

Задача поиска шаблонов в данных

Все в мире повторяется.

Ф. Бэкон

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

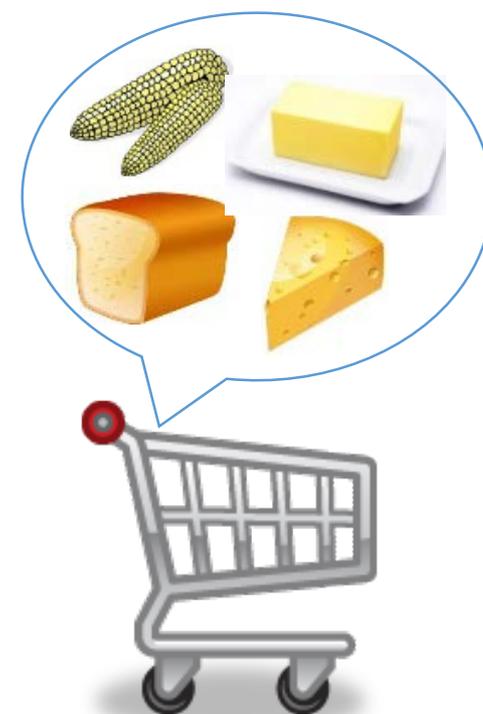
- Частые наборы и шаблоны
- Алгоритм Apriori поиска частых наборов
- Меры полезности шаблонов

Пример: анализ рыночной корзины

- Какие наборы товаров в супермаркете часто покупают совместно?



...



Пример: поиск побочных эффектов лекарств

- Какие симптомы у пациентов часто встречаются совместно с принимаемыми лекарствами?



жар



тошнота
рвота



тошнота
жар



тошнота
жар
сыпь



...

Пример: выявление заимствований

- Какие группы авторов часто используют одинаковые конструкции?

Диплом_Бендер

Курсовая_Михельсон

Диссертация_Берлага

в настоящее время
актуальной является
проблема плагиата

Диплом_Бендер

Реферат_Полыхаев

Отчет_Корейко

одним из эффективных
решений данной
проблемы является



Диссертация_Берлага

Диплом_Бендер

Курсовая_Михельсон

Отчет_Корейко

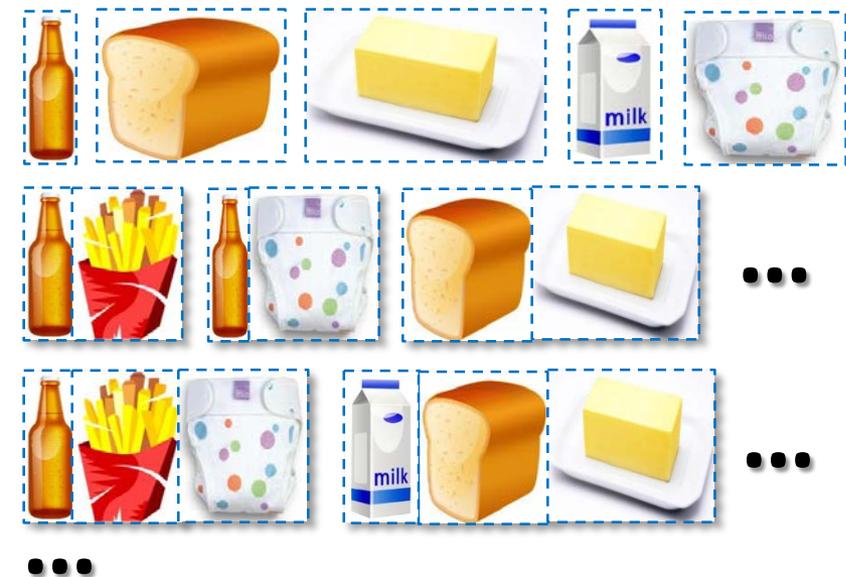
в итоге исследования
получены следующие
основные результаты

