

Mineração de Dados

Regras de Associação



- 1 Introdução
- 2 Mineração de *Itemsets* Frequentes
- 3 Regras de Associação
- 4 Algoritmo Força Bruta pra Minerar *Itemsets*
- 5 *Apriori*
- 6 *Eclat*
- 7 *dEclat*
- 8 Gerando Regras de Associação

Introdução

Mineração de Itens Frequentes

- ▶ Em muitas situações, é importante conhecer as relações entre dois ou mais objetos da base de dados
 - ▶ Esses conjuntos de itens são chamados de *itemsets*
- ▶ Análises de cestas de compras é uma aplicação tradicional
 - ▶ Itens que são comprados juntos em mercados
 - ▶ Um caso popular é a compra de fraldas e cervejas
- ▶ A mineração de *itemsets* frequentes é uma tarefa exploratória básica
 - ▶ Busca por co-ocorrências
 - ▶ Fornece uma estimativa da probabilidade conjunta
- ▶ Uma vez determinado os *itemsets* frequentes, pode-se extrair regras de associação desses conjuntos
 - ▶ Fornece informação da ocorrência condicional dos itens
- ▶ Este conteúdo é baseado no material de Zaki & Meira Jr., *Data Mining and Analysis*

Mineração de Itens Frequentes: Terminologia

▶ *Itemsets*

- ▶ Seja $I = \{x_1, x_2, \dots, x_m\}$ um conjunto de elementos chamados itens
- ▶ Um conjunto $X \subseteq I$ é chamado de *itemset* e um *itemset* de cardinalidade k é chamado de *k-itemset*
- ▶ $I^{(k)}$ é o conjunto de todos os *k-itemsets*

▶ *Tidsets*

- ▶ Seja $T = \{t_1, t_2, \dots, t_n\}$ um conjunto de identificadores de transações ou *tids*
- ▶ Um conjunto $\mathbf{T} \subseteq T$ é chamado de *tidset*

▶ Na prática, *itemsets* e *tidsets* são mantidos ordenados

▶ Transações

- ▶ Uma transação é uma tupla $\langle t, X \rangle$ onde $t \in T$ e X é um *itemset*

▶ Uma base de dados binária D é um conjunto que relaciona *tids* e itens, ou seja, $D \subseteq T \times I$

Representação de Bases de Dados

- ▶ Bases de dados binárias podem ser representadas por uma base transacional/horizontal ou vertical
- ▶ $\mathbf{i}(t)$ é o conjunto de itens do *tid* $t \in T$
- ▶ $\mathbf{t}(x)$ é o conjunto de *tids* que contem o item x

Bases Transacionais e Verticais

D	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

<i>t</i>	i(t)
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

(b) Transaction database

<i>x</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
t(x)	1	1	2	1	1
	3	2	4	3	2
	4	3	5	5	3
	5	4	6	6	4
		5			5
		6			

(c) Vertical database

- ▶ A base de dados *D* tem 5 itens ($I = \{A, B, C, D, E\}$) e 6 *tids* ($T = \{1, 2, 3, 4, 5, 6\}$)
- ▶ A primeira transação é $\langle 1, \{A, B, D, E\} \rangle$ ou $\langle 1, ABDE \rangle$

Mineração de *Itemsets* Frequentes

Suporte e *Itemsets* Frequentes

- ▶ Suporte de um *itemset* X é a quantidade de transações em D que contem X
$$sup(X) = |\{t | t \in D \text{ e } X \subseteq i(t)\}| = |t(X)|$$
- ▶ O suporte relativo pode ser definido como
$$rsup(X) = \frac{sup(X)}{|D|}$$
- ▶ O suporte relativo é uma estimativa da probabilidade conjunta dos itens em X
- ▶ X é chamado de *itemset* frequente quando $sup(X) \geq minsup$, onde $minsup$ é um limiar de suporte mínimo (definido pelo usuário)
- ▶ O conjunto F denota o conjunto de todos os *itemsets* frequentes e $F^{(k)}$ é o conjunto de k -*itemsets* frequentes

Itemsets Frequentes

► $minsup = 3$

t	$i(t)$
1	<i>ABDE</i>
2	<i>BCE</i>
3	<i>ABDE</i>
4	<i>ABCE</i>
5	<i>ABCDE</i>
6	<i>BCD</i>

