

# A Note on Depth-First Mining of Frequent Itemsets

Jef Wijsen

February 22, 2008

## Abstract

This note may help you in understanding Section 6.6 in the textbook [1].

## 1 Vertical Database Layout

Let  $(\mathcal{I}, \prec)$  be a linearly ordered, finite set of *items*. Let  $\mathcal{T}$  be a set of transaction identifiers. Under the *vertical database layout*, a *transaction database* is a total function  $\text{cover} : \mathcal{I} \rightarrow 2^{\mathcal{T}}$ . All the following definitions are relative to such a fixed transaction database  $\text{cover}$ .

**Example 1** Assume  $a \prec b \prec c \prec d \prec e$ . The following transaction database is shown in Fig. 6.24 on page 364 of [1].

$$\begin{aligned}\text{cover}(a) &= \{1, 3, 4, 5, 6, 7, 8, 9\} \\ \text{cover}(b) &= \{1, 2, 5, 6, 8, 9, 10\} \\ \text{cover}(c) &= \{2, 3, 5, 6, 8, 10\} \\ \text{cover}(d) &= \{2, 3, 4, 6, 9\} \\ \text{cover}(e) &= \{3, 4, 10\}\end{aligned}$$

An *itemset* is a sequence  $a_1 a_2 \dots a_n$  where  $n \geq 1$  and  $a_1 \prec a_2 \prec \dots \prec a_n$  (the technical treatment can be easily extended to deal with the empty itemset). The function  $\text{cover}$  naturally extends to itemsets:

$$\text{cover}(a_1 a_2 \dots a_n) = \bigcap_{i=1}^n \text{cover}(a_i)$$

**Example 2**

$$\text{cover}(abc) = \{5, 6, 8\}$$

The (left) concatenation operator  $\cdot$  is defined as usual: if  $s = a_1 a_2 \dots a_n$  is an itemset and  $a_0 \prec a_1$ , then  $a_0 \cdot s$  denotes the itemset  $a_0 a_1 \dots a_n$ .

## 2 Conditional Cover

Let  $s$  be an itemset. We define the “conditional cover” as follows:

$$\begin{aligned}\text{cocov}[s] = \{(a, T) \mid &a \in \mathcal{I} \text{ and} \\ &a \text{ precedes (w.r.t. } \prec \text{) the leftmost symbol of } s \text{ and} \\ &T = \text{cover}(a \cdot s)\}.\end{aligned}$$

If  $b$  be the leftmost symbol of  $s$ , then  $\text{cocov}[s]$  is a transaction database with domain  $\{a \in \mathcal{I} \mid a \prec b\}$ .

**Example 3**

$$\begin{aligned}\text{cocov}[de](a) &= \{3, 4\} \\ \text{cocov}[de](b) &= \{\} \\ \text{cocov}[de](c) &= \{3\}\end{aligned}$$

**Property 1** If  $b \cdot s$  is an itemset and  $a \prec b$ , then

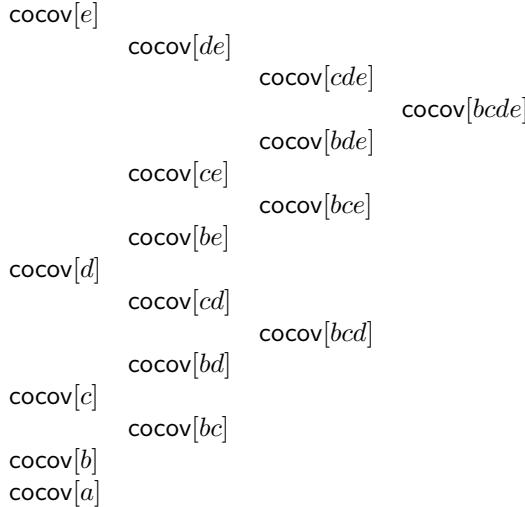
$$\text{cocov}[b \cdot s](a) = \text{cocov}[s](b) \cap \text{cocov}[s](a) .$$

**Proof.**

$$\begin{aligned}\text{cocov}[b \cdot s](a) &= \text{cover}(a \cdot b \cdot s) \\ \text{cover}(a \cdot b \cdot s) &= \text{cover}(a \cdot s) \cap \text{cover}(b \cdot s) \\ \text{cover}(a \cdot s) &= \text{cocov}[s](a) \\ \text{cover}(b \cdot s) &= \text{cocov}[s](b)\end{aligned}$$

□

This property suggests a recursive procedure for computing frequent itemsets with a support count greater than  $\sigma$ . Compute recursively the following sequence, while outputting frequent itemsets and ignoring infrequent items.



**Example 4** Assume  $\sigma = 2$ .

- $\text{cocov}[e] \boxed{\begin{matrix} a & \{3, 4\} \\ b & \{10\} \\ c & \{3, 10\} \\ d & \{3, 4\} \end{matrix}} \rightsquigarrow \{a, e\}$  In the *horizontal database layout*, this yields (omitting the infrequent item  $b$ ):  $\{(3, acd), (4, ad), (10, c)\}$ . The FP-tree is shown in Fig. 6.27 (b) on page 368 of [1].
- $\text{cocov}[de] \boxed{\begin{matrix} a & \{3, 4\} \\ c & \{3\} \end{matrix}} \rightsquigarrow \{a, d, e\}$  In the *horizontal database layout*, this yields (omitting the infrequent item  $c$ ):  $\{(3, a), (4, a)\}$ . The FP-tree is shown in Fig. 6.27 (d) on page 368 of [1].
- $\text{cocov}[ce] \boxed{a \quad \{3\}} \times$
- $\text{cocov}[d] \boxed{\begin{matrix} a & \{3, 4, 6, 9\} \\ b & \{2, 6, 9\} \\ c & \{2, 3, 6\} \end{matrix}} \rightsquigarrow \{a, d\}$
- $\text{cocov}[cd] \boxed{\begin{matrix} a & \{3, 6\} \\ b & \{2, 6\} \end{matrix}} \rightsquigarrow \{a, c, d\}$
- \*  $\text{cocov}[bcd] \boxed{a \quad \{6\}} \times$
- $\text{cocov}[bd] \boxed{a \quad \{6, 9\}} \rightsquigarrow \{a, b, d\}$
- $\text{cocov}[c] \boxed{\begin{matrix} a & \{5, 6, 8\} \\ b & \{2, 5, 6, 8, 10\} \end{matrix}} \rightsquigarrow \{a, c\}$
- $\text{cocov}[bc] \boxed{a \quad \{5, 6, 8\}} \rightsquigarrow \{a, b, c\}$
- $\text{cocov}[b] \boxed{a \quad \{1, 5, 6, 8, 9\}} \rightsquigarrow \{a, b\}$

## References

- [1] P.-N. Tan, M. Steinbach, and V. Kumar. *Introduction to Data Mining*. Addison Wesley, 2006.