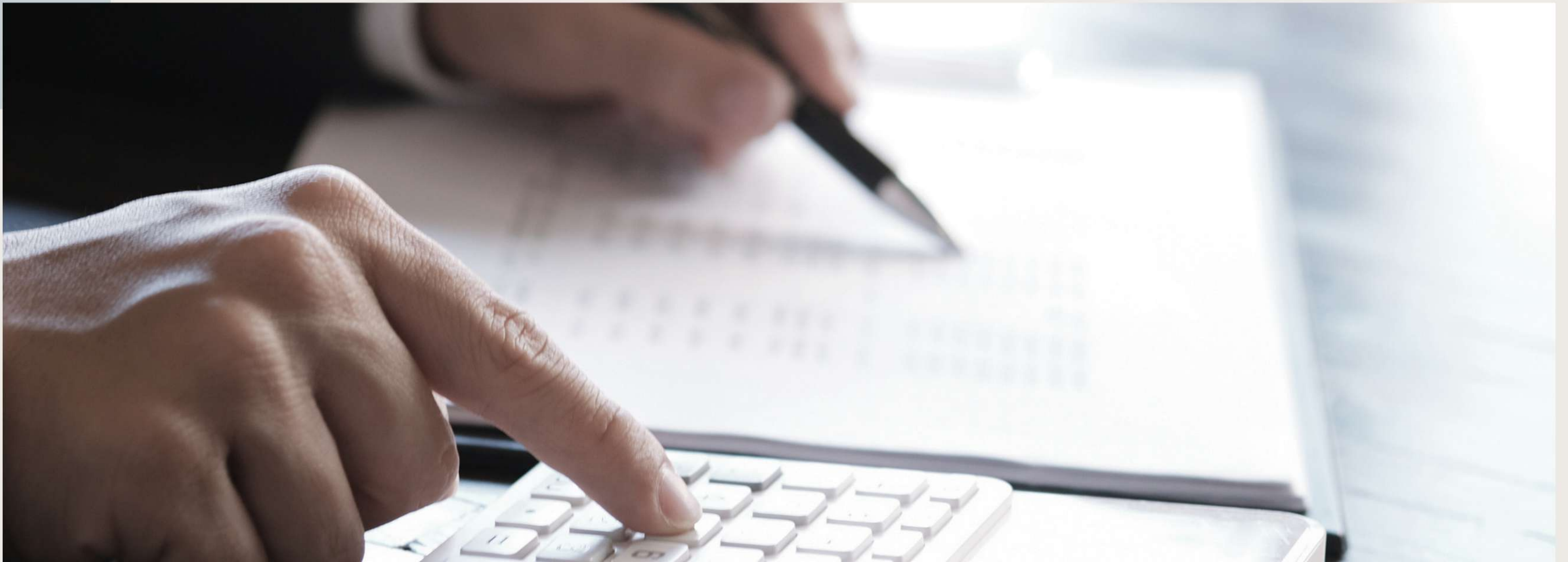


DATA MINING

GROUP 8



NỘI DUNG



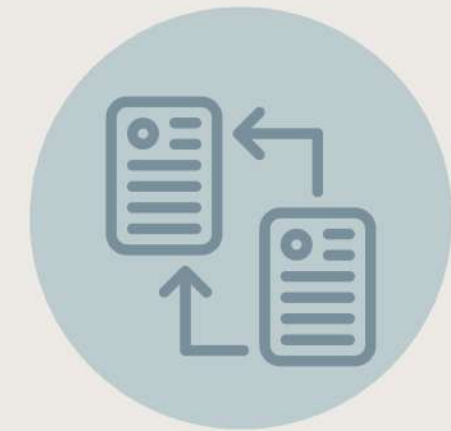
I

**ASSOCIATION
RULE**



II

**APRIORI
ALGORITHM**

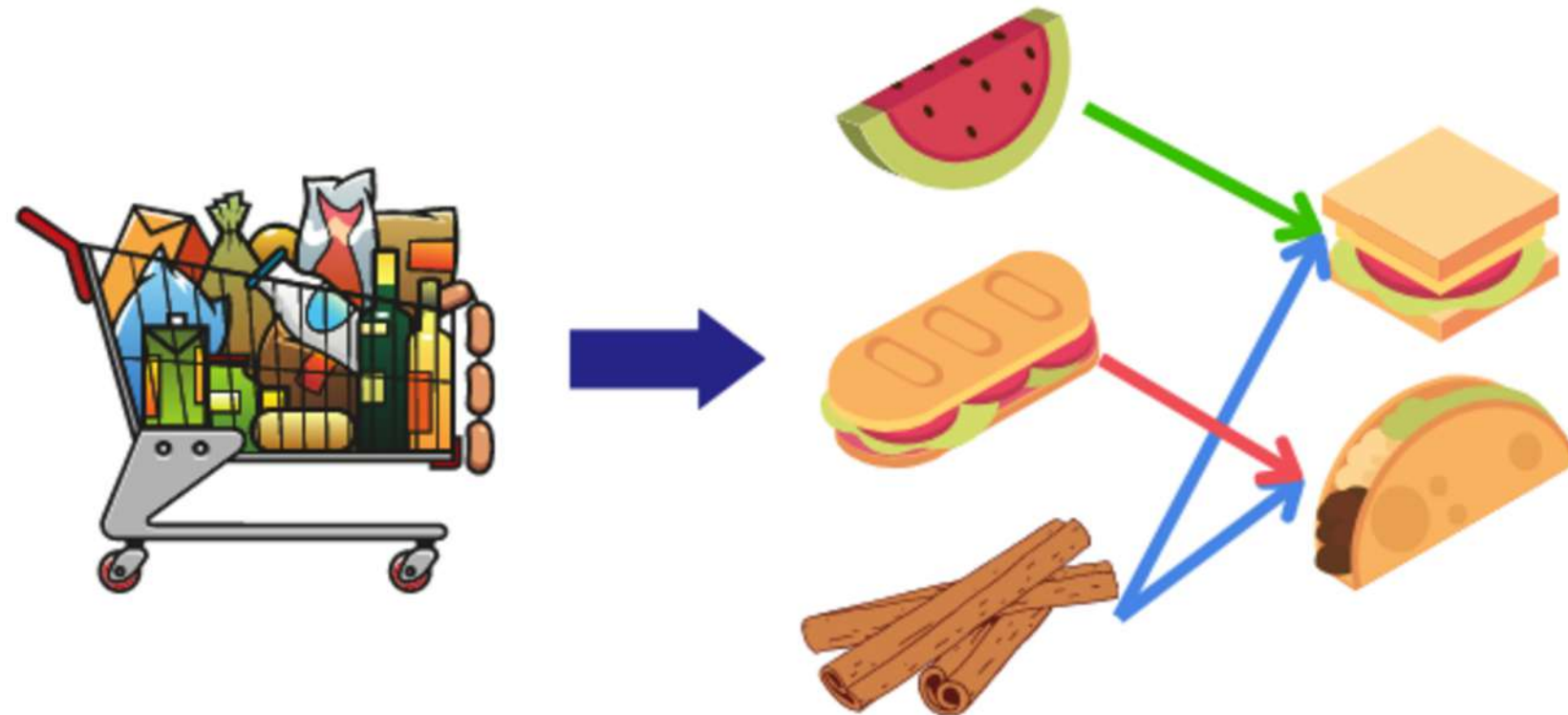


III

**FP-GROWTH
ALGORITHM**

I - ASSOCIATION RULE

Association Rule Learning



*"93% of people who purchased item A
also purchased item B"*



1. TẬP PHỔ BIẾN

Khái niệm: Tập phổ biến là một khái niệm quan trọng trong lĩnh vực khai phá dữ liệu. Nó là một tập hợp các mẫu (items) xuất hiện thường xuyên trong một tập dữ liệu.

Mục đích: Tìm ra các quy luật kết hợp giữa các mục (items) trong tập dữ liệu. Các quy luật kết hợp này có thể được sử dụng để phân tích hành vi tiêu dùng, tìm kiếm các mối liên hệ giữa các sản phẩm hoặc dịch vụ, đưa ra các khuyến nghị sản phẩm, và cải thiện hiệu quả kinh doanh.

MỘT SỐ KHÁI NIỆM

- **Cơ sở dữ liệu giao dịch:** lưu trữ các giao dịch được thực hiện trong một khoảng thời gian nhất định, trong đó mỗi giao dịch là một danh sách các mục (items) mà khách hàng đã mua hoặc sử dụng.
