Link Analysis and Web Search

Part III: Dynamics — Episode 10

Baochun Li Department of Electrical and Computer Engineering University of Toronto

Required reading: "Networks, Crowds, and Markets," Chapter 13.1 — 13.4, 14.1 — 14.3, 14.6

Information networks and the Web

- Logical relationships among pieces of information
 - Best example: the Web
- 1991: Tim Berners-Lee at CERN (Switzerland) created the Web
 - provided an easy way to make documents web pages for the world to see
 - view these pages using browsers
 - It is based on the idea of connecting these pages using links



The idea of links is both inspired and non-obvious

There are many ways of organizing information: classification (library), series of folders (files), or just alphabetically (phone book)





Modeling the web as a directed graph

Objective: create a "map" of the web But how?







Strongly Connected Components

Strongly connected component: a subset of nodes such that (1) every node in the subset has a path to every other; and (2) the subset is not part of some larger set in which every node can reach every other.





Now we can build a global map of the Web, using strongly connected components (Broder *et al.* [1999])

WIKIPEDIA: CORNELL UNIVERSITY WIKIPEDIA: CORNELL UNIVERSITY WIKIPEDIA: CORNELL WIKIPEDIA:

A giant strongly connected component



Can there be a second giant strongly connected component?







Not really — it's too fragile





11

The bow-tie structure of the web SCC Upstream **Downstream Disconnected Components**







How do we find web pages using search?

- (search)
 - librarians
 - patent attorneys
- were searching for were written by professionals
 - research articles
 - court documents
 - U.S. patents

Up through the 1980s, very few people cared about information retrieval

They are trained to formulate effective queries, and the documents they

1	$ \Delta $	
		Г

The Web is entirely different

- Both search users and web page authors are amateurs
- Scale is really large
- Highly dynamic nature of the content to be searched
 - Some of the authors may even optimize their content for a search engine
 - An industry called "Search Engine Optimization"
 - Millions of dollars on the line

1	5)
	\sim	ſ

When I search for a key phrase, what do I need?



- University of Wisconsin
 - University of Windsor
 - University of Winnipeg
 - University of Wyoming



Q



Basic idea: let the links "vote"





Using links as more than simple "votes"





Searching for "good museum"





Link-based ranking with hubs and authorities

Key idea: voting again and again





Link-based ranking with hubs and authorities

Key idea: principle of repeated improvement





But why stop here?





Let's make it more formal

- Two kinds of quality measures for web pages
 - Authority score Auth(p): level of endorsement
 - Hub score Hub(p) : quality as a list

scores of all pages that link to p.

scores of all pages that p points to.

Divide all scores so that they add to 1.

- Authority update rule: Auth(p) = sum of hub
- Hub update rule: Hub(p) = sum of authority



Using adjacency matrices to represent a graph



view a set of n pages, 1, 2, ...n, as a set of nodes in a directed graph

• n x n matrix M: M_{ij} is equal to 1 if thee is a link from node i to node j





Hub update rule as matrix multiplication $h_i \leftarrow M_{i1}a_1 + M_{i2}a_2 + \cdots + M_{in}a_n$ $h \leftarrow Ma$ $\begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 6 \\ 4 \\ 2 \\ 4 \end{bmatrix} = \begin{bmatrix} 9 \\ 7 \\ 2 \\ 4 \end{bmatrix}$ node 1 node 2 node 3 node 4





The authority update rule as matrix multiplication

 $a \leftarrow M^T h$

 $a_i \leftarrow M_{1i}h_1 + M_{2i}h_2 + \cdots + M_{ni}h_n$



Unwinding the k-step hub-authority updates $a^{\langle 1 \rangle} = M^T h^{\langle 0 \rangle}$

 $h^{\langle 1 \rangle} = M a^{\langle 1 \rangle} = M M^T h^{\langle 0 \rangle}$

 $a^{\langle 2 \rangle} = M^T h^{\langle 1 \rangle} = M^T M M^T h^{\langle 0 \rangle}$

 $a^{\langle k \rangle} = (M^{\prime})$ $h^{\langle k \rangle}$

 $h^{\langle 2 \rangle} = M a^{\langle 2 \rangle} = M M^T M M^T h^{\langle 0 \rangle} = (M M^T)^2 h^{\langle 0 \rangle}$

$$(MM^{T})^{k-1}M^{T}h^{\langle 0 \rangle}$$



PageRank: the core of Google search

Indirect endorsement



Direct endorsement





Basic PageRank update rule









Basic PageRank update rule

- Assign each page p a PageRank value
 - start with 1/n, n is the number of pages
- Basic PageRank update rule
 - them across outbound links
 - And then use the principle of repeated improvement
- Fact: If the network is strongly connected, then there are unique equilibrium PageRank values

Each node divides its current PageRank into equal shares, and then pass







Basic PageRank updates as matrix multiplication

- PageRank that should be passed to j in one single step
 - start with $1/L_i$, where L_i is the number of links out of i



• Each entry in the adjacency matrix N_{ij} specifies the portion of i's





Basic PageRank updates as matrix multiplication

- Vector r: the PageRanks of all the nodes

 $r_i \leftarrow N_{1i}r_1 + N_{2i}r_2 + \cdots + N_{ni}r_n$.

 $r \leftarrow N^T r$.



One major problem with the basic update rule





One major problem with the basic update rule







Scaled PageRank update rule

 Fact: Repeatedly applying the scaled PageRank update rule converges to a unique equilibrium set of PageRank values for all networks.





 $r \leftarrow N$

$$sN_{ij} + (1 - s)/n$$

+ $\tilde{N}_{2i}r_2 + \dots + \tilde{N}_{ni}r_n$
 $\tilde{N}^T r. \quad r^{\langle k \rangle} = (\tilde{N}^T)^k r^{\langle 0 \rangle}.$



Required reading: "Networks, Crowds, and Markets," Chapter 13.1 — 13.4, 14.1 — 14.3, 14.6

Final Examination

- December 15, 2023, BA 1180, 1:10pm 3pm, 110 minutes
- Covers all the lectures, but not the critique papers
- Sample examination questions on the course website
 - Length: a bit longer than the sample exams
 - Format: the same a large question divided into 2-3 small ones
- Special Office Hour: December 11, Monday, 2-5pm, BA 4118



Preparing for the final examination

- Strongly recommended: read the corresponding chapters in the textbook (in the course website)
 - Exercises at the end of each chapter will be very helpful
- If this is not feasible, understand the examples in the lecture slides
- Video recordings may help as well (though two were not available)

42

Sneak preview: my new graduate course next Fall

- Tentatively titled "Performant Software Systems"
- Covers how modern and high-performance software systems are built using the Rust programming language
- Still pending review and approval



Course Evaluations