

Was ist los?

Erhebung und Visualisierung von empirischen Daten verschiedener Internet-Zeitungen

Halbjahresarbeit 2016/2107 von Jaro Habiger

Contents

1	Idee	1
2	Die Auswahl der Zeitungen	2
	2.1 Ein Wort über Bild.de	2
3	Konzeptioneller Aufbau	2
4	Umsetzung	3
	4.1 Daten sammeln	3
	4.2 Speicherung	7
	4.3 MapReduce	7
	4.4 Datenbank	9
	4.5 Middleware	10
	4.6 Visualisierung	10
	4.7 Zusammenführung	11
5	Beobachtungen	13
	5.1 Evaluation des Verfahrens	13
	5.2 Verschiedene Zeitungen im Vergleich	13
6	Fazit	13
7	Quellen	13

1 Idee

- ich -> Programmieren
- ich -> politisch
- eig TR website
- ich -> Zeitung -> klo -> spiegel -> wortwolke
- wie funktionieren Zeitungen?
- anspruchsvoll
- -> was denken medien
- anwendungsmöglichkeiten

In meiner Freizeit programmiere ich gerne. Ich mag es, Probleme mit Computern oder anderen technischen Geräten zu lösen. Auch gewinne ich hierbei oft neue Erkenntnisse, da man sich beim Programmieren unkonventionell mit den Dingen beschäftigen muss: Man muss mehr in Strukturen denken, muss versuchen einem Computer, der kein Verständnis von Sinn oder Inhalt hat, das Problem nur in logischen Zusammenhängen zu beschreiben. Somit kann mithilfe eines Computerprogramms fast nie eine wirkliche Lösung für ein Problem oder eine Antwort auf eine Frage gefunden werden, allerdings können Computer einem Menschen helfen, eine bestimmte Aufgabe besser, schneller oder einfacher zu erledigen.

Besonders interessant finde ich es immer dann, diese Macht, die wir uns mit der Programmierung von Computern geben können, zu verwenden, wenn wir sie nutzen um neue Einblicke in unsere Umwelt zu erlangen. Auch hier gilt zwar wieder, dass der Computer eigentlich viel dümmer als wir ist, und keine Zusammenhänge begreifen kann, allerdings stumpfe Aufgaben viel schneller und ausdauernder zu erledigen vermag. Diese Kombination aus Eigenschaften ermöglicht es, große Datensätze mit mathematischen Methoden zu analysieren und die Ergebnisse für Menschen ansprechend aufzubereiten. Bei dieser Aufbereitung bietet sich eine optische Darstellung an, da die Augen des Menschen das Sinnesorgan sind, das am schnellsten Daten erfassen kann [spiegelminig]. Es bietet sich also an, empirische Daten grafisch aufzubereiten, sprich Datenvisualisierung zu betreiben.

Relativ schnell war mir also klar, dass 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 fand ich sofort sehr ansprechend.

