

Systeme II

9. Die Struktur des Webs

Christian Schindelhauer
Technische Fakultät
Rechnernetze und Telematik
Albert-Ludwigs-Universität Freiburg
Version 15.07.2013



Abschlussveranstaltung am Tunisee

Termin

3

- Dienstag, 30.07.2012, 16:00 Uhr am Tunisee
- oder 15:30 mit Fahrrad am Gebäude 051

Plan

 Moderner Dreikampf: Baden, Grillen, Trinken

Versorgung

- Getränke werden gestellt (vorher im Forum bestellen)
- Grill und Grillkohle wird bereit gestellt
- Essen durch Selbstorganisation





Der Webgraph

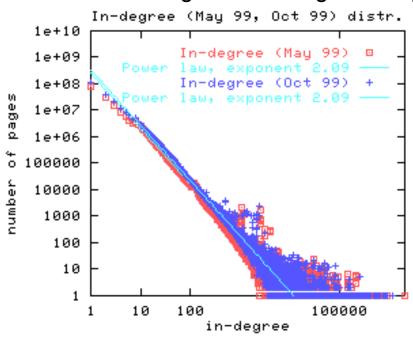
- G_{WWW}:
 - Statische HTML-Seiten sind Knoten
 - Link bezeichnen gerichtete Kanten
- Ausgrad eines Knoten:
 - Anzahl der Links einer Web-Seite
- Eingrad eines Knoten
 - Anzahl der Links, die auf eine Web-Seite zeigen
- Gerichteter Pfad von Knoten u nach v
 - Folge von Web-Seiten, indem man den Links folgt
- Ungerichteter Pfad (u=w₀,w₂,...,w_{m-1},v=w_m) von Seite u nach v
 - Für alle i gibt es entweder einen Link von winach with oder umgekehrt
- Starke (schwache) Zusammenhangskomponente
 - Maximale Knotenmenge in der zwischen allen Knoten dieser Menge ein (un) gerichteter Pfad besteht

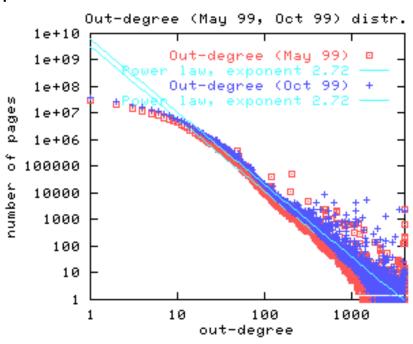




Ein- und Ausgradverteilung

- Der Ein- und Ausgrad gehorchen einem Potenzgesetz (power law)
 - d.h. die Häufigkeit von Eingrad i ist proportional zu ~ 1/iα





Ergebnisse von

- Kumar et al 97: 40 Millionen Web-Seiten
- Barabasi et al 99: Domain *.nd.edu + Web-Seiten in Abstand 3
- Broder et al 00: 204 million Web-Seiten (Scan Mai und Oktober '99)





Pareto Verteilung = Verteilung nach Potenzgesetz

Diskrete Pareto (power law) Verteilung für x ∈ {1,2,3,...}

$$\mathbf{P}[X=x] = \frac{1}{\zeta(\alpha) \cdot x^{\alpha}}$$

mit konstanten Faktor

$$\zeta(\alpha) = \sum_{i=1}^{\infty} \frac{1}{i^{\alpha}}$$

auch bekannt als Riemannsche Zeta-Funktion

- "Heavy tail"-Eigenschaft
 - nicht alle Momente E[Xk] sind definiert
 - Der Erwartungswert existiert genau dann wenn α>2
 - Varianz und E[X²] existieren genau dann wenn α>3
 - E[X^k] ist genau dann definiert wenn α>k+1
- Dichtefunktion der kontinuierlichen Pareto-Verteilung für x>x₀

$$f(x) = \frac{\alpha - 1}{x_0} \left(\frac{x_0}{x}\right)^{\alpha}$$



Spezialfall: Zipf-Verteilung

- George Kinsley Zipf behauptete, dass die Frequenz des n-häufigsten Word mit Häufigkeit f(n) erscheint,
 - wobei f(n) n = c
- Zipf-Wahrscheinlichkeitsverteilung für x ∈ {1,2,3,...}

$$\mathbf{P}[X=x] = \frac{c}{x}$$

mit konstantem Faktor c ist nur definiert auf einem endlichen Abschnitt, da

$$\ln n \le \sum_{i=1}^{n} \frac{1}{i} \le 1 + \ln n$$

für große n gegen unendlich wächst

- Zipf Verteilungen beziehen sich auf Ränge
 - Der Zipf-Exponent α kann wie die Pareto-Verteilung größer sein als 1, d.h. $f(n) = c/n^{\alpha}$
- Pareto-Verteilungen beziehen sich auf die absolute Größe
 - e.g. Anzahl der Einwohner, Anzahl der Links, etc.



