

# Lecture Three

---

## Decision Tree Example

# Example

age	income	student	credit_rating	sys_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

# Tree induction example

---

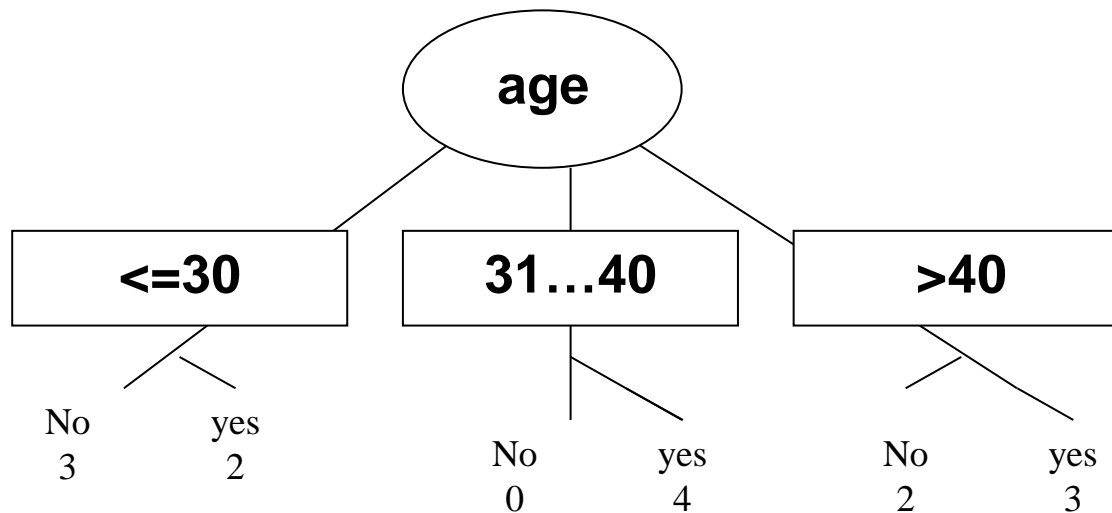
$$\bullet \text{info}(\mathbf{D}) = - p^+ \log_2 p^+ - p^- \log_2 p^-$$

$$\begin{aligned} \text{Info}(\mathbf{D}) &= -((9/14)\log_2(9/14)) - ((5/14)\log_2(5/14)) \\ &= 0.940 \end{aligned}$$

$D[9+, 5-]$   
age

$\leq 30$  [2+,3-]  
 $31 \dots 40$  [4+,0-]  
 $> 40$  [3+,2-]

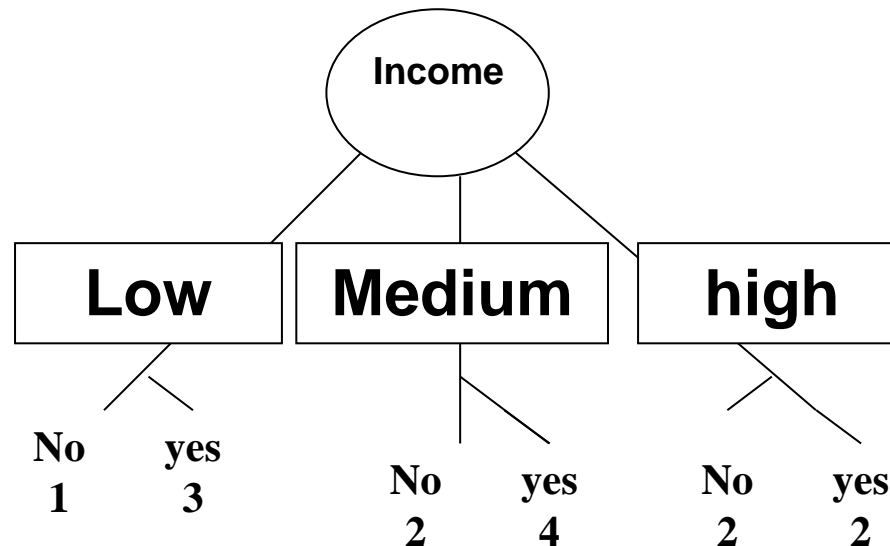
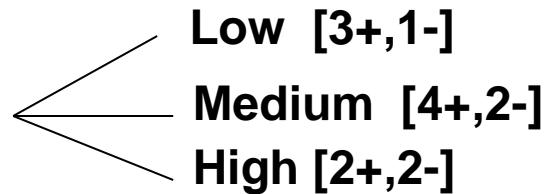
---



