

Ассоциативные правила

Виктор Китов

v.v.kitov@yandex.ru

Решаемая задача

Транзакции в магазине

<i>TID</i>	Items
1	{Bread, Milk}
2	{Bread, Diapers, Beer, Eggs}
3	{Milk, Diapers, Beer, Cola}
4	{Bread, Milk, Diapers, Beer}
5	{Bread, Milk, Diapers, Cola}

Пример правила:

$\{\text{Diapers}\} \longrightarrow \{\text{Beer}\}.$

Применения: анализ рыночных корзин (market basket analysis), медицинская диагностика, анализ данных в др. областях.

● Проблемы:

- вычислительная сложность перебора всех правил.
- статистическая значимость правил
- совместная встречаемость не означает причинно следственной связи (correlation vs. causation)
 - (например, товары лежат рядом на одной полке)

Пример генерации правил в политике

Данные по голосованиям ("товары"):

- | | |
|-------------------------------------|---------------------------------------|
| 1. Republican | 18. aid to Nicaragua = no |
| 2. Democrat | 19. MX-missile = yes |
| 3. handicapped-infants = yes | 20. MX-missile = no |
| 4. handicapped-infants = no | 21. immigration = yes |
| 5. water project cost sharing = yes | 22. immigration = no |
| 6. water project cost sharing = no | 23. synfuel corporation cutback = yes |
| 7. budget-resolution = yes | 24. synfuel corporation cutback = no |
| 8. budget-resolution = no | 25. education spending = yes |
| 9. physician fee freeze = yes | 26. education spending = no |

Определение принадлежности к партии (предпочтения)

Association Rule	Confidence
{budget resolution = no, MX-missile=no, aid to El Salvador = yes } → {Republican}	91.0%
{budget resolution = yes, MX-missile=yes, aid to El Salvador = no } → {Democrat}	97.5%
{crime = yes, right-to-sue = yes, physician fee freeze = yes} → {Republican}	93.5%
{crime = no, right-to-sue = no, physician fee freeze = no} → {Democrat}	100%

Бинарное представление

- $I = \{i_1, i_2, \dots, i_d\}$ - все товары, ищем наборы из этих товаров
- $T = \{t_1, t_2, \dots, t_N\}$ - все транзакции, t_i - подмножество I
- Поддержка (support)

$$\sigma(X) = |\{t_i | X \subseteq t_i, t_i \in T\}|$$

Ассоциативное правило

- Ассоциативное правило $X \rightarrow Y$, где X, Y - наборы товаров, $X \cap Y = \emptyset$

Свойства ассоциативного правила:

$$\text{Support, } s(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{N};$$

$$\text{Confidence, } c(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}.$$

- $X \cup Y$ - наборы X и Y одновременно встретились (пересечение).
- Примеры расчетов.
- Правило с низкой поддержкой может появиться случайно, нет смысла делать промоакции для редких правил.
- Уверенность измеряет, насколько правило надежно.

Мера lift

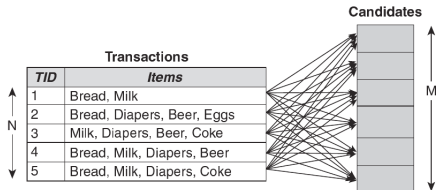
Мера lift показывает, полезнее ли правило случайного угадывания:

$$\begin{aligned} \text{lift}(A \rightarrow B) &= \frac{\text{support}(A \cup B)}{\text{support}(A) \cdot \text{support}(B)} = \\ &= \frac{|\{A \cup B \subseteq x \mid x \in X_{\text{train}}\}| \cdot |X_{\text{train}}|}{|\{A \subseteq x \mid x \in X_{\text{train}}\}| \cdot |\{B \subseteq x \mid x \in X_{\text{train}}\}|} \end{aligned}$$

Нахождение правил

- Нахождение правил состоит из 2х этапов:
 - 1 генерация частых наборов ($support \geq minsup$)
 - 2 генерация уверенных правил по наборам ($confidence \geq minconf$)
- Первая задача вычислительно сложнее.
 - полный перебор: сложность $O(N(2^k - 1)w)$ для наборов длины k , #транзакций N и средней длины транзакции w .
 - решения: $\downarrow k$ (Apriori), \downarrow число сравнений с транзакциями (FP-growth)

Полный перебор ($M = 2^k - 1$)



