



## Data Mining Nguyen Manh Tuan

khai thác dữ liệu và ứng dụng (Trường Đại học Công nghiệp Thành phố Hồ Chí Minh)



Scan to open on Studocu

Người gửi : Nguyễn Mạnh Tuấn  
MSSV: CH0601093  
Ngày gửi : 30/06/2007

---

## DATA MINING

### Bài tập 1:

- Cho tập các hoá đơn  $O = \{o1, o2, o3, o4, o5\}$ , mỗi hóa đơn chứa các mặt hàng như sau:  
 $o1 = \{i1, i3, i4\}$  ;  $o2 = \{i1, i3, i4\}$  ;  $o3 = \{i3, i5\}$  ;  $o4 = \{i4, i5\}$  ;  $o5 = \{i2, i3, i5\}$   
Cho ngưỡng phổ biến tối thiểu  $\text{minsupp} = 0,4$  hãy:
  - Tìm các tập phổ biến tối đại theo ngưỡng  $\text{minsupp} = 0,4$
  - Tìm tất cả các luật kết hợp có độ phổ biến tối thiểu là 0,4 và độ tin cậy tối thiểu là 0,8
- Sử dụng cây định danh để tìm các luật phân lớp từ bảng quyết định sau đây:

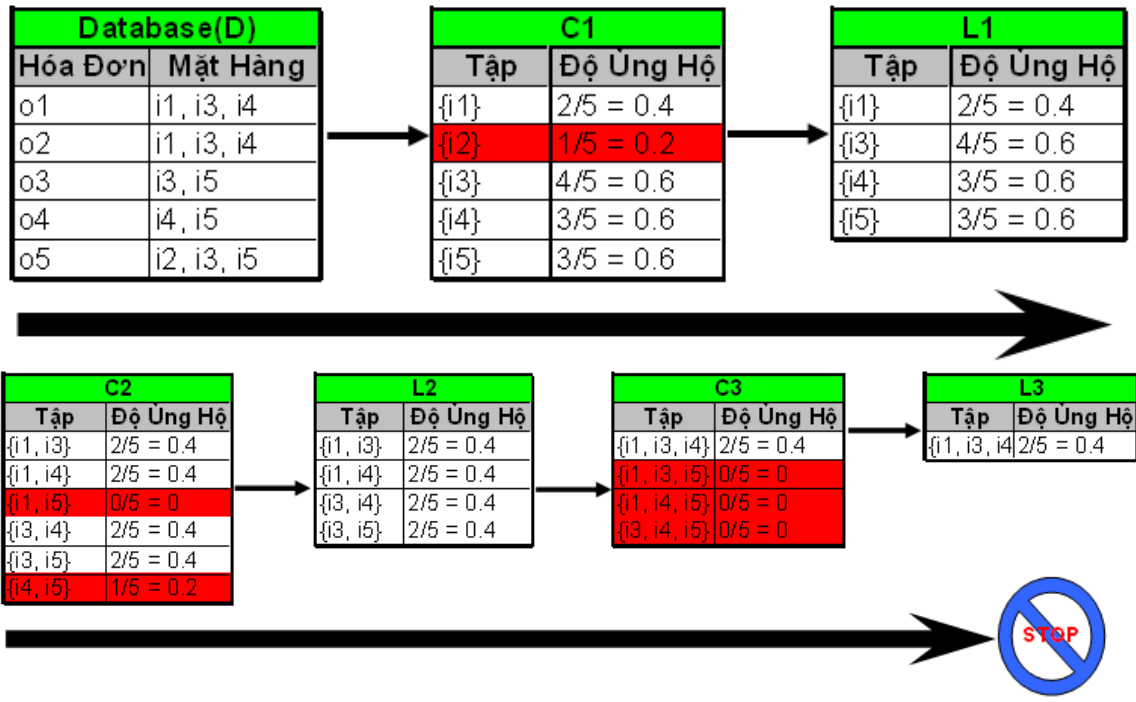
#	Trời	Áp Suất	Gió	Kết quả
1	Trong	Cao	Bắc	Không mưa
2	Mây	Cao	Nam	Mưa
3	Mây	Trung bình	Bắc	Mưa
4	Trong	Thấp	Bắc	Không mưa
5	Mây	Thấp	Bắc	Mưa
6	Mây	Cao	Bắc	Mưa
7	Mây	Thấp	Nam	Không mưa
8	Trong	Cao	Nam	Không mưa

Bạn có suy nghĩ gì về việc dùng luật kết hợp để làm luật phân lớp.  
Bảng dữ liệu lúc đó sẽ có các cột <Trời, Trong>, <Trời, mây>, <Áp suất, Cao>  
<Áp suất, trung bình>, <Áp suất, Thấp>

**Giải bài tập 1:**

**1. Luật kết hợp**

**a. Tìm các tập phổ biến tối đại theo ngưỡng minsupp=0.4**



Vậy các tập phổ biến thu được là :

