

# CHƯƠNG 4: PHÂN LỚP DỮ LIỆU

## 4.1. KHÁI NIỆM VỀ PHÂN LỚP DỮ LIỆU

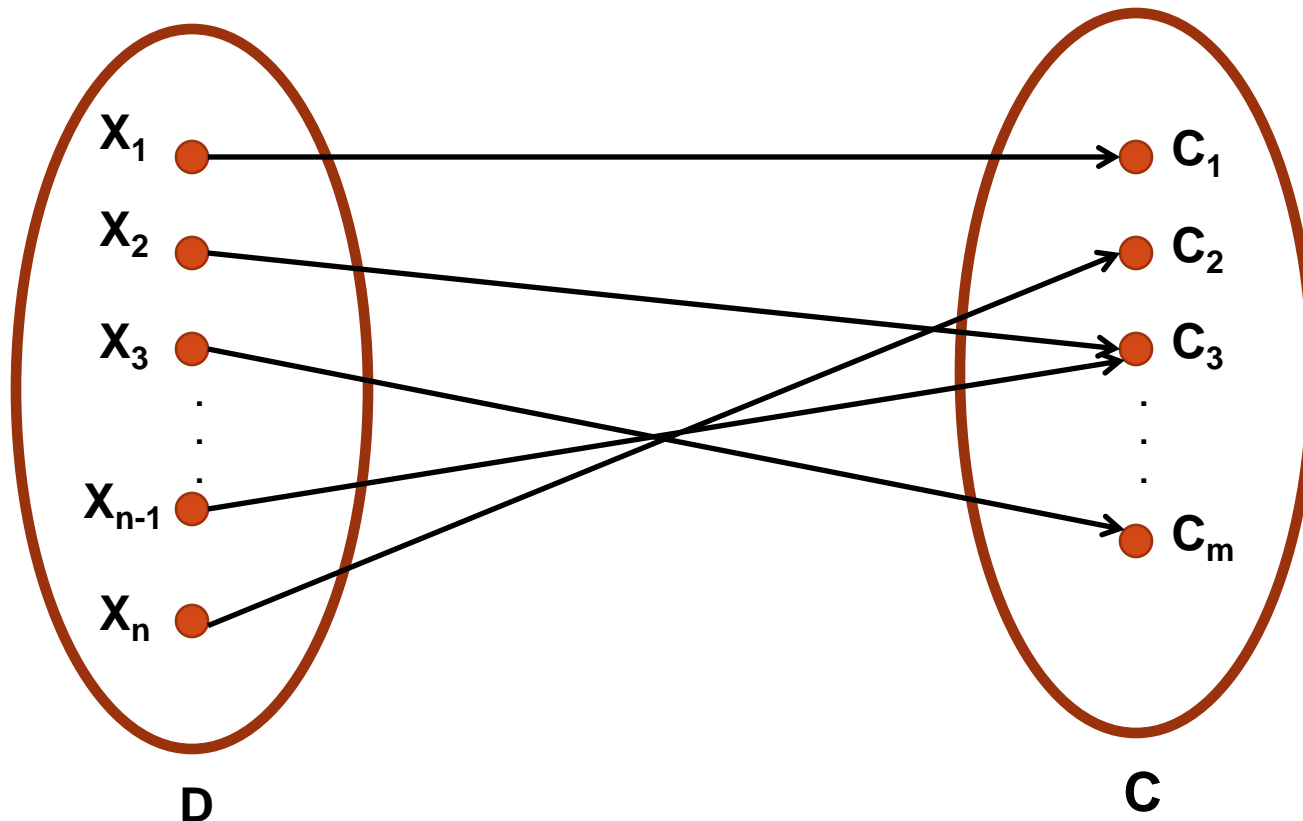
## 4.2. PHÂN LỚP DỰA TRÊN XÁC SUẤT CÓ ĐIỀU KIỆN (*Phân lớp Bayes – Naive Bayesian Classification*)

## 4.3. PHÂN LỚP DỰA TRÊN CÂY QUYẾT ĐỊNH

## 4.1. KHÁI NIỆM VỀ PHÂN LỚP DỮ LIỆU

Cho tập các lớp  $C = \{C_1, C_2, \dots, C_m\}$  và tập dữ liệu  $D = \{X_1, X_2, \dots, X_n\}$

Phân lớp dữ liệu là sự phân chia các đối tượng dữ liệu vào các lớp.  
Về bản chất đây quá trình ánh xạ mỗi đối tượng  $X_j \in D$  tương ứng với một lớp  $C_i \in C$ .

