

Semantische Suche auf einem Web-Korpus

Philipp Bausch

April 25, 2014

- 1 Einleitung
- 2 Die Daten
- 3 Verarbeitung
- 4 Fazit und Daten

- Verwendet Wikipedia als Korpus der Textsuche

- Verwendet Wikipedia als Korpus der Textsuche
- Was findet man hier nicht?

Was Wikipedia ist



Link



- *triviale* Informationen

- *triviale* Informationen
 - Verlauf von Ereignissen
 - Details von Ereignissen

- *triviale* Informationen
 - Verlauf von Ereignissen
 - Details von Ereignissen
- Meinungen von Autoren
- Debatten

- 1 Einleitung
- 2 Die Daten**
- 3 Verarbeitung
- 4 Fazit und Daten

- ClueWeb12

- ClueWeb12
 - **nicht** Kostenlos
 - nur englische Texte
 - nur kleine Texte ($< 10\text{MB}$)
 - 2.820.500 URLs

- ClueWeb12
- Commoncrawl Web Corpus 2012

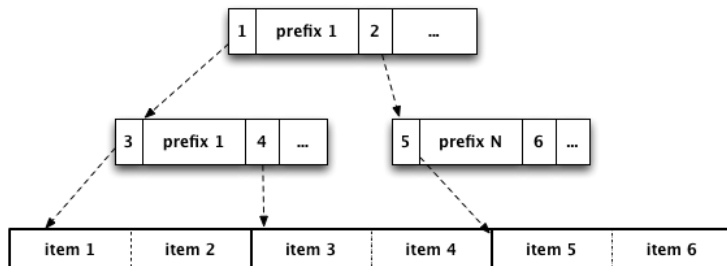
- ClueWeb12

- Commoncrawl Web Corpus 2012
 - Kostenlos
 - 3.8 Milliarden Web Dokumente
 - mehr als 100TB Datenmenge
 - 61 Millionen Domains

- Web Korpus von 2012 auf AmazonS3

- Web Korpus von 2012 auf AmazonS3
 - per EC2 Hadoop/HTTP abrufbar

- Web Korpus von 2012 auf AmazonS3
 - per EC2 Hadoop/HTTP abrufbar
 - URL-Index für ARC Dateien (Scott Robertson)



URL-Index Aufbau [IMG01]

- Web Korpus von 2012 auf AmazonS3
 - per EC2 Hadoop/HTTP abrufbar
 - URL-Index für ARC Dateien (Scott Robertson)
 - ARC Dateien enthalten vollen HTTP-Response

<http://www.spiegel.de/international/0,1518,druck-345720,00.html> 195.71.11.67 20120214055058 text/html 24686
HTTP/1.0 200 OK
Date:Tue, 14 Feb 2012 05:48:36 GMT
Server:Apache-Coyote/1.1
X-Powered-By:Servlet 2.4; JBoss-4.0.3SP1 (build: CVSTag=JBoss_4_0_3_SP1 date=200510231054)/Tomcat-5.5
Cache-Control:max-age=120
Expires:Tue, 14 Feb 2012 05:50:37 GMT
X-Host:Inxp-2885
X-Robots-Tag:noindex, nofollow, noarchive
Content-Type:text/html;charset=ISO-8859-1
Vary:Accept-Encoding
Content-Encoding:gzip
Content-Length:8341
X-Cache:MISS from Inxp-3958.srv.mediaways.net
X-Cache-Lookup:MISS from Inxp-3958.srv.mediaways.net:100
Via:1.1 www.spiegel.de, 1.0 Inxp-3958.srv.mediaways.net (squid/3.1.4)
Connection:close
x-commoncrawl-DetectedCharset:ISO-8859-1

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen

Eingabe:

www.cnn.com

www.news.google.com

...

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen
- 2 Alle zutreffenden **locations** merken

Ergebnis bis hier:

(ARCDatei1, Offset1, Größe1)
(ARCDatei2, Offset2, Größe2)

...

