

Phát hiện các luật kết hợp trong cơ sở dữ liệu

Nguyễn Hồng Phương
Bộ môn Hệ thống thông tin
Viện CNTT&TT – trường ĐHBK Hà Nội
phuongnh@soict.hut.edu.vn
<http://is.hut.edu.vn/~phuongnh>



1

Nội dung trình bày

- ❑ 1. Tổng quan
- ❑ 2. Phát hiện luật kết hợp trong cơ sở dữ liệu giao dịch
- ❑ 3. Phát hiện luật kết hợp trong cơ sở dữ liệu quan hệ
- ❑ 4. Một số vấn đề khác

Phát hiện luật kết hợp trong cơ sở dữ liệu

2

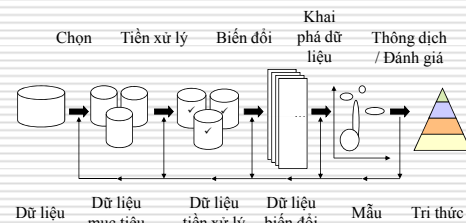
1. Tổng quan

- ❑ Khai phá dữ liệu và phát hiện tri thức
- ❑ Luật kết hợp: Bài toán “Cái giỏ hàng”
- ❑ Một số ứng dụng khác
- ❑ Các khái niệm cơ bản

Phát hiện luật kết hợp trong cơ sở dữ liệu

3

Khai phá dữ liệu và phát hiện tri thức



Phát hiện luật kết hợp trong cơ sở dữ liệu

4

Luật kết hợp: Bài toán “Cái giỏ hàng”

- ❑ Phân tích thói quen mua hàng của khách hàng: tìm sự kết hợp và tương quan giữa các mặt hàng khác nhau mà khách hàng đặt vào trong “giỏ hàng” của họ

Sữa, trứng, đường,
bánh mì



Khách hàng 1

Sữa, trứng, ngũ cốc, bánh mì



Khách hàng 2

Trứng, đường



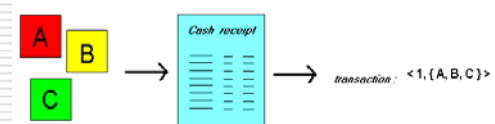
Khách hàng 3

Phát hiện luật kết hợp trong cơ sở dữ liệu

5

Phân tích bài toán “Cái giỏ hàng”

- ❑ Cho cơ sở dữ liệu gồm các giao dịch của khách hàng, mỗi giao dịch là một tập các mặt hàng
- ❑ Tìm các nhóm mặt hàng thường được mua cùng nhau



Phát hiện luật kết hợp trong cơ sở dữ liệu

6

Một số ứng dụng khác

□ Viễn thông

- Mỗi khách hàng là một giao dịch gồm một tập các cuộc gọi của khách hàng đó

□ Hiện tượng khí quyển

- Mỗi khoảng thời gian quan sát là một giao dịch chứa một tập các sự kiện quan sát được (mưa, gió, mây,...)

Phát hiện luật kết hợp trong cơ sở dữ liệu

7

Các khái niệm cơ bản

Giao dịch:

Dạng quan hệ

$\langle \text{Tid}, \text{item} \rangle$

$\langle 1, \text{item1} \rangle$

$\langle 1, \text{item2} \rangle$

$\langle 2, \text{item3} \rangle$

Dạng thu gọn

$\langle \text{Tid}, \text{itemset} \rangle$

$\langle 1, \{\text{item1}, \text{item2}\} \rangle$

$\langle 2, \{\text{item3}\} \rangle$

Mục (Item): phần tử đơn,

Tập mục (Itemset): Tập các mục

Độ hỗ trợ của 1 tập mục X - $\text{sup}(X)$: Số giao dịch chứa X

Độ hỗ trợ tối thiểu **minsup** : ngưỡng của độ hỗ trợ

Tập mục thường xuyên : độ hỗ trợ \geq **minsup**.

