

Prof. Dr. Katharina Morik,
Prof. Dr. Claus Weihs,
Dr. Wouter Duivesteijn,
M.Sc. Sarah Schnackenberg,
B.Sc. Melanie Dagge

Dortmund, 30.04.14
Abgabe: bis Do, 07.05.2015, an
wouter.duivesteijn@tu-dortmund.de
und/oder in den Briefkasten "Duivesteijn"
im OH12, R4.005

Übungen zur Vorlesung
Wissensentdeckung in Datenbanken
Sommersemester 2015

Blatt 3

Wiederholung Im Rahmen dieses Blattes sollten Sie einige Inhalte aus der Vorlesung wiederholen können. Dabei sollten Sie folgende Fragen beantworten können:

1. Was ist ein Verband? Geben Sie einen Beispiel-Verband anhand von Mengen mit der Teilmengenrelation als partieller Ordnung an!
2. Was ist die *Monotonie*-Eigenschaft im Bezug auf häufige Mengen und den *Apriori*-Algorithmus?
3. Geben Sie die zentrale Idee des FP-Growth-Algorithmus wieder und beschreiben Sie den Algorithmus kurz.
4. Was bewirkt die Sortierung der Itemsets der Transaktionen nach deren Häufigkeit?

Hinweis: Der FP-Growth-Algorithmus ist möglicherweise nicht leicht zu verstehen. Bei Verständnisproblemen sei auf das Originalpapier "*Mining Frequent Patterns without Candidate Generation*" (Han et. al., 1999) verwiesen, das auf der Web-Seite der Übungsblätter zu finden ist: www-ai.cs.uni-dortmund.de/LEHRE/VORLESUNGEN/KDD/SS15/exercises.html.

Aufgabe 3.1 (3 Punkte)

In der Vorlesung wurden mit Hilfe des Apriori-Algorithmus die häufigen Mengen in einer Transaktionsdatenbank gefunden. Gegeben sei die nachfolgende Aufstellung von Filmen, die von Zuschauern z_1, \dots, z_{10} besucht worden sind.

1. (1 Punkt) Formen Sie die Tabelle in eine Transaktionsdatenbank um!
2. (2 Punkte) Bestimmen Sie mit dem Apriori-Algorithmus die häufigen Mengen mit minimalem Support von $\frac{2}{5}$ und $\frac{3}{5}$. Geben Sie für jeden Schritt die Kandidatenmenge sowie die Menge der *large itemsets* (d.h. diejenigen Mengen, die den minimalen Support erfüllen) an.

| Titel | Jahr | z_1 | z_2 | z_3 | z_4 | z_5 | z_6 | z_7 | z_8 | z_9 | z_{10} |
|-----------------------------|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|
| Star Wars | 1977 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 |
| E.T. der Ausserirdische | 1982 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 |
| Indiana Jones | 1989 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 |
| Otto - der Ausserfriesische | 1989 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 |
| Wayne's World | 1992 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
| Bang Boom Bang | 1999 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 |
| Bridget Jones | 2001 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| Simpsons (Film) | 2007 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |

Aufgabe 3.2 (5 Punkte)

Diese Aufgabe behandelt den in der Vorlesung vorgestellten Algorithmus FP-Growth. Als Grundlage dient wieder die Datenbank aus Aufgabe 1. Es sei ein minimaler Support von $\frac{2}{5}$ gegeben, für den nun die häufigen Mengen in der Datenbank gefunden werden sollen.

- (1 Punkt) Geben Sie die Transaktionstabelle mit nach Häufigkeit sortierten Items (innerhalb der Transaktionen) an!
- (1 Punkt) Bestimmen Sie die Header-Tabelle sowie den *FP-Tree* aus der angegebenen Transaktionstabelle.
- (1 Punkt) Bestimmen Sie alle *conditional pattern bases* zum *FP-Tree*.
- (1 Punkt) Bestimmen Sie nun zu den *conditional pattern bases* die *conditional FP-Trees*.
- (1 Punkt) Bestimmen Sie anhand der *conditional FP-Trees* rekursiv die *frequent patterns*. Zeigen Sie die Erfassung der *frequent patterns* jeweils an der Entwicklung der *conditional pattern bases* sowie den *conditional FP-Trees*.

Aufgabe 3.3 (2 Punkte)

Die in den Übungen vorgestellte Software *RapidMiner* enthält sowohl eine Implementierung des *Apriori*-Algorithmus (WEKA-Extension) als auch einen Operator für das *FP-Growth* Verfahren. Die obige Datenbank-Tabelle finden sie als CSV-Datei unter:

<http://www-ai.cs.uni-dortmund.de/LEHRE/VORLESUNGEN/KDD/SS15/kino.csv>

- Installieren Sie die WEKA-Extension (Help/Updates and Extension)
- Laden Sie die CSV-Datei herunter und importieren Sie die Daten in ihr *RapidMiner* Repository!
- Laden Sie die Daten in ein *RapidMiner* Experiment und starten Sie das Experiment. Betrachten Sie die Daten in der Ergebnisansicht von *RapidMiner* und überlegen Sie sich, welche Attribute für den *Apriori*-Algorithmus benötigt werden.
- (2 Punkte) Erstellen Sie ein Experiment, dass die Daten liest und den Apriori-Algorithmus auf die Daten anwendet.

Hinweis: *Beachten Sie, dass der Apriori-Operator nur binäre Attribute verwenden kann! Sie benötigen dafür einen Operator, der Attribute in das gewünschte Ziel-Format konvertieren kann.*