- **Hạng mục (item):** là một đối tượng hoặc một thực thể riêng biệt trong một danh sách hoặc một tập hợp có thể là sản phẩm, mặt hàng (sữa, trứng, bánh mì,..., dịch vụ, hoặc bất kỳ đối tượng nào khác trong tập dữ liệu được phân tích.
- **Tập hạng mục (itemset):** có thể là một tập hợp các hạng mục đơn hoặc một tập hợp các hạng mục kết hợp trong tập dữ liệu. (các hạng mục trong giỏ hàng)

BÀI TOÁN KHAI THÁC TẬP PHỔ BIẾN

- Là bài toán tìm tất cả các hạng mục (itemset) X có độ phổ biến (support) thỏa mãn độ phổ biến tối thiểu (minsupport)

$$\text{Sup}(X) \geq \text{minsup}$$

ĐỘ PHỔ BIẾN?

Cho tập $I = \{i_1, i_2, \dots, i_n\}$

- Hạng mục (item): i_1, i_2, \dots, i_n
- Tập hạng mục (itemset): Tập $X \subseteq I$

Cho tập $D = \{T_1, T_2, \dots, T_n\}$ T_i :

- Một giao dịch (transaction)
- D : CSDL giao dịch
- $|D|$: Số giao dịch trong D

Độ phổ biến: $\text{sup}(X) = C(X) / |D|$

$X = \{A, C\}$

$\text{sup}(X) = 3/5$ (60%)

- $\text{sup}(\{A\}) = 5/5$
- $\text{sup}(\{C\}) = 3/5$

$I = \{A, B, C, D, E\}$

Ex. tập hạng mục: $X = \{A, E, D\}$

$|D| = 5$

$X = \{A, C\}$

$\text{sup}(X) = 3/5$ (60%)