$$\begin{aligned}
 Info_{age}(D) = & 5/14((-2/5)\log_2(2/5)) - ((3/5)\log_2(3/5)) \\
 & + 4/14((-4/4)\log_2(4/4)) - ((0/4)\log_2(0/4)) \\
 & + 5/14((-3/5)\log_2(3/5)) - ((2/5)\log_2(2/5)) \\
 = & 0.694 \text{ bits}
 \end{aligned}$$

$$Gain(\mathbf{age}) = Info(D) - Info_{age}(D) = 0.940 - 0.694 = 0.246 \text{ bits.}$$

D[9+, 5-]  
Income



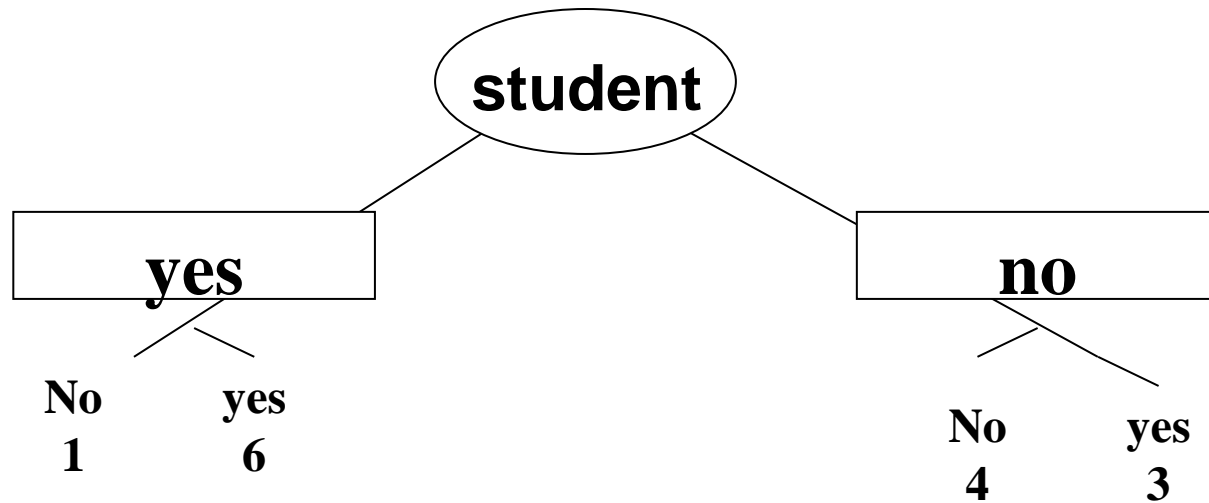
Info gain Income ( D) =

$$\begin{aligned}
 & 4/14((-3/4) \log_2(3/4)) - ((1/4) \log_2(1/4)) \\
 & + 6/14((-4/6) \log_2(4/6)) - ((2/6) \log_2(2/6)) \\
 & + 4/14((-2/4) \log_2(2/4)) - ((2/4) \log_2(2/4)) \\
 & = \mathbf{0.91104 \text{ bits}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain Income} &= \text{Info Gain (D)} - \text{Info gain Income (D)} \\
 &= 0.940 - 0.91194 = \mathbf{0.029 \text{ bits}}
 \end{aligned}$$

D[9+, 5-] student 
 $\swarrow$  No [3+, 4-]  
 $\searrow$  Yes [6+, 1-]

