

KHAI THÁC TẬP PHỔ BIẾN (Frequent Itemset Mining)

DATA MINING

HCMUS - 2024



Nội dung

1. Các khái niệm
2. Khai thác tập phổ biến
3. Khai thác tập phổ biến đóng
4. Khai thác tập phổ biến tối đại
5. Nhận xét

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

1

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

2

1. Các khái niệm

- ❖ Hạng mục (*item*): Cho I là một tập các thuộc tính nhị phân. Cho $I = \{I_1, I_2, \dots, I_m\}$, mỗi I_m là một item.
- ❖ Tập hạng mục (*itemset*): Một tập $X \subseteq I$ là một tập các hạng mục.
- ❖ Một CSDL giao tác là một tập gồm nhiều *itemset*, mỗi *itemset* là một giao tác được định danh bởi một giá trị duy nhất là mã giao tác (*tid*).

TID	Items
1	A, C, T, W
2	C, D, W
3	A, C, T, W
4	A, C, D, W
5	A, C, D, T, W
6	C, D, T

Cho CSDL giao tác D như sau.

- ❖ Độ hỗ trợ (*support*) của tập hạng mục X trong cơ sở dữ liệu D , $sup(X)$, là phần trăm số giao tác trong D có chứa X .
- ❖ Ví dụ:
 - $sup(A) = 4/6 * 100 = 66.67\%$
 - $sup(ACD) = 2/6 * 100 = 33.3\%$

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

3

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

4

1.1 Tập phổ biến

Cho một tập hạng mục X và cơ sở dữ liệu D .

- ❖ Tập X là phổ biến trong D nếu $sup(X) \geq minsup$, với $minsup$ là ngưỡng hỗ trợ tối thiểu do người dùng đặt.
- ❖ Ví dụ: $minsup = 70\%$
 $sup(A) = 66.67\% < minsup$
 $sup(C) = 100\% > minsup$
 - A không là tập phổ biến.
 - C là tập phổ biến.

TID	Items
1	A, C, T, W
2	C, D, W
3	A, C, T, W
4	A, C, D, W
5	A, C, D, T, W
6	C, D, T

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

5

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

6

1.2 Tập phổ biến đóng

Cho $I = \{i_1, i_2, \dots, i_m\}$ - là tập các items

Cho $T = \{t_1, t_2, \dots, t_m\}$ - là tập các giao tác.

❖ Kết nối Galois

Cho quan hệ hai ngôi $\delta \subseteq I \times T$ chứa CSDL cần khai thác.

Với: $X \subseteq I$ và $Y \subseteq T$, ta định nghĩa hai ánh xạ giữa $P(I)$ và $P(T)$ như sau:

- $t: P(I) \rightarrow P(T), t(X) = \{y \in T | \forall x \in X, x\delta y\}$
- $i: P(T) \rightarrow P(I), i(Y) = \{x \in I | \forall y \in Y, x\delta y\}$

1.2 Tập phổ biến đóng

Ánh xạ (1): $\iota(X)$ lấy tất cả tid của giao tác có chứa tập hạng mục X .

Ánh xạ (2): $i(Y)$ lấy tất cả item tồn tại trong tất cả giao tác Y .

❖ Toán tử đóng: $c = i \circ \iota$

❖ Tập hạng mục X là tập đóng nếu $c(X) = X$.

⇒ **Tập phổ biến đóng: là tập hạng mục đóng thỏa ngưỡng $minsup$ cho trước.**

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

7

1.2 Tập phổ biến đóng

❖ Ví dụ: Cho cơ sở dữ liệu D với $minsup = 30\%$.

Kiểm tra AW , CD có phải là tập phổ biến đóng?

TID	Items
1	A, C, T, W
2	C, D, W
3	A, C, T, W
4	A, C, D, W
5	A, C, D, T, W
6	C, D, T

Sử dụng toán tử đóng:

$c(AW) = i(\iota(AW)) = i(\{1345\}) = ACW$

$c(CD) = i(\iota(CD)) = i(\{2456\}) = CD$

Vậy CD là tập phổ biến đóng, AW **không** là tập phổ biến đóng.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

8

1.2 Tập phổ biến đóng

❖ Tóm tắt định nghĩa: Tập phổ biến đóng là tập phổ biến mà không có tập nào bao nó có cùng độ phổ biến.

