

Was ist Los?

Erhebung und Visualisierung von empirischen Daten verschiedener Internet-Zeitungen

Halbjahresarbeit 2016/2107 von Jaro Habiger

Inhaltsverzeichnis

1	Idee.	1
2	Umsetzung	1
	2.1 Daten sammeln	1
	2.2 Speicherung	4
	2.3 map/reduce	4
	2.4 Datenbank	7
	2.5 Middleware	7
	2.6 Visualisierung	7
	2.7 Zusammenführung	8
3	Beobachtungen	8
	3.1 Evaluation des Verfahrens	8
	3.1.1 Ergebniss	8
	3.1.2 Technisch	8
	3.2 Verschiedene Zeitungen im Vergleich	8
4	Fazit	8

1 Idee

Relativ schnell war mir klar, das ich mich mit Datenvisualisierung beschäftigen will. Anfangs war allerdings nicht ganz klar, welchen Datensatz ich auf welche Aspekte hin untersuchen will. Eines tages nahm ich dann eine Ausgabe des Spiegels zur Hand und sah auf der ersten seite eine sogenannte Wortwolke. Diese form der Darstellung

Am Anfang dieser Halbjahresarbeit stand die Idee, die aktuelle Nachrichtenlage visuell sichtbar zu machen. Dies sollte möglichst intuitiv und aufschlussreich sein. Außerdem

2 Umsetzung

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

2.1 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 verarbeitet werden können.

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

Der letztendliche Datensammler ist ein Programm, welches ich in Python geschrieben habe. Python ist eine Programmiersprache, die auf Einfachheit

SPIEGEL ONLINE

DER SPIEGEL SPIEGEL TV

Anmelden

PANORAMA

- Justiz
- Leute
- Gesellschaft
- Multimedia-Reportagen
- Heimwerkerblog

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

```
<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>
      Bundeslands 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>
    <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: Ostfrösche suchen Reeperbahn in Braunschweig
      </title>
      <link>
        http://www.spiegel.de/panorama/gesellschaft/ostfrösche-sucht-reeperbahn-in-
        braunschweig-a-1125427.html#ref=rss
      </link>
      <description>
        Er lag um fast 200 Kilometer daneben: Ein schwer alkoholisierter 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>
    </item>
    <item>
      <title>
        http://www.spiegel.de/panorama/gesellschaft/ostfrösche-sucht-reeperbahn-in-
        braunschweig-a-1125427.html
      </title>
      <description>
        Er lag um fast 200 Kilometer daneben: Ein schwer alkoholisierter Mann aus
        Ostfriesland hat die Hamburger Reeperbahn in Braunschweig gesucht. Die
        Polizei half zumindest bei der Ausnüchterung.
      </description>
    </item>
  </channel>
</rss>
```

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

und Flexibilität optimiert ist. Das Datensammelprogramm lädt sich den RSS-Feed einer einzelnen Nachrichtenquelle periodisch herunter und verarbeitet ihn weiter. Ich habe mich dazu entschieden, jeden RSS-Feed alle 10s neu zu analysieren. 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 Author des Textes, der Titel, oder eben der eigentliche Inhalt des Textes. Diese Textuelle Repräsentation einer 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 aus wie in ??ir).

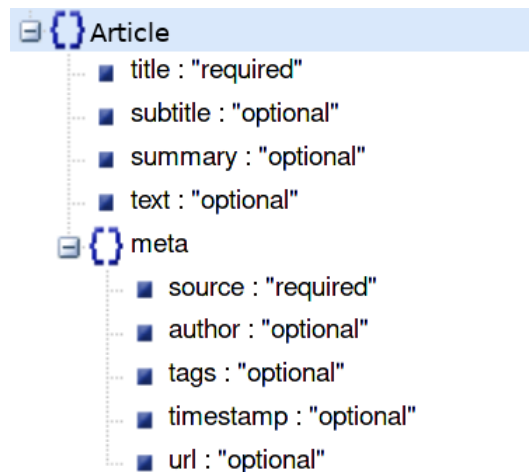


Fig. 3: Die Zwischenrepräsentation der Artikel

Diese umwandlung ist nötig, da der RSS-Standard zwar die grobe Struktur und ihre representation als XML-Datei spezifiziert, letztendlich jeder Nachrichtenanbieter allerdings doch nicht ganz Kompatible Feeds ausliefern. Diese kleinen Unterschiede werden hier also angeglichen, damit die weitere Verarbeitung leichter von statten gehen kann, und keine Unterschiede mehr beachtet werden müssen. Zusätzlich hierzu enthält der RSS-Feed nur Kurzfassungen der Artikel. In der Umwandlung werden zusätzlich hierzu noch die gesamten Artikel von der Internetseite heruntergeladen und der Volltext aus dieser extrahiert. Hierzu wird zuerst das HTML Dokument der Seite geparsed und die relevanten Teile mithilfe sogenannter "CSS-Selektoren" identifiziert. Hiernach werden alle HTML-Elemente dieser Bereiche in reinen Text umgewandelt.