Transaction Database

sup	itemsets
6	<i>B</i>
5	<i>E, BE</i>
4	<i>A, C, D, AB, AE, BC, BD, ABE</i>
3	<i>AD, CE, DE, ABD, ADE, BCE, BDE, ABDE</i>

Frequent Itemsets

The 19 frequent itemsets shown in the table comprise the set \mathcal{F} . The sets of all frequent k -itemsets are

$$\mathcal{F}^{(1)} = \{A, B, C, D, E\}$$

$$\mathcal{F}^{(2)} = \{AB, AD, AE, BC, BD, BE, CE, DE\}$$

$$\mathcal{F}^{(3)} = \{ABD, ABE, ADE, BCE, BDE\}$$

$$\mathcal{F}^{(4)} = \{ABDE\}$$

Regras de Associação

Regras de Associação

- ▶ Uma regra de associação é definida como $X \rightarrow Y$, onde X e Y são *itemsets* disjuntos
- ▶ Seja o suporte da regra o número de transações em que X e Y ocorrem conjuntamente ($X \cup Y$ ou XY)
$$\text{sup}(X \rightarrow Y) = \text{sup}(XY) = |t(XY)|$$

- ▶ A fração em que XY ocorre é o suporte relativo da regra
$$\text{rsup}(X \rightarrow Y) = \frac{\text{sup}(XY)}{|D|} = P(X \wedge Y)$$

- ▶ A confiança de uma regra é a probabilidade condicional das transações conterem Y tal que contém X
$$\text{conf}(X \rightarrow Y) = P(Y|X) = \frac{P(X \wedge Y)}{P(X)} = \frac{\text{sup}(XY)}{\text{sup}(X)}$$

Algoritmo Força Bruta pra Minerar *Itemsets*

Algoritmo Força Bruta

- ▶ Determina todos os *itemsets*, computa seus valores de suporte e guarda os *itemsets* frequentes
- ▶ Complexidade computacional $O(|I| |D| 2^{|I|})$

Algoritmo Força Bruta

BRUTEFORCE ($\mathbf{D}, \mathcal{I}, \text{minsup}$):

```
1  $\mathcal{F} \leftarrow \emptyset$  // set of frequent itemsets
2 foreach  $X \subseteq \mathcal{I}$  do
3    $\text{sup}(X) \leftarrow \text{COMPUTESUPPORT}(X, \mathbf{D})$ 
4   if  $\text{sup}(X) \geq \text{minsup}$  then
5      $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X, \text{sup}(X))\}$ 
6 return  $\mathcal{F}$ 
```

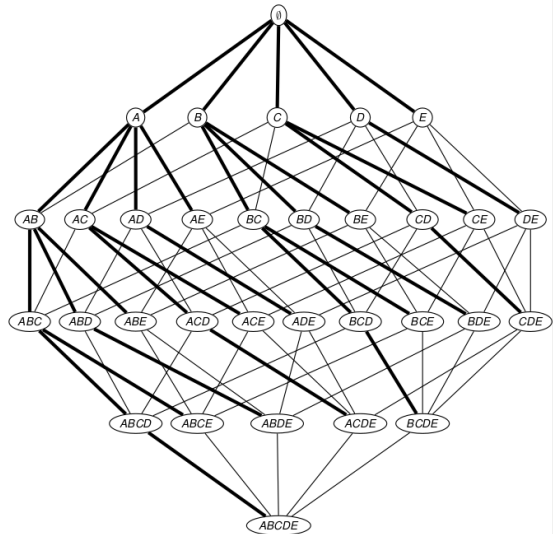
COMPUTESUPPORT (X, \mathbf{D}):

```
1  $\text{sup}(X) \leftarrow 0$ 
2 foreach  $\langle t, \mathbf{i}(t) \rangle \in \mathbf{D}$  do
3   if  $X \subseteq \mathbf{i}(t)$  then
4      $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 
5 return  $\text{sup}(X)$ 
```

Algoritmo Força Bruta

Itemset search space is a lattice where any two itemsets X and Y are connected by a link iff X is an *immediate subset* of Y , that is, $X \subseteq Y$ and $|X| = |Y| - 1$.

Frequent itemsets can be enumerated using either a BFS or DFS search on the *prefix tree*, where two itemsets X, Y are connected by a link iff X is an immediate subset and prefix of Y . This allows one to enumerate itemsets starting with an empty set, and adding one more item at a time.

