Authoritative Sources in a Hyperlinked Environment

Jon M. Kleinberg Presenter: Zhe Zhao

Overview

- Background and Motivation
- Approach Authorities & Hubs
 - Construct a focused subgraph based on query
 - Computing ``hubs'' and ``authorities''
 - Iterative Algorithm and its convergence
- Expansions:
 - Similar-Page Queries
 - Multiple Set of Hubs and Authorities
- Related Work
- Conclusions

Types of Queries

- Three Types of Queries
 - Specific queries
 - Does Netscape support the JDK 1.1 code-signing API?
 - Broad-topic queries
 - Find information about the Java programming language.
 - Similar-page queries
 - Find pages `similar' to java.sun.com.

Types of Queries

- Three Types of Queries
 - Specific queries
 - Does Netscape support the JDK 1.1 code-signing API?
 - Broad-topic queries Abundance Problem!

- Find information about the Java programming language.
- Similar-page queries
 - Find pages `similar' to java.sun.com.

Background and Motivation

• Hard to imagine no ranking algorithms in search engine.

Authoritative Sources in a Hyperlinked Environment Page 49 of about 39,000 results (0.35 seconds) Review - Authoritative Sources in a Hyperlinked Environment. - Pu... www.pubzone.org/dblp/journals/dr/Mendelzon00 +1 Publication Info · Discussion / Material · Links · Rating · Subscribe. Review -Authoritative Sources in a Hyperlinked Environment. ... Webmining Techniques for Program Comprehension Andy Zaidman To... www.docstoc.com/.../Webmining-Techniques-for-Program-Compreh... Apr 15, 2009 - Authoritative sources in a hyperlinked environment. Journal of the ACM, 1999 Hubs and Authorities n n n Recursive definition of hub and ... source - DBLP dblp.cloudmining.net/search?..... +1 Jon M. Kleinberg: Authoritative Sources in a Hyperlinked Environment. ... authoritative (1) environment (1) hyperlinked (1) source (1) ... xls:Authoritative Sources ina Hyperlinked Environment JonM ... www.searchuu.com/.../Authoritative+Sources+ina+Hyperlinked+Envi... eBook: Authoritative Sources ina Hyperlinked Environment JonM. Kleinbergy. 0 Results. ©2011 www.searchuu.com. All rights reserved.

Ranking algorithms in web search

- Modern search engines may return millions of pages for a single query. This amount is prohibitive to preview for human users.
- Ranking algorithms will process the search results and only show the most useful information to the search engine user.

Ranking algorithms in web search

Ļ

Authoritative Sources in a Hyperlinked Environment

About 39,000 results (0.16 seconds)

Scholarly articles for Authoritative Sources in a Hyperlinked Environment Authoritative sources in a hyperlinked environment - Kleinberg - Cited by 6005 ... for topic distillation in a hyperlinked environment - Bharat - Cited by 908 Automatic resource compilation by analyzing hyperlink ... - Chakrabarti - Cited by 805

[PDF] Authoritative Sources in a Hyperlinked Environment - Cornell ...

www.cs.cornell.edu/home/kleinber/auth.pdf +1 File Format: PDF/Adobe Acrobat - Quick View by JM Kleinberg - Cited by 6005 - Related articles HITs is a link-structure analysis algorithm which ranks pages by "authorities" (pages which have many incoming links and provide the best **source** of information ...

Jon Kleinberg's Homepage

www.cs.cornell.edu/home/kleinber/ +1

Web Analysis and Search: Hubs and Authorities. J. Kleinberg. Authoritative ...

Show more results from cornell.edu

Authoritative sources in a hyperlinked environment

dl.acm.org/citation.cfm?id=324140

Ranking algorithms in web search

- To find a small set of most ``authoritative'' pages relevant to the query.
- Authority
 - Most useful/relevant/helpful results of a query.
 - ``java'' java.com
 - ``harvard'' harvard.edu
 - ``search engine'' powerful search engines.

Challenge of content-based ranking

- Most useful webpage don't have the keyword
 - Query: ``Harvard''
 - 49 ``Harvard'' in <u>www.harvard.edu</u>
 - 357 ``Harvard'' in <u>http://en.wikipedia.org/wiki/Harvard University</u>
- Pages are not sufficiently descriptive

– ``automobile manufacturers'' in Honda or Toyota

Analysis of Link Structure

- Hyperlinks encode human latent judgment
- Reasons:
 - Navigation:
 - Back to top...
 - Relevant:
 - Webpage discussing java link to java.com
 - Popular:
 - www.yahoo.com www.google.com
 - Advertisement

Analysis of Link Structure

- Hyperlinks encode human latent judgment
- Reasons:
 - Navigation:
 - Back to top...
 - Relevant:
 - Webpage discussing java link to java.com
 - Popular:
 - <u>www.yahoo.com</u> <u>www.google.com</u>
 - Advertisement
- Some of the Reasons may be very helpful to find authoritative results.