Phát hiện luật kết hợp trong cơ sở dữ liệu

8

Tập mục thường xuyên

ID giao dịch	Các mặt hàng đã mua
1	Sữa, trứng, đường, bánh mì
2	Sữa, trứng, ngũ cốc, bánh mì
3	Trứng, đường

- $\text{Sup}(\{\text{Sữa, trứng, bánh mì}\}) = 2$ (66.6%)
- $\text{Sup}(\{\text{Trứng, đường}\}) = 2$ (66.6%)
- $\text{Sup}(\{\text{Ngũ cốc, bánh mì}\}) = 1$ (33.3%)
- Nếu **minsup** = 50% thì $\{\text{Sữa, trứng, bánh mì}\}$ và $\{\text{Trứng, đường}\}$ là các tập mục thường xuyên còn $\{\text{Ngũ cốc, bánh mì}\}$ thì không phải.

Phát hiện luật kết hợp trong cơ sở dữ liệu

9

Luật kết hợp

- A, B là tập các mục trong tập mục I

- Luật $r = A \Rightarrow B$

- Độ hỗ trợ của r: $\text{sup}(r) = \text{sup}(A \cup B)$

- Độ tin cậy của r:

- $\text{conf}(r) = \text{sup}(A \cup B) / \text{sup}(A)$

- r được gọi là luật kết hợp nếu $\text{sup}(r) \geq$ **minsup** và $\text{conf}(r) \geq$ **minconf**

Độ hỗ trợ tối thiểu

Độ tin cậy tối thiểu

Phát hiện luật kết hợp trong cơ sở dữ liệu

10

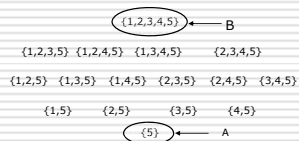
Hai tính chất cơ bản

□ Tính chất 1:

- Nếu một tập mục là không thường xuyên thì các siêu tập của nó cũng không thường xuyên

□ Tính chất 2:

- Nếu một tập mục là thường xuyên thì các tập con của nó cũng thường xuyên



Phát hiện luật kết hợp trong cơ sở dữ liệu

11

2. Phát hiện luật kết hợp trong CSDL giao dịch

- Phát hiện các tập mục thường xuyên

- Kiểu Apriori
- Sử dụng FP-tree

- Phát hiện các luật kết hợp

- Khai phá luật kết hợp đa mức

Phát hiện luật kết hợp trong cơ sở dữ liệu

12

Phát hiện các tập mục thường xuyên

- ❑ Giải thuật Apriori
- ❑ Sử dụng FP-tree

Phát hiện luật kết hợp trong cơ sở dữ liệu

13

Giải thuật Apriori

Đầu vào: Cơ sở dữ liệu các giao dịch D và s_{min}

Đầu ra: Tập Answer chứa tất cả các tập mục thường xuyên của D

Giải thuật:

```

1)  $L_1 = \{\text{large 1-itemsets}\};$ 
2) for( $k=2; L_{k-1} \neq \emptyset; k++$ ) do begin
3)    $C_k = \text{AprioriGen}(L_{k-1});$  // New candidate
4)   forall transactions  $t \in D$  do begin
5)      $C_t = \text{Subset}(C_k, t);$  // Candidates contained in t
6)     forall candidates  $c \in C_t$  do
7)        $c.\text{count}++$ 
8)   end
9)    $L_k = \{c \in C_k \mid c.\text{count} \geq s_{min}\}$ 
10) end
11) Answer =  $\cup_k L_k;$ 
    
```

