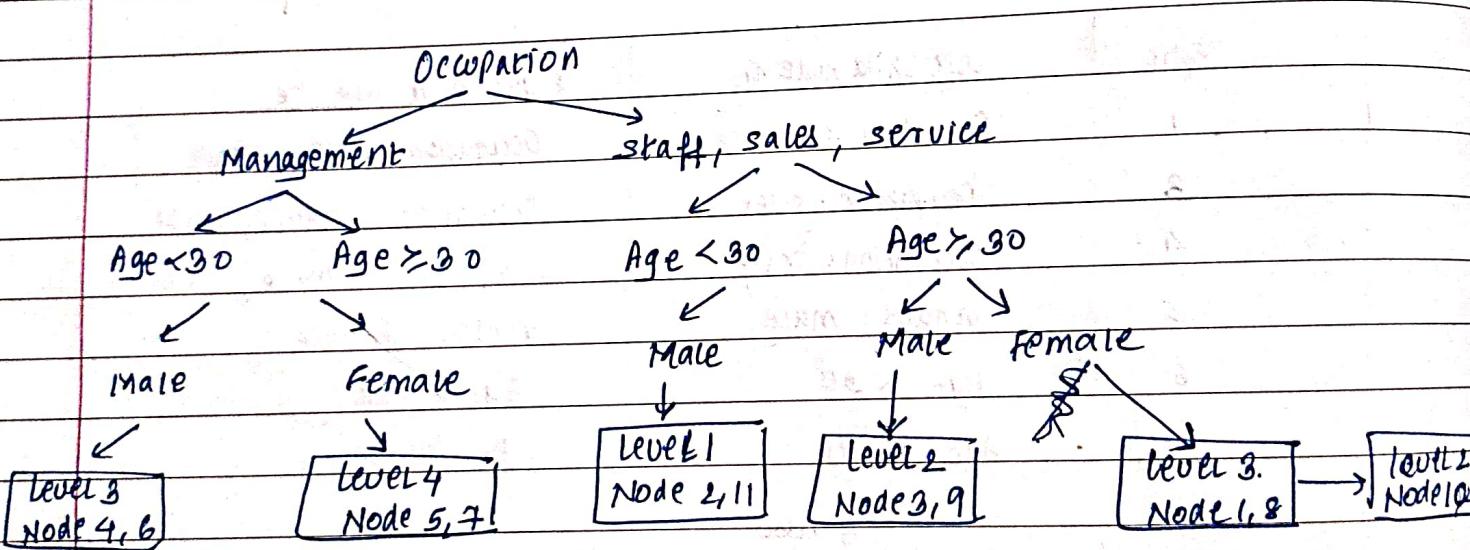
Calculate : $\Phi(S|t) = 2PLK \sum_{j=1}^{\# \text{classes}} |P(j|t_L) - P(j|t_R)|$

P_L	P_R	$P(j t_L)$	$P(j t_R)$	$\Phi(S t)$	$2PLPK$	$\Phi(S t')$
1. $3/7 = 0.428$	$4/7 = 0.5714$	$L_1 = 1/3 = 0.33$ $L_2 = 0.33$ $L_3 = 0.33$ $L_4 = 0$	$L_1 = 1/4 = 0.25$ $L_2 = 2/4 = 0.5$ $L_3 = 1/4 = 0.25$ $L_4 = 0$	0.332	0.489	0.1625
2. $2/7 = 0.2857$	$5/7 = 0.7142$	$L_1 = 0$ $L_2 = 1/2 = 0.5$ $L_3 = 1/2 = 0.5$ $L_4 = 0$	$L_1 = 2/5 = 0.4$ $L_2 = 2/5 = 0.4$ $L_3 = 1/5 = 0.2$ $L_4 = 0$	0.8	0.408	0.3264

