

Frequent itemsets, association rules

Jiří Kléma

Department of Cybernetics,
Czech Technical University in Prague



<http://ida.felk.cvut.cz>

-

Association rules

- Association Rules (ARs)

- Definition

- simple event co-occurrence assertions based on data,
- probabilistic character – co-occurrence is not strict,

- Notation and meaning

- if **Ant** then **Suc**,
- another notation: $\text{Ant} \Rightarrow \text{Suc}$,
- antecedent (Ant) and succedent/consequent (Suc) define general events observable in data,
- event – a binary phenomenon, it either occurs or not,
- an extensive representation (data) transformed into a concised and understandable description (knowledge).

- Association rules, examples

- book store recommendations (Amazon):
 $\{\text{Castaneda: } \textbf{Teachings of Don Juan}\}$
 $\Rightarrow \{\text{Hesse: } \textbf{Der Steppenwolf} \ \& \ \text{Ruiz: } \textbf{The Four Agreements}\}$
- relation among risk factors and diseases in medicine (Stulong):
 $\{\textbf{beer} \geq 1\text{litre/day} \ \& \ \textbf{liquors}=0\} \Rightarrow \{\text{not}(\textbf{heart disease})\}$

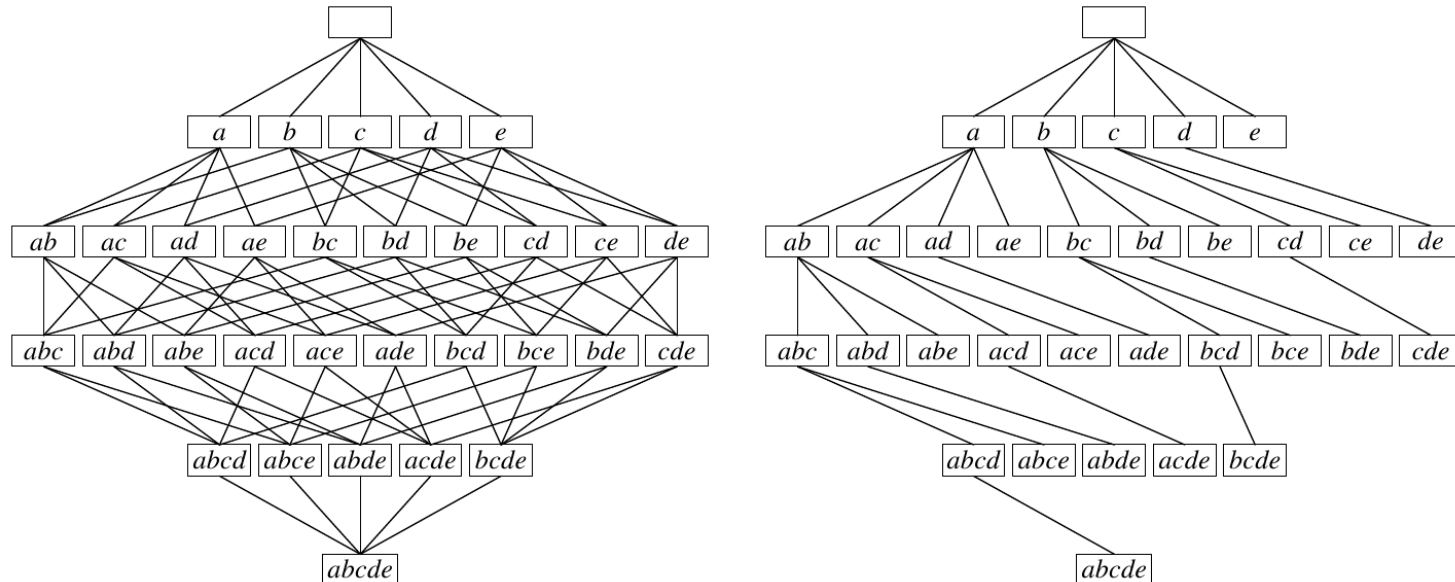
- A4M33SAD

-
- A horizontal progress bar consisting of 25 small square icons. The first 6 squares are solid black, and the remaining 19 squares are white with a black outline.