Phát hiện luật kết hợp trong cơ sở dữ liệu

14

Hàm AprioriGen

Đầu vào: Một tập L_{k-1} chứa tất cả các $(k-1)$ -tập mục thường xuyên

Đầu ra: Tập C_k ứng cử là một siêu tập chứa tất cả các k -tập mục thường xuyên

Giải thuật:

```

1) Function AprioriGen( $L_{k-1}$ : tập  $(k-1)$ -tập mục thường xuyên):tập  $k$ -tập mục thường xuyên
2)   // Pha kết nối
3)   insert into  $C_k$ 
4)   select  $p.\text{item}_1, p.\text{item}_2, \dots, p.\text{item}_{k-1}, q.\text{item}_{k-1}$ 
5)   from  $L_{k-1} p, L_{k-1} q$ 
6)   where  $p.\text{item}_1 = q.\text{item}_1, \dots, p.\text{item}_{k-2} = q.\text{item}_{k-2}, p.\text{item}_{k-1} < q.\text{item}_{k-1}$ 
7)   // Pha cắt tỉa
8)   forall itemsets  $c \in C_k$  do
9)     forall  $(k-1)$ -subsets  $s$  of  $c$  do
10)      if ( $s \notin L_{k-1}$ ) then delete  $c$  from  $C_k;$ 
11)   return  $C_k;$ 
12) end;
    
```

Phát hiện luật kết hợp trong cơ sở dữ liệu

15

Vấn đề của giải thuật kiểu Apriori

- ❑ Chi phí cho việc kiểm soát một số lượng lớn các tập mục ứng cử
 - 10^4 1-tập mục thường xuyên sẽ sinh 10^7 tập ứng cử kích thước 2
 - Lặp nhiều lần việc duyệt CSDL để kiểm tra các tập ứng cử

➡ Tránh việc sinh quá nhiều tập ứng cử

➡ Sử dụng cấu trúc cây mẫu thường xuyên

Phát hiện luật kết hợp trong cơ sở dữ liệu

16

Xây dựng cây mẫu thường xuyên

- ❑ FP-tree (Frequent Pattern tree)
- ❑ Các bước xây dựng:
 - Duyệt DB lần 1 để sinh ra danh sách L

TID	Items		Item frequency
100	f, a, c, d, g, i, m, p	→	f 4
200	a, b, c, f, l, m, o		c 4
300	b, f, h, j, o		a 3
400	b, c, k, s, p		b 3
500	a, f, c, e, l, p, m, n		m 3
			p 3

Phát hiện luật kết hợp trong cơ sở dữ liệu

17

Xây dựng cây mẫu thường xuyên

- Duyệt DB lần 2, sắp xếp lại các giao dịch theo danh sách L

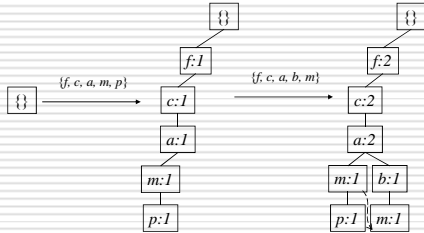
TID	Items	Các mục đã sắp xếp
100	f, a, c, d, g, i, m, p	f, c, a, m, p
200	a, b, c, f, l, m, o	f, c, a, b, m
300	b, f, h, j, o	f, b
400	b, c, k, s, p	c, b, p
500	a, f, c, e, l, p, m, n	f, c, a, m, p

Phát hiện luật kết hợp trong cơ sở dữ liệu

18

Xây dựng cây mẫu thường xuyên

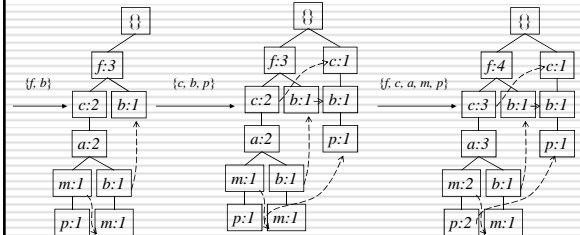
Tiến hành xây dựng cây



Phát hiện luật kết hợp trong cơ sở dữ liệu

19

Xây dựng cây mẫu thường xuyên

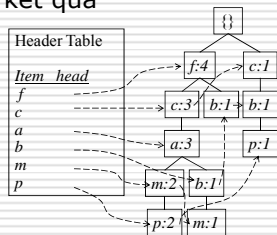


Phát hiện luật kết hợp trong cơ sở dữ liệu

20

Xây dựng cây mẫu thường xuyên

Cây kết quả



Khai phá mẫu thường xuyên?