(1) $\text{Minsup} = 60\%$

(2) $\text{Minsup} = 70\%$

T_1	$\{A, B, C, D\}$
T_2	$\{A, C, E\}$
T_3	$\{A, E\}$
T_4	$\{A, B, E\}$
T_5	$\{A, B, C, D, E\}$

T_1	$\{A, B, C, D\}$
T_2	$\{A, C, E\}$
T_3	$\{A, E\}$
T_4	$\{A, B, E\}$
T_5	$\{A, B, C, D, E\}$

$\text{Minsup} = 60\%$

Tính chất tập phổ biến

- Tập con của tập phổ biến là tập phổ biến.
- Tập bao của tập KHÔNG phổ biến là tập KHÔNG phổ biến.
- Các tính chất của độ phổ biến trong data mining bao gồm:
tính mô tả, tính ứng dụng, tính tương quan

VÍ DỤ

T ₁	{A, B, C, D}
T ₂	{A, C, E}
T ₃	{A, E}
T ₄	{A, B, E}
T ₅	{A, B, C, D, E}

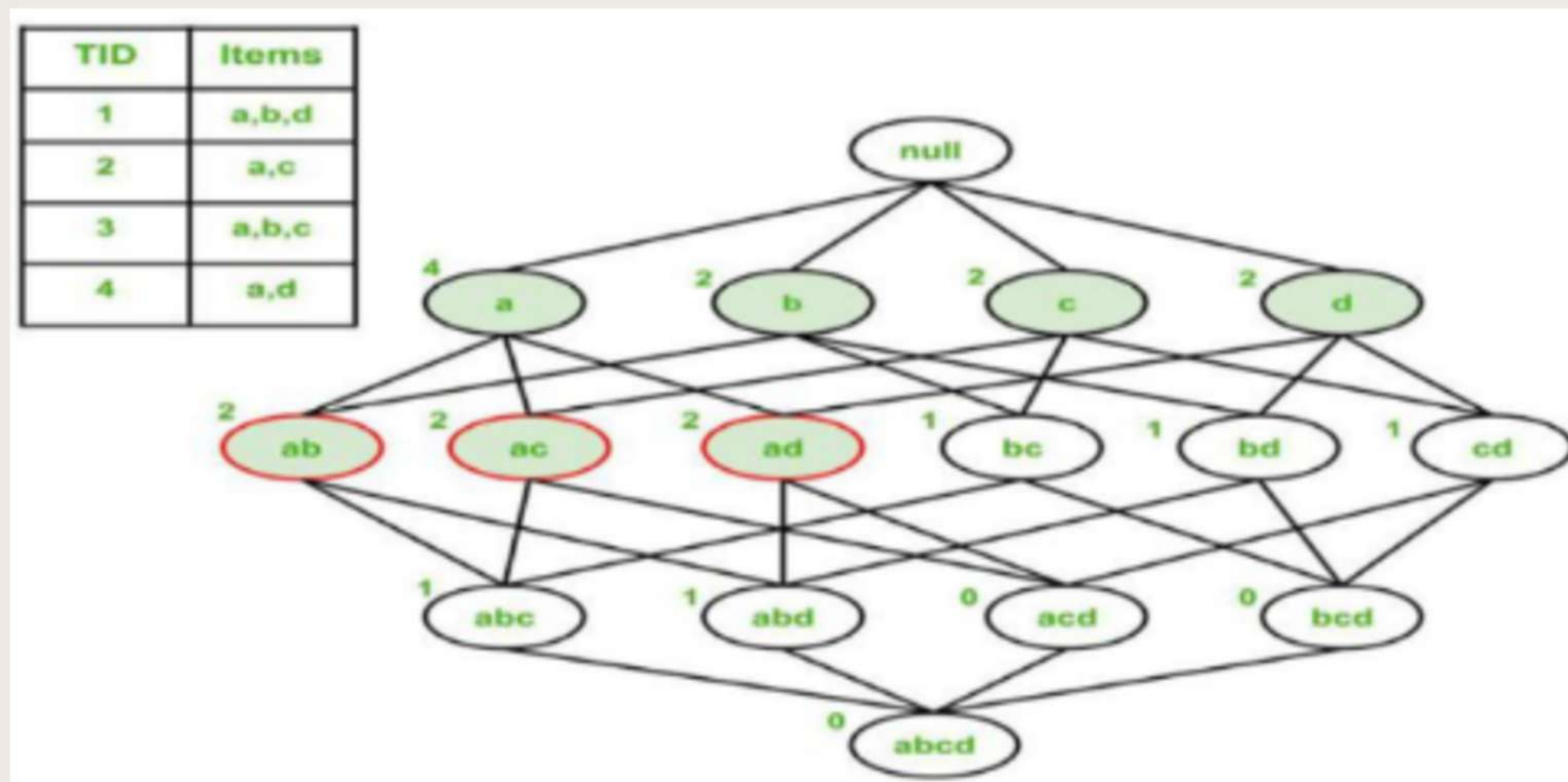
Minsup=60%

Nếu tập {A,D} KHÔNG phổ biến, thì các tập chứa nó như {A,B,C,D} cũng KHÔNG phổ biến ($\text{sup}(\{D\})=2/5 < \text{minsup}$: tập không phổ biến).

=> {A, D}: không phải là tập phổ biến.