-

Frequent itemset mining – method categorization (1)

- a set of all the itemsets is **partially ordered** (makes a poset)
 - can be depicted as an acyclic graph – Hasse diagram,
 - nodes = itemsets, an edge $I \rightarrow J$ iff $I < J$ and there is no $K : I < K < J$,
 - when depicting all the subsets it also makes a **lattice**,
 - efficient to reduce on a **tree** (each node needs to be visited and tested only once),
- methods for the itemset lattice/tree search
 - breath-first – level-wise, each level concerns itemsets of a certain length,
 - depth-first – traverse the itemsets with an identical prefix.



Frequent itemset mining – method categorization (2)

- transaction set/database representation

- horizontal – transactions as the main units, transaction \approx a list/array of items
 - * a natural way,
- vertical – items as the main units, a transaction list is stored for each item
 - * advantage: efficient (recursive) access to the transaction list of an itemset,
 - * the transaction list for a pair of items is the intersection of the transaction lists of the individual items.

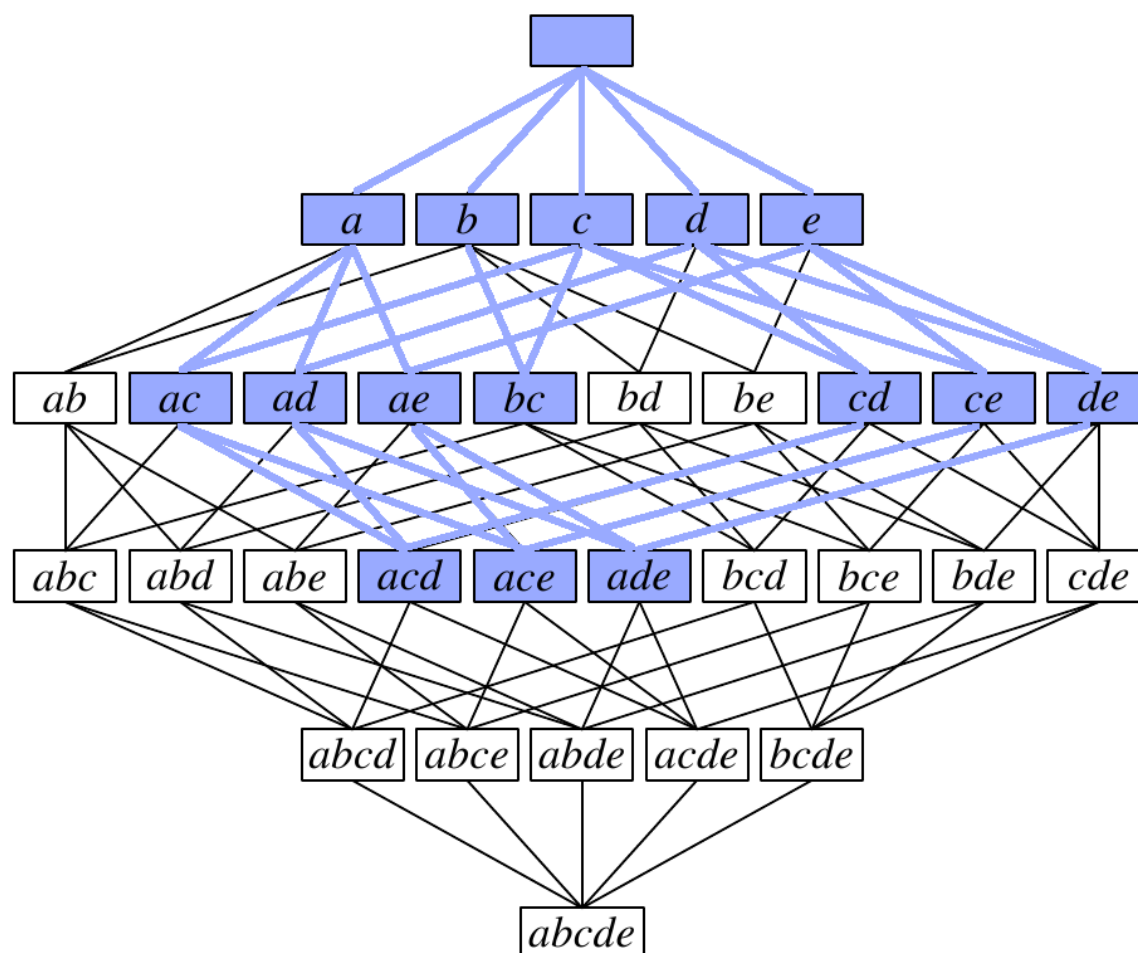
Transactions	Items
t_1	a, d, e
t_2	b, c, d
t_3	a, c, e
t_4	a, c, d, e
t_5	a, e
t_6	a, c, d
t_7	b, c
t_8	a, c, d, e
t_9	b, c, e
t_{10}	a, d, e