$L1 = \{i1\}, \{i3\}, \{i4\}, \{i5\}$

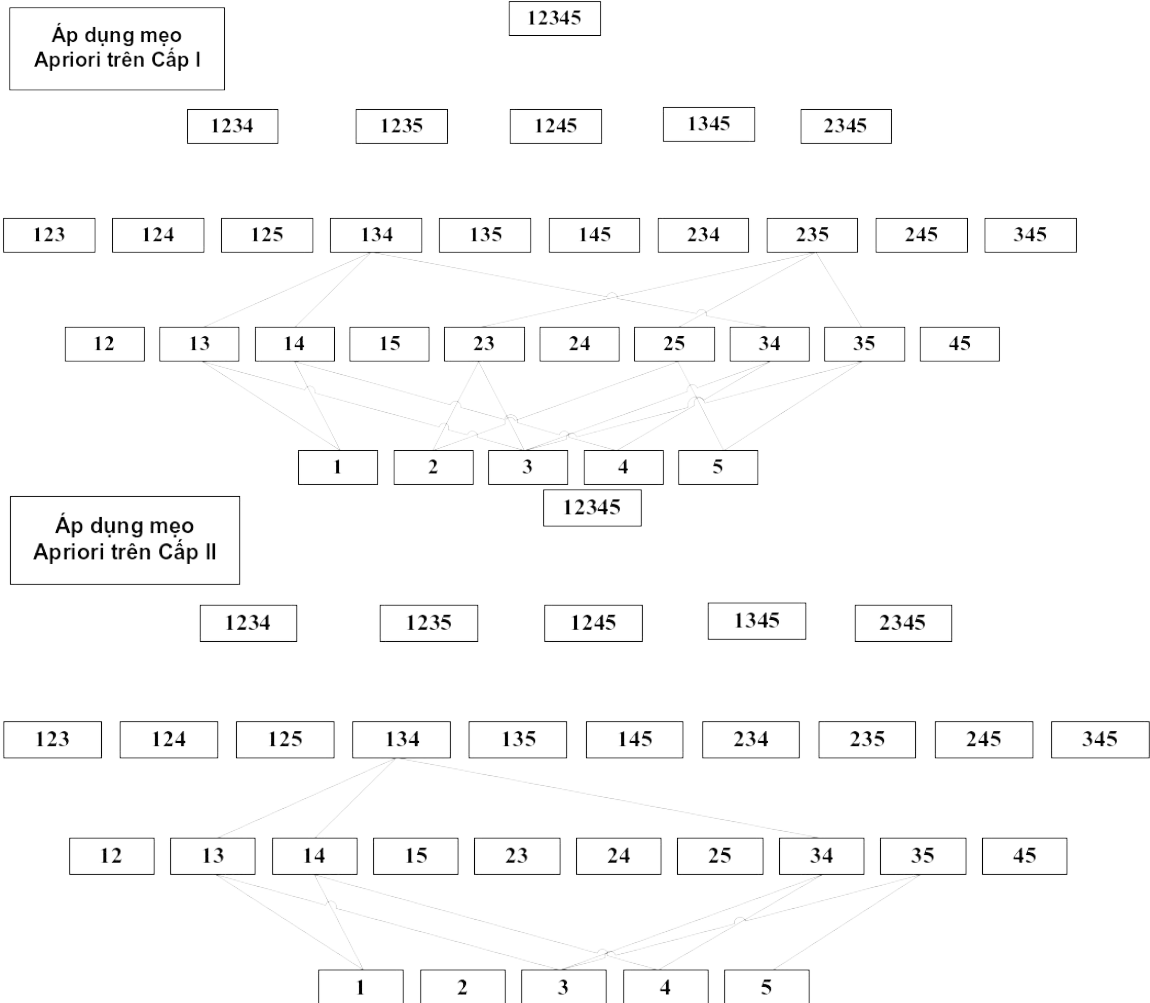
$L2 = \{\{i1, i3\}, \{i1, i4\}, \{i3, i4\}, \{i3, i5\}\}$

$L3 = \{\{i1, i3, i4\}\}$

Vậy tập phổ biến tối đại là:  $\{\{i3, i5\}, \{i1, i3, i4\}\}$

**Chú ý:** chúng ta có thể áp dụng heuristic tại bước  $L2 \rightarrow C3$ , vì tạo ra  $C3$  phải có 3 phần tử nên ta chỉ cần quan tâm đến 2 tập 3 phần tử xuất tập  $D$ :  $\{i1, i3, i4\}$  và  $\{i2, i3, i5\}$ , cho nên tại bước này ta chỉ cần chọn tập  $\{i1, i3, i4\}$  để tìm tập phổ biến.

Hay chúng ta dùng mẹo Apriori:



**b. Tìm tất cả các luật kết hợp có độ phổ biến tối thiểu là 0,4 và độ tin cậy tối thiểu là 0,8**

Tập Phổ Biến (S)	Luật	Độ Tin Cậy
{i1, i3}	I1 => i3	$\text{Supp}(S)/\text{Supp}(i1) = (2/5)/(2/5) = 1$
	I3 => i1	$\text{Supp}(S)/\text{Supp}(i3) = (2/5)/(4/5) = 0.5$
{i1, i4}	I1 => i4	$\text{Supp}(S)/\text{Supp}(i1) = (2/5)/(2/5) = 1$
	I4 => i1	$\text{Supp}(S)/\text{Supp}(i4) = (2/5)/(3/5) = 0.67$
{i3, i4}	I3 => i4	$\text{Supp}(S)/\text{Supp}(i3) = (2/5)/(4/5) = 0.5$
	I4 => i3	$\text{Supp}(S)/\text{Supp}(i4) = (2/5)/(3/5) = 0.67$
{i3, i5}	I3 => i5	$\text{Supp}(S)/\text{Supp}(i3) = (2/5)/(4/5) = 0.5$
	I5 => i3	$\text{Supp}(S)/\text{Supp}(i5) = (2/5)/(3/5) = 0.67$
{i1, i3, i4}	I1, i3 => i4	$\text{Supp}(S)/\text{Supp}(i1, i3) = (2/5)/(2/5) = 1$
	I1, i4 => i3	$\text{Supp}(S)/\text{Supp}(i1, i4) = (2/5)/(2/5) = 1$
	I3, i4 => i1	$\text{Supp}(S)/\text{Supp}(i3, i4) = (2/5)/(2/5) = 1$
	I1 => i3, i4	$\text{Supp}(S)/\text{Supp}(i1) = (2/5)/(2/5) = 1$

I3 => i1, i4	$\text{Supp}(S)/\text{Supp}(i3) = (2/5)/(4/5) = 0.5$
I4 => i1, i3	$\text{Supp}(S)/\text{Supp}(i4) = (2/5)/(3/5) = 0.67$

Tất cả các luật kết hợp có độ phổ biến tối thiểu là 0,4 và độ tin cậy tối thiểu là 0,8:

- Luật 1 : i1 => i3
- Luật 2 : i1 => i4
- Luật 3: i1,i3 => i4
- Luật 4: i1,i4 => i3
- Luật 5: i3,i4 => i1
- Luật 6: i1 => i3, i4

## 2. Luật Phân Lớp

#	Trời	Áp Suất	Gió	Kết quả
1	Trong	Cao	Bắc	Không mưa
2	Mây	Cao	Nam	Mưa
3	Mây	Trung bình	Bắc	Mưa
4	Trong	Thấp	Bắc	Không mưa
5	Mây	Thấp	Bắc	Mưa
6	Mây	Cao	Bắc	Mưa
7	Mây	Thấp	Nam	Không mưa
8	Trong	Cao	Nam	Không mưa

Dùng thuật toán ID3 để phân hoạch:

**Ký hiệu:** P: Mưa; N: Không mưa

$$I(P, N) = -P/(P+N) \cdot \log_2(P/(P+N)) - N/(P+N) \cdot \log_2(N/(P+N))$$

Ta có: P = 4 và N = 4

$$\Rightarrow I(4, 4) = 1$$

- Tính entropy cho thuộc tính [Trời]:

Trời	pi	ni	I(pi, ni)
Trong	0	3	0
Mây	4	1	0.72
<b>E (Trời) = 3/8*0 + 5/8*0.72 = 0.45</b>			
<b>Gain(Trời) = I(4,4) - E (Trời) = 0.55</b>			

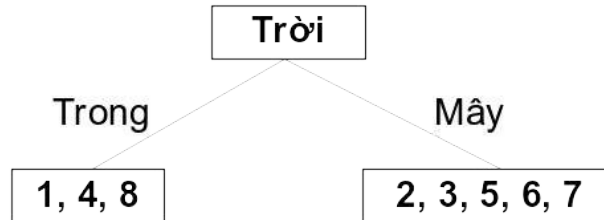
- Tính entropy cho thuộc tính [Áp Suất]:

Áp Suất	pi	ni	I(pi, ni)
Cao	2	2	1.00
Trung Bình	1	0	0.00
Thấp	1	2	0.92
<b>E (Áp Suất) = 4/8*1 + 1/8*0 + 3/8*0.92 = 0.84</b>			
<b>Gain(Áp Suất) = I(4,4) - E (Áp Suất) = 0.155</b>			

- Tính entropy cho thuộc tính [Gió]:

Gió	pi	ni	I(pi, ni)
Bắc	3	2	0.97
Nam	1	2	0.92
<b>E (Gió) = 4/8*0.97 + 3/8*0.92 = 0.95</b>			
<b>Gain(Trời) = I(4,4) - E (Gió) = 0.05</b>			

Ta nhận thấy Gain của thuộc tính [Trời] là lớn nhất, nên ta dùng thuộc tính [Trời] để phân lớp:



- Phân lớp nhánh [Trời - Trong]:

- Bảng dữ liệu của nhánh:

#	Áp Suất	Gió	Kết quả
1	Cao	Bắc	Không mưa
4	Thấp	Bắc	Không mưa
8	Cao	Nam	Không mưa

Với kết quả này ta không cần phân lớp nữa cho lớp [Trời - Trong].

- Phân lớp nhánh [Trời - Mây]:

- Bảng dữ liệu của nhánh:

#	Áp Suất	Gió	Kết quả
2	Cao	Nam	Mưa
3	Trung bình	Bắc	Mưa
5	Thấp	Bắc	Mưa
6	Cao	Bắc	Mưa
7	Thấp	Nam	Không mưa

$$I(4,1) = 0.72$$

- Tính entropy cho thuộc tính [Áp Suất]:

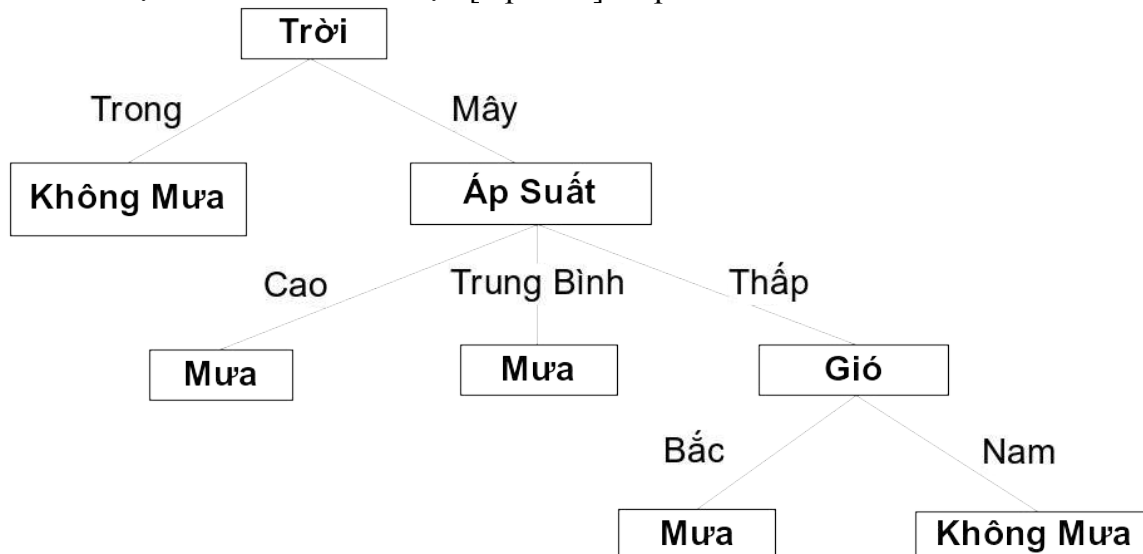
Áp Suất	Pi	ni	I(pi, ni)
Cao	2	0	0.00
Trung Bình	1	0	0.00
Thấp	1	1	1.0
<b>E (Áp Suất) = 2/5*0 + 1/5*0 + 2/5*1 = 0.4</b>			
<b>Gain(Áp Suất) = I(4,1) - E (Áp Suất) = 0.32</b>			

- Tính entropy cho thuộc tính [Gió]:

Gió	pi	Ni	I(pi, ni)
Bắc	3	0	0.00
Nam	1	1	1.0

$E(\text{Gió}) = 3/5 \cdot 0 + 2/5 \cdot 1 = 0.4$
$\text{Gain}(\text{Trời}) = I(4,1) - E(\text{Gió}) = 0.32$

Thuộc tính [Gió] và [Áp Suất] có Gain bằng nhau, vì [Áp Suất] có nhiều thuộc tính hơn nên ta chọn [Áp Suất] để phân tích:



Vậy, ta có các luật sau:

- Trời trong => Không mưa
- Trời mây, Áp Suất Cao => Mưa
- Trời mây, Áp Suất Trung Bình => Mưa
- Trời mây, Áp Suất Thấp, Gió Bắc => Mưa
- Trời mây, Áp Suất Thấp, Gió Nam => Không mưa

## Bài Tập 2: Luật Phân Lớp

Dùng thuật toán ID3 và Naïve Bayes để tìm luật phân lớp trong bảng sau đây.

TT	Màu Tóc	Chiều Cao	Cân Nặng	Dùng Thuốc?	Kết Quả
1	Đen	Tầm thước	Nhẹ	Không	Bị râm
2	Đen	Cao	Vừa phải	Có	Không
3	Râm	Thấp	Vừa phải	Có	Không
4	Đen	Thấp	Vừa phải	Không	Bị râm
5	Bạc	Tầm thước	Nặng	Không	Bị râm
6	Râm	Cao	Nặng	Không	Không
7	Râm	Tầm thước	Nặng	Không	Không
8	Đen	Thấp	Nhẹ	Có	Không

So sánh kết quả.

