

link to the user is bucket



support is count of buckets

数值: support confidence interest

below confidence  $\rightarrow$  指的是 threshold.

$$\begin{array}{ccc} m \rightarrow b & \frac{m \vee b}{m} & \frac{cm \ b}{cm} \\ (m \rightarrow b) & & \end{array}$$

1. 购物篮模型: 描述一种多对多关系. 其中一类对象是 item, 另一类是 basket, 也叫 transaction <sup>交易</sup>

item: super market 所有 goods 的集合, 很大.

basket: 包含一部分 goods, 是 item 的子集, 一般数量比较多.

2. 几个概念:

→ 通常用分数表示, 即 包含 I 的 basket 数目 / 总数目.

① 支持度 support: I 是一个序列, 包含 I 的 basket 的数目为 support. 规则中所有商品

② 频繁项集: I 的支持度  $\geq$  support threshold. <sup>事务比例</sup>

③ association rules:  $\{i_1, i_2, \dots, i_k\} \rightarrow j$  means: 如果包含  $i_1, i_2, \dots, i_k$ , 则很可能包含 j.

④ 置信度  $\text{conf}(I \rightarrow j)$ : association rules 的置信度

$$\text{conf}(I \rightarrow j) = \frac{\text{support}(I \cup j)}{\text{support}(I)}$$

前提成立时, 结论也成立的. 可能性

可靠

⑤ 兴趣度 (interest): association rules 防止热门产品影响表达式所引入.

$$\text{interest} = \text{conf}(I \rightarrow j) - \text{pr}[j] \rightarrow \text{包含 } j \text{ 的 basket 比例.}$$

$\{a, b, c, d\}$  是频繁的, 那么  $\{a, b, c\}$  则一定也是频繁的.

if  $A, B, C \rightarrow D$  的  $\text{conf}$  小于  $x$ , 则  $A, B \rightarrow C, D$  一定也小于.

⑥ maximal frequent itemsets:  $\{a, b, c\}$  是 frequent, 但任意加一个  $\{a, b, c, x\}$  就不再 frequent.

⑦ closed itemsets:  $\{a, b, c\}$  和任意  $x$  的  $\{a, b, c, x\}$  有不同的 count.

用 the number of passes an algorithm makes (读取数据的遍数) 计算  $\text{count}$ .

⑧ lift: 提升度:  $\text{lift}(A \rightarrow B) = \frac{\text{support}(A \cup B)}{\text{support}(A) * \text{support}(B)}.$

>1时: A与B有正相关性, 买A同时更可能买B. 越大表示关联程度越强.

=1时 无关联.

<1时: 负相关, 买A不太可能买B, 越小表示关联越弱.

### 3. find 频繁项:

Counting Pairs in memory:

① approach 1: 用矩阵

若数字用4 bytes 则  $\rightarrow$  一个 pair 用 4 bytes

② 用三元表:  $[i, j, c]$

$\rightarrow$  一个 pair 用 12 bytes

for approach 1: 按列优先  
按  $\{1,2\}, \{1,3\}, \dots, \{1,n\}, \{2,3\}, \{2,4\}, \dots, \{2,n\}$  存储  
 $i < j$ , 存上三角.  $\{i, j\}$  对位于  $(i-1)(n-i/2) + j - 1$  的位置.



数据总数  $\frac{n(n-1)}{2}$ , 总字节数  $2n^2$

当 less than  $\frac{1}{3}$  of possible pairs 时 2 更好.

上述方法为 memory 足够时, 但若 memory 不多呢?

Frequent Pairs.

A-Priori 算法: key-idea: monotonicity - 单调性.

如果集合 I 出现不少于  $s$  次, 则 I 的子集也是.

如果集合 I 并未出现在  $s$  个 basket 中, 则没有包含 I 的集合出现在  $s$  个 basket 中.

step 1: Read baskets, 计算每个 item 出现的次数, 把大于  $s$  的 item 称为频繁 item.

这一步只需正比于 item 数目的内存.  $\rightarrow$  当然可以根据需要变动.

step 2: 再读篮子, 在内存中只计算 2 个元素或是频繁 item 的 basket.

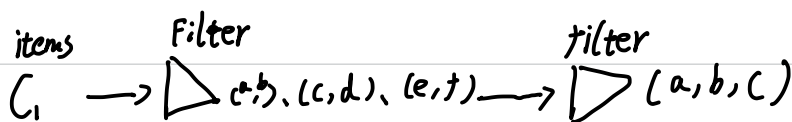
需要与频繁项的平方成正比的内存 (计数).

频繁的 pairs 列表.

$\rightarrow$  用三角矩阵存储的话,  $n = \text{frequent item 数目}$ .

一般  $k=2$  时内存要求最多.  
即第 2 步.

目前得到 Pairs ( , ) count, 作为  $k-1$  向  $k$  传递, 每层都计算该层所有集合的 count, 删去不重要的再传入下一层。



更大数据集时用 PCY 算法 (Park-Chen-Yu):

来源于 A-Priori 中第一步有大部分内存闲置, 用第一步闲置的内存减少第2步所需。

Pass 1 的 PCY: 除项目计数外, 维护一个 hash 表, 要有尽可能多 bucket → hash 冲突时保存内容  
为每一个散列桶保存一个计数, 桶只保留计数而不保留实际项目对。

在第一步, 不仅要得到每个 item 的 count, 还要对每个 basket 中产生所有项目对 Pairs, 创一张 hash 表。桶计数小于  $S$  的, 它的所有项目对都不是频繁的。

含频繁 Pair 的桶一定频繁。但频繁的桶可能内容全是不频繁的。

↓  
因为多个不同内容 hash 到一起。

Pass 2: 只计算 hash 频繁桶里面的。

用 bit 图代替 hash table.    1-频繁桶.    0-无.