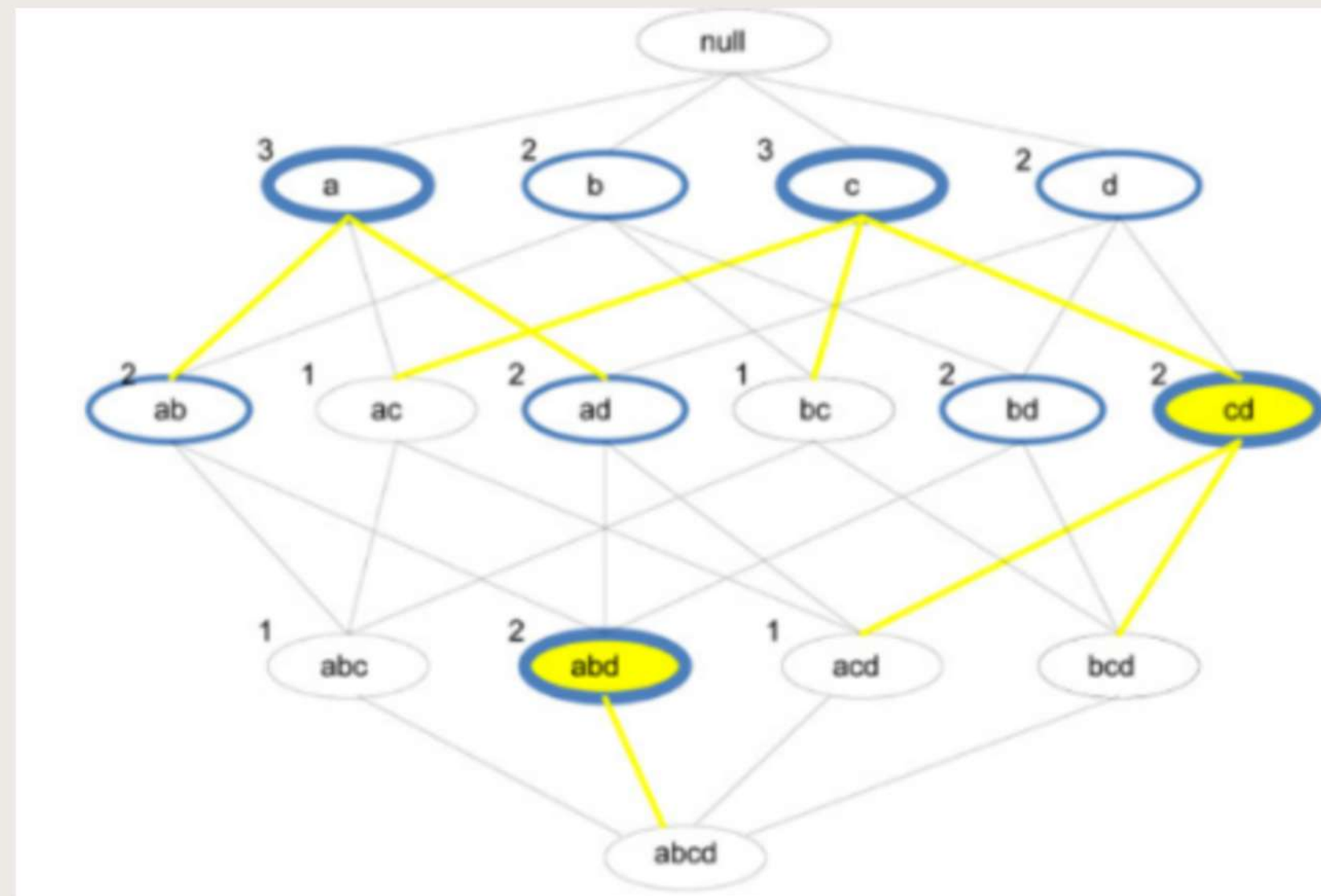
Tập phổ biến tối đại (Maximal frequent itemset)

Là tập phổ biến và KHÔNG tồn tại tập bao nó là tập phổ biến



Tập bao phổ biến (Closed frequent itemset)

Là tập phổ biến và không tồn tại tập nào bao nó có cùng độ phổ biến.



SO SÁNH CÁC TẬP

TẬP PHỔ BIẾN (PATTERN BASE)	TẬP PHỔ BIẾN TỐI ĐA (MAXIMAL PATTERN BASE)	TẬP BAO PHỔ BIẾN (CLOSED PATTERN BASE)
<p>Tập phổ biến là tập hợp các mục (items) xuất hiện trong Luật kết hợp của một tập dữ liệu.</p> <p>VD: Nếu tập dữ liệu có các giao dịch mua hàng, tập PB có thể chứa các mục như {sữa, bánh mì}, {sữa, trái cây}, {sữa, thịt}.</p>	<p>Tập phổ biến tối đa là một tập phổ biến mà không thể mở rộng thêm bằng cách thêm bất kỳ mục nào khác.</p> <p>VD: Nếu tập PB gốc chứa các tập {A, B} và {A, B, C}, thì chỉ tập {A, B} được coi là tối đa.</p>	<p>Tập bao phổ biến là một tập phổ biến mà không có tập PB khác chứa nó nhưng có cùng hỗ trợ (support).</p> <p>VD: Nếu tập PB gốc chứa các tập {A, B}, {A, B, C}, và {A, B, D}, thì chỉ tập {A, B} được coi là tập bao PB.</p>



2. LUẬT KẾT HỢP

Khái niệm: Luật kết hợp là mối quan hệ giữa các tập thuộc tính trong cơ sở dữ liệu. Luật kết hợp là phương tiện hữu ích để khám phá các mối liên kết trong dữ liệu. Có 2 độ đo quan trọng đối với luật kết hợp: Độ hỗ trợ (support) và độ tin cậy (confidence).