---



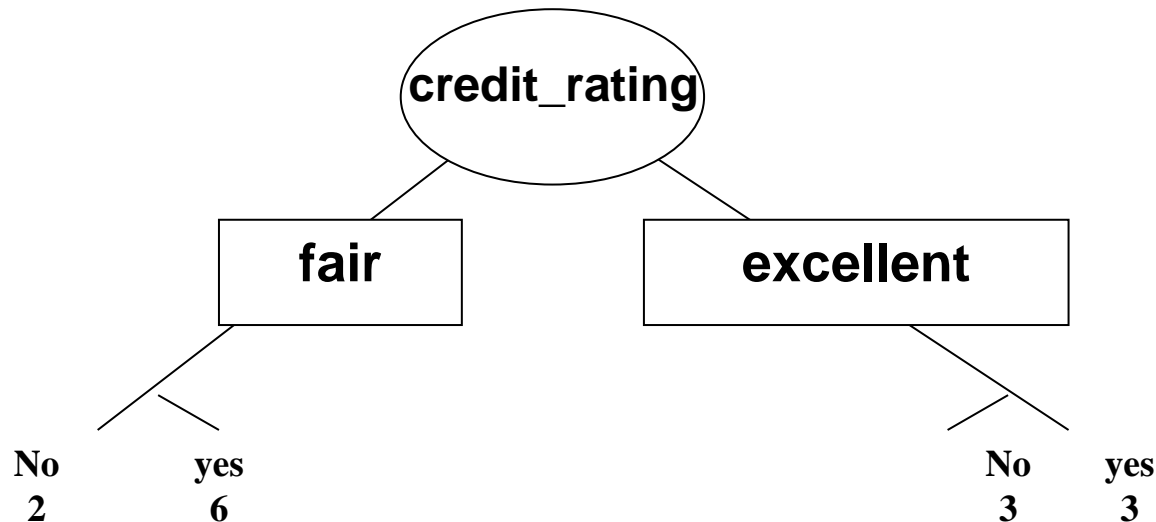
$$\begin{aligned}
 \text{Info gain student}(D) &= \frac{7}{14}((-4/7)\log_2(4/7)) - ((3/7)\log_2(3/7)) \\
 &\quad + \frac{7}{14}((-6/7)\log_2(6/7)) - ((1/7)\log_2(1/7)) \\
 &= \mathbf{0.7884 \text{ bits.}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain student} &= \text{InfoGain}(D) - \text{InfoGain student}(D) \\
 &= \mathbf{0.940} - 0.7884 = \mathbf{0.151 \text{ bits.}}
 \end{aligned}$$

D[9+, 5-] credit\_rating

- Fair [6+, 2-]
- excellent[3+, 3-]

---



$$\begin{aligned} \text{Info gain credit\_rating (D)} &= 8/14((-6/8)\log_2(6/8)) - ((2/8)\log_2(2/8)) \\ &\quad + 6/14((-3/6)\log_2(3/6)) - ((3/6)\log_2(3/6)) \\ &= 0.892 \text{ bits.} \end{aligned}$$

$$\begin{aligned} \text{Gain} &= \text{InfoGain (D)} - \text{Info gain credit\_rating (D)} \\ &= 0.940 - 0.892 = 0.048 \text{ bits.} \end{aligned}$$

age	income	student	credit_rating	comp
<=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

Gain      = 0.246      = 0.029      = 0.151      = 0.048



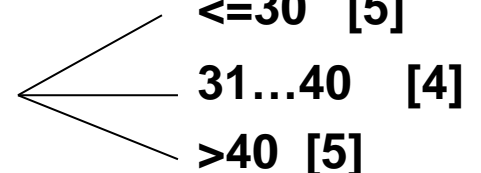
## □ Gain Ratio for Attribute Selection (C4.5)

---

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$


split info    بشوف عندي كام متغير في كل عمود وبحسب علي اساسه

$$SplitInfo_{Age}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

**D[14] age** 

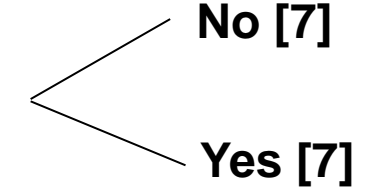
$$SplitInfo_{Age}(D) = - \frac{5}{14} \times \log_2 \left( \frac{5}{14} \right) - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{5}{14} \times \log_2 \left( \frac{5}{14} \right) = 1.5774$$

$$SplitInfo_{Income}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

**D[14] Income** 

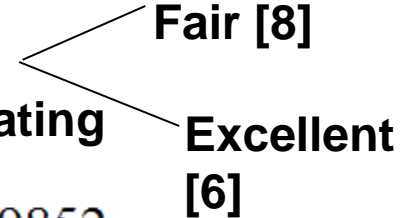
$$SplitInfo_{income}(D) = - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) - \frac{4}{14} \times \log_2 \left( \frac{4}{14} \right) = 1.557$$

$$SplitInfo_{student}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

**D[14] student** 

$$SplitInfo_{student}(D) = - \frac{7}{14} \times \log_2 \left( \frac{7}{14} \right) - \frac{7}{14} \times \log_2 \left( \frac{7}{14} \right) = 1$$

$$SplitInfo_{credit\_rating}(D) = - \sum_{i=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