Pareto-Verteilung

- Beispiele für Power-Laws (= Pareto Verteilungen)
 - Pareto 1897: Einkommensverteilung in der Bevölkerung
 - Yule 1944: Word-Häufigkeit in Sprachen
 - Zipf 1949: Größe von Städten
 - Länge von Molekülketten
 - Dateilängen von Unix-Dateien
 - Initialmasse von Sternen
 - Zugriffshäufigkeit auf Web-Seiten
 - Länge von Telefonanrufen

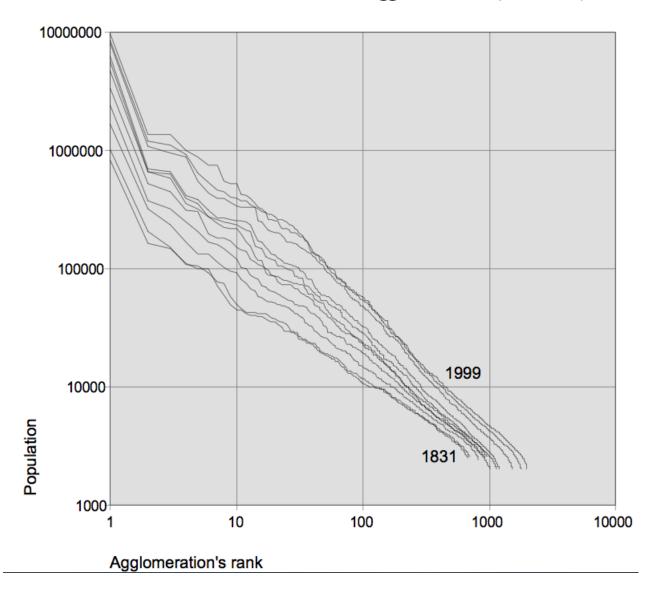
- ...



Größe von Städten (I)

Scaling Laws and Urban Distributions, Denise Pumain, 2003

Figure 1 The hierarchical differentiation in urban systems: Rank-size distribution of French agglomerations (1831-1999)

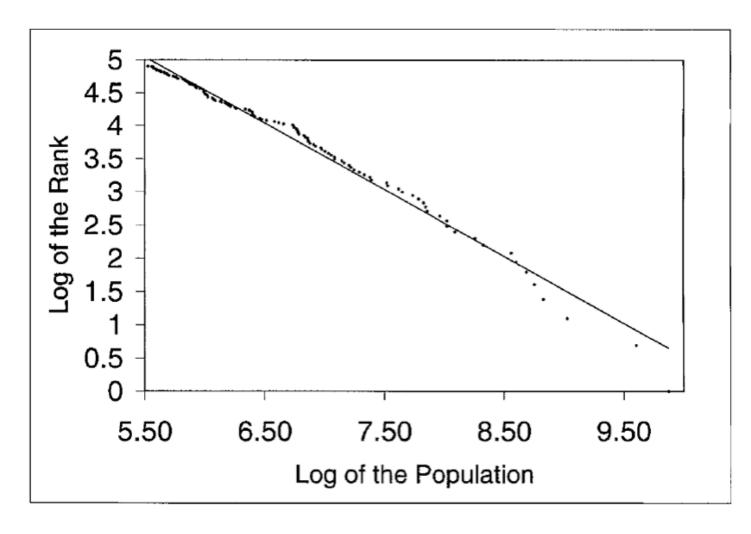


Zipf Verteilung





Größe von Städten (II)



Zipf Verteilung

FIGURE I
Log Size versus Log Rank of the 135 largest U. S. Metropolitan Areas in 1991
Source: Statistical Abstract of the United States [1993].