ĐỘ PHỔ BIẾN

$$\text{sup}(X \rightarrow Y) = P(X \cup Y) = \text{sup}(X \cup Y)$$

Độ phổ biến của một luật kết hợp $X \Rightarrow Y$ là tỷ lệ giữa số lượng các bản ghi chứa tập hợp $X \Rightarrow Y$, so với tổng số các bản ghi trong D - Ký hiệu $\text{sup}(X \Rightarrow Y)$.

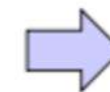
ĐỘ TIN CẬY

$$\text{Conf}(X \rightarrow Y) = P(Y|X) = \text{sup}(X \cup Y) / \text{sup}(X)$$

Độ tin cậy (confidence) của luật $X \rightarrow Y$ trong D , ký hiệu $\text{conf}(X \rightarrow Y, D)$, là tỷ lệ giữa số giao dịch chứa cả X và Y trên số giao dịch chỉ chứa X .

VÍ DỤ

Transaction ID	Items bought
1	Tea, Cake, Cold Drink
2	Tea, Coffee, Cold Drink
3	Eggs, Tea, Cold Drink
4	Cake, Milk, Eggs
5	Cake, Coffee, Cold Drink, Milk, Eggs



$X \rightarrow Y$	sup	conf
Tea \rightarrow Cold Drink	60%	100%
Cake \rightarrow Milk	40%	67%
Tea \rightarrow Milk	0%	0%

QUY TRÌNH KHAI THÁC LUẬT KẾT HỢP

B1: Tìm tất cả các tập PB (theo ngưỡng minsup)

B2: Tạo ra các luật từ các tập PB

- Mỗi tập phổ biến X , tạo các tập con khác rỗng của X

• Đối với mỗi tập con (khác rỗng) A của X , $B=X-A$

$A \rightarrow B$ là LKH nếu $\text{conf}(A \rightarrow B) \geq \text{minconf}$

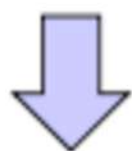
$\text{conf}(A \rightarrow B) = \frac{\text{sup}(A \cup B)}{\text{sup}(A)} = \frac{\text{sup}(X)}{\text{sup}(A)}$

VÍ DỤ

T ₁	{A, B, C, D}
T ₂	{A, C, E}
T ₃	{A, E}
T ₄	{B, C, E}
T ₅	{A, B, C, D, E}

*minsup=60%,
minconf=80%*

minsup=60%,



Frequent itemset	Support (%)
{A}, {B}, {C}, {E}	80, 60, 80, 80
{A, C}, {A, E}, {B, C}, {C, E}	60, 60, 60, 60