Bài giải:

- Sử dụng thuật toán ID3:

**Ký hiệu:** P: Bị Rám; N: Không

$$I(P, N) = -P/(P+N) \cdot \log_2(P/(P+N)) - N/(P+N) \cdot \log_2(N/(P+N))$$

Ta có: P = 3 và N = 5

$$\Rightarrow I(3, 5) = 0.95$$

- Tính entropy cho thuộc tính **[Màu Tóc]**:

Màu Tóc	pi	ni	I(pi, ni)
Đen	2	2	1.00
Râm	0	3	0.00
Bạc	1	0	0.00
<b>E (Màu Tóc) = 4/8*1 + 3/8*0 + 1/8*0 = 0.5</b>			
<b>Gain(Màu Tóc) = I(3,5) - E (Màu Tóc) = 0.45</b>			

- Tính entropy cho thuộc tính **[Chiều Cao]**:

Chiều Cao	Pi	ni	I(pi, ni)
Cao	0	2	0.00
Tầm thước	2	1	0.92
Thấp	1	2	0.92
<b>E (Chiều Cao) = 2/8*0 + 3/8*0.92 + 3/8*0.92 = 0.69</b>			
<b>Gain(Chiều Cao) = I(3,5) - E (Chiều Cao) = 0.26</b>			

- Tính entropy cho thuộc tính **[Cân Nặng]**:

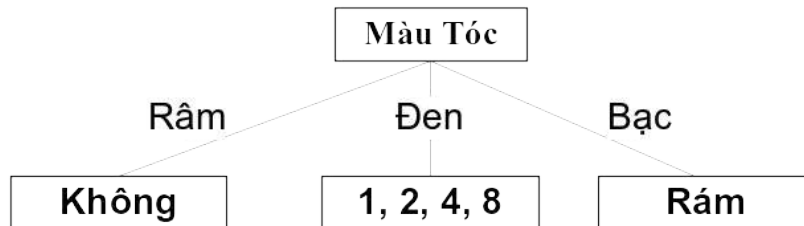
Cân Nặng	pi	ni	I(pi, ni)
Nặng	1	2	0.92
Vừa phải	1	2	0.92
Nhẹ	1	1	1.00
<b>E (Cân Nặng) = 3/8*0.92 + 3/8*0.92 + 2/8*1 = 0.94</b>			
<b>Gain(Cân Nặng) = I(3,5) - E (Cân Nặng) = 0.01</b>			

- Tính entropy cho thuộc tính **[Dùng Thuốc]**:

Dùng Thuốc	Pi	ni	I(pi, ni)
Không	3	2	0.97
Có	0	3	0.00
<b>E (Dùng Thuốc) = 5/8*0.97 + 3/8*0 = 0.6</b>			
<b>Gain(Dùng Thuốc) = I(3,5) - E (Dùng Thuốc) = 0.35</b>			

Ta nhận thấy Gain của thuộc tính **[Màu Tóc]** là lớn nhất, nên ta dùng thuộc tính **[Màu Tóc]** để phân lớp:





- Phân lớp nhánh [Màu Tóc - Đen]:

- Bảng dữ liệu của nhánh:

TT	Chiều Cao	Cân Nặng	Dùng Thuốc?	Kết Quả
1	Tầm thước	Nhẹ	Không	Bị râm
2	Cao	Vừa phải	Có	Không
4	Thấp	Vừa phải	Không	Bị râm
8	Thấp	Nhẹ	Có	Không

$$I(2,2) = 1$$

- Tính entropy cho thuộc tính [Chiều Cao]:

Chiều Cao	Pi	Ni	I(pi, ni)
Cao	0	1	0.00
Tầm thước	1	0	0.00
Thấp	1	1	1.00
<b>E (Chiều Cao) = <math>\frac{1}{4} \cdot 0 + \frac{1}{4} \cdot 0 + \frac{2}{4} \cdot 1.00 = 0.5</math></b>			
<b>Gain(Chiều Cao) = I(2,2) - E (Chiều Cao) = 0.5</b>			

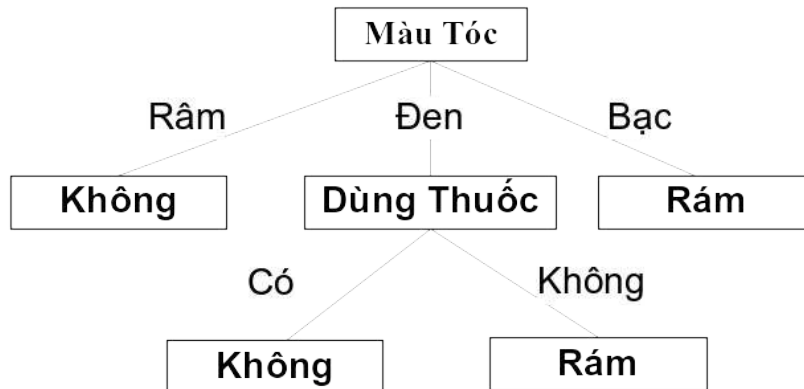
- Tính entropy cho thuộc tính [Cân Nặng]:

Cân Nặng	pi	Ni	I(pi, ni)
Nặng	0	0	0.00
Vừa phải	1	1	1.00
Nhẹ	1	1	1.00
<b>E (Cân Nặng) = <math>\frac{0}{4} \cdot 0 + \frac{2}{4} \cdot 1 + \frac{2}{4} \cdot 1 = 1</math></b>			
<b>Gain(Cân Nặng) = I(2,2) - E (Cân Nặng) = 0</b>			

- Tính entropy cho thuộc tính [Dùng Thuốc]:

Dùng Thuốc	Pi	Ni	I(pi, ni)
Không	2	0	0.00
Có	0	2	0.00
<b>E (Dùng Thuốc) = <math>\frac{2}{4} \cdot 0 + \frac{2}{4} \cdot 0 = 0</math></b>			
<b>Gain(Dùng Thuốc) = I(2,2) - E (Dùng Thuốc) = 1</b>			

Ta nhận thấy Gain của thuộc tính [Dùng Thuốc] là lớn nhất, nên ta dùng thuộc tính [Dùng Thuốc] để phân lớp:



Vậy, ta có các luật sau:

- Màu Tóc râm => Không Rám
- Màu Tóc bạc => Rám
- Màu Tóc đen, dùng thuốc => Không Rám
- Màu Tóc đen, không dùng thuốc => Rám

#### • Sử dụng thuật toán Naive Bayes

- $P(P) = 3/8 = 0.375$
- $P(N) = 5/8 = 0.625$

Màu Tóc	
$P(\text{Đen}/P) = 2/3$	$P(\text{Đen}/N) = 2/5$
$P(\text{Râm}/P) = 0/3$	$P(\text{Râm}/N) = 3/5$
$P(\text{Bạc}/P) = 1/3$	$P(\text{Bạc}/N) = 0/5$

Chiều Cao	
$P(\text{Cao}/P) = 0/3$	$P(\text{Cao}/N) = 2/5$
$P(\text{Tầm Thước}/P) = 2/3$	$P(\text{Tầm Thước}/N) = 1/5$
$P(\text{Thấp}/P) = 1/3$	$P(\text{Thấp}/N) = 2/5$

Cân Nặng	
$P(\text{Nặng}/P) = 1/3$	$P(\text{Nặng}/N) = 2/5$
$P(\text{Vừa Phái}/P) = 1/3$	$P(\text{Vừa Phái}/N) = 2/5$
$P(\text{Nhẹ}/P) = 1/3$	$P(\text{Nhẹ}/N) = 1/5$

Dùng Thuốc	
$P(\text{Không}/P) = 3/3$	$P(\text{Không}/N) = 2/5$
$P(\text{Có}/P) = 0/3$	$P(\text{Có}/N) = 3/5$

