

Böööhze Links

Wie Suchmaschinen die Gumm nudeln in der Linksuppe finden...



Prof. Dr. Mario Fischer
Fakultät für Informatik & Wirtschaftsinformatik
Schwerpunkt Electronic Commerce
Hochschule für angewandte Wissenschaft
University of Applied Sciences
Würzburg

7.-8. März 2009 - SEO-Campixx Berlin

Alle hier gezeigten Webseiten/Screenshots sind zufällig zur Demo ausgewählt.
Sie haben keinerlei Bezug zu den Aussagen der Folien!

Über welche Dimensionen sprechen wir?



- Index: Mehrere hundert Milliarden Dokumente
- eine Billion Links
- ca. 200 Datacenter weltweit, ca. 500.000 Server und 10 Exabyte Speichervolumen
- Speicherung aller gesetzten (und entfernten) Links
- Eine Suchanfrage beschäftigt mehrere hundert Server, braucht > 10 Mrd. Prozessorzyklen und fragt > 100 MB Daten ab
Antwortzeit: 0,2 Sekunden – incl. Realtime-Filter
- ca. 15.000 Suchabfragen weltweit / Sekunde
- Datenverarbeitung führt Google etwa zehn Mal effizienter / günstiger als andere IT-Unternehmen aus
- (interne) PageRank-Berechnung wichtiger Sites mehrmals pro Tag (näherungsweise)
- Geschätzter Zeitraum, bis Maschinen verstehen, was Menschen meinen...?

Warum werden Links gesetzt?

- Hinweis auf guten Content
- Gefälligkeit / Geschäftsbeziehungen
- Linktausch
- Linkkauf
- Affiliate
- Zum Behumsen des Rankings



Wie kann *gut* und *schlecht* (bööööhze) unterschieden werden?

- Hinweis auf guten Content
- Gefälligkeit / Geschäftsbeziehungen
- Linktausch
- Linkkauf
- Affiliate
- Zum Behumsen des Rankings



Unnütze Links: Herkömmliche Methoden

Kernfrage:

Wer ist mit wem verwandt?

- Gleicher IP D-Block
- Gleicher Admin C
- Gleiche Faxnummer
- Identische Google-Accounts (Webmastertools, Analytics)
- Gleicher Adsense-Code
- Gleiche Adresse im Impressum
- ...

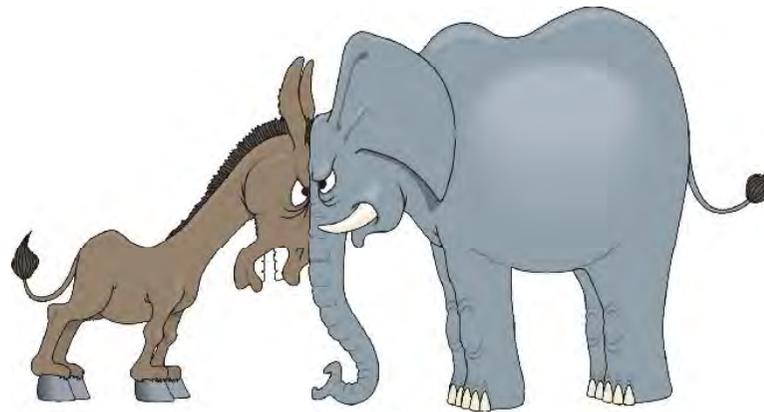


Unnütze Links: Weitere Methoden



- Aufstellen von sog. Blacklists 
- Einsatz von Heuristiken 
- Nachfolgende Datenanalyse 

Problem:
So leicht wie es sich anhört,
ist es bei weitem nicht...



