Information Networks

Link Analysis Ranking Lecture 9





PageRank algorithm [BP98]

- § Good authorities should be pointed by good authorities
- § Random walk on the web graph
 - § pick a page at random
 - § with probability 1- α jump to a random page
 - § with probability α follow a random outgoing link
- § Rank according to the stationary distribution

§
$$PR(p) = \alpha \sum_{q \to p} \frac{PR(q)}{|F(q)|} + (1 - \alpha) \frac{1}{n}$$



- 1. Red Page
- 2. Purple Page
- 3. Yellow Page
- 4. Blue Page
- 5. Green Page



A PageRank algorithm

§ Performing vanilla power method is now too expensive – the matrix is not sparse

$$q^{0} = v$$

$$t = 1$$

repeat

$$q^{t} = (P^{\prime})^{T} q^{t-1}$$

$$\delta = \left\|q^{t} - q^{t-1}\right\|$$

$$t = t + 1$$

until $\delta < \epsilon$

Efficient computation of $y = (P'')^T x$

$$y = \mathbf{a} \mathbf{P}^{\mathsf{T}} \mathbf{x}$$
$$\boldsymbol{\beta} = \|\mathbf{x}\|_{1} - \|\mathbf{y}\|_{1}$$
$$\mathbf{y} = \mathbf{y} + \boldsymbol{\beta} \mathbf{v}$$

P = normalized adjacency matrix P' = P + dv^T, where d_i is 1 if i is sink and 0 o.w. P'' = α P' + (1- α)uv^T, where u is the vector of all 1s



Hubs and Authorities [K98]

- § Authority is not necessarily transferred directly between authorities
- § Pages have double identity
 - § hub identity
 - § authority identity
- § Good hubs point to good authorities
- § Good authorities are pointed by good hubs





HITS Algorithm

- § Initialize all weights to 1.
- § Repeat until convergence
 - § O operation : hubs collect the weight of the authorities

$$h_i = \sum_{j:i \to j} a_j$$

§ I operation: authorities collect the weight of the hubs

$$a_i = \sum_{j: j \to i} h_j$$

§ Normalize weights under some norm



- § ...in the beginning...
- § previous work
- § some more algorithms
- § some experimental data
- § a theoretical framework

§ Problems with HITS

- § multiple links from or to a single host
 - view them as one node and normalize the weight of edges to sum to 1
- § topic drift: many unrelated pages
 - prune pages that are not related to the topic
 - weight the edges of the graph according the relevance of the source and destination
- § Other approaches?



§ Perform a random walk alternating between hubs and authorities





- § Start from an authority chosen uniformly at random
 - § e.g. the red authority





- § Start from an authority chosen uniformly at random
 - § e.g. the red authority
- § Choose one of the in-coming links uniformly at random and move to a hub
 - § e.g. move to the yellow authority with probability 1/3





- § Start from an authority chosen uniformly at random
 - § e.g. the red authority
- § Choose one of the in-coming links uniformly at random and move to a hub
 - § e.g. move to the yellow authority with probability 1/3
- § Choose one of the out-going links uniformly at random and move to an authority
 - § e.g. move to the blue authority with probability 1/2





- § In matrix terms
 - § A_c = the matrix A where columns are normalized to sum to 1
 - § A_r = the matrix A where rows are normalized to sum to 1
 - § p = the probability state vector
- § The first step computes § $y = A_c p$
- § The second step computes § $p = A_r^T y = A_r^T A_c p$
- § In MC terms the transition matrix § $P = A_r A_c^T$



 $y_2 = 1/3 p_1 + 1/2 p_2$ $p_1 = y_1 + 1/2 y_2 + 1/3 y_3$



- § The SALSA performs a random walk on the authority (right) part of the bipartite graph
 - § There is a transition between two authorities if there is a BF path between them





§ Stationary distribution of SALSA

§ authority weight of node i =

fraction of authorities in the hub-authority community of i

fraction of links in the community that point to node i

§ Reduces to InDegree for single community graphs



X



- § Rank a node according to the reachability of the node
- § Create the neighborhood
 by alternating between
 Back and Forward steps
- § Apply exponentially decreasing weight as you move further away



VV =



- § Rank a node according to the reachability of the node
- § Create the neighborhood
 by alternating between
 Back and Forward steps
- § Apply exponentially decreasing weight as you move further away