2 Die Auswahl der Zeitungen

2.1 Ein Wort über Bild.de

3 Konzeptioneller Aufbau

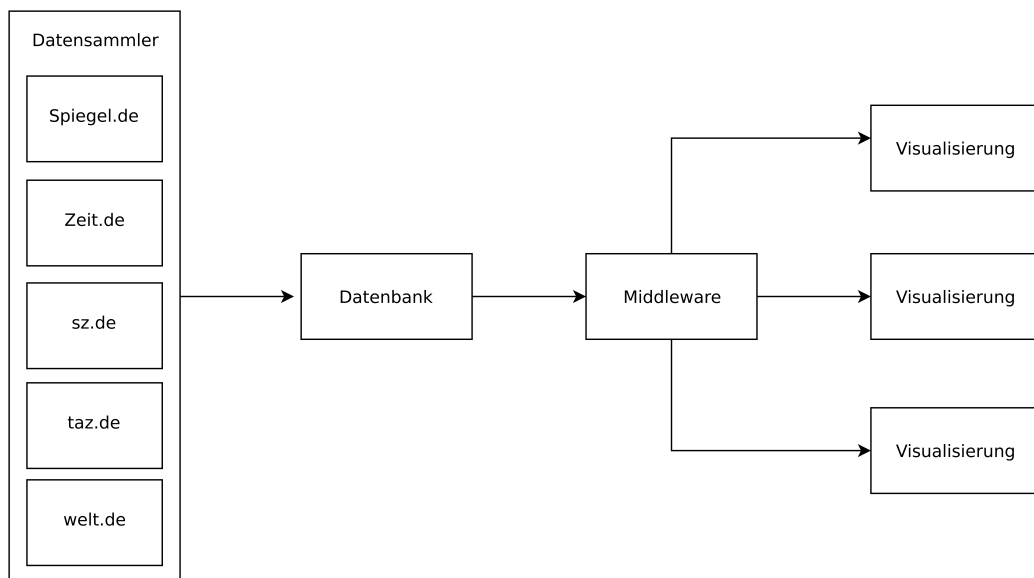


Abb. 1: Der Aufbau des Gesamtsystems

Nun gilt es sich ein System auszudenken, welches die Daten der Ausgewählten Zeitungen periodisch herunterlädt und die gesammelten Ergebnisse für spätere Analyse abspeichert. Dieser schritt ist ein sehr wichtiger, der gut bedacht sein muss, da

more

Dieses System sieht am Ende aus wie in Abb. 1.

4 Umsetzung

Der nun konzipierte Aufbau 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. Die einzelnen Schritte sind teilweise eher schwierig zu verstehen, da ihr Sinn wenig greifbar ist. Dies liegt daran, dass viele verschiedene Repräsentationen und Strukturen von Daten ineinander umgewandelt werden müssen.

4.1 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 sind also die gleichen Daten wie auf der normalen Internetseite abrufbar, mit dem Unterschied, dass diese einfacher mit Hilfe von Programmen verarbeitet werden können (Siehe Abb. 2 und 3).

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 Datensammelprogramm ist eine Software, welche ich in Python geschrieben habe. Python ist eine Programmiersprache, die auf Einfachheit 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 einmal pro Minute 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