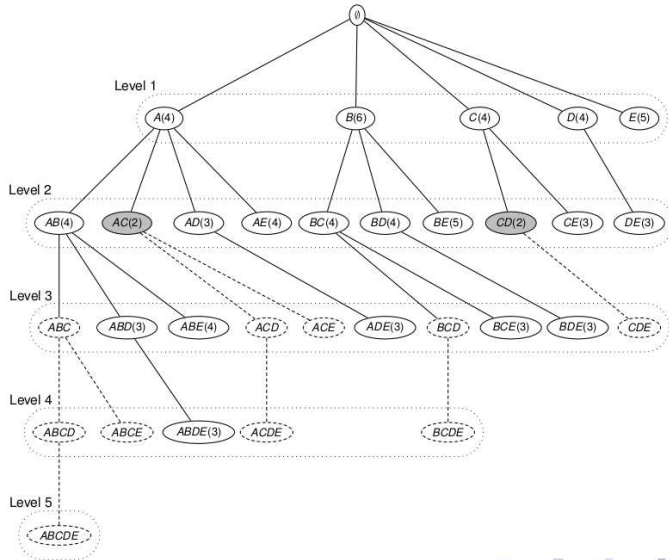


Apriori

Apriori

- ▶ Se $X \subseteq Y$ então $\text{sup}(X) \geq \text{sup}(Y)$ e, assim
 - ▶ se X é um *itemset* frequente então qualquer subconjunto $Z \subseteq X$ também é frequente
 - ▶ se X não é frequente então qualquer superconjunto $Z \supseteq X$ não é frequente
- ▶ O algoritmo *Apriori* explora essas propriedades
 - ▶ Reduz as buscas por evitar os candidatos a *itemsets* infrequentes
 - ▶ Subconjuntos infrequentes não compõem *itemsets* frequentes

Apriori



Apriori

APRIORI (\mathbf{D} , \mathcal{I} , $minsup$):

```

1  $\mathcal{F} \leftarrow \emptyset$ 
2  $\mathcal{C}^{(1)} \leftarrow \{\emptyset\}$  // Initial prefix tree with single items
3 foreach  $i \in \mathcal{I}$  do Add  $i$  as child of  $\emptyset$  in  $\mathcal{C}^{(1)}$  with  $sup(i) \leftarrow 0$ 
4  $k \leftarrow 1$  //  $k$  denotes the level
5 while  $\mathcal{C}^{(k)} \neq \emptyset$  do
6   COMPUTESUPPORT ( $\mathcal{C}^{(k)}$ ,  $\mathbf{D}$ )
7   foreach leaf  $X \in \mathcal{C}^{(k)}$  do
8     if  $sup(X) \geq minsup$  then  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X, sup(X))\}$ 
9     else remove  $X$  from  $\mathcal{C}^{(k)}$ 
10   $\mathcal{C}^{(k+1)} \leftarrow \text{EXTENDPREFIXTREE} (\mathcal{C}^{(k)})$ 
11   $k \leftarrow k + 1$ 
12 return  $\mathcal{F}^{(k)}$ 

```

Apriori

COMPUTESUPPORT ($\mathcal{C}^{(k)}$, **D**):

```

1 foreach  $\langle t, i(t) \rangle \in \mathbf{D}$  do
2   foreach  $k$ -subset  $X \subseteq i(t)$  do
3     if  $X \in \mathcal{C}^{(k)}$  then  $sup(X) \leftarrow sup(X) + 1$ 
```