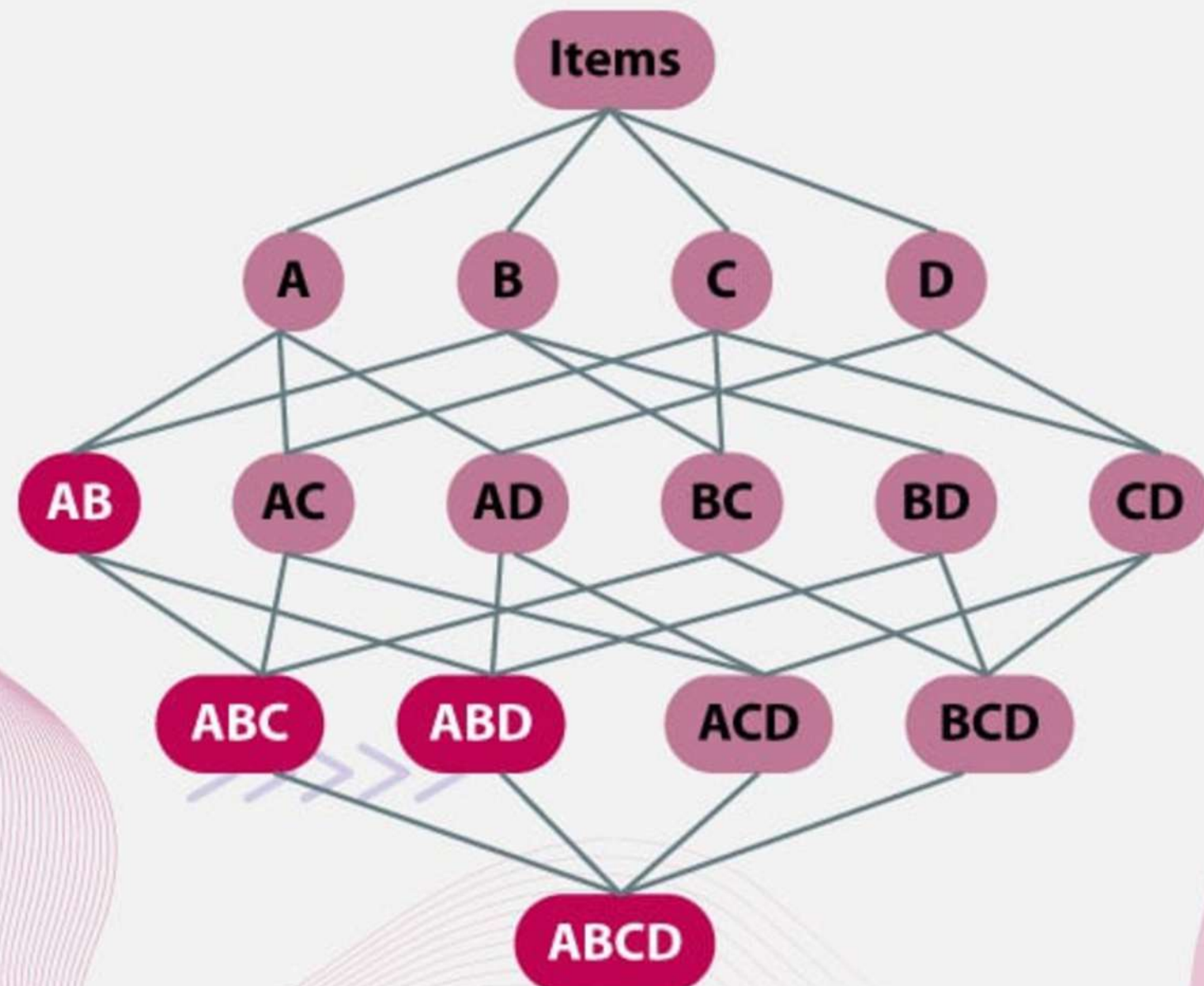
minconf=80%

Luật	conf(X→Y)	Chọn/Loại
A → C	$\text{sup}(\{A\} \cup \{C\}) / \text{sup}(\{A\}) = 3/4 = 75\%$	Loại
C → A	$\text{sup}(\{C\} \cup \{A\}) / \text{sup}(\{C\}) = 3/4 = 75\%$	Loại
A → E	$\text{sup}(\{A\} \cup \{E\}) / \text{sup}(\{A\}) = 3/4 = 75\%$	Loại
E → A	$\text{sup}(\{E\} \cup \{A\}) / \text{sup}(\{E\}) = 3/4 = 75\%$	Loại
B → C	$\text{sup}(\{B\} \cup \{C\}) / \text{sup}(\{B\}) = 3/3 = 100\%$	Chọn
C → B	$\text{sup}(\{C\} \cup \{B\}) / \text{sup}(\{C\}) = 3/4 = 75\%$	Loại
C → E	$\text{sup}(\{C\} \cup \{E\}) / \text{sup}(\{C\}) = 3/4 = 75\%$	Loại
E → C	$\text{sup}(\{E\} \cup \{C\}) / \text{sup}(\{E\}) = 3/4 = 75\%$	Loại

II - APRIORI ALGORITHM

Apriori Algorithm for

Association Rule Mining



Khái niệm: Thuật toán Apriori là một thuật toán được sử dụng trong khai phá dữ liệu để tìm kiếm các tập hợp mục thường xuyên và học các quy tắc kết hợp trong cơ sở dữ liệu quan hệ.

Cách thức hoạt động:

1. Xác định ngưỡng hỗ trợ và ngưỡng tin cậy
2. Tìm tất cả các tập phổ biến 1 hạng mục
3. Tạo các tập ứng viên k+1 hạng mục từ tập phổ biến có k hạng mục
4. Kiểm tra độ phổ biến của các ứng viên
5. Dừng khi không tạo được tập phổ biến hoặc tập ứng viên

T_1	{A, B, C, D}	minsup= 60% ; minconf = 80%
T_2	{A, C, E}	A: Bánh mì
T_3	{B, C, E}	B: Sữa
T_4	{B, C, E}	C: Bơ
T_5	{A, B, C, D, E}	D: Táo
		E: Kem

$minsup=60\%$

T ₁	{A, B, C, D}
T ₂	{A, C, E}
T ₃	{B, C, E}
T ₄	{B, C, E}
T ₅	{A, B, C, D, E}

1st scan

C_1	Itemset	Sup (%)
	{A}	60
	{B}	80
	{C}	100
	{D}	40
	{E}	80

L_1

Itemset	Sup (%)
{A}	60
{B}	80
{C}	100
{E}	80

2nd scan

C_2	Itemset	Sup (%)
	{A, B}	40
	{A, C}	60
	{A, E}	40
	{B, C}	80
	{B, E}	60
	{C, E}	80

