

Universität Potsdam
Institut für Informatik
Lehrstuhl Maschinelles Lernen



Suche im Web

Tobias Scheffer

WWW

- 1990 am CERN von Tim Berners Lee zum besseren Zugriff auf Papers entwickelt.
 - ◆ HTTP, URLs, HTML, Webserver.
- Verbindet FTP mit der Idee von Hypertext.
- Multilingual (ca. 75% englisch, je 5% japanisch und deutsch).
- Groß, verteilt.
- Volatil: Dokumente erscheinen und verschwinden.
- Unstrukturiert, heterogen.

Entwicklung des WWW

- Hypertext: 1960er Jahre.
- ARPANET / DARPA NET / Internet: 1970er.
- FTP: 80er Jahre.
 - ◆ Archie: Crawler sammelt Dateinamen von FTP-Servern und ermöglicht Regex-Suche.
 - ◆ Gopher.
- Web-Browser: Anfang der 90er von Doktoranden.
 - ◆ Mosaic: Marc Andreessen & Eric Bina, UIUC NCSA.
 - ◆ Netscape und IE aus Mosaic hervorgegangen.
- Web-Directories: Verzeichnis „Yahoo“, 1994, zwei Stanford-Doktoranden.

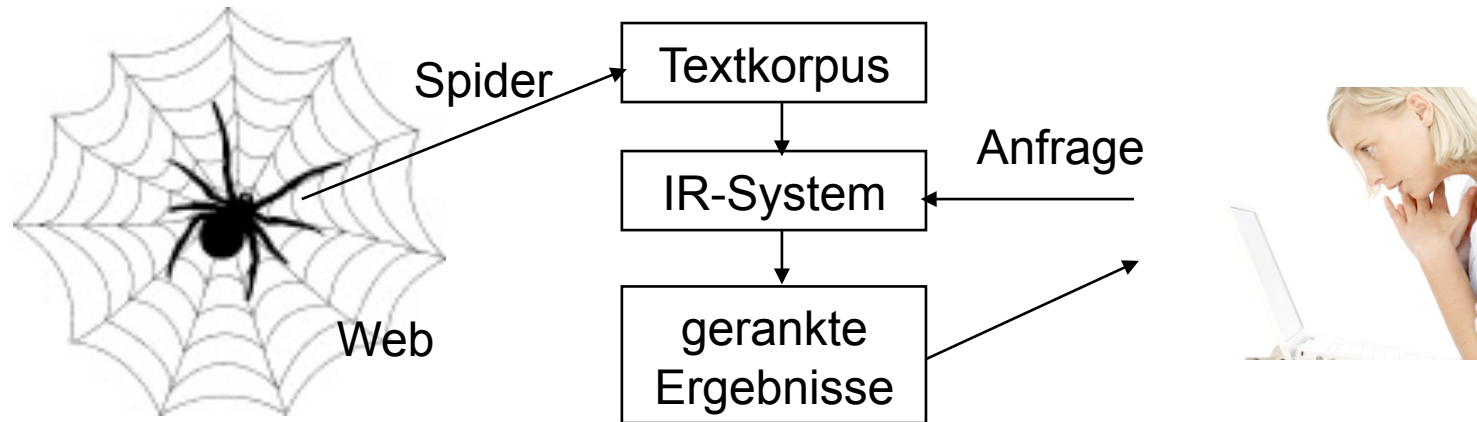
Suche im WWW

- WebCrawler: Programmier-Projekt, Uni Washington, 1994. → Excite, AOL.
- 1994: Doktorand an CMU entwickelt Lycos.
- 1995: DEC entwickelt Altavista.
- 1998: Doktoranden aus Stanford entwickeln Google. Idee: Analyse der Linkstruktur, relevantere Ergebnisse.

