

Extensions: Topic Specific (aka Personalized) PageRank

Mining of Massive Datasets
Leskovec, Rajaraman, and Ullman
Stanford University



Topic-Specific PageRank

- Instead of generic popularity, can we measure popularity within a topic?
- **Goal:** Evaluate Web pages not just according to their popularity, but by how close they are to a particular topic, e.g. “sports” or “history”
- **Allows search queries to be answered based on interests of the user**
 - **Example:** Query “Trojan” wants different pages depending on whether you are interested in sports, history and computer security

Topic-Specific PageRank

- Random walker has a small probability of teleporting at any step
- **Teleport can go to:**
 - **Standard PageRank:** Any page with equal probability
 - To avoid dead-end and spider-trap problems
 - **Topic Specific PageRank:** A topic-specific set of “relevant” pages (**teleport set**)
- **Idea: Bias the random walk**
 - When walker teleports, she pick a page from a set S
 - S contains only pages that are relevant to the topic
 - E.g., Open Directory (DMOZ) pages for a given topic/query
 - For each teleport set S , we get a different vector r_S

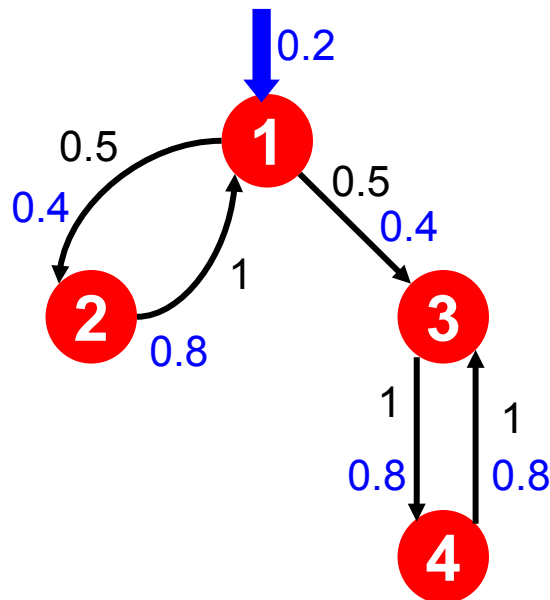
Matrix Formulation

- To make this work all we need is to update the teleportation part of the PageRank formulation:

$$A_{ij} = \begin{cases} \beta M_{ij} + (1 - \beta)/|S| & \text{if } i \in S \\ \beta M_{ij} & \text{otherwise} \end{cases}$$

- A is stochastic!
- We weighted all pages in the teleport set S equally
 - Could also assign different weights to pages!
- **Random Walk with Restart:** S is a single element
- **Compute as for regular PageRank:**
 - Multiply by M , then add a vector
 - Maintains sparseness

Example: Topic-Specific PageRank



Suppose $S = \{1\}$, $\beta = 0.8$

Node	Iteration			
	0	1	2	stable
1	0.25	0.4	0.28	0.294
2	0.25	0.1	0.16	0.118
3	0.25	0.3	0.32	0.327
4	0.25	0.2	0.24	0.261

$S=\{1\}$, $\beta=0.90$:

$r=[0.17, 0.07, 0.40, 0.36]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

$S=\{1\}$, $\beta=0.70$:

$r=[0.39, 0.14, 0.27, 0.19]$

$S=\{1,2,3,4\}$, $\beta=0.8$:

$r=[0.13, 0.10, 0.39, 0.36]$

$S=\{1,2,3\}$, $\beta=0.8$:

$r=[0.17, 0.13, 0.38, 0.30]$

$S=\{1,2\}$, $\beta=0.8$:

$r=[0.26, 0.20, 0.29, 0.23]$

$S=\{1\}$, $\beta=0.8$:

$r=[0.29, 0.11, 0.32, 0.26]$

Discovering the Topic Vector S

- **Create different PageRanks for different topics**
 - The 16 DMOZ top-level categories:
 - arts, business, sports,...
- **Which topic ranking to use?**
 - User can pick from a menu
 - Classify query into a topic
 - Can use the **context** of the query
 - E.g., query is launched from a web page talking about a known topic
 - History of queries e.g., “basketball” followed by “Jordan”
 - User context, e.g., user’s bookmarks, ...