	P_L	P_R	$P(L E_L)$	$P(L R_L)$	$\phi(G1L)$	$2PLR$	$\phi(G1R)$
4.	$2/7 = 0.2857$	$5/7 = 0.7142$	$L_1 = 1/2 = 0.5$ $L_2 = 1/2 = 0.5$ $L_3 = 0$ $L_4 = 0$	$L_1 = 1/5 = 0.2$ $L_2 = 2/5 = 0.4$ $L_3 = 2/5 = 0.4$ $L_4 = 0$	0.8	0.408	0.3264
5.	$6/11 = 0.5454$	0.4285	$L_1 = 0.5$ $L_2 = 0.5$ $L_3 = 0$ $L_4 = 0$	$L_1 = 0$ $L_2 = 0.66$ $L_3 = 0.66$ $L_4 = 0$	1.333 0.4896	0.4896 0.4896	0.6514
→ 6	0.2857	0.7142	$L_1 = 1$ $L_2 = 0$ $L_3 = 0$ $L_4 = 0$	$L_1 = 0$ $L_2 = 0.6$ $L_3 = 0.4$ $L_4 = 0$	2	0.408	0.816
7	0.5714	0.4285	$L_1 = 0.5$ $L_2 = 0.5$ $L_3 = 0$ $L_4 = 0$	$L_1 = 0$ $L_2 = 0.83$ $L_3 = 0.66$ $L_4 = 0$	1.333 0.4896	0.4896 0.4896	0.6514

Here in split 2nd: we see that split 6 is having $\phi(G1R) = 0.816$ highest among other values.

The decision tree looks like:



Decision Rules :-

If Occupation = Management and Age < 30 and gender = Male then L_3 (Node 4, 6)

If Occupation = Management and Age ≥ 30 and gender = Female then L_4 (Node 5, 7)

If Occupation = staff, sales, service and Age < 30 and gender = male then L_1 (Node 2, 11)

If occupation = staff, sales, service and Age ≥ 30 and gender = male	then L2 (Node 3, 9)
If occupation = staff, sales, service and Age ≥ 30 and gender = female	then L3 (Node 1, 8)
If occupation = staff, sales, service and Age ≥ 30 and gender = Female	then L2 (Node 10)

G. Construct a C4.5 decision tree to classify salary based on the other variables,
do as much as you can by hand, before turning to the software

split	child nodes (left)	child node (right)
1	occupation: service, occupation: management	occupation: sales, staff
2	Gender = male	Gender = Female
3	Age ≤ 25	Age > 25
4	Age ≤ 35	Age > 35
5	Age ≤ 45	Age > 45

$$H(T) = \text{Entropy}$$

levels	salary (target)	Probability
level 1	less than 85,000 \$	2/11
level 2	Between $\geq 85,000$ and $< 45,000$	8/11
level 3	Between $\geq 45,000$ and $< 55,000$	4/11
level 4	$\geq 55,000$ \$	2/11

$$\text{Entropy } H(X) = - \sum_j p_j \log_2(p_j)$$

Variable X has K values with probabilities $p_1, p_2 \dots p_K$
for variables with several outcomes, use weighted sum of $-\log_2(p)$
to transmit result

$$= -\frac{2}{11} \log_2\left(\frac{2}{11}\right) - \frac{8}{11} \log_2\left(\frac{8}{11}\right) - \frac{4}{11} \log_2\left(\frac{4}{11}\right) - \frac{2}{11} \log_2\left(\frac{2}{11}\right)$$

$$H(T) = -p_j \log_2(p_j) = 1.927 \text{ bits}$$

split on Occupation:

$$P(\text{Service}) = 3/11$$

$$P(\text{Management}) = 4/11$$

$$P(\text{Sales}) = 2/11$$

$$P(\text{Staff}) = 2/11$$

$$\text{Occupation : Service} : P(< 55,000) = 1/3 = P(> 55,000) = 0$$

$$\text{Occupation : Management} : P(< 55,000) = 1/2 = P(> 55,000) = 1/2$$

$$\text{Occupation : Sales} : P(< 55,000) = 1/2 = P(> 55,000) = 0$$

$$\text{Occupation : Staff} : P(< 55,000) = 0 = P(> 55,000) = 0$$

$$\text{Entropy of 4 branches + service, Management, Sales, Staff} \\ H(\text{Service}) = -1/3 \log(1/3) - 1/3 \log(1/3) - 1/3 \log(1/3) - 0 \log(0) = 1.58$$