Graph-Struktur des WWW

- Große zentrale, stark verknüpfte Komponente.
 - ◆ Jede Seite von jeder Seite erreichbar.
- Randgebiete die nur
 - ◆ Auf die zentrale Komponente verweisen.
 - ◆ Von zentraler Komponente aus verwiesen werden.
 - ◆ Isolierte Komponenten.
- Zipfs Gesetz (spezielles „power law“):
 - ◆ Wenn n natürlich auftretende Zahlen z_i der Größe nach geordnet werden, dann ist $z_i \sim 1/i$.
 - ◆ Beschreibt im Web Länge der Seiten, Anzahl eingehender und ausgehender Kanten, Anzahl Hits etc.

Suche im Web



- Unterschied zu anderen IR-Systemen:
 - ◆ Spider / crawler sammelt Texte für Indexstruktur.
 - ◆ Analyse der Linkstruktur liefert Relevanzinformation.
 - ◆ Suchmaschinen sind stark verteilte Systeme.

Crawling

- Beginne mit URLs aus zentraler Komponente,
- Folge Links rekursiv und füge neue Seiten in Index.
- Suchstrategien
 - ◆ Breadth first, depth first,
 - ◆ Hoher PageRank first.
 - ◆ Crawlhäufigkeit in Abhängigkeit von Änderungshäufigkeit und PageRank (WallstreetJournal.com vs. HottesSockenShop.de)
 - ◆ Crawlhäufigkeit in Abhängigkeit der Zugriffshäufigkeit auf die Seite.
 - ◆ Crawling zu den richtigen Zeitpunkten bei zyklischen Updates.

Crawling

- Vor Ablegen testen, ob Seite schon abgelegt ist.
 - ◆ Web ist kein Baum sondern Graph.
 - ◆ Trie, Hash-Tabelle, schneller Zugriff auf URL.
 - ◆ Mehrere (unendlich viele) URLs für dieselben Seiten durch Zyklen in Verzeichnisstrukturen.
 - ◆ Seiten mit dynamischem Inhalt.
 - ◆ Seiteninhalt von Cookie abhängig.

Crawling

- Gecrawlte Seiten müssen parsiert werden.
 - ◆ Finden der Links, die dann wieder rekursiv verfolgt werden.
- Ankertexte werden als Indexterme für die verlinkte Seite abgelegt.
- Beschränkung des crawlbaren Bereichs durch Benutzer möglich.
 - ◆ Robots.txt, Robots META Tag



Crawling

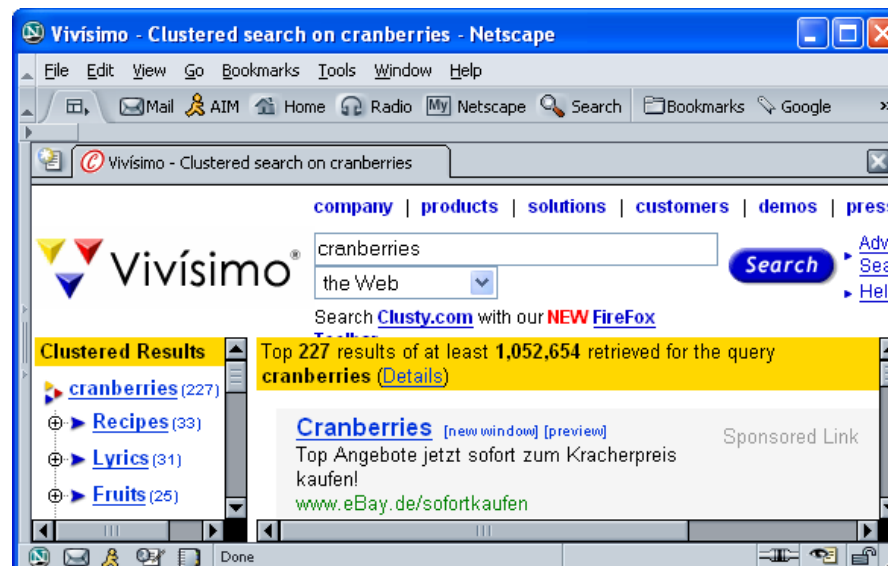
- Crawlen muss in parallelen threads erfolgen.
 - ◆ Antwortzeit auf jeden HTTP-Request relativ hoch,
 - ◆ Ein Crawler-Prozess muss einige hundert Request-Threads verwalten.
- Aber: zuviele Anfragen an denselben Server vermeiden
 - ◆ „Denial of Service Attack“.