Phát hiện luật kết hợp trong cơ sở dữ liệu

21

Phát hiện các luật kết hợp

Giải thuật đơn giản để sinh các luật

Đầu vào: Tập tất cả các tập mục thường xuyên có nhiều hơn một mục

$$\bar{F} = \bigcup_{k \geq 2} F_k = F \setminus F_1$$

Đầu ra: Tất cả các luật kết hợp

Phương pháp:

- for all $f_k \in \bar{F}$ do
- GenRules(f_k, f_k);

Phát hiện luật kết hợp trong cơ sở dữ liệu

22

Phát hiện các luật kết hợp

Thủ tục GenRules

Đầu vào: Hai tập mục thường xuyên f_k và l_m , và một ngưỡng độ tin cậy c_{min}
Đầu ra: Các luật kết hợp với nhiều nhất m-1 mục ở phần đầu luật (m>2)
Phương pháp:

```

1) procedure GenRules( $f_k$ : k-tập mục thường xuyên,  $l_m$ : m-tập mục thường xuyên)
2)    $L \leftarrow \{ \text{các } (m-1)\text{-tập mục } l_{m-1} | l_{m-1} \subset l_m \}$ 
3)   forall  $l_{m-1} \in L$  do begin
4)      $c \leftarrow s(f_k) / s(l_{m-1})$ ; // Độ chắc chắn của luật
5)     if  $c \geq c_{min}$  then begin
6)       output luật  $l_{m-1} \Rightarrow (f_k \setminus l_{m-1})$ ;
7)       if  $m-1 \geq 1$  then
8)         GenRules( $f_k, l_{m-1}$ );
9)     end;
10)  end;
11) end;

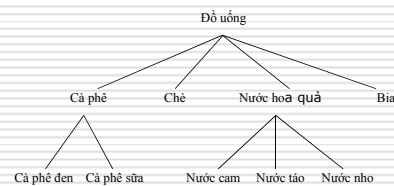
```

Phát hiện luật kết hợp trong cơ sở dữ liệu

23

Vấn đề khai phá luật kết hợp đa mức

Phân cấp khái niệm trên các mục của CSDL



Phát hiện luật kết hợp trong cơ sở dữ liệu

24

Thuật toán cơ bản khai phá luật kết hợp đa mức

```

 $L_1 := \{\text{các 1-tập mục thường xuyên}\};$ 
 $k := 2;$ 
while ( $L_{k-1} \neq \emptyset$ ) do
  begin
     $C_k := \text{các ứng cử viên mới kích thước } k \text{ được sinh ra từ } L_{k-1}$ 
    forall giao dịch  $t \in D$  do
      begin
        Thêm tất cả các tổ tiên của từng mục trong  $t$  vào  $t$ , loại bỏ sự trùng lặp
        Tăng bộ đếm của tất cả các ứng viên trong  $C_k$  mà có mặt trong  $t$ 
      end
       $L_k := \text{Tất cả ứng viên trong } C_k \text{ đạt độ hỗ trợ tối thiểu}$ 
       $k := k + 1;$ 
    end
  end
  Câu trả lời :=  $\bigcup_k L_k$ 

```

Phát hiện luật kết hợp trong cơ sở dữ liệu

25

3. Phát hiện luật kết hợp trong CSDL quan hệ

□ CSDL quan hệ: các quan hệ thường chứa các thuộc tính định lượng, phạm trù.