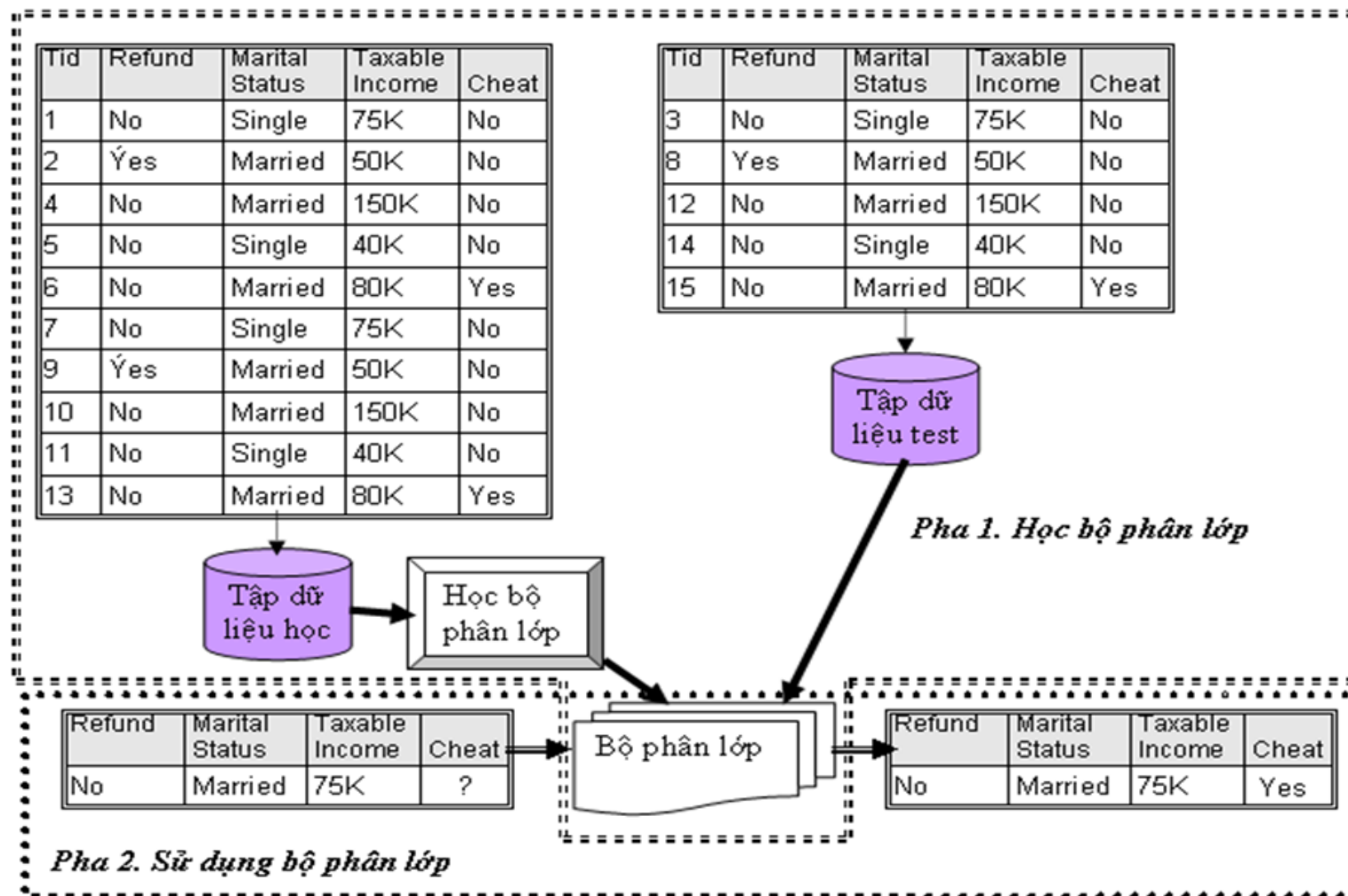


$$f: D \rightarrow C \text{ hay } c = f(X) \text{ (với } X \in D \text{ và } c \in C)$$

- Mỗi ánh xạ được gọi là một mô hình phân lớp (Classification Model).

⇒ **Làm sao để xây dựng mô hình phân lớp?**

**Thông qua quá trình huấn luyện dựa trên tập dữ liệu học (học có giám sát – supervised learning)**



## ***Xây dựng mô hình***

**B1:** Chọn một tập ví dụ mẫu (gồm các đối tượng đã được phân lớp):

$$\mathbf{D}_{\text{exam}} = \mathbf{D}_1 \cup \mathbf{D}_2 \cup \dots \cup \mathbf{D}_m \text{ trong đó } \mathbf{D}_i = \{X | (X \in \mathbf{D}) \wedge (X \rightarrow \mathbf{C}_i)\} \quad i=1, \dots, m$$

**B2:** Tách  $\mathbf{D}_{\text{exam}}$  thành 02 tập:

❖ Tập dữ liệu học  $\mathbf{D}_{\text{train}}$

❖ Tập dữ liệu kiểm tra  $\mathbf{D}_{\text{test}}$

Hiển nhiên  $\mathbf{D}_{\text{exam}} = \mathbf{D}_{\text{train}} \cup \mathbf{D}_{\text{test}}$  và thường thì người ta tách sao cho:

$$|\mathbf{D}_{\text{train}}| = \frac{2}{3} |\mathbf{D}_{\text{exam}}| \quad |\mathbf{D}_{\text{test}}| = \frac{1}{3} |\mathbf{D}_{\text{exam}}|$$

**B3:** Dùng  $\mathbf{D}_{\text{train}}$  để xây dựng mô hình (xác định tham số). Có nhiều loại mô hình phân lớp như: Bayes, cây quyết định, luật phân lớp,...

**B4:** Dùng  $\mathbf{D}_{\text{test}}$  để kiểm tra, đánh giá mô hình xây dựng được.

**B5:** Chọn mô hình có chất lượng nhất.

## ***Sử dụng mô hình***

Cho  $X \in \mathbf{D}$  (là tập dữ liệu chưa phân lớp)  $\Rightarrow$  Xác định lớp của  $X$

## 4.2. PHÂN LỚP DỰA TRÊN XÁC SUẤT CÓ ĐIỀU KIỆN (*Naive Bayes Classifier*)

### 4.2.1. Xác suất có điều kiện và công thức Bayes

- Gọi  $X$  là một bộ dữ liệu (data tuple). Theo ngôn ngữ xác suất,  $X$  được xem là một biến cố (evidence).
- Gọi  $H$  là một giả thuyết (hypothesis): bộ  $X$  thuộc về lớp  $C_i$ .

⇒ **Cần xác định  $P(H|X)$** : xác suất xảy ra  $H$  khi đã xuất hiện  $X$  (hay xác suất để  $X$  thuộc về lớp  $C_i$  nếu như đã biết các thuộc tính của  $X$ ).