Частый набор vs. шаблон



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

Шаблоны (ассоциативные правила)
антецедент \rightarrow *консеквент*



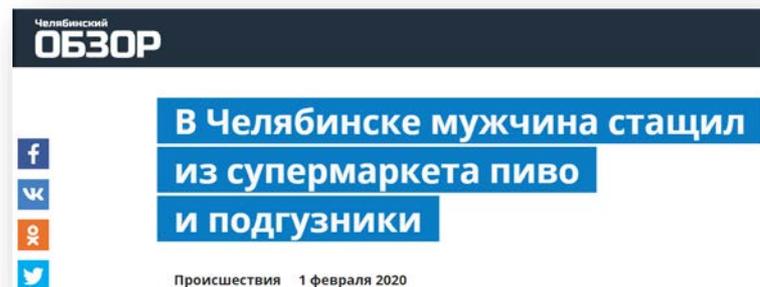
ЕСЛИ		ТО	

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 (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 



Сотрудники Росгвардии в Челябинске задержали мужчину, подозреваемого в краже из супермаркета. «Уловом» грабителя стали три упаковки подгузников и три банки пива.

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

Объекты, наборы, транзакции

- Объекты («товары»):

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

- Наборы:

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

- Транзакции:

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

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

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

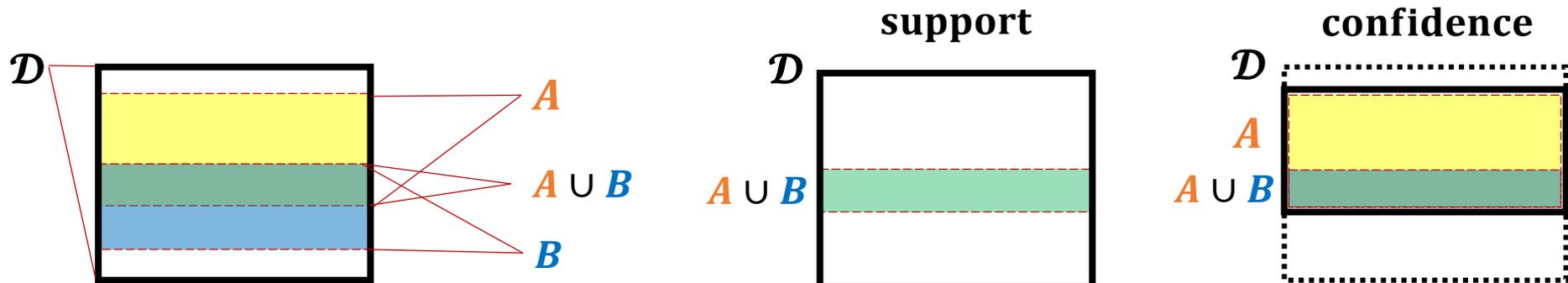
Шаблон, поддержка, достоверность

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

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

- Достоверность