Apriori

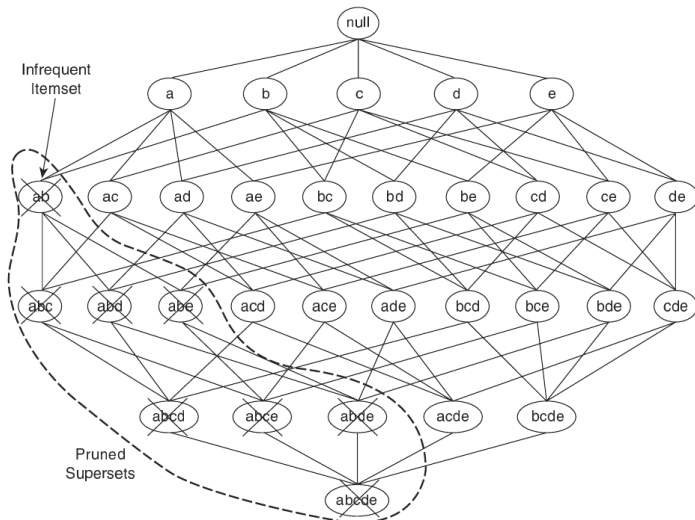
- У следующих правил - одинаковая поддержка:

$$\begin{aligned} \{Beer, Diapers\} &\longrightarrow \{Milk\}, & \{Beer, Milk\} &\longrightarrow \{Diapers\}, \\ \{Diapers, Milk\} &\longrightarrow \{Beer\}, & \{Beer\} &\longrightarrow \{Diapers, Milk\}, \\ \{Milk\} &\longrightarrow \{Beer, Diapers\}, & \{Diapers\} &\longrightarrow \{Beer, Milk\}. \end{aligned}$$

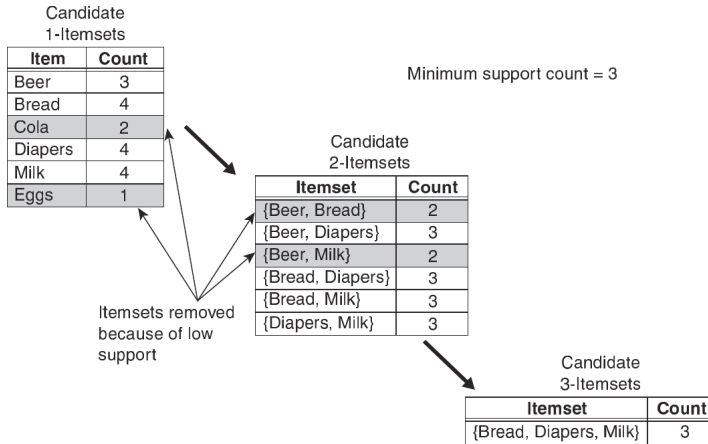
- Идея алгоритма Apriori - если набор частый, то все его поднаборы тоже частые.
- Математически это свойство антимонотонности поддержки:

$$X' \subset X \Rightarrow \sigma(X') \geq \sigma(X)$$

Идея Apriori



Поиск, использующий антимонотонность



Экономия вычислений

- Экономия перебора:

$$\binom{6}{1} + \binom{6}{2} + \binom{6}{3} \rightarrow \binom{6}{1} + \binom{4}{2} + 1$$

- Алгоритм:

- проход наборам одного товара, выбор частых F_1
- для $k = 1, 2, 3, \dots$ пока $F_k \neq \emptyset$:
 - генерация F_{k+1} комбинацией $\{F_k\}_k$

Генерация F_k

- Не должны генерироваться кандидаты, содержащие нечастые k поднаборы.
- Перебор кандидатов должен быть полным (среди потенциально подходящих).
- Сгенерированные кандидаты не должны дублироваться.
 - $\{a, b, c, d\} = \{a, b\} + \{c, d\} = \{a, b, c\} + \{d\} = \{a\} + \{b, c, d\} = \dots$
- Подходы генерации:
 - полный перебор, сложность C_k^d , $d = \# \text{товаров}$, k -длина набора

Подходы генерации

- Подходы генерации

- $F_{k-1} \times F_1$: комбинация всех F_{k-1} и F_1 , сложность $O(|F_{k-1}| \times |F_1|)$
- могут генерироваться повторения: $\{\text{Bread}, \text{Diapers}\} + \{\text{Milk}\} = \{\text{Bread}, \text{Milk}\} + \{\text{Diapers}\} = \dots$
 - решение: объединять только наборы по возрастанию в лексикографическом порядке:
 $\{\text{Bread}, \text{Diapers}\} + \{\text{Milk}\}, \{\text{Bread}, \text{Milk}\} + \{\text{Diapers}\}$
 - $F_{k-1} \times F_{k-1}$: объединяются только при условии

$$a_i = b_i, i = 1, 2, \dots, k-2; a_{k-1} \neq b_{k-1}$$

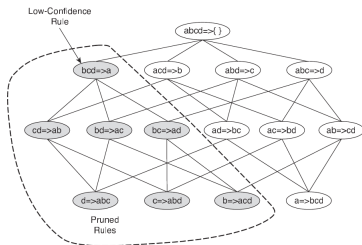
- перебор полный и \uparrow эффективность от лексикографического упорядочивания.

Генерация правил

- Для набора длины k существует $2^k - 2$ правил
 - по набору X генерируем $Y \rightarrow X - Y$
 - X -частый, значит и поднаборы $Y, X - Y$ -частые.
 - игнорируем $\emptyset \rightarrow Y$ и $Y \rightarrow \emptyset$
- Оптимизация перебора правил: если $X \rightarrow Y - X$ малой уверенности (confidence), то любое $X' \rightarrow Y - X', X' \subset X$ - тоже малой уверенности, т.к. $\sigma(X') \geq \sigma(X)$
 - $conf(X \rightarrow Y - X) = \frac{\sigma(Y)}{\sigma(X)} \geq \frac{\sigma(Y)}{\sigma(X')} = conf(X' \rightarrow Y - X')$
 - т.е. если $conf(bcd \Rightarrow a)$ мало, то $conf(cd \Rightarrow ab)$ еще меньше.