- Với F là tập hợp gồm tất cả tập phổ biến.

$$F = \{X \mid X \subseteq I \text{ và } sup(X) \geq minsup\}$$

- Gọi C là tập hợp gồm tất cả tập phổ biến đóng.

$$\Rightarrow C = \{X \mid X \in F \text{ và } \nexists Y \supset X \text{ mà } sup(X) = sup(Y)\}$$

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

9

1.3 Tập phổ biến tối đại

❖ Định nghĩa: Tập phổ biến tối đại là tập phổ biến mà không có tập nào bao nó là phổ biến.

$$M = \{X \mid X \in F \text{ và } \nexists Y \supset X \text{ mà } Y \in F\}$$

❖ Ví dụ: Cho 3 tập phổ biến $\{A,B\}$, $\{A,C\}$, $\{A,B,D\}$

- $\{A,C\}$ và $\{A,B,D\}$ là **tập phổ biến tối đại**.

- $\{A,B\}$ **không** phải là tập phổ biến tối đại. Do $\{A,B\}$ là tập con của $\{A,B,D\}$.

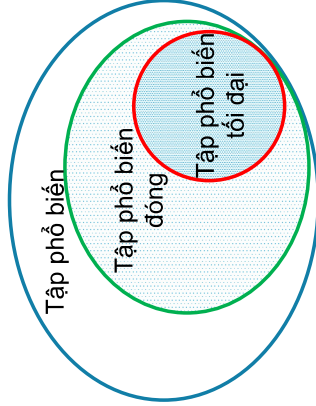
3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

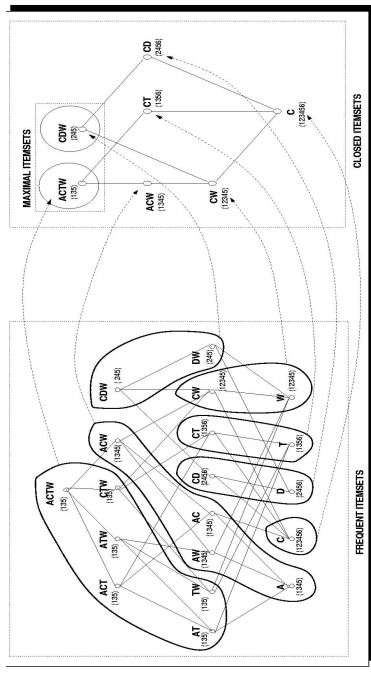
10

1.4 So sánh tập phổ biến

❖ Số lượng tập phổ biến phát sinh trong quá trình khai thác.



Tập phổ biến, đóng và tối đại



3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

11

BỘ MÔN KHOA HỌC MÁY TÍNH

12

2. Khai thác tập phổ biến

- **Input:** Tập các giao dịch T , với tập itemsets I
- **Output:** Tất cả các itemsets chứa trong I thỏa:
 - $\text{support} \geq \text{minsup}$

• Tham số:

- $N = |T|$: số lượng giao dịch
- $d = |I|$: số lượng itemsets riêng biệt.
- w : số lượng tối đa items của 1 giao dịch.
- Có bao nhiêu itemsets có thể có ?

• Quy mô của vấn đề:

- WalMart bán 100,000 mặt hàng và có thể lưu trữ hàng tỉ giỏ hàng.
- The Web có hàng tỉ từ và hàng tỉ trang

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

13

2. Khai thác tập phổ biến

• Quy tắc Apriori :

- Nếu một tập là phổ biến, thì tất cả tập con của nó phải phổ biến.
- Nếu 1 tập không phổ biến thì tất cả tập chứa nó không phổ biến.

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

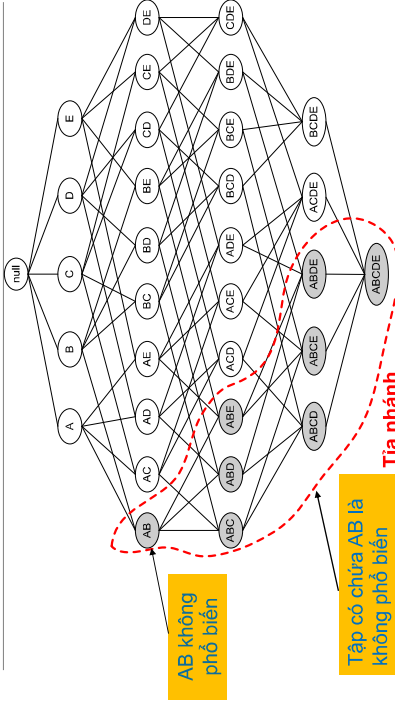
- Độ hỗ trợ của 1 tập không bao giờ vượt quá độ hỗ trợ các tập con của nó.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