~~$$H(\text{Management}) = -0 \log(0) - 0 \log(0) - 1/2 \log(1/2) - 1/2 \log(1/2) = 1$$~~

~~$$H(\text{Sales}) = -0 \log(0) - 1/2 \log(1/2) - 1/2 \log(1/2) - 0 \log(0) = 1$$~~

~~$$H(\text{Staff}) = -1/2 \log(1/2) - 1/2 \log(1/2) - 0 \log(0) - 0 \log(0) = 1$$~~

~~H_{total}~~

Combining Entropy of all 4 branches along with their proportion p_i

$$p_i = \frac{3}{11} * (1.58) + \frac{4}{11} * (1) + \frac{2}{11} * (1) + \frac{2}{11} * (1)$$

$$= 1.1572 \text{ bits}$$

$$K H_0(T) = p_i H_0(T_i) i=1$$

\therefore From the split on Occupation we have the attribute.

$$H(T) - T = 1.927 - 1.1572 = 0.7698 H_0 \text{ bits}$$

Now split for Gender!

$$P(\text{Male}) = 6/11 \quad P(\text{Female}) = 5/11$$

Gender \rightarrow Male

$$P(< 35,000) = 1/3$$

Gender \rightarrow Female

$$P(< 35,000) = 0$$

$$P(\geq 35,000 \& < 45,000) = 1/3$$

$$P(\geq 35,000 \& < 45,000) = 1/5$$

$$P(\geq 45,000 \& < 55,000) = 1/3$$

$$P(\geq 45,000 \& < 55,000) = 2/5$$

$$P(> 55,000) = 0$$

$$P(> 55,000) = 2/5$$

From above the Entropy

$$H(\text{Male}) = -\frac{1}{3} \log(1/3) - \frac{1}{3} \log(1/3) - \frac{1}{3} \log(1/3) - 0 \log(0) = 1.58$$

$$H(\text{Female}) = -0 \log(0) - \frac{1}{5} \log(1/5) - \frac{2}{5} \log(2/5) = 1.52$$

Combining Entropy of 2 branches along with their proportion p_i

$$\frac{6}{11} \times (1.58) + \frac{5}{11} (1.52) = 1.55 \text{ bits}$$

$$H_G(T) = p_i H_G(T_i) \quad i=1$$

$$\therefore H(\text{Gain}) = H(T) - T = 1.927 - 1.55 = 0.377 \text{ bits}$$

SPLIT ON AGE :-

$$P(Age \leq 35) = 7/11 \Rightarrow P(Age > 35) = 4/11$$

$$\boxed{\text{Age } \leq 35} \implies P(Age \leq 35, Income) = 2/7$$

$$P(Age \leq 35, Income \in [35,000] \cup [45,000]) = 2/7$$

$$P(Age \leq 35, Income \in [45,000] \cup [55,000]) = 2/7$$

(i) Age ≥ 35

$$P(Age \geq 35, Income) = 0$$

$$P(Age \geq 35, Income \in [35,000] \cup [45,000]) = 1/4$$

$$P(Age \geq 35, Income \in [45,000] \cup [55,000]) = 2/4$$

$$P(Age \geq 35, Income \in [55,000]) = 1/4$$

Entropy of from the above 2 branches Age ≤ 35 & Age > 35

$$H(Age \leq 35) = -\frac{2}{7} \log\left(\frac{2}{7}\right) - \frac{2}{7} \log\left(\frac{2}{7}\right) - \frac{2}{7} \log\left(\frac{2}{7}\right) - \\ 1/7 \log(1/7) = 1.93$$