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



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

$minsup = 0.6$

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

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

$minsup = 0.1$

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

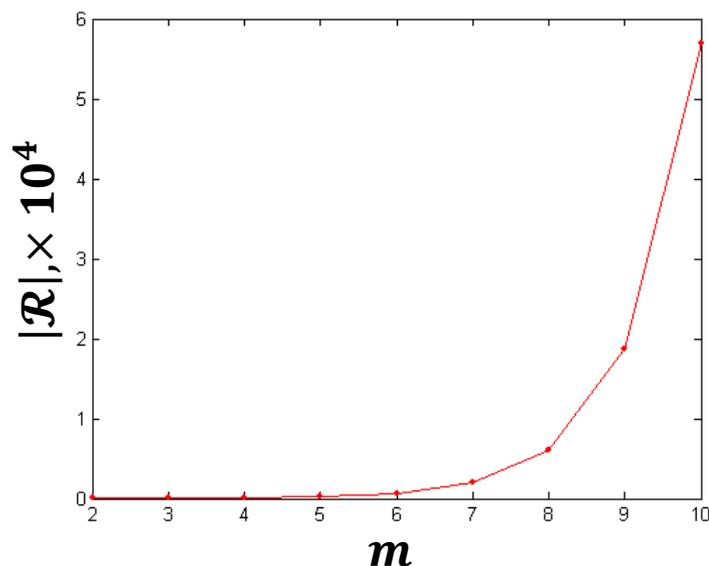
Устойчивый шаблон

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

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 шаблонов: алгоритм

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

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

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

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

$$\text{conf} := \frac{\text{sup}(I)}{\text{sup}(S)}$$

if $\text{conf} \geq \text{minconf}$ **then**

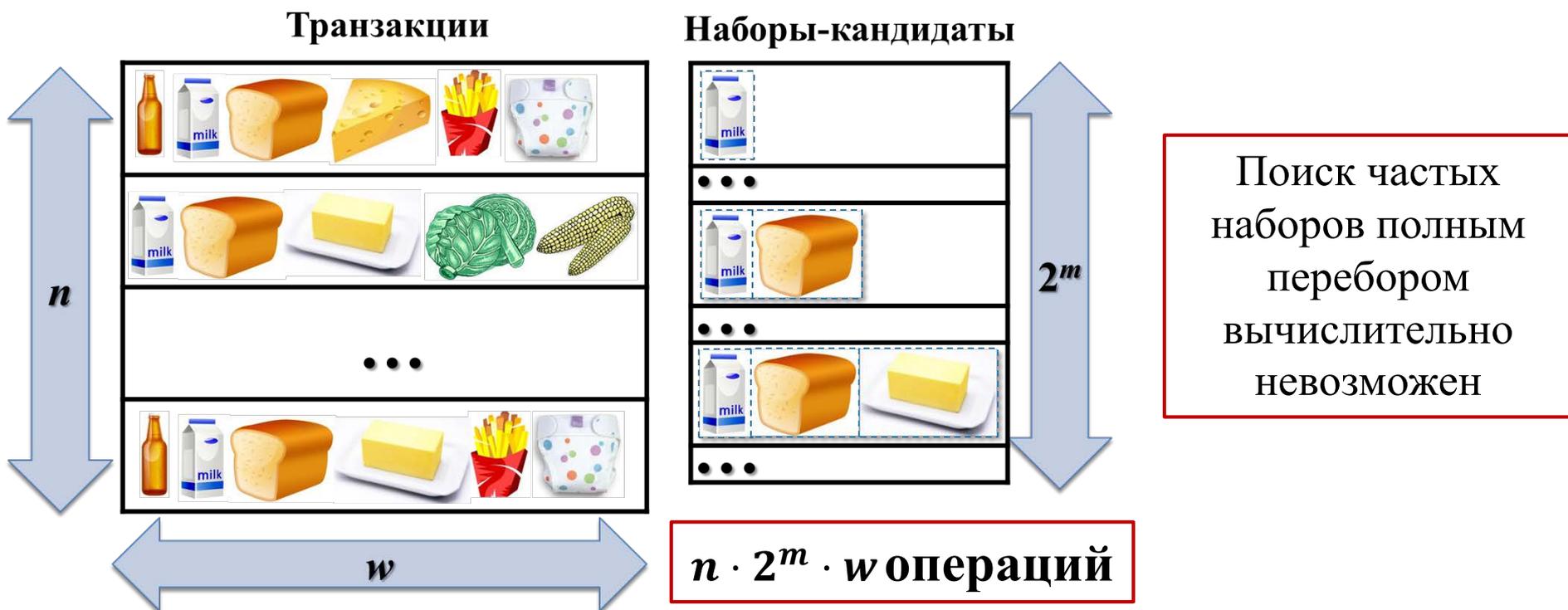
$$\text{pattern} := "S \rightarrow I \setminus S"$$

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

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

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 $\forall c \in \mathcal{C}: \text{sup}(c) < \text{minsup}$



