

We want to find all **strong rules**. These are rules $a \rightarrow b$ such that:

$$\text{Supp}(a \cup b) \geq \theta, \text{ and } \text{Conf}(a \rightarrow b) \geq \text{minconf}.$$

Here θ is called the **minimum support threshold**.

The support has a monotonicity property called *downward closure*:

$$\text{If } \text{Supp}(a \cup b) \geq \theta \text{ then } \text{Supp}(a) \geq \theta \text{ and } \text{Supp}(b) \geq \theta.$$

That is, if $a \cup b$ is a frequent item set, then so are a and b .

$$\begin{aligned} \text{Supp}(a \cup b) &= \# \text{times } a \text{ and } b \text{ are purchased} \\ &\leq \# \text{times } a \text{ is purchased} = \text{Supp}(a). \end{aligned}$$

Apriori finds all frequent itemsets (a such that $\text{Supp}(a) \geq \theta$). We can use Apriori's result to get all strong rules $a \rightarrow b$ as follows:

- For each frequent itemset ℓ :
 - Find all nonempty subsets of ℓ
 - For each subset a , output $a \rightarrow \{\ell \setminus a\}$ whenever

$$\frac{\text{Supp}(\ell)}{\text{Supp}(a)} \geq \text{minconf}.$$

Now for Apriori. Use the downward closure property: generate all k -itemsets (itemsets of size k) from $(k - 1)$ -itemsets. It's a breadth-first-search.

Example:

$\theta = 10$

	<i>apples</i>	<i>bananas</i>	<i>cherries</i>		<i>elderberries</i>		<i>grapes</i>
1-itemsets:	a	b	c	d	e	f	g
supp:	25	20	30	45	29	5	17
2-itemsets:	{a,b}	{a,c}	{a,d}	{a,e}	...	{e,g}	
supp:	7	25	15	23		3	
3-itemsets:	{a,c,d}	{a,c,e}	{b,d,g}	...			
supp:	15	22	15				
4-itemsets:	{a,c,d,e}						
supp:	12						

Apriori Algorithm:

Input: Matrix M

$L_1 = \{\text{frequent 1-itemsets}; i \text{ such that } \text{Supp}(i) \geq \theta\}$.

For $k = 2$, while $L_{k-1} \neq \emptyset$ (while there are large $k - 1$ -itemsets), $k++$

- $C_k = \text{apriori_gen}(L_{k-1})$ generate candidate itemsets of size k
- $L_k = \{c : c \in C_k, \text{Supp}(c) \geq \theta\}$ frequent itemsets of size k (loop over transactions, scan the database)

end

Output: $\bigcup_k L_k$.

The subroutine `apriori_gen` joins L_{k-1} to L_{k-1} .

`apriori_gen` Subroutine:

Input: L_{k-1}

Find all pairs of itemsets in L_{k-1} where the first $k - 2$ items are identical.