Benutzerschnittstellen

- HTML-Webinterface
- <form>-Umgebung ermöglicht Übergabe des Suchstrings als Parameter an Suchmaschine.
- Suchmaschine generiert Rückgabeseite,
 - ◆ CGI
 - ◆ Java Servlet.
- Session Tracking:
 - ◆ Speicherung von Präferenzen, Anfrage-Verfeinerung.
 - ◆ HTTP ist zustandsloses Protokoll, Tracking nur über Cookies oder dynamische URLs mit Session-ID.
- Meist nur kurze Anfragen (2.35 Wörter).

Benutzerschnittstellen

- Advanced Search:
 - ◆ Wird fast nie benutzt.
 - ◆ Ähnliche Seiten finden,
 - ◆ Links auf Seite finden,
 - ◆ Maschinelle Übersetzung, cross-language retrieval.
- Clusterung von Seiten



Retrieval

- Suchterme \rightarrow Index \rightarrow Liste von URLs.
- Texte mit mehreren Suchtermen durch Mengenoperationen.
- Aufbau eines Ähnlichkeitsvektors mit verschiedenen Merkmalen

$$\begin{matrix}
 \text{sim}(x, y) \text{ im Vektorraummodell} \\
 \tau_{\#} \text{ gleicher Terme in k1-Tags} \\
 \tau_{\#} \\
 \tau_{\#} \dots \\
 \tau_{\#} \\
 \tau_{\#} \text{ PageRank}(x) \\
 \tau_{\#} \\
 \tau_{\#}
 \end{matrix}$$

- Rückgabe der Seiten $\text{argmax}_x(w \cdot x)$
- Manuelle Konstruktion von w .
- Lernen von w aus Klickdaten.

Ranking

- Ende der 90er wurde alle Suchmaschinen in Web-Portale mit breitem Informations- und Unterhaltungsangebot umgebaut.
- Gegenannahme von Google: nur gute Suchfunktion ist wichtig, der Rest interessiert niemanden.
- Idee von Kleinberg und Page & Brin: Verweisstruktur zeigt, welche Webseiten irgend jemanden interessieren.
- Suchergebnisse werden nach PageRank sortiert.

PageRank

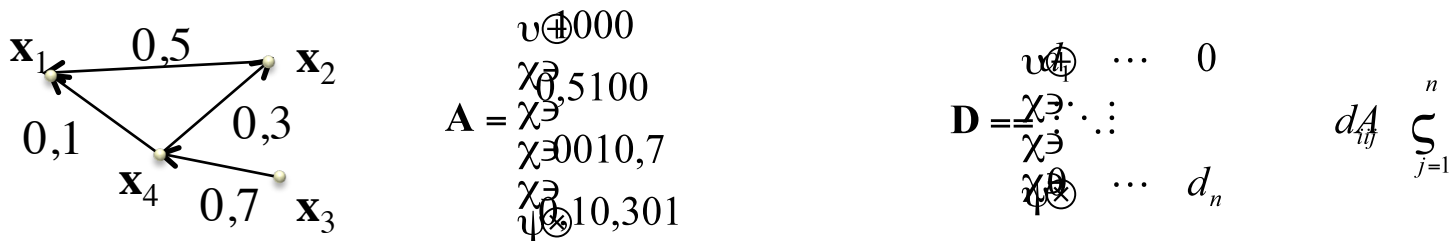
- Keine Unterscheidung in Hubs und Authorities.
- Ranking der Seiten nur nach Authority.
- PageRank ist umso höher, je mehr andere Seiten auf die Seiten verweisen, und je höher deren PageRank ist.
- Wird auf das ganze Web angewandt, nicht nur auf Ergebnismenge einer Suchmaschine.