Nhãn phân lớp

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

$X$  được xác định  
thông qua tập  
giá trị của các  
thuộc tính

- **$P(H|X)$  là xác suất có điều kiện của H đối với X (xác suất xảy ra H khi biết X xảy ra).**

Ví dụ:  $X = (\text{age}=35 \text{ years old}, \text{income}=\$40,000),$   
 $H = (\text{buy\_computer}=\text{Yes})$

$$P(H|X) = P(\text{buy\_computer}=\text{yes} \mid \text{age}=35 \text{ years old}, \text{income}=\$40,000)$$

$\Rightarrow$  Xác suất để một người 35 tuổi có thu nhập \$40,000 mua máy tính

- **$P(X|H)$  là xác suất có điều kiện của X đối với H (xác suất xảy ra X khi biết H xảy ra).**

Ví dụ:

$$P(X|H) = P(\text{age}=35 \text{ years old}, \text{income}=\$40,000 \mid \text{buy\_computer}=\text{yes})$$

$\Rightarrow$  Xác suất để một người mua máy tính có độ tuổi là 35 và thu nhập là \$40,000.

- **$P(X)$  là xác suất tiên nghiệm của X.**

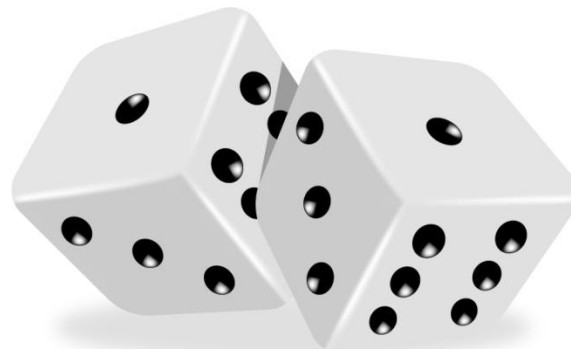
Ví dụ:  $P(X) = P(\text{age}=35 \text{ years old}, \text{income}=\$40,000) \Rightarrow$  Xác suất để tìm thấy trong tập dữ liệu đang xét một người có độ tuổi là 35 và thu nhập là \$40,000.

- **$P(H)$  là xác suất tiên nghiệm của H.**

Ví dụ:  $P(H) = P(\text{buy\_computer}=\text{Yes}) \Rightarrow$  Xác suất mua máy tính của khách hàng nói chung (không quan tâm đến độ tuổi hay thu nhập)



**Thomas Bayes**  
(1702 – 1761)



**Công thức Bayes:**

$$P(H | X) = \frac{P(X | H)P(H)}{P(X)}$$

*T. Bayes.*

## 4.2.2. Phân lớp dữ liệu dựa trên xác suất có điều kiện (phân lớp Bayes)

Bộ phân lớp Bayes hoạt động như sau:

1. Cho  $D$  là tập dữ liệu học gồm các bộ và nhãn lớp tương ứng (đã được phân lớp). Mỗi bộ được biểu diễn bởi một vector  $n$  chiều  $X = (x_1, x_2, \dots, x_n)$  trong đó  $x_i$  là giá trị tương ứng với thuộc tính  $A_i$  ( $i = 1, 2, \dots, n$ ). Tập  $D_i = \{X | (X \in D) \wedge (X \rightarrow C_i)\}$  là tập các bộ trong  $D$  thuộc về lớp  $C_i$ .
2. Giả sử có  $m$  lớp  $C_1, C_2, \dots, C_m$ . Bộ  $X$  được dự đoán là thuộc về lớp  $C_i$  khi và chỉ khi:  $P(C_i|X) > P(C_j|X)$  với mọi  $j \neq i$  và  $1 \leq j \leq m$  ( $X$  thuộc về lớp mà xác suất có điều kiện khi biết  $X$  là lớn nhất)  $\Rightarrow$  **Đi tìm lớp  $C_i$  trong số  $m$  lớp sao cho  $P(C_i|X)$  là lớn nhất.**
3.  $P(X)$  là giống nhau với tất cả các lớp nên theo công thức Bayes thì  $P(C_i|X)$  lớn nhất tương ứng với tích  $P(X|C_i)P(C_i)$  lớn nhất  $\Rightarrow$  **Đi tìm  $C_i$  sao cho tích  $P(X|C_i)P(C_i)$  là lớn nhất ( $i = 1, 2, \dots, m$ ).**

4. Ta có thể tính:

$$P(C_i) = \frac{|D_i|}{|D|}$$

và nếu coi  $n$  thuộc tính của  $X$  là độc lập thì:

$$P(X | C_i) = \prod_{k=1}^n P(x_k | C_i) = P(x_1 | C_i)P(x_2 | C_i) \dots P(x_n | C_i)$$

$$P(x_k | C_i) = \frac{|\{X' | (X'(A_k) = x_k) \wedge (X' \in D_i)\}|}{|D_i|}$$



### Chú ý:

- ❖ Nếu không tính được  $P(C_i)$  thì có thể coi  $P(C_1) = P(C_2) = \dots = P(C_m)$  và bài toán quy về  **tìm lớp  $C_i$  trong số  $m$  lớp sao cho  $P(X|C_i)$  có giá trị lớn nhất.**
- ❖ Nếu tồn tại  $P(x_k|C_i) = 0$  thì có thể áp dụng hiệu chỉnh Laplace và công thức tính của  $P(x_k|C_i)$  được hiệu chỉnh như sau:

$$P(x_k | C_i) = \frac{\left| \{X' \mid (X'(A_k) = x_k) \wedge (X' \in D_i)\} \right| + 1}{|D_i| + q}$$

***q: số giá trị khác nhau của  $A_k$***

## Ví dụ:

Cho tập dữ liệu học gồm các bộ dữ liệu đã được phân lớp như sau:

<i>RID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>Class: buys_computer</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