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

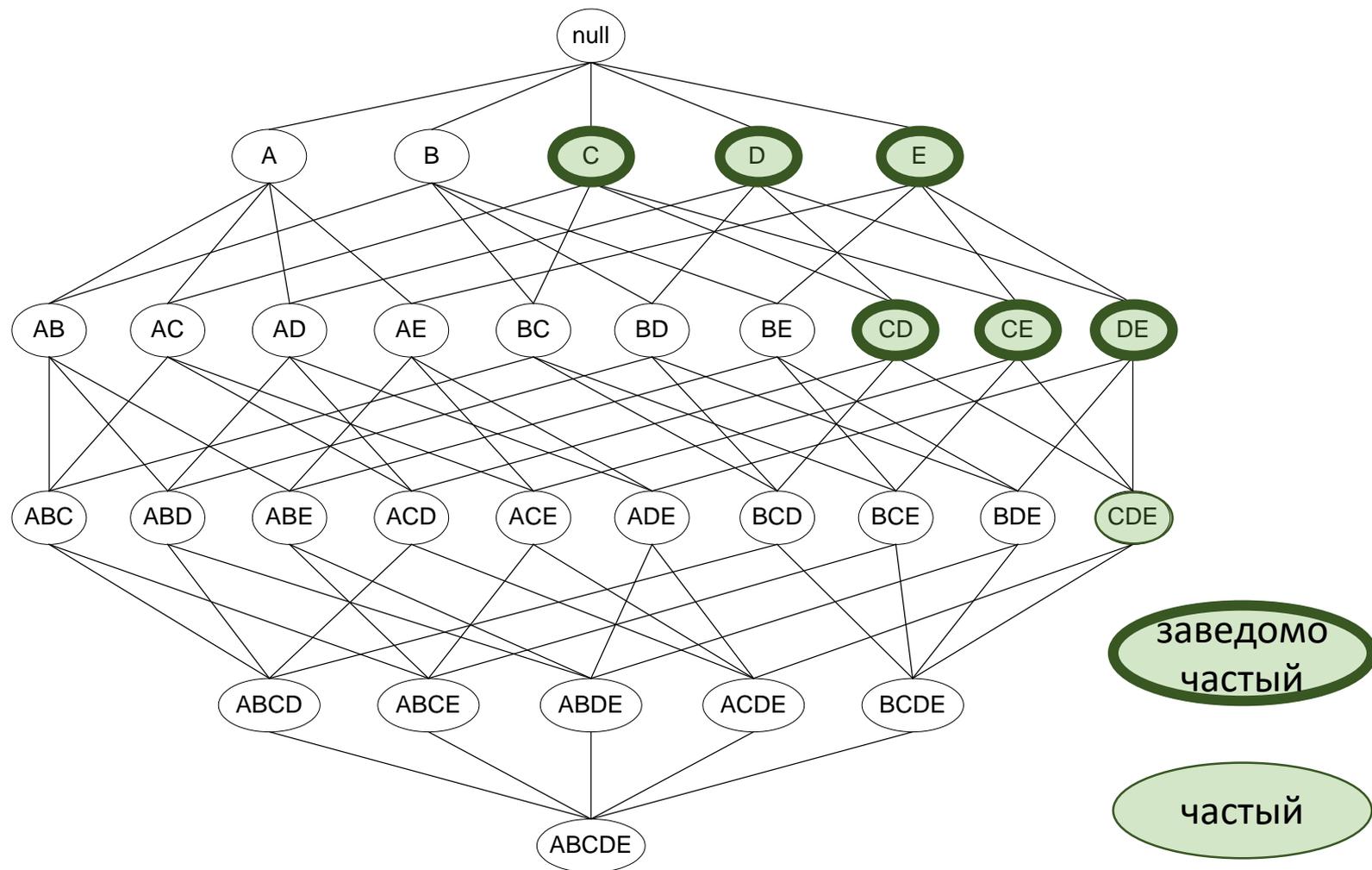
- любое подмножество частого набора является частым набором