a	b	c	d	e
1	2	2	1	1
3	7	3	2	3
4	9	4	4	4
5		6	6	5
6		7	8	8
8		8	10	9
10		9		10

APRIORI algorithm – the basic idea

- pioneering, the most well-known, but not the most efficient,
- based on the elemental characteristic of any frequent itemset:
Each subset of a frequent itemset is frequent.
- as we proceed bottom-up from subsets to supersets
the logical **contraposition** principle
 $(p \Rightarrow q) \Leftrightarrow (\neg q \Rightarrow \neg p)$
- the anti-monotone property transformed to a monotone property, consequence:
No superset of an infrequent itemset can be frequent.
- candidate itemsets
 - potentially frequent – all the subsets are known to be frequent.
- APRIORI categories: breath-first search, horizontal transaction representation.

Transakce	Položky
t_1	a, d, e
t_2	b, c, d
t_3	a, c, e
t_4	a, c, d, e
t_5	a, e
t_6	a, c, d
t_7	b, c
t_8	a, c, d, e
t_9	b, c, e
t_{10}	a, d, e

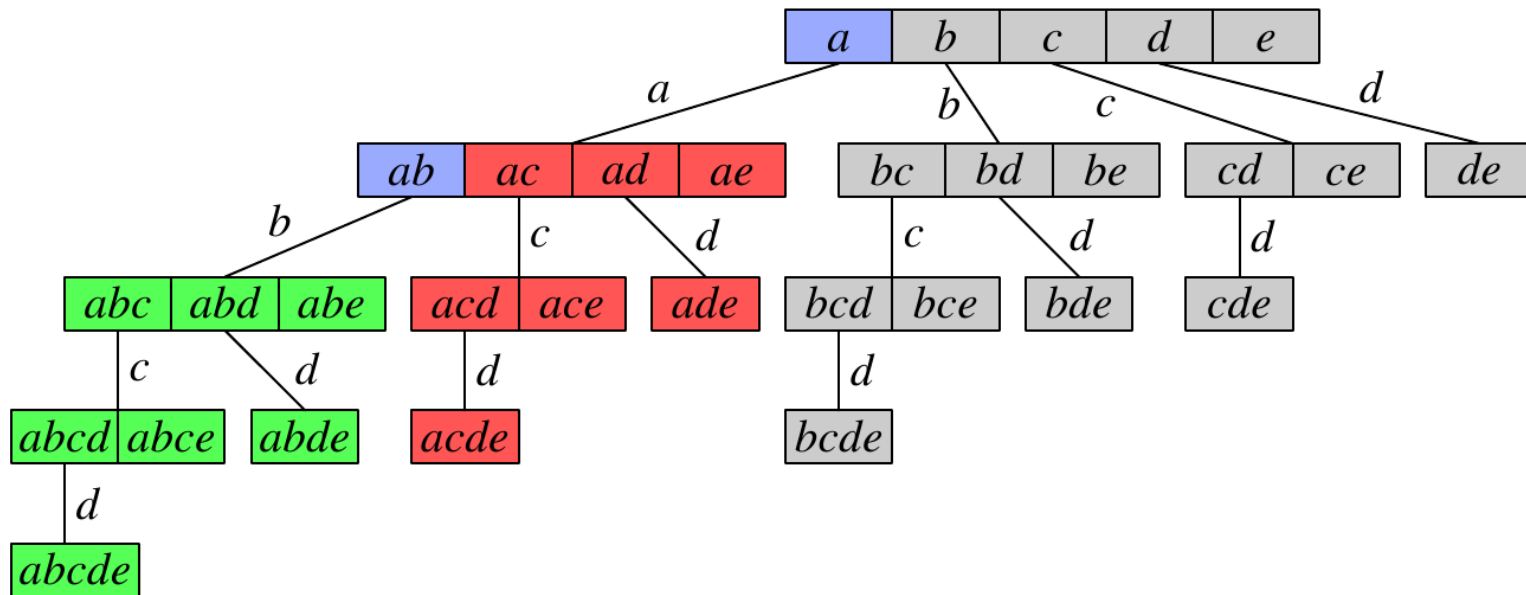


ECLAT algorithm (Zaki et al., 1997) – the basic idea