- Or Hypertext-Induced Topic Search(HITS) developed by Jon Kleinberg, while visiting IBM Almaden
- IBM expanded HITS into Clever.
- Authorities
 - pages that are relevant and are linked to by many other pages
- Hubs

pages that link to many related authorities

- Intuitive Idea to find authoritative results using link analysis:
 - Not all hyperlinks related to the conferral of authority.
 - Find the pattern authoritative pages have:
 - Authoritative Pages share considerable overlap in the sets of pages that point to them.

- Intuitive Idea to find authoritative results using link analysis:
 - Not all hyperlinks related to the conferral of authority.
 - Find the pattern authoritative pages have:
 - Authoritative Pages share considerable overlap in the sets of pages that point to them.





• First Step:

Constructing a focused subgraph of the WWW based on query

- Second Step
 - Iteratively calculate authority weight and hub weight for each page in the subgraph

- Why not find authorities on the entire WWW?
 - The algorithm is non-trivial.
 - not necessary when there is a query.
- Objective: S $_{\sigma}$
 - S $_{\sigma}$ is relatively small.
 - S $_{\sigma}$ is rich in relevant pages.
 - S $_{\sigma}$ contains most (or many) of the strongest authorities
- Solution:
 - Generate a Root Set $Q\sigma$ from text-based search engine
 - Expand the root set

Subgraph (σ, εt,d)

σ : a query string
ε : a text-based search engine.
t, d: natural numbers.
Let R denote the top t results of ε on σ

Set S := R For each page $p \in R$ Let $\Gamma^+(p)$ denote the set of all pages p points to. Let $\Gamma^-(p)$ denote the set of all pages pointing to p. Add all pages in $\Gamma^+(p)$ to S. If $(\Gamma^-(p)) < d$ then Add all pages in $\Gamma(p)$ to S. Else Add an arbitrary set of d pages from $\Gamma^-(p)$ to S End Root Set

Subgraph (σ, εt,d)

σ : a query string
E : a text-based search engine.
t, d: natural numbers.
Let R denote the top t results of E on σ

Set S := R For each page $p \in R$ Let $\Gamma^+(p)$ denote the set of all pages p points to. Let $\Gamma^-(p)$ denote the set of all pages pointing to p. Add all pages in $\Gamma^+(p)$ to S. If $(\Gamma^-(p)) < d$ then Add all pages in $\Gamma(p)$ to S. Else Add an arbitrary set of d pages from $\Gamma^-(p)$ to S

End



Subgraph (σ, εt,d)

σ : a query string
ε : a text-based search engine.
t, d: natural numbers.
Let R denote the top t results of ε on σ

Set S := R For each page $p \in R$ Let $\Gamma^+(p)$ denote the set of all pages p points to. Let $\Gamma^-(p)$ denote the set of all pages pointing to p. Add all pages in $\Gamma^+(p)$ to S. If $(\Gamma^-(p)) < d$ then Add all pages in $\Gamma(p)$ to S. Else

Add an arbitrary set of *d pages from* $\Gamma^-(p)$ to S



Computing Hubs and Authorities

- Rules:
 - A good hub points to many good authorities.
 - A good authority is pointed to by many good hubs.
 - Authorities and hubs have a mutual reinforcement relationship.



Computing Hubs and Authorities

- Let authority score of the page i be x(i), and the hub score of page i be y(i).
- mutual reinforcing relationship:
- I step:

$$x(i) = \sum_{(j,i)\in E} y(j)$$

• O step:

$$y(i) = \sum_{(i,j)\in E} x(j)$$

• 1st Iteration



- 1st Iteration
- I Step



- 1st Iteration
- I Step
- O Step



- 2nd Iteration
- I Step



- 2nd Iteration
- I Step
- O Step



- 2nd Iteration
- I Step
- O Step

••••

•••



Iterate(G,k)

G: a collection of n linked pages

k: a natural number

Let z denote the vector $(1, 1, 1, \ldots, 1) \in \mathbf{R}^n$.

Set $x_0 := z$.

Set $y_0 := z$.

For i = 1, 2, ..., k

Apply the \mathcal{I} operation to (x_{i-1}, y_{i-1}) , obtaining new *x*-weights x'_i . Apply the \mathcal{O} operation to (x'_i, y_{i-1}) , obtaining new *y*-weights y'_i . Normalize x'_i , obtaining x_i . Normalize y'_i , obtaining y_i .

Initialization

End

 $\operatorname{Iterate}(G,k)$

G: a collection of n linked pagesk: a natural numberLet z denote the vector $(1, 1, 1, \dots, 1) \in \mathbb{R}^n$.Set $x_0 := z$.Set $y_0 := z$.For $i = 1, 2, \dots, k$ I Step

Apply the \mathcal{I} operation to (x_{i-1}, y_{i-1}) , obtaining new x-weights x'_i .

Apply the \mathcal{O} operation to (x'_i, y_{i-1}) , obtaining new *y*-weights y'_i . Normalize x'_i , obtaining x_i .

Normalize y'_i , obtaining y_i .

