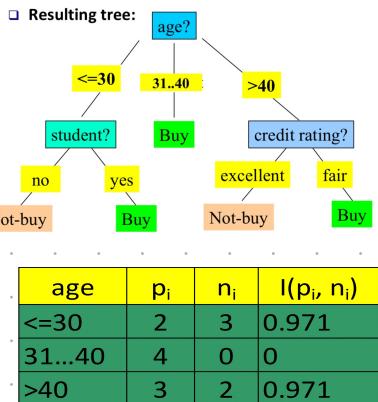


## Attribute selection with Information gain



Training data set: Who buys computer?

age	income	student	credit_rating	buys_computer
$\leq 30$	high	no	fair	no
$\leq 30$	high	no	excellent	no
$31..40$	high	no	fair	yes
$>40$	medium	no	fair	yes
$>40$	low	yes	fair	yes
$>40$	low	yes	excellent	no
$31..40$	low	yes	excellent	yes
$\leq 30$	medium	no	fair	no
$\leq 30$	low	yes	fair	yes
$>40$	medium	yes	fair	yes
$\leq 30$	medium	yes	excellent	yes
$31..40$	medium	no	excellent	yes
$31..40$	high	yes	fair	yes
$>40$	medium	no	excellent	no

class P: buys\_computer = "yes"  
 class N: buys\_computer = "no"  
 Yes : 9  
 No : 5

Age  
 $\leq 30$  yes: 2 / No: 3  
 $31..40$  yes: 4 / No: 0  
 $>40$  yes: 3 / No: 2

income  
 high yes: 2 / No: 2  
 medium yes: 4 / No: 2  
 low yes: 3 / No: 1

student  
 Yes yes: 6 / No: 1  
 No yes: 3 / No: 4

Credit  
 Fair yes: 6 / No: 2  
 Excellent yes: 3 / No: 3

$$\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) = 0.940$$

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$= \frac{5}{14} \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) \right) + \frac{4}{14} \left( -\frac{4}{4} \log_2 \frac{4}{4} \right) + \frac{5}{14} \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) = 0.694$$

$$\text{Info}_{\text{income}}(D) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1)$$

$$= \frac{4}{14} \left( -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) + \frac{6}{14} \left( -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) + \frac{4}{14} \left( -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) = 0.911$$

$$\text{Info}_{\text{student}}(D) = \frac{1}{14} I(6,1) + \frac{7}{14} I(3,4)$$

$$= \frac{1}{14} \left( -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right) + \frac{7}{14} \left( -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) = 0.989$$

$$\text{Info}_{\text{credit}}(D) = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$

$$= \frac{8}{14} \left( -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) + \frac{6}{14} \left( -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right) = 0.892$$

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246$$

$$\text{Gain}(\text{income}) = \text{Info}(D) - \text{Info}_{\text{income}}(D) = 0.940 - 0.911 = 0.029$$

$$\text{Gain}(\text{student}) = \text{Info}(D) - \text{Info}_{\text{student}}(D) = 0.940 - 0.989 = 0.151$$

$$\text{Gain}(\text{credit\_rating}) = \text{Info}(D) - \text{Info}_{\text{credit\_rating}}(D) = 0.940 - 0.892 = 0.048$$

L = 30

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
<=30	medium	yes	excellent	yes

$$\text{Sol}^n \quad \text{Info}(D) = I(2,3) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.991$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) \\ &= \frac{2}{5} \left( -\frac{2}{2} \log_2 \frac{2}{2} \right) + \frac{2}{5} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{1}{5} \left( -\frac{1}{1} \log_2 \frac{1}{1} \right) \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}}(D) &= \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) \\ &= \frac{2}{5} \left( -\frac{2}{2} \log_2 \frac{2}{2} \right) + \frac{3}{5} \left( -\frac{3}{3} \log_2 \frac{3}{3} \right) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{credit}} &= \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1) \\ &= \frac{3}{5} \left( -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) + \frac{2}{5} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) \\ &= 0.951 \end{aligned}$$

$$\text{Gain}_{\text{income}}(D) = \text{Info}(D) - \text{Info}_{\text{income}}(D) = 0.991 - 0.4 = 0.591 \quad \text{Gain zinsen}$$

$$\text{Gain}_{\text{student}}(D) = \text{Info}(D) - \text{Info}_{\text{student}}(D) = 0.991 - 0 = 0.991$$

$$\text{Gain}_{\text{credit\_rating}}(D) = \text{Info}(D) - \text{Info}_{\text{credit\_rating}}(D) = 0.991 - 0.951 = 0.02$$

