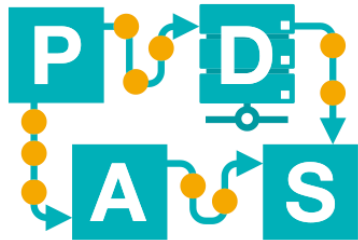


# Association Rules and Sequence Mining

*Instruction 8*

IDS



Chair of Process  
and Data Science

**RWTH**AACHEN  
UNIVERSITY

# Recap

**Association rule mining can be seen as two steps:**

**1. Find all frequent item sets using:**

- Apriori algorithm
- FP-Growth algorithm

**2. Generate strong association rules with the given criteria from the frequent item sets:**

- By definition, this rules must satisfy minimum support and minimum confidence.

# Why FP-Growth algorithm?

## Disadvantages of Apriori algorithm:

- Find candidate sets in an expensive way. If frequent items are large in amount, so the combination would be huge and it would be an expensive operation.

**So Apriori algorithm is a slow algorithm.**

# FP-Growth Algorithm

TID	Item sets
T <sub>1</sub>	{I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> , I <sub>6</sub> }
T <sub>2</sub>	{I <sub>7</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> , I <sub>6</sub> }
T <sub>3</sub>	{I <sub>1</sub> , I <sub>8</sub> , I <sub>4</sub> , I <sub>5</sub> }
T <sub>4</sub>	{I <sub>1</sub> , I <sub>9</sub> , I <sub>10</sub> , I <sub>4</sub> , I <sub>6</sub> }
T <sub>5</sub>	{I <sub>10</sub> , I <sub>2</sub> , I <sub>4</sub> , I <sub>11</sub> , I <sub>5</sub> }

Item	Support count
I <sub>1</sub>	3
I <sub>2</sub>	3
I <sub>3</sub>	2
I <sub>4</sub>	5
I <sub>5</sub>	4
I <sub>6</sub>	3
I <sub>7</sub>	1
I <sub>8</sub>	1
I <sub>9</sub>	1
I <sub>10</sub>	2
I <sub>11</sub>	1

# FP-Growth Algorithm

- **Min-support = 3**

Item	Support count
$l_1$	3
$l_2$	3
$l_4$	5
$l_5$	4
$l_6$	3

Write in  
descending  
order



L

item	Support count
$l_4$	5
$l_5$	4
$l_1$	3
$l_2$	3
$l_6$	3

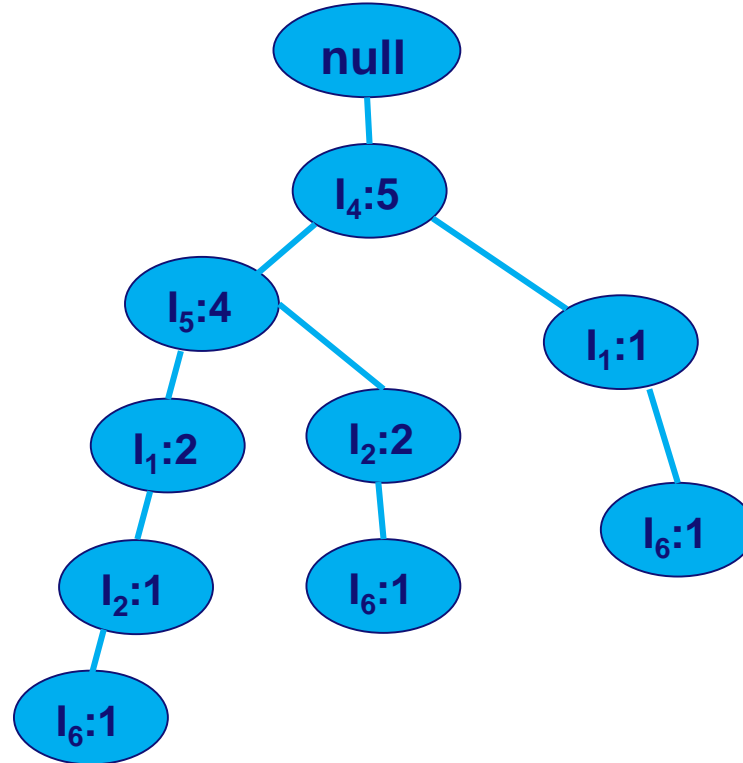
# FP-Growth Algorithm

item	Support count
$l_4$	5
$l_5$	4
$l_1$	3
$l_2$	3
$l_6$	3

TID	Item sets	Ordered item set
$T_1$	$\{l_1, l_2, l_3, l_4, l_5, l_6\}$	$\{l_4, l_5, l_1, l_2, l_6\}$
$T_2$	$\{l_7, l_2, l_3, l_4, l_5, l_6\}$	$\{l_4, l_5, l_2, l_6\}$
$T_3$	$\{l_1, l_8, l_4, l_5\}$	$\{l_4, l_5, l_1\}$
$T_4$	$\{l_1, l_9, l_{10}, l_4, l_6\}$	$\{l_4, l_1, l_6\}$
$T_5$	$\{l_{10}, l_2, l_4, l_{11}, l_5\}$	$\{l_4, l_5, l_2\}$