**D[14] credit\_rating** 

$$SplitInfo_{credit\_rating}(D) = - \frac{8}{14} \times \log_2 \left( \frac{8}{14} \right) - \frac{6}{14} \times \log_2 \left( \frac{6}{14} \right) = 0.9852$$

age	income	student	credit_rating	correct
<=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

Gain            = 0.246            = 0.029            = 0.151            = 0.048  
 split Info   = 1.5774            = 1.577            = 1            = 0.9852

# $\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$

---

	Age	Income	Student	credit rating
Gain	0.246	0.029	0.151	0.048
Split Info	1.5774	1.577	1	0.9852

$$\text{Gain Ratio}(\text{Age}) = \text{Gain}(\text{age}) / \text{Split Info}(\text{age})$$

$$\text{Gain Ratio}(\text{Age}) = 0.246 / 1.5774 = 0.1559$$

$$\text{Gain Ratio}(\text{Income}) = \text{Gain}(\text{income}) / \text{Split Info}(\text{income})$$

$$\text{Gain Ratio}(\text{income}) = 0.029 / 1.577 = 0.019$$

# $\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$

---

	Age	Income	Student	credit rating
Gain	0.246	0.029	0.151	0.048
Split Info	1.5774	1.577	1	0.9852

$$\text{Gain Ratio}(\text{Student}) = \text{Gain}(\text{Student}) / \text{Split Info}(\text{Student})$$

$$\text{Gain Ratio}(\text{Student}) = 0.151 / 1 = 0.151$$

$$\text{Gain Ratio}(\text{credit\_rating}) = \text{Gain}(\text{credit\_rating}) / \text{Split Info}(\text{credit\_rating})$$

$$\text{Gain Ratio}(\text{credit\_rating}) = 0.048 / 0.9852 = 0.048782$$

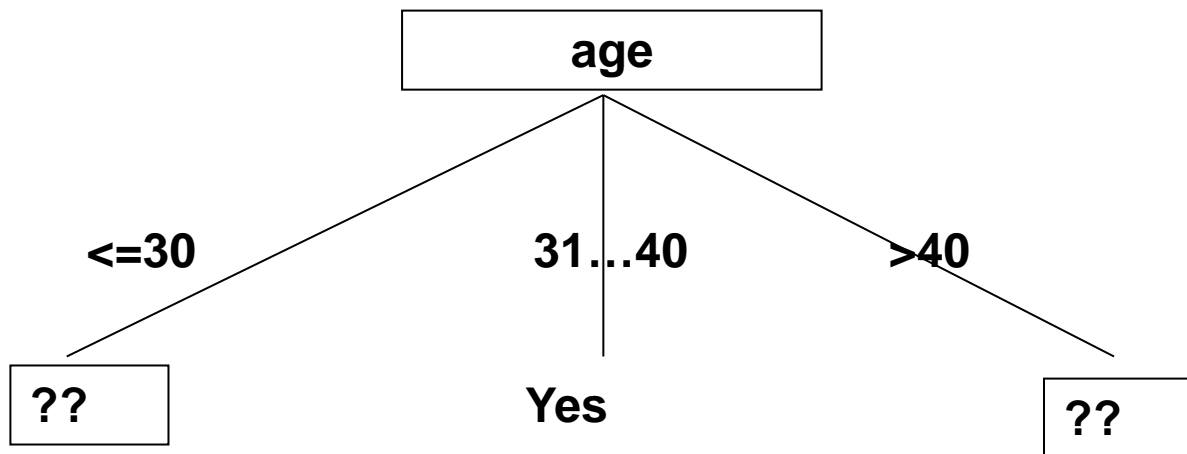
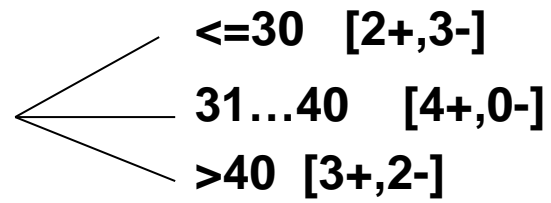
# $\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$

---

	Age	Income	Student	credit rating
Gain	0.246	0.029	0.151	0.048
Split Info	1.5774	0.926	1	0.9852
Gain Ratio	0.1559	0.019	0.151	0.048782