Zipf's Law and the Internet

Lada A. Adamic, Bernardo A. Huberman, 2002

a)

b)

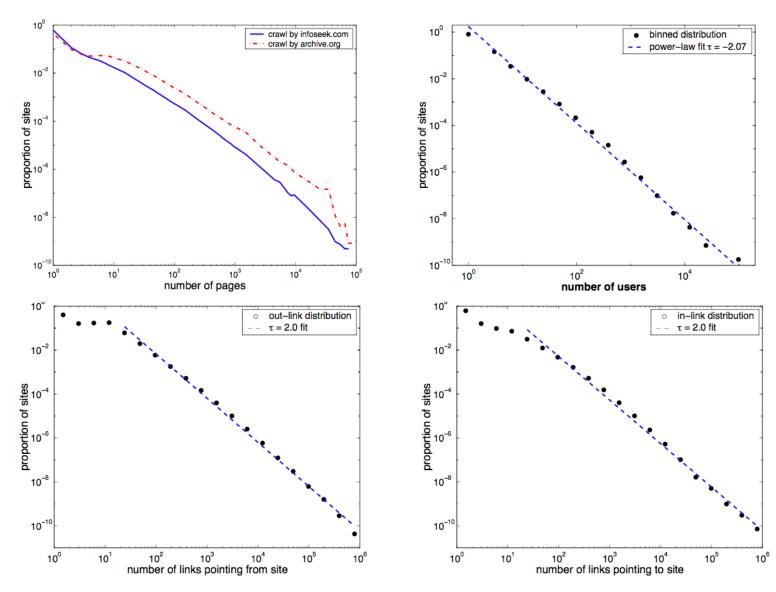


Figure 1. Fitted power law distributions of the number of site a) pages, b) visitors, c) out links, and d) in links, measured in 1997.

Pareto Verteilung





Zipf's Law and the Internet

Lada A. Adamic, Bernardo A. Huberman, 2002

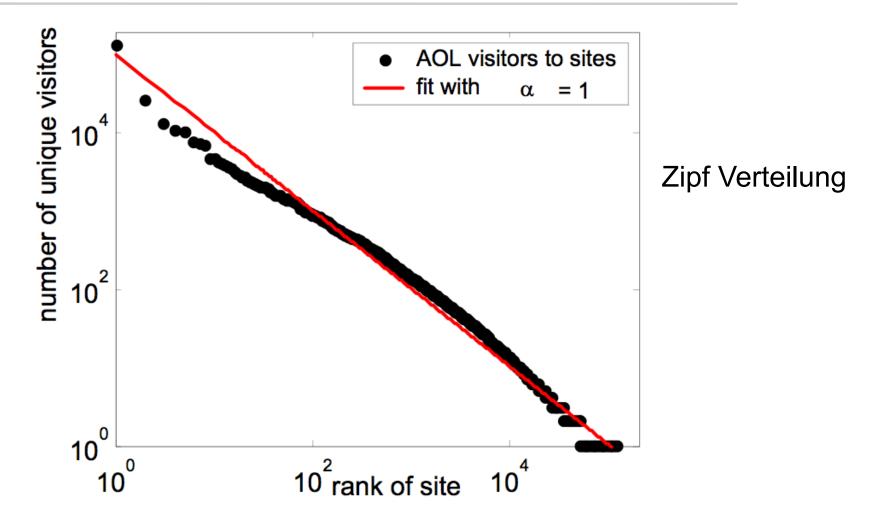
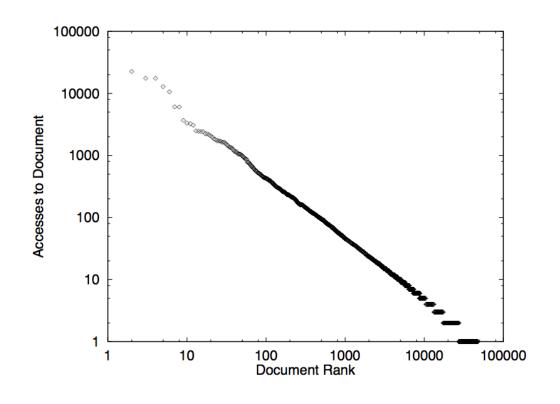
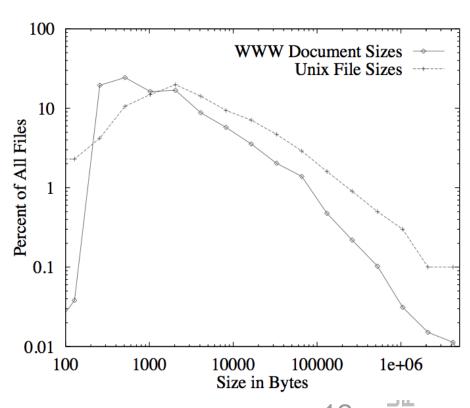