W = 3 * 1



- § Rank a node according to the reachability of the node
- § Create the neighborhood
 by alternating between
 Back and Forward steps
- § Apply exponentially decreasing weight as you move further away



VV = 3 + (1/2) * 0



- § Rank a node according to the reachability of the node
- § Create the neighborhood
 by alternating between
 Back and Forward steps
- § Apply exponentially decreasing weight as you move further away



VV = 3 + (1/4) * 1



Implicit properties of the HITS algorithm

§ Symmetry

- § both hub and authority weights are defined in the same way (through the sum operator)
- § reversing the links, swaps values

§ Equality

§ the sum operator assumes that all weights are equally important



A bad example



- § The red authority seems better than the blue authorities.
 - § quantity becomes quality



- § Is the hub quality the same as the authority quality?
 - § asymmetric definitions
 - § preferential treatment



- § Small authority weights should not contribute to the computation of the hub weights
- § Repeat until convergence
 - § O operation : hubs collect the k highest authority weights

$$h_i = \sum_{j:i \to j} a_j : a_j \in F_k(i)$$

§ *I* operation: authorities collect the weight of hubs

$$a_i = \sum_{j: j \to i} h_j$$

§ Normalize weights under some norm



Norm(p) algorithm

- § Small authority weights should contribute less to the computation of the hub weights
- § Repeat until convergence
 - § O operation : hubs compute the p-norm of the authority weight vector

$$h_{i} = \left(\sum_{j:i \to j} a_{j}^{p}\right)^{1/p} = \left\|\overline{F(i)}\right\|_{p}$$

§ *I* operation: authorities collect the weight of hubs

$$a_i = \sum_{j: j \to i} h_j$$

§ Normalize weights under some norm



- § A hub is as good as the best authority it points to
- § Repeat until convergence § *O* operation : hubs collect the highest authority weight $h_i = \max_{j:i \to j} a_j$ § *I* operation: authorities collect the weight of hubs $a_i = \sum_{j:j \to i} h_j$ § Normalize weights under some norm
- § Special case of AT(k) (for k=1) and Norm(p) (p= ∞)



§ Discrete Dynamical System: The repeated application of a function g on a set of weights

Initialize weights to w^{0} For t=1,2,... $w^{t}=g(w^{t-1})$

- § LAR algorithms: the function g propagates the weight on the graph G
- § Linear vs Non-Linear dynamical systems
 - § eigenvector analysis algorithms (PageRank, HITS) are linear dynamical systems
 - § AT(k), Norm(p) and MAX are non-linear



- § Notoriously hard to analyze not well understood
 - § we cannot easily prove convergence
 - § we do not know much about stationary weights
- § Convergence is important for an LAR algorithm to be well defined.
- § The MAX algorithm converges for any initial configuration

















after normalization with the max weight





The hubs are mapped to the seed node

before normalization w=3after normalization with the max weight w=1

normalization factor = 3





weight of blue node





weight of yellow node W = (1 + W)/3











- § ...in the beginning...
- § previous work
- § some more algorithms
- § some experimental data [BRRT05]
- § a theoretical framework



- § 34 different queries
- § user relevance feedback
 - § high relevant/relevant/non-relevant
- § measures of interest
 - § "high relevance ratio"
 - § "relevance ratio"
- § Data (and code?) available at

http://www.cs.toronto.edu/~tsap/experiments/journal (or /thesis)


Aggregate Statistics

	AVG HR	STDEV HR	AVG R	STDEV R
HITS	22%	24%	45%	39%
PageRank	24%	14%	46%	20%
In-Degree	35%	22%	58%	29%
SALSA	35%	21%	59%	28%
MAX	38%	25%	64%	32%
BFS	43%	18%	73%	19%



Aggregate Statistics

	AVG HR	STDEV HR	AVG R	STDEV R
HITS	22%	24%	45%	39%
PageRank	24%	14%	46%	20%
In-Degree	35%	22%	58%	29%
SALSA	35%	21%	59%	28%
MAX	38%	25%	64%	32%
BFS	43%	18%	73%	19%



Aggregate Statistics

	AVG HR	STDEV HR	AVG R	STDEV R
HITS	22%	24%	45%	39%
PageRank	24%	14%	46%	20%
In-Degree	35%	22%	58%	29%
SALSA	35%	21%	59%	28%
MAX	38%	25%	64%	32%
BFS	43%	18%	73%	19%



HITS and the TKC effect

§



"recipes"

