

Was ist Los?

Erhebung und Visualisierung von empirischen Daten verschiedener Internet-Zeitungen

Halbjahresarbeit 2016/2107 von Jaro Habiger

Inhaltsverzeichnis

Idee	1
Umsetzung	1
Daten Sammeln	1
Speicherung	3
Stichpunkte	4
Vorgehen	5
Probleme	5
Leztendliche Umsetzung	5
Beobachtungen	5
Evaluation des Verfahrens	5
Verschiedene Zeitungen im Vergleich	5
Fazit	5

Idee

Umsetzung

Die Umsetzung lässt sich sehr gut in einzelne Teilprobleme unterteilen. Hierbei ist es am einfachsten, den Datenfluss von den verschiedenen Nachrichtenquellen zur fertigen Visualisierung zu betrachten. In den nachfolgenden Abschnitten wird dieser Verarbeitungsprozess beschrieben.

Daten Sammeln

Als erstes müssen Daten zur weiteren Verwertung von den verschiedenen Nachrichtenquellen gesammelt werden. Dies geschieht über die sogenannten "RSS-Feeds". Bei diesen handelt es sich um ein standardisiertes Format, über das Nachrichtenanbieter ihre Artikel, inklusive Metadaten wie z.B. den Zeitpunkt der Veröffentlichung, in maschinenlesbarer Form bereitstellen. Im Prinzip werden also die gleichen Daten bereitgestellt, wie auf der normalen Internetseite, mit dem Unterschied, dass sie einfacher mit Programmen verarbeitbar sind.

Diese Darstellung als RSS-Feed ermöglicht es, die Artikel verschiedener Online-Zeitungen zu betrachten, ohne für jede einen komplett neuen Datensammler programmieren zu müssen.

Der letztendliche Datensammler ist ein Programm, welches ich in Python geschrieben habe. Es lädt sich den RSS-Feed einer einzelnen Nachrichtenquelle periodisch herunter und verarbeitet ihn weiter. Aktuell wird jeder RSS-Feed alle 10s neu analysiert. Im ersten Schritt der Verarbeitung lädt das Programm den Feed als Datei von den Servern der jeweiligen Online-Zeitung herunter. Diese Datei ist eine sogenannte XML-Datei. In ihr werden Daten als Baumstruktur abgebildet. Hierbei muss man sich die Datei als "Stamm" des Baums vorstellen. Die ersten "Äste" der Baumstruktur beinhalten die Metadaten, wie z.B. den Erstellungszeitpunkt und den Herausgeber des Feeds. Die nachfolgenden "Äste" enthalten jeweils einen Artikel. Diese Artikel wiederum beinhalten verschiedene Unterelemente, d.h. die "Äste" verzweigen sich in weitere, kleinere "Ästchen". In diesen stehen nun z.B. der Autor des Textes, der Titel, oder eben der eigentliche Inhalt des Textes. Diese textuelle Repräsentation einer

SPIEGEL ONLINE

DER SPIEGEL

SPIEGEL TV

Q

Anmelden

PANORAMA

Justiz

Leute

Gesellschaft

Multimedia-Reportagen

Heimwerkerblog



Gut zwei Promille

Ostfriesen sucht Reeperbahn in Braunschweig

Er lag um fast 200 Kilometer daneben: Ein schwer alkoholisierten Mann aus Ostfriesland hat die Hamburger Reeperbahn in Braunschweig gesucht. Die Polizei half zumindest bei der Ausnüchterung.

mehr...

Fig. 1: Ein Artikel einmal in der normalen Darstellung...