■ Xử lý những thuộc tính định lượng:

□ Phân vùng rõ

→ Khai phá luật kết hợp định lượng

□ Phân vùng mờ

→ Khai phá luật kết hợp mờ

Phát hiện luật kết hợp trong cơ sở dữ liệu

26

Khai phá luật kết hợp định lượng

□ Phân vùng Equi-Depth: Các vùng có kích thước như nhau

■ Phân vùng dựa trên các giá trị có thể có của thuộc tính. Ví dụ: nếu kiểu thuộc tính có giá trị từ 1 đến 15 và depth $d=3$ thì sinh ra các khoảng [1,3], [4,6], [7,9], [10,12], [13,15]

■ Phân vùng dựa trên các giá trị có thực trong CSDL: d giá trị đầu được đặt vào khoảng thứ nhất, d giá trị tiếp theo được đặt vào khoảng thứ hai,...

Phát hiện luật kết hợp trong cơ sở dữ liệu

27

Khai phá luật kết hợp định lượng

□ Phân vùng dựa trên khoảng cách: có xem xét tính chất định lượng và ngữ nghĩa của dữ liệu. Khoảng cách giữa các điểm dữ liệu càng nhỏ thì chúng càng nên thuộc về 1 nhóm

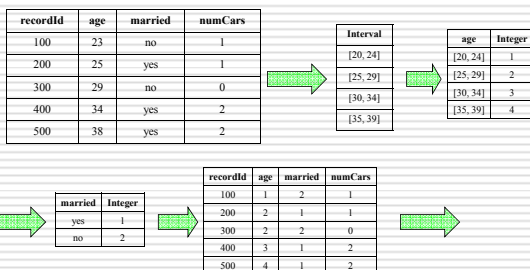
Lương	equi-depth	distance-based
18	[18, 30]	[18, 18]
30		[30, 31]
31	[31, 80]	
80		[80, 82]
81		
82		

Phát hiện luật kết hợp trong cơ sở dữ liệu

28

Khai phá luật kết hợp định lượng

□ Các bước:



Phát hiện luật kết hợp trong cơ sở dữ liệu

29

Khai phá luật kết hợp định lượng

recordId	<age, [20,29]>	<age, [30,39]>	<married, yes>	<married, no>	<numCars, 0>	<numCars, 1>	<numCars, 2>
100	1	0	0	1	0	1	0
200	1	0	1	0	0	1	0
300	1	0	0	1	1	0	0
400	0	1	1	0	0	0	1
500	0	1	1	0	0	0	1

Itemset	Support
{<age, [20, 29]>}	3
{<age, [30, 39]>}	2
{<married, yes>}	3
{<married, no>}	2
{<numCars, [0, 1]>}	3
{<age, [30, 39]>, <married, yes>}	2

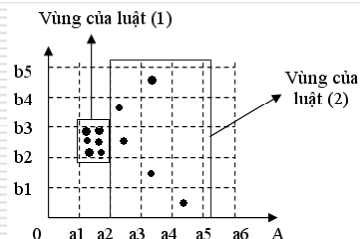
Rule	Support	Confidence
{<age, [30, 39]>, <married, yes>} → {<numCars, [2, 2]>}	0.40	1.00
{<age, [20, 29]>} → {<numCars, [0, 1]>}	0.60	0.67

Phát hiện luật kết hợp trong cơ sở dữ liệu

30

Khai phá luật kết hợp định lượng

□ Cách tiếp cận khối dày đặc



Phát hiện luật kết hợp trong cơ sở dữ liệu

31

Khai phá luật kết hợp mờ

□ Khái niệm luật kết hợp mờ

- Nếu $X = \{x_1, x_2, \dots, x_p\}$ là $A = \{a_1, a_2, \dots, a_p\}$ thì $Y = \{y_1, y_2, \dots, y_q\}$ là $B = \{b_1, b_2, \dots, b_q\}$
- X, Y là tập các thuộc tính
- $x_1, x_2, \dots, y_1, y_2, \dots$ là các thuộc tính
- A, B là tập các tập mờ
- $a_1, a_2, \dots, b_1, b_2, \dots$ là các tập mờ