Authority-Ranking: PageRank

- Gegeben: Trainingsdaten $\{x_1, x_2, \dots, x_n\}$ mit gegebenen lokalen Authority Scores (Kompetenz-Bewertungen):
 - ◆ Authority Score A_{ij} gibt an wie „kompetent“ x_j aus Sicht von x_i ist.
 - ◆ Link oder kein Link, Link-Position. Authority Matrix
- Gesucht: Modell $f: A \in R^{n \times n} \mapsto s \in R^n$ welches für Instanzen $\{x_1, x_2, \dots, x_n\}$ globale Authority Scores (Ranking) s_i liefert.
 - Annahme: Je kompetenter x_i und je höher der Authority Score A_{ij} , desto kompetenter x_j .

Authority-Ranking: PageRank

- Modellierung der (nicht-symmetrischen) Kompetenz-Bewertungen als ungerichteter Graph:
 - ◆ Instanzen sind Knoten, Authority Scores sind Kanten-Gewichte
 - ⇒ Authority Matrix = Adjazenzmatrix A .



- ◆ Intuition: (normierter) Graph beschreibt mit welcher Wahrscheinlichkeit Knoten x_i Knoten x_j als „Experten“ nennen würde.
- ◆ Beispiel: x_4 hält x_2 für 3-mal so kompetent wie x_1 , x_3 ist aus seiner Sicht kein Experte.

Authority-Ranking: PageRank

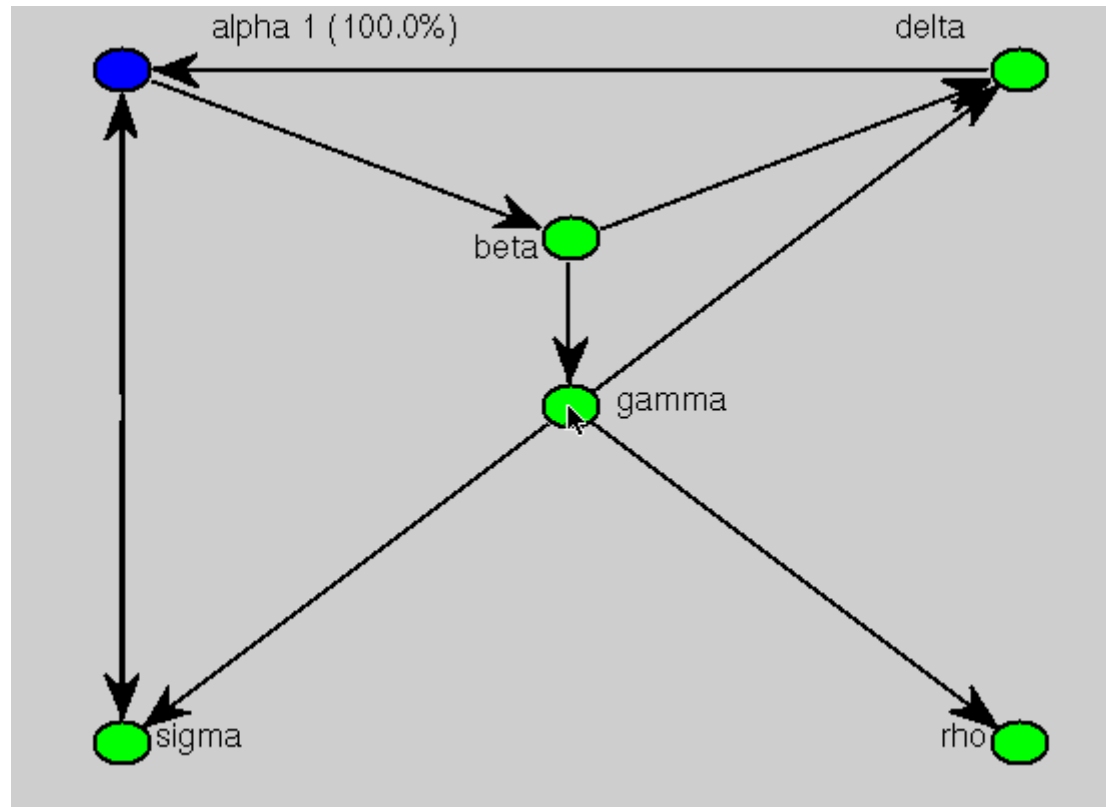
- Random Surfer:
 - ◆ Beginnend bei einem beliebigen Knoten folgt der Surfer mit Wahrscheinlichkeit $1-\varepsilon$ einem zufälligen Link.
 - ◆ Mit Wahrscheinlichkeit ε startet er neu an Zufallsknoten.
 - ◆ Wahrscheinlichkeit für Aufenthalt an einem Knoten = globaler Authority Score.
- Hubs & Authorities (HITS):
 - ◆ Jeder Knoten besitzt...
 - ★ *Hub Score*: Wie gut sind seine Verweise auf „kompetente“ Knoten
= ausgehende Kanten.
 - ★ *Authority Score*: Wie kompetent ist der Knoten, d.h. wie viele kompetente Knoten verweisen auf ihn = eingehende Kanten.