End

 $\operatorname{Iterate}(G,k)$

G: a collection of n linked pages k: a natural number Let z denote the vector $(1, 1, 1, \ldots, 1) \in \mathbf{R}^n$. Set $x_0 := z$. Set $y_0 := z$. For i = 1, 2, ..., kApply the \mathcal{I} operation to (x_{i-1}, y_{i-1}) , obtaining new x-weights x'_i . Apply the \mathcal{O} operation to (x'_i, y_{i-1}) , obtaining new y-weights y'_i . Normalize x'_i , obtaining x_i . O Step Normalize y'_i , obtaining y_i . End

 $\operatorname{Iterate}(G,k)$

G: a collection of n linked pages k: a natural number Let z denote the vector $(1, 1, 1, \ldots, 1) \in \mathbf{R}^n$. Set $x_0 := z$. Set $y_0 := z$. For i = 1, 2, ..., kApply the \mathcal{I} operation to (x_{i-1}, y_{i-1}) , obtaining new x-weights x'_i . Apply the \mathcal{O} operation to (x'_i, y_{i-1}) , obtaining new *y*-weights y'_i . Normalize x'_i , obtaining x_i . Normalization Normalize y'_i , obtaining y_i . End

Proof of Convergence

- A Matrix Perspective:
 - Denote A as adjacent matrix of the subgraph
 - -Istep: $x(i) = \sum_{(j,i)\in E} y(j) \longrightarrow x = A^T y$

- O step:

$$y(i) = \sum_{(i,j)\in E} x(j) \longrightarrow y = Ax$$

• Converge to eigenvector.

A Statistical View of HITS

- 1st Eigenvalue of AA^T = singular value of A
- 1st Eigenvector of AA^{T} = transform vector to the 1st principal component.
- Principal Component:
 - Matrix A \rightarrow a set of vectors.
 - The dimension where vectors significantly distributed
 ^{0.2}
 ^{1.2}
 ^{1.2}



A Statistical View of HITS

- The weight of authority equals the contribution of transforming the dataset to first principal component.
 - Importance of this vector for the distribution of whole dataset.
- From the statistical view:
 - HITS can be implemented by PCA
 - HITS is different from clustering using dimensionality reduction.
 - The number of samples of PCA is limited.

Example of Results:

Query "censorship" : Authorities

- .378 http://www.e.org/
- .344 http://www.e.org/blueribbon.html
- .238 http://www.cdt.org/
- .235 http://www.vtw.org/
- .218 http://www.aclu.org/

EFFweb The Electronic Frontier Foundation The Blue Ribbon Campaign for Online Free Speech The Center for Democracy and Technology Voters Telecommunications Watch ACLU: American Civil Liberties Union

Query "search engines" : Authorities

.346 <u>http://www.yahoo.com/</u>
.291 <u>http://www.excite.com/</u>
.239 <u>http://www.mckinley.com/</u>
.231 <u>http://www.lycos.com/</u>
.231 http://www.altavista.digital.com/

Yahoo! Excite Welcome to Magellan! Lycos Home Page AltaVista: Main Page

Expansions: Similar Page Queries

• Similar-page queries

Find pages `similar' to <u>www.honda.com</u>

- Applying HITS on Similar-Page Queries
 - Find t pages pointing to p as root set

Query "www.honda.com" : Authorities

.202 http://www.toyota.com/ .199 http://www.honda.com/ .192 http://www.ford.com/ .173 http://www.ford.com/ .162 http://www.bmwusa.com/ .162 http://www.volvocars.com/ .158 http://www.saturncars.com/ .155 http://www.nissanmotors.com/ .145 http://www.audi.com/ Welcome to @Toyota Honda Ford Motor Company BMW of North America, Inc. VOLVO Welcome to the Saturn Web Site NISSAN - ENJOY THE RIDE Audi Homepage

Expansions: Similar Page Queries

- Why it works?
 - Does this mean that toyota.com offers a friendly hyperlink to honda.com?

Expansions: Similar Page Queries

- Why it works?
 - Does this mean that toyota.com offers a friendly hyperlink to honda.com?
 - Hubs from the root set make it possible.



Expansions: Multiple set of Hubs and Authorities

- Varies of Reasons for this:
 - The query string may have several very different meanings.
 - The current algorithm cannot find all the meanings.
 - Hubs of different meanings may not have overlap.
 - Only one type of hubs and authorities won out after iterations of mutually reinforcing.
- ``Natural'' solution:
 - Use other eigenvectors.

-- use other principal components

Connections to Related Work

- Standing, Impact, and Influence
 - Social Network
 - Scientific Citations
- Hypertext and WWW rankings
- Clustering of Link Structures.

PageRank v.s. HITS

- PageRank
 - Computed for all web pages stored prior to the query
 - Computes authorities only
 - Fast to compute

- HITS
 - Performed on the subset generated by each query.
 - Computes authorities and hubs
 - Easy to compute, real-time execution is hard.

Which one is more suitable for large scale data set??

Conclusion

- Motivation
 - Ranking is necessary.
 - Hyperlink information is useful
- Authorities & Hubs.
 - Find authoritative pages.
 - Construct subgraph
 - Mutually reinforcing relationship
 - Iterative algorithm
- Compare to PageRank