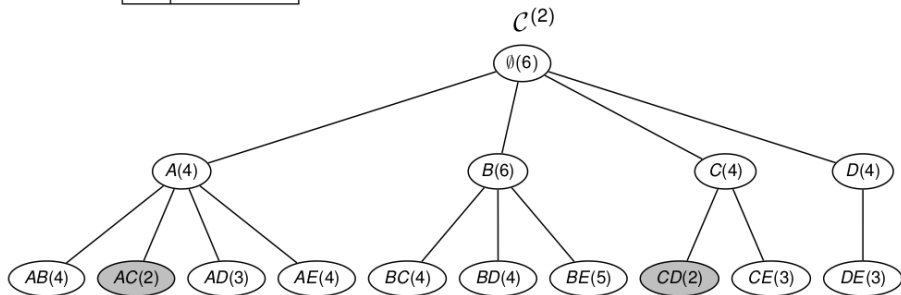
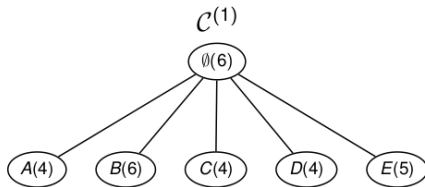
EXTENDPREFIXTREE ($\mathcal{C}^{(k)}$):

```

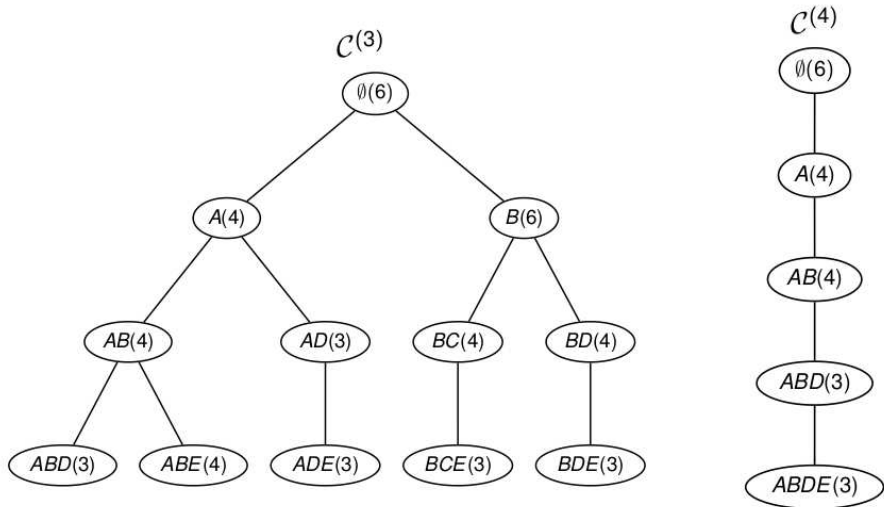
1 foreach leaf  $X_a \in \mathcal{C}^{(k)}$  do
2   foreach leaf  $X_b \in \text{SIBLING}(X_a)$ , such that  $b > a$  do
3      $X_{ab} \leftarrow X_a \cup X_b$ 
      // prune candidate if there are any infrequent
      subsets
4     if  $X_j \in \mathcal{C}^{(k)}$ , for all  $X_j \subset X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
5       Add  $X_{ab}$  as child of  $X_a$  with  $sup(X_{ab}) \leftarrow 0$ 
6   if no extensions from  $X_a$  then
7     remove  $X_a$ , and all ancestors of  $X_a$  with no extensions, from  $\mathcal{C}^{(k)}$ 
8 return  $\mathcal{C}^{(k)}$ 
```

Apriori

D	
<i>t</i>	<i>i(t)</i>
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD



Apriori

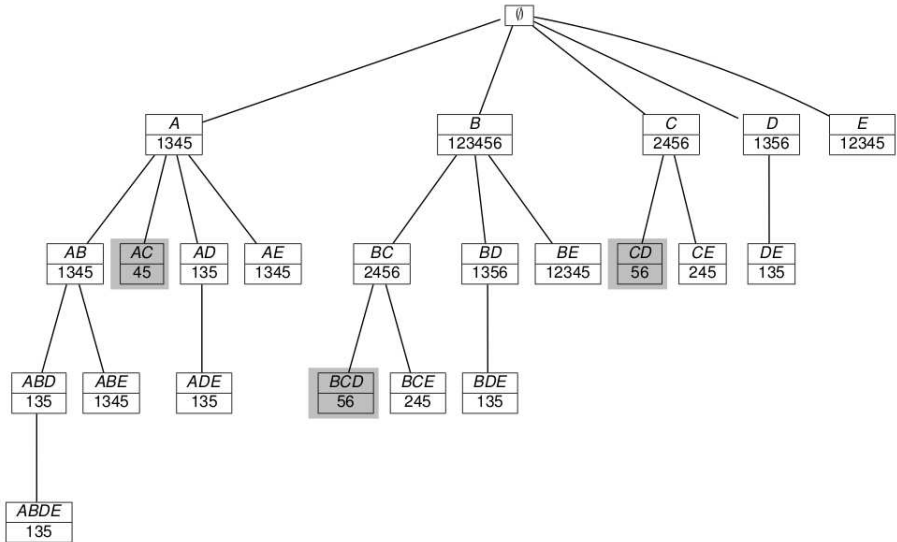


Eclat

Eclat

- ▶ Melhora o desempenho da contagem do suporte
- ▶ Método baseado na intersecção de *tidsets* e os indexa na estrutura de dados
- ▶ O suporte de um *itemset* candidato pode ser computado pela intersecção dos *tidsets* dos subconjuntos
 - ▶ Dados $t(X)$ e $t(Y)$, então $t(XY) = t(X) \cap t(Y)$
 - ▶ $sup(XY) = |t(XY)|$