15

2. Khai thác tập phổ biến

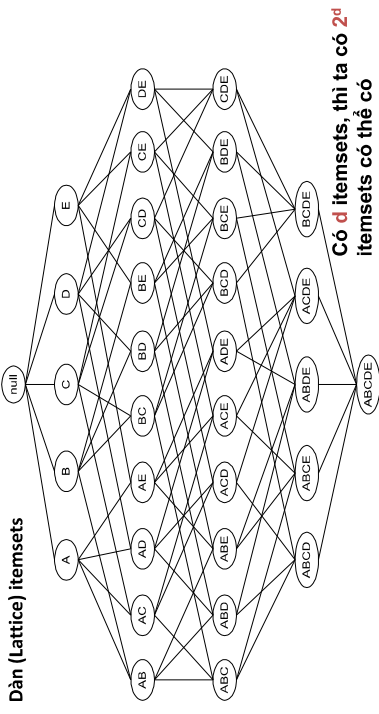


3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

17

2. Khai thác tập phổ biến

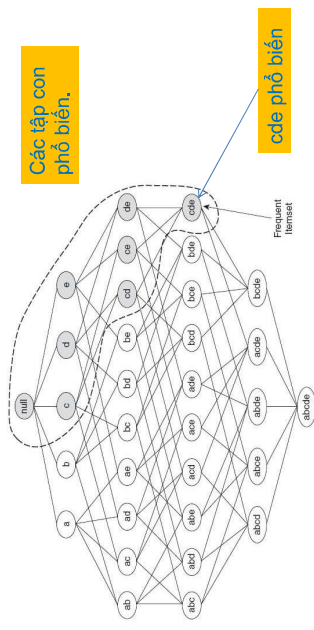


3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

14

2. Khai thác tập phổ biến



3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

16

2. Khai thác tập phổ biến

- ❖ Thuật toán **Apriori** (*state-of-the art*) được đề xuất bởi R. Agrawal và R. Srikant vào năm 1994 để khai thác tập phổ biến.
 - ❖ Gọi C_k là các tập có k hạng mục. Thuật toán thực hiện như sau: $k = 1$. F là tập hợp các tập phổ biến.
 - Bước 1: Đếm độ hỗ trợ của từng tập trong C_k .
 - Bước 2: Phát sinh ứng viên C_{k+1} dựa trên C_k .
 - Bước 3: Loại bỏ các ứng viên C_{k+1} chứa tập con C_k không phổ biến.
 - Bước 4: Thêm các tập C_k thỏa ngưỡng minsup vào F .
- Thuật toán lặp lại đến khi tất cả tập phổ biến được phát sinh.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

18

2.2 Thuật toán Eclat

Đầu vào:

- P , các tập 1-hạng mục cùng tidset.
- $minsup$, ngưỡng support tối thiểu.

Kết quả:

- F , tập các itemset phổ biến.

0. $Eclat([P])$:

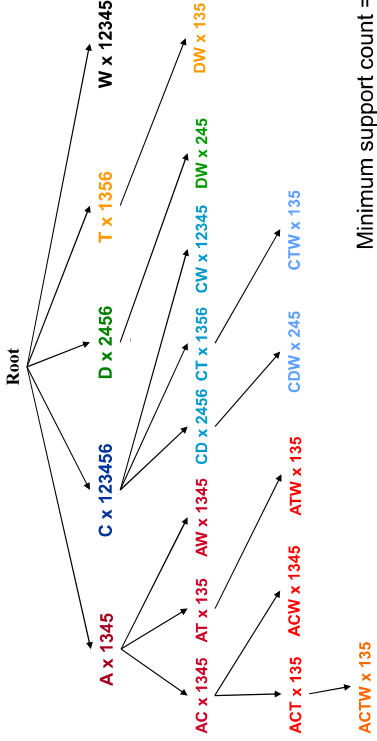
- for all $X_i \in [P]$ do
- $T_i = \emptyset$
- for all $X_j \in [P]$, with $j > i$ do
 $R = X_i \cup X_j$;
- $t(R) = t(X_j) \cap t(X_i)$;
- if $\sigma(R) \geq minsup$ then
 $T_i = T_i \cup \{R\}$; $F_{|R|} = F_{|R|} \cup \{R\}$;
- for all $T_i \neq \emptyset$ do $Eclat(T_i)$;