- 1. (1.000) HonoluluAdvertiser.com URL: http://www.hawaiisclassifieds.com
- § 2. (0.999) Gannett Company, Inc. URL: http://www.gannett.com
- § 3. (0.998) AP MoneyWire URL: http://apmoneywire.mm.ap.org
- ß 4. (0.990) e.thePeople : Honolulu Advertiser URL: http://www.e-thepeople.com/
- 5. (0.989) News From The Associated Press Ş URL: http://customwire.ap.org/
- § 6. (0.987) Honolulu Traffic URL: http://www.co.honolulu.hi.us/
- § 7. (0.987) News From The Associated Press URL: http://customwire.ap.org/
- 8. (0.987) <u>News From The Associated Press</u> URL: http://customwire.ap.org/ Ş
 - 9. (0.987) News From The Associated Press URL: http://customwire.ap.org/
 - **10.** (0.987) News From The Associated Press URL: http://customwire.ap.org/



MAX - "net censorship"

- § 1. (1.000) EFF: Homepage URL: http://www.eff.org
- § 2. (0.541) Internet Free Expression Alliance URL: http://www.ifea.net
- § **3.** (0.517) <u>The Center for Democracy and Technology</u> URL: http://www.cdt.org
- § 4. (0.517) <u>American Civil Liberties Union</u> URL: http://www.aclu.org
- § **5.** (0.386) <u>Vtw Directory Page</u> URL: http://www.vtw.org
- § 6. (0.357) <u>PEACEFIRE</u> URL: http://www.peacefire.org
- § **7.** (0.277) <u>Global Internet Liberty Campaign Home Page</u> URL: http://www.gilc.org
- § 8. (0.254) <u>libertus.net: about censorship and free speech</u> URL: http://libertus.net
- § 9. (0.196) EFF Blue Ribbon Campaign Home Page URL: http://www.eff.org/blueribbon.html
- § **10.** (0.144) <u>The Freedom Forum</u> URL: http://www.freedomforum.org



MAX – "affirmative action"

- § **1.** (1.000) Copyright Information URL: http://www.psu.edu/copyright.html
- § **2.** (0.447) **PSU Affirmative Action** URL: http://www.psu.edu/dept/aaoffice
- § **3.** (0.314) <u>Welcome to Penn State's Home on the Web</u> URL: http://www.psu.edu
- § **4.** (0.010) <u>University of Illinois</u> URL: http://www.uiuc.edu
- § **5.** (0.009) <u>Purdue University-West Lafayette, Indiana</u> URL: http://www.purdue.edu
- § 6. (0.008) <u>UC Berkeley home page</u> URL: http://www.berkeley.edu
- § 7. (0.008) <u>University of Michigan</u> URL: http://www.umich.edu
- § **8.** (0.008) <u>The University of Arizona</u> URL: http://www.arizona.edu
- § 9. (0.008) <u>The University of Iowa Homepage</u> URL: http://www.uiowa.edu
- § **10.** (0.008) <u>Penn: University of Pennsylvania</u> URL: http://www.upenn.edu



PageRank

- § **1.** (1.000) WCLA Feedback URL: http://www.janeylee.com/wcla
- § **2.** (0.911) **Planned Parenthood Action Network** URL: http://www.ppaction.org/ppaction/
- § **3.** (0.837) <u>Westchester Coalition for Legal Abortion</u> URL: http://www.wcla.org
- § **4.** (0.714) <u>Planned Parenthood Federation</u> URL: http://www.plannedparenthood.org
- § **5.** (0.633) <u>GeneTree.com Page Not Found</u> URL: http://www.qksrv.net/click
- § 6. (0.630) <u>Bible.com Prayer Room</u> URL: http://www.bibleprayerroom.com
- § **7.** (0.609) <u>United States Department of Health</u> URL: http://www.dhhs.gov

8. (0.538) <u>Pregnancy Centers Online</u> URL: http://www.pregnancycenters.org

- § 9. (0.517) <u>Bible.com Online World</u> URL: http://bible.com
- § **10.** (0.516) <u>National Organization for Women</u> URL: http://www.now.org



link-spam structure



- § ...in the beginning...
- § previous work
- § some more algorithms
- § some experimental data
- § a theoretical framework



Theoretical Analysis of LAR algorithms [BRRT05]

- § Why bother?
 - § Plethora of LAR algorithms: we need a formal way to compare and analyze them
 - § Need to define properties that are useful
 - sensitivity to spam
 - § Need to discover the properties that characterize each LAR algorithm



- § A Link Analysis Ranking Algorithm is a function that maps a graph to a real vector $A:G_n \to \mathbb{R}^n$
- § G_n : class of graphs of size n
- § LAR vector the output A(G) of an algorithm A on a graph G
- § G_n: the class of all possible graphs of size
 n



$W_1 = \begin{bmatrix} 0.9 & 1 & 0.7 & 0.6 & 0.8 \end{bmatrix}$

§ How close are the LAR vectors w_1 , w_2 ?