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen
- 2 Alle zutreffenden **locations** merken
- 3 Auf die **location** Zugreifen

Ergebnis bis hier:

```
http://www.spiegel.de/intern...  
HTTP/1.0 200 OK  
Date: Tue, 14 Feb 2012 05:48:...  
...
```

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen
 - 2 Alle zutreffenden **locations** merken
 - 3 Auf die **location** Zugreifen
- vorhanden? lokal laden
- ansonsten herunterladen und speichern

Ergebnis bis hier:

```
http://www.spiegel.de/intern...
HTTP/1.0 200 OK
Date: Tue, 14 Feb 2012 05:48:...
...
```

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen
- 2 Alle zutreffenden **locations** merken
- 3 Auf die **location** Zugreifen
vorhanden? lokal laden
ansonsten herunterladen und speichern
- 4 **Whitespaces** entfernen

Ergebnis bis hier:

Wie gerade nur ohne
Whitespaces

Schritte:

- 1 Für alle gegebenen URLs durch den Index gehen
- 2 Alle zutreffenden **locations** merken
- 3 Auf die **location** Zugreifen
vorhanden? lokal laden
ansonsten herunterladen und speichern
- 4 **Whitespaces** entfernen
- 5 TSV-File erstellen

Ergebnis bis hier:

Eine TSV Datei mit einem Dokument pro Zeile:
`URL<TAB>TITEL<TAB>INHALT<NL>`

- 1 Einleitung
- 2 Die Daten
- 3 Verarbeitung**
- 4 Fazit und Daten

- HTML → inhaltlichen Texten
- Nutzt Features wie durchschnittliche Satzlänge

- HTML → inhaltlichen Texten
- Nutzt Features wie durchschnittliche Satzlänge
- Entfernt fast die gesamte Navigation
- Schnell [KFN10]

Beim Verwenden der Commoncrawl Daten:

- Probleme mit dem vorhandenen Entitätenerkennung:

Beim Verwenden der Commoncrawl Daten:

- Probleme mit dem vorhandenen Entitätenerkennung:
 - Titel Entitäten nur bei genauer Übereinstimmung
 - arbeitet mit [Sektionen](#) und [Links](#) (nur auf wikipedia)

Beim Verwenden der Commoncrawl Daten:

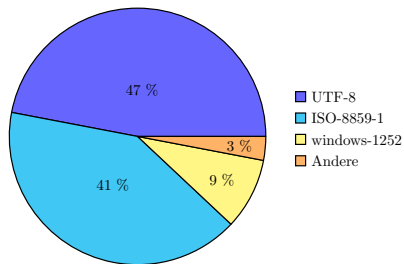
- Probleme mit dem vorhandenen Entitätenerkennung:
 - Titel Entitäten nur bei genauer Übereinstimmung
 - arbeitet mit [Sektionen](#) und [Links](#) (nur auf wikipedia)
- Probleme mit Kontexten:

Beim Verwenden der Commoncrawl Daten:

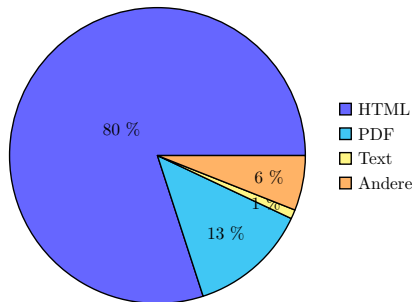
- Probleme mit dem vorhandenen Entitätenerkennung:
 - Titel Entitäten nur bei genauer Übereinstimmung
 - arbeitet mit [Sektionen](#) und [Links](#) (nur auf wikipedia)
- Probleme mit Kontexten:
 - manche Seiten erzeugen zu große Kontexte (keine Sätze)

- *verallgemeinerter* Entitätenerkennung (mit und ohne POS-Tags)
- maximale Kontextgröße im Decomposer
- Parser zum Einlesen der TSV-Datei
- Boilerpipe zum Extrahieren des Textes

- 1 Einleitung
- 2 Die Daten
- 3 Verarbeitung
- 4 Fazit und Daten**



Verteilung der Codierungen



Verteilung der Typen

Siehe Browser

- KFN10** Kohlschütter, Christian ; Fankhauser, Peter ; Nejd, Wolfgang:
Boilerplate detection using shallow text features. In:
Proceedings of the third ACM international conference on
Web search and data mining ACM, 2010, S. 441–450
- IMG01** https://raw.githubusercontent.com/trivio/common_crawl_index/master/docs/tree.png

- Seite die als HTML erkannt wurde.
- Komische Seite (Vorsicht!!)
- NULL-Byte im Titel
- RSA-Key
- RNA