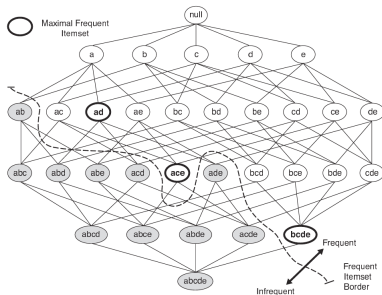
Оптимизация перебора правил в Apriori

Использование принципа для оптимизации перебора правил:



Компактное представление частых наборов

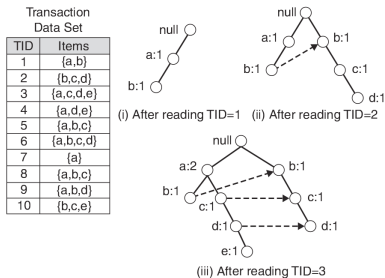
- Для компактного представления частых наборов достаточно знать только нерасширяемые частые наборы (maximal frequent itemsets).
- Все поднаборы нерасширяемые частые наборов - частые. Но теряется информация о поддержке.



FP-growth

FP-growth алгоритм использует структуру данных для компактного представления наборов и их поддержек.

дерево FP-growth:



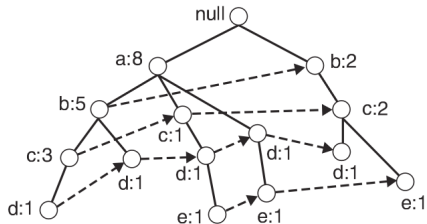
- Узлы - частые товары, упорядоченные по \downarrow поддержки.
- Узлы, отвечающие одинаковым товарам соединены указателем (уровень дерева)

Полное FP-дерево

Полное FP-дерево

Transaction
Data Set

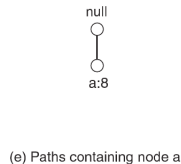
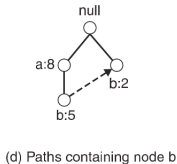
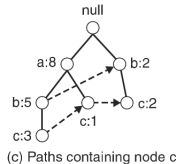
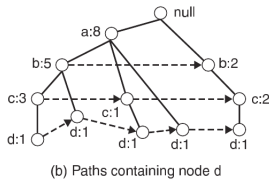
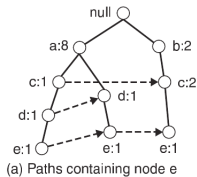
TID	Items
1	{a,b}
2	{b,c,d}
3	{a,c,d,e}
4	{a,d,e}
5	{a,b,c}
6	{a,b,c,d}
7	{a}
8	{a,b,c}
9	{a,b,d}
10	{b,c,e}



(iv) After reading TID=10

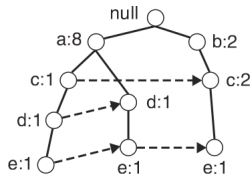
Поиск частых наборов

Можем найти частые наборы, оканчивающиеся на заданный суффикс, например "e":

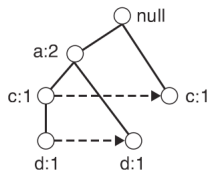


Извлеченные наборы

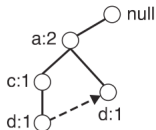
- Находим поддержку $\{e\}$: сумируем поддержки всех путей, заканчивающихся на e , получим 3.



(a) Prefix paths ending in e



(b) Conditional FP-tree for e



(c) Prefix paths ending in de



(d) Conditional FP-tree for de

Извлеченные наборы

- Т.к. $\{e\}$ - частый, то ищем частые поднаборы с расширенным суффиксом: de,ce,be,ae.
- Рассмотрим de. Строим условное FP-дерево - при условии окончания на e.
 - 1 Обновляем счетчики: например, путь $null \rightarrow b:2 \rightarrow c:2 \rightarrow e:1$ содержит $\{b,c\}$. Т.к. условие=окончание на e, то пересчитываем поддержку:

$$null \rightarrow b:1 \rightarrow c:1 \rightarrow e:1$$

- 2 некоторые товары могут перестать быть частыми в условном дереве, например "b", т.к. $\#\{b,e\}=1$.
- 3 убираем узлы с "e" (т.к. счетчики выше уже учитывают наличие "e" в конце)

Заключение

- Ассоциативные правила - эффективный метод извлечения интерпретируемых зависимостей в дискретных данных.
 - анализ покупательских корзин, политика, экология и др.
- Нужно быстро
 - искать частотные наборы (Apriori, FP-growth)
 - правила (Apriori)
- FP-growth быстрее Apriori.
- Развернутый обзор темы.