Architektonisch ist der Datensammler ein Programm, in das weitere kleinere Module "hereingesteckt" werden, die die einzelnen Nachrichtenseiten ansprechen.

Am Ende der Datenaquirierung werden die gesammelten Daten an die nächste Stufe weitergegeben: Die Speicherung.

2.2 Speicherung

Die nun gesammelten Daten müssen vor der weiteren Verwendung erst einmal strukturiert zwischengespeichert werden. Dies geschieht in einer Datenbank. Ich habe mich dafür entschieden, eine sogenannte "noSQL-Datenbank" zu verwenden, da diese flexiblere Abfragemethoden ermöglichen. "noSQL" steht hierbei dafür, dass die Datenbank nicht über die Datenbanksprache SQL angesprochen wird, was normalerweise der Fall ist, sondern eine andere Schnittstelle bietet. Relativ früh habe ich mich dafür entschieden, dass ich meine Datenbankabfragen als "map/reduce" Funktionen formulieren möchte.

2.3 map/reduce

Das map/reduce Verfahren hat einige große Stärken, wie z.B. die hohe Parallelisierbarkeit und die damit verbundene hohe Geschwindigkeit, bei gleichzeitig hoher Flexibilität. MapReduce ist ein von der Firma Google eingeführtes Programmiermodell, welches wie folgt funktioniert:

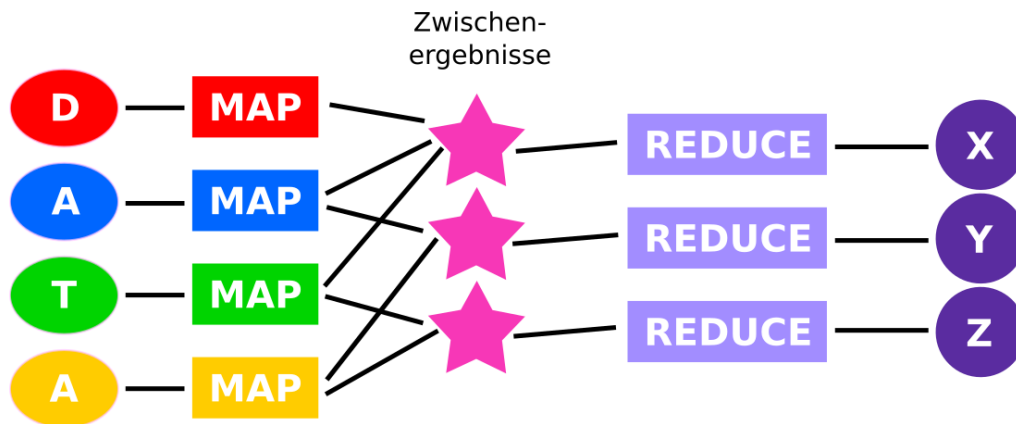


Fig. 4: Der map/reduce prozess als Grafik

1. Am Anfang des Prozesses steht eine Menge aus N Eingangsdatensätzen. In meinem Fall sind das die Zeitungsartike, die als Objekte mit der oben beschriebenen Datenstruktur vorliegen. Für jeden dieser Artikel wird jetzt die map funktion ausgeführt. Diese ordnet jedem Artikel N Schlüssel-Wert Paare zu. In dem Fall, das wir die zu jedem vorkommenen Wort die Häufigkeit bestimmen wollen, wordnet die map Funktion also jedem Artikel die Menge der Darin enthaltenen Wörter zu. Da diese Funktion auf jeden Artikel angewnd wird, kann man sich in diesem Fall den gesamten map Prozess als eine Funktion vorstellen, in die die Menge aller Artikel der Menge der darin enthaltenen Worte zuordnet. Der hierzu zugehörige Code der map funktion ist:

```
function map() {
  var text = this.text.replace(/[~A-Za-zÄäÖöÜüß ]/g, " "); // entferne alle ü
  var words = text.split(" "); // Teile den text in Wörter
  for (var i = 0; i < words.length; i++) { // für jedes Wort
    var word = words[i];
    if(word) {
      emit(word, 1); // füge der Ausgabemenge das Wort hinzu
    }
  }
}
```