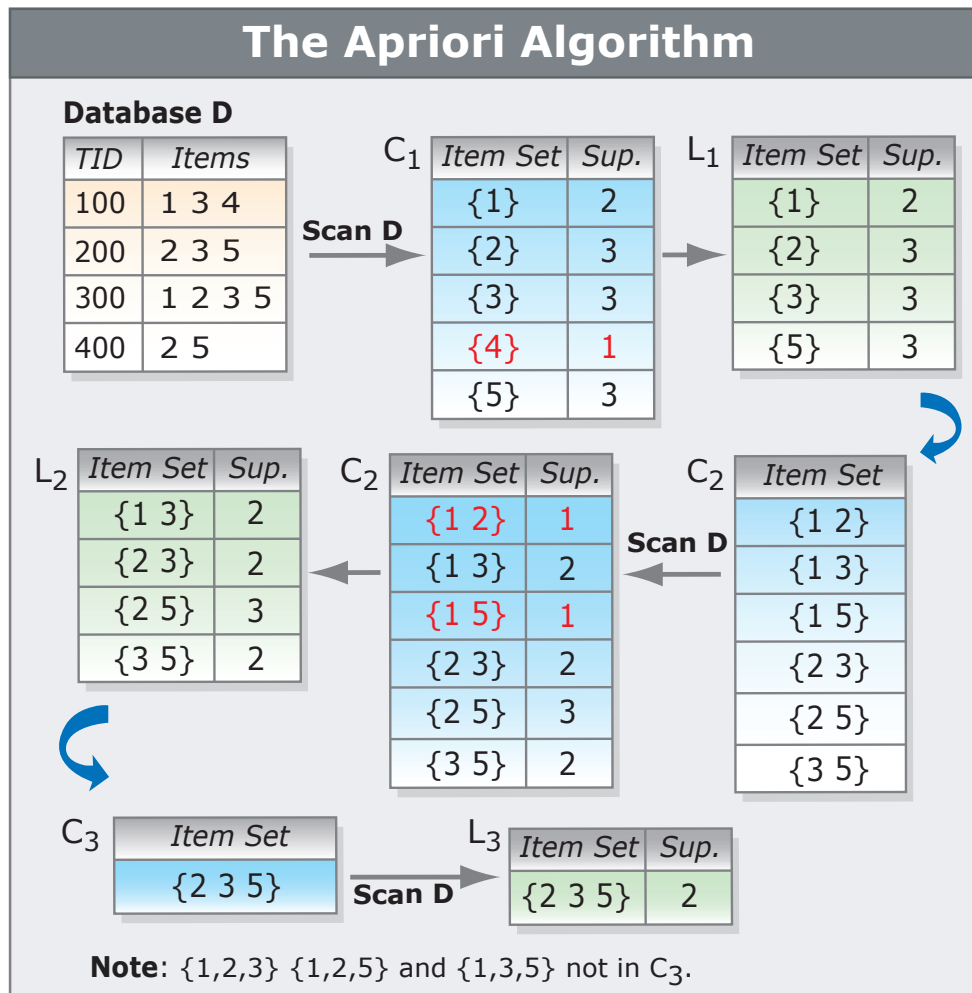
Union them (lexicographically) to get $C_k^{\text{too big}}$,

e.g., $\{a, b, c, d, e, f\}, \{a, b, c, d, e, g\} \rightarrow \{a, b, c, d, e, f, g\}$

Prune: $C_k = \{c \in C_k^{\text{too big}}, \text{all } (k - 1)\text{-subsets } c_s \text{ of } c \text{ obey } c_s \in L_{k-1}\}$.

Output: C_k .

Example of Prune step: consider $\{a, b, c, d, e, f, g\}$ which is in $C_k^{\text{too big}}$, and I want to know whether it's in C_k . Look at $\{a, b, c, d, e, f, g\}, \{a, \cancel{b}, c, d, e, f, g\}, \{a, b, \cancel{c}, d, e, f, g\}, \{a, b, c, \cancel{d}, e, f, g\}$, etc. If any are not in L_6 , then prune $\{a, b, c, d, e, f, g\}$ from L_7 .



-
- Apriori scans the database at most how many times?
 - Huge number of candidate sets. ☹
 - Spawned huge number of apriori-like papers.
-

What do you do with the rules after they're generated?

- Information overload (give up)
- Order rules by “interestingness”

– Confidence

$$\hat{P}(b|a) = \frac{\text{Supp}(a \cup b)}{\text{Supp}(a)}$$

– “Lift” / “Interest”

$$\frac{\hat{P}(b|a)}{\hat{P}(b)} = \frac{\text{Supp}(b)}{1 - \frac{\text{Supp}(a \cup b)}{\text{Supp}(a)}}$$

:

– Hundreds!

Research questions:

- mining more than just itemsets (e.g, sequences, trees, graphs)
- incorporating taxonomy in items
- boolean logic and “logical analysis of data”
- Cynthia’s questions: Can we use rules within ML to get good predictive models?

MIT OpenCourseWare
<http://ocw.mit.edu>

15.097 Prediction: Machine Learning and Statistics
Spring 2012

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.