

## Big Data / Data Mining

### Assoziationsanalyse

Bei der Assoziationsanalyse (*pattern mining*) wird in (häufig sehr großen) Datenmengen nach Strukturen gesucht, die ein häufiges Auftreten aufweisen, z.B.:

**häufige Artikelmengen** (*frequent itemsets*) bei Darstellung der Daten als Menge (Warenkörbe, ...)

**häufige Sequenzen** (*frequent sequences, frequent strings*) bei Darstellung der Daten als Sequenz (natürlichsprachige Texte, Ausleihhistorie in Bibliothek, ...)

**häufige Bäume/Graphen** (*frequent trees/graphs*) bei Darstellung der Daten als Baum/Graph (Verbreitung von Informationen, ...)

Aus den identifizierten Mustern lassen sich daraufhin **Assoziationsregeln** ableiten. Diese können verwendet werden, um Vorhersagen zu treffen oder Empfehlungen zu generieren. Die Warenkorbanalyse als eine der ersten Anwendungen der Assoziationsanalyse hat die Terminologie maßgeblich geprägt (*item, transaction, ...*).

#### Beispiel:

$$\{Chips, Windeln\} \rightarrow \{Bier\}$$

*Wer Chips und Windeln kauft, kauft häufig auch Bier.*

Eine Assoziationsregel hat also die Form  $A \rightarrow B$ , wobei  $A$  Prämisse oder Regelrumpf und  $B$  Konsequenz oder Regelkopf genannt werden. Gegeben ist eine Menge von Items  $I$  und eine Menge von Transaktionen  $T$ , die als *Datenbasis*  $D = (T_1, \dots, T_n)$  mit  $T_j \subset I$  zusammengefasst werden. Ein Item kann beispielsweise als ein Artikel und eine Transaktion als ein Einkauf/Warenkorb in einem Shop verstanden werden (vgl. Cleve und Lämmel 2014, S. 73 f., 223 f.; vgl. „Assoziationsanalyse“ 2008, S. 261).

#### Beispiel:

$$I = \{Bier, Chips, Cola, Wurst, Käse, Brot\}$$
$$D = \left\{ \begin{array}{l} (1, \{Bier, Cola\}), (2, \{Wurst, Chips\}), (3, \{Bier, Cola\}), (4, \{Brot, Wurst, Käse\}), \\ (5, \{Brot, Wurst, Käse, Bier\}), (6, \{Chips, Cola, Bier\}), (7, \{Brot, Wurst, Käse\}) \end{array} \right\}$$

Assoziationsregeln sind meist nur zu einem bestimmten Prozentsatz korrekt, da i.d.R. nicht *alle* Kunden entsprechend dieser Regel einkaufen. Wir benötigen daher **Interessantheitsmaße** um die **Qualität einer Regel zu bestimmen**.

**Definition:** Der Support einer Itemmenge  $X \subset I$  ist definiert als die relative Häufigkeit dieser Itemmenge in der Datenbasis:

$$\text{supp}(X) = \frac{|\{T \in D : X \subset T\}|}{|D|}$$

Wenn  $X$  mechte Teilmenge von  $T$  ist und  $T$  Element der Datenbasis ist

(1)

**Definition:** Der Support einer Assoziationsregel ist gleich dem Support der Vereinigung von Prämisse und Konsequenz der Regel:

$$\text{supp}(X \rightarrow Y) = \text{supp}(X \cup Y) = \frac{|\{T \in D : (X \cup Y) \subset T\}|}{|D|} \quad (2)$$

Somit ergibt sich der Support einer Assoziationsregel  $X \rightarrow Y$  als Anteil der Transaktionen  $T$  der Datenbasis  $D$ , die diese Regel erfüllen. Damit wird zum Ausdruck gebracht, ob eine Kombination von Items zur Bildung einer Assoziationsregel überhaupt von Bedeutung ist. Die Konfidenz (confidence) einer Regel gibt dann Aufschluss über die Stärke des Zusammenhangs bzw. das Ausmaß der Gültigkeit einer Regel  $X \rightarrow Y$ .

**Definition:** Die Konfidenz ist definiert als Anteil der Transaktionen, die sowohl  $X$  als auch  $Y$  beinhalten, an der Menge der Transaktionen, die die Prämisse  $X$  erfüllen.

$$\text{conf}(X \rightarrow Y) = \frac{\text{supp}(X \rightarrow Y)}{\text{supp}(X)} \quad (3)$$

Um die wertvollen von den weniger wertvollen Assoziationsregeln zu trennen, werden Schwellwerte für Konfidenz und Support eingeführt. Eine Assoziationsregel ist i.d.R. umso wertvoller, je größer Support und Konfidenz sind. Die Schwellwerte für minimale Konfidenz und minimalen Support werden i.d.R. vom Nutzer festgelegt.