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

25

2.2 Thuật toán Eclat



Minimum support count = 3

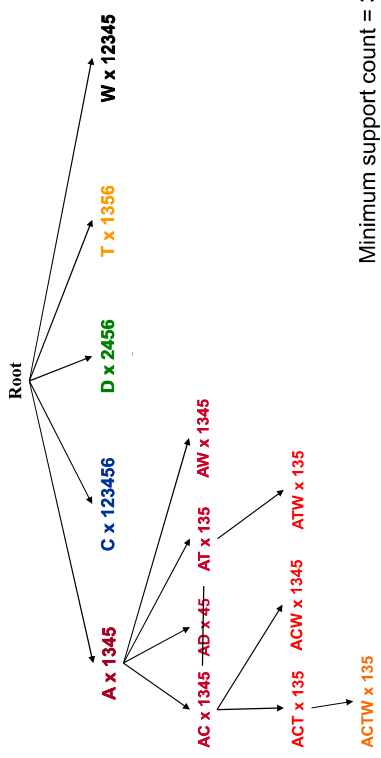
27

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

28

2.2 Thuật toán Eclat



Minimum support count = 3

26

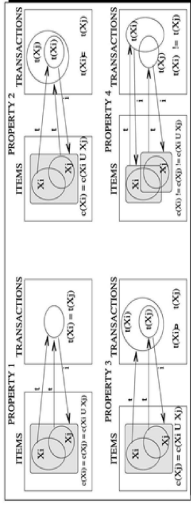
3. Khai thác tập phổ biến đóng

❖ M. J. Zaki cùng đồng sự đề xuất Thuật toán CHARM để khai thác những mẫu phổ biến đóng.

❖ Thuật toán sử dụng *tidset* và *duyet theo chiều sâu trước* tương tự như thuật toán Eclat.

❖ Thuật toán áp dụng một số cải tiến để cắt tỉa bớt các tập hạng mục không phổ biến và tìm tập đóng bằng phương pháp dự trên mối quan hệ của các tập hạng mục.

3. Thuật toán Charm



Định lý 1: Đặt $X_i \times t(X_i)$ và $X_j \times t(X_j)$ là hai thành viên bất kỳ của một lớp $[P]$. Bốn thuộc tính sau là:

- Nếu $t(X_i) = t(X_j)$, thì $c(X_i) = c(X_j) = c(X_i \cup X_j)$.
- Nếu $t(X_i) \subset t(X_j)$, thì $c(X_i) \neq c(X_j)$, nhưng $c(X_i) = c(X_i \cup X_j)$.
- Nếu $t(X_i) \supset t(X_j)$, thì $c(X_i) \neq c(X_j)$, nhưng $c(X_j) = c(X_i \cup X_j)$.
- Nếu $t(X_i) \not\subset t(X_j)$ và $t(X_j) \not\subset t(X_i)$, thì $c(X_i) \neq c(X_j) \neq c(X_i \cup X_j)$.

Đầu vào: CSDL D , $minsup$

Kết quả: tập FCI gồm tất cả các tập phổ biến đóng của CSDL

CHARM (D , $minsup$):

- $[\emptyset] = \{l_i \times t(l_i) : l_i \in I \wedge \sigma(l_i) \geq minsup\}$
 - CHARM-EXTEND ($[\emptyset]$, $C = \emptyset$)
 - return C //tất cả itemset đóng
- CHARM-EXTEND ($[P]$, C):
- for each $l_i \times t(l_i)$ in $[P]$ do
 - $P_i = P \cup l_i$ and $[P_i] = \emptyset$
 - for each $l_j \times t(l_j)$ in $[P_i]$, with $j > i$ do
 $X = l_i$ and $Y = t(l_i) \cap t(l_j)$
 - CHARM-PROPERTY ($X \times Y$, l_i , l_j , P_i , $[P_i]$, $[P]$)
 - SUBSUMPTION-CHECK (C , P_i)
 - CHARM-EXTEND ($[P_i]$, C)
 - delete $[P_i]$

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

29

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

30

3. Thuật toán Charm

CHARM-PROPERTY ($X \times Y, X_i, X_j, P_i, [P_i], [P_i]$)

12: if ($\sigma(X) \geq \text{minsup}$) then

13: if ($i(X_i) = i(X_j)$) then (1)

14: remove X_j from $[P_i]$

15: $P_i = P_i \cup X_j$

16: else if ($i(X_i) \subset i(X_j)$) then (2)

17: $P_i = P_i \cup X_j$

18: else if ($i(X_i) \supset i(X_j)$) then (3)

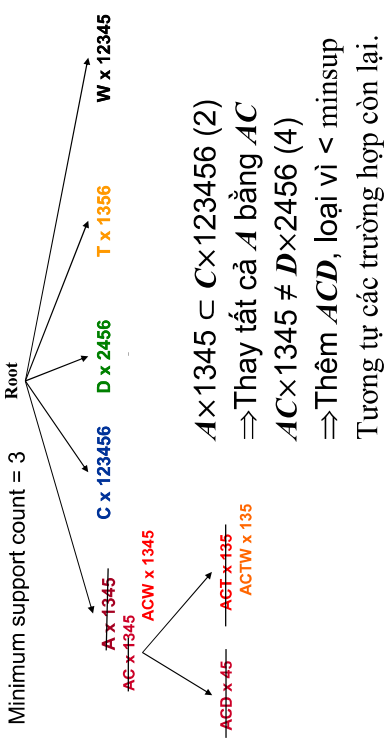
19: remove X_j from $[P_i]$

20: Add $X \times Y$ to $[P_i]$

21: else if ($i(X_i) \neq i(X_j)$) then (4)