**Áp dụng phân lớp Bayes hãy dự đoán bộ dữ liệu**

$X = (age = youth, income = medium, student = yes, credit\_rating = fair)$

**thuộc lớp nào?**

Có 02 lớp dữ liệu tương ứng với **buys\_computer = yes** và **buys\_computer = no**

$$P(\text{buys\_computer} = \text{yes}) = 9/14 = 0.643$$

$$P(\text{buys\_computer} = \text{no}) = 5/14 = 0.357$$

$$P(\text{age} = \text{youth} \mid \text{buys\_computer} = \text{yes}) = 2/9 = 0.222$$

$$P(\text{age} = \text{youth} \mid \text{buys\_computer} = \text{no}) = 3/5 = 0.600$$

$$P(\text{income} = \text{medium} \mid \text{buys\_computer} = \text{yes}) = 4/9 = 0.444$$

$$P(\text{income} = \text{medium} \mid \text{buys\_computer} = \text{no}) = 2/5 = 0.400$$

$$P(\text{student} = \text{yes} \mid \text{buys\_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{student} = \text{yes} \mid \text{buys\_computer} = \text{no}) = 1/5 = 0.200$$

$$P(\text{credit\_rating} = \text{fair} \mid \text{buys\_computer} = \text{yes}) = 6/9 = 0.667$$

$$P(\text{credit\_rating} = \text{fair} \mid \text{buys\_computer} = \text{no}) = 2/5 = 0.400$$

**Suy ra:**  $P(X \mid \text{buys\_computer} = \text{yes}) = P(\text{age} = \text{youth} \mid \text{buys\_computer} = \text{yes}) \times$   
 $P(\text{income} = \text{medium} \mid \text{buys\_computer} = \text{yes}) \times$   
 $P(\text{student} = \text{yes} \mid \text{buys\_computer} = \text{yes}) \times$   
 $P(\text{credit\_rating} = \text{fair} \mid \text{buys\_computer} = \text{yes})$   
 $= 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044.$

**Tương tự:**  $P(X \mid \text{buys\_computer} = \text{no}) = 0.600 \times 0.400 \times 0.200 \times 0.400 = 0.019$

$$P(X \mid \text{buys\_computer} = \text{yes})P(\text{buys\_computer} = \text{yes}) = 0.044 \times 0.643 = 0.028$$

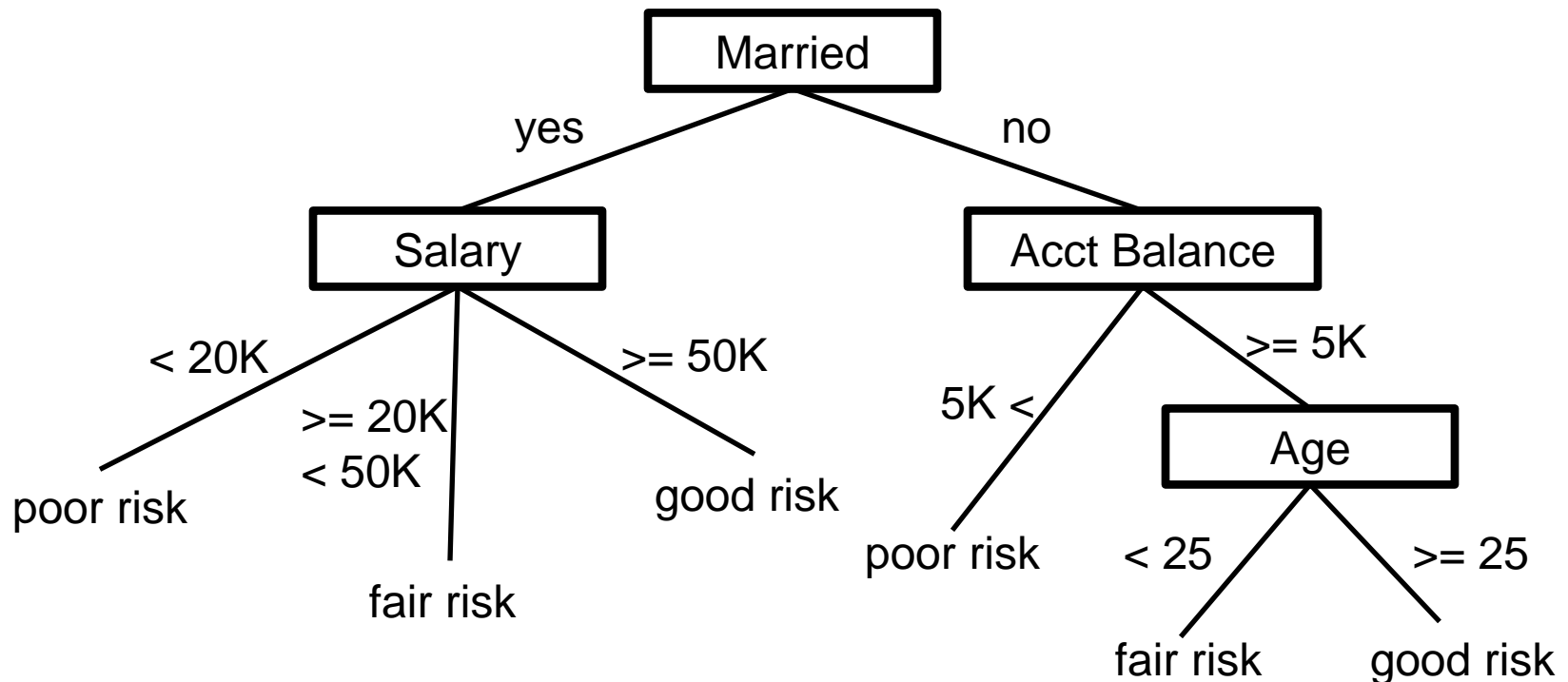
$$P(X \mid \text{buys\_computer} = \text{no})P(\text{buys\_computer} = \text{no}) = 0.019 \times 0.357 = 0.007$$