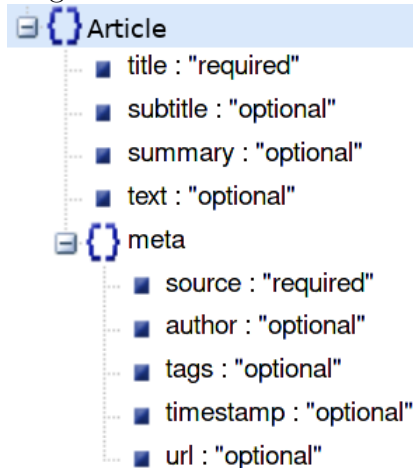
```

<?xml version="1.0" encoding="UTF-8" ?>
<rss xmlns:content="http://purl.org/rss/1.0/modules/content/" version="2.0">
  <channel>
    <title>SPIEGEL ONLINE - Schlagzeilen</title>
    <link>http://www.spiegel.de/</link>
    <description>
      Deutschlands führende Nachrichtenseite. Alles Wichtige aus Politik,
      Wirtschaft, Sport, Kultur, Wissenschaft, Technik und mehr.
    </description>
    <language>de</language>
    <pubDate>Sun, 11 Dec 2016 20:15:15 +0100</pubDate>
    <lastBuildDate>Sun, 11 Dec 2016 20:15:15 +0100</lastBuildDate>
  </channel>
  <image>
    <title>SPIEGEL ONLINE</title>
    <link>http://www.spiegel.de/</link>
    <url>http://www.spiegel.de/static/sys/logo_128x61.gif</url>
  </image>
  <item>
    <title>
      Gut zwei Promille: Ostfriesen sucht Reeperbahn in Braunschweig
    </title>
    <link>
      http://www.spiegel.de/panorama/gesellschaft/ostfriesen-sucht-reeperbahn-in-
      braunschweig-a-1125427.html#ref=RSS
    </link>
    <description>
      Er lag um fast 200 Kilometer daneben: Ein schwer alkoholisierten Mann aus
      Ostfriesland hat die Hamburger Reeperbahn in Braunschweig gesucht. Die
      Polizei half zumindest bei der Ausnüchterung.
    </description>
    <category>Panorama</category>
    <pubDate>Sun, 11 Dec 2016 19:55:39 +0100</pubDate>
    <guid>
      http://www.spiegel.de/panorama/gesellschaft/ostfriesen-sucht-reeperbahn-in-
      braunschweig-a-1125427.html
    </guid>
    <content:encoded>
      <![CDATA[
        Er lag um fast 200 Kilometer daneben: Ein schwer alkoholisierten Mann aus
        Ostfriesland hat die Hamburger Reeperbahn in Braunschweig gesucht. Die
        Polizei half zumindest bei der Ausnüchterung.
      ]]>
    </content:encoded>
  </item>
</rss>

```

Fig. 2: ... und einmal als Teil eines RSS-Feeds

Baumstruktur gilt es nun zunächst in eine einfacher verwendbare native representation im Speicher des Python Programms umzuwandeln. Hierfür wird eine Programmbibliothek verwendet, die einen sogenannten XML-Parser beinhaltet. Nach diesem schritt können die Einzelnen Artikel betrachtet werden. Hierbei wird zu aller erst überprüft, welche Artikel schn einmal verarbeitet wurden. Diese werden Verworfen. Die übriggebliebenen, also neuen Artikel werden in eine Allgemeinere Repräsentation für Neuigkeiten und ihre Metadaten gebracht, die ich mir überlegt habe. Diese ist auch wieder eine Baumstruktur



und sieht wie folgt aus:

Diese umwandlung ist nötig, da der RSS-Standart zwar die grobe Struktur und ihre representation als XML-Datei spezifiziert, letztenendes jeder Nachrichtenanbieter allerdings doch nicht ganz Kompatible Feeds ausliefern. Diese Kleinen unterscheide werden Hier also angeglichen, damit die weitere Verarbeitung leichter von statten gehen kann, und keine unterschiede mehr beachtet werden müssen.

Am ende der Datenaquirierung werden die gesammelten Daten an die nächste stufe weitergegeben: Die Speicherung.

Speicherung

Die nun gesammelten Daten müssen vor der weiteren verwendung ersteinmal strukturiert zwischengespeichert werden. Dies geschieht in einer Datenbank. Ich habe mich dafür entschieden, eine sogenannte "noSQL-Datenbank" zu verwenden, da

Stichpunkte

- Docker!!!1!
- build + tests
- Klare Unterteilung:
- Datensammler
 - In Python geschrieben
 - Crawlen hauptsächlich rss-feeds
- MongoDB
 - noSQL-Datenbank
 - Map/Reduce-Querys
 - zuerst Mongo, dann Couch, dann Mongo
- Http Middleware
 - stellt MongoDB über HTTP bereit
 - Clinets können eigene Map/Reduce Anfragen an die Datenbank stellen
 - * langsam
 - * flexibel
 - -> kleine datenmengen, daher ok
- Frontend
 - 3d.js visualisierung
 - einfache API zur DB
 - einfache umstrukturierbarkeit

Vorgehen

Probleme

Leztendliche Umsetzung

Beobachtungen

Evaluation des Verfahrens

Verschiedene Zeitungen im Vergleich

Fazit