

8. Übung „Knowledge Discovery“

Wintersemester 2008/2009

1 Hierarchische Assoziationsregeln

1. Seien A, A_1, A_2, B Items und A eine Verallgemeinerung von A_1 und A_2 . Beweisen oder widerlegen Sie (durch Angabe eines Gegenbeispiels) die folgenden Behauptungen:

a) $\text{support}(A \rightarrow B) = \text{support}(A_1 \rightarrow B) + \text{support}(A_2 \rightarrow B)$

Diese Implikation gilt nicht. Ein Gegenbeispiel bilden die Transaktionen $D = \{(A_1, A_2, B), (A_1, B), (A_2, B)\}$ für die gilt:

$$\text{support}(A \rightarrow B) = 1 \neq 4/3 = 2/3 + 2/3 = \text{support}(A_1 \rightarrow B) + \text{support}(A_2 \rightarrow B)$$

b) Falls $\text{support}(A_1 \rightarrow B) > \text{minsupport}$, dann ist $\text{support}(A \rightarrow B) > \text{minsupport}$.

Gilt, denn

$$\begin{aligned} \text{support}(A \rightarrow B) &= \frac{|(A \cup B)|}{|D|} = \frac{|(A_1 \cup A_2 \cup B)|}{|D|} \\ &\geq \frac{|(A_1 \cup B)|}{|D|} = \text{support}(A_1 \rightarrow B) > \text{minsupport} \end{aligned}$$

c) Falls $\text{support}(A \rightarrow B) > \text{minsupport}$, dann ist $\text{support}(A_1 \rightarrow B) > \text{minsupport}$.

Diese Richtung gilt nicht, wie man für $\text{minsupport} = 0.8$ und das in 1a gegebene D leicht sieht. Dort ist $\text{support}(A \rightarrow B) = 1$ aber $\text{support}(A_1 \rightarrow B) = 2/6 < 0.8$.

2. Wie berechnet sich die erwartete Konfidenz einer Regel $X \rightarrow Y$?

$$\text{erwarteteKonfidenz}(X \rightarrow Y) = \frac{\text{support}(X, Y)}{\text{support}(X)}$$

wären X und Y statistisch unabhängig voneinander, so wäre

$$\text{support}(X, Y) = \text{support}(X) \text{support}(Y)$$

und damit

$$\text{erwarteteKonfidenz}(X \rightarrow Y) = \frac{\text{support}(X) \text{support}(Y)}{\text{support}(X)}$$

$$= \text{support}(Y)$$

3. Gegeben folgende Menge von Transaktionen:

TID	Items
1	D,F,G
2	C,D,E,G
3	A,D,E
4	C,D,H
5	B,C,D,F,G
6	C,F,G
7	C,D,F
8	E,G

a) Geben Sie für einen Support von 25% die vom Apriori-Algorithmus generierten Mengen C_k und L_k an.

C_k und L_k :

- i. $L_1 = \{C, D, E, F, G\}$
- ii. $C_2 = \{CD, CE, CF, CG, DE, DF, DG, EF, EG, FG\}$
- iii. $L_2 = \{CD, CF, CG, DE, DF, DG, EG, FG\}$
- iv. $C_3 = \{CDF, CDG, CFG, \del{DEF}, DEG, DFG\}$
- v. $L_3 = \{CDF, CDG, CFG, DFG\}$
- vi. $C_4 = \{CDFG\}$
- vii. $L_4 = \{\}$

b) Bestimmen Sie Support und Konfidenz der Assoziationsregel $DF \rightarrow G$.

Support = 1/4, Konfidenz = 2/3

- c) Nehmen Sie an, dass die Items E, F und G zum Item X generalisiert sind. Bestimmen Sie den Support von $X \rightarrow C$. Geben Sie an, ob $X \rightarrow C$ und $E \rightarrow C$ R -interessant mit $R = 2$ sind. Begründen Sie ihre Antworten.

- $\text{Support}(X \rightarrow C) = 1/2$
- $X \rightarrow C$ ist R -interessant, da keine generellere Regel existiert
- $E \rightarrow C$: Support von $X = 7/8$, Support von $E = 3/8$, zu erwarten ist für $E \rightarrow C$ ein Support von $3/14$ ($= 1/2 * 3/7 = \text{support}(X, C) \frac{\text{support}(E)}{\text{support}(X)}$); tatsächlich beträgt der Support $1/8$, so dass die Regel nicht R -interessant ist.

2 TITANIC

Folgend ist der TITANIC-Algorithmus in einer Form wiedergegeben, die der Form des Apriori-Algorithmus aus der Vorlesung ähnelt.

```

1  TITANIC ( $I, D, \text{minsup}$ )
2    support( $\{\emptyset\}$ );
3     $L_0 := \{\emptyset\}$ ;
4     $k := 1$ ;
5    forall  $m \in M$  do  $\{m\}.p.s := \emptyset.s$ ;
6     $C := \{\{m\} \mid m \in M\}$ ; // Singletons
7    loop begin
8      support( $C$ ); // Support berechnen
9      forall  $X \in L_{k-1}$  do  $X.closure := \text{closure}(X)$ ;
10      $L_k := \{X \in C \mid X.s \neq X.p.s \wedge X.s \geq \text{minsup}\}$ ; // Pruning
11     if  $L_k = \emptyset$  then exit loop;
12      $k ++$ ;
13      $C := \text{Apriori-Kandidatengenerierung}(L_{k-1})$ ;
14 end loop;
15 return  $\bigcup_{i=0}^{k-1} \{X.closure \mid X \in L_i\}$ ;

```

1. Berechnen Sie für den Kontext auf Seite 12 der Vorlesungs-Folien (Kapitel 4, Teil 2 Begriffsverbände) die Begriffsinhalte mit einem Mindest-Support von 1 (also $33\frac{1}{3}\%$).

$\{\emptyset, \{a, c\}, \{b, d\}, \{c\}, \{b, c, d\}\}$

2. Stellen Sie den Zusammenhang zwischen Ihren Zwischenergebnissen und den auf der Folie 12 genannten Begriffen und Formeln her.

Die „dicken schwarzen Kästchen“ sind die minimalen Erzeuger, welche im Algorithmus in den Mengen L_k enthalten sind (in diesem Fall $\{\emptyset\}, \{\{a\}, \{b\}, \{c\}, \{d\}\}, \{\{b, c\}, \{c, d\}\}$). Die Begriffsinhalte ergeben sich als die Hüllen der minimalen Erzeuger.

3. Vergleichen Sie die einzelnen Schritte (Pruning, Kandidatengenerierung, Hüllen) des TITANIC mit denen des Apriori-Algorithmus.

Pruning: bei Apriori Entfernung von Teilmengen, die nicht minimalen Support haben, bei TITANIC zusätzlich Entfernung nicht minimaler Generatoren

Kandidatengenerierung: Join: identisch, Pruning: bei TITANIC zusätzlich Berechnung von $X.p_s := \min(X.p_s, S.s)$, wobei S $(k - 1)$ -elementige Teilmenge von X

Hüllen: bei TITANIC zusätzlich Berechnung der Hüllen (Closure)

4. Diskutieren Sie für jeden der angegebenen Teilschritte die möglichen Performance-Gewinne und -Verluste.

Da bei TITANIC durch Entfernen nicht minimaler Generatoren (10) vor der Kandidatengenerierung (13) zunächst meist weniger Kandidaten generiert werden, kann Rechenzeit bei der Kandidatengenerierung eingespart werden. Aus dem gleichen Grund terminiert TITANIC meist nach weniger Durchläufen (11) als Apriori.