

Задача поиска шаблонов в данных

*Мыслить по шаблону –
вернейший способ завалить дело.*

Дж.Ф. Энрайт

Челябинский
ОБЗОР

**В Челябинске мужчина стащил
из супермаркета пиво
и подгузники**

Происшествия 1 февраля 2020

On the left side of the screenshot, there are four social media icons: Facebook (f), VK (VK), Odnoklassniki (OK), and Twitter (bird).

<https://obzor174.ru/v-chelyabinske-muzhchina-stashchil-iz-supermarketa-pivo-i-podguzniki>

Содержание

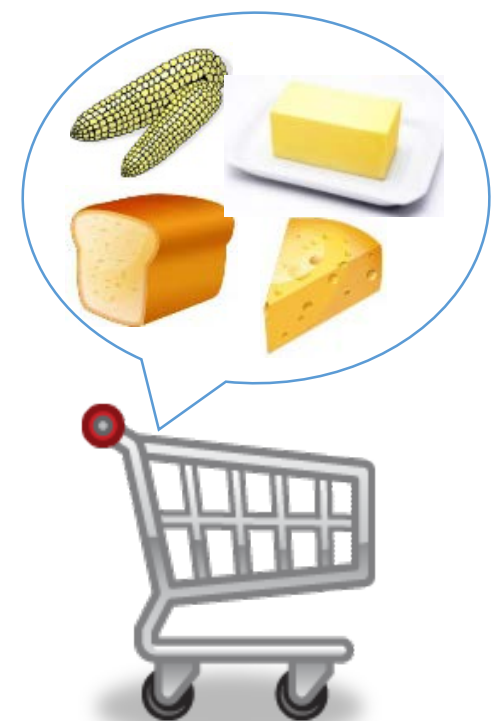
- **Постановка задачи**
 - Частые наборы
 - Ассоциативные правила
- **Основные алгоритмы поиска частых наборов**
 - Apriori
 - Eclat
 - FP-Growth
- Меры полезности шаблонов:
 - support
 - confidence
 - lift и др.
- Поиск шаблонов в Больших данных
 - Компактное представление частых наборов
 - Иерархии наборов
 - Фрагментация и сэмплинг для поиска частых наборов

Задача анализа рыночной корзины

- Какие наборы товаров в супермаркете часто покупают совместно?



...

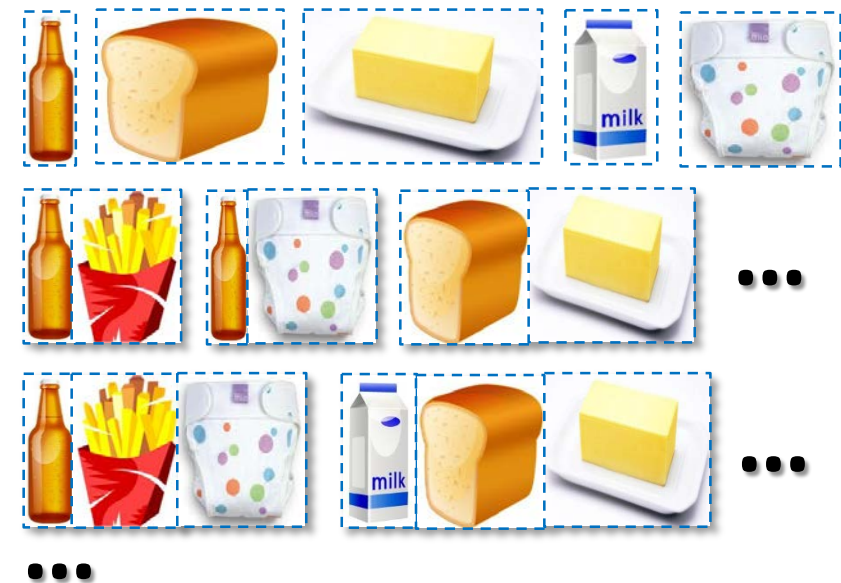


Частый набор vs. шаблон






Частые k -наборы
 $1 \leq k \leq k_{max}$

Шаблоны (ассоциативные правила)
антецедент \rightarrow *консеквент*



ЕСЛИ		ТО	

А при чем здесь  и  ?

- В 1992 г. в США был проведен анализ 1.2 млн. рыночных корзин в 25 магазинах формата «у дома» компании Osco, который выявил частый набор  для покупок в рабочие дни с 5 до 7 час. вечера
- Руководство Osco не стало ставить эти товары рядом на полках, поскольку им были неясны причины такого частого набора
- *Объяснение:* молодая семья приходит с работы домой, жена отправляет мужа в ближайший магазин купить для ребенка , а муж дополнительно покупает себе 

Объекты, наборы, транзакции

- Объекты (items):

$$\mathcal{I} = \{i_1, \dots, i_m\}$$

- Набор (itemset):

$$I \subseteq \mathcal{I}, I \neq \emptyset; k\text{-набор: } I \subseteq \mathcal{I}, |I| = k$$

- Транзакции (transaction database):

$$\mathcal{D} = \{(tid; I) \mid I \subseteq \mathcal{I}\}$$

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Cola, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Cola, Diaper, Milk

$$\mathcal{I} = \{\text{Beer, Bread, Cola, Diaper, Milk}\}$$

Шаблон, поддержка, достоверность

- Шаблон: $A \rightarrow B, A \neq \emptyset, B \neq \emptyset, A \cap B = \emptyset$

- Поддержка

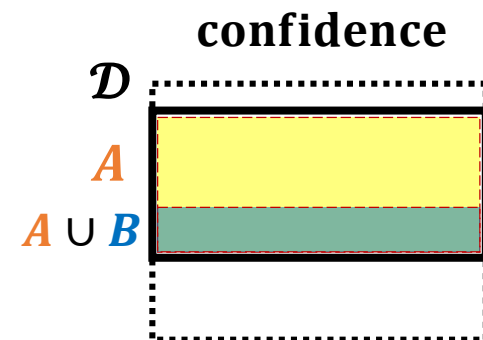
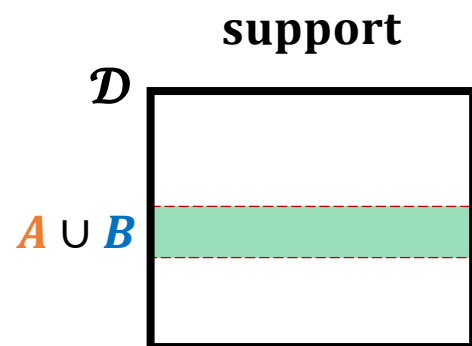
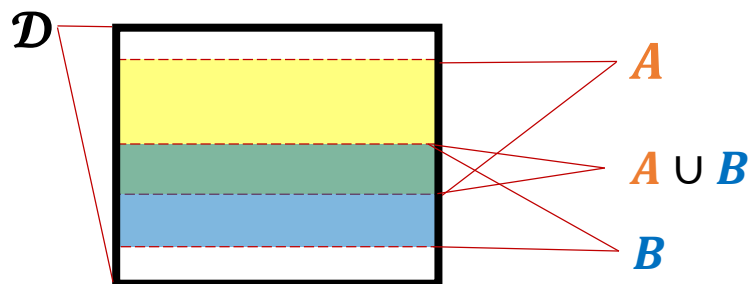
$$\text{sup}(A \rightarrow B) = P(A \cup B) = \frac{|\{t \in \mathcal{D} \mid (A \cup B) \subseteq tI\}|}{|\mathcal{D}|}$$

Более точно:
 $P(E_A \cap E_B)$

- Достоверность

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

Более точно:
 $P(E_B | E_A)$



Поддержка и достоверность шаблона

- Поддержка:

$$\text{sup}(A \rightarrow B) = P(A \cup B)$$

- Достоверность:

$$\text{conf}(A \rightarrow B) = P(B|A)$$

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Cola, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Cola, Diaper, Milk

{Diaper, Milk} → Beer

$$\text{sup} = \frac{2}{5} = 0.4, \quad \text{conf} = \frac{2}{3} = 0.67$$

Частый набор

- $minsup$ – порог поддержки (параметр)
- $I \subseteq \mathcal{I}$ – частый $\Leftrightarrow sup(I) \geq minsup$
- Множество всех частых наборов: $\mathcal{L} = \bigcup_{k=1}^{k_{max}} \mathcal{L}_k$,
 \mathcal{L}_k – множество частых k -наборов

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Cola, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Cola, Diaper, Milk

$minsup = 0.6$

$\mathcal{L}_1 = \{\text{Beer, Bread, Diaper, Milk}\}$

$\mathcal{L}_2 = \left\{ \begin{array}{l} \{\text{Beer, Diaper}\}, \{\text{Bread, Diaper}\}, \\ \{\text{Diaper, Milk}\} \end{array} \right\}$

$minsup = 0.1$

$\mathcal{L}_3 = \left\{ \begin{array}{l} \{\text{Beer, Bread, Diaper}\}, \\ \{\text{Beer, Diaper, Milk}\}, \\ \{\text{Bread, Diaper, Milk}\} \end{array} \right\}$

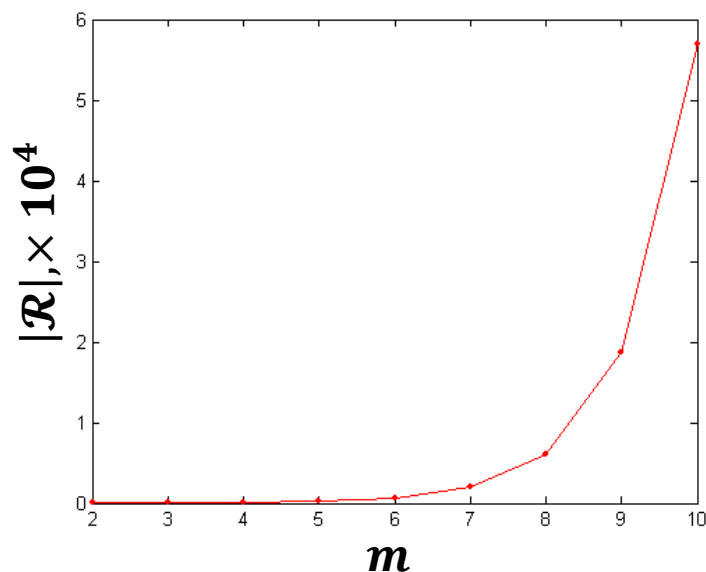
Устойчивый шаблон

- $minconf$ – порог достоверности (параметр)
- Шаблон $A \rightarrow B$ устойчив, если
$$\sup(A \rightarrow B) \geq minsup \wedge \text{conf}(A \rightarrow B) \geq minconf$$
- **Постановка задачи поиска шаблонов**
 - Дано:
$$\mathcal{I} = \{i_1, \dots, i_m\}, \mathcal{D} = \{t_1, \dots, t_n\}, minsup, minconf$$
 - Найти:
$$\mathcal{R} = \{A \rightarrow B \mid A, B \subseteq \mathcal{I}, A \neq \emptyset, B \neq \emptyset, A \cap B = \emptyset, \sup(A \rightarrow B) \geq minsup, \text{conf}(A \rightarrow B) \geq minconf\}$$

Поиск шаблонов: полный перебор

1. Сгенерируем $\mathcal{R} = \{A \rightarrow B \mid A, B \subseteq \mathcal{I}, A \neq \emptyset, B \neq \emptyset, A \cap B = \emptyset\}$
2. $\forall r \in \mathcal{R}$ вычислим $\text{sup}(r)$, $\text{conf}(r)$
3. Отбросим $\forall r \in \mathcal{R}: \text{sup}(r) < \text{minsup}, \text{conf}(r) < \text{minconf}$

- $|\mathcal{R}| = \sum_{k=1}^{m-1} \left[C_k^m \cdot \sum_{i=1}^{m-k} C_i^{m-k} \right] = 3^m - 2^{m+1} + 1$



m	$ \mathcal{R} $
6	602
10	57 002

Поиск устойчивых шаблонов
полным перебором
вычислительно невозможен

Поиск шаблонов: идеи

TID	Items
1	Bread, Milk
2	Beer, Bread, Diaper, Eggs
3	Beer, Cola, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Bread, Cola, Diaper, Milk

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$ ($s=0.4, c=0.67$)

$\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$ ($s=0.4, c=1.0$)

$\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$ ($s=0.4, c=0.67$)

$\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$ ($s=0.4, c=0.67$)

$\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$ ($s=0.4, c=0.5$)

$\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$ ($s=0.4, c=0.5$)

- Шаблоны примера – разбиения одного набора
- Разбиения одного набора имеют равную поддержку, но могут иметь разную достоверность
- Для поиска шаблонов обработку поддержки и достоверности можно отделить друг от друга

Поиск шаблонов: алгоритм

1. Найдем все частые наборы
(с поддержкой не ниже *minsup*)
2. Сгенерируем шаблоны, выполняя разбиение
каждого частого набора; устойчивыми будут
шаблоны с достоверностью не ниже *minconf*

Поиск частых наборов: полный перебор

1. Сгенерируем кандидатов в частые наборы: $\mathcal{C} = \mathcal{P}(\mathcal{I}) \setminus \emptyset$
 2. $\forall c \in \mathcal{C}$ вычислим $\text{sup}(c)$
 3. Отбросим $\forall c \in \mathcal{C}: \text{sup}(c) < \text{minsup}$
- Сложность: $O(n \cdot 2^m \cdot w)$,
где w – средняя длина транзакции