Figure 2. Sites ranked by the number of unique AOL visitors they received Dec. 1, 1997. AOL (America Online) is the largest Internet service provider in the United States. The fit is a Zipf distribution $n_r \sim r^{-1}$

Heavy-Tailed Probability Distributions in the World Wide Web Mark Crovella, Murad, Taqqu, Azer Bestavros, 1996

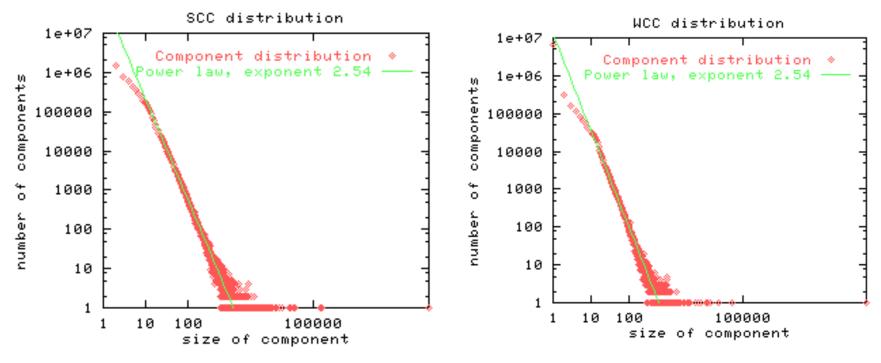






Größe der Zusammenhangskomponenten

Starke und schwache Zusammenhangskomponenten unterliegen einer Pareto-Verteilung

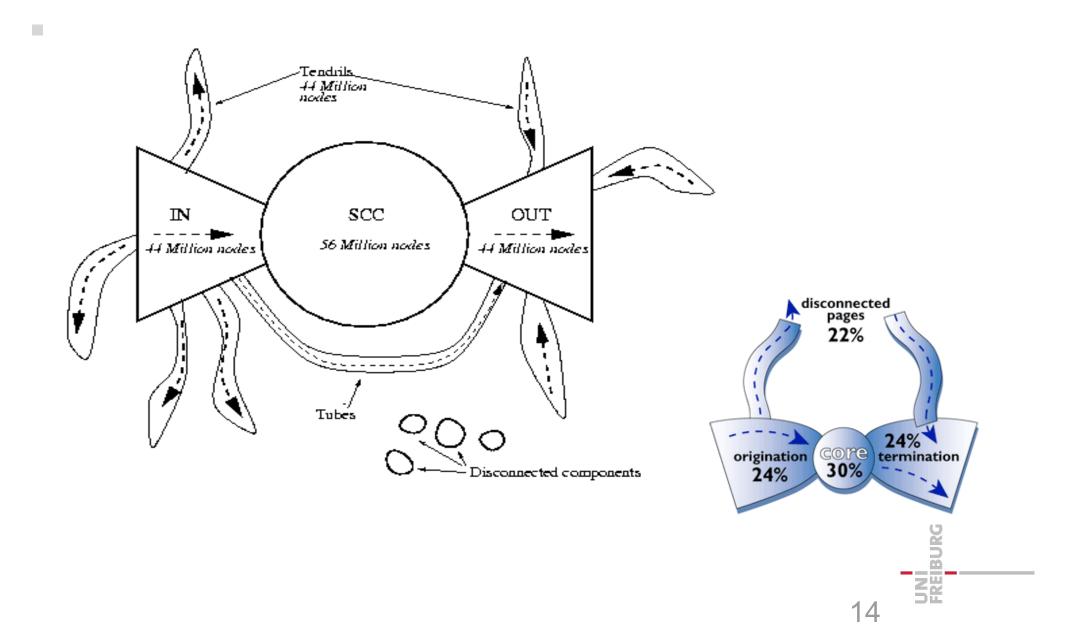


- A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener. "Graph Structure in the Web: Experiments and Models." In Proc. of the 9th World Wide Web Conference, pp. 309—320. Amsterdam: Elsevier Science, 2000.
 - Größte schwache Zusammenhangskomponente hat 91% aller Web-Seiten
 - Größte starke Zusammenhangskomponente hat 28% aller Webseiten
 - Durchmesser ist ≥ 28





Der Web-Graph (1999)





Power-Laws im Internet!?!