Eclat



Eclat

```

// Initial Call:  $\mathcal{F} \leftarrow \emptyset, P \leftarrow \{\langle i, \mathbf{t}(i) \rangle \mid i \in \mathcal{I}, |\mathbf{t}(i)| \geq \text{minsup}\}$ 
ECLAT ( $P, \text{minsup}, \mathcal{F}$ ):
1 foreach  $\langle X_a, \mathbf{t}(X_a) \rangle \in P$  do
2    $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X_a, \text{sup}(X_a))\}$ 
3    $P_a \leftarrow \emptyset$ 
4   foreach  $\langle X_b, \mathbf{t}(X_b) \rangle \in P$ , with  $X_b > X_a$  do
5      $X_{ab} = X_a \cup X_b$ 
6      $\mathbf{t}(X_{ab}) = \mathbf{t}(X_a) \cap \mathbf{t}(X_b)$ 
7     if  $\text{sup}(X_{ab}) \geq \text{minsup}$  then
8        $P_a \leftarrow P_a \cup \{\langle X_{ab}, \mathbf{t}(X_{ab}) \rangle\}$ 
9   if  $P_a \neq \emptyset$  then ECLAT ( $P_a, \text{minsup}, \mathcal{F}$ )

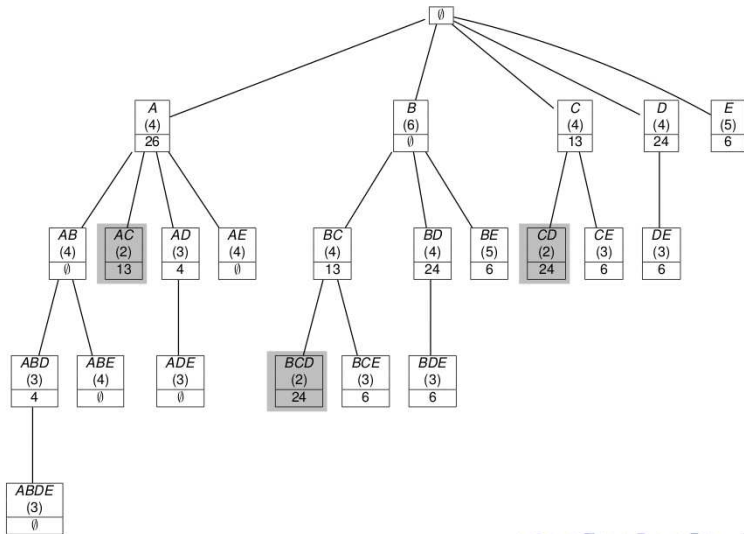
```

dEclat

dEclat

- ▶ Melhora o desempenho do *Eclat* reduzindo os *tidsets* intermediários
- ▶ Guarda as diferenças ao invés de todos os *tidsets*
- ▶ Sejam $X_a = \{x_1, \dots, x_j, x_a\}$ e $X_b = \{x_1, \dots, x_j, x_b\}$, então $X_{ab} = X_a \cup X_b = \{x_1, \dots, x_j, x_a, x_b\}$
- ▶ No processo de busca, pode-se guardar apenas as diferenças $d(X_{ab}) = t(X_a) \setminus t(X_{ab}) = t(X_a) \setminus t(X_b)$
- ▶ Pode-se determinar $d(X_{ab})$ com base nas diferenças como $d(X_{ab}) = d(X_b) \setminus d(X_a)$
- ▶ $sup(X_{ab}) = sup(X_a) - |d(X_{ab})|$
- ▶ Assim, é possível substituir as operações de intersecção do *Eclat* por operações de diferença entre conjuntos