Поиск частых наборов
полным перебором
вычислительно невозможен

Отбрасывание заведомо редких наборов

- Антимонотонность поддержки

$$\forall X, Y: X \subseteq Y \Leftrightarrow \text{sup}(X) \geq \text{sup}(Y)$$

- поддержка набора всегда не больше поддержки любого своего подмножества

- Принцип Априори

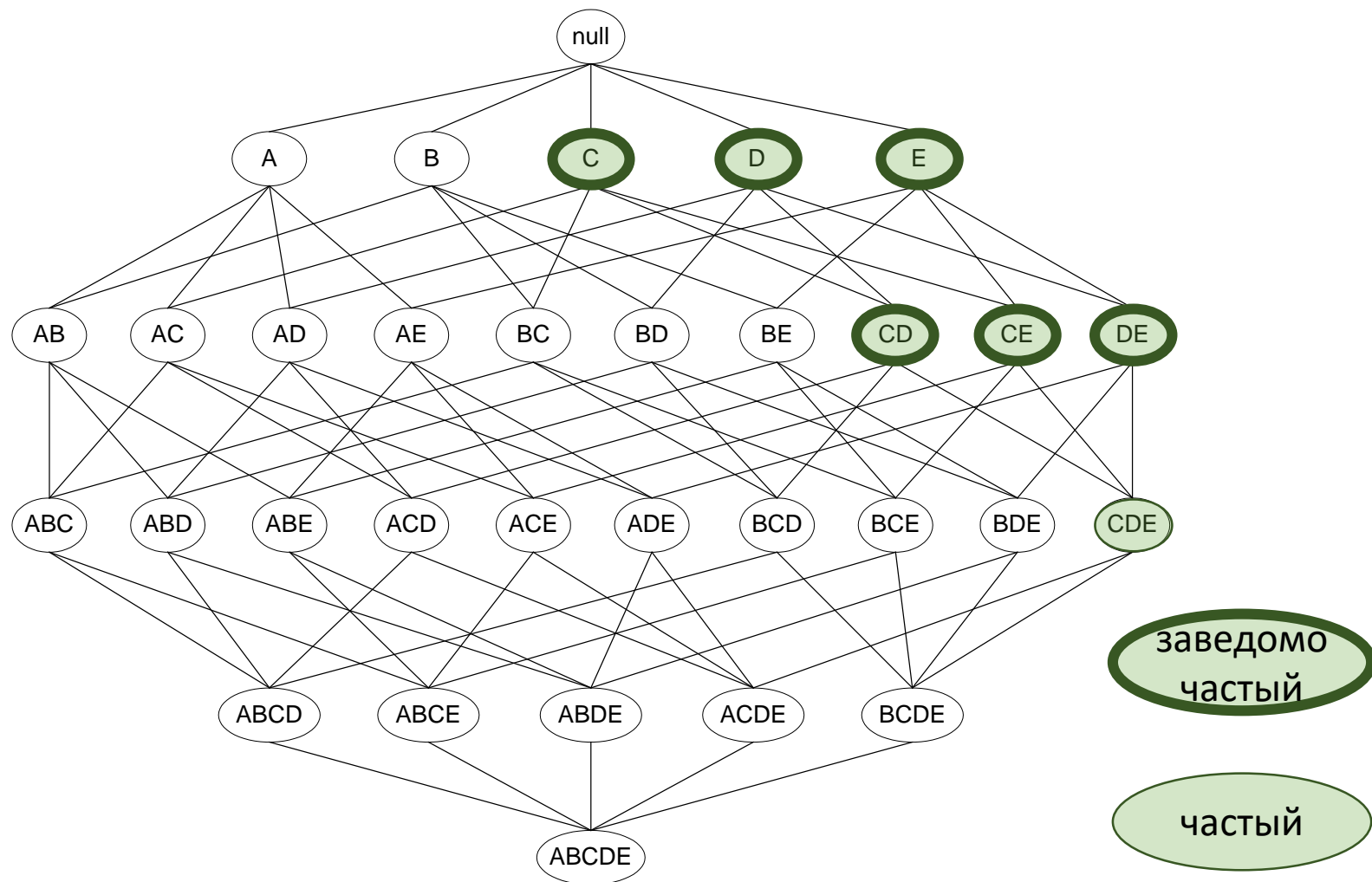
- любое подмножество частого набора является частым набором

$$\text{sup}(Y) \geq \text{minsup} \Leftrightarrow \forall X \subseteq Y \text{ sup}(X) \geq \text{minsup}$$

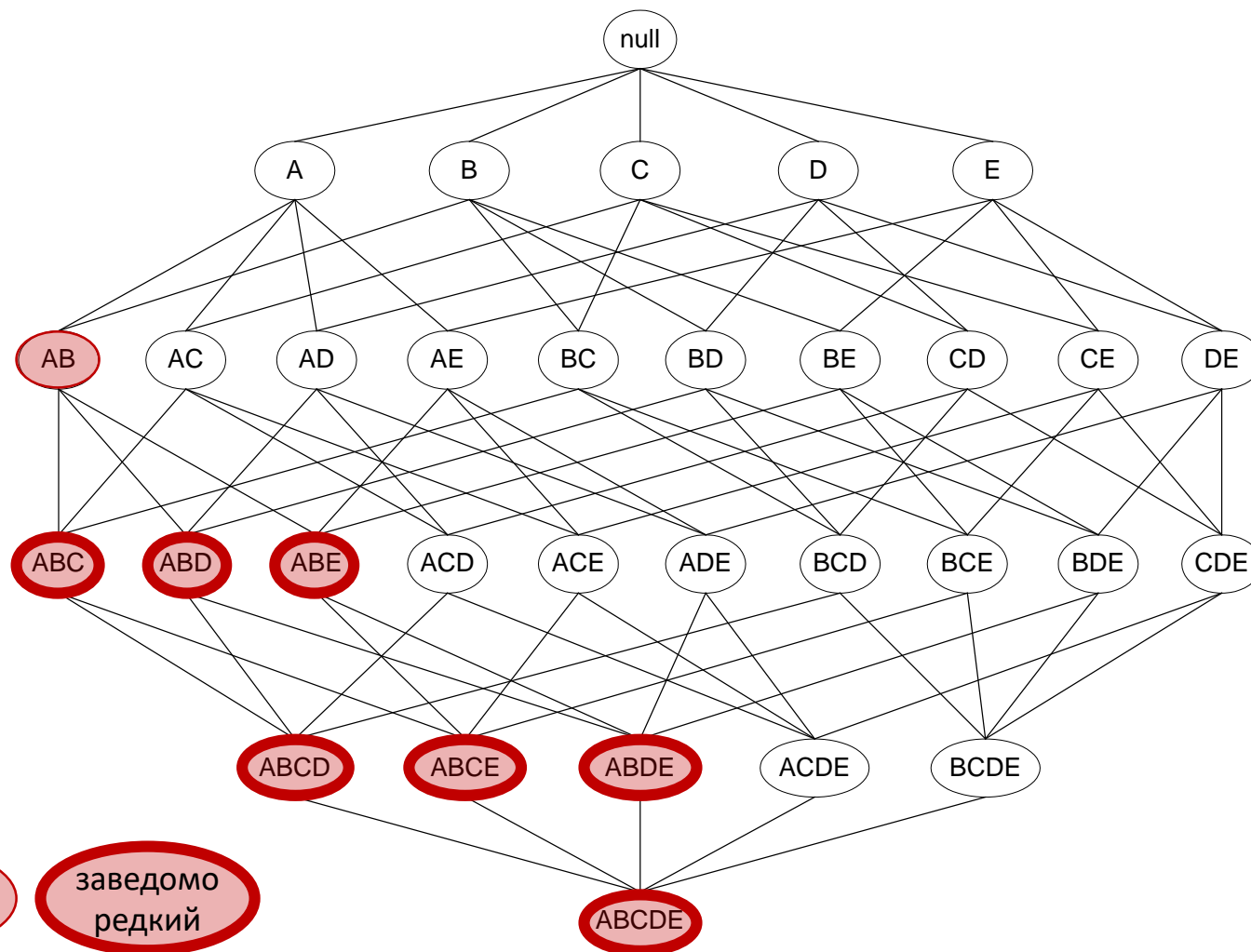
- если некое подмножество набора является редким набором, то набор является редким

$$\text{sup}(Y) < \text{minsup} \Leftrightarrow \exists X \subseteq Y \text{ sup}(X) < \text{minsup}$$

Принцип Априори



Принцип Априори



Генерация наборов-кандидатов

- $k = 1: C_1 := \mathcal{I}$ (множество 1-наборов)
- $k = 2: C_2 := \mathcal{L}_1 \times \mathcal{L}_1$
 - Декартово произведение множества частых 1-наборов на себя
- $k \geq 3$:
 - Соединение множества частых $(k - 1)$ -наборов с самим собой
 - $C_k := \mathcal{L}_{k-1} \bowtie_{\Theta} \mathcal{L}_{k-1}$
 $X = (x_1, \dots, x_{k-1}), Y = (y_1, \dots, y_{k-1}), X, Y \in \mathcal{L}_{k-1}$,
элементы в наборах лексикографически упорядочены
 $\Theta = (\bigwedge_{i=1}^{k-2} x_i = y_i) \wedge (x_{k-1} < y_{k-1})$
 - $X \bowtie_{\Theta} Y = (x_1, \dots, x_{k-2}, x_{k-1}, y_{k-1})$
 - Отбрасывание заведомо редких наборов из C_k без вычисления поддержки по принципу Априори
 - если $I \subseteq C_k$, но $I \notin \mathcal{L}_{k-1}$, то отбросить I

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$

Декартово произведение

$$C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$$

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$

Декартово произведение

$$C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$$

Вычисление поддержки

$$\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \del{DE}\}$$

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$

Декартово произведение

$$C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$$

Вычисление поддержки

$$\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \emptyset E\}$$

- $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE\}$

Самосоединение

$$C_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$$

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$

Декартово произведение

$$C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$$

Вычисление поддержки

$$\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \cancel{DE}\}$$

- $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE\}$

Самосоединение

$$C_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$$

Вычисление поддержки

$$\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, \cancel{ACE}, \cancel{ADE}, BCD, \cancel{BCE}, BDE, CDE\}$$

- $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$

Самосоединение

$$C_4 = \{ABCD, ABCE, ABDE\}$$

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$
 Декартово произведение
 $C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$
 Вычисление поддержки
 $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \emptyset E\}$
- $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE\}$
 Самосоединение
 $C_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$
 Вычисление поддержки
 $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$
- $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$
 Самосоединение
 $C_4 = \{ABCD, ABCE, ABDE\}$
 Отбрасывание
 $C_4 = \{ABCD, ABCE, ABDE\}$ и $BCE \notin \mathcal{L}_3$, $ADE \notin \mathcal{L}_3$

Генерация наборов-кандидатов: пример

- $C_1 = \mathcal{I} = \{A, B, C, D, E\}$, $\mathcal{L}_1 = \{A, B, C, D, E\}$
 Декартово произведение
 $C_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, DE\}$
 Вычисление поддержки
 $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \cancel{DE}\}$
- $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE\}$
 Самосоединение
 $C_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$
- $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$
 Самосоединение
 $C_4 = \{ABCD, ABCE, ABDE\}$
 Отбрасывание
 $C_4 = \{ABCD, ABCE, ABDE\} \Rightarrow C_4 = \{ABCD\}$

Алгоритм Априори

Input:

- D , a database of transactions;
- min_sup , the minimum support count threshold.

Output: L , frequent itemsets in D .

Method:

```
(1)  $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;  
(2) for ( $k = 2; L_{k-1} \neq \phi; k++$ ) {  
(3)    $C_k = \text{apriori\_gen}(L_{k-1})$ ;  
(4)   for each transaction  $t \in D$  { // scan  $D$  for counts  
(5)      $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates  
(6)     for each candidate  $c \in C_t$   
(7)        $c.\text{count}++$ ;  
(8)   }  
(9)    $L_k = \{c \in C_k \mid c.\text{count} \geq min\_sup\}$   
(10) }  
(11) return  $L = \cup_k L_k$ ;
```

Алгоритм Априори

procedure `apriori_gen`(L_{k-1} :frequent $(k - 1)$ -itemsets)

```
(1)  for each itemset  $l_1 \in L_{k-1}$ 
(2)    for each itemset  $l_2 \in L_{k-1}$ 
(3)      if  $(l_1[1] = l_2[1]) \wedge (l_1[2] = l_2[2])$ 
           $\wedge \dots \wedge (l_1[k-2] = l_2[k-2]) \wedge (l_1[k-1] < l_2[k-1])$  then {
(4)         $c = l_1 \bowtie l_2$ ; // join step: generate candidates
(5)        if has_infrequent_subset( $c, L_{k-1}$ ) then
(6)          delete  $c$ ; // prune step: remove unfruitful candidate
(7)        else add  $c$  to  $C_k$ ;
(8)      }
(9)  return  $C_k$ ;
```

procedure `has_infrequent_subset`(c : candidate k -itemset;

L_{k-1} : frequent $(k - 1)$ -itemsets); // use prior knowledge

```
(1)  for each  $(k - 1)$ -subset  $s$  of  $c$ 
(2)    if  $s \notin L_{k-1}$  then
(3)      return TRUE;
(4)  return FALSE;
```

```
(1)   $L_1 = \text{find\_frequent\_1-itemsets}(D)$ ;
(2)  for  $(k = 2; L_{k-1} \neq \phi; k++)$  {
(3)     $C_k = \text{apriori\_gen}(L_{k-1})$ ;
(4)    for each transaction  $t \in D$  { // scan  $D$  for counts
(5)       $C_t = \text{subset}(C_k, t)$ ; // get the subsets of  $t$  that are candidates
(6)      for each candidate  $c \in C_t$ 
(7)         $c.\text{count}++$ ;
(8)    }
(9)     $L_k = \{c \in C_k \mid c.\text{count} \geq \text{min\_sup}\}$ 
(10) }
(11) return  $L = \cup_k L_k$ ;
```

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



$C_1 = \mathcal{I}$

Items	SUP
I1	
I2	
I3	
I4	
I5	
I6	

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



$C_1 = \mathcal{I}$

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2
I6	1

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



C_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2
I6	1



\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2



$C_2 = \mathcal{L}_1 \times \mathcal{L}_1$

Items	SUP
I1,I2	
I1,I3	
I1,I4	
I1,I5	
I2,I3	
I2,I4	
I2,I5	
I3,I4	
I3,I5	
I4,I5	

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2



$C_2 = \mathcal{L}_1 \times \mathcal{L}_1$

Items	SUP
I1,I2	4
I1,I3	4
I1,I4	1
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2
I3,I4	0
I3,I5	1
I4,I5	0

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2



$C_2 = \mathcal{L}_1 \times \mathcal{L}_1$

Items	SUP
I1,I2	4
I1,I3	4
I1,I4	1
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2
I3,I4	0
I3,I5	1
I4,I5	0



\mathcal{L}_2

Items	SUP
I1,I2	4
I1,I3	4
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_2