Maschinelles Lernen

- Mustererkennung
 - Automatisches Training
 - Gesteuertes Training



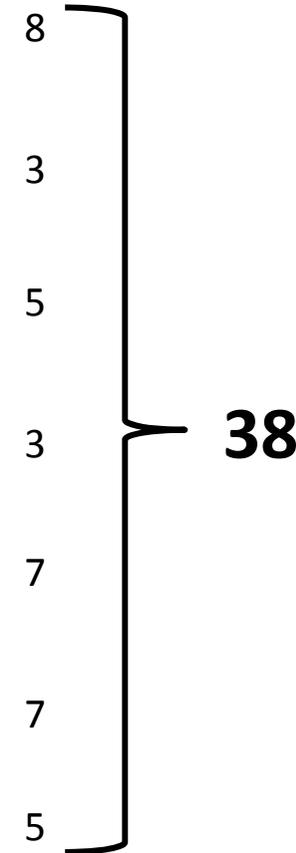
Maschinelles Lernen: Mustererkennung Beispiel

Im Prüfset von Domains stellt der selbstlernende Algorithmus z. B. fest:

- 1 Über 2 Ausgehende Links am Seitenende. Auf >75% aller Seiten einer Domain
- 2 >60% aller eingehenden Links kommen aus einem D-Class-Netz
- 3 Die Linkmuster (extern) von >50% der anlinkenden Domains sind zu 80% identisch
- 4 >40% der anlinkenden Domains sind direkt untereinander verlinkt
- 5 >60% der anlinkenden Seiten haben andere Keywordtopics als die verlinkten Seiten
- 6 Es fehlen Links von vertrauenswürdigen Domains
- 7 Am Linknetzwerk hängen Sites mit bestimmten Keywords wie z. B. „Suchmaschinenoptimierung“

Entscheidungstabelle: ab 20 De-Ranking
 ab 30 P4 Penalty
 ab 40 tschüss

PenaltyPoints



Test,
Feinjustierung,
Collateral-Prüfung



Alle Zahlen beispielhaft! Um Gottes willen!

Maschinelles Lernen: Erkennungsraten

Wie gut sind die Maschinen heute?

Wahrscheinlichkeit, eine **Linkfarm** zu entdecken:

95,1%

Wahrscheinlichkeit, einen **Spam-Linkhub** (Kern) zu entdecken:

~80%

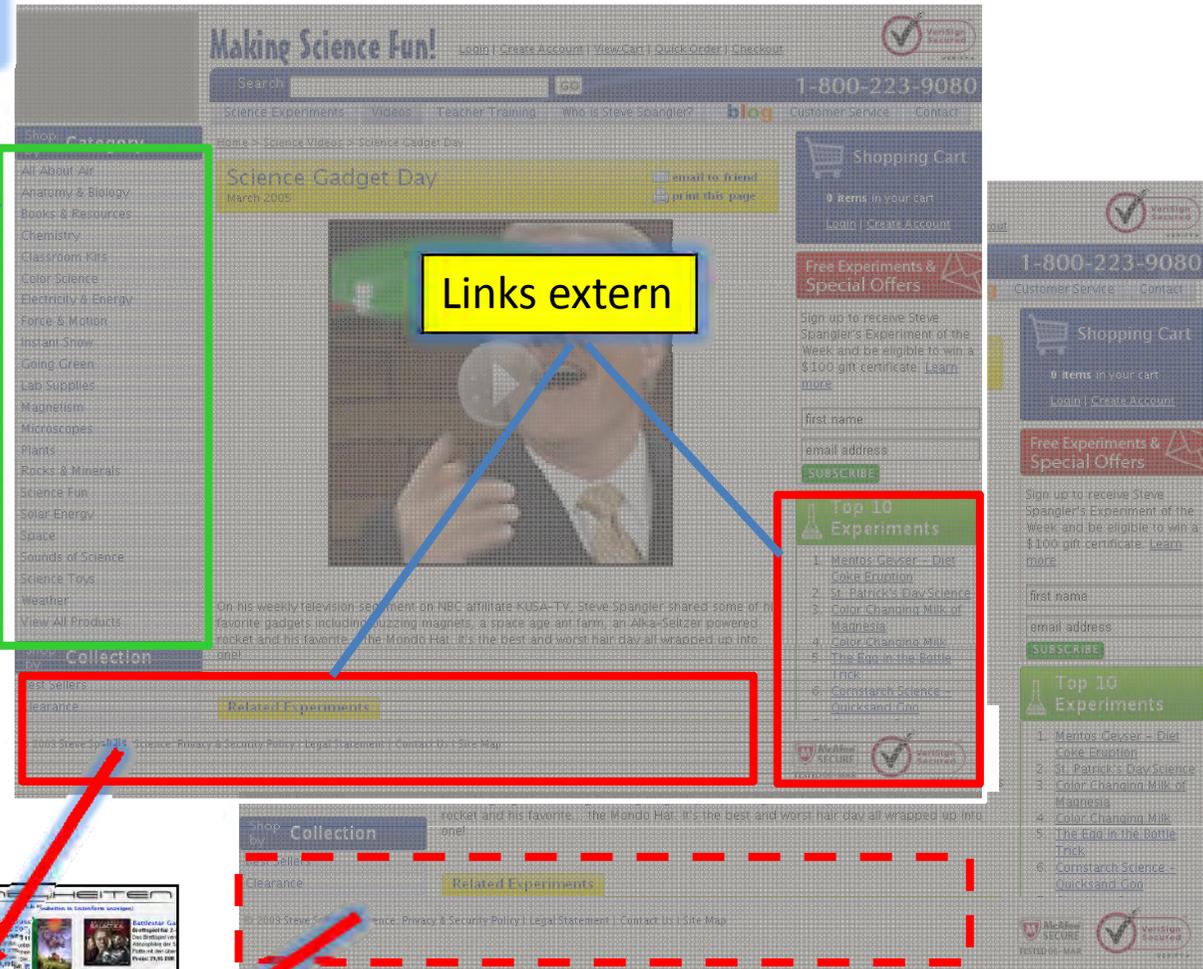
Maschinelles Lernen: Mustererkennung Beispiel