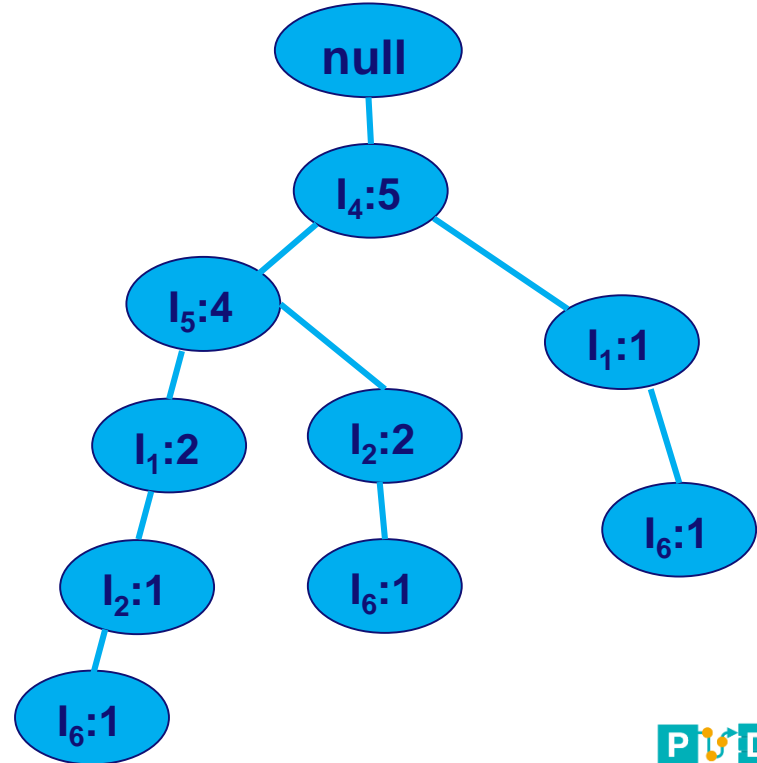
# FP-Growth Algorithm

TID	Ordered item set
T <sub>1</sub>	{I <sub>4</sub> , I <sub>5</sub> , I <sub>1</sub> , I <sub>2</sub> , I <sub>6</sub> }
T <sub>2</sub>	{I <sub>4</sub> , I <sub>5</sub> , I <sub>2</sub> , I <sub>6</sub> }
T <sub>3</sub>	{I <sub>4</sub> , I <sub>5</sub> , I <sub>1</sub> }
T <sub>4</sub>	{I <sub>4</sub> , I <sub>1</sub> , I <sub>6</sub> }
T <sub>5</sub>	{I <sub>4</sub> , I <sub>5</sub> , I <sub>2</sub> }



# FP-Growth Algorithm

Item	Conditional pattern base
$l_6$	$(\{l_4, l_5, l_1, l_2\}:1, \{l_4, l_5, l_2\}:1, \{l_4, l_1\}:1)$
$l_2$	$(\{l_4, l_5, l_1\}:1, \{l_4, l_5\}:2)$
$l_1$	$(\{l_4, l_5\}:2, \{l_4\}:1)$
$l_5$	$(\{l_4\}:4)$
$l_4$	-





# FP-Growth Algorithm

Item	Conditional pattern base	Conditional FP tree	Frequent pattern generated
$l_6$	$(\{l_4, l_5, l_1, l_2\}:1, \{l_4, l_5, l_2\}:1, \{l_4, l_1\}:1)$	$[l_4:3]$	$\langle l_4, l_6: 3 \rangle$
$l_2$	$(\{l_4, l_5, l_1\}:1, \{l_4, l_5\}:2)$	$[l_4, l_5:3]$	$\langle l_4, l_2:3 \rangle \langle l_5, l_2:3 \rangle \langle l_2, l_5, l_4:3 \rangle$
$l_1$	$(\{l_4, l_5\}:2, \{l_4\}:1)$	$[l_4:3]$	$\langle l_4, l_1:3 \rangle$
$l_5$	$(\{l_4\}: 4)$	$[l_4:4]$	$\langle l_4, l_5:4 \rangle$
$l_4$	-	-	-

**Next step: association rule mining... what is confidence and support for  $l_4 \Rightarrow l_6$  and  $l_6 \Rightarrow l_4$  ?**



# Association Rules

- **T** is the set of transactions.
- **I** is the set of all possible item sets composed by items in **T**.
- **A**  $\subseteq$  **I** and **B**  $\subseteq$  **I** are two item sets/sub-item sets from **T**.
- **A**  $\Rightarrow$  **B** is an association rule.