22: Add $X \times Y$ to $[P_i]$

3. Thuật toán Charm



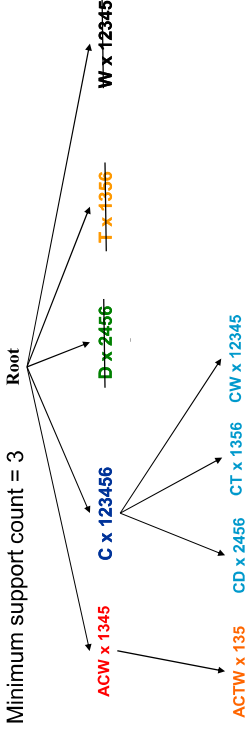
3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

31

32

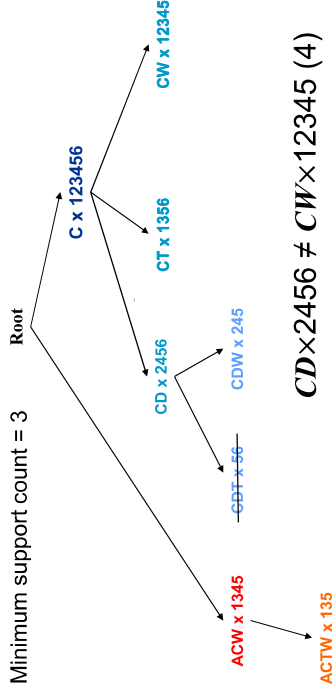
3. Thuật toán Charm



33

34

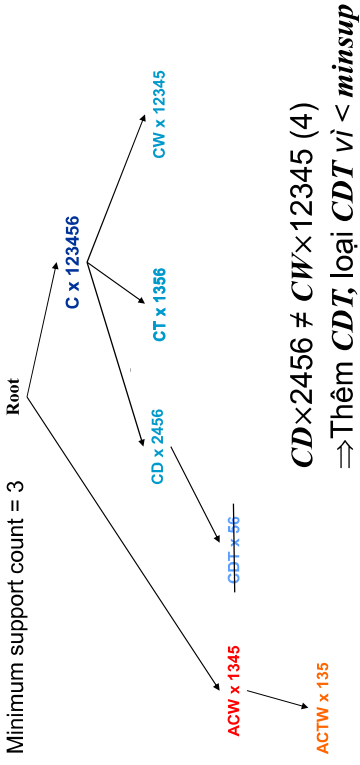
3. Thuật toán Charm



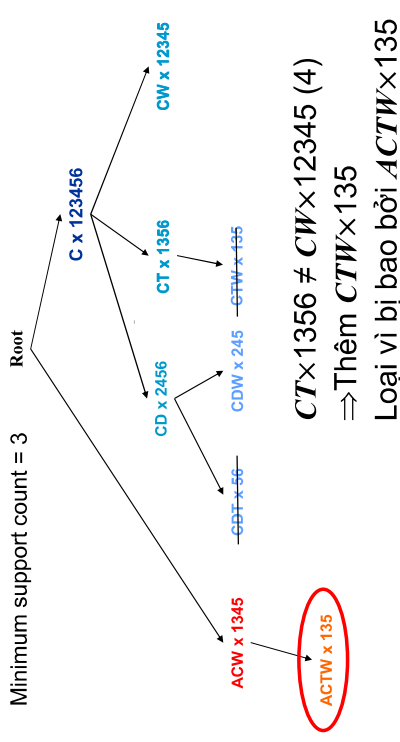
35

36

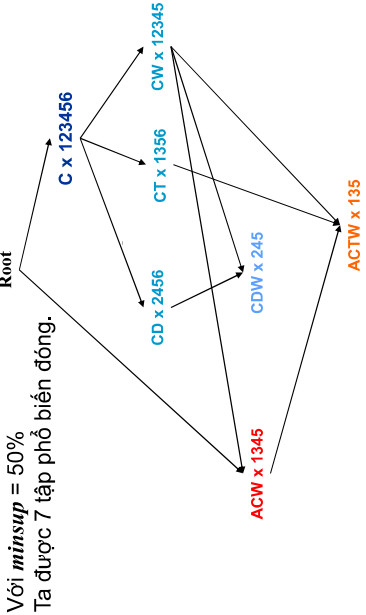
3. Thuật toán Charm



3. Thuật toán Charm



3. Thuật toán Charm



37

BỘ MÔN KHOA HỌC MÁY TÍNH

3/18/2024

38

4. Thuật toán GenMax

- ❖ K. Gouda và M. J.Zaki đề xuất Thuật toán GenMax để tìm tập phổ biến tối đại dựa trên tiến trình backtrack.
- ❖ Thuật toán cũng sử dụng tidset, và cách duyệt cây tương tự như Eclat.
- ❖ Từng hạng mục sẽ được lấy ra những hạng mục có thể kết hợp với nó (tập kết hợp thỏa *minsup*).