- Erste Arbeit zu Power-Laws im Internet sagte JA
 - C. Faloutsos, P. Faloutsos, and M. Faloutsos, "On Power-Law Relationships of the Internet Topology," in Proceedings of the ACM SIGCOMM, Sept. 1999.
- Diese Ergebnisse wurden 2002 angezweifelt:
 - Chen, Chang, Govindan, Jamin, Shenker, Willinger, The Origin of Power Laws in Internet Topologies Revisited, Infocom 2002
 - Gründe: Datenmenge zu gering, veraltet und spezielle Teilmengen ausgewählt
- Mittlerweile
 - Die Struktur des Internets gehorcht einem Power-Law, wenn man die Struktur berücksichtigt



Zipf's Law and the Internet

Lada A. Adamic, Bernardo A. Huberman, 2002

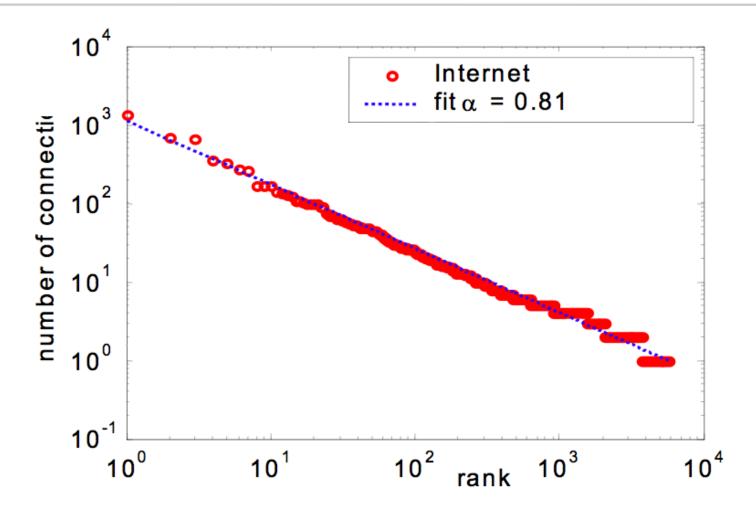


Figure 3. The connectivity of the internet backbone at the autonomous system (AS level). Each AS is itself a network corresponding to a single ISP, business entity or educational institution.



Die Struktur des Internets

Aus

 "Comparing the structure of power-law graphs and the Internet AS graph", Sharad Jaiswal, Arnold L. Rosenberg, Don Towsley, INCP 2004

Stub AS

- verbinden zwei Hosts; Sackgassen
- typischerweise Universitäten oder große Unternehmen
- vertrauen auf Transit ASes f
 ür Anbindung ans Internet.

Transit ASes

- Service providers
- typisch regionaler oder nationaler Internet
 Service Provider oder Backbone Network

Tier (oder Ebene/Level)

- bezeichnet Verbindungshierarchie

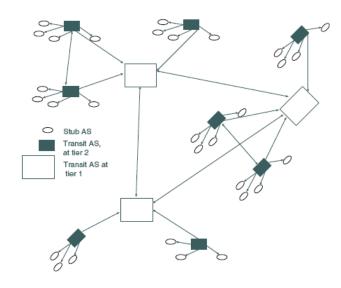


Fig. 1. Structure of the AS graph



Identifikation der Ebenen

- Die Identifikation der Ebenen ist nicht trivial und mehrdeutig
- Alternative Methoden zur Bestimmung der Ebenen
 - Nach Grad der Knoten
 - Mit Hilfe von PageRank-Algorithmus
 - Bestimmt "wichtigsten" Knoten
 - Anzahl der kürzesten Wege durch einen Knoten
 - Bestimmt Flaschenhälse
- All diese Methoden führen zu vergleichbaren Ergebnissen

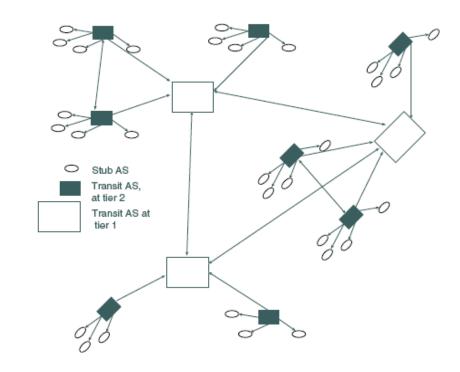


Fig. 1. Structure of the AS graph

"Comparing the structure of power-law graphs and the Internet AS graph", Sharad Jaiswal, Arnold L. Rosenberg, Don Towsley, INCP 2004