**$\Rightarrow X$  thuộc lớp dữ liệu tương ứng với **buys\_computer = yes****

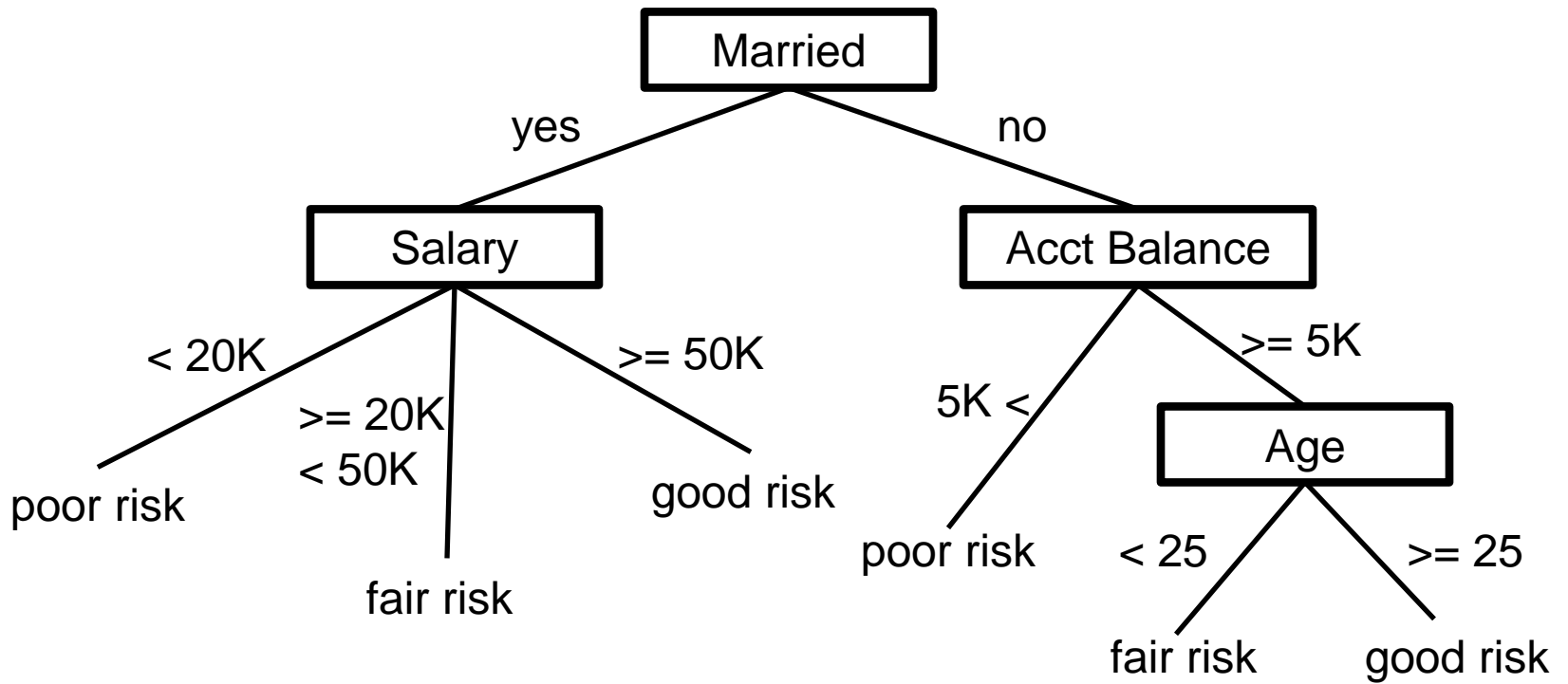
## 4.3. PHÂN LỚP DỰA TRÊN CÂY QUYẾT ĐỊNH

### 4.3.1. Mô hình phân lớp cây quyết định

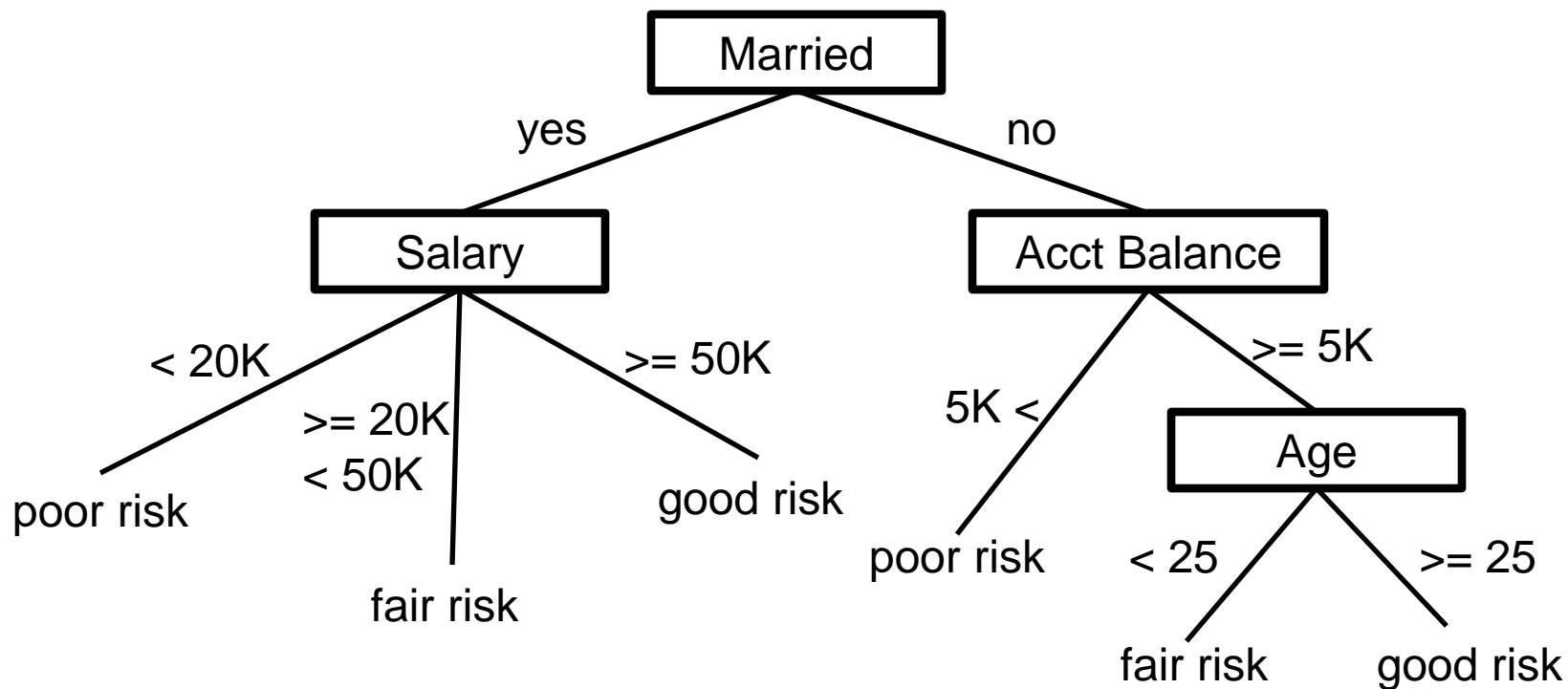
- Cây quyết định (decision tree) là một mô hình phân lớp điển hình.
- Cây quyết định bao gồm:
  - ❖ **Các nút trong**: biểu diễn cho một thuộc tính được kiểm thử (test).
  - ❖ **Các nút lá**: nhãn/mô tả của một lớp (class label).
  - ❖ **Nhánh**: xuất phát từ một nút trong, phản ánh kết quả của một phép thử trên thuộc tính tương ứng.