Items	SUP
I1,I2	4
I1,I3	4
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2



$\mathcal{C}_3 = \mathcal{L}_2 \bowtie \mathcal{L}_2$

Items	SUP
I1,I2,I3	
I1,I2,I5	
I1,I3,I5	
I2,I3,I4	
I2,I3,I5	
I2,I4,I5	

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_2

Items	SUP
I1,I2	4
I1,I3	4
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2



$\mathcal{C}_3 = \mathcal{L}_2 \bowtie \mathcal{L}_2$

Items	SUP
I1,I2,I3	
I1,I2,I5	
I1, I3,I5	
I2, I3,I4	
I2, I3,I5	
I2, I4,I5	

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_2

Items	SUP
I1,I2	4
I1,I3	4
I1,I5	2
I2,I3	4
I2,I4	2
I2,I5	2



$C_3 = \mathcal{L}_2 \bowtie \mathcal{L}_2$

Items	SUP
I1,I2,I3	2
I1,I2,I5	2



\mathcal{L}_3

Items	SUP
I1,I2,I3	2
I1,I2,I5	2

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



\mathcal{L}_3

Items	SUP
I1, I2, I3	2
I1, I2, I5	2



$C_4 = \mathcal{L}_3 \bowtie \mathcal{L}_3$

Items	SUP
I1, I2, I3, I5	



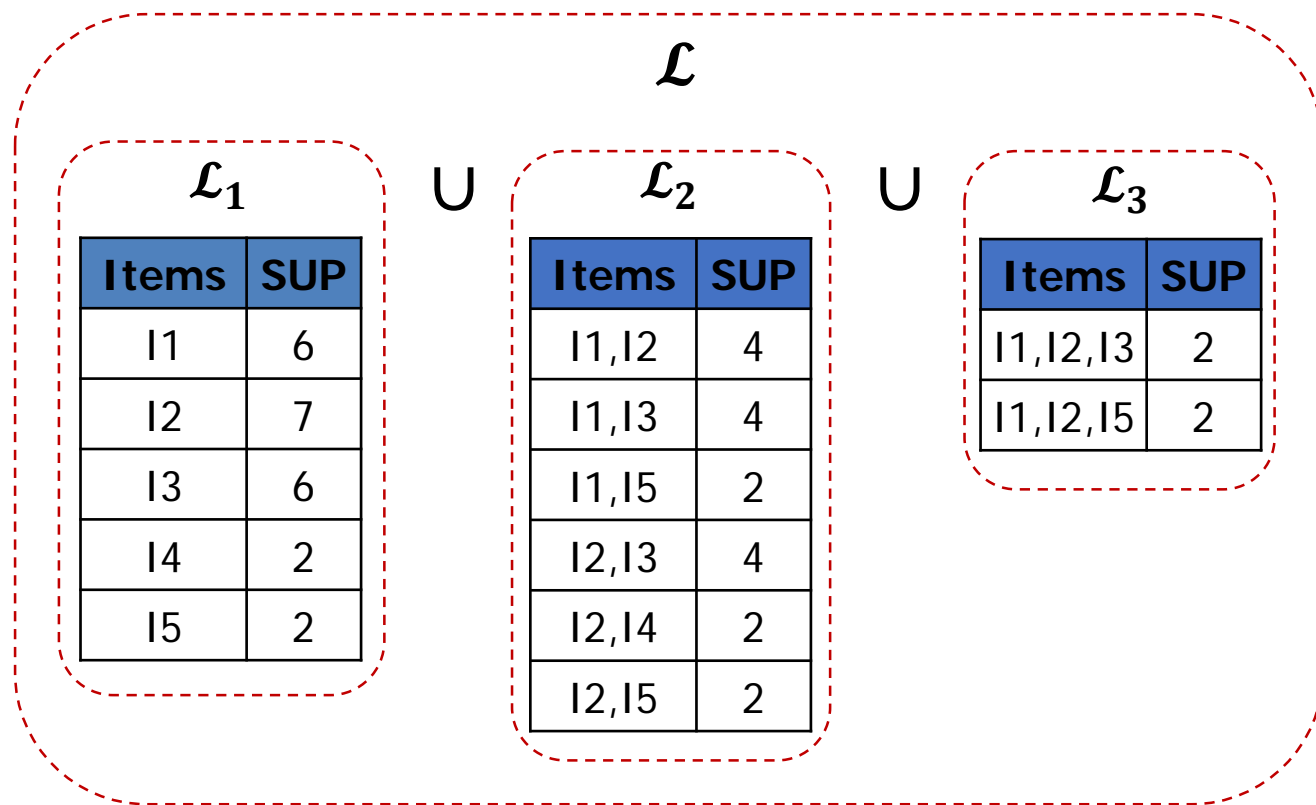
$\mathcal{L}_4 = \emptyset$

Items	SUP

Алгоритм Apriori: пример

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



Поиск шаблонов: алгоритм

$$\mathcal{L} := \bigcup_{k=1}^{k_{max}} \{I \subseteq \mathcal{I} \mid |I| = k, \sup(I) \geq \text{minsup}\}$$

$$\mathcal{R} := \emptyset$$

for all $I \in \mathcal{L}$

for all $S \in \mathcal{P}(I) \setminus \emptyset$ **do**

if $\frac{\sup(I)}{\sup(S)} \geq \text{minconf}$ **then**

$pattern := "S \rightarrow I \setminus S"$

$\mathcal{R} := \mathcal{R} \cup pattern$

$\mathcal{P}(I)$ – множество всех подмножеств I

Генерация шаблонов: пример

$D,$

$minsup = 0.2$

$minconf = 0.7$

\mathcal{L}_3

Шаблоны

для $\{I1, I2, I5\}$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



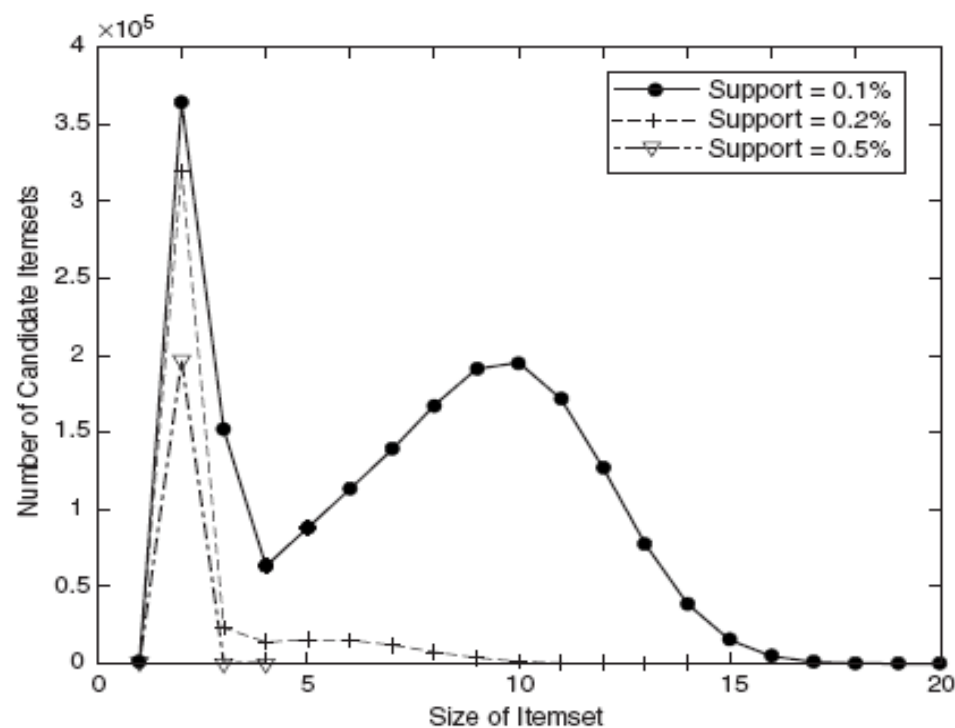
Items	SUP
I1, I2, I3	2
I1, I2, I5	2



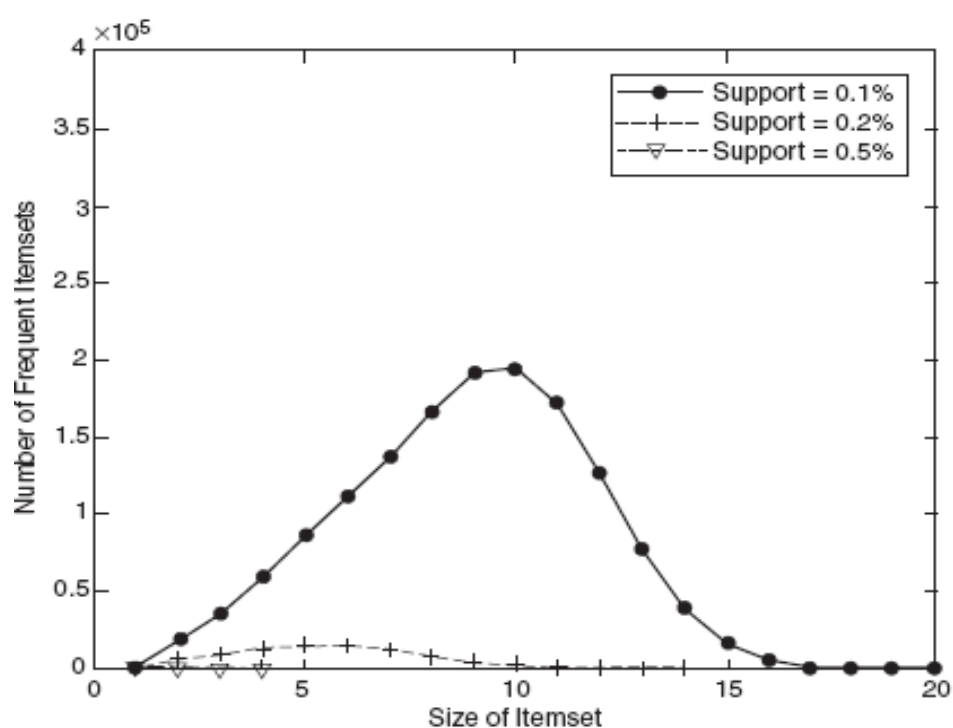
Rule	CONF
$\{I1, I2\} \rightarrow I5$	2/4
$\{I1, I5\} \rightarrow I2$	2/2
$\{I2, I5\} \rightarrow I1$	2/2
$I1 \rightarrow \{I2, I5\}$	2/6
$I2 \rightarrow \{I1, I5\}$	2/7
$I5 \rightarrow \{I1, I2\}$	2/2

Факторы, влияющие на сложность Apriori

minsup (порог поддержки)



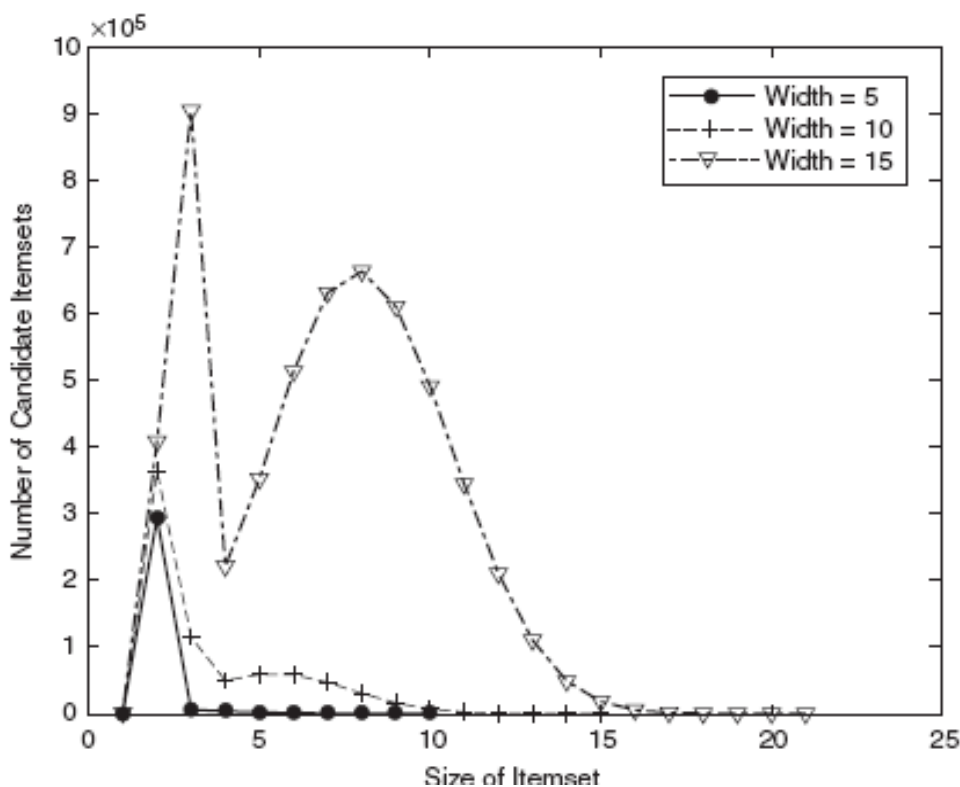
$$\bigcup_k |C_k|$$



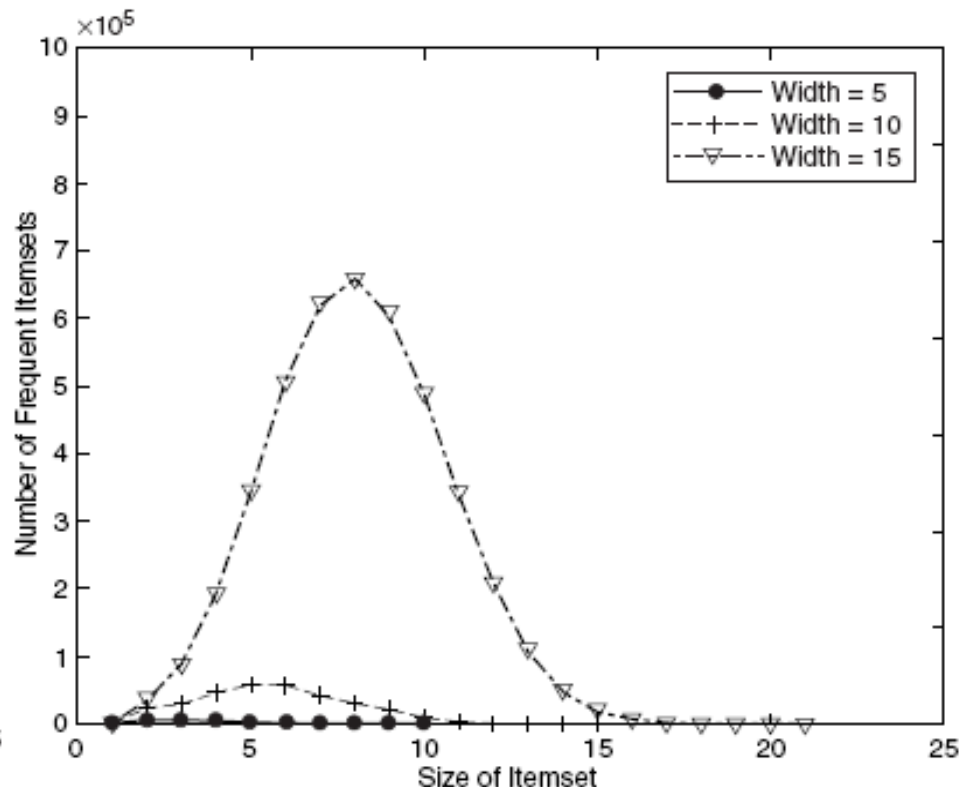
$$\bigcup_k |\mathcal{L}_k|$$

Факторы, влияющие на сложность Apriori

$|\bar{T}|$ (средняя длина транзакции)



$$\bigcup_k |C_k|$$



$$\bigcup_k |\mathcal{L}_k|$$

Недостатки Apriori

- Потенциально большое число генерируемых кандидатов
 - $C_1 = \mathcal{I}$
 - $C_2 = \mathcal{L}_1 \times \mathcal{L}_1 \Rightarrow |C_2| = |\mathcal{L}_1|^2$
 - $k \geq 3: C_k = \mathcal{L}_{k-1} \bowtie_{\Theta} \mathcal{L}_{k-1}$
- Многократное сканирование D и проверка вхождения каждого кандидата в каждую транзакцию
 - $O(|D| \cdot |C_k| \cdot |\bar{T}|)$

Алгоритм ECLAT

- Вертикальный формат данных
 - $TIDlist(A) ::= \{T \in D \mid A \subseteq T\}$
 - $sup(A) = |TIDlist(A)|$
 - $sup(A, B) = |TIDlist(A) \cap TIDlist(B)|$