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

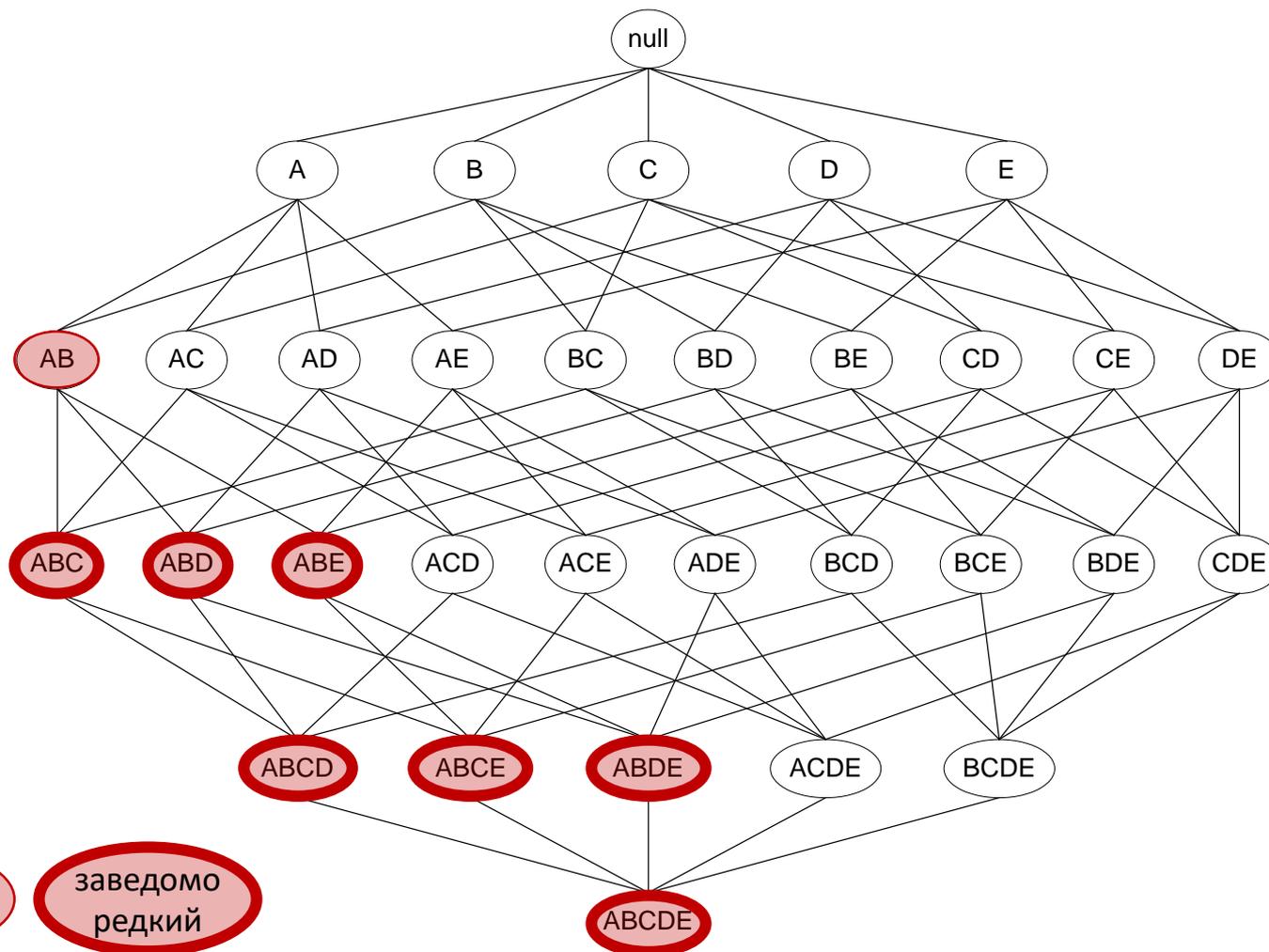
- если некое подмножество набора является редким набором, то набор является редким

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

Принцип Априори



Принцип Априори



редкий

заведомо редкий

ABCDE

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 в C_k по принципу Априори
 - если $I \subseteq C_k$, но $I \notin \mathcal{L}_{k-1}$, то отбросить I

Генерация наборов-кандидатов: пример

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

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
 $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE, \cancel{DE}\}$
- $\mathcal{L}_2 = \{AB, AC, AD, AE, BC, BD, BE, CD, CE\}$
 Самосоединение
 $C_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$
 Вычисление поддержки
 $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, ACE, ADE, BCD, BCE, BDE, CDE\}$
- $\mathcal{L}_3 = \{ABC, ABD, ABE, ACD, BCD, BDE, CDE\}$
 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