- Có thể dễ dàng chuyển đổi từ mô hình **cây quyết định** sang mô hình **luật phân lớp** bằng cách: đi từ nút gốc cho tới nút lá, mỗi đường đi tương ứng với một luật phân lớp.



1. **If** (Married = yes) **And** (Salary > 20K) **Then** Class = poor risk
2. **If** (Married = yes) **And** (50K > Salary >= 20K) **Then** Class = fair risk
3. **If** (Married = yes) **And** (Salary >= 50K) **Then** Class = good risk
4. **If** (Married = no) **And** (Acct Balance < 5K) **Then** Class = poor risk
5. **If** (Married = no) **And** (Acct Balance >= 5K) **And** (Age < 25) **Then** Class = fair risk
6. **If** (Married = no) **And** (Acct Balance >= 5K) **And** (Age >= 25) **Then** Class = good risk



Name	Age	Married	Salary	Acct Balance	Class
Alice	19	yes	30K	6K	?
Pike	28	no	60K	7K	?
Tom	35	yes	10K	10K	?
Peter	24	no	20K	8K	?
Lucas	40	no	20K	3K	?



Name	Age	Married	Salary	Acct Balance	Class
Alice	19	yes	30K	6K	fair risk
Pike	28	no	60K	7K	good risk
Tom	35	yes	10K	10K	poor risk
Peter	24	no	20K	8K	fair risk
Lucas	40	no	20K	3K	poor risk

### 4.3.2. Các độ đo sử dụng trong phân lớp

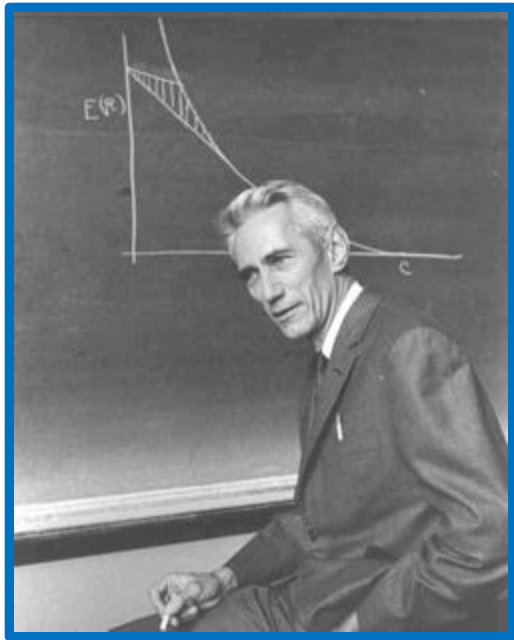
#### A. Entropy của tập dữ liệu

*Là lượng thông tin cần để phân loại một phần tử trong tập dữ liệu  $D$ .  
Ký hiệu là  $Infor(D)$ .*

Gọi:

$p_i$ : xác suất để một phần tử bất kỳ trong  $D$  thuộc về lớp  $C_i$  ( $i=1, 2, \dots, m$ ).

$D_i$ : Tập các phần tử trong  $D$  thuộc về lớp  $C_i$ .



**Claude Elwood Shannon**  
(1916 – 2001)

$$p_i = \frac{|D_i|}{|D|}$$

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

## B. Entropy của dữ liệu ứng với một thuộc tính

*Là lượng thông tin cần để phân loại một phần tử trong tập dữ liệu  $D$  dựa trên thuộc tính  $A$ . Ký hiệu là  $Infor_A(D)$ .*

- ❖ Thuộc tính  $A$  dùng để phân tách  $D$  thành  $v$  phân hoạch (tập con) là  $D_1, D_2, \dots, D_v$ .
- ❖ Mỗi phân hoạch  $D_j$  có  $|D_j|$  phần tử.
- ❖ Lượng thông tin này sẽ cho biết mức độ trùng lặp giữa các phân hoạch, nghĩa là một phân hoạch chứa các phần tử từ một hay nhiều lớp khác nhau.