- sorts items lexicographically → canonical itemset representation
 - $\{a, b, c\} \approx abc < bac < bca < \dots < cba$,
 - an itemset encoded by its lexicographically smallest (largest) code word,
- the tree is searched in a **depth-first** way
 - owing to the canonical representation it is a prefix tree,
- uses purely **vertical** transaction set representation,
- can generate more candidate itemsets than APRIORI
 - decides when support of any subset is not necessarily available.

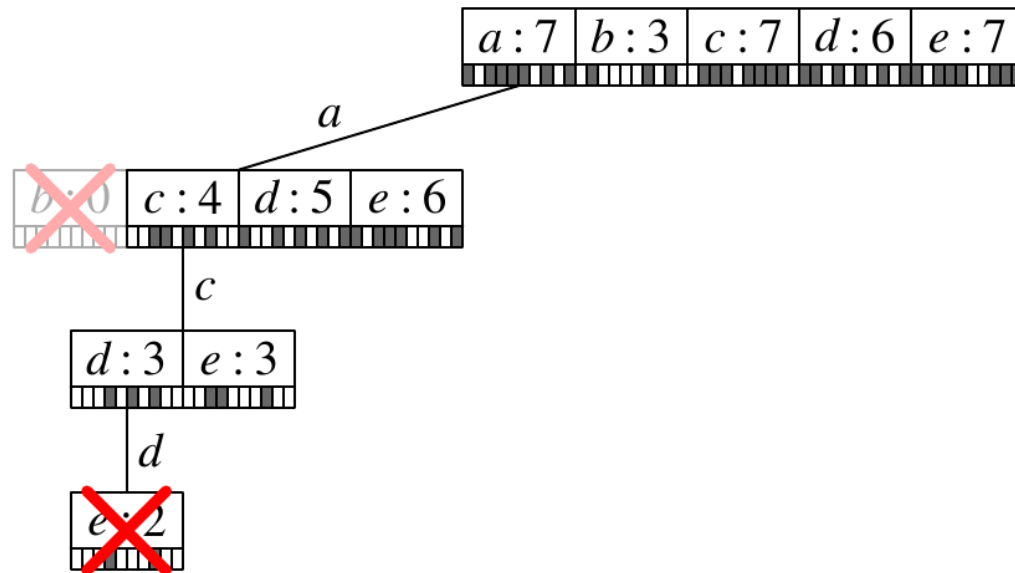
Conditional transaction database – depth-first search

- depth-first search the prefix tree, **divide and conquer** strategy
 - find all the frequent itemsets with the given prefix first,
 - do the same for the rest of itemsets,
 - leads to transaction set splits (transactions with/without the given prefix),
- node colors
 - the prefix in blue, itemsets having the prefix in green, itemsets without the prefix in red,
 - a recursive procedure, the previous α -step needs also to be concerned.