Age attribute with the maximum gain ratio is selected as the splitting attribute

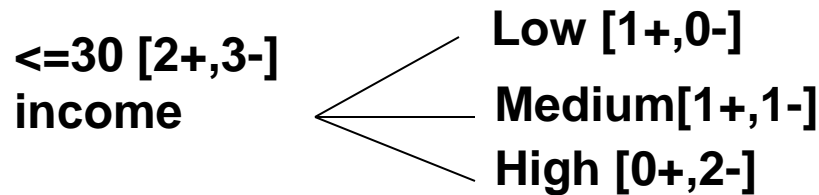
D[9+, 5-]  
age



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

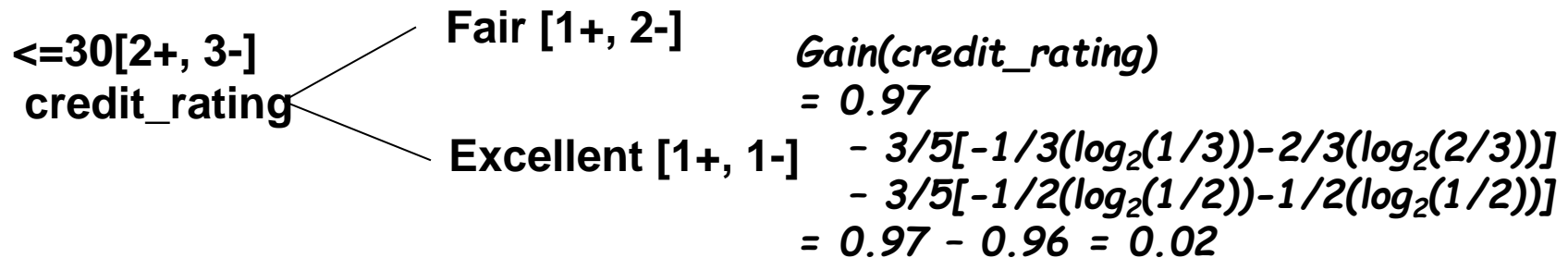
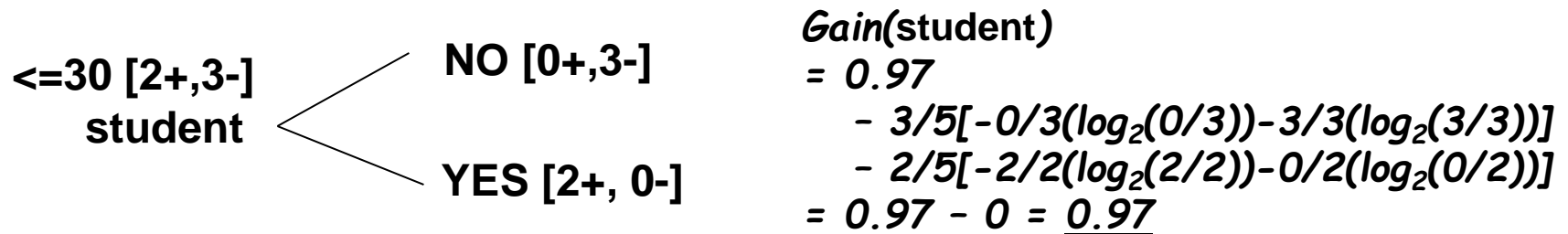
$$\begin{aligned}
 \text{Info}(\leq 30) &= -2/5(\log_2(2/5)) \\
 &\quad -3/5(\log_2(3/5)) \\
 &= 0.97
 \end{aligned}$$





---


$$\begin{aligned}
 \text{Gain}(\text{income}) &= 0.97 - 1/5[-1/1(\log_2(1/1)) - 0/1(\log_2(0/1))] \\
 &\quad - 2/5[-1/2(\log_2(1/2)) - 1/2(\log_2(1/2))] \\
 &\quad - 2/5[-0/2(\log_2(0/2)) - 2/2(\log_2(2/2))] \\
 &= 0.97 - 0.4 = 0.57
 \end{aligned}$$



age	income	student	credit_rating	correct
<=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

Gain

= **0.57**

= **0.97**

= **0.02**

$$SplitInfo_{Income \leq 30}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} \mathbf{D[5]} \\ \mathbf{Income} \end{array} \begin{array}{l} \text{Low [1]} \\ \text{Medium [2]} \\ \text{High [2]} \end{array}$$