Sites erkennen, die Links verkaufen:

Mögliche Kriterien:

Links intern

- Linkverhältnis intern/extern Domain und/oder Seite
- Wechselnde Linktexte bei gleichem Linkziel
- Linktextwechsel über die Zeit
- Verhältnis Deeplinks zu Startseitenlinks
- Anzeichen von PR-Sculpting
- Fehlen von Domainnamen in den Linktexten
- Externe Links ohne inhaltlichen Zusammenhang zur Site
- Unterschiedliche Topics der extern verlinkten Sites
- Werthaltige/“verdächtige“ Keywords auf den Linktexten
- Alter des Linkbestehens
- Linkwechsel in bestimmten Intervallen (z. B. monatlich)
- Keine/wenig Intext-Links
- Nicht natürlich wirkende Geschwindigkeit beim Linkaufbau



Im Quellcode zusammenstehende, externe Links zu
Immer den gleichen Domains

Maschinelles Lernen: Mustererkennung Beispiel

Entscheidungsbaum des C4.5 Algorithmus von Brian D. Davidson:

```
FromPath contains tilde = 0:
|   From contains Simple gTLD = 0:
|   |   Domains are same = 0: -1 (35.0)
|   |   Domains are same = 1: 1 (30.0/1.0)
|   From contains Simple gTLD = 1:
|   |   FromPage has > 200 links = 1: 1 (812.0)
|   |   FromPage has > 200 links = 0:
|   |   |   Same DNS servers = 0:
|   |   |   |   Same contact email = 1: 1 (5.0)
|   |   |   |   Same contact email = 0:
|   |   |   |   |   Pages share > 10% links = 0: -1 (62.0/7.0)
|   |   |   |   |   Pages share > 10% links = 1: 1 (3.0/1.0)
|   |   |   |   Same DNS servers = 1:
|   |   |   |   |   Complete hostnames are same = 1: 1 (374.0)
|   |   |   |   |   Complete hostnames are same = 0:
|   |   |   |   |   |   Domains are same = 0: 1 (68.0/1.0)
|   |   |   |   |   |   Domains are same = 1:
|   |   |   |   |   |   |   Pages share > 20% links = 1: 1 (18.0)
|   |   |   |   |   |   |   Pages share > 20% links = 0:
|   |   |   |   |   |   |   |   FromPage has > 100 links = 1: 1 (14.0)
|   |   |   |   |   |   |   |   FromPage has > 100 links = 0:
|   |   |   |   |   |   |   |   |   FromPage has <= 5 links = 1: 1 (2.0)
|   |   |   |   |   |   |   |   |   FromPage has <= 5 links = 0:
|   |   |   |   |   |   |   |   |   |   1 path components are same = 1: -1 (5.0/2.0)
|   |   |   |   |   |   |   |   |   |   1 path components are same = 0:
|   |   |   |   |   |   |   |   |   |   |   FromPage has <= 10 links = 0: 1 (21.0/6.0)
|   |   |   |   |   |   |   |   |   |   |   FromPage has <= 10 links = 1: -1 (5.0/2.0)
FromPath contains tilde = 1:
|   1 path components are same = 0: -1 (56.0)
|   1 path components are same = 1: 1 (26.0)
```