Abb. 2: 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>
      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>
    ▼<image>
      <title>SPIEGEL ONLINE</title>
      <link>http://www.spiegel.de</link>
      <url>http://www.spiegel.de/static/sys/logo_120x61.gif</url>
    </image>
    ▼<item>
      ▼<title>
        Gut zwei Promille: Ostfrieze sucht Reeperbahn in Braunschweig
      </title>
      ▼<link>
        http://www.spiegel.de/panorama/gesellschaft/ostfrieze-sucht-reeperbahn-in-
        braunschweig-a-1125427.html#ref=rss
      </link>
    </item>
  </channel>
</rss>

```

Abb. 3: ... und einmal als Teil eines RSS-Feeds

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. Autor, Titel, oder eine Kurzfassung des Textes. Außerdem ist ein Link zur Vollversion des Artikels gegeben. Diese textuelle Repräsentation einer Baumstruktur gilt es nun in eine einfacher verwendbare Repräsentation im Speicher des Python Programms umzuwandeln, also zu “parsen”. Das Wort “parsen” ist sehr schwierig ins Deutsche zu übersetzen. Es beschreibt den Sachverhalt, die textuelle Repräsentation eines strukturierten Datensatzes in die “native” Repräsentation einer Programmiersprache für diesen Datensatz zu bringen. Die “native” Repräsentation eines Datensatzes in einer Programmiersprache ist die Standardform, diesen Datensatz abzubilden. Diese hat oft bestimmte Funktionen, die das Programmieren vereinfachen, wie z.B. die Zugriffsmöglichkeit auf einzelne Elemente bei Listen. Hierfür wird eine Programmibibliothek verwendet, die einen sogenannten “XML-Parser” beinhaltet, der genau dies erreicht. Nach diesem Schritt können die einzelnen Artikel betrachtet werden. Hierbei wird zuallererst überprüft, welche Artikel schon 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 Abb. 4.

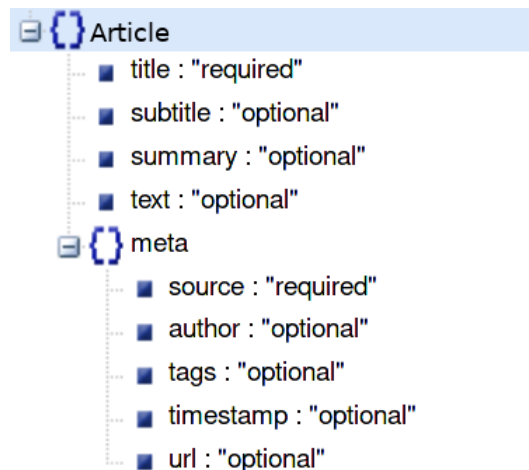


Abb. 4: Die Zwischenrepräsentation der Artikel