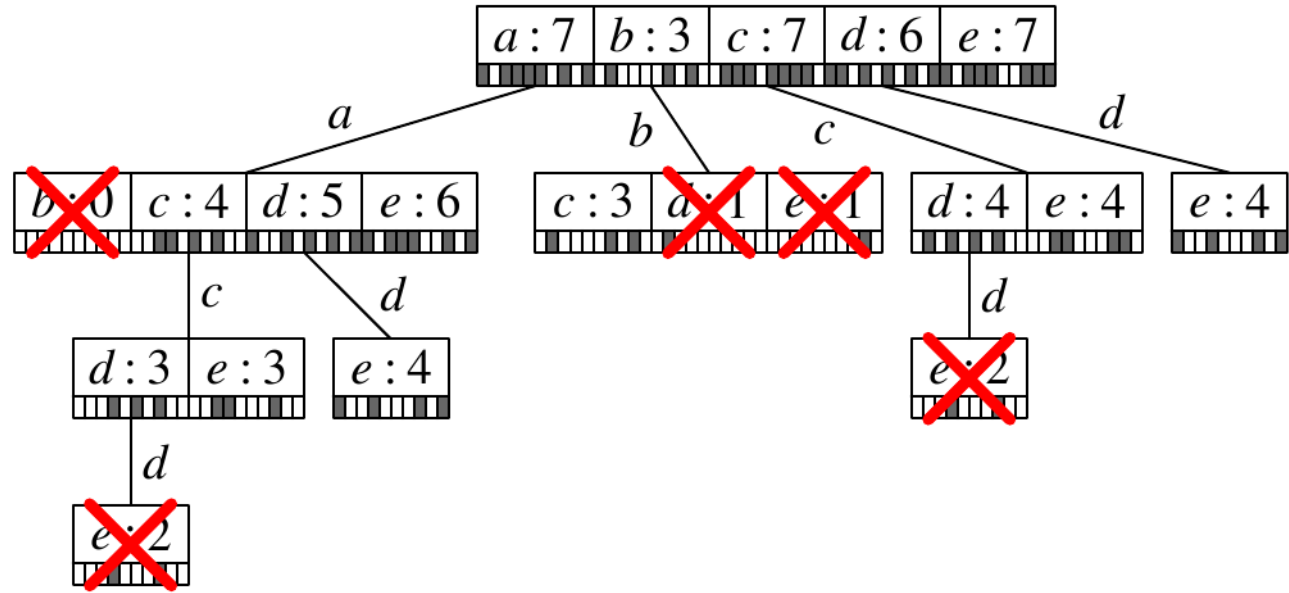
ECLAT example, $s_{min} = 3$ (Borgelt: Frequent Pattern Mining)

t_1	a, d, e
t_2	b, c, d
t_3	a, c, e
t_4	a, c, d, e
t_5	a, e
t_6	a, c, d
t_7	b, c
t_8	a, c, d, e
t_9	b, c, e
t_{10}	a, d, e



- preprocessing: vertical representation by the bit vector (grey/white – item in/out of transaction)
 - the only transaction database scan, intersections follow exclusively,
- step 1: the conditional transaction database for a item,
- step 2: $\{a, b\}$ infrequent – prune,
- step 3: the conditional transaction database for $\{a, c\}$ itemset,
- step 4: the conditional transaction database for $\{a, c, d\}$ itemset and prune $\{a, c, d, e\}$.

t_1	a, d, e
t_2	b, c, d
t_3	a, c, e
t_4	a, c, d, e
t_5	a, e
t_6	a, c, d
t_7	b, c
t_8	a, c, d, e
t_9	b, c, e
t_{10}	a, d, e



- the whole tree shown, the outcome (certainly) identical with APRIORI,
- APRIORI might prune $\{a, c, d, e\}$ without counting its support,
 - knowing that $s_{\{c,d,e\}} = 2 \leq s_{min} = 3$,
- in contrary, APRIORI needs more transaction database scans.

Reducing the output – the pruned sets of frequent itemsets

- the number of frequent itemsets can be prohibitive
 - the output is not comprehensible, a user can be interested in long patterns only,
 - leads to a notion of **maximal** itemset
 - * frequent and none of its proper supersets is frequent,
 - * the set of maximal itemsets:

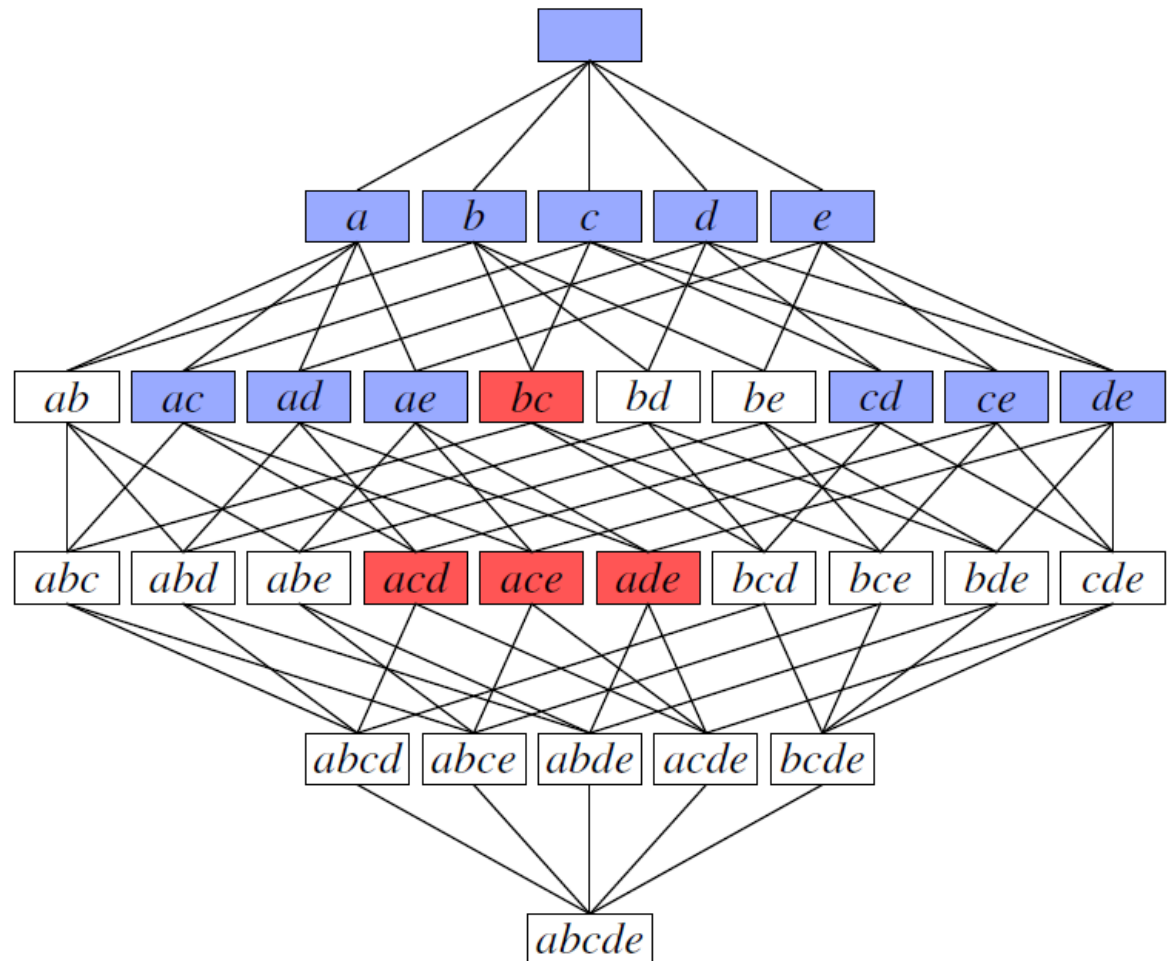
$$M_D(s_{\min}) = \{J \subseteq I \mid s_D(J) \geq s_{\min} \wedge \forall K \supset J : s_D(K) < s_{\min}\},$$

- the set of frequent itemsets is redundant
 - all the information about it can be preserved in a smaller set (subset),
 - leads to a notion of **closed** itemset
 - * frequent and none of its proper supersets has the same support,
 - * the set of closed itemsets:

$$C_D(s_{\min}) = \{J \subseteq I \mid s_D(J) \geq s_{\min} \wedge \forall K \supset J : s_D(K) < s_D(J)\},$$

- obvious relations
 - all maximal itemsets and all closed itemsets are frequent,
 - any maximal itemset is necessarily closed.