Maschinelles Lernen: Mustererkennung Beispiel

Spammer sind Menschen und Menschen machen Fehler: Footprints

- Identische Titel/Description
- Vergessene Comments im Quelltext
- Im Pfad stehen gleiche Begriffe wie /user, /home /~ oder ähnliche, identische Zeichenketten (gleiche Software!)
- Gleiche E-Mail-Adresslinks
- Domain A hat > x% identische Ausgangslinks wie Domain B

```
<div id=showdiv style='visibility:hidden;position:absolute;' onclick='showDiv();'>
<!-- insert your hidden text here. Do not forget to <H1></H1> your keywords -->
<h1>Music box</h1> - A mechanical music device, can be disc music box or cylinder r
<!-- end insert here -->
```

Maschinelles Lernen: Positivmuster erkennen

Set vertrauenswürdiger Sites



DVD Player Shop

Startseite



Philips DVP 3260 DVD Player (DivX-zertifiziert)
 EUR 53,45



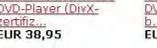
Philips DVD Player
 EUR 79,...



Western Digital WD TV/HDTV Media Player
 EUR 101,00



Samsung DVD-P 181 DVD-Player (DivX-zertifiziert)
 EUR 38,95



LG GH22NP20 AUBA10B DVD-Brenner 22x P-ATA
 EUR 23,95



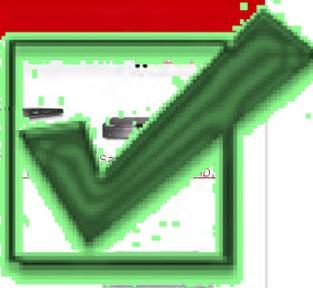
Philips DVP 3120 / 12 DVD-Player silber
 EUR 39,97



Panasonic DMP BD 35 EG K Blu-ray Player
 EUR 238,72



Samsung SE-S2240/EURN externer DVD-Brenner
 EUR 54,95



Inbound-Linkmuster:

- Trustrank 6,4
- 1.250 Links, davon
 - 920 diff. Domains
 - 633 diff. IP-Class
 - 540 mit Domainnamen

Maschinelles Lernen: Semantisch: Hilltop / HITS

Identifikation von Expertenseiten (Hubs) und Autoritäten (Authorities)

ab PR 4



„Experten“

müssen in den Top 200 sein und zu mindesten 5 **unabhängigen***, themenrelevanten Seiten verweisen

ab PR 4



ab PR 4

ab PR 4



„Autorität“

>2 **unabhängige*** „Experten“-Links

*nach der IP-Nummer und Domainnamen!

Kandidat für die Top 10

Berechnung eines **Relevanzwertes**, der vom Verlinkungsgrad der Expertenseiten abhängt.

mögl. wenig Links nach außen

Derzeit wahrscheinlich nur bei allgemeinen, werthaltigen Begriffen einsetzbar, da rechenintensiv.
Hilltop: Entwickelt von Krishna Bharat und George A. Mihaila, Google hat das Patent erworben ;-)

Quality-Rank / Vertrauenswürdigkeit: Beispiel



Login

Home Blog **Quality Rank Tool** sc

<http://www.website-boosting.de>

Crawled on: Mar 6, 2009 06:03am EST

65.15% trustworthy

Alexa Ranking

1.66 / 2 points

Alexa rating is your current ranking amongst built an unparalleled database of information this information can be found on Alexa's Site

Tips For Improvement

To get access Tips for Improver

Backlink Profile

10.00 / 10 points

Your backlink profile takes into consideration compares it to the total volume of links across index of your site, it makes your backlink profile directly to content as opposed to the index.