Алгоритм ECLAT

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6



C_1

Item	TID_list	SUP
I1	10, 40, 50, 70, 80, 90	6
I2	10, 20, 30, 40, 60, 80, 90	7
I3	30, 50, 60, 70, 80, 90	6
I4	20, 40	2
I5	10, 80	2
I6	100	1



\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2

Алгоритм ECLAT

 \mathcal{L}_1

Item	TID_list	SUP
I1	10, 40, 50, 70, 80, 90	6
I2	10, 20, 30, 40, 60, 80, 90	7
I3	30, 50, 60, 70, 80, 90	6
I4	20, 40	2
I5	10, 80	2



$$C_2 = \mathcal{L}_1 \times \mathcal{L}_1$$

$$TIDlist(A, B) = TIDlist(A) \cap TIDlist(B)$$


 C_2

Items	TID_list	SUP
I1, I2	10, 40, 80, 90	4
I1, I3	50, 70, 80, 90	4
I1, I4	40	1
I1, I5	10, 80	2
I2, I3	30, 60, 80, 90	4
I2, I4	20, 40	2
I2, I5	10, 80	2
I3, I5	80	1

Алгоритм ECLAT

 \mathcal{L}_2

Items	TID_list	SUP
11, 12	10, 40, 80, 90	4
11, 13	50, 70, 80, 90	4
11, 15	10, 80	2
12, 13	30, 60, 80, 90	4
12, 14	20, 40	2
12, 15	10, 80	2

$$C_3 = \mathcal{L}_2 \bowtie \mathcal{L}_2$$

 C_3

Items
11, 12, 13
11, 12, 15
11, 13, 15
12, 13, 14
12, 13, 15
12, 14, 15



Apriori

$$TIDlist(A, B) = TIDlist(A) \cap TIDlist(B)$$

 \mathcal{L}_3

Items	TID_list	SUP
11, 12, 13	80, 90	2
11, 12, 15	10, 80	2

 \mathcal{L}_1

Item	TID_list	SUP
11	10, 40, 50, 70, 80, 90	6
12	10, 20, 30, 40, 60, 80, 90	7
13	30, 50, 60, 70, 80, 90	6
14	20, 40	2
15	10, 80	2

Алгоритм FP-Growth

- Поиск частых наборов без генерации кандидатов
- Длинный частый набор строится из коротких, используя только локально частые наборы
 - Пусть abc частый набор
 - Выполнить проектирование \mathcal{D} на abc : $\mathcal{D}|abc$ (получить транзакции, содержащие abc)
 - Если d – частый набор в $\mathcal{D}|abc$, то $abcd$ – частый набор

Алгоритм FP-Growth: предобработка

$D, \text{minsup} = 0.2$

TID	Items
10	I1, I2, I5
20	I2, I4
30	I2, I3
40	I1, I2, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I1, I2, I3, I5
90	I1, I2, I3
100	I6