Transactions	Items
t_1	a, d, e
t_2	b, c, d
t_3	a, c, e
t_4	a, c, d, e
t_5	a, e
t_6	a, c, d
t_7	b, c
t_8	a, c, d, e
t_9	b, c, e
t_{10}	a, d, e



- filter the set of frequent itemsets
 - * reasonable when the set of frequent itemsets is needed anyway,
- direct search with earlier and more efficient pruning
 - * a compact representation accelerates search,
 - * specialized algorithms derived from classical ones – MaxMiner, Closet, Charm, GenMax,
 - * among other properties, for any closed itemset it holds
 - the closed itemset matches the intersection of all the transactions that contain it,
 - it also explains why $\{d, e\}$ is not closed:

Transactions	Items
t_1	a, d, e
t_4	a, c, d, e
t_8	a, c, d, e
t_{10}	a, d, e
\cap	a, d, e

Generate rules from frequent itemsets – step 2

Inputs:

$$I, D, L, \alpha_{min};$$

Output :

```
R; % pravidla splňující  $s_{min}$  a  $\alpha_{min}$ 
```

AR-Gen:

$$R = \emptyset;$$

```
for  $\forall l \in L$  do:
```

for $\forall x \subset l$ such that $x \neq \emptyset$ and $x \neq l$ do:

if $s(l)/s(x) \geq \alpha_{min}$, then $R = R \cup \{x \Rightarrow (1-x)\}$

(apply the property: $s(l)/s(x) < \alpha_{min} \Rightarrow \forall x' \subset x \ s(l)/s(x') < \alpha_{min}$)

- Example: market basket analysis

- Inputs: $L=\{\text{Bread, Butter}\}$ (generated for $s_{min}=30\%$), $\alpha_{min}=50\%$
- Output: $R=\{\text{Bread} \Rightarrow \text{Butter: } s=60\%, \alpha=75\%, \text{Butter} \Rightarrow \text{Bread: } s=60\%, \alpha=100\%\}$

Transactions	Items
t_1	RZN
t_2	VI, SAD, AU
t_3	PAH, AU
t_4	PAH, VI, AU
t_5	PAH, MAS
t_6	VI, AU
t_7	PAH, SAD
t_8	PAH, VI, MAS, AU
t_9	PAH
t_{10}	PAH, VI, AU

Transactions	Items
t_{11}	AU
t_{12}	RZN, PAH, VI, SAD, AU
t_{13}	PAH, VI, MAS, AU
t_{14}	VI, SAD, AU
t_{15}	PAH, AU
t_{16}	SAD, AU
t_{17}	RZN, PAH, SAD
t_{18}	PAH, VI, MAS, AU
t_{19}	PAH
t_{20}	PAH, VI, MAS, AU

APRIORI step – $s_{min}=20\%$, resp. 4

i	C_i	L_i
1	$\{\text{RZN}\}, \{\text{PAH}\}, \{\text{VI}\}$ $\{\text{MAS}\}, \{\text{SAD}\}, \{\text{AU}\}$	$\{\text{PAH}\}, \{\text{VI}\}, \{\text{MAS}\}$ $\{\text{SAD}\}, \{\text{AU}\}$
2	$\{\text{PAH}, \text{VI}\}, \{\text{PAH}, \text{MAS}\}, \{\text{PAH}, \text{SAD}\}$ $\{\text{PAH}, \text{AU}\}, \{\text{VI}, \text{MAS}\}, \{\text{VI}, \text{SAD}\}$ $\{\text{VI}, \text{AU}\}, \{\text{MAS}, \text{SAD}\}, \{\text{MAS}, \text{AU}\}$ $\{\text{SAD}, \text{AU}\}$	$\{\text{PAH}, \text{VI}\}, \{\text{PAH}, \text{MAS}\}$ $\{\text{PAH}, \text{AU}\}, \{\text{VI}, \text{MAS}\}$ $\{\text{VI}, \text{AU}\}, \{\text{MAS}, \text{AU}\}$ $\{\text{SAD}, \text{AU}\}$
3	$\{\text{PAH}, \text{VI}, \text{MAS}\}, \{\text{PAH}, \text{VI}, \text{AU}\}$ $\{\text{PAH}, \text{MAS}, \text{AU}\}, \{\text{PAH}, \text{SAD}, \text{AU}\}$ $\{\text{VI}, \text{MAS}, \text{AU}\}, \{\text{VI}, \text{SAD}, \text{AU}\}$ $\{\text{MAS}, \text{SAD}, \text{AU}\}$	$\{\text{PAH}, \text{VI}, \text{MAS}\}$ $\{\text{PAH}, \text{VI}, \text{AU}\}$ $\{\text{PAH}, \text{MAS}, \text{AU}\}$ $\{\text{VI}, \text{MAS}, \text{AU}\}$
4	$\{\text{PAH}, \text{VI}, \text{MAS}, \text{AU}\}$	$\{\text{PAH}, \text{VI}, \text{MAS}, \text{AU}\}$
5	\emptyset	\emptyset