Diese Umwandlung ist nötig, da der RSS-Standard zwar die grobe Struktur und ihre Repräsentation als XML-Datei spezifiziert, letzten Endes jeder Nachrichtenanbieter allerdings doch nicht ganz kompatible Feeds ausliefert. Diese kleinen Unterschiede werden hier also angeglichen, damit die weitere Verarbeitung leichter vonstatten gehen kann, und in der weiteren Verarbeitung keine Unterschiede mehr beachtet werden müssen.

Ein weiteres Problem der RSS-Feeds ist, dass sie nur Kurzfassungen der Artikel enthalten. In der Umwandlung müssen also noch die vollständigen Artikel von der Internetseite des Nachrichtenanbieters heruntergeladen und der Volltext aus dieser extrahiert werden. Hierzu wird zuerst das HTML-Dokument der Seite “geparsed”. Aus dem “geparsten” HTML-Dokument werden nun die relevanten Teile, zu denen z.B. nicht die Kopf- oder Fußzeile gehören, 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 (Siehe Abb. 5).

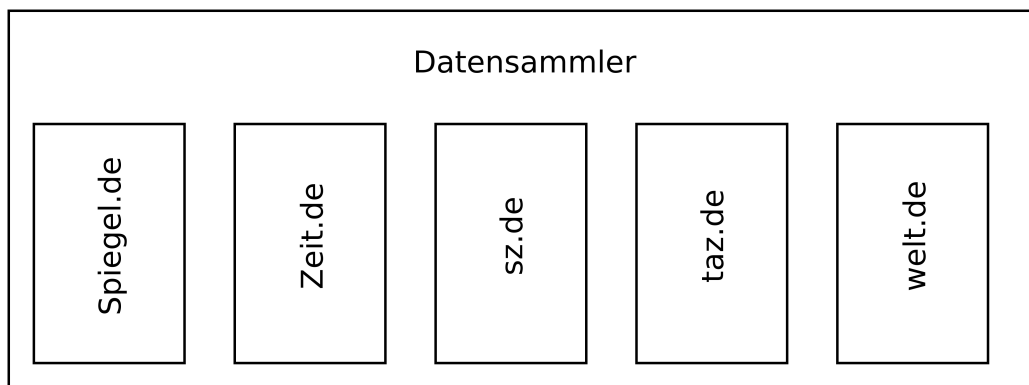


Abb. 5: Der Datensammler und seine Module

Am Ende der Datenaquirierung werden die gesammelten Daten an die nächste Stufe weitergegeben, in der die Daten gespeichert werden.

adapt
to ac-
tual
used
news-
papers

4.2 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, wie es normalerweise der Fall ist, sondern eine andere Abfragemöglichkeit bietet. Relativ früh habe ich mich dafür entschieden, dass ich meine Datenbankabfragen als “MapReduce”-Funktionen formulieren möchte.

4.3 MapReduce

Das MapReduce-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 (Vgl. Abb. 6):

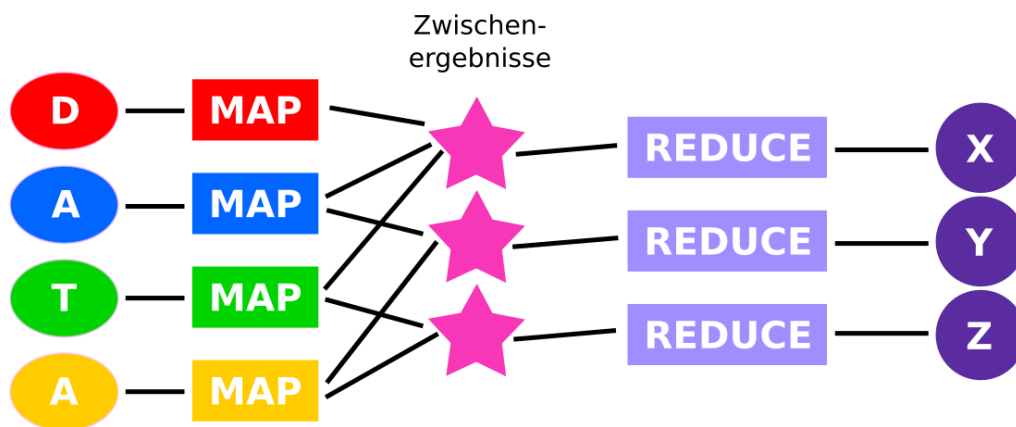


Abb. 6: Der MapReduce prozess als Grafik

1. Am Anfang des Prozesses steht eine Menge aus n Eingangsdatensätzen. In meinem Fall sind das die Zeitungsartikel, 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, dass wir zu jedem vorkommenden Wort die Häufigkeit bestimmen wollen, ordnet die Map-Funktion also jedem Artikel die Menge der darin enthaltenen Wörter zu. Da diese Funktion auf jeden Artikel angewandt wird, kann man sich in diesem Fall den gesamten Map-Prozess als eine Funktion vorstellen, die die Menge aller Artikel der Menge der darin enthaltenen Worte zuordnet. Der hierzu zugehörige Code dieser Map-Funktion ist:

```
1 function map() {
2   var text = this.text.replace(/[~A-Za-zÄäÖöÜüß ]/g, " "); //
   entferne alle überflüssigen Satzzeichen, wie .,?!-
3   var words = text.split(" "); // Teile den text in Wörter
4   for (var i = 0; i < words.length; i++) { // für jedes Wort
5     var word = words[i];
6     if(word) {
7       emit(word, 1); // füge der Ausgabemenge das Wort
   hinzu
8     }
9   }
10 }
```

2. Im nächsten Schritt werden Elemente mit gleichen Schlüsseln gruppiert. Nach diesem Schritt liegt also eine Menge aus Schlüssel-Wertmengenpaaren vor (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 [`@asdf`]). 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, was stark an eine Strichliste erinnert:

```
1 [
2   "der": {1; 1; 1; 1; 1; 1; 1; 1; 1},
3   "die": {1; 1; 1; 1},
4   ...
5 ]
```

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 nur die Elemente der Menge aufsummieren muss:

```
1 reduce(key, values) {  
2   return values.reduce((previousValue, currentValue) =>  
   currentValue + previousValue);  
3 }
```

4. Der vierte und letzte Schritt gehört eigentlich nicht mehr zum MapReduce-Verfahren. Dieser wird nach dem MapReduce-Verfahren ausgeführt und dient dazu, die Ergebnisse zu sortieren und eventuell Feinheiten zu verbessern. So ist es zum Beispiel möglich, selten genutzte Worte oder sogenannte Stoppworte in diesem Schritt auszusortieren. Stoppworte sind häufig auftretende Worte, die keine Relevanz für die Erfassung des Dokumenteninhaltes haben. Hierfür verwende ich verschiedene Stoppwortlisten, die von Sprachforschern erstellt worden sind, zusammen mit eigenen Ergänzungen. Der hierzu gehörige Code, der die Wortliste filtert, sieht wie folgt aus:

```
1 function filter(data) {  
2   return data.filter(word =>  
   stopwords.indexOf(word["_id"].toLowerCase()) < 0)  
3 }
```

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, dass die Wörter “Trump” und “Trumps” zusammengezählt werden. Hierbei wird immer das kürzeste Wort behalten, da dies meist die Grundform ist.

4.4 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 arbeiten, störte mich aber sehr stark an deren Komplexität. 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 Datenbank erfüllte die meisten meiner Anforderungen relativ gut, weshalb ich meinen gesamten, bis dahin existierenden Code an CouchDB anpasste. Nach einiger Zeit des Testens stellte sich allerdings heraus, dass CouchDB wahrscheinlich zu langsam sein würde und ein unpassendes Rechteverwaltungssystem mitbringt, was die Arbeitersparnis im Gegensatz zu MongoDB zunichte machen würde. In Folge dieser Erkenntnis entschied ich mich dazu, meine Software zurück auf MongoDB zu portieren und für diese eine HTTP-Schnittstelle zu entwickeln.

4.5 Middleware

Diese HTTP-Schnittstelle ist gewissermaßen das Bindeglied zwischen der Datenbank und der Visualisierung, weshalb es im nachfolgenden “Middleware” genannt wird. Sie nimmt HTTP-Anfragen von der Visualisierung (“Frontend”) entgegen, leitet diese an die Datenbank weiter und schickt die Antwort der Datenbank an das Frontend zurück. Insofern könnte man sie auch als Übersetzer zwischen verschiedenen Protokollen verstehen. Eine Übersetzung ist notwendig, da ich die Visualisierung im Webbrowser implementiert habe, und im Browser nur sehr wenige Schnittstellen verfügbar sind. HTTP ist eine dieser wenigen Möglichkeiten und bietet sich an, da es verhältnismäßig einfach zu implementieren ist.

4.6 Visualisierung

Ich habe mich entschieden, die Visualisierung im Browser umzusetzen. Dies lag hauptsächlich daran, dass ich bereits einige der verwendeten Technologien kannte, und somit die Einstiegshürde niedriger war. Außerdem gibt es für den Browser und die damit verbundenen Technologien bereits sehr gute Bibliotheken zur Datenvisualisierung, auf die ich für meine Arbeit zurückgreifen konnte. Ein weiterer, nicht unwichtiger Punkt, ist die einfache Portabilität von Programmen, welche im Browser laufen: Sie können auf fast jedem Computer, unabhängig von Betriebssystem oder Prozessorarchitektur, innerhalb von Sekunden aufgerufen und ausgeführt werden.

Eines der Hauptziele bei der Konzeption und Entwicklung der gesamten Analysesoftware war die Flexibilität: Es sollte möglichst schnell und einfach möglich sein, verschiedene Analysen über den Datensatz durchzuführen. Dies bezog sich nicht nur auf vorgefertigte, bei der Entwicklung bedachte Analysen, sondern auch zukünftige Ansätze. So sollte es einfach möglich sein, eigene neue Analysen über den Datensatz durchzuführen. Dies führte zu der konzeptionellen Entscheidung, dass die Analyse durch den Quellcode des Visualisierungsmoduls vorgegeben sein sollte und alle anderen Zwischenschritte, wie die Middleware oder die Datenbank nur auf die Anfragen des Analysemodul hören sollten. Dies hat zur Folge, dass ein großer Teil der Komplexität des Codes im Frontend liegt.

Jede Visualisierung fragt am Anfang die benötigten Daten bei der Middleware über MapReduce-Anfragen an. Die einzelnen Funktionen, die hierzu benötigt werden, werden also in dem Code der Visualisierung definiert und hängen davon ab, welche Daten aus dem Datensatz benötigt werden. Dieser Vorgang kann mehrere Male wiederholt werden, um verschiedene Informationen über verschiedene Aspekte des Datensatzes zu bekommen. Hiernach werden die nun vorliegenden Daten weiterverarbeitet. Sie können also gefiltert, sortiert oder miteinander verrechnet werden. Liegen die Daten nun in der gewünschten Form vor, kann die eigentliche Visualisierung beginnen.

Um die nun aufbereiteten Daten zu visualisieren, gibt es viele verschiedene Möglichkeiten. Ich habe mich dazu entschieden, die Softwarebibliothek “d3.js” zu verwenden, welche es ermöglicht, Verbindungen zwischen Datenmengen und Elementen im Browser herzustellen. Das Aussehen dieser Elemente kann nun mithilfe von gängigen Techniken, wie “CSS” angepasst werden. Für komplexere Visualisierungen werden auch andere Bibliotheken verwendet, um z.B. die einzelnen Positionen der Worte zu berechnen.

read;
konkreter;
bild

4.7 Zusammenführung

Da die gesamte Software modular aufgebaut ist, besteht sie aus vielen kleinen Komponenten, die spezialisiert und alleine nicht sehr nützlich sind. Diese kleinen Komponenten nennt man in der Softwareentwicklung “Microservices”. Um aus diesen kleinen Microservices eine funktionierende Software zu erhalten, muss man diese jeweils in einer passenden Umgebung starten und richtig miteinander verbinden. Um dies zu bewerkstelligen, habe ich das Software-

read

Paket “Docker” verwendet. Dieses erlaubt es zu allererst, die einzelnen Umgebungen für jeden Microservice zu beschreiben. Diese enthalten zum Beispiel die verschiedenen Programmbibliotheken, die ich in den einzelnen Teilen verwandt habe, oder einen Webserver, um den Visualisierungscode an den Webbrowser auszuliefern. Eine solche Umgebung mit allen darin enthaltenen Programmen nennt man “Container”. Diese Container erstellt man immer auf der Basis eines anderen “Basiscontainers”. Dies sorgt dafür, dass man sich nicht immer um alles innerhalb eines Containers kümmern muss, also zum Beispiel nicht immer selber das Betriebssystem installieren muss, sondern nur noch die eigenen Komponenten hinzufügen muss. Dies geschieht automatisiert über eine Datei, die jedem Microservice beigelegt ist, dem sogenannten “Dockerfile”. Auch ist es möglich komplett, vorgefertigte Container zu verwenden, was ich z.B. für die Datenbank nutze.

Die einzelnen Container, die nun die komplette Umgebung enthalten, die die jeweiligen Microservices zum laufen benötigen, müssen nun nur noch alle richtig verbunden werden. Dies geschieht mit einem Programm namens “Docker-compose”. Diesem gibt man eine Datei, die den Zustand des gesamten gewünschten Systems beschreibt. Automatisiert wird aus diesem und allen Dockerfiles der Microservices dann genau das beschriebene System erstellt und gestartet. Die Datei, die dieses System abbildet, sieht wie folgt aus:

```
1 version: '2'
2
3 services:
4   data-collectors:
5     build: data-collectors/
6     links:
7       - mongodb
8   backend:
9     build: backend/
10    links:
11      - mongodb
12
13   frontend:
14     build: frontend/
15     links:
16       - backend
```

```
17     ports:
18         - "80:80"
19
20     mongodb:
21         image: "mongo:latest"
```

Auch hier sieht man wieder sehr schön die Struktur und die Modularität der Software und ihre einzelnen Dienste.

Dieses System, das alle ständig laufenden Komponenten enthält, wird nun auf einem Server gestartet und ermöglicht es, sowohl den Code der Visualisierung als auch Auskünfte über den Datensatz zu bekommen.

5 Beobachtungen

5.1 Evaluation des Verfahrens

- Map/Reduce bei jedem Query skaliert nicht :(

5.2 Verschiedene Zeitungen im Vergleich

6 Fazit

7 Quellen