dEclat



dEclat

```

// Initial Call:  $\mathcal{F} \leftarrow \emptyset$ ,
 $P \leftarrow \{ \langle i, \mathbf{d}(i), \text{sup}(i) \rangle \mid i \in \mathcal{I}, \mathbf{d}(i) = \mathcal{T} \setminus \mathbf{t}(i), \text{sup}(i) \geq \text{minsup} \}$ 
DECLAT ( $P$ ,  $\text{minsup}$ ,  $\mathcal{F}$ ):
1 foreach  $\langle X_a, \mathbf{d}(X_a), \text{sup}(X_a) \rangle \in P$  do
2    $\mathcal{F} \leftarrow \mathcal{F} \cup \{ (X_a, \text{sup}(X_a)) \}$ 
3    $P_a \leftarrow \emptyset$ 
4   foreach  $\langle X_b, \mathbf{d}(X_b), \text{sup}(X_b) \rangle \in P$ , with  $X_b > X_a$  do
5      $X_{ab} = X_a \cup X_b$ 
6      $\mathbf{d}(X_{ab}) = \mathbf{d}(X_b) \setminus \mathbf{d}(X_a)$ 
7      $\text{sup}(X_{ab}) = \text{sup}(X_a) - |\mathbf{d}(X_{ab})|$ 
8     if  $\text{sup}(X_{ab}) \geq \text{minsup}$  then
9        $P_a \leftarrow P_a \cup \{ \langle X_{ab}, \mathbf{d}(X_{ab}), \text{sup}(X_{ab}) \rangle \}$ 
10  if  $P_a \neq \emptyset$  then DECLAT ( $P_a$ ,  $\text{minsup}$ ,  $\mathcal{F}$ )

```

Gerando Regras de Associação

Gerando Regras de Associação

- ▶ Para cada *itemset* frequente $Z \in F$, deve-se gerar as regras $X \rightarrow Y$, onde $Y = Z \setminus X$
- ▶ A regra XY deve ser frequente ($\text{sup}(XY) > \text{minsup}$)
- ▶ Depois computa-se a confiança $\text{conf}(X \rightarrow Y) = \frac{\text{sup}(XY)}{\text{sup}(X)}$
 - ▶ $\text{conf}(X \rightarrow Y) \geq \text{minconf}$, ou seja, as confianças devem atender a um limiar de confiança mínima
 - ▶ Se $\text{conf}(X \rightarrow Y) < \text{minconf}$, então $\text{conf}(W \rightarrow Z \setminus W) < \text{minconf} \forall W \subset X$ pois $\text{sup}(W) \geq \text{sup}(X)$
 - ▶ Não é necessário investigar os subconjuntos de X

Gerando Regras de Associação

```

ASSOCIATIONRULES ( $\mathcal{F}$ , minconf):
1 foreach  $Z \in \mathcal{F}$ , such that  $|Z| \geq 2$  do
2    $\mathcal{A} \leftarrow \{X \mid X \subset Z, X \neq \emptyset\}$ 
3   while  $\mathcal{A} \neq \emptyset$  do
4      $X \leftarrow$  maximal element in  $\mathcal{A}$ 
5      $\mathcal{A} \leftarrow \mathcal{A} \setminus X$  // remove  $X$  from  $\mathcal{A}$ 
6      $c \leftarrow \text{sup}(Z) / \text{sup}(X)$ 
7     if  $c \geq \text{minconf}$  then
8       | print  $X \longrightarrow Y$ ,  $\text{sup}(Z)$ ,  $c$ 
9     else
10    |  $\mathcal{A} \leftarrow \mathcal{A} \setminus \{W \mid W \subset X\}$ 
    | // remove all subsets of  $X$  from  $\mathcal{A}$ 
  
```

Lift

- ▶ Uma vez encontradas as regras de associação que atendem às restrições impostas, essas podem ser ordenadas segundo um critério de interesse
- ▶ O *Lift* é uma medida de correlação
- ▶ $lift(X \rightarrow Y) = \frac{conf(X \rightarrow Y)}{sup(Y)} = \frac{sup(XY)}{sup(X)sup(Y)}$
 - ▶ se $lift(X \rightarrow Y) > 1$ então X e Y são positivamente correlacionados e a ocorrência de um implica na ocorrência do outro
 - ▶ se $lift(X \rightarrow Y) < 1$ então a ocorrência de X é negativamente correlata a Y
 - ▶ se $lift(X \rightarrow Y) = 1$, então X e Y são independentes ($P(XY) = P(X)P(Y)$)