2. Im nächsten Schritt werden Elemente mit Gleichen Schlüsseln Gruppirt.

Nach diesem Schritt also eine Menge aus Schlüssel-Wertmengenpaaren (Auch wenn ich immer wieder das Wort Menge verwende, meine ich eigentlich Multimengen, da es relevant ist, wie oft ein bestimmtes Element in der Menge vorkommt). Dieser Schritt ist, anders als die anderen beiden Schritte bei allen anderen Abfragen gleich. Da in unserem Beispiel das Wort als Schlüssel verwendet wurde sähe eine beispielhafte Menge nach diesem schritt wie folgt aus:

```
{
  "der": {1; 1; 1; 1; 1; 1; 1; 1; 1; 1};
  "die": {1; 1; 1; 1};
  ...
}
```

3. Im dritten und letzten Schritt wird für jeden Schlüssel die sogenannte “reduce” funktion angewandt. Diese reduziert die Mengen, die den Schlüsseln zugeordnet sind auf einen Wert. In unserem Beispiel fällt die reduce-funktion relativ einfach aus, da sie einfach nur die elemente der Menge aufsummieren muss:

```
reduce(key, values) {
  return values.reduce((previousValue, currentValue) => currentValue + previousValue)
}
```

4. Der vierte und letzte schritt gehört eigentlich nicht mehr zum map/reduce verfahren. Dieser wird nach diesem ausgeführt und dient dazu die ergebnisse zu sortieren und eventuell feinheiten zu verbessern. So ist es zum beispiel möglich selten genutzte- oder sogenannte stoppworte in diesem schritt auszusortieren:

```
function filter(data) {
  return data.filter(word => stopwords.indexOf(word["_id"].toLowerCase()) < 0)
}
```

Des weiteren wird in diesem schritt versucht, Wörter mit dem gleichen Stamm, und somit mit der gleichen Bedeutung zusammenzuführen, auch wenn diese unterschiedliche endungen haben. Ein beispiel hierfür ist, das die Wörter “Trump” und “Trumps” zusammengezählt werden. Hierbei wird immer das kürzeste der worte behalten, da dies meist die Grundform ist.

2.4 Datenbank

Nachdem klar war, wie die Abfragen formuliert werden sollten, habe ich verschiedene Datenbanken in betracht gezogen. Zuerst habe ich angefangen mit der Datenbank “MongoDB” zu arbeite, störte mich aber shr stark an der Komplexität von dieser. Außerdem bietet MongoDB keine möglichkeit, diese über das HTTP Protokoll anzusprechen, was für die Visualisierung allerdings sehr wichtig ist. Also sah ich mich nach anderen Alternativen um und fand “CouchDB”, welche vom Apache Projekt entwickelt wird. Diese erfüllte die meisten meiner anforderungen relativ gut, weshalb ich meinen gesamten bis dahin existierenden Code an CoucDB anpasste. Nach einiger Zeit des Testens stellte sich allerdings heraus, das CouchDB wahrchienlich zu langsam sein würde und ein unpassendes Rechteverwaltungssystem mitbringt, was die Arbeitersparniss im gegensatz zu MongoDB zu nichte machen würde. In Folge dieser Erkenntniss entschied ich mich dazu, meine Software zurück auf MongoDB zu portieren und für diese Eine HTTP schnitschtelle zu schreiben.

2.5 Middleware

Diese HTTP API ist in gewissermaßen das Bindeglied zwischen der Datenbank und der Visualisierung, weshalb es im nacholgenden “Middleware” genannt wird. Sie nimmt HTTP anfragen vom Frontend entgegen, leitet diese an die Datenbank weiter und schickt die Antwort an das Frontend zurück. In sofern könnte man sie auch als übersetzer zwischen verschiedenen Protokollen verstehen. Diese ist notwendig, da ich die Visualisierung im Browser implementieren möchte und im Browser nur sehr wenige Protokolle verfügbar sind. HTTP ist eines dieser wenigen Möglichkeiten und bietet sich an, da es verhältnissmäßig einfach zu implementieren ist.

2.6 Visualisierung

Ich habe mich entschieden die Visualisierung im Browser zu schreiben. Dies

- d3.js
- Flexibel

2.7 Zusammenführung

Um die einzelnen Teilkomponenten oder auch “Microservices” zusammenzuführen habe ich Docker verwendet.

- Mehr details

3 Beobachtungen

3.1 Evaluation des Verfahrens

3.1.1 Ergebniss

3.1.2 Technisch

- Map/Reduce bei jedem Query skaliert nicht :(

3.2 Verschiedene Zeitungen im Vergleich

4 Fazit