Authority-Ranking: PageRank

- PageRank: Random Surfer Modell zum Ranking von Webseiten.
 - ◆ (Ursprünglicher) Ranking-Algorithmus von Google.
 - ◆ Abhängig von Query werden relevante Webseiten gefunden und nach ihrem globalen Authority Score sortiert.
- Annahmen:
 - ◆ Link auf Webseite x_i verweist auf „kompetente“ Webseite x_j , d.h. Adjazenzmatrix des Webgraphen = Authority Matrix mit $A_{ij} = 1$.
 - ◆ Mit Wahrscheinlichkeit $1 - \varepsilon$ folgt der Nutzer (Random Surfer) einem Link auf der Webseite.
 - ◆ Mit Wahrscheinlichkeit ε wechselt er auf eine zufällige Webseite.

Authority-Ranking: PageRank

- Beispiel:



Authority-Ranking: PageRank

- Gegeben: Wahrscheinlichkeit dafür, dass der Nutzer von Webseite x_i zu Webseite x_j wechselt ist

$$P_{ij} = \frac{A_{ij}}{\sum_{k=1}^n A_{ik}} \cdot \frac{1}{n}$$

Transitionswahrscheinlichkeit

und somit $P_{ij} = \frac{A_{ij}}{\sum_{k=1}^n A_{ik}} \cdot \frac{1}{n}$ mit $U_{ij} = \frac{1}{n}$.

- Gesucht: Wahrscheinlichkeit dafür, dass man auf Webseite x_i ist, d.h. $s_i = \dots$.

Aufenthaltswahrscheinlichkeit

Authority-Ranking: PageRank

- Algorithmus: Beginnend mit initialem Ranking Scores s iterativ neue Scores bestimmen mit

$$s \mathbf{P} s \frac{1}{c} \quad \text{wobei} \quad c = \|\mathbf{P}^T s\| \quad \|\mathbf{1}\| = 1$$

- Konvergenz von PageRank bei $s \mathbf{P} s = s$, sodass gilt

$$s \mathbf{P} s \frac{1}{\lambda} \quad \text{d.h. } s \text{ ist ein Eigenvektor von } \mathbf{P}^T \text{ mit Eigenwert } \lambda.$$

- Man kann zeigen, dass s der Eigenvektor mit größtem Eigenwert λ ist.

