

7.2 Learning Decision Trees (continued)

➤ 3 parameters are computed:

- a) **Information Content** or **Entropy** regarding the classes in the samples
- b) **Remainder** or Expected **Remaining Entropy** of an **Attribute Test**
- c) **Information Gain** of an **Attribute Test**

a) $I(P(C_1), P(C_2), \dots, P(C_m)) = \sum_{i=1:m} -P(C_i) \times \log_2(P(C_i))$,
where $P(C_i)$ – probability of the class C_i ,
 m – number of classes.

For our example,

$$I(P(\text{positive}), P(\text{negative})) = \\ -P(\text{positive}) \times \log_2(P(\text{positive})) - P(\text{negative}) \times \log_2(P(\text{negative}))$$

<i>ID</i>	<i>Age</i>	<i>Income</i>	<i>Student</i>	<i>Credit Rating</i>	<i>Decision/ Class/ Label</i>
1	≤ 30	high	no	fair	negative
2	≤ 30	high	no	excellent	negative
3	31...40	high	no	fair	positive
4	> 40	medium	no	fair	positive
5	> 40	low	yes	fair	positive
6	> 40	low	yes	excellent	negative
7	31...40	low	yes	excellent	positive
8	≤ 30	Medium	no	fair	negative
9	≤ 30	low	yes	fair	positive
10	> 40	medium	yes	fair	positive
11	≤ 30	medium	yes	excellent	positive
12	31...40	medium	no	excellent	positive
13	31...40	high	yes	fair	positive
14	> 40	medium	no	excellent	negative

➤ That is, $I(9/14, 5/14) =$
 $- 9/14 \times \log_2 (9/14) -$
 $5/14 \times \log_2 (5/14) \approx$
 0.940 (bits)

b. Remainder (A) =

$\sum_{i=1:v} P(\text{Samples with } i\text{th value of } A) \times I(\text{Classes in samples with } i\text{th value of } A),$
 where $v = \text{number of distinct values of } A.$

$$\begin{aligned}
 \text{i. Remainder(Age)} &= 5/14 \times I(2/5, 3/5) && [\leq 30] \\
 &+ 4/14 \times I(4/4, 0/4) && [31 \dots 40] \\
 &+ 5/14 \times I(3/5, 2/5) && [> 40]
 \end{aligned}$$

<i>ID</i>	<i>Age</i>	<i>Income</i>	<i>Student</i>	<i>Credit Rating</i>	<i>Decision/ Class/ Label</i>
1	≤ 30	high	no	fair	negative
2	≤ 30	high	no	excellent	negative
3	31...40	high	no	fair	positive
4	> 40	medium	no	fair	positive
5	> 40	low	yes	fair	positive
6	> 40	low	yes	excellent	negative
7	31...40	low	yes	excellent	positive
8	≤ 30	Medium	no	fair	negative
9	≤ 30	low	yes	fair	positive
10	> 40	medium	yes	fair	positive
11	≤ 30	medium	yes	excellent	positive
12	31...40	medium	no	excellent	positive
13	31...40	high	yes	fair	positive
14	> 40	medium	no	excellent	negative

➤ That is,
 $\text{Remainder}(\text{Age}) = 2 \times 5/14 \times I(2/5, 3/5) \approx 0.694 \text{ (bits)}$

ii) $\text{Remainder}(\text{Income}) = ?$

iii) $\text{Remainder}(\text{Student}) = ?$

iv) $\text{Remainder}(\text{Credit Rating}) = ?$

c) $\text{Gain}(A) = I(\text{Classes in all samples in the table}) - \text{Remainder}(A)$

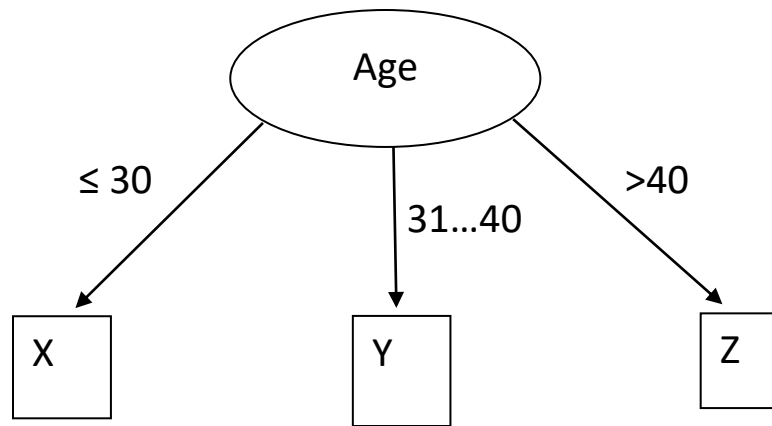
$\text{Gain}(\text{Age}) = 0.940 - 0.694 = 0.246 \text{ (bits)}$

Check that, $\text{Gain}(\text{Income}) = 0.029 \text{ bits,}$

$\text{Gain}(\text{Student}) = 0.151 \text{ bits, and}$

$\text{Gain}(\text{Credit Rating}) = 0.048 \text{ bits.}$

So, we have the attribute 'Age' with the highest Gain, and this leads to what follows.