Age

yes: 2 / No: 3

Income

high yes: 0 / No: 2

Medium yes: 1 / No: 1

low yes: 1 / No: 0

student

yes: 2 / No: 3

Credit

fair yes: 1 / No: 2

excellent yes: 1 / No: 1

age	income	student	credit_rating	buys_computer
31...40	high	no	fair	yes
31...40	low	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes

## Age

yes: 4 / No: 0

## Income

High yes: 2 / No: 0

Medium yes: 1 / No: 0

Low yes: 1 / No: 0

## Student

Yes

yes: 2 / No: 0

No

yes: 2 / No: 0

## Credit

Fair

yes: 2 / No: 0

Excellent:

yes: 2 / No: 0

Training data set: Who buys computer?

age	income	student	credit	rating	buys_computer
>40	medium	no	fair		yes
>40	low	yes	fair		yes
>40	low	yes	excellent		no
>40	medium	yes	fair		yes
>40	medium	no	excellent		no

Age  
yes: 3 / No: 2

Income

High yes: 0 / No: 0  
Medium yes: 2 / No: 1  
Low yes: 1 / No: 1

Student

Yes  
yes: 2 / No: 1  
No  
yes: 1 / No: 1

Credit

Fair  
yes: 3 / No: 0  
Excellent  
yes: 0 / No: 2

Sol<sup>n</sup>  $\text{Info}(D) = I(3,2) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.991$

$$\begin{aligned}\text{Info}_{\text{income}}(D) &= \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) \\ &= \frac{3}{5} \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{2}{5} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) \\ &= 0.951\end{aligned}$$

$$\begin{aligned}\text{Info}_{\text{student}}(D) &= \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) \\ &= \frac{3}{5} \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{2}{5} \left( -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) \\ &= 0.951\end{aligned}$$

$$\begin{aligned}\text{Info}_{\text{credit}}(D) &= \frac{3}{5} I(3,0) + \frac{2}{5} I(0,2) \\ &= \frac{3}{5} \left( -\frac{2}{3} \log_2 \frac{2}{3} \right) + \frac{2}{5} \left( -\frac{1}{2} \log_2 \frac{1}{2} \right) \\ &= 0\end{aligned}$$

$$\text{Gain}_{\text{income}}(D) = \text{Info}(D) - \text{Info}_{\text{income}}(D) = 0.991 - 0.951 = 0.02$$

$$\text{Gain}_{\text{student}}(D) = \text{Info}(D) - \text{Info}_{\text{student}}(D) = 0.991 - 0.951 = 0.02$$

$$\text{Gain}_{\text{credit\_rating}}(D) = \text{Info}(D) - \text{Info}_{\text{credit\_rating}}(D) = 0.991 - 0 = 0.991$$

## Training data set: Who buys computer?

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
<=30	medium	yes	excellent	yes

age	income	student	credit_rating	buys_computer
31...40	high	no	fair	yes
31...40	low	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes

age	income	student	credit_rating	buys_computer
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
>40	medium	yes	fair	yes
>40	medium	no	excellent	no

Gain student = 0.991  
" តាន់ Gain នៅក្នុង "

student	buys_computer
no	no
no	no
no	no
yes	yes
yes	yes

No  
Not buy

Yes  
Buy

Buy

Gain credit-rating = 0.991  
" តាន់ Gain នៅក្នុង "

credit_rating	buys_computer
fair	yes
fair	yes
excellent	no
fair	yes
excellent	no

Excellent  
Not buy

fair  
Buy