$$H(Age > 35) = -1/4 \log(1/4) - 1/4 \log(1/4) - 2/4 \log(2/4) = 1.5$$

Combining entropy of above 2 branches, p_i

$$= \frac{7}{11} \times (1.93) + \frac{4}{11} (1.5) = 1.77 \text{ bits}$$

K

$$H_{\text{age}}(T) = p_i H_{\text{age}}(T_i) \geq$$

\therefore from the info gained by split on gender attribute is

$$H(T) - H_{\text{Gender}}(T) = 1.927 - 1.77$$

$$= \boxed{0.157 \text{ bits}}$$

split on Age: $P(C \leq 45)$

$$= 10/11 \Rightarrow P(C > 45) = 1/11$$

$$\text{Age: } C \leq 45 = P(C < 35,000) = 2/10$$

$$P(C \geq 35,000 \text{ and } < 45,000) = 2/10$$

$$P(C \geq 45,000 \text{ and } < 55,000) = 4/10$$

$$P(C \geq 55,000) = 2/10$$

Age: $C > 45$

$$P(C \leq 35,000) = 0$$

$$P(C \geq 35,000 \text{ and } < 45,000) = 1/11$$

$$P(C \geq 45,000 \text{ and } < 55,000) = 0$$

$$P(C \geq 55,000) = 0$$

Entropy for above two: Age ≤ 45 and Age > 45

$$H(\text{Age} \leq 45) = -2/10 \log(2/10) - 2/10 \log(2/10) - 4/10 \log(4/10)$$

$$-2/10 \log(2/10) = 1.92$$

$$H(\text{Age} > 45) = -1/11 \log(1/11) = 0.3136$$

Combining the Entropy of 2 branches, along with their proportion P_i

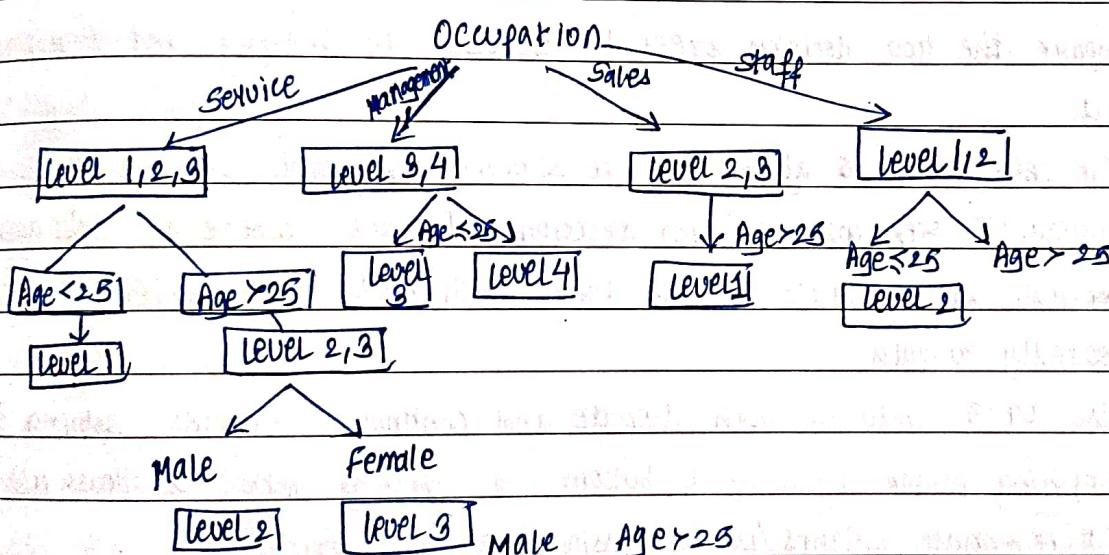
$$= 10/11 * 1.92 + 1/11 (0.3136) = 1.7735 \text{ bits}$$