\mathcal{L}_1

Items	SUP
I1	6
I2	7
I3	6
I4	2
I5	2



Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



D

TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

Редкие наборы
не рассматриваются

Построение FP-дерева

$D, \text{minsup} = 0.2$ sorted \mathcal{L}_1

TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



I2:7

I1:6

I3:6

I4:2

I5:2

NULL

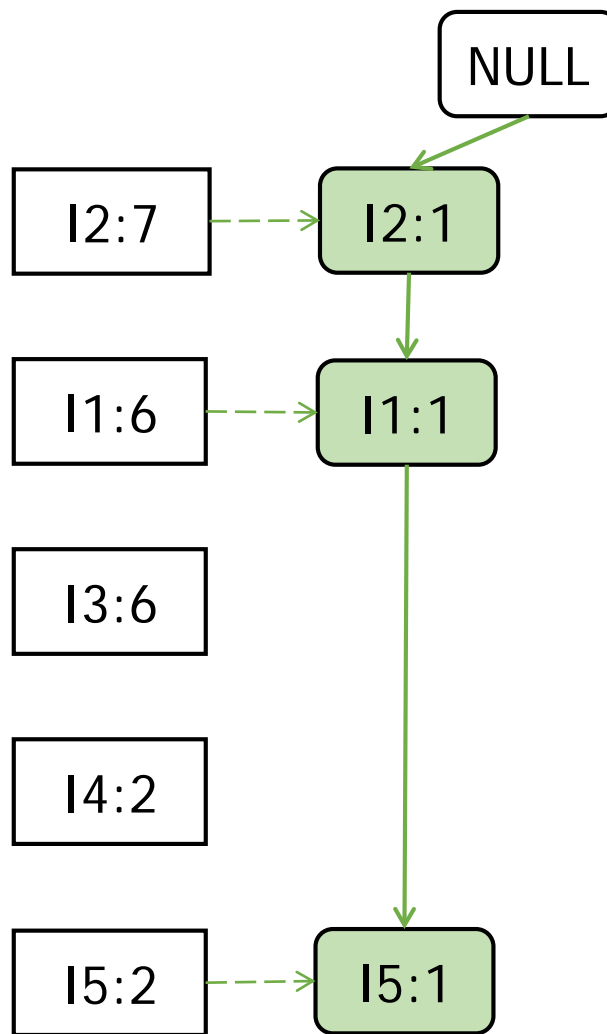
Построение FP-дерева

$D, \text{minsup} = 0.2$ sorted \mathcal{L}_1



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



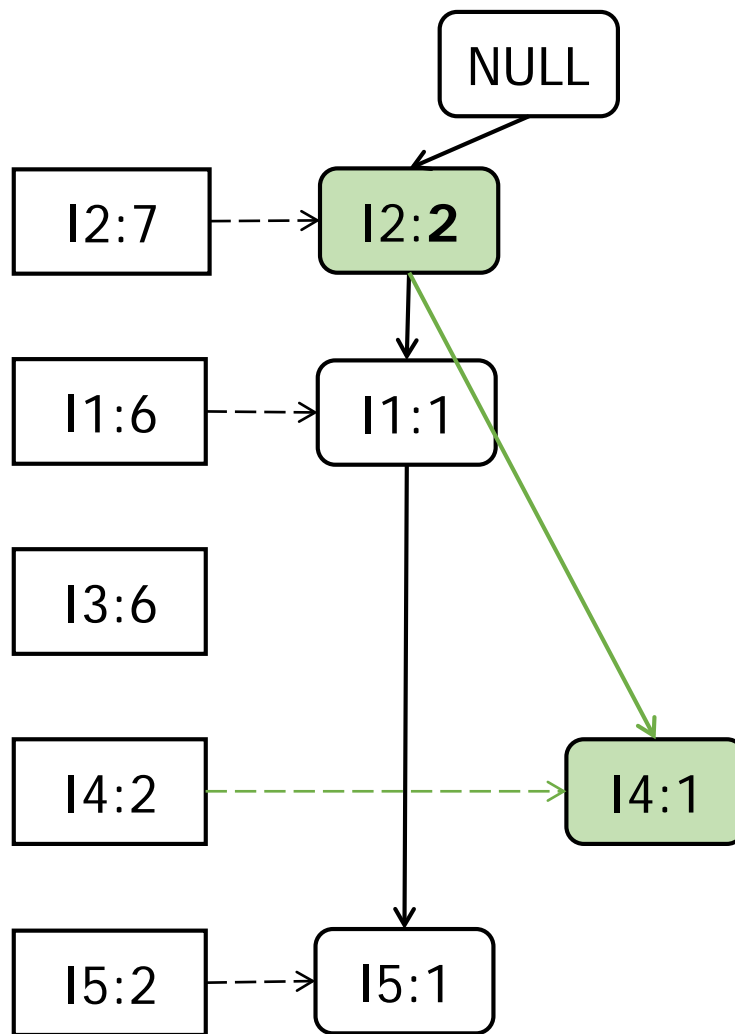
Построение FP-дерева

$D, \text{minsup} = 0.2$ sorted \mathcal{L}_1



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



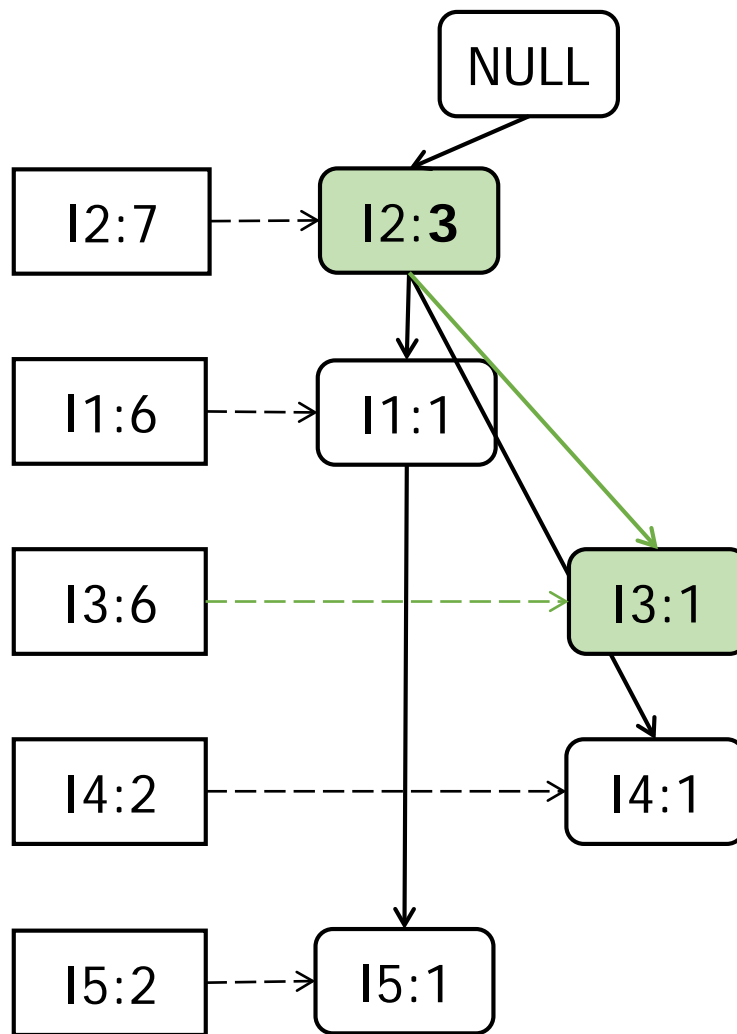
Построение FP-дерева

$D, \text{minsup} = 0.2$ sorted \mathcal{L}_1



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



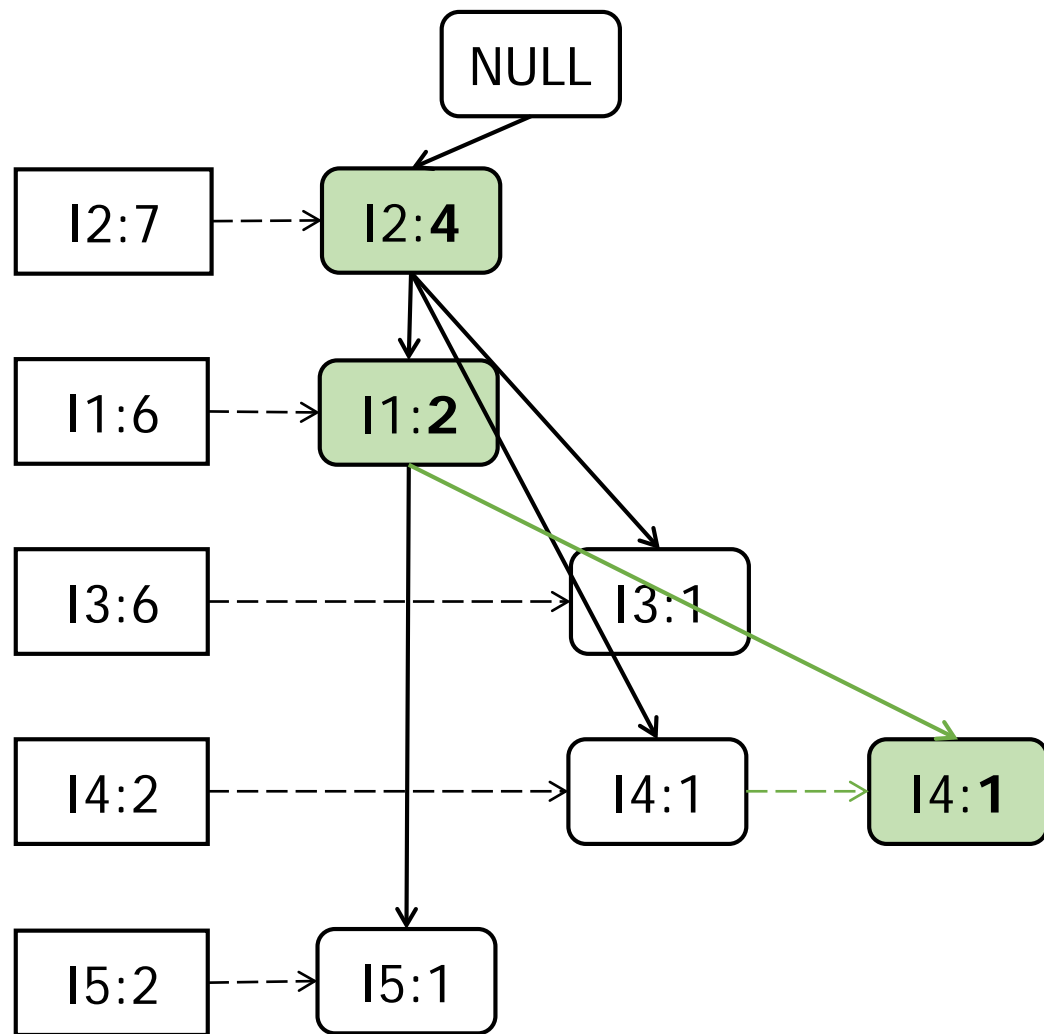
Построение FP-дерева

$D, minsup = 0.2$ sorted \mathcal{L}_1



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

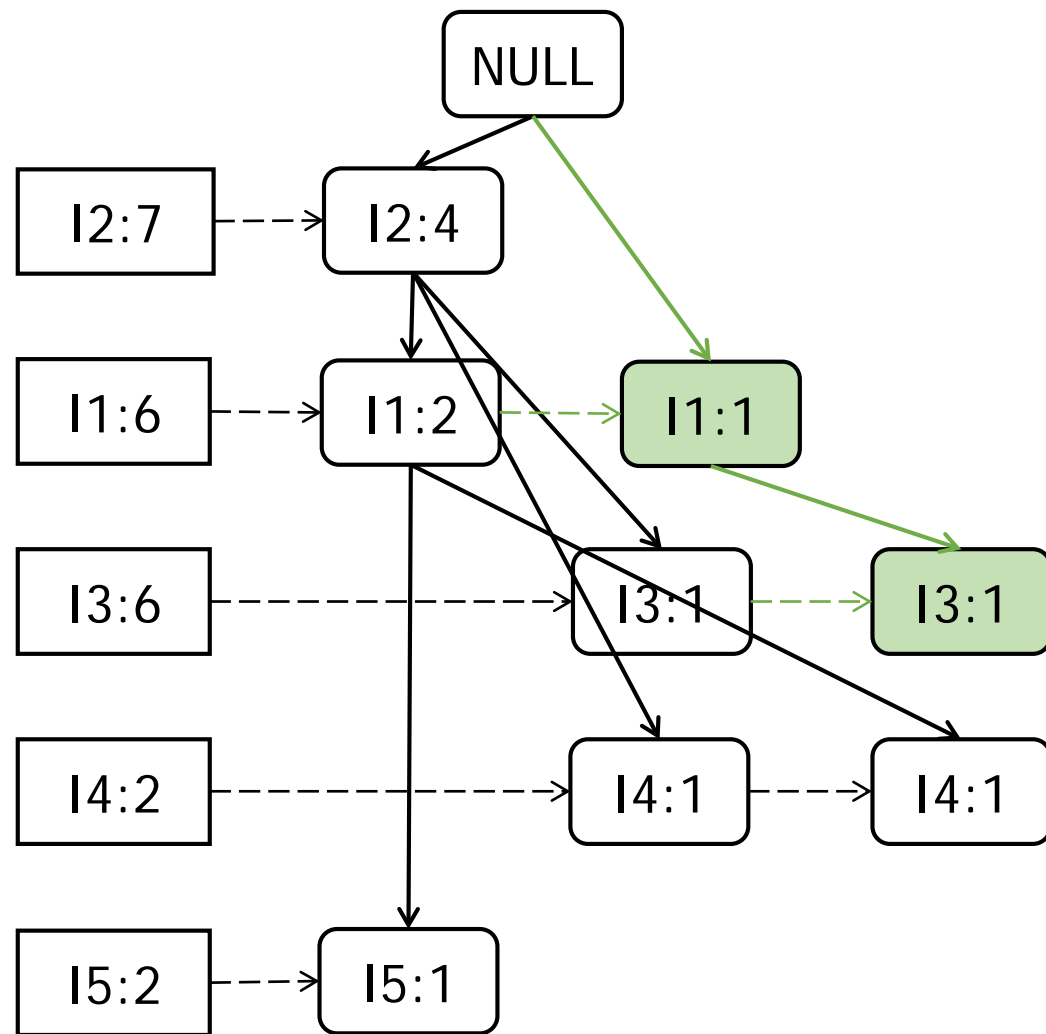
Items	SUP
I2	7
I1	6
I3	6
I4	2
I5	2



Построение FP-дерева

$D, \text{minsup} = 0.2$

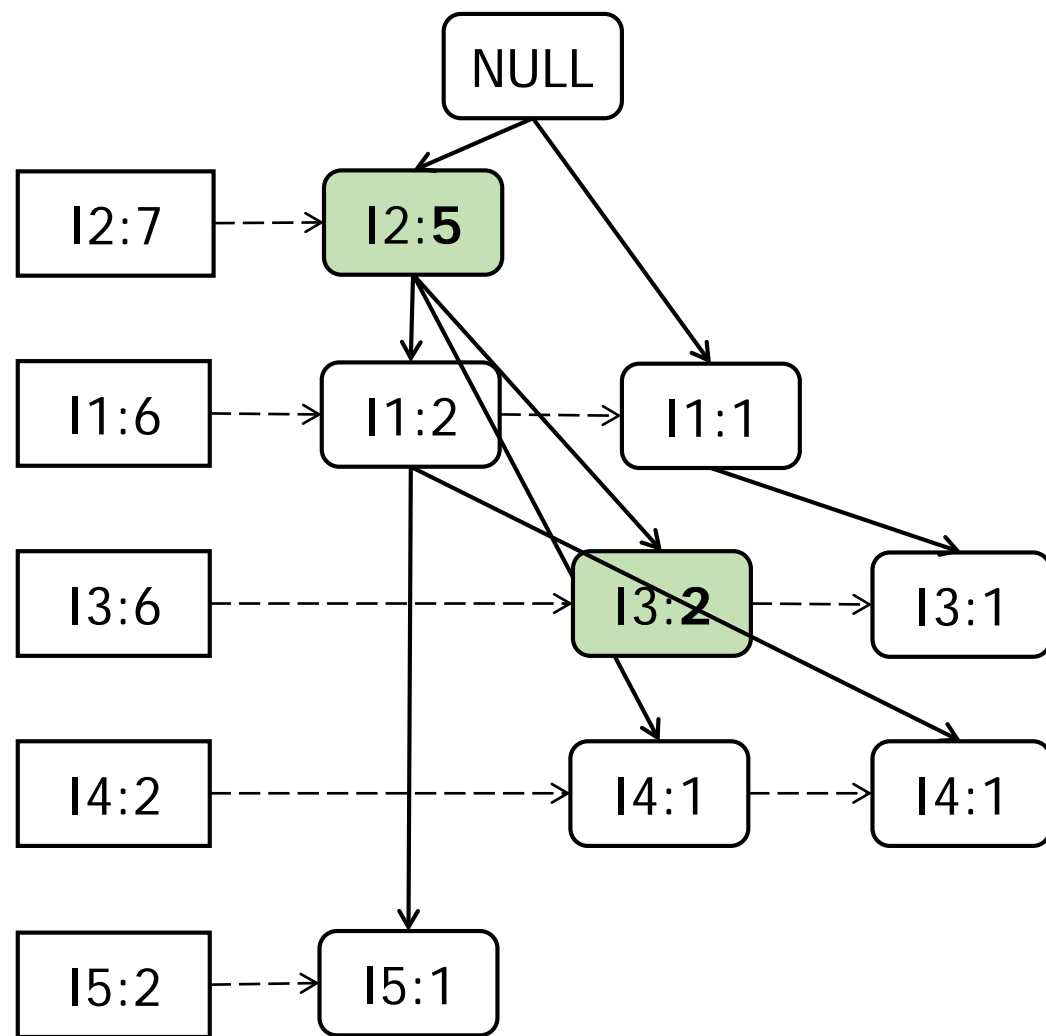
TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6



Построение FP-дерева

$D, \text{minsup} = 0.2$

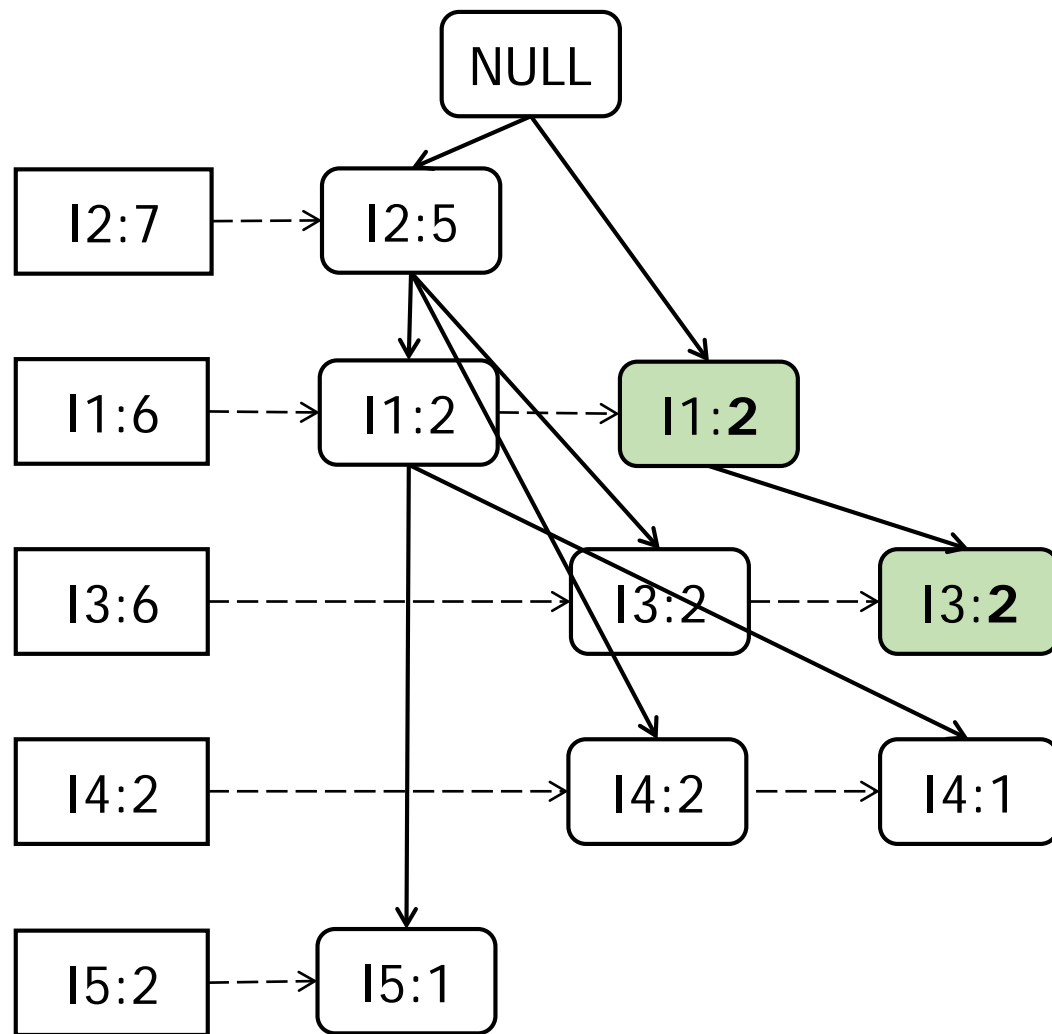
TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6