**Beispiel (fort.):** Wir legen den minimalen Support  $S_{\min}$  auf 3 fest. Die folgenden Frequent Itemsets genügen  $S_{\min}$ :

0 Items	1 Item	2 Items	3 Items
$\emptyset: 7$	$\{\text{Bier}\}: 4$ $\{\text{Cola}\}: 3$ $\{\text{Wurst}\}: 4$ $\{\text{Käse}\}: 3$ $\{\text{Brot}\}: 3$	$\{\text{Brot, Käse}\}: 3$ $\{\text{Brot, Wurst}\}: 3$ $\{\text{Wurst, Käse}\}: 3$ $\{\text{Bier, Cola}\}: 3$	$\{\text{Brot, Wurst, Käse}\}: 3$

Die Konfidenz der Assoziationsregel  $\{\text{Brot, Käse}\} \rightarrow \{\text{Bier}\}$  z.B. beträgt:

$$\text{conf}(\{\text{Brot, Käse}\} \rightarrow \{\text{Bier}\}) = \frac{\frac{1}{7}}{\frac{3}{7}} = \frac{1}{3}$$

## Aufgabe 1

- a) Warum ist eine Assoziationsregel mit großem Support und großer Konfidenz nicht immer wertvoll bzw. interessant? Finden Sie ein Beispiel.
- b) Welche weiteren Interessantheitsmaße können verwendet werden?
- c) Welche Arten von Assoziationsregeln gibt es?
- d) Welche wichtigen Algorithmen zum Erzeugen von Assoziationsregeln gibt es?

## Aufgabe 2

Gegeben sind die folgenden 4 Items/Produkte (Milch, Bier, Cola, Saft) und die folgenden 8 Baskets/Warenkörbe:

B1 = {Saft}	B5 = {Milch, Saft, Bier}
B2 = {Saft, Milch, Bier}	B6 = {Bier, Milch}
B3 = {Bier, Cola}	B7 = {Milch, Saft}
B4 = {Cola, Saft, Milch, Bier}	B8 = {Bier, Saft}

- a) Der minimale Support  $S_{min}$  beträgt 3 Baskets. Geben Sie alle Frequent Itemsets an, die der Bedingung von  $S_{min}$  genügen.
- b) Wie groß ist die Konfidenz und der Interest für folgende Assoziationsregeln?:

R1 = {Saft} $\rightarrow$ {Bier}	R3 = {Milch, Saft} $\rightarrow$ {Bier}
R2 = {Milch, Bier} $\rightarrow$ {Cola}	R4 = {Cola} $\rightarrow$ {Bier}

## Aufgabe 3

Assoziationsanalyse mit RapidMiner

Laden Sie den Datensatz *Transactions* aus den Beispieldaten des Repository.

- a) Reichern Sie diesen Datensatz mit weiteren beschreibenden Informationen aus den Beispieldaten an und verbinden Sie entsprechende Tabellen.
- b) Erzeugen Sie mit einem geeigneten Algorithmus Assoziationsregeln aus diesen Datensätzen.

# Lösungen

1 a)  $\left[ \begin{array}{l} 85\% \text{ der Männer sind verheiratet} = \text{supp}(B) \\ 75\% \text{ der verh. Männer sind regelm. Bierkonsumenten} = \text{conf}(A \rightarrow B) \end{array} \right]$

$\left[ \begin{array}{l} 26\% \text{ der Männer sind schwere Raucher} = \text{supp}(B) \\ 38\% \text{ der Männer, die schwere Raucher sind, sind regelm. Bierkons.} = \text{conf}(A \rightarrow B) \end{array} \right]$

↳ Frage ist, Risiko für Männer mittleren Alters für atherosclerosis  $\rightarrow$  verheiratet sein wirklich entscheidend?  
 $\rightarrow$  echter Zusammenhang?

b) Lift  $\rightarrow \text{lift}(x \rightarrow y) = \frac{\text{supp}(x \cup y)}{\text{supp}(x) \cdot \text{supp}(y)}$

$\rightarrow L() > 1 \hat{=}$  Abhängigkeit

nach mehr:  $\rightarrow$  [https://michael.hahsler.net/research/association\\_rules/measures.html](https://michael.hahsler.net/research/association_rules/measures.html)