\therefore from above values gained by split on gender $H(T) - H_{\text{Age}}(T)$

$$= 1.927 - 1.7735 = \boxed{0.1465 \text{ bits}}$$

Condition split	child node	Information gain
1	Occupation : Service, Management, Sales, Staff	0.7689 bits
2	Gender = Male, Female	0.347 bits
3	Age ≤ 25 & Age > 25	0.547 bits
4	Age ≤ 35 and Age > 35	0.157 bits
5	Age ≤ 45 and Age > 45	0.1465 bits

4.5 Decision Tree :-



Decision Rules :-

Antecedent	Consequent	Support	Confidence
1 if occupation = service ↳ Age ≤ 25	Then L1	3/11	1
2 if occupation = service and Age > 25 and male	Then L2	1/3	1
3 if occupation = service and Age > 25 and female	Then L3	1/3	1

Antecedent	Consequent	Support	Confidence
4. if occupation = Management and Age ≤ 25	Then L9	3/11	1
5. if occupation = Management and Age > 25	Then L4	8/11	1
6. if occupation = sales and Age ≥ 25 and Male	Then L2	1/2	1
7. if occupation = sales and Age > 25 and Female	Then L3	1/2	1
8. if occupation = staff and Age ≤ 25	Then L1	3/11	1
9. if occupation = staff and Age > 25	Then L2	8/11	1

7. Compare the two decision trees & discuss the benefits and drawbacks of each.

→ The CART and C4.5 algorithms are similar, however C4.5 uses an intermediate step when building decision rules set, whereas CART uses numerical split criteria to split data which needs to be applied repeatedly to data.

→ Also C4.5 includes both discrete and continuous features, solving the overfitting problem by using a bottom-up pruning method & allows users to differentiate weights/values assigned to each attribute.

8. Generate the full-set of decision rules for the CART decision tree

- | | |
|---|---------------------|
| 1. if occupation = Management and Age < 30
and Male | then L9 (Node 4, 6) |
| 2. If occupation = Management and Age ≥ 30
and female | then L4 (Node 5, 7) |

if occupation = staff, sales, service & Age
and Age ≥ 30 and Male

then L1 (Node 2, 11)

then L2 (Node 3, 9)

if occupation = staff, sales, service
and Age ≥ 30 & Female

then L3 (Node 1, 8)

if occupation = staff, sales, service
and Age ≥ 30 and female and staff

then L2 (Node 10)

9. Generate the full set of decision rules for the C4.5 decision tree.

Antecedent	Consequent	Support	confidence
1. if occupation = Service and Age ≤ 25	Then L1	3/11	1
2. if occupation = Service and Age > 25 & Male	Then L2	4/3	1
3. if occupation = Service and Age > 25 & Female	Then L3	1/3	1
4. if occupation = Mgmt and Age > 25	Then L2	3/11	1
5. if occupation = Mgmt & Age > 25	Then L4	8/11	1
6. if occupation = Sales & Age > 25 & Male	Then L2	1/2	1
7. if occupation = sales Age > 25 & Female	Then L3	1/2	1
8. if occupation = staff Age ≤ 25	Then L1	3/11	1
9. if occupation = staff & Age > 25	Then L2	8/11	1

10. Evaluate the advantages & disadvantages of the two sets of decision rules

CART advantages :-

- (i) It builds binary trees, in which each internal node has exactly two outgoing branches
- (ii) Using the two criteria, splits are chosen & the resulting tree is pruned using cost complexity method
- (iii) CART is able to handle variables that are both numerical & categorical & it does so with outliers (works fine).
- (iv) It has more general decision making rules

CART disadvantages :-

- (i) It can only split based on one variable.
- (ii) The decision tree growth can be unstable.

C4.5 advantages :-

- (i) The algorithm inherently employs single pass pruning process to mitigate Overfitting
- (ii) It can work with both Discrete & continuous data.
- (iii) C4.5 can handle the issue of incomplete data very well.

C4.5 disadvantages :-

- 1. The algorithm suffers from Overfitting
- 2. It has poor attribute split technique, inability to handle continuous valued & missing valued attributes with high learning cost.

— x-x — END OF ASSIGNMENT — x —