

# Anwendung von Frequent Itemset Mining auf nutzergenerierte Geodaten

Christian Sengstock  
Universität Heidelberg  
Institut für Informatik  
Lehrstuhl für Datenbanksysteme  
sengstock@informatik.uni-heidelberg.de

Michael Gertz  
Universität Heidelberg  
Institut für Informatik  
Lehrstuhl für Datenbanksysteme  
gertz@informatik.uni-heidelberg.de

## 1. EINFÜHRUNG

Geodaten werden zunehmend durch Nutzer generiert. Openstreetmap (OSM), Google Earth, Qype sowie mit geographischen Attributen annotierte Nutzerinhalte wie Wikipedia-Artikel, Blogs und Tweets sind einige Beispiele. Die Datenformate und -modelle sind oft standardisiert (GML, KML) oder es handelt sich um einfache flache Modelle (GeoJSON, OSM-Format, zeilenbasierte Formate). Die Modelle erlauben dabei die freie Definition von Attributen (Metadaten, Sachdaten) in unterschiedlicher Mächtigkeit.

Wir stellen einen Ansatz und ein Framework vor, das die statistische Analyse, das Data Mining und die Integration solcher Datenmengen vereinfacht. Hierzu werden die Daten in ein flaches *Key-Value-Modell* überführt, wobei jedes *Key-Value-Paar* eine *Eigenschaft* darstellt. Auf Basis dieser *eigenschaftsbasierten* Sichtweise lassen sich Verfahren des *Frequent Itemset Mining* auf die Daten anwenden, um automatisiert Schema-Modelle zu identifizieren und zu analysieren.

Dieser Ansatz bildet eine Grundlage für Anwendungen und Verfahren wie Qualitätsanalyse, Schema-Integration, Clustering und Klassifikation, Visualisierung und die automatisierte Generierung von Ontologien.

In dieser Arbeit verwenden wir den genannten Ansatz, um die Qualität und die Konsistenz der frei editierbaren OSM-Daten automatisiert zu analysieren. Hierzu wird auf die ermittelten Schemata eine Konsistenzprüfung der Wertebereiche durchgeführt sowie ein Qualitätsmaß berechnet und auf Karten visualisiert.

## 2. VERWANDTE FORSCHUNGSTHEMEN

Die Analyse der Eigenschaften von Geodaten wird auf unterschiedlichen Ebenen untersucht. Im *frequent geographic pattern mining* wird das Auftreten von geographischen Mustern auf Basis geographischer Eigenschaften (*contains, intersects, near-by*) analysiert [2]. Suchmöglichkeiten auf Basis semantisch-geographischer Eigenschaften sind eine Forschungsfrage in Bezug auf ein *Geospatial Semantic Web* [3]. Speziell mit der Analyse und Bewertung von nutzergenerierten Annotationen im Rahmen von Verschlagwortungen haben sich Golder et al. [4] beschäftigt. In unserer Arbeit konzentrieren wir uns dagegen nicht auf die semantische Analyse, sondern auf Formen der automatisierten Schemaanalyse, um gemeinsam verwendete Schemata und Konzeptualisierungen in von Nutzern generierten Geodaten zu finden.

## 3. FREQUENT ITEMSET MINING

Das *Frequent Itemset Mining* wurde Anfang der 90er Jahre entwickelt, um Zusammenhänge von großen Warenkorb Datensätzen zu analysieren [1]. Hierbei wird von einer Transaktion ausgegangen (typischerweise eine Kauftransaktion), der eine Menge von Waren (*Items*) zugeordnet sind. Durch das Key-Value-Modell kann man ein Objekt auf das Transaktions-Waren-Modell abbilden. Das Objekt (identifiziert durch eine ObjektID) ist die Transaktion. Die Key-Value-Paare ent-

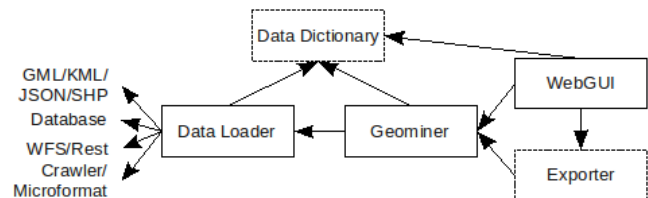


Abb. 1: Komponenten des Frameworks

sprechen den zugeordneten Waren.

Um das im Folgenden beschriebene Frequent Itemset Mining auf Geodaten anwenden zu können, wird jedes Geoobjekt, bestehend aus einer Geometrie und einer Menge an Metadaten in ein Transaktionsobjekt mit zugeordneter Geometrie und einer Menge aus Key-Value-Paaren (Eigenschaften) überführt:

Geo-Transaktionsobjekt(*Geometrie, Liste((Key,Value))*)

Die Abbildung des *Frequent Itemset Mining* auf die Attribute von Datenmodellen nennen wir *Frequent Property Mining*, eine gefundene häufige Menge von Eigenschaften ein *Frequent Property Set*.

Um häufig auftretende Mengen von Key-Value-Paaren zu ermitteln, wird ein *Frequent Itemset Mining* auf die Geo-Transaktionsobjekte mit den zugeordneten Key-Werten ausgeführt. Hierfür haben wir einen leicht angepassten *Apriori-Algorithmus* [1] verwendet. Das Ergebnis sind Mengen aus Keys (*Frequent Property Set*), wobei jeder Menge seine Häufigkeit (*Support*) im Datensatz *D* zugeordnet ist ( $Anzahl/|D|$ ). Beispielsweise:

```
(highway) 0.3
(highway name) 0.25
(highway access) 0.2
(highway access name) 0.18
```

Die zu einem Frequent Property Set gehörenden Geoobjekte lassen sich als individuelle Datenmengen weiterverarbeiten, visualisieren und exportieren. Für das im Folgenden vorgestellte Framework bildet die Generierung der Frequent Property Sets die Grundlage für weitere Analysen.

## 4. FRAMEWORK KOMPONENTEN

Der *Data Loader* (siehe Abb. 1) liest die Geodaten und überführt sie in das Key-Value-Modell. 1:N Beziehungen müssen dabei auf Listen abgebildet werden. Um unterschiedliche Bezeichnungen von Keys zu unterstützen, können alternative Bezeichner in einem *Data Dictionary* abgelegt werden. Das *Data Dictionary* kann als Ontologie weiterentwickelt werden, um semantische Beziehungen zwischen den Attributen abbilden zu können.

Der *Geominer* liest die im Key-Value-Modell verfügbaren Daten und führt zunächst ein einfaches Frequent Property Mining durch. Dieses bildet die Basis für weitere Verfahren, wie statistische Analysen (insbesondere Histogramme und Wertebereiche) oder Clustering-Verfahren. Der *Geominer* kann auf das *Data Dictionary* zugreifen, um alternative Bezeich-