$\rightarrow$  gibt ein Maß der Abhängigkeit zwischen den Größen an

Correlation  $\rightarrow \text{corr}(A, B) = P(A \cup B) / P(A) \cdot P(B)$

Conviction  $\rightarrow \text{conv}(x \rightarrow y) = \frac{1 - \text{supp}(y)}{1 - \text{conf}(x \rightarrow y)} \Rightarrow$  Wkt.  $x$  erscheint ohne  $y$

c) single dimensional, multi dimensional, boolean, quantitativ

d) AIS,

↓  
 ineffizient  
 1) erzeugt frequent itemsets  
 2) erzeugt confidente und frequent assoziationsregeln

SETM

↓  
 wie AIS,  
 bsp auf SQL Basis

, Apriori

↓  
 erzeugt nur Assoziations-  
 regeln, wenn diese über einem  
 Threshold (supp, conf) liegen

, Apriori-TD, AprioriHybrid Quelle:

<http://www.ijsrp.org/research-paper-0513/ijrsrp-p17133.pdf>

1) erzeugt frequent itemsets (supp > Thresh)  
 2) Erzeugung der Assoziationsregeln (conf > Thresh)

2) a)

Items = { Saft, Milch, Bier, Cola }  $D = \{ B_1, \dots, B_8 \}$   $|D| = 8$

$\text{supp}(\text{Saft}) = \frac{6}{8} = \frac{3}{4} > S_{\min} = \frac{3}{8} \checkmark$

$\text{supp}(\text{Milch}) = \frac{5}{8} > S_{\min} = \frac{3}{8} \checkmark$

$\text{supp}(\text{Bier}) = \frac{6}{8} = \frac{3}{4} > S_{\min} = \frac{3}{8} \checkmark$

$\text{supp}(\text{Cola}) = \frac{2}{8} = \frac{1}{4} < S_{\min} = \frac{3}{8} \times$

b)

$\text{conf}(R_1) = \frac{\text{supp}(x \cup y)}{\text{supp}(x)} = \frac{\frac{4}{8}}{\frac{6}{8}} = \frac{1}{2} \cdot \frac{4}{3} = \frac{2}{3} = 66,6\%$

$\text{conf}(R_2) = \frac{\frac{1}{8}}{\frac{5}{8}} = \frac{1}{5} \cdot \frac{8 \cdot 2^1}{4 \cdot 1} = \frac{1}{4} = 25\%$

$\text{conf}(R_3) = \frac{\frac{3}{8}}{\frac{4}{8}} = \frac{3}{4} \cdot \frac{2^1}{1} = \frac{3}{4} = 75\%$

$\text{conf}(R_4) = \frac{\frac{2}{8}}{\frac{2}{8}} = 1 = 100\%$

$\text{lift}(R_1) = \frac{\text{supp}(x \cup y)}{\text{supp}(x) \cdot \text{supp}(y)} = \frac{\frac{4}{8}}{\frac{6}{8} \cdot \frac{5}{8}} = \frac{1}{2} < 1 \Rightarrow$  geringe Abhängigkeit

$\text{lift}(R_2) = \frac{\frac{1}{8}}{\frac{4}{8} \cdot \frac{2}{8}} = \frac{1}{1} = 1 \Rightarrow$  Abhängigkeit

$\text{lift}(R_3) = \frac{\frac{3}{8}}{\frac{4}{8} \cdot \frac{4}{8}} = \frac{3}{1} \cdot \frac{16}{8} = 2 > 1 \Rightarrow$  große Abhängigkeit

$\text{lift}(R_4) = \frac{\frac{2}{8}}{\frac{1}{4} \cdot \frac{2}{4}} = \frac{1}{1} \cdot \frac{16}{3} = 1,3 > 1 \Rightarrow$  große Abhängigkeit

# Literatur

„Assoziationsanalyse“ (2008). In: *Datenanalyse und Statistik*. Wiesbaden: Gabler, S. 261–272. ISBN: 978-3-8349-0434-8 978-3-8349-9654-1. DOI: 10.1007/978-3-8349-9654-1\_17. URL: [http://link.springer.com/10.1007/978-3-8349-9654-1\\_17](http://link.springer.com/10.1007/978-3-8349-9654-1_17) (besucht am 12.01.2020).

Cleve, Jürgen und Uwe Lämmel (2014). *Data mining*. Studium. München: De Gruyter Oldenbourg. 306 S. ISBN: 978-3-486-72034-1 978-3-486-71391-6.