- Phân lớp X: Xét các mẫu chưa tìm thấy

- $X1 = \{\text{Đen, Tầm thước, Nhẹ, Có}\}$
- $P(X1,P).P(P) = P(\text{Đen},P)*P(\text{Tầm Thước},P)*P(\text{Nhẹ},P)*P(\text{Có},P)*P(P) = 2/3*2/3*1/3*0/3*0.375 = 0$
- $P(X1,N).P(N) = P(\text{Đen},N)*P(\text{Tầm Thước},N)*P(\text{Nhẹ},N)*P(\text{Có},N)*P(N) = 2/5*1/5*1/5*3/5*0.625 = 0.006$   
Mẫu X2 được phân vào lớp N
- $X2 = \{\text{Đen, Tầm thước, Vừa Phải, Không}\}$
- $P(X2,P).P(P) = P(\text{Đen},P)*P(\text{Tầm Thước},P)*P(\text{Vừa Phải},P)*P(\text{Không},P)*P(P) = 2/3*2/3*1/3*3/3*0.375 = 0.05555$
- $P(X2,N).P(N) = P(\text{Đen},N)*P(\text{Tầm Thước},N)*P(\text{Vừa Phải},N)*P(\text{Không},N)*P(N) = 2/5*1/5*2/5*2/5*0.625 = 0.008$   
Mẫu X2 được phân vào lớp P
- $X3 = \{\text{Đen, Tầm thước, Vừa Phải, Có}\}$
- $P(X3,P).P(P) = P(\text{Đen},P)*P(\text{Tầm Thước},P)*P(\text{Vừa Phải},P)*P(\text{Có},P)*P(P) = 2/3*2/3*1/3*0/3*0.375 = 0$
- $P(X3,N).P(N) = P(\text{Đen},N)*P(\text{Tầm Thước},N)*P(\text{Vừa Phải},N)*P(\text{Có},N)*P(N) = 2/5*1/5*2/5*3/5*0.625 = 0.012$   
Mẫu X3 được phân vào lớp N
- $X4 = \{\text{Đen, Tầm thước, Nặng, Không}\}$
- $P(X4,P).P(P) = P(\text{Đen},P)*P(\text{Tầm Thước},P)*P(\text{Nặng},P)*P(\text{Không},P)*P(P) = 2/3*2/3*1/3*3/3*0.375 = 0.055$
- $P(X4,N).P(N) = P(\text{Đen},N)*P(\text{Tầm Thước},N)*P(\text{Nặng},N)*P(\text{Không},N)*P(N) = 2/5*1/5*2/5*2/5*0.625 = 0.008$   
Mẫu X4 được phân vào lớp P
- $X5 = \{\text{Đen, Tầm thước, Nặng, Có}\}$  phân vào lớp N vì  $P(\text{Có},P) = 0$
- $X6 = \{\text{Đen, Cao, ..., ...}\}$  phân vào lớp N vì  $P(\text{Cao}/P) = 0$
- $X7 = \{\text{Đen, Thấp, Nặng, Không}\}$
- $P(X7,P).P(P) = P(\text{Đen},P)*P(\text{Thấp},P)*P(\text{Nặng},P)*P(\text{Không},P)*P(P) = 2/3*1/3*1/3*3/3*0.375 = 0.028$
- $P(X7,N).P(N) = P(\text{Đen},N)*P(\text{Thấp},N)*P(\text{Nặng},N)*P(\text{Không},N)*P(N) = 2/5*2/5*2/5*2/5*0.625 = 0.016$   
Mẫu X7 được phân vào lớp P
- $X8 = \{\text{Đen, Thấp, Nặng, Có}\}$  phân vào lớp N vì  $P(\text{Có},P) = 0$
- $X9 = \{\text{Đen, Thấp, Nhẹ, Không}\}$

- $P(X_{10}, P).P(P) = P(\text{Đen}, P) * P(\text{Thấp}, P) * P(\text{Nhẹ}, P) * P(\text{Không}, P) * P(P) = 2/3 * 1/3 * 1/3 * 3/3 * 0.375 = 0.028$
- $P(X_{10}, N).P(N) = P(\text{Đen}, N) * P(\text{Thấp}, N) * P(\text{Nhẹ}, N) * P(\text{Không}, N) * P(N) = 2/5 * 2/5 * 1/5 * 2/5 * 0.625 = 0.008$   
Mẫu X10 được phân vào lớp P
- $X_{11} = \{\text{Đen}, \text{Thấp}, \text{Nhẹ}, \text{Có}\}$  phân vào lớp N vì  $P(\text{Có}, P) = 0$
- $X_{12} = \{\text{Râm}, \dots, \dots, \dots\}$  phân vào lớp N vì  $P(\text{Râm}, P) = 0$
- $X_{12} = \{\text{Bạc}, \dots, \dots, \dots\}$  phân vào lớp P vì  $P(\text{Bạc}, N) = 0$
- Rút ra các phân lớp:
  - Màu tóc Râm  $\Rightarrow$  Không Rám
  - Màu tóc Bạc  $\Rightarrow$  Rám
  - Chiều cao Cao  $\Rightarrow$  Không Rám
  - Có dùng Thuốc  $\Rightarrow$  Không Rám
  - Màu tóc Đen, Tầm Thước, Vừa Phải, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Tầm Thước, Nhẹ, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Tầm Thước, Nặng, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Thấp, Vừa Phải, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Thấp, Nhẹ, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Thấp, Nặng, Không  $\Rightarrow$  Rám
- Rút gọn các phân lớp:
  - Màu tóc Râm  $\Rightarrow$  Không Rám
  - Màu tóc Bạc  $\Rightarrow$  Rám
  - Chiều cao Cao  $\Rightarrow$  Không Rám
  - Có dùng Thuốc  $\Rightarrow$  Không Rám
  - Màu tóc Đen, Tầm Thước, Không  $\Rightarrow$  Rám
  - Màu tóc Đen, Thấp, Không  $\Rightarrow$  Rám

**Kết Luận:** dùng phương pháp Naive Bayes phức tạp hơn phương pháp ID3, tuy nhiên có khai phá ra được nhiều luật mới hơn ID3.