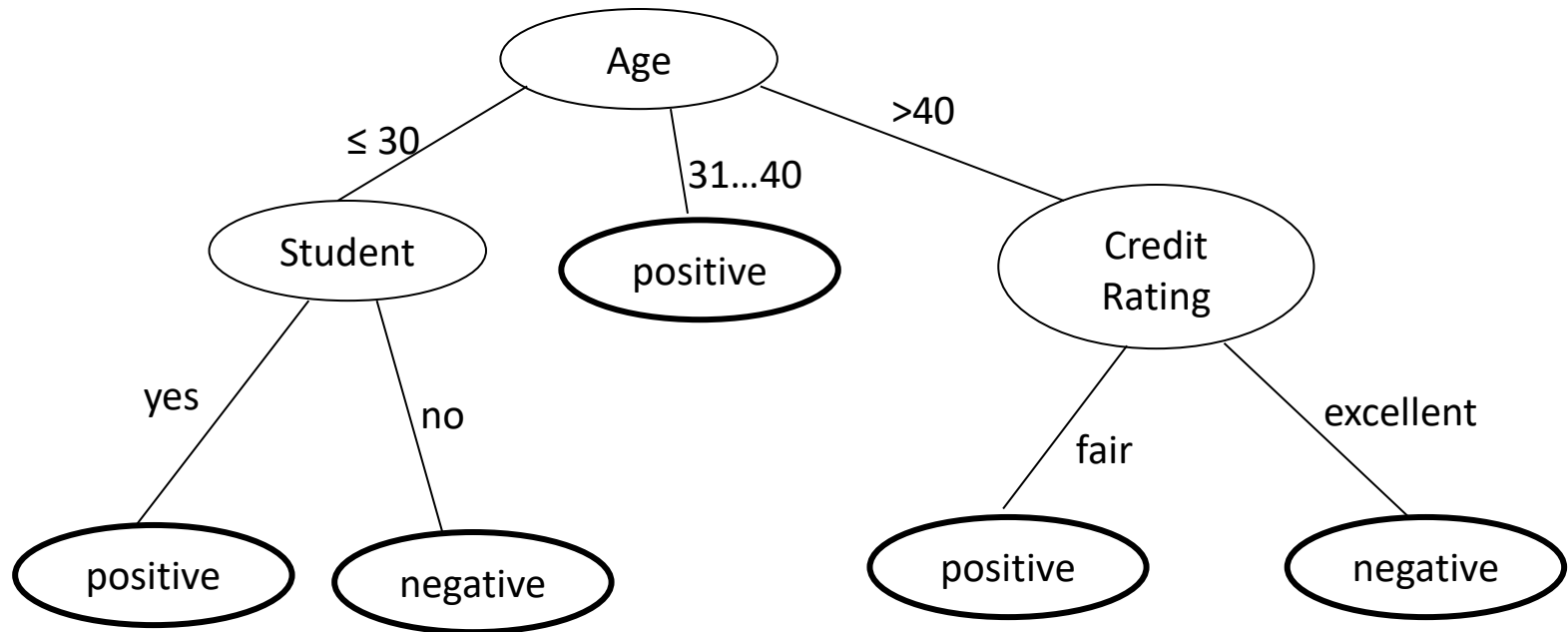
X =

<i>Age</i>	Income	Student	Credit Rating	<i>Decision/ Class/ Label</i>
≤ 30	high	no	fair	negative
≤ 30	high	no	excellent	negative
≤ 30	medium	no	fair	negative
≤ 30	low	yes	fair	positive
≤ 30	medium	yes	excellent	positive

Y = ?

Z = ?

Finally, we have the following tree. [Self study]



And it means that we have **learned the following 5 rules**.

1. If 'Age' = '≤ 30' and 'Student' = 'yes', then 'Class' = 'Buys a computer'.
2. If 'Age' = '≤ 30' and 'Student' = 'no', then 'Class' = 'Does not buy a computer'.
3.
4.
5. If 'Age' = '>40' and 'Credit Rating' = 'excellent', then 'Class' = 'Does not buy a computer'.