Tips For Improvement

To get access Tips for Improver

Directories

7 / 10 points

Directories do not hold the value they once did are still considered to hold value. These directories constant fluctuation and bring stability to your

Social Bookmarking

1.50 / 3 points

If you have quality content, and a decent user larger Social Bookmarking sites. For most content provided it is not spam.

Tips For Improvement

To get access Tips for Improver

Title Search

4.92 / 5 points

Title search takes the first four words in your these words combined.

Domain Age

0.00 / 10 points

Domain age is important as spammers tend to Yahoo trust sites that are registered long enough the long haul. Also, having your domain registered

Tips For Improvement

To get access Tips for Improver

External Links

15 / 15 points

External links checks to see how many outgoing dofollow links on any one page. Link exchanges are considered against Google's webmaster keep outbound homepage links on your homepage

For more information, see [Google's webmaster](#)

Tips For Improvement

To get access Tips for Improver

Google Page Rank

3.44 / 5 points

While pagerank means a lot less than it once that are PR5+ still tend to be very strong in

Tips For Improvement

To get access Tips for Improver

Inbound Links

28.15 / 50 points

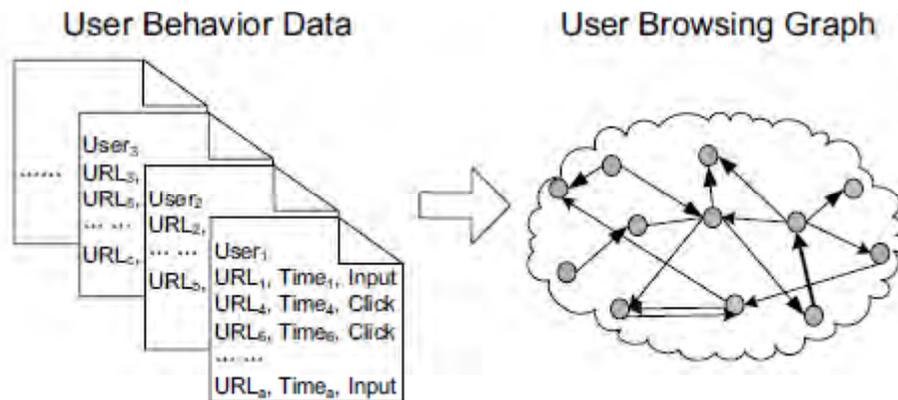
Inbound links is the total volume of site with volume is important for more competitive ranking

Hinweis: Bei all diesen Werkzeugen:
Werte immer mit Vorsicht / Skepsis betrachten

Guckst Du: www.linkvoodoo.com/score_site.php

BrowseRank / Wo klicken die Menschen?

- Von Microsoft Research veröffentlicht
- Soll die Schwächen des Linkgraphen-basierten PageRanks beheben
- Bewegungen zwischen und Verweildauer auf Webseiten der Nutzer werden - anonym- via Toolbars und Phone-Home (?) erhoben – auch URL-Eingabe oder Linkklick



| URL | TIME | TYPE |
|------------------------------|----------------------|-------|
| http://aaa.bbb.com/ | 2007-04-12, 21:33:05 | INPUT |
| http://aaa.bbb.com/1.htm | 2007-04-12, 21:34:11 | CLICK |
| http://ccc.ddd.org/index.htm | 2007-04-12, 21:34:52 | CLICK |
| http://eee.fff.edu/ | 2007-04-12, 21:39:03 | INPUT |
| ... | ... | ... |

Quelle: <http://research.microsoft.com/en-us/people/tyliu/fp032-liu.pdf>, S. 3

$$\pi_i = \frac{\bar{\pi}_i}{q_{ii}} \quad - \frac{q_{ij}}{q_{ii}} = \begin{cases} \alpha \frac{\bar{w}_{ij}}{\sum_{k=1}^{N+1} \bar{w}_{ik}} + (1-\alpha)\sigma_j, & i \in V, j \in \tilde{V}^* \\ \sigma_j, & i = N+1, j \in V \end{cases}$$

- **Problematisch: Verstärkung von Trampelpfaden, Neue (interessante) Seiten/Sites werden in der Suchmaschine schlecht gefunden, Signifikanz „kleiner“ Sites möglicherweise zu niedrig**
 - -> Nur wo viele sind bzw. schon waren, ist es gut?
 - -> Braucht man zum Finden von Powersites überhaupt eine SuMa?
 - -> Wie PageRank auch manipulierbar, z. B. über Botnetze

Quelle: <http://research.microsoft.com/en-us/people/tyliu/fp032-liu.pdf>, S. 2

* Formel hat keinen Bezug zur Abbildung, sieht aber schön kompliziert aus!