⇒ **Mong đợi:**  $Infor_A(D)$  càng nhỏ càng tốt.

$$Infor_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} Infor(D_j)$$





### C. Độ lợi thông tin (Information Gain)

- Mục tiêu: Tối thiểu hóa lượng thông tin cần thiết để phân lớp các các mẫu dữ liệu (tối thiểu hóa số lượng các điều kiện kiểm tra cần thiết để phân lớp một bản ghi mới).*

Độ lợi thông tin ứng với thuộc tính A (ký hiệu  $Gain(A)$ ) chính là độ sai biệt giữa Entropy ban đầu của tập dữ liệu (trước phân hoạch) và Entropy của dữ liệu ứng với thuộc tính A (sau khi phân hoạch bởi A).

$$Gain(A) = Infor(D) - Infor_A(D)$$



### 4.3.3. Giải thuật ID3 xây dựng cây quyết định

**Input:** Tập dữ liệu học Records gồm m đối tượng (bản ghi)  $R_1, R_2, \dots, R_m$ .

Tập thuộc tính Attributes gồm m thuộc tính  $A_1, A_2, \dots, A_n$ .

**Output:** Mô hình cây quyết định.

**procedure** Build\_tree(Records, Attributes)

**begin**

    Tạo nút N;

**if** (tất cả các bản ghi thuộc về một lớp  $C_i$  nào đó) **then**

**begin**

            N.Label =  $C_i$ ;

            return N;

**end;**

**if** (Attributes =  $\emptyset$ ) **then**

**begin**

            Tìm lớp  $C_j$  mà phần lớn các bản ghi  $r \in$  Records thuộc về lớp đó.

            N.Label =  $C_j$ ;

            return N;

**end;**

    Chọn  $A_i \in$  Attribute sao cho  $\text{Gain}(A_i) \rightarrow \max$ ;

    N.Label =  $A_i$ ;

**for each** giá trị  $v_i$  đã biết của  $A_i$  **do**

**begin**

            Thêm một nhánh mới vào nút N ứng với  $A_i = v_j$ ;

$S_j =$  Tập con của Records có  $A_i = v_j$ ;

**if** ( $S_j = \emptyset$ ) **then**

                Thêm một nút lá L với nhãn là lớp mà phần lớn các bản ghi  $r \in$  Records thuộc về lớp đó;

                Return L;

**else**

                Thêm vào nút được trả về bởi Build\_Tree( $S_j$ , Attribute  $\setminus \{A_i\}$ );

**end ;**

**end;**

## Phương pháp lựa chọn thuộc tính

Dùng heuristic để chọn tiêu chí rẽ nhánh tại một nút: Phân hoạch tập dữ liệu học  $D$  thành các phân hoạch con với các nhãn phù hợp:

- Xếp hạng mỗi thuộc tính.
- Thuộc tính được chọn để rẽ nhánh là thuộc tính có trị số điểm (score) là lớn nhất.
- Độ đo để chọn thuộc tính phân tách (splitting attribute) là Information Gain (được xây dựng dựa trên lý thuyết thông tin của Claude Elwood Shannon).

Cụ thể: **Thuộc tính có giá trị Information Gain lớn nhất sẽ được chọn làm thuộc tính phân nhánh cho nút  $N$ .**

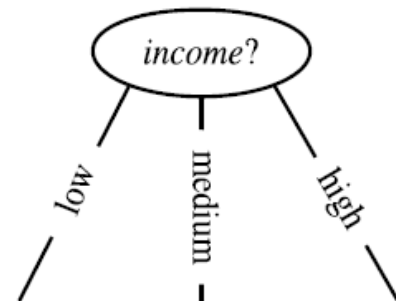
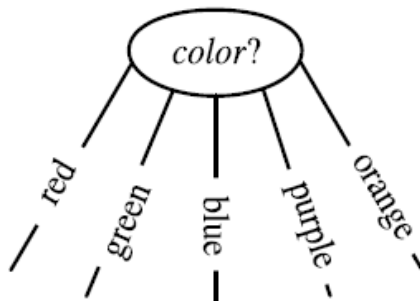
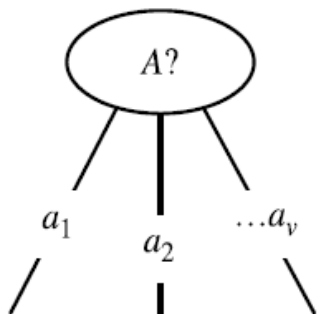
- ❖ Nút  $N$  là nút hiện tại cần phân hoạch các phần tử trong  $D$ .
- ❖ Thuộc tính phân hoạch đảm bảo sự trùng lặp ngẫu nhiên ít nhất giữa các phân hoạch tạo được.

**⇒ Giúp tối thiểu số phép thử (test) cần để phân loại một phần tử.**

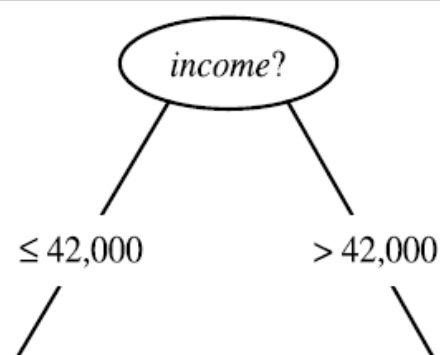
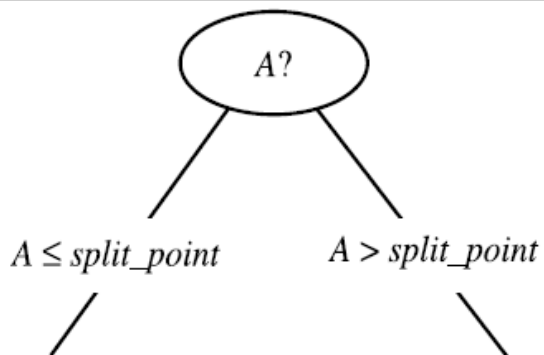
# Partitioning Scenarios