### Bài Tập 3: Tập Thô

Hãy rút gọn bảng quyết định sau đây:

**Bảng :** Số liệu quan sát về hiện tượng râm nắng.

TT	Tên người	Màu tóc	Chiều cao	Cân nặng	Dùng thuốc?	Kết quả
		M	C	N	T	R
1	Hoa	Đen	Tầm thước	Nhẹ	Không	Bị râm
2	Lan	Đen	Cao	Vừa phải	Có	Không
3	Xuân	Râm	Thấp	Vừa phải	Có	Không
4	Hạ	Đen	Thấp	Vừa phải	Không	Bị râm
5	Thu	Bạc	Tầm thước	Nặng	Không	Bị râm

6	Đông	Râm	Cao	Nặng	Không	Không
7	Mơ	Râm	Tầm thước	Nặng	Không	Không
8	Đào	Đen	Thấp	Nhẹ	Có	Không

Khi dùng cây quyết định, ta có các luật:

- IF Tóc đen  
Người đó dùng thuốc  
THEN Không sao cả
- IF Người tóc đen  
Không dùng thuốc  
THEN Họ bị râm nắng
- IF Người tóc bạc  
THEN Bị râm nắng
- IF Người tóc râm  
THEN Không sao cả

Sau khi rút gọn ( để có các reducts) , bạn có luật gì ???  
So sánh với kết quả tạo luật từ cây quyết định

Giải:

Ta có ma trận phân biệt:

	Hoa	Lan	Xuân	Hạ	Thu	Đông	Mơ
Lan	C, N, T						
Xuân	M,C,N,T						
Hạ	$\lambda$	C, T	M,T				
Thu	$\lambda$	M,C,N,T	M,C,N,T				
Đông	M, C, N	$\lambda$	$\lambda$	M,C,N	M,C		
Mơ	M,N	$\lambda$	$\lambda$	M,C,N	M	$\lambda$	
Đào	C,T	$\lambda$	$\lambda$	N,T	M,C,N,T	$\lambda$	$\lambda$

$$\Rightarrow F(M,C,N,T) = (C^{\vee}N^{\vee}T) \wedge (M^{\vee}C^{\vee}N^{\vee}T) \wedge (M^{\vee}N^{\vee}C) \wedge (M^{\vee}N) \wedge (C^{\vee}T) \wedge (M^{\vee}T) \wedge (N^{\vee}T) \wedge (M^{\wedge}C) \wedge (M) = M \wedge (C^{\vee}N^{\vee}T) \wedge (C^{\vee}T) \wedge (N^{\vee}T) = (M^{\wedge}T)^{\vee} (M^{\wedge}C^{\wedge}N)$$

○ Ta có rút gọn sau:

$$\text{Reduct1} = \{M, T\}$$

$$\text{Reduct2} = \{M, C, N\}$$

$$\text{Core} = \{\text{Kết quả}\} = \{M, C, N\} \cap \{M, T\} = \{M\}$$

○ Tìm các luật

○ Với  $C = \{\text{Kết quả}\}$  ta có các lớp tương đương sau:

- $X1 = \{1,4,5\}$  Bị rám
- $X2 = \{2,3,6,7,8\}$  Không bị rám
- Xét  $\text{Reduct1} = \{M, T\}$  với  $\{M, T\} \Rightarrow K$ , ta có các lớp tương đương:
  - $Z1 = \{1, 4\}$  Đen, Không dùng thuốc
  - $Z2 = \{2, 8\}$  Đen, Có dùng thuốc
  - $Z3 = \{3\}$  Râm, Có dùng thuốc
  - $Z4 = \{5\}$  Bạc, Không dùng thuốc
  - $Z5 = \{6, 7\}$  Râm, Không dùng thuốc
- Xét phân lớp  $X1 \{1,4,5\}$  Bị Rám:
  - $X1 \cap Z1 = \{1, 4\} \neq \emptyset$  và  $Z1 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Không dùng thuốc thì Bị rám.
  - $X1 \cap Z2 = \emptyset \Rightarrow$  không có luật phân lớp.
  - $X1 \cap Z3 = \emptyset \Rightarrow$  không có luật phân lớp.
  - $X1 \cap Z4 = \{5\} \neq \emptyset$  và  $Z4 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Bạc, Không dùng thuốc thì Bị rám.
  - $X1 \cap Z5 = \emptyset \Rightarrow$  không có luật phân lớp.

Vậy ta có các luật sau :

L1: Nếu tóc Đen, Không dùng thuốc thì Bị rám.

L5: Nếu tóc Bạc, Không dùng thuốc thì Bị rám.

- Xét phân lớp  $X2 \{2,3,6,7,8\}$  Không Bị Rám:
  - $X1 \cap Z1 = \emptyset \Rightarrow$  không có luật phân lớp.
  - $X1 \cap Z2 = \{2, 8\} \neq \emptyset$  và  $Z2 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Có dùng thuốc thì Không Bị rám.
  - $X1 \cap Z3 = \{3\} \neq \emptyset$  và  $Z3 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Râm, Có dùng thuốc thì Không Bị rám.
  - $X1 \cap Z4 = \emptyset \Rightarrow$  không có luật phân lớp.
  - $X1 \cap Z5 = \{6,7\} \neq \emptyset$  và  $Z5 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Râm, Không dùng thuốc thì Không Bị rám.
  - $X1 \cap Z5 = \emptyset \Rightarrow$  không có luật phân lớp.

Vậy ta có các luật sau :

L2: Nếu tóc Đen, Có dùng thuốc thì Không Bị rám.

L3: Nếu tóc Râm, Có dùng thuốc thì Không Bị rám.

L6: Nếu tóc Râm, Không dùng thuốc thì Không Bị rám.

Vậy đối với trường hợp  $\{M, T\} \rightarrow \{K\}$  . Ta có các luật sau :

L1: Nếu tóc Đen, Không dùng thuốc thì Bị rám.

L5: Nếu tóc Bạc, Không dùng thuốc thì Bị rám.

L2: Nếu tóc Đen, Có dùng thuốc thì Không bị râm.  
 L3: Nếu tóc Râm, Có dùng thuốc thì Không bị râm.  
 L6: Nếu tóc Râm, Không dùng thuốc thì Không bị râm.

Rút gọn luật L3 và L6. Ta có các luật còn lại:

L1: Nếu tóc Đen, Không dùng thuốc thì Bị râm.  
 L5: Nếu tóc Bạc, Không dùng thuốc thì Bị râm.  
 L2: Nếu tóc Đen, Có dùng thuốc thì Không bị râm.  
 L3: Nếu tóc Râm thì Không bị râm.

- Xét Reduct2 = {M, C, N} với {M, C, N}  $\Rightarrow$  K, ta có các lớp tương đương:

Z1 = {1} Đen, Tầm thước, Nhẹ  
 Z2 = {2} Đen, Cao, Vừa phải  
 Z3 = {3} Râm, Thấp, Vừa phải  
 Z4 = {4} Đen, Thấp, Vừa phải  
 Z5 = {5} Bạc, Tầm thước, Nặng  
 Z6 = {6} Râm, Cao, Nặng  
 Z7 = {7} Râm, Tầm thước, Nặng  
 Z8 = {8} Đen, Thấp, Nhẹ

- Xét phân lớp X1 {1,4,5} Bị Râm:

- $X1 \cap Z1 = \{1\} \neq \emptyset$  và  $Z1 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Tầm thước, Nhẹ thì Bị râm.
- $X1 \cap Z2 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z3 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z4 = \{4\} \neq \emptyset$  và  $Z4 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Thấp, Vừa phải thì Bị râm.
- $X1 \cap Z5 = \{5\} \neq \emptyset$  và  $Z5 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Bạc, Tầm thước, Nặng thì Bị râm.
- $X1 \cap Z6 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z7 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z8 = \emptyset \Rightarrow$  không có luật phân lớp.