§ Geometric distance: how close are the numerical weights of vectors w₁, w₂?

$$d_{1}(w_{1}, w_{2}) = \sum |w_{1}[i] - w_{2}[i]|$$

$$w_{1} = [1.0 \ 0.8 \ 0.5 \ 0.3 \ 0.0]$$

$$w_{2} = [0.9 \ 1.0 \ 0.7 \ 0.6 \ 0.8]$$

$$d_{1}(w_{1}, w_{2}) = 0.1 + 0.2 + 0.2 + 0.3 + 0.8 = 1.6$$



§ Rank distance: how close are the ordinal rankings induced by the vectors w₁, w₂?
§ Kendal's T distance

 $d_r(w_1, w_2) = \frac{\text{pairs ranked in a different order}}{\text{total number of distinct pairs}}$









Rank distance of partial rankings

- § Extreme value p = 1
 - § charge for every potential conflict
- § Extreme value p = 0
 - § charge only for inconsistencies
 - § problem: not a metric
- § Intermediate values 0 < p < 1
 - § Details [FMNKS04] [T04]
 - § Interesting case p = 1/2
- § We will use whatever gives a stronger result



- § Intuition: a small change on a graph should cause a small change on the output of the algorithm.
- § Definition: Link distance between graphs G=(P,E) and G'=(P,E')

 $d_{g}(G,G') = |E \cup E'| - |E \cap E'|$





§ $C_k(G)$: set of graphs G' such that $d_{\ell}(G,G') \leq k$

§ Definition: Algorithm A is stable if

 $\lim_{n\to\infty}\max_{G}\max_{G'\in C_k(G)}d_1(A(G), A(G'))=0$

§ Definition: Algorithm A is rank stable if $\lim_{n\to\infty} \max_{G} \max_{G'\in C_k(G)} d_r(A(G), A(G')) = 0$



Stability: Results

- § InDegree algorithm is stable and rank stable on the class G_n
- § HITS, Max are neither stable nor rank stable on the class G_n



Instability of HITS





Stability of HITS

§ HITS is stable if $\sigma_1 - \sigma_2 \rightarrow \infty$ [NZJ01] § The two strongest linear trends are well

- separated
- § What about the converse?



§ PageRank is unstable



§ PageRank is rank unstable [Lempel Moran 2005]



§ Perturbations to unimportant nodes have small effect on the PageRank values [NZJ01][BGS03]

$$d_1(A(G), A(G')) \le \frac{2a}{1-2a} \sum_{i \in P} A(G)[i]$$



Stability of PageRank

§ Lee Borodin model [LB03]

- § upper bounds depend on authority and hub values
- § PageRank, Randomized SALSA are stable
- § HITS, SALSA are unstable
- § Open question: Can we derive conditions for the stability of PageRank in the general case?



- § Definition: Two algorithms A_1 , A_2 are similar if $\lim_{n \to \infty} \frac{\max_{G \in G_n} d_1(A_1(G), A_2(G))}{\max_{W_1, W_2} d_1(W_1, W_2)} = 0$
- § Definition: Two algorithms A₁, A₂ are rank similar if

$$\lim_{n\to\infty}\max_{G\in G_n}d_r(A_1(G), A_2(G))=0$$

§ Definition: Two algorithms A₁, A₂ are rank equivalent if

$$\max_{G \in G_n} d_r(A_1(G), A_2(G)) = 0$$



Similarity: Results

§ No pairwise combination of InDegree, SALSA, HITS and MAX algorithms is similar, or rank similar on the class of all possible graphs G_n



Product Graphs

- § Latent authority and hub vectors a, h
 - $h_i = probability of node i being a good hub$
 - $\frac{1}{3}$ = probability of node j being a good authority
- § Generate a link $i \rightarrow j$ with probability $h_i a_j$ $W[i, j] = \begin{cases} 1 & \text{with probability } h_i a_j \\ 0 & \text{with probability } 1 - h_i a_j \end{cases}$