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

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	



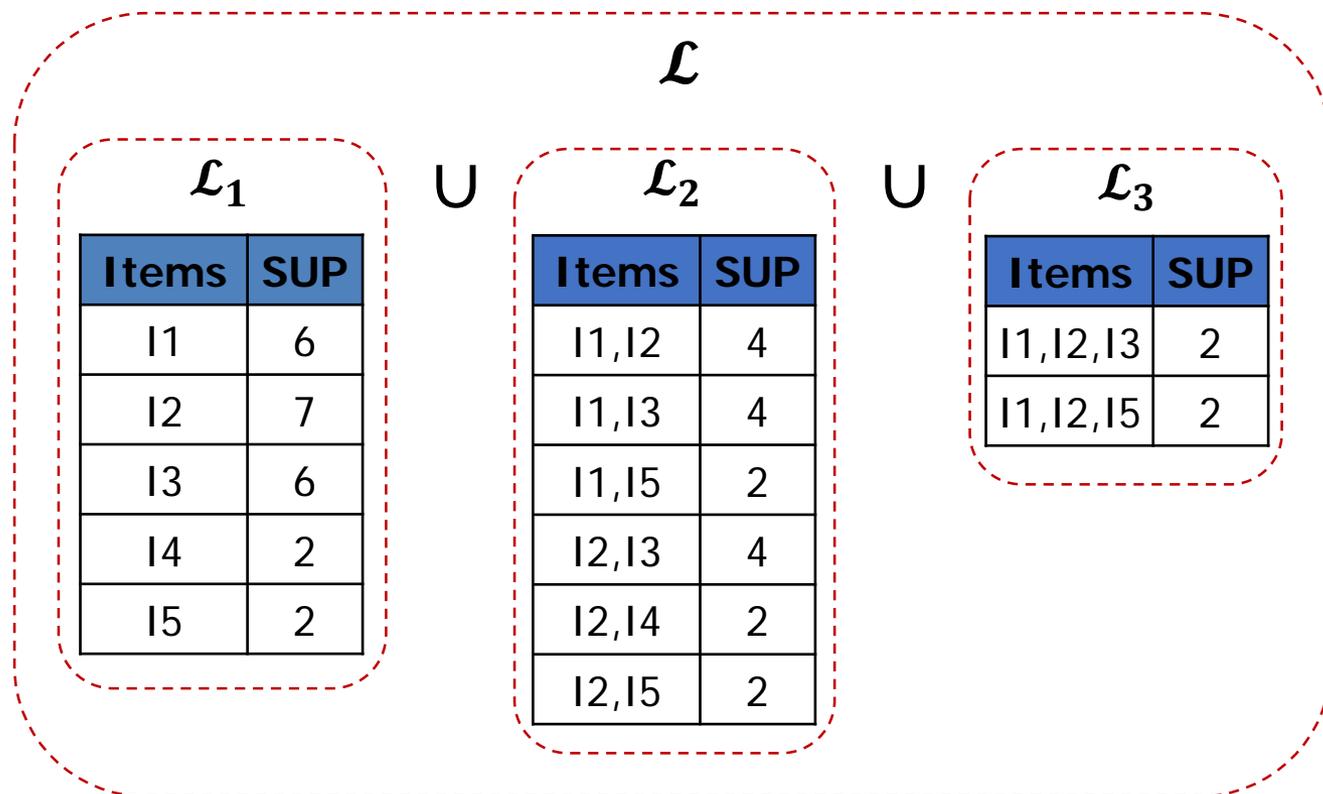
$\mathcal{L}_4 = \emptyset$

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



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



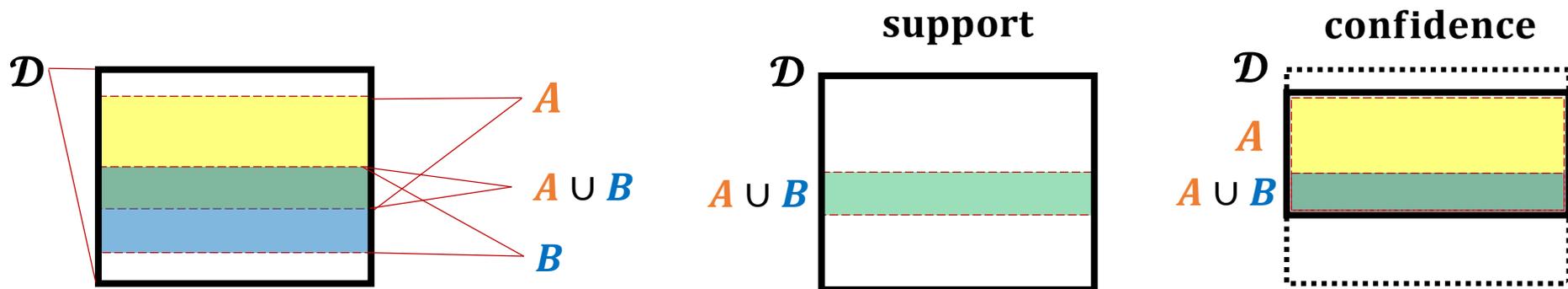
Items	SUP
I1, I2, I3	2
I1, I2, I5	2



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

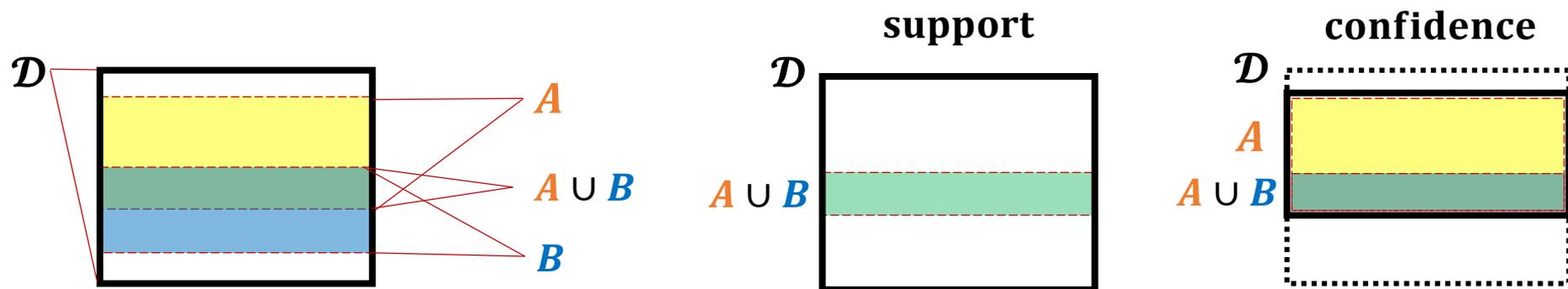
Полезность шаблонов

№	Шаблон	sup	conf	Польза
1	хлеб → масло	👍	👍	?
2	хлеб → авокадо	👍	👍	?
3	молоко → тунец	👍	👎	?
4	масло → стир. порошок	👎	👎	?
5	водка → икра красная	👎	👍	?
6	водка → икра баклажанная	👎	👍	?



Полезность шаблонов

№	Правило	sup	conf	Польза
1	хлеб → масло	👍	👍	👎
2	хлеб → авокадо	👍	👍	👍
3	молоко → тунец	👍	👎	👎
4	масло → стир. порошок	👎	👎	👎
5	водка → икра красная	👎	👍	👎
6	водка → икра баклажанная	👎	👍	👍



Пример: шаблон устойчив, но бесполезен

- Таблица сопряженности

	<i>coffee</i>	\overline{coffee}	Σ
<i>tea</i>	4000	3500	7500
\overline{tea}	2000	500	2500
Σ	6000	4000	10000

- $coffee \rightarrow tea$ [$sup = \frac{4000}{10000} = 0.4$, $conf = \frac{4000}{6000} = 0.66$]
- $P(tea) = \frac{7500}{10000} = 0.75 > P(tea|coffee) = conf$
- Проще угадать наличие *tea*, чем предсказать это с помощью шаблона

Пример: шаблон неустойчив, но полезен