Построение FP-дерева

$D, \text{minsup} = 0.2$

TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

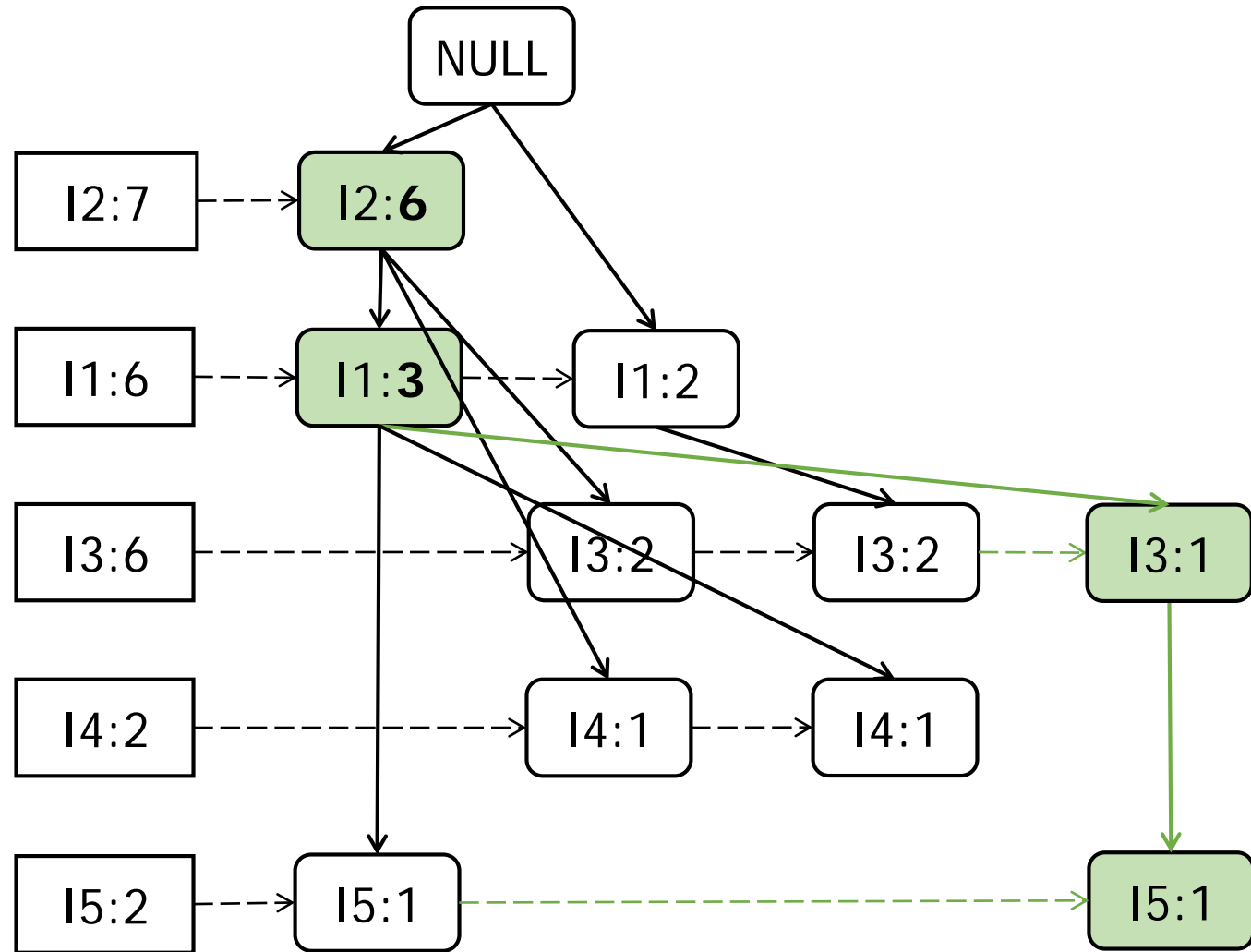


Построение FP-дерева

$D, \text{minsup} = 0.2$



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

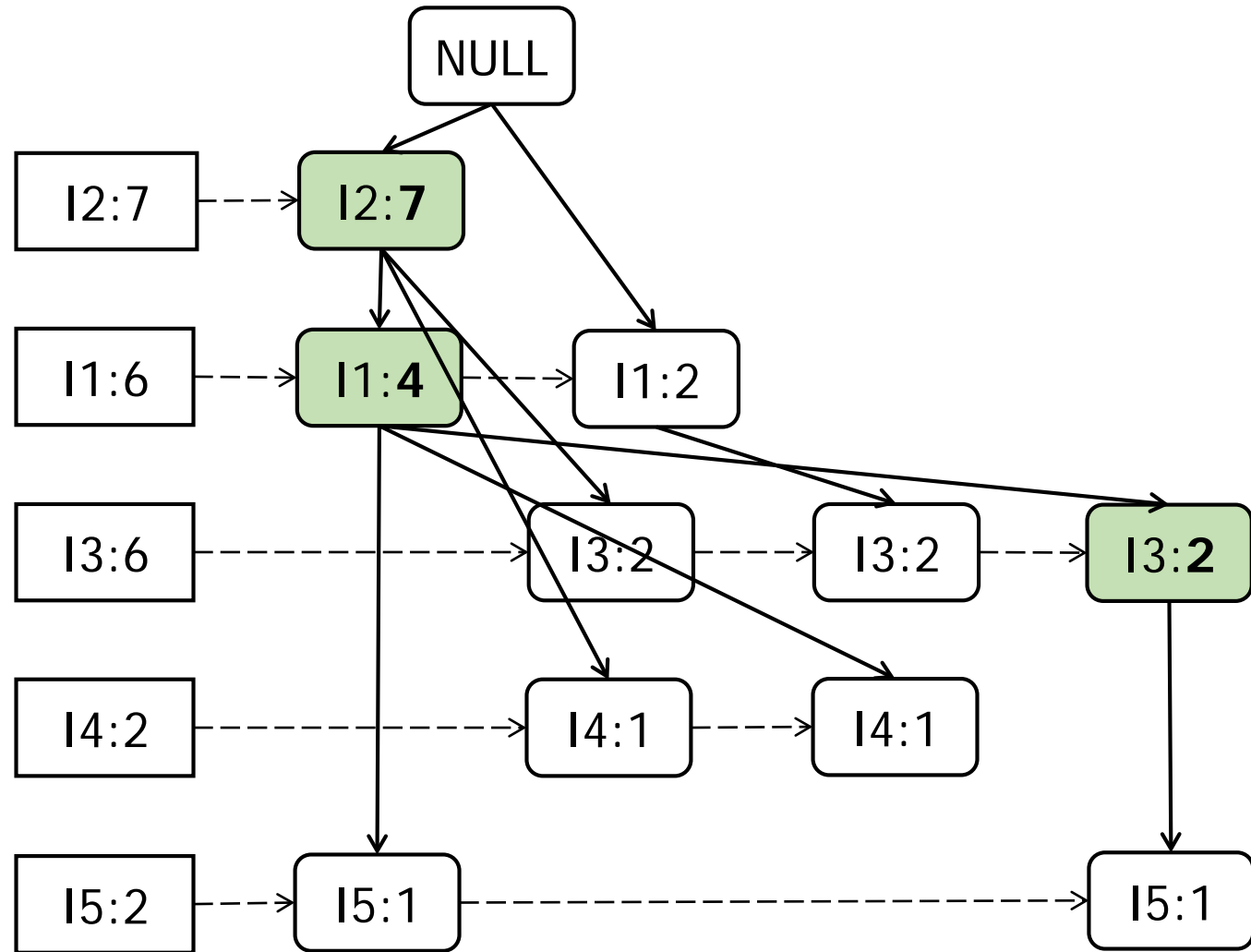


Построение FP-дерева

$D, \text{minsup} = 0.2$



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

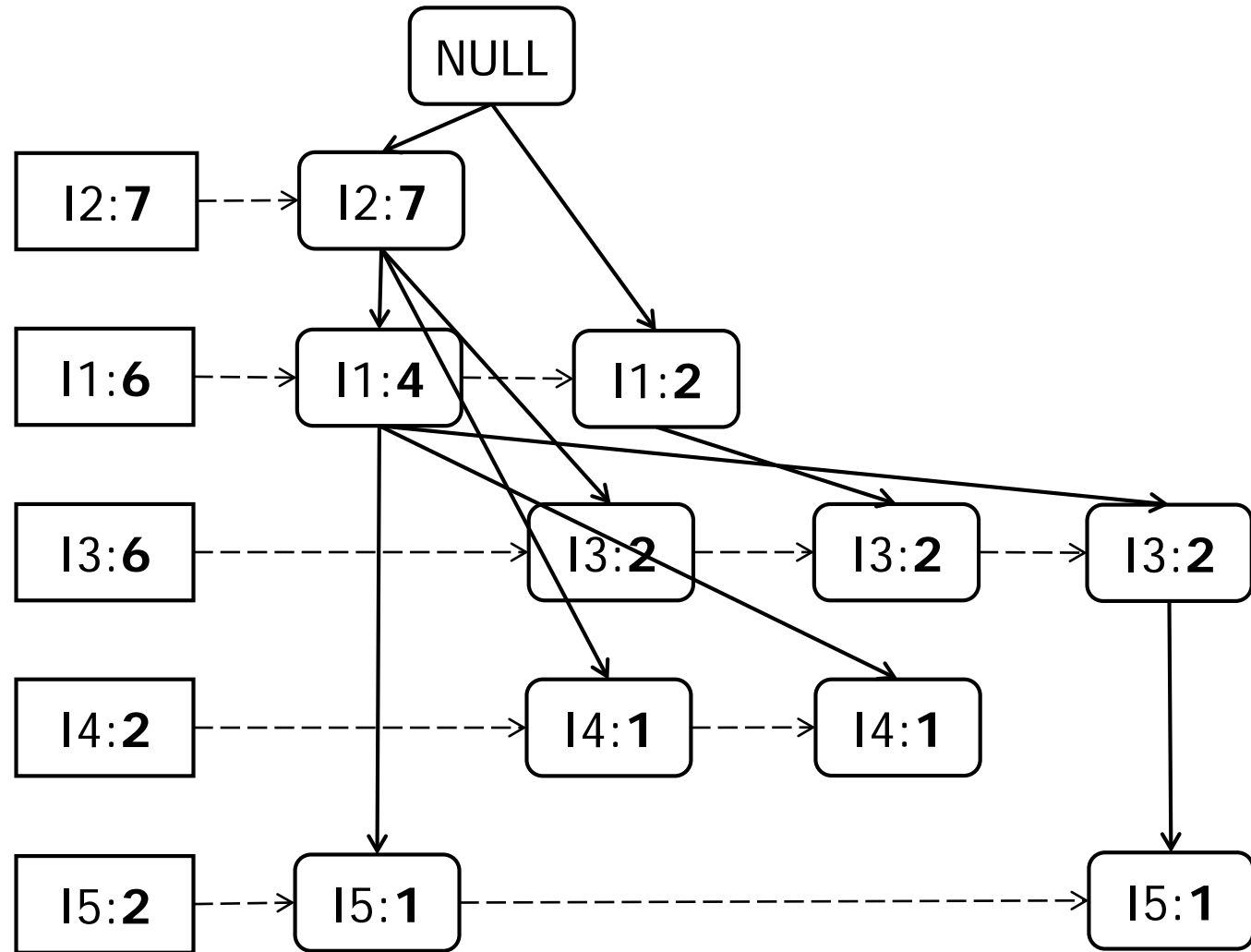


Построение FP-дерева

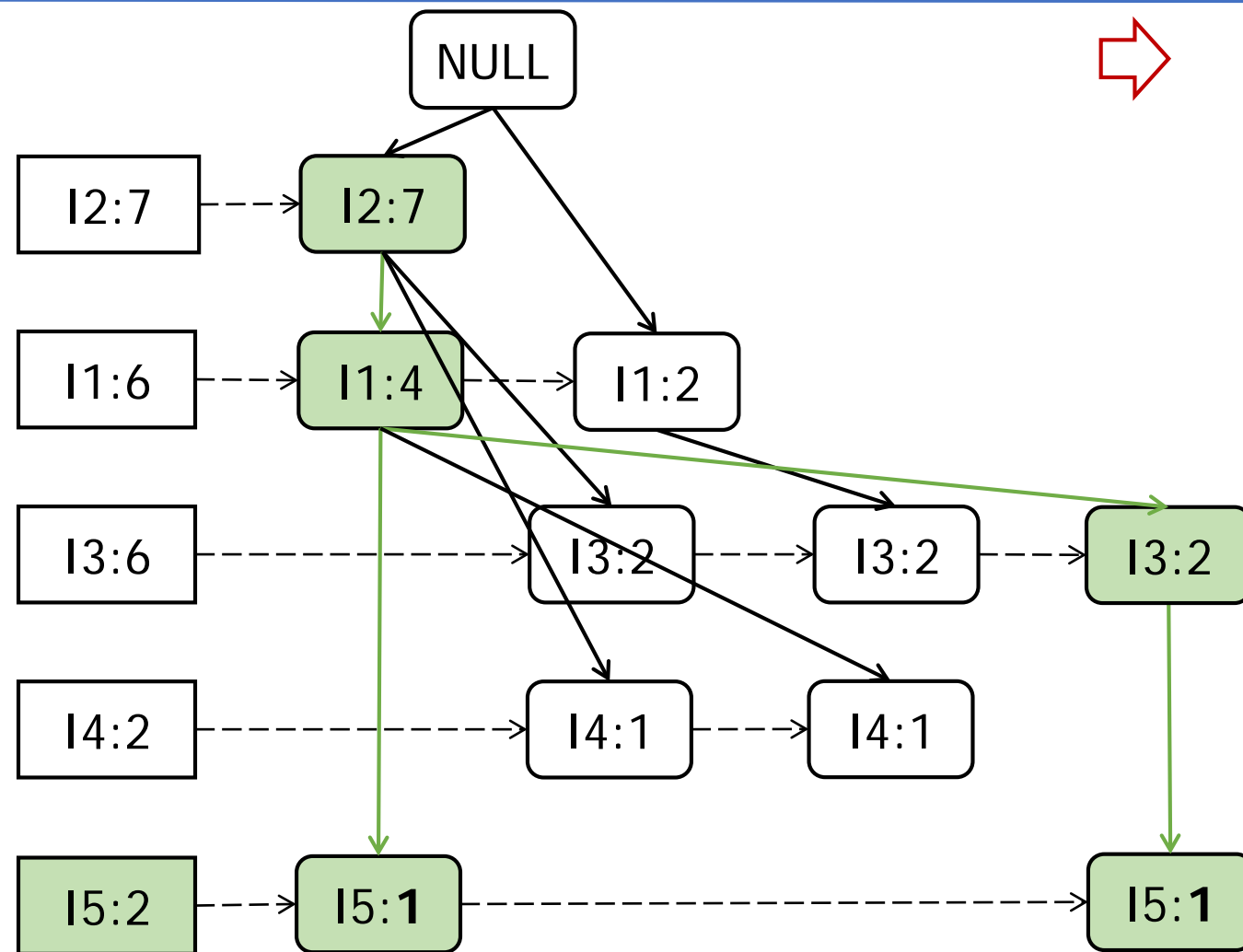
$D, \text{minsup} = 0.2$



TID	Items
10	I2, I1, I5
20	I2, I4
30	I2, I3
40	I2, I1, I4
50	I1, I3
60	I2, I3
70	I1, I3
80	I2, I1, I3, I5
90	I2, I1, I3
100	I6