Vậy ta có các luật sau :

L'1: Nếu tóc Đen, Tầm thước, Nhẹ thì Bị râm.  
 L'4: Nếu tóc Đen, Thấp, Vừa phải thì Bị râm  
 L'5: Nếu tóc Bạc, Tầm thước, Nặng thì Bị râm.

- Xét phân lớp X2 {2,3,6,7,8} Không Bị Râm:

- $X2 \cap Z1 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X2 \cap Z2 = \{2\} \neq \emptyset$  và  $Z2 \subseteq X2$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Cao, Vừa phải thì Không Bị râm.

- $X1 \cap Z3 = \{3\} \neq \emptyset$  và  $Z3 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Râm, Thấp, Vừa phải thì Không Bị râm.
- $X1 \cap Z4 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z5 = \emptyset \Rightarrow$  không có luật phân lớp.
- $X1 \cap Z6 = \{6\} \neq \emptyset$  và  $Z6 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Râm, Cao, Nặng thì Không Bị râm.
- $X1 \cap Z7 = \{7\} \neq \emptyset$  và  $Z7 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Râm, Tầm thước, Nặng thì Không Bị râm.
- $X1 \cap Z8 = \{8\} \neq \emptyset$  và  $Z8 \subseteq X1$  nên ta có luật phân lớp đúng chính xác 100%. Vậy ta có luật : Nếu tóc Đen, Thấp, Nhẹ thì Không Bị râm.

Vậy ta có các luật sau :

- L'2: Nếu tóc Đen, Cao, Vừa phải thì Không bị râm.
- L'3: Nếu tóc Râm, Thấp, Vừa phải thì Không bị râm.
- L'6: Nếu tóc Râm, Cao, Nặng thì Không bị râm.
- L'7: Nếu tóc Râm, Tầm thước, Nặng thì Không bị râm.
- L'8: Nếu tóc Đen, Thấp, Nhẹ thì Không bị râm.

Vậy đối với trường hợp  $\{M, T\} \rightarrow \{K\}$  . Ta có các luật sau :

- L'1: Nếu tóc Đen, Tầm thước, Nhẹ thì Bị râm.
- L'4: Nếu tóc Đen, Thấp, Vừa phải thì Bị râm
- L'5: Nếu tóc Bạc, Tầm thước, Nặng thì Bị râm.
- L'2: Nếu tóc Đen, Cao, Vừa phải thì Không bị râm.
- L'3: Nếu tóc Râm, Thấp, Vừa phải thì Không bị râm.
- L'6: Nếu tóc Râm, Cao, Nặng thì Không bị râm.
- L'7: Nếu tóc Râm, Tầm thước, Nặng thì Không bị râm.
- L'8: Nếu tóc Đen, Thấp, Nhẹ thì Không bị râm.

**Kết hợp  $\{M, T\} \Rightarrow \{K\}$  và  $\{M, C, N\} \Rightarrow \{K\}$  , ta có tổng cộng các luật sau:**

- L1: Nếu tóc Đen, Không dùng thuốc thì Bị râm.
- L5: Nếu tóc Bạc, Không dùng thuốc thì Bị râm.
- L2: Nếu tóc Đen, Có dùng thuốc thì Không bị râm.
- L3: Nếu tóc Râm thì Không bị râm.
- L'1: Nếu tóc Đen, Tầm thước, Nhẹ thì Bị râm.
- L'4: Nếu tóc Đen, Thấp, Vừa phải thì Bị râm
- L'5: Nếu tóc Bạc, Tầm thước, Nặng thì Bị râm.
- L'2: Nếu tóc Đen, Cao, Vừa phải thì Không bị râm.
- L'3: Nếu tóc Râm, Thấp, Vừa phải thì Không bị râm.



L'6: Nếu tóc Râm, Cao, Nặng thì Không bị râm.  
 L'7: Nếu tóc Râm, Tầm thước, Nặng thì Không bị râm.  
 L'8: Nếu tóc Đen, Thấp, Nhẹ thì Không bị râm.

Loại bỏ các luật thừa , ta có :

L1: Nếu tóc Đen, Không dùng thuốc thì Bị râm.  
 L5: Nếu tóc Bạc, Không dùng thuốc thì Bị râm.  
 L2: Nếu tóc Đen, Có dùng thuốc thì Không bị râm.  
 L3: Nếu tóc Râm thì Không bị râm.  
 L'5: Nếu tóc Bạc, Tầm thước, Nặng thì Bị râm.

Kết hợp L5 và L'5 , ta có các luật sau cùng:

**L1: Nếu tóc Đen, Không dùng thuốc thì Bị râm.**  
**L5: Nếu tóc Bạc thì Bị râm.**  
**L2: Nếu tóc Đen, Có dùng thuốc thì Không bị râm.**  
**L3: Nếu tóc Râm thì Không bị râm.**

**Kết luận:** Kết quả tạo luật từ cây quyết định và kết quả tạo luật từ rút gọn các reducts thì giống nhau.

### Bài Tập 3: Episodes

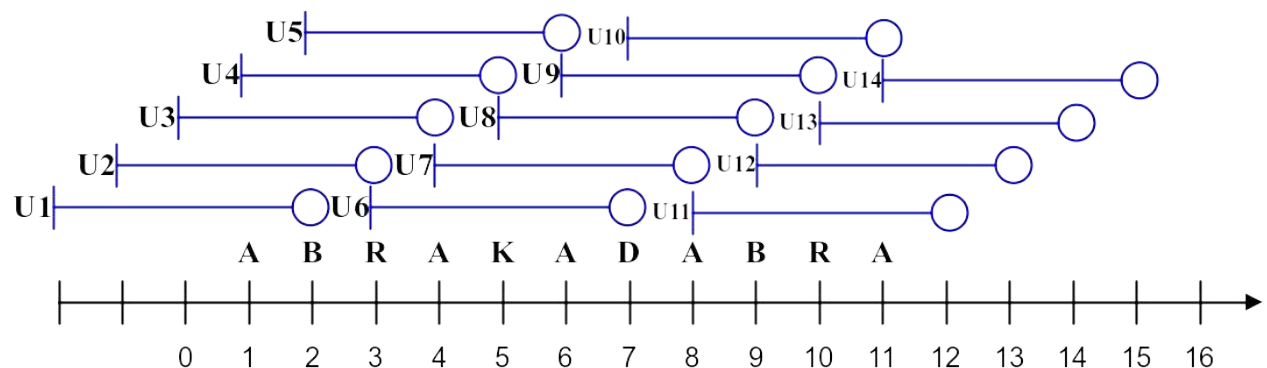
Cho chuỗi sự kiện sau đây:

**A B R A K A D A B R A**

- 1) Có bao nhiêu cửa sổ có bề rộng là 4 sự kiện được xử lý để tìm các episodes phổ biến theo tiếp cận WINEPI ?
- 2) Giả sử ngưỡng min\_fr là 0.3. Tìm các episode phổ biến tuần tự và song song trong chuỗi sự kiện nêu trên trên ?
- 3) Tìm các episode tối đại ?