§ Azar, Fiat, Karlin, McSherry Saia 2001
 § The class of product graphs G^p_n



§ Theorem: HITS and InDegree are similar with high probability on the class of product graphs, G^p_n (subject to some assumptions)



§ Monotonicity: Algorithm A is strictly monotone if for any nodes x and y

 $B_N(x) \subset B_N(y) \Leftrightarrow A(G)[x] < A(G)[y]$





§ Locality: An algorithm A is strictly rank local if, for every pair of graphs G=(P,E) and G'=(P,E'), and for every pair of nodes x and y, if B_G(x)=B_{G'}(x) and B_G(y)=B_{G'}(y) then

$A(G)[x] < A(G)[y] \Leftrightarrow A(G')[x] < A(G')[y]$

§ the relative order of the nodes remains the same



§ The InDegree algorithm is strictly rank local



- § Label Independence: An algorithm is label independent if a permutation of the labels of the nodes yields the same permutation of the weights
 - § the weights assigned by the algorithm do not depend on the labels of the nodes



Axiomatic characterization of the InDegree algorithm [BRRT05]

§ Theorem: Any algorithm that is strictly rank local, strictly monotone and label independent is rank equivalent to the InDegree algorithm



§ Consider two nodes i and j with d(i) > d(j)
§ Assume that w(i) < w(j)</p>





§ Remove all links except to i and j § w₁(i) < w₁(j) (from locality)





§ Add links from C and L to node k § w₂(i) < w₂(j) (from locality) § w₂(k) < w₂(i) (from monotonicity)

§ $w_2(k) < w_2(j)$




§ Remove links from L to i and add links from R to i

 $w_3(k) < w_3(j)$ (from locality)





§ Graphs G₂ and G₃ are the same up to a label permutation







§ Graphs G₂ and G₃ are the same up to a label permutation







- § We now have
 - § $w_2(j) < w_2(k)$ and $w_3(j) < w_3(k)$ (shown before) § $w_2(j) = w_3(k)$ and $w_2(k) = w_3(j)$ (label independ.) § $w_2(j) > w_2(k)$ CONTRADICTION!





§ All three properties are needed

- § locality
 - PageRank is also strictly monotone and label independent
- § monotonicity
 - consider an algorithm that assigns 1 to nodes with even degree, and 0 to nodes with odd degree
- § label independence
 - consider and algorithm that gives the more weight to links that come from some specific page (e.g. the Yahoo page)



References

- § [BP98] S. Brin, L. Page, The anatomy of a large scale search engine, WWW 1998
- § [K98] J. Kleinberg. Authoritative sources in a hyperlinked environment. Proc. 9th ACM-SIAM Symposium on Discrete Algorithms, 1998.
- § [HB98] Monika R. Henzinger and Krishna Bharat. Improved algorithms for topic distillation in a hyperlinked environment. Proceedings of the 21'st International ACM SIGIR Conference on Research and Development in IR, August 1998.
- § [BRRT05] A. Borodin, G. Roberts, J. Rosenthal, P. Tsaparas, Link Analysis Ranking: Algorithms, Theory and Experiments, ACM Transactions on Internet Technologies (TOIT), 5(1), 2005
- § R. Lempel, S. Moran. The Stochastic Approach for Link-Structure Analysis (SALSA) and the TKC Effect. 9th International World Wide Web Conference, May 2000.
- § A. Y. Ng, A. X. Zheng, and M. I. Jordan. Link analysis, eigenvectors, and stability. International Joint Conference on Artificial Intelligence (IJCAI), 2001.
- § Ronny Lempel, Shlomo Moran: Rank-Stability and Rank-Similarity of Link-Based Web Ranking Algorithms in Authority-Connected Graphs. Inf. Retr. 8(2): 245-264 (2005)
- § P. Tsaparas, Using Non-Linear Dynamical Systems for Web Searching and Ranking Principles of Database Systems (PODS), Paris, 2004
- § Azar, Fiat, Karlin, McSherry, and Saia, Spectral Analysis of Data, STOC, 2001
- § [FKMSV04] Ron Fagin, Ravi Kumar, Mohammad Mahdian, D. Sivakumar, Erik Vee, Comparing and aggregating rankings with ties, PODS 2004