Pareto-Verteilung des Grades von ASen im Internet

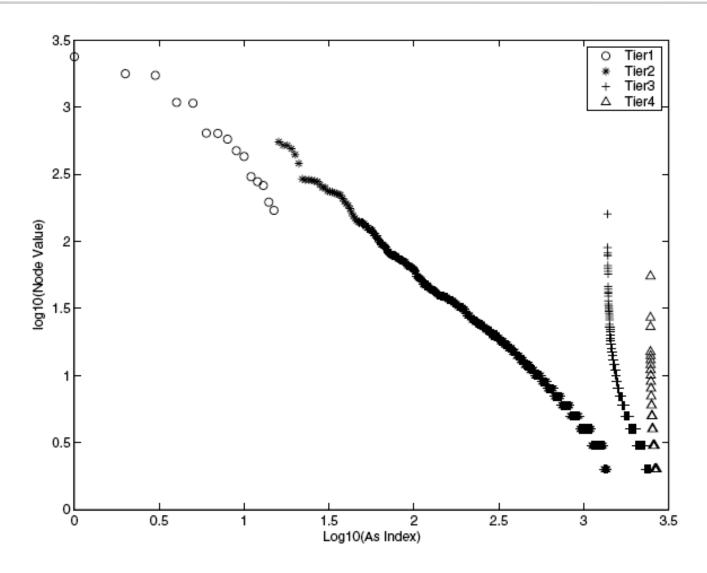


Fig. 2. Degree of ASes in different tiers

"Comparing the structure of power-law graphs and the Internet AS graph", Sharad Jaiswal, Arnold L.
Rosenberg, Don Towsley, INCP 2004



Webseitensuche

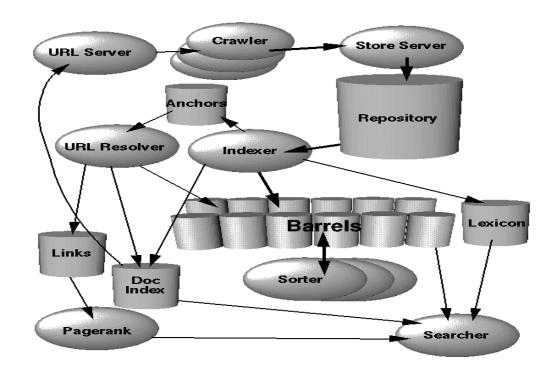
- PageRank [Brin&Page 98]
 - Vergibt jeder Web-Seite einen absoluten Rang (rank)/Autorität
 - Rang berücksichtigt Eingrad und Autorität des Eingrads
 - Idee: Seiten sind wichtig, wenn wichtige Seite auf sie zeigen

REIBURG



Die Anatomie einer Web Search Maschine

- "The Anatomy of a Large-Scale Hypertextual Web Search Engine", Sergey Brin and Lawrence Page, Computer Networks and ISDN Systems, Vol. 30, 1-6, S. 107-117, 1998
- Design des Prototyps von Google
 - Stanford University 1998
- Hauptkomponenten
 - Web Crawler
 - Indexer
 - Pagerank
 - Searcher
- Hauptunterschied zwischen Google und anderen Suchmaschinen (1998)
 - Der Pagerank Algorithmus



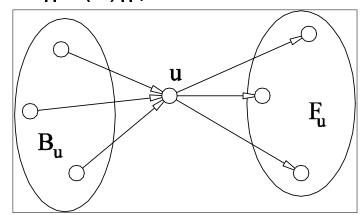


Der vereinfachte PageRank-Algorithmus

- Vereinfachter PageRank-Algorithmus
 - Rang einer Web-seite R(u) ∈ [0,1]
 - Wichtige Seiten übergeben ihre Gewicht an verlinkte Seite

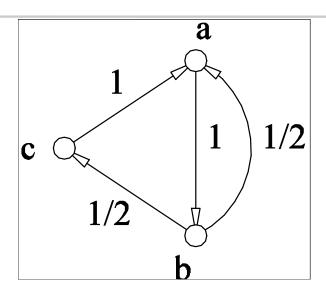
$$R(u) \leftarrow c \sum_{v \in B_u} \frac{R(v)}{|F_v|}$$

- c ist Normalisierungsfaktor so dass $||R(u)||_1 = 1$, d.h.
 - die Summe der Pageranks ist 1
- Vorgänger-Knoten B_u
- Nachfolger-Knoten F_u