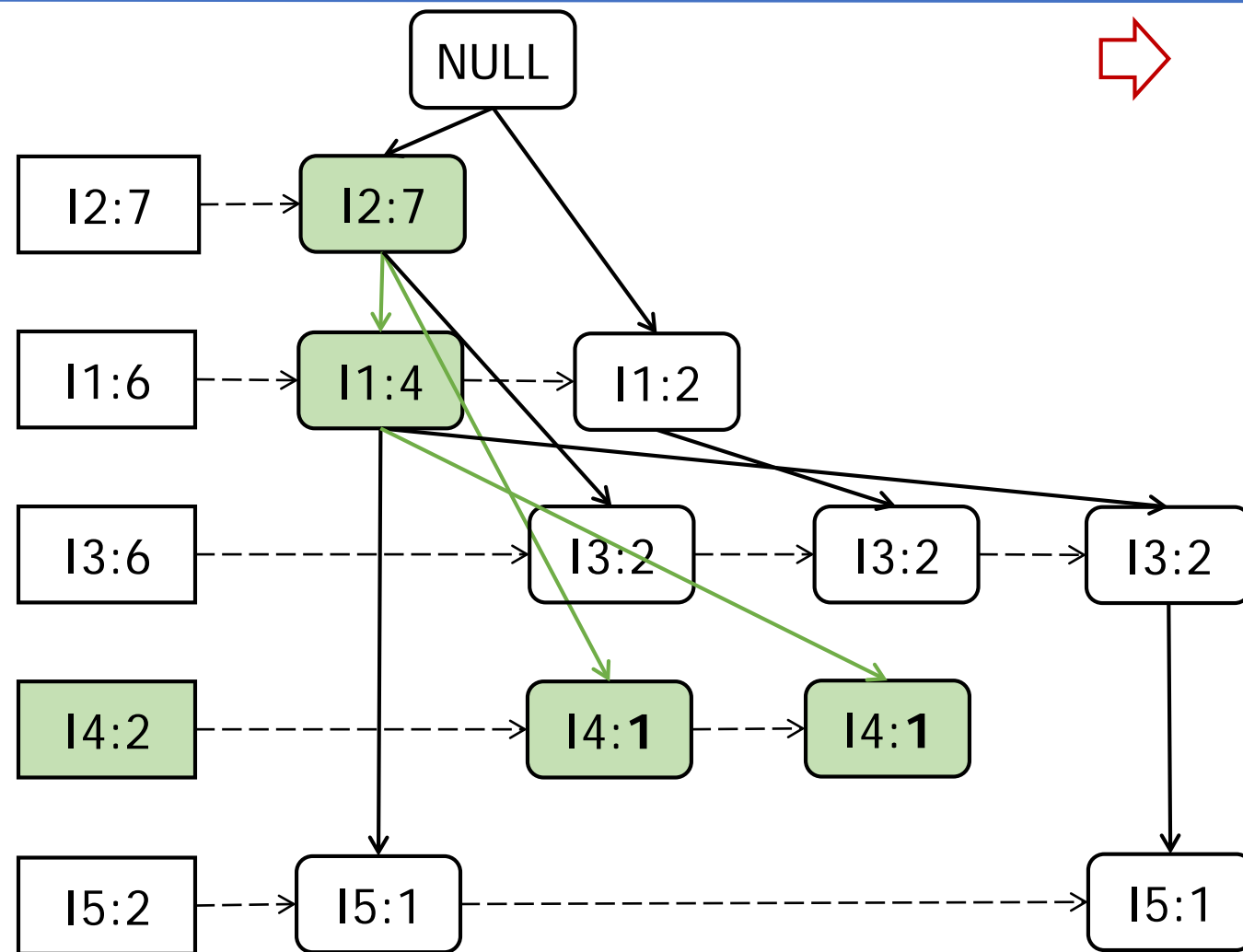


Обход FP-дерева: условные базисы



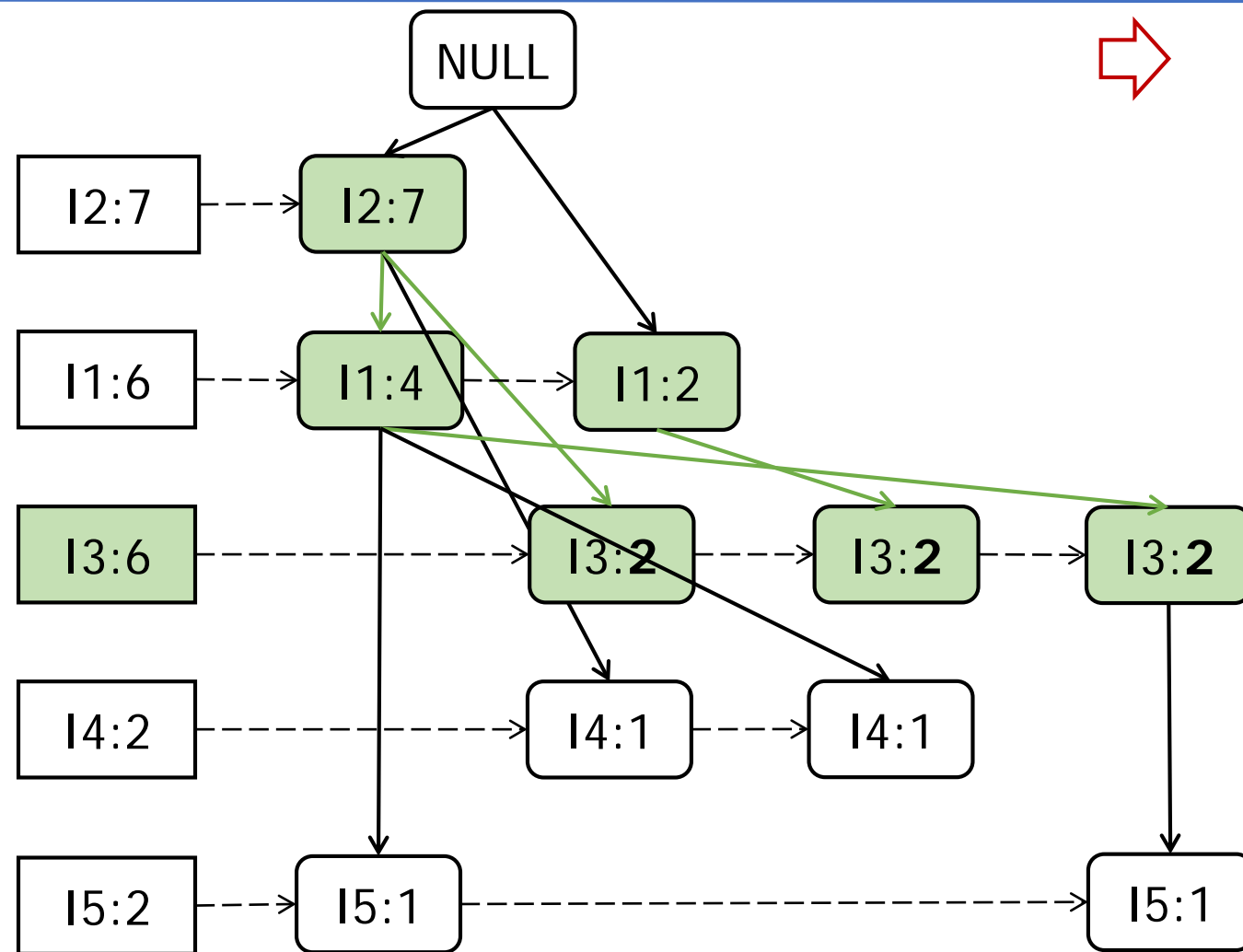
Item	Conditional pattern base
I2	∅
I1	
I3	
I4	
I5	{I2,I1:1}, {I2,I1,I3:1}

Обход FP-дерева: условные базисы



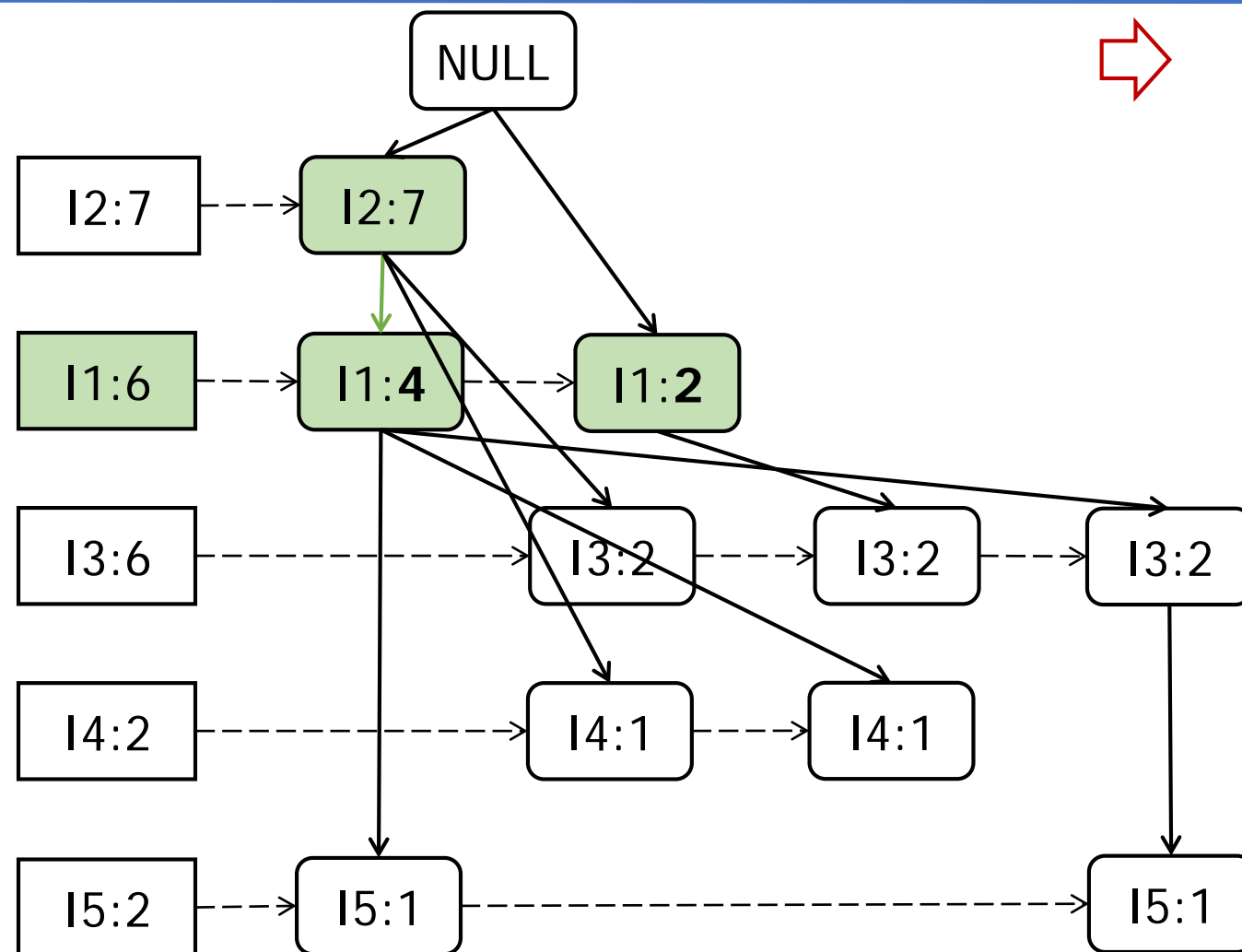
Item	Conditional pattern base
12	\emptyset
11	
13	
14	{12,11:1}, {12:1}
15	{12,11:1}, {12,11,13:1}

Обход FP-дерева: условные базисы



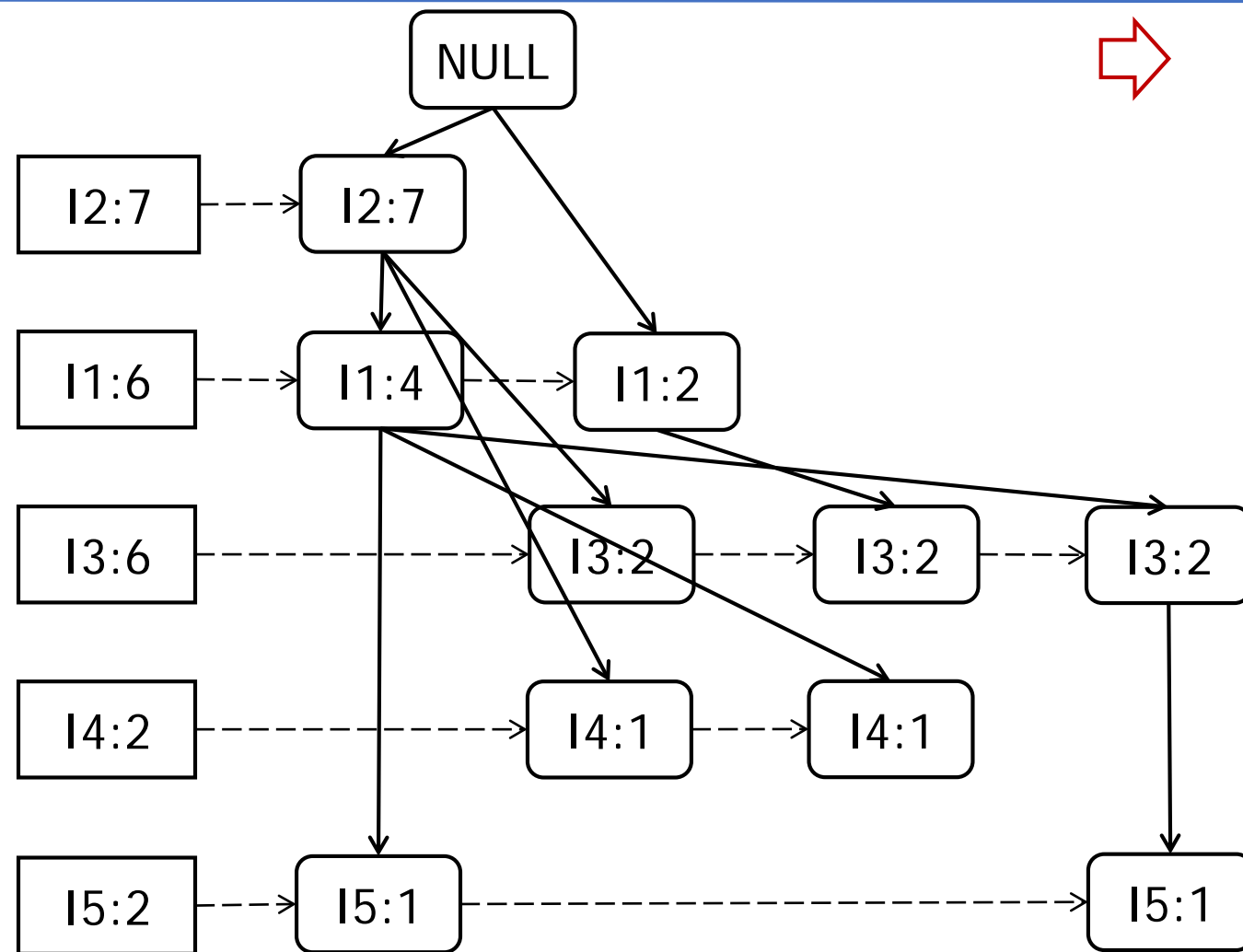
Item	Conditional pattern base
I2	\emptyset
I1	
I3	{I2, I1:2}, {I2:2}, {I1:2}
I4	{I2, I1:1}, {I2:1}
I5	{I2, I1:1}, {I2, I1, I3:1}

Обход FP-дерева: условные базисы



Item	Conditional pattern base
I2	\emptyset
I1	{I2:I4}
I3	{I2,I1:I2}, {I2:I2}, {I1:I2}
I4	{I2,I1:I1}, {I2:I1}
I5	{I2,I1:I1}, {I2,I1,I3:I1}

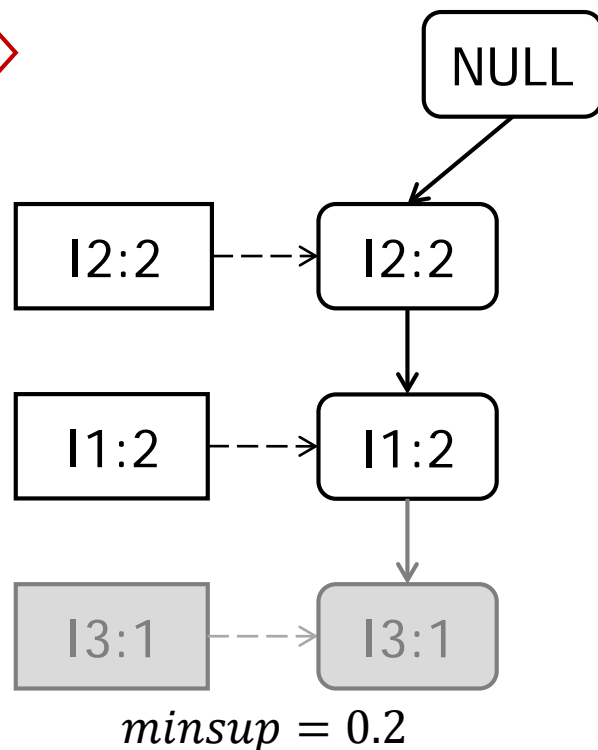
Обход FP-дерева: условные базисы



Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}

Обход FP-дерева: условные FP-деревья

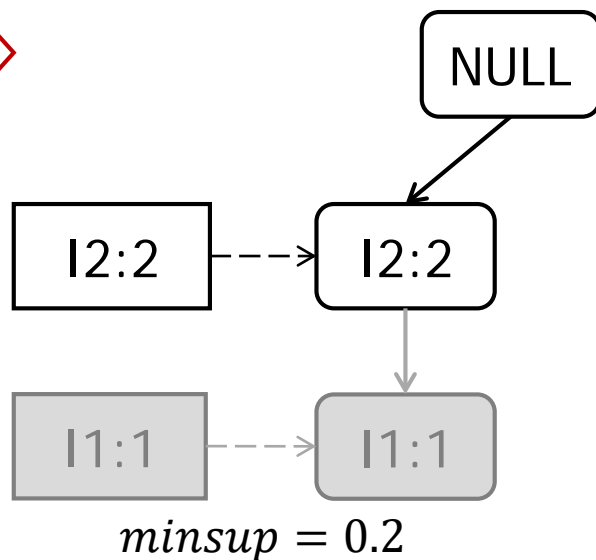
Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I2:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I2	\emptyset
I1	
I3	
I4	
I5	{I2:2, I1:2}

Обход FP-дерева: условные FP-деревья

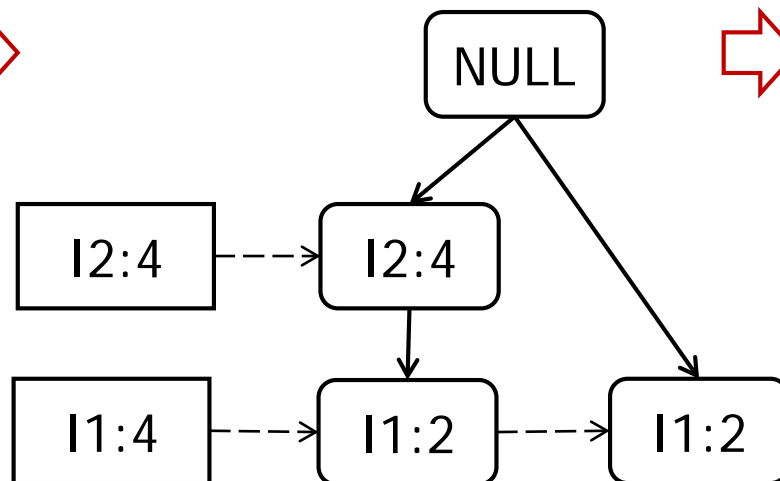
Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I2:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I2	\emptyset
I1	
I3	
I4	{I2:2}
I5	{I2:2}, {I1:2}