Usually, we would like to discover the association rule **A**  $\Rightarrow$  **B** of which the support and confidence are above certain levels.

# Association Rules

- $\text{support}(A \Rightarrow B) = \text{support}(A \cup B) = \frac{\text{support}_{\text{count}}(A \cup B)}{|T|}$
- $\text{confidence}(A \Rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support}(A)} = \frac{\text{support}_{\text{count}}(A \cup B)}{\text{support}_{\text{count}}(A)}$

- Min\_sup represents minimum support and min\_conf represents minimum confidence.
- $A \Rightarrow B$  is a desired association rule if:

$\text{support}(A \Rightarrow B) \geq \text{min\_sup}$  and  $\text{confidence}(A \Rightarrow B) \geq \text{min\_conf}$

# Association Rules

- Set **min\_sup** to 0.5 and **min\_conf** to 0.7. Is **{bread} => {meat}** from **D** a desired association rule?

Set of transactions D

TID	Set of items
0	bread, meat, wine
1	bread, meat
2	pizza, wine
3	bread, meat, pizza, wine

# Association Rules

- Set **min\_sup** to 0.5 and **min\_conf** to 0.7. Is **{bread} => {meat}** from **D** a desired association rule?

Set of transactions D	
TID	Set of items
0	bread, meat, wine
1	bread, meat
2	pizza, wine
3	bread, meat, pizza, wine

$$\text{support}(\{\text{bread}\} \Rightarrow \{\text{meat}\}) = \frac{\text{support}_{\text{count}}(\{\text{bread, meat}\})}{|D|} = \frac{3}{4} = 0.75 > \text{min\_sup}$$

$$\text{confidence}(\{\text{bread}\} \Rightarrow \{\text{meat}\}) = \frac{\text{support}_{\text{count}}(\{\text{bread, meat}\})}{\text{support}_{\text{count}}(\{\text{bread}\})} = \frac{3}{3} = 1 > \text{min\_conf}$$



**{bread} => {meat}** is a desired association rule

# Association Rules

We use **lift** to evaluate the quality of the discovered association rule **A => B**.

$$\text{lift}(A \Rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support}(A) \cdot \text{support}(B)} = \frac{P(A \cup B)}{P(A) \cdot P(B)}$$

If  $\text{lift}(A \Rightarrow B) \approx 1$  then  $A$  and  $B$  are independent

If  $\text{lift}(A \Rightarrow B) \ll 1$  then  $A$  and  $B$  are negatively correlated

If  $\text{lift}(A \Rightarrow B) \gg 1$  then  $A$  and  $B$  are positively correlated

# Association Rules

Evaluate the quality of the association rule **{bread} => {meat}** by using **lift**:

Set of transactions D

TID	Set of items
0	bread, meat, wine
1	bread, meat
2	pizza, wine
3	bread, meat, pizza, wine

$$\text{lift}(\{bread\} \Rightarrow \{meat\}) = \frac{\text{support}(\{bread, meat\})}{\text{support}(\{bread\}) \cdot \text{support}(\{meat\})} = \frac{(3/4)}{(3/4) \cdot (3/4)} = 1.33$$

# Association Rules

**Exercise 3:** Judge if  $\{A, B\} \Rightarrow \{E\}$ ,  $\{A\} \Rightarrow \{B\}$  and  $\{A\} \Rightarrow \{C\}$  are the desired association rules under minimum support 0.5 and minimum confidence 0.75. Also evaluate the quality of the desired rules.

Example data set S

**TID   Data items**

1	A, B, E
2	C, A, D
3	C, B, D
4	C, A, B, E



# Association Rules

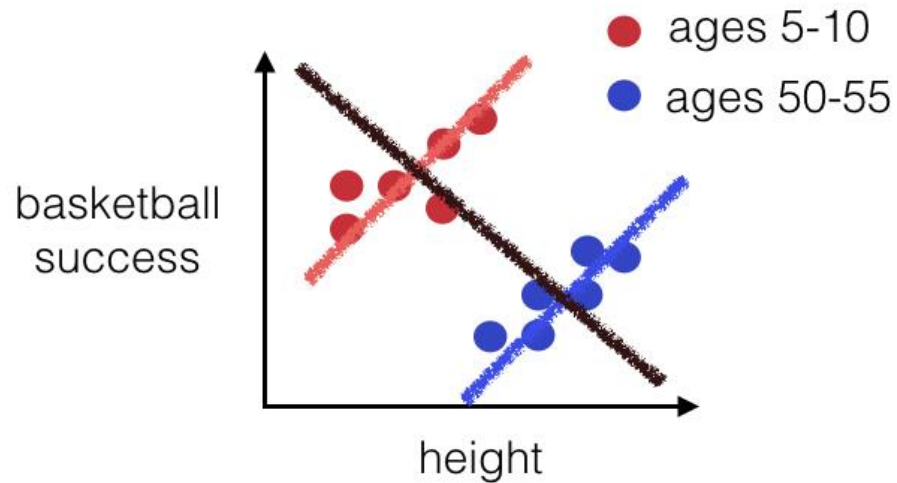
Judge if  $\{A, B\} \Rightarrow \{E\}$ ,  $\{A\} \Rightarrow \{B\}$  and  $\{A\} \Rightarrow \{C\}$ :