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

39

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

40

4. Thuật toán GenMax

Thuật toán FI-backtrack

Đầu vào:

- I_ℓ tập các itemsets có độ dài ℓ .
- C_ℓ tập những items có thể kết hợp với I_ℓ .
- ℓ là độ dài của itemset.

Kết quả: itemset phổ biến

FI-backtrack (I_ℓ, C_ℓ, ℓ)

- 1: **for each** $x \in C_\ell$
- 2: $I_{\ell+1} = I_\ell \cup \{x\}$ //đồng thời thêm $I_{\ell+1}$ vào **FI**
- 3: $P_{\ell+1} = \{y: y \in C_\ell \text{ and } y > x\}$
- 4: $C_{\ell+1} = \text{FI-combine}(I_{\ell+1}, P_{\ell+1})$
- 5: **FI-backtrack** ($I_{\ell+1}, C_{\ell+1}, \ell+1$)

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

41

4. Khai thác tập phổ biến tối đại

- ❖ K. Gouda và M. J.Zaki đề xuất Thuật toán GenMax để tìm tập phổ biến tối đại dựa trên tiến trình backtrack.
- ❖ Thuật toán cũng sử dụng tidset, và cách duyệt cây tương tự như Eclat.
- ❖ Từng hạng mục sẽ được lấy ra những hạng mục khả kết hợp với nó (tập kết hợp thỏa *minsup*).

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

38

4. Thuật toán GenMax

Mục tiêu của tiến trình backtrack là:

- ❖ Lấy ra những tập khả kết hợp với tập hạng mục đang xét.
- ❖ Kết hợp tập hạng mục với các tập khả kết hợp với nó để tạo tập $k+1$ -hạng mục tiếp theo.
- ❖ Thực hiện đệ quy đến khi tất cả tập hạng mục phổ biến được rút trích.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

40

4. Thuật toán GenMax

Hàm FI-combine: dùng kết hợp các hạng mục lại với nhau.

FI-combine ($I_{\ell+1}, P_{\ell+1}$)

- 1: $C = \emptyset$
- 2: **for each** $y \in P_{\ell+1}$
- 3: **if** $I_{\ell+1} \cup \{y\}$ là phổ biến
- 4: $C = C \cup \{y\}$ //sắp xếp lại C nếu cần
- 5: **return** C ;

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

42

4. Thuật toán GenMax

- Để tìm tập phổ biến tối đại, chỉ cần áp dụng điều kiện loại bỏ đi những tập phổ biến không tối đại.

```

MFI-backtrack ( $I_t, C_t, I$ )
1: for each  $x \in C_t$ 
2:    $I_{t+1} = I_t \cup \{x\}$ 
3:    $P_{t+1} = \{y; y \in C_t \text{ and } y > x\}$ 
4:   if  $I_{t+1} \cup P_{t+1}$  có tập bao nó trong MFI
5:     return //tất cả nhánh con bị cắt tĩa
6:    $C_{t+1} = \text{FI-combine}(I_{t+1}, P_{t+1})$ 
7:   if  $C_{t+1}$  is empty
8:     if  $I_{t+1}$  không có tập nào bao nó trong MFI
9:        $\text{MFI} = \text{MFI} \cup I_{t+1}$ 
10:    else MFI-backtrack ( $I_{t+1}, C_{t+1}, I+1$ )

```

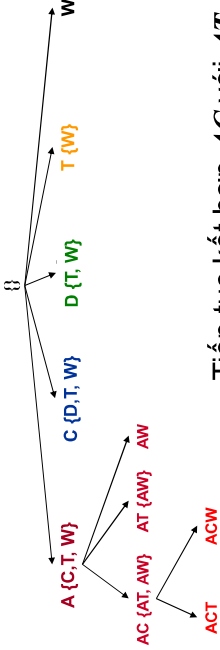
3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

43

4. Thuật toán GenMax

- Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $\text{minsup} = 50\%$

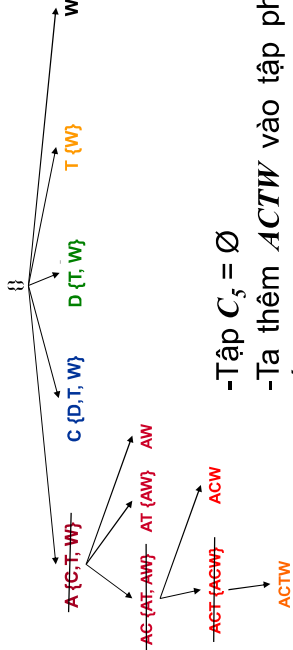


- Tiếp tục kết hợp AC với AT, AW.
- Tập $C_3 = \{ACT, ACW\} \neq \emptyset$
- Gọi đệ quy hàm backtrack.