## Examples

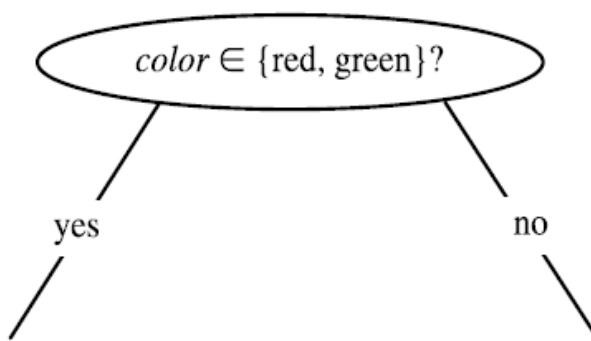
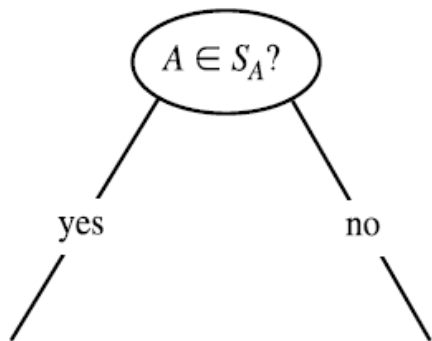
a)



b)



c)



### Ví dụ 1: Cho tập dữ liệu học:

RID	age	income	student	credit_rating	Class: buys_computer
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

**Tính toán tương tự:**

$$Gain(income) = 0.029 \text{ bits}$$

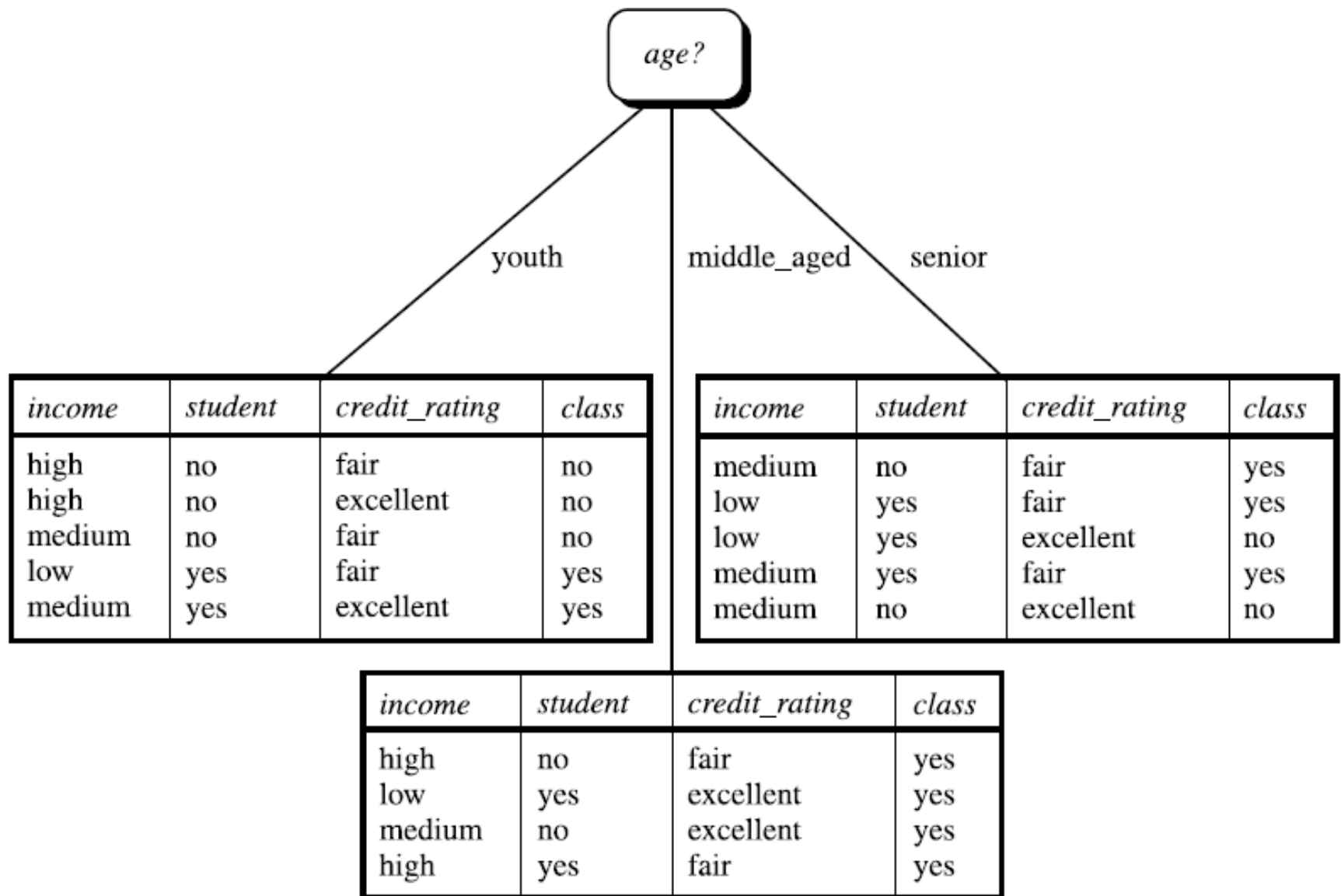
$$Gain(student) = 0.151 \text{ bits}$$

$$Gain(credit\_rating) = 0.048 \text{ bits}$$

**⇒ Chọn age là thuộc tính phân tách**

$$\begin{aligned} Info_{age}(D) &= \frac{5}{14} \times \left( -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) \\ &\quad + \frac{4}{14} \times \left( -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right) \\ &\quad + \frac{5}{14} \times \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\ &= 0.694 \text{ bits.} \end{aligned}$$

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



**Q & A**