

Web Mining: Übung 2

Lösungsvorschlag

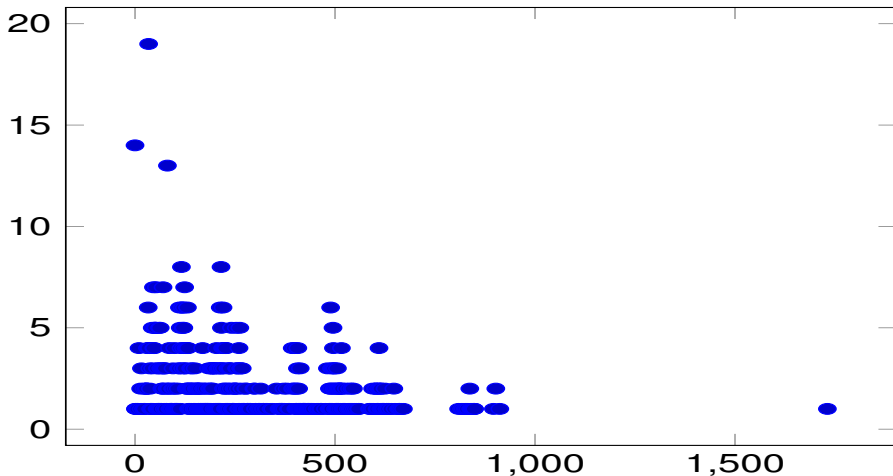


TECHNISCHE
UNIVERSITÄT
DARMSTADT



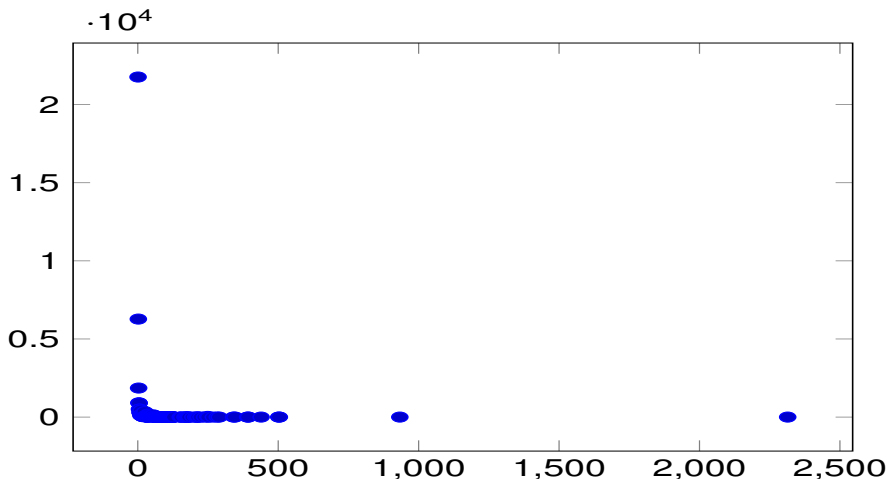
Aufgabe 2: Entwicklung eines Crawlers

Histogramm: Es gibt y Seiten mit x Links



Aufgabe 2: Entwicklung eines Crawlers

Histogramm: Es gibt x Links, die y mal auftraten



Aufgabe 2: Entwicklung eines Crawlers

Erkennung wiederkehrender Links



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Entfernung aller Listen:

```
1 <ul>
2     <li><a href="start.html">Start</a></li>
3     <li><a href="news.html">News</a></li>
4     <li><a href="impressum.html">Impressum</a></li>
5 </ul>
```

- Entfernung mit CSS-Selektor: [class*='nav']:

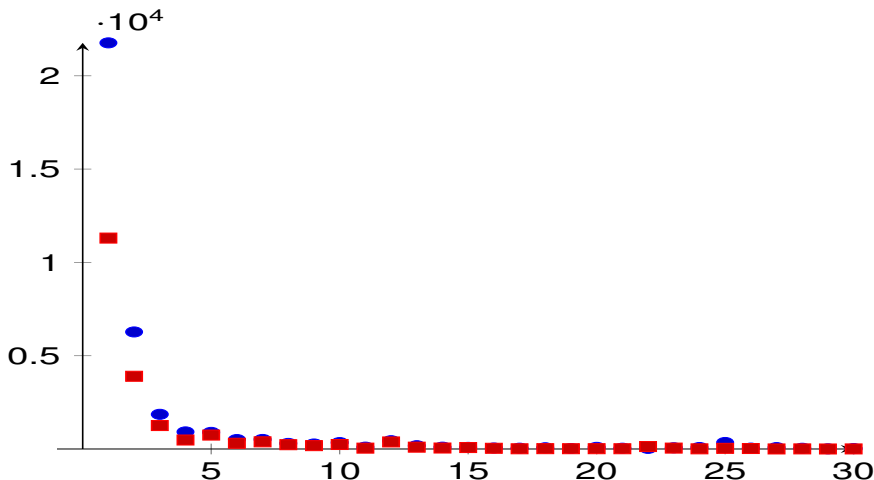
```
1 <div class="main-navigation">
2     ...
3 </div>
```

- Speichern der HTML-Dateien als HASH(Kanonisierte URL).html

Aufgabe 2: Entwicklung eines Crawlers

Wirksamkeit der Duplikaterkennung

Es gibt x Links, die y mal auftraten



Aufgabe 2: Entwicklung eines Crawlers

Wie oft wurden Hosts besucht?

(Subdomains subsumiert, Seiten: 1.107)





tbd



- ▶ Suchstrategie funktioniert; die besuchten Hosts streuen
- ▶ Darstellung der Verlinkungen zwischen den Hosts wäre spannend („Wie vernetzt sind Hochschulen?“)
- ▶ Extrem häufig verlinkte Social Media Seiten erschweren solche Analysen. Bsp.: Anteil Twitter in offenen Links zu einem fortgeschrittenen Zeitpunkt:

$$\frac{8736}{24877} \approx 35\%$$

- ▶ Technische Lösung skaliert erwartungsgemäß nicht. Datenverwaltung nimmt mehr Zeit in Anspruch als das eigentliche Crawlern
 - ▶ Zu besuchende Links und Statistiken in Textdateien
 - ▶ Speichern eines jeden Datums direkt nach Erhebung (Persistenz)