45

4. Thuật toán GenMax

- Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $\text{minsup} = 50\%$

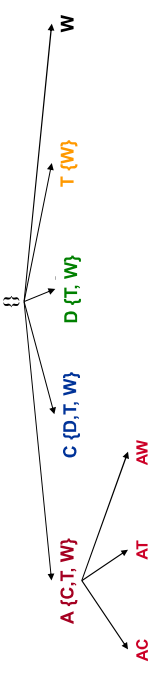


- Tập $C_5 = \emptyset$
- Ta thêm ACTW vào tập phổ biến tối đại.

47

4. Thuật toán GenMax

- Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $\text{minsup} = 50\%$

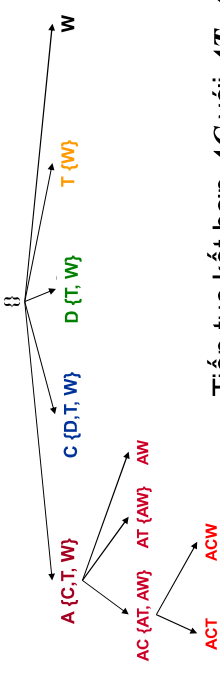


- Đầu tiên kết hợp A lần lượt C, T, W.
- Tập $C_2 = \{AC, AT, AW\} \neq \emptyset$
- Gọi đệ quy hàm backtrack.

44

4. Thuật toán GenMax

- Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $\text{minsup} = 50\%$

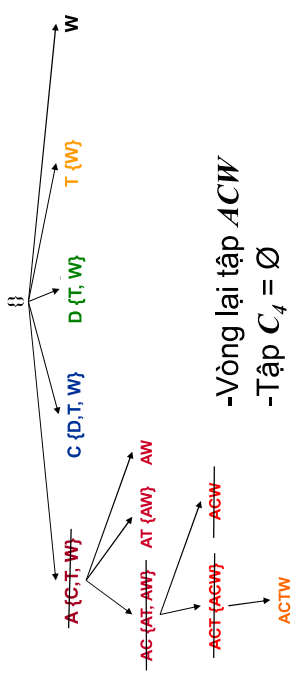


- Tiếp tục kết hợp AC với AT, AW.
- Tập $C_3 = \{ACT, ACW\} \neq \emptyset$
- Gọi đệ quy hàm backtrack.

46

4. Thuật toán GenMax

- Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $\text{minsup} = 50\%$

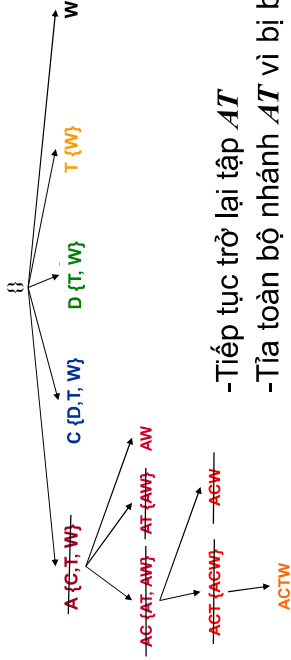


- Vòng lại tập ACW
- Tập $C_4 = \emptyset$
- Loại ACW vì bị bao bởi ACTW

48

4. Thuật toán GenMax

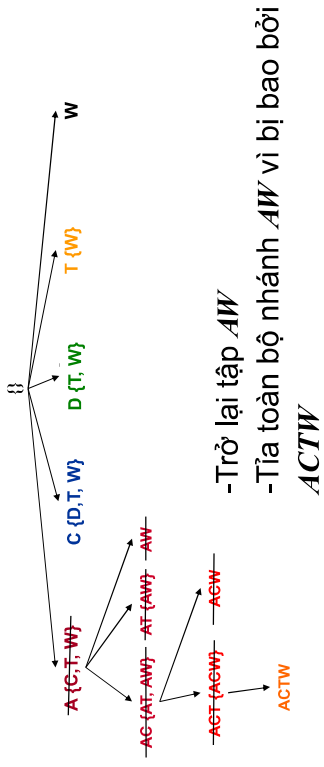
Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $minsup = 50\%$



49

4. Thuật toán GenMax

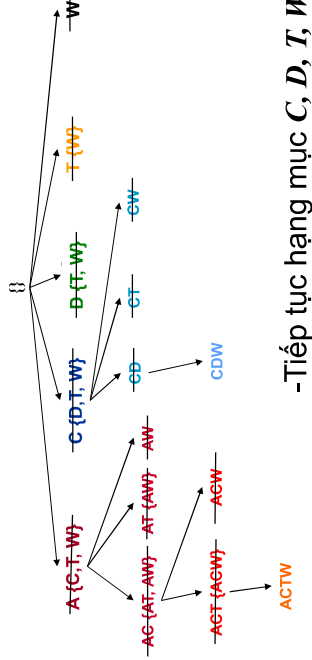
Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $minsup = 50\%$