BrowseRank / Wo klicken die Menschen?

Unterschiede PageRank / TrustRank / BrowseRank in der Bewertung

| No. | PageRank | TrustRank | BrowseRank |
|-----|------------------------|------------------------|------------------------|
| 1 | adobe.com | adobe.com | <i>myspace.com</i> |
| 2 | passport.com | yahoo.com | msn.com |
| 3 | msn.com | google.com | yahoo.com |
| 4 | microsoft.com | msn.com | <i>youtube.com</i> |
| 5 | yahoo.com | microsoft.com | live.com |
| 6 | google.com | passport.net | <i>facebook.com</i> |
| 7 | mapquest.com | ufindus.com | google.com |
| 8 | miibeian.gov.cn | <i>sourceforge.net</i> | ebay.com |
| 9 | w3.org | <i>myspace.com</i> | <i>hi5.com</i> |
| 10 | godaddy.com | <i>wikipedia.org</i> | <i>bebo.com</i> |
| 11 | statcounter.com | phpbb.com | <i>orkut.com</i> |
| 12 | apple.com | yahoo.co.jp | aol.com |
| 13 | live.com | ebay.com | <i>friendster.com</i> |
| 14 | xbox.com | nifty.com | <i>craigslist.org</i> |
| 15 | passport.com | mapquest.com | google.co.th |
| 16 | <i>sourceforge.net</i> | cafepress.com | microsoft.com |
| 17 | amazon.com | apple.com | <i>comcast.net</i> |
| 18 | paypal.com | infoseek.co.jp | <i>wikipedia.org</i> |
| 19 | aol.com | miibeian.gov.cn | <i>pogo.com</i> |
| 20 | <i>blogger.com</i> | <i>youtube.com</i> | <i>photobucket.com</i> |

Wissenschaftliche Quellen

- Zum Hilltop-Algorithmus: <http://ftp.cs.toronto.edu/pub/reports/csrq/405/hilltop.html>
- Zum HITS-Algorithmus, Univ. of Missouri-Columbia: <http://www2002.org/CDROM/refereed/643/>
- Deutsche Seminararbeit zum HITS-Algorithmus:
<http://www-dbs.informatik.uni-heidelberg.de/teaching/ss2005/seminar/Ausarbeitung/1-2.pdf>
- Weiterentwicklung von HITS zum ACR-Algorithmus (mehr Themenbezug der Links + Gewichtung):
<http://www.cs.cornell.edu/home/kleinber/www98-arc.pdf>
- Stanford-Publikation zu ACR: <http://theory.stanford.edu/people/raghavan/www7/181.html>
- Zum Trust-Rank-Algorithmus, Univ. Stanford: <http://www.vldb.org/conf/2004/RS15P3.PDF>
- Zum BrowseRank von Microsoft: <http://research.microsoft.com/en-us/people/tyliu/fp032-liu.pdf>
- PageRank-Patent von 2006 mit Nutzerverhalten: <http://tinyurl.com/y9mnts>
- Linknachbarschaften C4.5-Algo: <http://www.cse.lehigh.edu/~brian/pubs/2000/aaais/aaai2000ws.pdf>
- Linknachbarschaften am unstimigen Spachmodell erkennen: <http://tinyurl.com/bty6f4>
- Zufallsgesteuert Linkfarmen finden: http://airweb.cse.lehigh.edu/2007/papers/paper_116.pdf
- Spamseiten über Content-Analyse finden: <http://research.microsoft.com/pubs/65140/www2006.pdf>
- Linkanalyse zur Spamererkennung (Univ. Rom und Yahoo!): <http://tinyurl.com/cambcv>

Bissi Eigenmarketing: Blog auf www.Website-Boosting.de und
<http://twitter.com/mariofischer>