Обход FP-дерева: условные FP-деревья

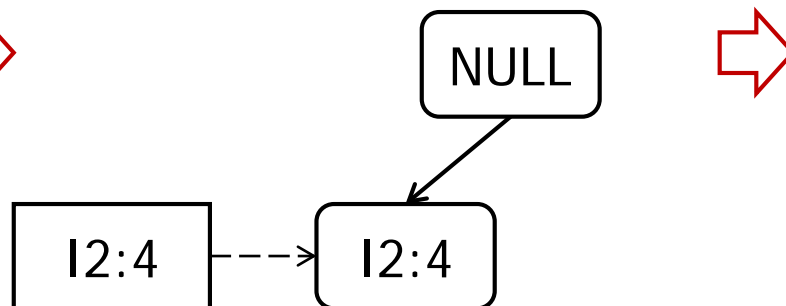
Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I2	\emptyset
I1	
I3	{I2:4,I1:2} {I1:2}
I4	{I2:2}
I5	{I2:2, I1:2}

Обход FP-дерева: условные FP-деревья


Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I2	\emptyset
I1	{I2:4}
I3	{I2:4,I1:2} {I1:2}
I4	{I2:2}
I5	{I2:2, I1:2}

Обход FP-дерева: условные FP-деревья

Item	Conditional pattern base
I2	\emptyset
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I2	\emptyset
I1	{I2:4}
I3	{I2:4,I1:2} {I1:2}
I4	{I2:2}
I5	{I2:2, I1:2}

Обход FP-дерева: частые наборы

Item	Conditional pattern base
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



Item	Conditional FP-tree
I1	{I2:4}
I3	{I2:4,I1:2} {I1:2}
I4	{I2:2}
I5	{I2:2, I1:2}



Prefix	Frequent Itemsets
I1	{I2,I1:4}
I3	{I2,I3:4}, {I1,I3:4}, {I2,I1,I3:2}
I4	{I2,I4:2}
I5	{I2,I5:2}, {I1,I5:2}, {I2,I1,I5:2}

Обход FP-дерева: частые наборы

Item	Conditional pattern base
I1	{I2:4}
I3	{I2,I1:2}, {I2:2}, {I1:2}
I4	{I2,I1:1}, {I2:1}
I5	{I2,I1:1}, {I2,I1,I3:1}



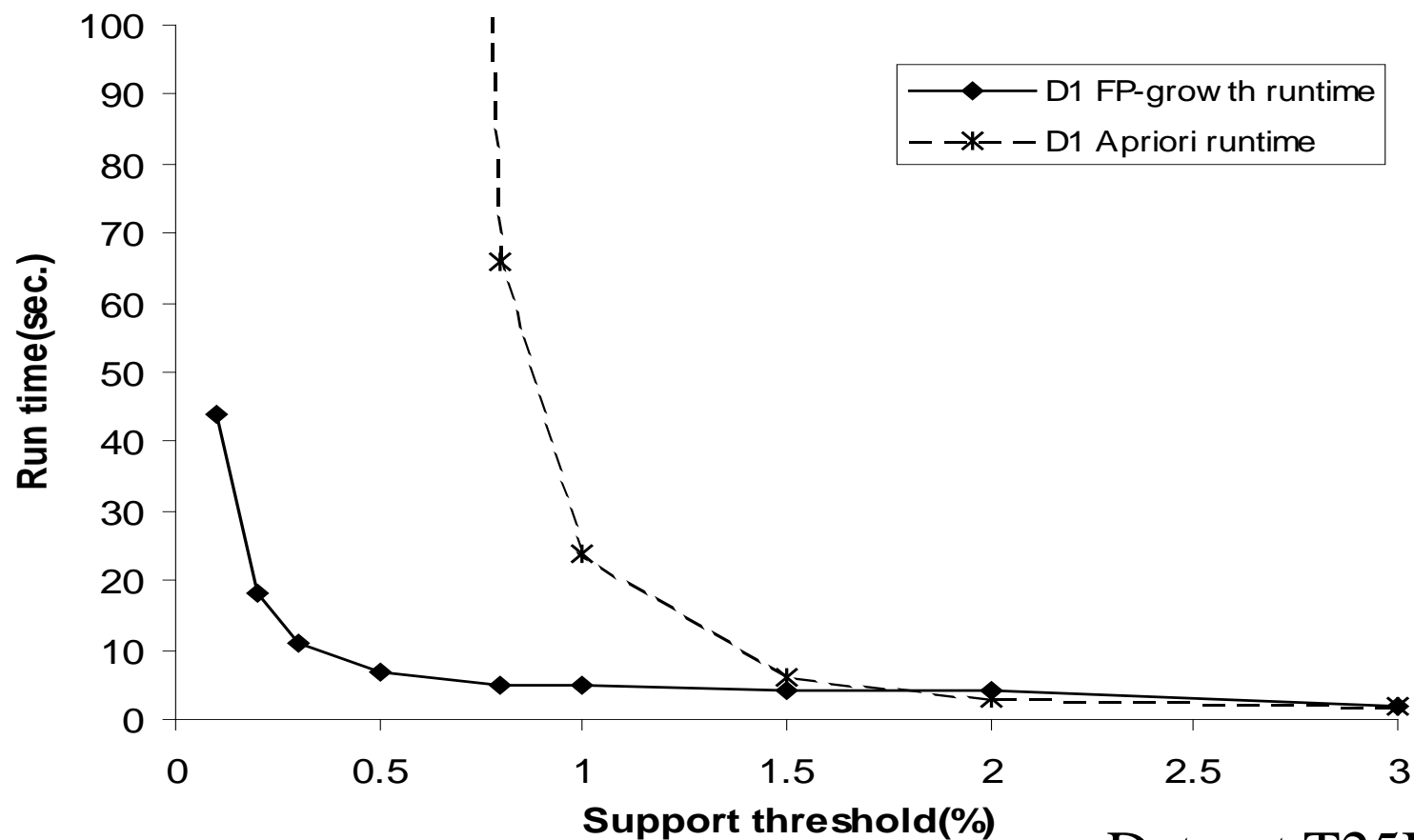
Item	Conditional FP-tree
I1	{I2:4}
I3	{I2:4,I1:2} {I1:2}
I4	{I2:2}
I5	{I2:2, I1:2}



Prefix	Frequent Itemsets
I1	{I2,I1:4}
I3	{I2,I3:4}, {I1,I3:4}, {I2,I1,I3:2}
I4	{I2,I4:2}
I5	{I2,I5:2}, {I1,I5:2}, {I2,I1,I5:2}

$$U \quad \mathcal{L}_1$$

FP-Growth vs. Apriori



Dataset T25I20D10K

Реализация алгоритмов поиска шаблонов

- Christian Borgelt, Prof. for Data Science at the Paris Lodron University of Salzburg

The screenshot shows a web browser displaying the page "Christian Borgelt's Web Pages" with the URL "borgelt.net/fpm.html". The page title is "Software for Frequent Pattern Mining". A navigation menu on the left includes "Home", "Publications", "Slides", "Teaching", "Software", "Frida", and "FPM". The "Software" section contains a table of software tools.

Name	Language	Description
FIMGUI	Java	Frequent Item Set Mining GUI and Viewer
ARuleGUI	Java	Association Rule Mining GUI and Viewer
Apriori	C	Frequent Item Set Mining (all, closed, maximal, generators) and Association Rule Induction
Eclat/LCM	C/Python	Frequent Item Set Mining (all, closed, maximal, generators) and Association Rule Induction
FPgrowth	C	Frequent Item Set Mining (all, closed, maximal, generators) and Association Rule Induction
RElim	C	Frequent Item Set Mining (all, closed, maximal, generators, fault-tolerant)
Sam	C	Frequent Item Set Mining (all, closed, maximal, generators, fault-tolerant)
SODIM	C	Frequent Item Set Mining (fault-tolerant)
IsTa	C	Frequent Item Set Mining (closed and maximal)
FPgrowth	C	Frequent Item Set Mining (closed and maximal)
PyFIM	C/Python	Frequent Item Set Mining for Python
JIM	C	Jaccard Item Set Mining / Cover Similarity
JaM	C/Java	Java API for Frequent Item Set Mining
FIMGUI	C/Java	Graphical User Interface for Frequent Item Set Mining
CoCoNAD	C	Continuous-time Closed Neuron Assembly Detection Frequent Pattern Mining in Point Processes
PyCoCo	C/Python	CoCoNAD for Python
CoCo4R	C/R	CoCoNAD for R
JNICoCo	C/Java	Java API for CoCoNAD
FIM4R	C/Java	Graphical User Interface for CoCoNAD + PSF + PSR
Seqwog	C	Frequent Sequence Mining
Sequola	C	Frequent Sequence Mining
MoSS	Java	Molecular Substructure Miner

License

On October 23, 2014, I decided to abandon the (L)GPL licenses and adopt the MIT license for my programs, in order to avoid problems some people see with using software that is licensed under the LGPL in other software (even though the LGPL actually permits use in proprietary programs, while the GPL does not). I hope to remove these and related problems by switching to the MIT license.

The transition to the new license will be accomplished with new versions that get published. The following applies:

For any version published on or after October 23, 2014:
(MIT license, or more precisely Expat License, to be found in the file `!t-1license.txt` in the directory `<prname>/doc` in the source package of the program, see also opensource.org and wikipedia.org)

© 1996-2020 Christian Borgelt

- IBM Synthetic Data Generator for Itemsets and Sequences

<https://github.com/zakimjz/IBMGenerator>

Литература

- Han J., Kamber M., Pei J. Data Mining: Concepts and Techniques. 3rd Edition. Morgan Kaufmann, 2011. 740 p. ISBN 978-0123814791
 - Chapter 6. Mining Frequent Patterns, Associations, and Correlations: Basic Concepts and Methods, pp. 243-278
- Tan P.-N., Steinbach M., Karpatne A., Kumar V. Introduction to Data Mining. 2nd Edition. Pearson, 2019. 839 p. ISBN-13: 978-0-13-312890-1
 - 5. Association Analysis: Basic Concepts and Algorithms, pp. 357-450

Отбрасывание заведомо неустойчивых шаблонов

- В общем случае достоверность не обладает свойством антимонотонности
- **Теорема.**

Пусть A, B – непустые наборы и $A \subseteq B$.

Если $\text{conf}(A \rightarrow B \setminus A) < \text{minconf}$, то

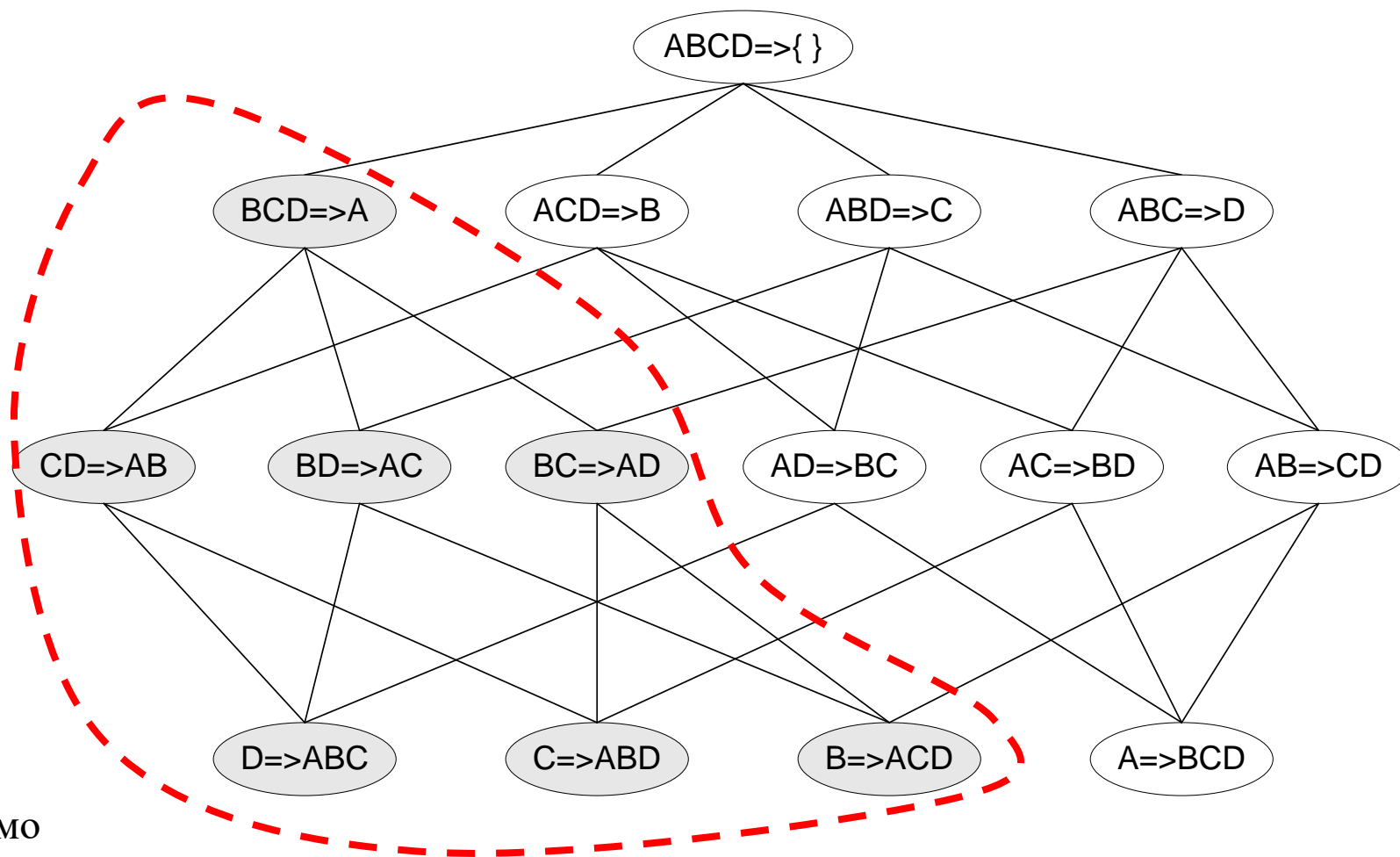
$\forall \tilde{A} \subseteq A \text{ conf}(\tilde{A} \rightarrow B \setminus \tilde{A}) < \text{minconf}$

Доказательство.

$$\text{conf}(A \rightarrow B \setminus A) = \frac{\text{sup}(B)}{\text{sup}(A)}, \text{ conf}(\tilde{A} \rightarrow B \setminus \tilde{A}) = \frac{\text{sup}(B)}{\text{sup}(\tilde{A})}.$$

$$\tilde{A} \subseteq A \Rightarrow \text{sup}(\tilde{A}) \geq \text{sup}(A).$$

Отбрасывание заведомо неустойчивых шаблонов



Заведомо
неустойчивые шаблоны