L_2

Itemset	Sup (%)
{A, C}	60
{B, C}	80
{B, E}	60
{C, E}	80

3th scan

C_3

Itemset	Sup (%)
{B, C, E}	60

L_3

Itemset	Sup (%)
{B, C, E}	60

T ₁	{A, B, C, D}
T ₂	{A, C, E}
T ₃	{B, C, E}
T ₄	{B, C, E}
T ₅	{A, B, C, D, E}

Frequent Itemset	Sup (%)
{A, C}	60
{B, C}	80
{B, E}	60
{C, E}	80
{B, C, E}	60

Từ các bước vừa rồi, ta được bảng Frequent Itemset từ các mục L2 và L3

KIỀU DOAN

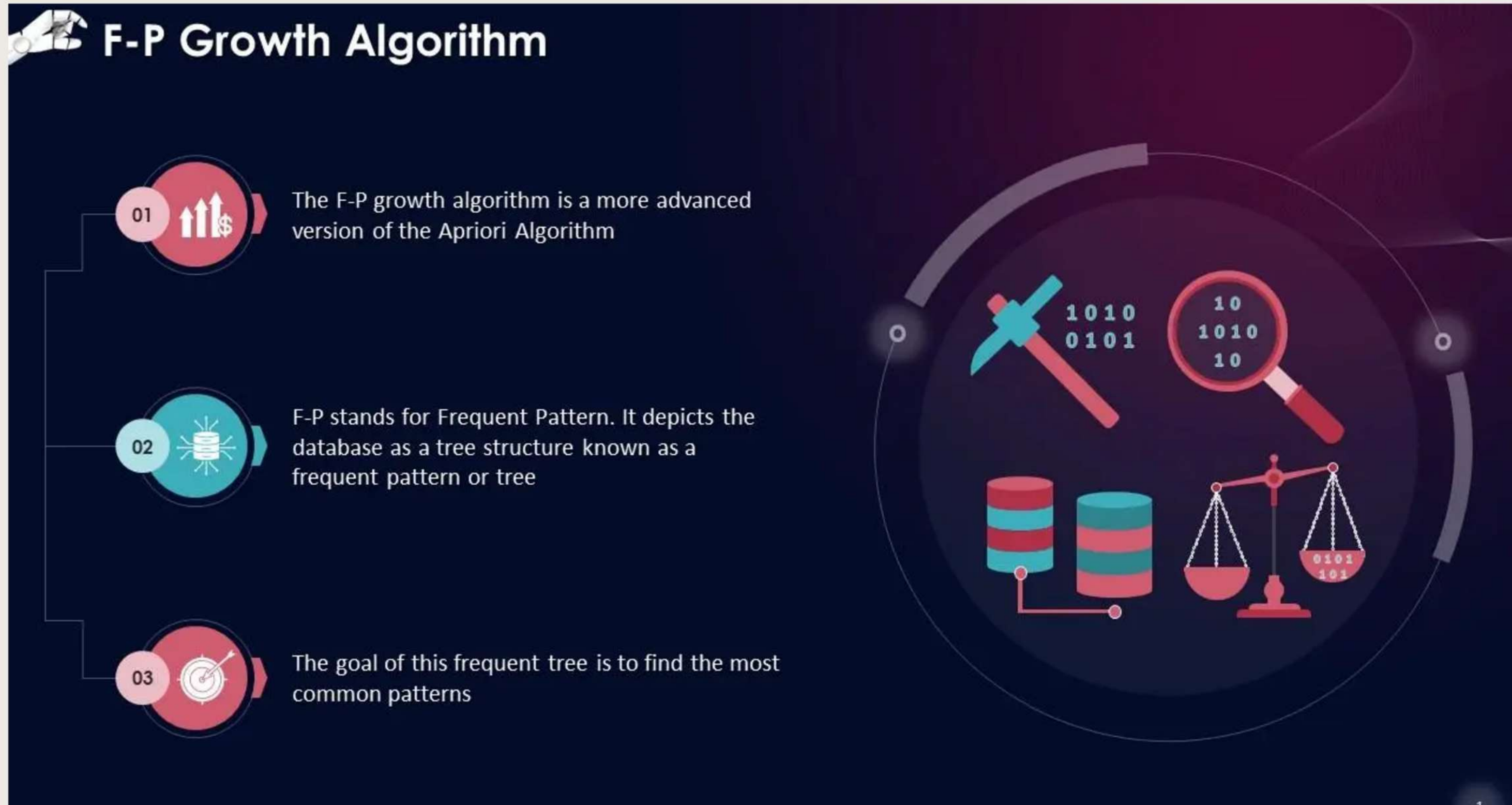
AR	conf (%)
$\{A\} \rightarrow \{C\}$	100
$\{C\} \rightarrow \{A\}$	60
$\{B\} \rightarrow \{C\}$	100
$\{C\} \rightarrow \{B\}$	80
$\{B\} \rightarrow \{E\}$	75
$\{E\} \rightarrow \{B\}$	75
$\{C\} \rightarrow \{E\}$	80
$\{E\} \rightarrow \{C\}$	100
$\{B\} \rightarrow \{C, E\}$	75
$\{C\} \rightarrow \{B, E\}$	60
$\{E\} \rightarrow \{B, C\}$	75
$\{B, C\} \rightarrow \{E\}$	75
$\{B, E\} \rightarrow \{C\}$	100
$\{C, E\} \rightarrow \{B\}$	75

Tính độ tin cậy và so với minconf =80% và sẽ có 6 luật kết hợp như hình.

$\{A\} \rightarrow \{C\}$ là Bánh mì \rightarrow Bơ: Khách hàng mua Bánh mì khả năng cao mua Bơ

$\{B, E\} \rightarrow \{C\}$ là Sữa, Kem \rightarrow Bơ: Khách hàng mua Sữa, Kem khả năng cao mua Bơ