- $\text{Support}(\{A, B\} \Rightarrow \{E\}) = 0.5$ ,  $\text{confidence}(\{A, B\} \Rightarrow \{E\}) = 1$ ,  $\text{lift}(\{A, B\} \Rightarrow \{E\}) = 2$ , it is a desired association rule, and lift is larger than 1.
- $\text{Support}(\{A\} \Rightarrow \{B\}) = 0.5$ ,  $\text{confidence}(\{A\} \Rightarrow \{B\}) = 0.67$  it is not a desired association rule
- $\text{Support}(\{A\} \Rightarrow \{C\}) = 0.5$ ,  $\text{confidence}(\{A\} \Rightarrow \{C\}) = 0.67$  it is not a desired association rule

# Simpson's Paradox

	Restaurant 1	Restaurant 2
Males	$50 \backslash 150 = 33.3\%$	$180 \backslash 360 = 50\%$
Females	$200 \backslash 250 = 80\%$	$36 \backslash 40 = 90\%$
General	$250 \backslash 400 = 62.5\%$	$216 \backslash 400 = 54\%$

# Simpson's Paradox



# Sequence Mining



# Support

A central concept in sequence mining is *support* (and *support count*): the frequency of appearance (relative or absolute) of a certain pattern within the database.

# Support

In this database, find the support count of:

(bc)(de)

b(de)

(bc)de

(ac)(bc)

(bc)(ac)

D = [  
<a(bc)d(eb)>,  
<(ac)(bc)de>,  
<(ac)b(cd)>,  
<ab(bc)(cde)>,  
<(bc)(bd)(bde)>,  
<(abc)(ac)(bc)de>,  
<a(bd)c(de)>  
]

# Support: Solutions

In this database, find the support of:

(bc)(de): **2**

b(de)

(bc)de

(ac)(bc)

(bc)(ac)

D = [  
<a(bc)d(eb)>,  
<(ac)(bc)de>,  
<(ac)b(cd)>,  
**<ab(bc)(cde)>**,  
**<(bc)(bd)(bde)>**,  
<(abc)(ac)(bc)de>,  
<a(bd)c(de)>  
]

# Support: Solutions

In this database, find the support of:

(bc)(de): **2**

b(de): **3**

(bc)de

(ac)(bc)

(bc)(ac)

D = [  
<a(bc)d(eb)>,  
<(ac)(bc)de>,  
<(ac)b(cd)>,  
**<ab(bc)(cde)>**,  
**<(bc)(bd)(bde)>**,  
<(abc)(ac)(bc)de>,  
**<a(bd)c(de)>**  
]



# Support: Solutions

In this database, find the support of:

(bc)(de): **2**

b(de): **3**

(bc)de: **4**

(ac)(bc)

(bc)(ac)

D = [  
    <a(bc)d(eb)>,  
    <(ac)(bc)de>,  
    <(ac)b(cd)>,  
    <ab(bc)(cde)>,  
    <(bc)(bd)(bde)>,  
    <(abc)(ac)(bc)de>,  
    <a(bd)c(de)>  
]

# Support: Solutions

In this database, find the support of:

(bc)(de): **2**

b(de): **3**

(bc)de: **4**

(ac)(bc): **2**

(bc)(ac)

D = [  
    <a(bc)d(eb)>,  
    <(ac)(bc)de>,  
    <(ac)b(cd)>,  
    <ab(bc)(cde)>,  
    <(bc)(bd)(bde)>,  
    <(abc)(ac)(bc)de>,  
    <a(bd)c(de)>  
]

# Support: Solutions

In this database, find the support of:

(bc)(de): **2**

b(de): **3**

(bc)de: **4**

(ac)(bc): **2**

(bc)(ac): **1**

D = [  
    <a(bc)d(eb)>,  
    <(ac)(bc)de>,  
    <(ac)b(cd)>,  
    <ab(bc)(cde)>,  
    <(bc)(bd)(bde)>,  
    <(abc)(ac)(bc)de>,  
    <a(bd)c(de)>  
]