---

$$= -\frac{1}{5} \times \log_2 \left( \frac{1}{5} \right) - \frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) - \frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) = \mathbf{0.993}$$

$$SplitInfo_{student \leq 30}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} \mathbf{D[5]} \\ \mathbf{student} \end{array} \begin{array}{l} \text{No [3]} \\ \text{Yes [2]} \end{array}$$

$$= -\frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) = \mathbf{0.970}$$

$$SplitInfo_{credit\_rating \leq 30}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} \mathbf{D[5]} \\ \mathbf{credit\_rating} \end{array} \begin{array}{l} \text{Fair [3]} \\ \text{Excellent [2]} \end{array}$$

$$= -\frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) = \mathbf{0.970}$$

age	income	student	credit_rating	correct
<=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

Gain                      = **0.57**                      = **0.97**                      = **0.02**  
 split Info                = **0.993**                      = **0.970**                      = **0.970**

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$


---

	Age	Income	Student	credit rating
Gain	<= 30	0.057	0.97	0.02
Split Info		0.993	0.97	0.97

$$\text{Gain Ratio}(\text{Income} \leq 30) = \text{Gain}(\text{income}) / \text{Split Info}(\text{income})$$

$$\text{Gain Ratio}(\text{Income} \leq 30) = 0.057 / 0.993 = 0.0574$$

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$


---

	Age	Income	Student	credit rating
Gain	<= 30	0.057	0.97	0.02
Split Info		0.993	0.97	0.97

$$\text{Gain Ratio}(\text{Student} \leq 30) = \text{Gain}(\text{Student}) / \text{Split Info}(\text{Student})$$

$$\text{Gain Ratio}(\text{Student} \leq 30) = 0.97 / 0.97 = 1$$

$$\text{Gain Ratio}(\text{credit\_rating} \leq 30) = \text{Gain}(\text{credit\_rating}) / \text{Split Info}(\text{credit\_rating})$$

$$\text{Gain Ratio}(\text{credit\_rating}) = 0.02 / 0.97 = 0.0206$$

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$

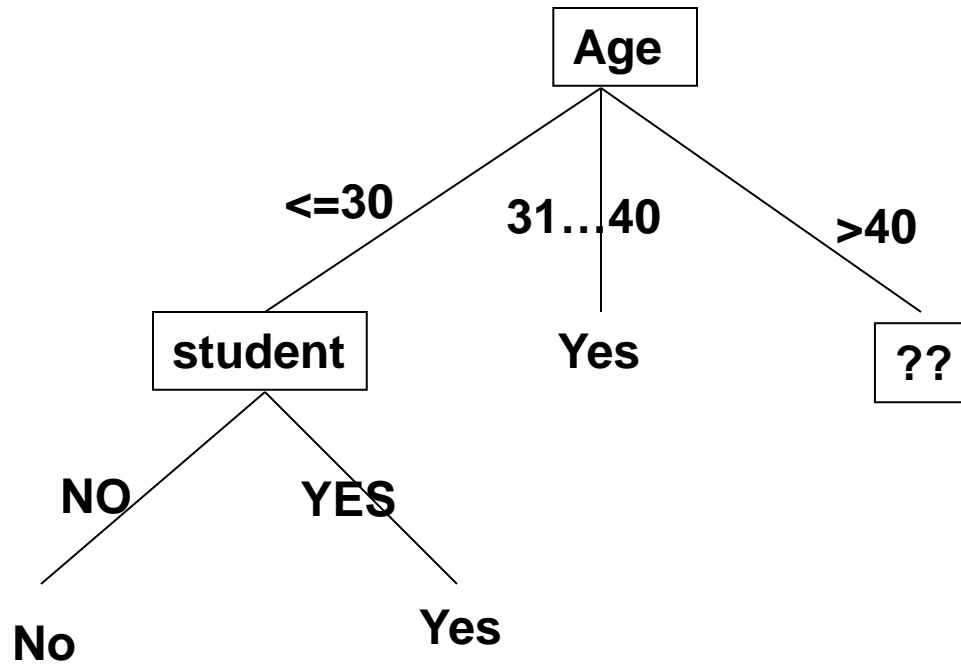

---

	Age	Income	Student	credit rating
Gain	<= 30	0.057	0.97	0.02
Split Info		0.993	0.97	0.97
GainRatio		0.0574	<u>1</u>	0.0206

