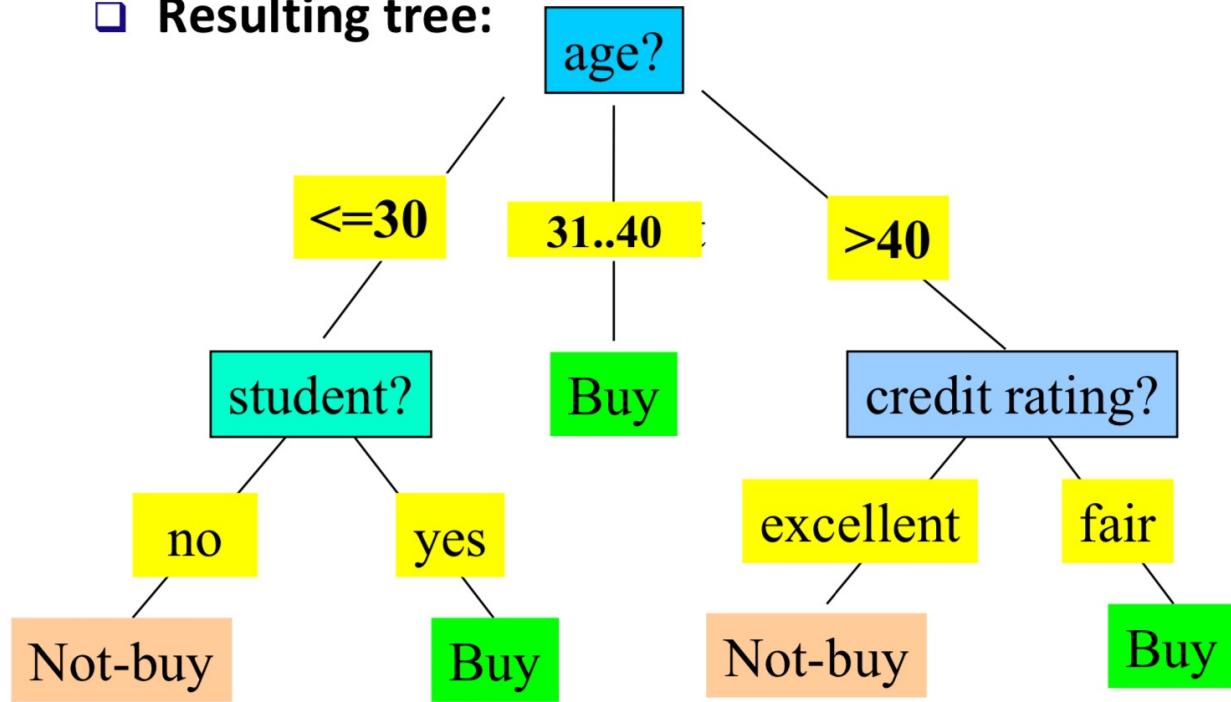


HW 5

□ Resulting tree:



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

ສືບຕໍ່ກໍາໄນກາງຄົງເກີດ

- Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$$

Class ດັວນທີ່

- Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

Feature ດັວນທີ່

- Information gained by branching on attribute A
ກິ່ານຈຳກັດໃນຫຼື້ສູ່ $Gain(A) = Info(D) - Info_A(D)$

1. $Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$

- Class P: buys_computer = "yes"

- Class N: buys_computer = "no"

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

Y N
class
n(ທຸນທັນຂອງ) → 14

$$2. \quad Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- age ແມ່ນອອກເປັນ 3 ກົ່ວມ ຕື່ອ $age \leq 30$, $31 > age \leq 40$, $age > 40$

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

- income ແມ່ນອອກເປັນ 3 ກົ່ວມ ຕື່ອ low, medium, high

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

- student ແມ່ນອອກເປັນ 2 ກົ່ວມ ຕື່ອ yes, no

$$Info_{student}(D) = \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4) = \frac{7}{14} \left(-\frac{6}{7} \log_2 \left(\frac{6}{7} \right) - \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right) + \frac{7}{14} \left(-\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right) = 0.788$$

- credit_rating ແມ່ນອອກເປັນ 2 ນໍາໃຈ fair, excellent

$$Info_{credit_rating}(D) = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3) = \frac{8}{14} \left(-\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right) + \frac{6}{14} \left(-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right) = 0.892$$

$$3. \quad Gain(A) = Info(D) - Info_A(D)$$

$$Gain(age) = 0.940 - 0.694 = 0.246$$

$$Gain(income) = 0.940 - 0.911 = 0.029$$

$$Gain(student) = 0.940 - 0.788 = 0.152$$

$$Gain(credit_rating) = 0.940 - 0.892 = 0.048$$

ລັງນັ້ນເສີ່ງ $Gain(age)$ ເປົ້າຮັກ ເພງກະນຳຂ່າຍກຳສຸດ

4. សម្រាប់លក់មុននៃ feature ទាំងអស់នៃវាក (root node)

4.1 ≤ 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

$$\text{Info}(D) = I(2,3) = 0.971$$

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

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

$$\text{Info}_{\text{credit}}(D) = \frac{2}{5} I(1,2) + \frac{2}{5} I(1,1) = 0.951$$

ការបាន Information Grain

$$\text{Grain}(\text{income}) = 0.971 - 0.4 = 0.571$$

$$\text{Grain}(\text{student}) = 0.971 - 0 = 0.971$$

$$\text{Grain}(\text{Credit_rating}) = 0.971 - 0.951 = 0.02$$

បែនពន្លេ Grain(student) នៃ node លាងទី ≤ 30

4.2 $31\dots 40$

age	income	student	credit_rating	buys_computer
$31\dots 40$	high	no	fair	yes
$31\dots 40$	medium	no	excellent	yes
$31\dots 40$	high	yes	fair	yes
$31\dots 40$	low	yes	excellent	yes

នៅលើ buys_computer ត្រូវបានពិនិត្យថាមព័ត៌មាន yes និង no

\therefore ដើម្បីកញ្ចប់ $31\dots 40$ ត្រូវរាយចក yes នៃ buys_computer តួចតាម

4.3 > 40

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

$$\text{Info}(D) = I(3,2) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.971$$

$$\text{Info}_{\text{income}}(D) = \frac{2}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$$

$$\text{Info}_{\text{student}}(D) = \frac{2}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$$

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

ការបាន Information Grain

$$\text{Grain}(\text{income}) = 0.971 - 0.951 = 0.02$$

$$\text{Grain}(\text{student}) = 0.971 - 0.951 = 0.02$$

$$\text{Grain}(\text{Credit_rating}) = 0.971 - 0 = 0.971$$

បែនពន្លេ Grain(Credit_rating) នៃ node លាងទី > 40

5. សម្រាប់រាយចក Decision Tree តួចតាម