□ Công thức tính độ hỗ trợ mờ

$$FS_{\langle X, A \rangle} = \frac{\sum_{t_i \in D} \prod_{x_j \in X} d_{x_j}(a_j, t_i, x_j)}{|D|}$$

Phát hiện luật kết hợp trong cơ sở dữ liệu

32

Khai phá luật kết hợp mờ

□ Công thức tính độ tin cậy mờ

$$FC_{\langle \langle X, A \rangle, \langle Y, B \rangle \rangle} = \frac{FS_{\langle Z, C \rangle}}{FS_{\langle X, A \rangle}} = \frac{\sum_{t_i \in D} \prod_{z_j \in Z} d_{z_j}(c_j, t_i, z_j)}{\sum_{t_i \in D} \prod_{x_j \in X} d_{x_j}(a_j, t_i, x_j)}$$

- Ví dụ: Có $X = \{\text{Balance, Income}\}$, $A = \{\text{medium, high}\}$, $Y = \{\text{Credit}\}$, $B = \{\text{high}\}$

Balance, medium	Credit, high	Income, high
0.5	0.6	0.4
0.8	0.9	0.4
0.7	0.8	0.7
0.9	0.8	0.3
0.9	0.7	0.6

$$FS_{\langle X, A \rangle} = 0.364$$

$$FC_{\langle \langle X, A \rangle, \langle Y, B \rangle \rangle} = 0.766$$

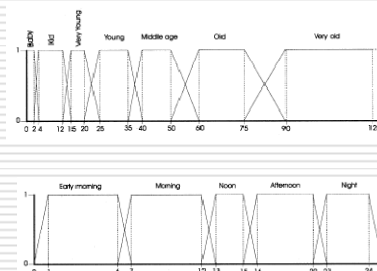
Phát hiện luật kết hợp trong cơ sở dữ liệu

33

Khai phá luật kết hợp mờ

□ Các bước

	Age	Hour
t_1	60	20:15
t_2	80	23:45
t_3	22	15:30
t_4	55	01:00
t_5	3	19:30
t_6	18	06:51



Phát hiện luật kết hợp trong cơ sở dữ liệu

34

Khai phá luật kết hợp mờ

	<Age, Baby>	<Age, Kid>	<Age, Very young>	<Age, Young>	<Age, Middle age>	<Age, Old>	<Age, Very old>	<Hour, Early Morning>	<Hour, Morning>	<Hour, Noon>	<Hour, After Noon>	<Hour, Night>
t_1	0	0	0	0	0	1	0	0	0	0	0.75	0.25
t_2	0	0	0	0	0	0.67	0.33	0	0	0	0	1
t_3	0	0	0.6	0.4	0	0	0	0	0	0.5	0.5	0
t_4	0	0	0	0	0.5	0.5	0	1	0	0	0	0
t_5	0.5	0.5	0	0	0	0	0	0	0	0	1	0
t_6	0	0	1	0	0	0	0	0.85	0.15	0	0	0

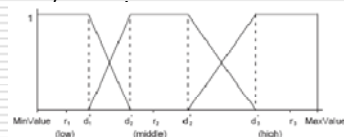
Phát hiện luật kết hợp trong cơ sở dữ liệu

35

Khai phá luật kết hợp mờ

- Phân vùng mờ miền thuộc tính?
- Attila Gyenesi giới thiệu kỹ thuật phân vùng mờ dựa trên chỉ số độ tốt (Goodness Index)

□ Tìm tâm, các cận các nhóm



□ Tính hàm độ thuộc

Phát hiện luật kết hợp trong cơ sở dữ liệu

36

4. Một số vấn đề khác

- Phát hiện luật có yếu tố thời gian
- Phát hiện luật trên nhiều quan hệ
- Phân loại luật kết hợp



Lời hay ý đẹp

Thành công, đó là cách khuyến khích ta cố gắng làm những việc lớn lao hơn nữa. Thất bại, đó là cách cổ vũ ta làm lại việc đã làm với nhiều hi vọng hơn.

Gabriel Palau