- Таблица сопряженности

	<i>honey</i>	\overline{honey}	Σ
<i>tea</i>	100	100	200
\overline{tea}	20	780	800
Σ	120	880	1000

- $tea \rightarrow honey$ [$sup = \frac{100}{1000} = 0.1$, $conf = \frac{100}{200} = 0.5$]
- $P(honey) = \frac{120}{1000} = 0.12 < P(honey|tea) = conf$
- $P(\overline{tea} \cup honey) = \frac{20}{1000} = 0.025$
- Отбрасывание шаблона не учитывает информацию о предпочтении *honey* при условии *tea*

Мера lift – дополнительная к sup и conf

- $\text{lift}(A, B) = \frac{P(A \cup B)}{P(A) \cdot P(B)}$
- Оценка
 - $\text{lift}(A, B) = 1$: A и B независимы
 - $\text{lift}(A, B) < 1$: A и B имеют отрицательную корреляцию (наличие A подразумевает отсутствие B)
 - $\text{lift}(A, B) > 1$: A и B имеют положительную корреляцию (наличие A подразумевает наличие B)
- $\text{lift}(A \rightarrow B) = \frac{P(B|A)}{P(B)} = \frac{\text{conf}(A \rightarrow B)}{\text{sup}(B)} = \frac{\text{sup}(A \rightarrow B)}{\text{sup}(A) \cdot \text{sup}(B)}$

Мера lift

- Пример:

- Таблица сопряженности

	A	\bar{A}	Σ
B	4000	3500	7500
\bar{B}	2000	500	2500
Σ	6000	4000	10000

- $A \rightarrow B$ [$sup = 0.4, conf = 0.66$]

- $lift(A, B) = \frac{P(A \cup B)}{P(A) \cdot P(B)} = \frac{0.4}{0.6 \cdot 0.75} = 0.89$

- A и B имеют отрицательную корреляцию

- $A \rightarrow B$ малополезно

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