Student attribute with the maximum gain ratio is selected as the splitting attribute

# Tree induction example

---



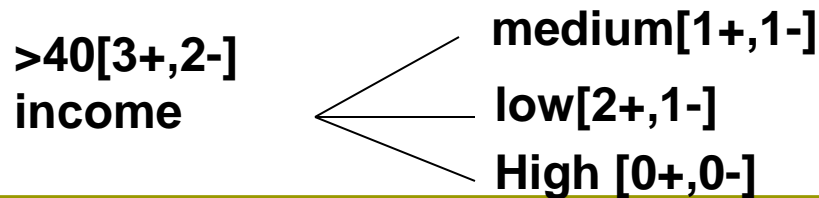


# Tree induction example

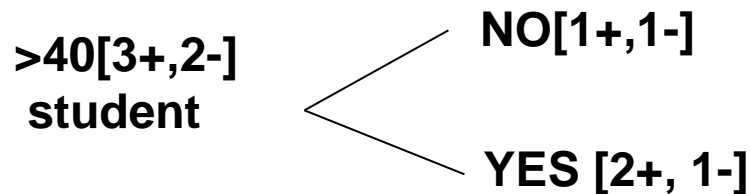
---

age	income	student	credit_rating	correct
>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

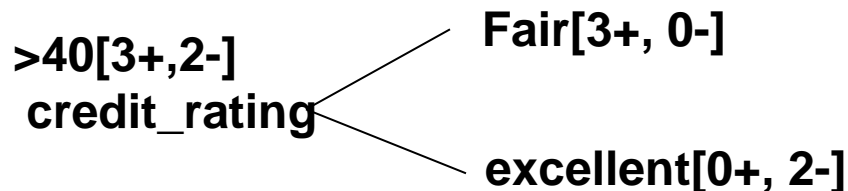
$$\begin{aligned} \text{Info}(>40) &= -3/5(\log_2(3/5)) - 2/5(\log_2(2/5)) \\ &= 0.97 \end{aligned}$$



$$\begin{aligned}
 \text{Gain}(\text{income}) &= 0.97 - 2/5[-1/2(\log_2(1/2)) - 1/2(\log_2(1/2))] \\
 &\quad - 3/5[-2/3(\log_2(2/3)) - 1/3(\log_2(1/3))] \\
 &\quad - 0/5[-0/0(\log_2(0/0)) - 0/0(\log_2(0/0))] \\
 &= 0.97 - 0.75 = 0.22
 \end{aligned}$$



$$\begin{aligned}
 \text{Gain}(\text{student}) &= 0.97 \\
 &\quad - 2/5[-1/2(\log_2(1/2)) - 1/2(\log_2(1/2))] \\
 &\quad - 3/5[-2/3(\log_2(2/3)) - 1/3(\log_2(1/3))] \\
 &= 0.97 - 0.43 = 0.54
 \end{aligned}$$



$$\begin{aligned}
 \text{Gain}(\text{credit\_rating}) &= 0.97 \\
 &\quad - 3/5[-3/3(\log_2(3/3)) - 0/3(\log_2(0/3))] \\
 &\quad - 2/5[-0/2(\log_2(0/2)) - 2/2(\log_2(2/2))] \\
 &= 0.97 - 0 = \underline{0.97}
 \end{aligned}$$

---

age	income	student	credit_rating	correct
>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            = **0.22**            = **0.54**            = **0.97**

$$SplitInfo_{Income}^{>40}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} >40[3+,2-] \\ \text{income} \end{array} \quad \begin{array}{l} \text{medium}[1+,1-] \\ \text{low}[2+,1-] \\ \text{High}[0+,0-] \end{array}$$


---

$$-\frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) = 0.970$$

$$SplitInfo_{student}^{>40}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} >40[3+,2-] \\ \text{student} \end{array} \quad \begin{array}{l} \text{NO}[1+,1-] \\ \text{YES}[2+,1-] \end{array}$$

$$-\frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) = 0.970$$

$$SplitInfo_{credit\_rating}^{>40}(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \quad \begin{array}{l} >40[3+,2-] \\ \text{credit\_rating} \end{array} \quad \begin{array}{l} \text{Fair}[3+,0-] \\ \text{excellent}[0+,2-] \end{array}$$

$$= -\frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \times \log_2 \left( \frac{2}{5} \right) = 0.970$$

age	income	student	credit_rating	comm
>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 = **0.57**      = **0.54**      = **0.97**  
 split Info = **0.970**      = **0.970**      = **0.970**