50

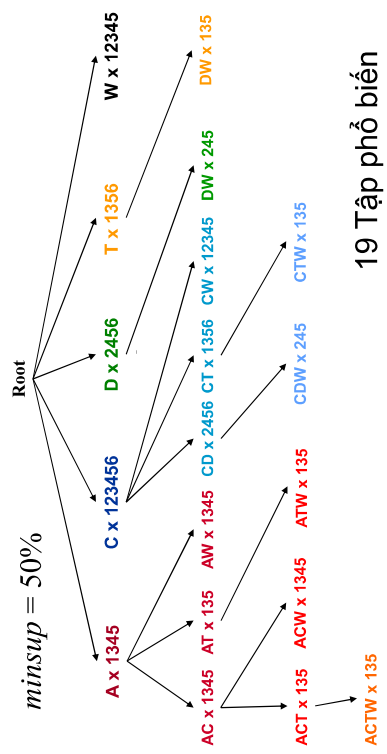
4. Thuật toán GenMax

Ví dụ: Đầu tiên ta có các tập khả kết hợp với hạng mục A với $minsup = 50\%$



51

5. Nhận xét

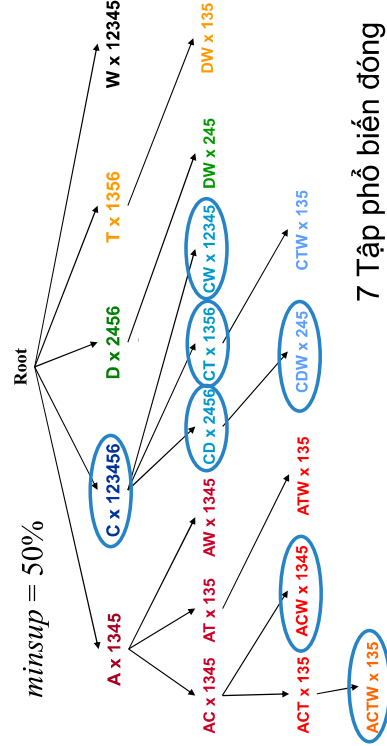


3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

52

5. Nhận xét

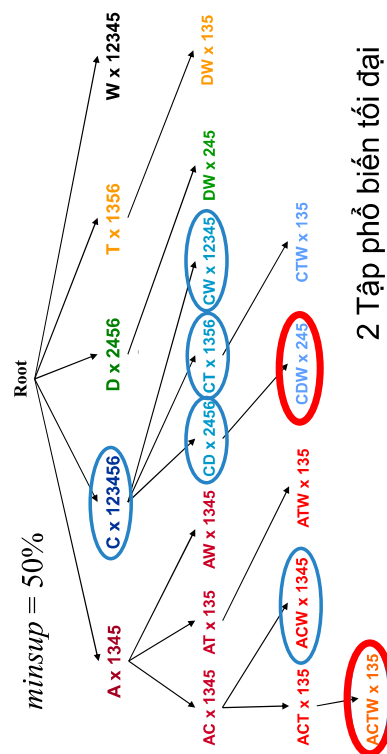


3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

53

5. Nhận xét



3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

54

5. Nhận xét

- ❖ Mọi quan hệ giữa các tập phổ biến như sau:
 $M \subseteq C \subseteq F$.
- ❖ Tập phổ biến đóng thể hiện đầy đủ thông tin của tất cả các tập phổ biến cùng với độ hỗ trợ chính xác của nó.
- ❖ Luật kết hợp rút trích từ tập phổ biến đóng sẽ nhỏ gọn hơn, dễ quản lý, phân tích.
- ❖ Khai thác tập phổ biến tối đại thích hợp với CSDL dây đặc, khi mà số lượng tập đóng cũng có thể rất lớn.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

55

Tài liệu tham khảo

- [1] M. J. Zaki, **Closed Itemset Mining And Non-redundant Association Rule Mining**, Computer Science Department, Rensselaer Polytechnic Institute.
- [2] M. J. Zaki, **Scalable Algorithms for Association Mining**, IEEE Transactions on Knowledge and Data Engineering, 12(3), May/Jun 2000, pp. 372-390.
- [3] M. J. Zaki and K. Gouda, **GenMax: An Efficient Algorithm for Mining Maximal Frequent Itemsets**, Data Mining and Knowledge Discovery: An International Journal, 11(3), 2005, pp .223-242.

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

56

Thanks for your listening !!
Q & A

3/18/2024

BỘ MÔN KHOA HỌC MÁY TÍNH

57