L₂

PAH, VI: PAH \Rightarrow VI $\alpha=50\%$, VI \Rightarrow PAH $\alpha=70\%$
(PAH & VI concurrently 7times, PAH 14times, VI 10times)

PAH, MAS: PAH \Rightarrow MAS 36%, **MAS \Rightarrow PAH 100%**
(PAH & MAS concurrently 5times, PAH 14times, MAS 5times)

L₃

PAH, VI, MAS: PAH & VI \Rightarrow MAS 57%, **PAH & MAS \Rightarrow VI 80%**, **VI & MAS \Rightarrow PAH 100%**
(PAH nor VI cannot make an antecedent, test MAS only)
MAS \Rightarrow PAH & VI 80%

$$\mathbf{L}_4$$

PAH, VI, MAS, AU: **PAH & VI & MAS \Rightarrow AU** 100%, PAH & VI & AU \Rightarrow MAS 57%,
PAH & MAS & AU \Rightarrow VI 100%, **VI & MAS & AU \Rightarrow PAH** 100%
 (the antecedents PAH & VI, PAH & AU, VI & AU without testing)
PAH & MAS \Rightarrow VI & AU 80%, **VI & MAS \Rightarrow PAH & AU** 100%,
MAS & AU \Rightarrow PAH & VI 100%
 (the antecedents PAH, VI a AU without testing)
MAS \Rightarrow PAH & VI & AU 80%

Four-fold table, quantifiers for the relation between Ant and Suc

- 4-fold table (4FT),

- $a, b, c, d \rightarrow$ the numbers of transactions meeting conditions.

4FT	Suc	\neg Suc	Σ
Ant	a	b	$r=a+b$
\neg Ant	c	d	$s=c+d$
Σ	$k=a+c$	$l=b+d$	$n=a+b+c+d$

- Confidence is not the only/always best quantifier

- its **implicative** nature is misleading for frequent succedents,
- independent itemsets can show a high confidence,
- 4-fold table example ($s=45\%, \alpha=90\%$, Ant and Suc independent):

450	50	500
450	50	500
900	100	1000

- | | | | | | | | | |
|-----|-----|------|-----|-----|------|-----|-----|------|
| 450 | 50 | 500 | 10 | 1 | 11 | 450 | 50 | 500 |
| 450 | 50 | 500 | 90 | 899 | 989 | 50 | 450 | 500 |
| 900 | 100 | 1000 | 100 | 900 | 1000 | 500 | 500 | 1000 |
- $s=0.45$, $\alpha=0.9$,
lift=1, **leverage=0**,
conviction=1
- $s=0.01$** , $\alpha=0.91$,
lift=9.09, **leverage=0.01**,
conviction=9.9
- $s=0.45$, $\alpha=0.9$,
lift=1.8, **leverage=0.2**,
conviction=5

Recommended reading, lecture resources

:: Reading

- Agrawal, Srikant: **Fast Algorithms for Mining Association Rules.**
 - the article that introduced the task and proposed APRIORI algorithm,
 - <http://rakesh.agrawal-family.com/papers/vldb94apriori.pdf>,
- Borgelt: **Frequent Pattern Mining.**
 - slides, a detailed course, including a formal notation,
 - <http://www.borgelt.net/teach/fpm/slides.html>,
- Hájek, Havránek: **Mechanizing Hypothesis Formation.**
 - a pioneering theory from 1966, decades before Agrawal,
 - <http://www.cs.cas.cz/hajek/guhabook/>.