Bài Giải :

**1\ Có bao nhiêu cửa sổ có bề rộng là 4 sự kiện được xử lý để tìm các episodes phổ biến theo tiếp cận WINEPI?**



- Bằng cách trượt cửa sổ, chúng ta có 14 cửa sổ, có bề rộng là 4 sự kiện:

Cửa Sổ $U_i$	Nội dung của $U_i$	Episodes song song xảy ra trong $U_i$
$U_{1,[-2,2]}$	$[-,-,A]$	$\{A\}$
$U_{2,[-1,3]}$	$[-,A,B]$	$\{A,B\}, \{AB\}$
$U_{3,[0,4]}$	$[-,A,B,R]$	$\{A,B,R\}, \{AB,AR,BR\}, \{ABR\}$
$U_{4,[1,5]}$	$[A,B,R,A]$	$\{A,B,R\}, \{AB,AR,BR\}, \{ABR\}$
$U_{5,[2,6]}$	$[B,R,A,K]$	$\{A,B,K,R\}, \{AB,AK,AR,BK,BR,KR\}, \{ABK,ABR,AKR,BKR\}, \{ABKR\}$
$U_{6,[3,7]}$	$[R,A,K,A]$	$\{A,K,R\}, \{AK,AR,KR\}, \{AKR\}$
$U_{7,[4,8]}$	$[A,K,A,D]$	$\{A,D,K\}, \{AD,AK,DK\}, \{ADK\}$
$U_{8,[5,9]}$	$[K,A,D,A]$	$\{A,D,K\}, \{AD,AK,DK\}, \{ADK\}$
$U_{9,[6,10]}$	$[A,D,A,B]$	$\{A,B,D\}, \{AB,AD,BD\}, \{ABD\}$
$U_{10,[7,11]}$	$[D,A,B,R]$	$\{A,B,D,R\}, \{AB,AD,AR,BD,BR,DR\}, \{ABD,ABR,ADR,BDR\}, \{ABDR\}$
$U_{11,[8,12]}$	$[A,B,R,A]$	$\{A,B,R\}, \{AB,AR,BR\}, \{ABR\}$
$U_{12,[9,13]}$	$[B,R,A,-]$	$\{A,B,R\}, \{AB,AR,BR\}, \{ABR\}$
$U_{13,[10,14]}$	$[R,A,-,-]$	$\{A,R\}, \{AR\}$
$U_{14,[11,15]}$	$[A,-,-,-]$	$\{A\}$

**2\ Giả sử ngưỡng min\_fr là 0.3. Tìm các episode phổ biến tuần tự và song song trong chuỗi sự kiện nêu trên trên ?**

**a\ Tìm các Episode song song**

- Xây dựng episode L1 chứa 1 phần tử:

Episode	Tần suất của Episode
A	$fr(A, S, W) = 14/14 = 1 > min\_fr$
B	$fr(B, S, W) = 8/14 = 0.57 > min\_fr$
D	$fr(D, S, W) = 4/14 = 0.29 < min\_fr$
K	$fr(K, S, W) = 4/14 = 0.29 < min\_fr$
R	$fr(R, S, W) = 8/14 = 0.57 > min\_fr$

**Tập Episode phổ biến song song có 1 sự kiện là  $L1 = \{A, B, R\}$**

- Xây dựng episode L2 chứa 2 phần tử:

Episode	Tần suất của Episode
AB	$fr(AB, S, W) = 8/14 = 0.57 > min\_fr$
AR	$fr(AR, S, W) = 8/14 = 0.57 > min\_fr$
BR	$fr(BR, S, W) = 6/14 = 0.43 > min\_fr$

**Tập Episode phổ biến song song có 2 sự kiện là  $L2 = \{AB, AR, BR\}$**

- Xây dựng episode L3 chứa 3 phần tử:

Episode	Tần suất của Episode
ABR	$fr(ABR, S, W) = 6/14 = 0.43 > \min\_fr$
Tập Episode phổ biến song song có 3 sự kiện là $L3 = \{ABR\}$	

- Không có Episode phổ biến song song có 4 sự kiện từ  $L3$

**Tập Episode phổ biến song song**

$$= L1 \cup L2 \cup L3$$

$$= \{A, B, R, AB, AR, BR, ABR\}$$

**Tập phổ biến song song tối đại chính là tập phổ biến  $L3 = \{ABR\}$**

#### b\ Tìm các Episode tuần tự

- Xây dựng episode  $L1$  chứa 1 phần tử:

Episode	Tần suất của Episode
A	$fr(A, S, W) = 14/14 = 1 > \min\_fr$
B	$fr(B, S, W) = 8/14 = 0.57 > \min\_fr$
D	$fr(D, S, W) = 4/14 = 0.29 < \min\_fr$
K	$fr(K, S, W) = 4/14 = 0.29 < \min\_fr$
R	$fr(R, S, W) = 8/14 = 0.57 > \min\_fr$
Tập Episode phổ biến song song có 1 sự kiện là $L1 = \{A, B, R\}$	

- Xây dựng episode  $L2$  chứa 2 phần tử:

Episode	Tần suất của Episode
AB	$fr(AB, S, W) = 6/14 = 0.43 > \min\_fr$
BA	$fr(BA, S, W) = 4/14 = 0.29 < \min\_fr$
AR	$fr(AR, S, W) = 4/14 = 0.29 < \min\_fr$
RA	$fr(RA, S, W) = 6/14 = 0.43 > \min\_fr$
BR	$fr(BR, S, W) = 6/14 = 0.43 > \min\_fr$
RB	$fr(RB, S, W) = 0/14 = 0.00 < \min\_fr$
Tập Episode phổ biến song song có 2 sự kiện là $L2 = \{AB, RA, BR\}$	

- Xây dựng episode  $L3$  chứa 3 phần tử:

Episode	Tần suất của Episode
ABR	$fr(ABR, S, W) = 4/14 = 0.29 < \min\_fr$
RAB	$fr(RAB, S, W) = 0/14 = 0.00 < \min\_fr$
BRA	$fr(BRA, S, W) = 4/14 = 0.29 < \min\_fr$
Tập Episode phổ biến song song có 3 sự kiện là $L3 = \{\}$	

- Không có Episode phổ biến song song có 4 sự kiện từ  $L3$

**Tập Episode phổ biến song song**

$$= L1 \cup L2$$

$$= \{A, B, R, AB, RA, BR\}$$

**Tập phổ biến song song tối đại chính là tập phổ biến  $L_2 = \{AB, RA, BR\}$**