Vereinfachter Pagerank-Algorithmus mit Beispiel



$$x \leftarrow 0 \cdot x + \frac{1}{2} \cdot y + 1 \cdot z$$

$$y \leftarrow 1 \cdot x + 0 \cdot y + 0 \cdot z$$

$$z \leftarrow 0 \cdot x + \frac{1}{2} \cdot y + 0 \cdot z$$

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} \leftarrow \begin{pmatrix} 0 & \frac{1}{2} & 1 \\ 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

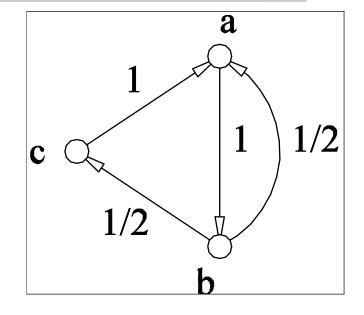
Runde	a	ъ	c
0	1	1	1
1	3/2	1	1/2
2	1	3/2	1/2
3	5/4	1/2	3/4
•	:	:	•
10	1,19	1,22	0,59
:	:	:	÷
20	1,20	1,20	0,60



Matrixdarstellung

 $\blacksquare R \leftarrow cMR$

- wobei R ein Vektor(R(1),R(2),... R(n)) ist und
- M die folgende n × n Matrix bezeichnet



$$M_{ij} := \left\{ egin{array}{l} rac{1}{|F_j|} \;, & ext{if } i \in F_j \ 0 \;, & ext{else.} \end{array}
ight.$$

$$M = \begin{pmatrix} 0 & \frac{1}{2} & 1 \\ 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{pmatrix}$$



Eigenvektor als Fixpunkt der Rekursion Stochastische Matrix

Die L1-Norm eines Vektors ist gegeben als $||x||_1 := \sum |x|$

$$||x||_1 := \sum_{i=1}^n |x|$$

chastisch, genau dann wenn $M_{ij} \geq 0$ und für Vektor x gilt: alle $j \in [n]$ gilt:

$$\sum_{i=1}^n M_{ij} = 1 ,$$

d.h. die Spaltensummen sind 1.

Definition 1 Eine $n \times n$ -Matrix M heißt sto- **Lemma 1** Für jede stochastische Matrix und

$$||M \ x||_1 \le ||x||_1$$

 $||M \ x||_1 = ||x||_1$ falls $x \ge 0$ oder $x \le 0$

 \Rightarrow Eigenwerte von M $|\lambda_i| \le 1$

Theorem 4 Für jede stochastische Matrix M gibt es einen Eigenvektor x mit Eigenwert 1, wobei $x \ge 0$ und $||x||_1 = 1$.



Nachteile des vereinfachten Pagerank-Algorithmus

- Der Web-graph hat Senken, d.h. Seiten ohne Links
 - M ist keine stochastische Matrix
- Der Web-Graph ist periodisch
 - Daher ist die Konvergenz unsicher
- Der Web-Graph ist nicht stark zusammenhängend
 - Er ist noch nicht einmal schwach zusammenhängend
 - Unterschiedliche Konvergenzvektoren
- Rank-Senken
 - Starke Zusammenhangskomponenten saugen alle Gewichte der Vorgänger auf
 - Alle Vorgänger dieser Web-Seiten verlieren ihr Gewicht



Der Pagerank-Algorithmus

- Füge einer Senke Links auf alle Web-Seiten hinzu
- Wähle gleichwahrscheinlich eine Web-Seite
 - Mit einer gewissen Wahrscheinlichkeit q < 1 führe Schritt des vereinfachten Pagerank-Algorithmus durch
 - Mit Wahrscheinlichkeit 1-q wähle gleichwahrscheinlich eine Web-Siete

$$R(u) \leftarrow \frac{1-q}{n} + q \cdot \left(\sum_{v \in B_u} \frac{R(v)}{|F_v|} + \frac{|\{v \mid F_u = \emptyset\}|}{n} \right)$$

Beachte: M ist stochastisch

tisch
$$M_{ij} := \left\{egin{array}{l} rac{1-q}{n} + rac{q}{|F_j|} \ rac{1-q}{n} \ , & ext{falls } F_j
eq \emptyset ext{ und } i
otin F_j \ rac{1}{n} \ , & ext{falls } F_j = \emptyset \end{array}
ight.$$