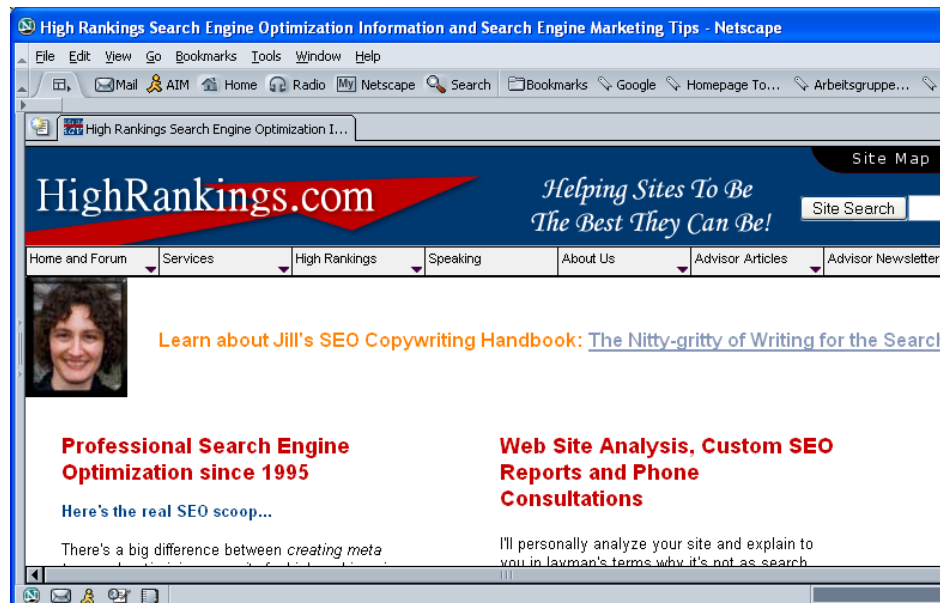
Authority-Ranking: PageRank

- Vorteile:
 - ◆ Leicht und effizient berechenbar.
 - ◆ Existenz & Eindeutigkeit der Lösung sowie Konvergenz des Algorithmus ist garantiert für $0 < \alpha < 1$.
- Nachteile:
 - ◆ Links können schlechte Indikatoren für Kompetenz-Bewertung sein:
 - ★ Kompetenz zweier Instanzen kann sich gegenseitig verstärken.
 - ★ Automatisch generierte Links haben kaum Aussagekraft.
 - ★ „Künstliche“ Links (z.B. Link-Spam) verändern Ranking.
 - ◆ Eigenschaften der Instanzen fließen nicht ins Ranking ein.

PageRank: Link Spam

- PageRank ist beeinflussbar.
- Link Farmen:
 - ◆ Felder künstlich generierter Seiten, deren PageRank $E(u)$ „geerntet“ und zu Zielseite geleitet wird.
- Guestbook Spam:
 - ◆ Generierte Einträge in Gästebüchern und Blogs, die Verweis auf Zielseite enthalten.
- Link Exchange Services:
 - ◆ Seiten mit Links auf (thematisch nicht verwandte) Seiten, meist gegen Geld.
- Partner-Programme: z.B. Amazon, Ebay.
 - ◆ Link auf Produktseite gegen Umsatzbeteiligung.

Link Spam



Link Spam: BadRank

- Wenn PageRank beeinflussbar wird, dann verliert er seine Korrelation zur Relevanz der Seiten.
- Link Spam sollte bei der Berechnung des Page Ranks nicht so berücksichtigt werden wie „natürliche“ Links.
- Suchmaschinenbetreiber haben „Blacklists“.
 - ◆ URLs von Link-Spam-Seiten.
 - ◆ Werden manuell erstellt.

BadRank

- Invertierter PageRank-Algorithmus, „bestraft“ Seiten, die auf Link Spam verweisen.
- Initialisierung: $B(u) = 1$, wenn u auf Blacklist.
- Für alle Seiten:

$$Bl(u) = \sum_{v \in \mathcal{R}(u)} \frac{B(v)}{N_u}$$

$$B(u) = \frac{Bl(u)}{|\mathcal{B}_u|}$$

- BadRank wirkt wie negativer PageRank.

Browsing

- Web-Verzeichnisse
 - ◆ Dmoz.org
 - ◆ Yahoo, Google.
- „Taxonomien des menschlichen Wissens“.

Browsing: Webverzeichnisse

- Dmoz.org,
- Yahoo,
- Google.