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$


---

	Age	Income	Student	credit rating
Gain	> 40	0.057	0.54	0.97
Split Info		0.970	0.97	0.97

$$\text{Gain Ratio}(\text{Income} > 40) = \text{Gain}(\text{income}) / \text{Split Info}(\text{income})$$

$$\text{Gain Ratio}(\text{Income} > 40) = 0.057 / 0.970 = 0.05876$$

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$


---

	Age	Income	Student	credit rating
Gain	> 40	0.057	0.54	0.97
Split Info		0.970	0.97	0.97

$$\text{Gain Ratio}(\text{Student} > 40) = \text{Gain}(\text{Student}) / \text{Split Info}(\text{Student})$$

$$\text{Gain Ratio}(\text{Student} > 40) = 0.54 / 0.97 = 0.5567$$

$$\text{Gain Ratio}(\text{credit\_rating} > 40) = \text{Gain}(\text{credit\_rating}) / \text{Split Info}(\text{credit\_rating})$$

$$\text{Gain Ratio}(\text{credit\_rating}) = 0.97 / 0.97 = 1$$

$$\text{Gain Ratio}(A) = \text{Gain}(A) / \text{Split Info}(A)$$


---

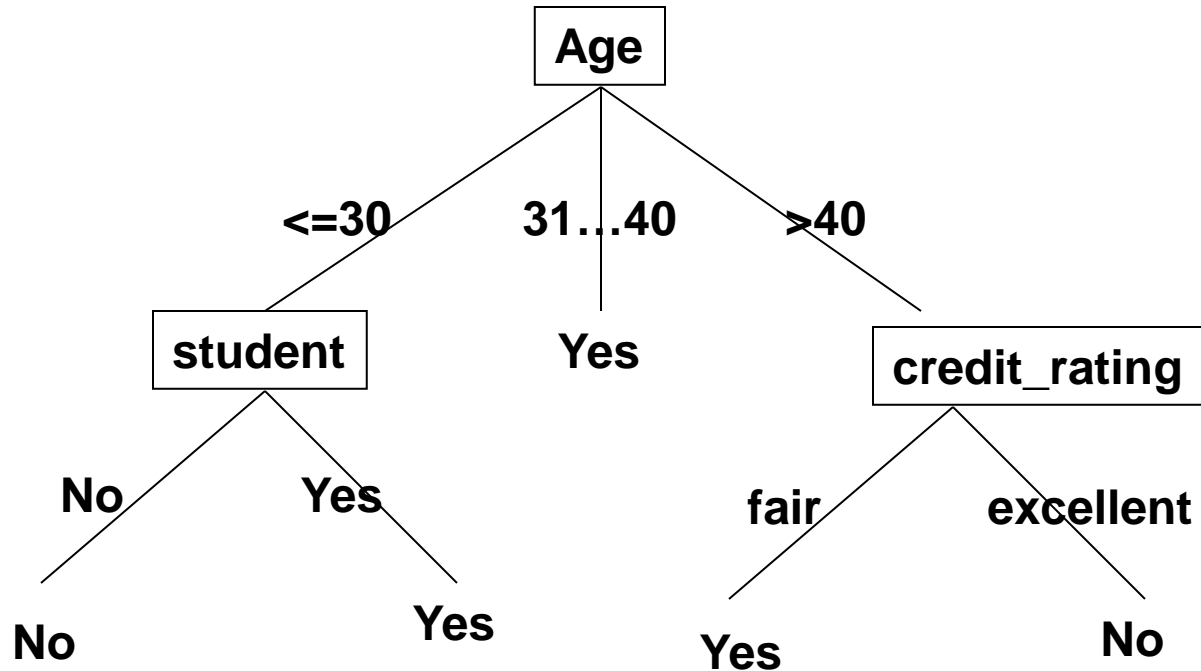
	Age	Income	Student	credit rating
Gain	> 40	0.057	0.54	0.97
Split Info		0.970	0.97	0.97
GainRatio		0.05876	0.5567	1

credit rating attribute with the maximum gain ratio is selected as the splitting attribute



# Tree induction example

---



---

*Thank you*