Eigenschaften des Pagerank-Algorithmus

- Der Graph der Matrix ist stark zusammenhängend
- Es gibt Rundwege der Länge 1
- Theorem
 - In nicht-periodischen Matrizen mit stark zusammenhängenden Graphen konvergiert der Markov-Prozess zu einem eindeutigen Eigenvektor mit Eigenwert 1

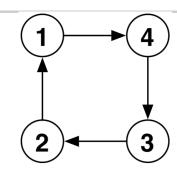
■ ⇒ Pagerank konvergiert zu diesem eindeutigen Eigenvektor



Was passiert beim Matrix-Auffüllen

Beispiel: q = 0.96, n = 4,

Graph:



0,01

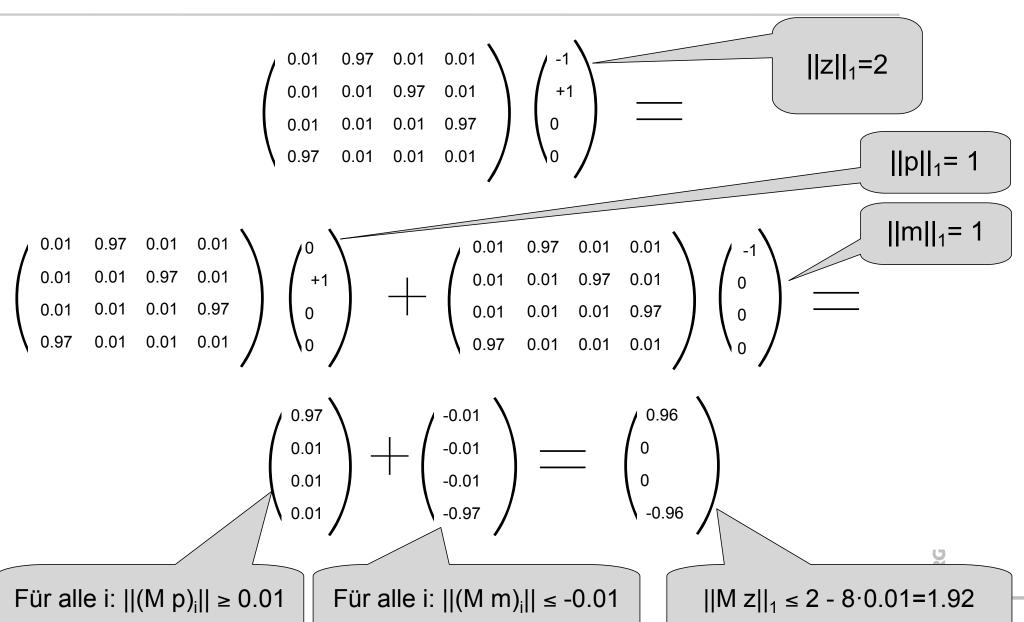
• Startvektor $x = (1,0,0,0)^T$

- Beobachtung:
 - Alle Einträge von M sind mindestens (1-q)/n (hier: 0,01)
 - Alle Einträge von Mx sind mindestens ||x||₁ (1-q)/n (hier: 0,01)
- Fakt: Für alle Vektoren x ≥ 0: (M x)_i ≥ -||x||₁(1-q)/n
- Fakt: Für alle Vektoren $x \le 0$: $(M x)_i \le ||x||_1(1-q)/n$





Was passiert mit gemischten Vektoren?



JU



Eigenschaften des Pagerank-Algorithm

- Es gibt einen eindeutigen (reellen) Eigenvektor mit Eigenwert 1 für die Matrix des Pagerank-Algorithmus
- Pagerank konvergiert zu einer (1+ε) Approximation des eindeutigen Eigenvektors in höchsten (-ln ε - ln n) / ln q Iterationen



Diskussion

- q = W'keit den vereinfachten Pagerank-Algorithmus zu benutzen
- Falls q klein ist
 - Pagerank konvergiert schneller
 - Kleinere Pfade sind relevanter
 - Weniger Strukturinformation wird benutzt
 - Die Pageranks sind ähnlicher
- Falls q groß ist
 - Pagerank konvergiert (möglicherweise) langsamer
 - Länge Wege spielen eine größere Rolle
 - Web-Senken sammeln mehr Gewicht auf
 - Daher löscht Google aus dem Web-Graph

Problem:

- Wie wählt man q
- Macht es Sinn, dass jede Web-Seite unabhängig vom Suchwort gewählt wird?





Systeme II

9. Die Struktur des Webs

Christian Schindelhauer
Technische Fakultät
Rechnernetze und Telematik
Albert-Ludwigs-Universität Freiburg