III - FP-GROWTH ALGORITHM



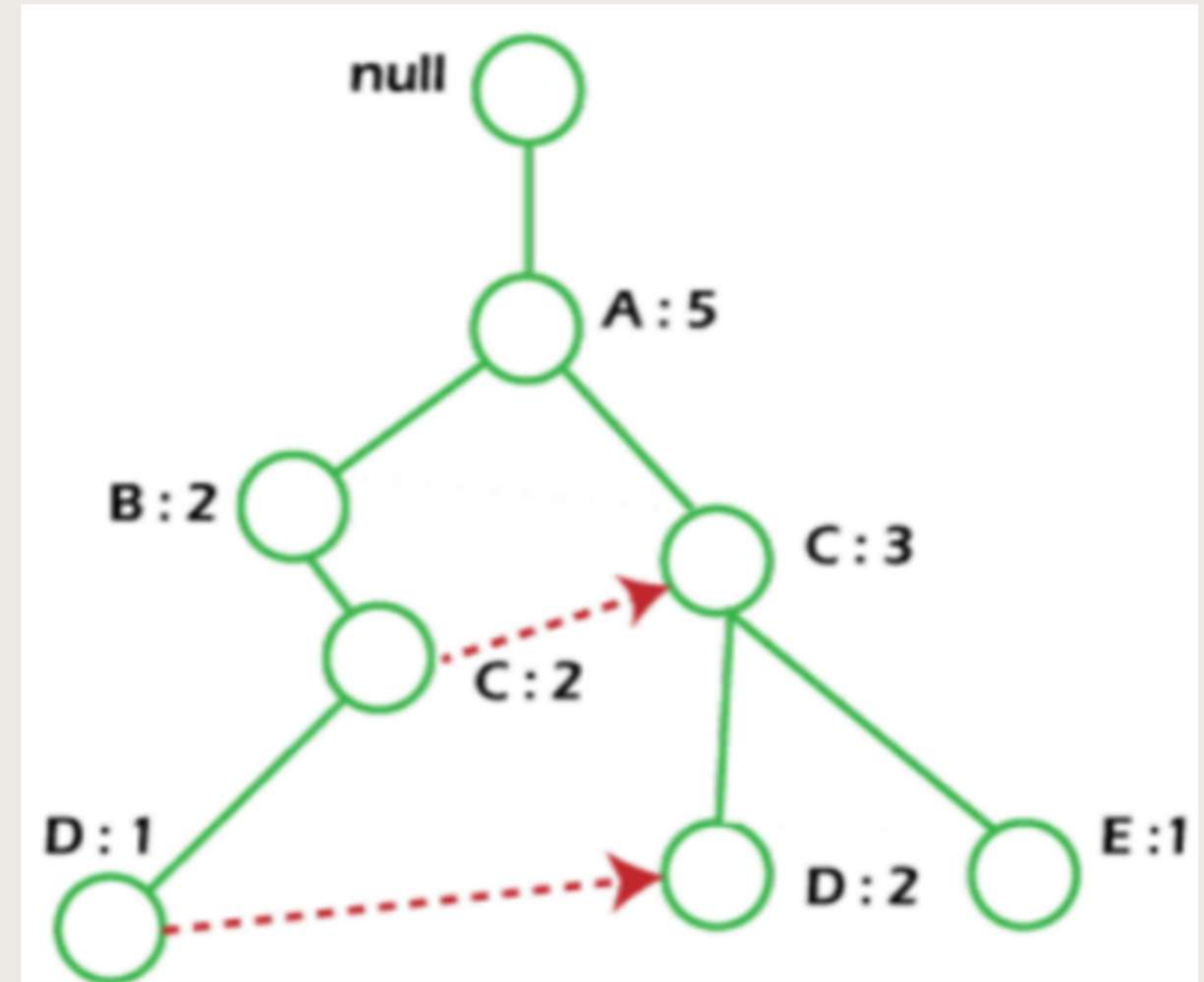
KHÁI NIỆM

- **Thuật toán tăng trưởng FP** là một cách khác để tìm các tập mục phổ biến mà không cần sử dụng các thể hệ ứng viên, do đó cải thiện hiệu suất. Đối với rất nhiều, nó sử dụng chiến lược chia để trị. Cốt lõi của phương pháp này là việc sử dụng một cấu trúc dữ liệu đặc biệt có tên là cây mẫu phổ biến (cây FP), giữ lại thông tin liên kết của tập mục.
- **Các hoạt động của thuật toán:**
- Đầu tiên, nó nén cơ sở dữ liệu đầu vào tạo ra một thể hiện cây FP để biểu diễn các mục thường xuyên.
- Sau bước đầu tiên này, nó chia cơ sở dữ liệu đã nén thành một tập hợp các cơ sở dữ liệu có điều kiện, mỗi cơ sở dữ liệu được liên kết với một mẫu phổ biến.
- Cuối cùng, mỗi cơ sở dữ liệu như vậy được khai thác riêng.

CÂY FP

- **Cây mẫu phổ biến (FP-tree)** là một cấu trúc dữ liệu nhỏ gọn lưu trữ thông tin định lượng về
- **Các mẫu phổ biến trong cơ sở dữ liệu.**
- **Cây mẫu phổ biến được tạo với các bộ mục ban đầu của cơ sở dữ liệu.**

Mục đích: Mục đích của cây FP là khai thác mẫu thường xuyên nhất. Mỗi nút của cây FP đại diện cho một mục của tập mục.



ƯU ĐIỂM

- Cần quét cơ sở dữ liệu hai lần so với Apriori (quét các giao dịch cho mỗi lần lặp lại).
- Việc ghép nối các mục không được thực hiện trong thuật toán này, làm cho nó nhanh hơn.
- Cơ sở dữ liệu được lưu trữ trong một phiên bản nhỏ gọn trong bộ nhớ.
- Nó hiệu quả và có khả năng mở rộng để khai thác cả các mẫu phổ biến dài và ngắn.

NHƯỢC ĐIỂM

- FP Tree công kênh và khó xây dựng hơn Apriori.
- Nó có thể đắt tiền.
- Thuật toán có thể không vừa với bộ nhớ dùng chung khi cơ sở dữ liệu lớn.

THANK YOU

FOR

ATTENTION