



Introduction to Data Mining

Link Analysis-2

U Kang
Seoul National University



In This Lecture

- Pagerank: Google formulation
 - Make the solution to converge
- Computing Pagerank for very large graphs
 - Pagerank vector and/or stochastic matrix do not fit in the memory
- Topic specific Pagerank



Outline

- ➔ **PageRank: Google Formulation**
- Computing PageRank
- Topic-Specific PageRank
- Measuring Proximity in Graphs



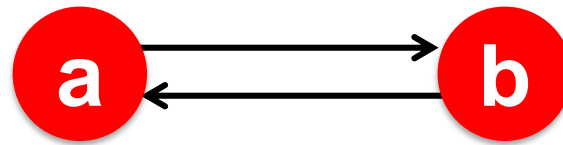
PageRank: Three Questions

$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i} \quad \text{or equivalently} \quad r = Mr$$

- Does this converge?
- Does it converge to what we want?
- Are results reasonable?



Does this converge?



$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

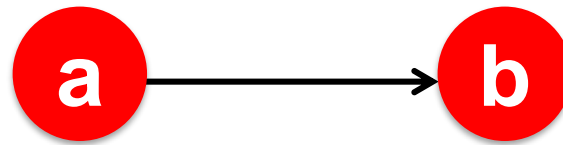
■ Example:

$$\begin{array}{l} r_a \\ r_b \end{array} = \begin{array}{cccc} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{array}$$

Iteration 0, 1, 2, ...



Does it converge to what we want?



$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

■ Example:

$$\begin{array}{l} \mathbf{r}_a \\ \mathbf{r}_b \end{array} = \begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{array}$$

Iteration 0, 1, 2, ...

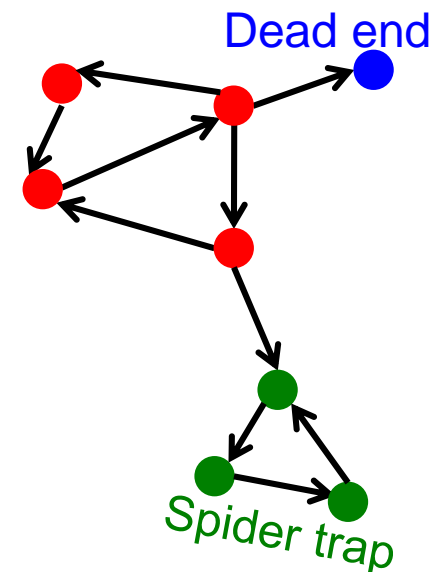


PageRank: Problems

2 problems:

- **(1)** Some pages are **dead ends** (have no out-links)
 - Random walk has “nowhere” to go to
 - Such pages cause importance to “leak out”

- **(2) Spider traps:**
(all out-links are within the group)
 - Random walked gets “stuck” in a trap
 - And eventually spider traps absorb all importance

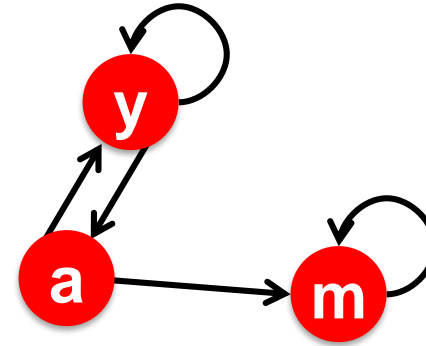




Problem: Spider Traps

■ Power Iteration:

- Set $r_j = 1$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
 - And iterate



m is a spider trap

	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	1

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

■ Example:

Iteration 0, 1, 2, ...

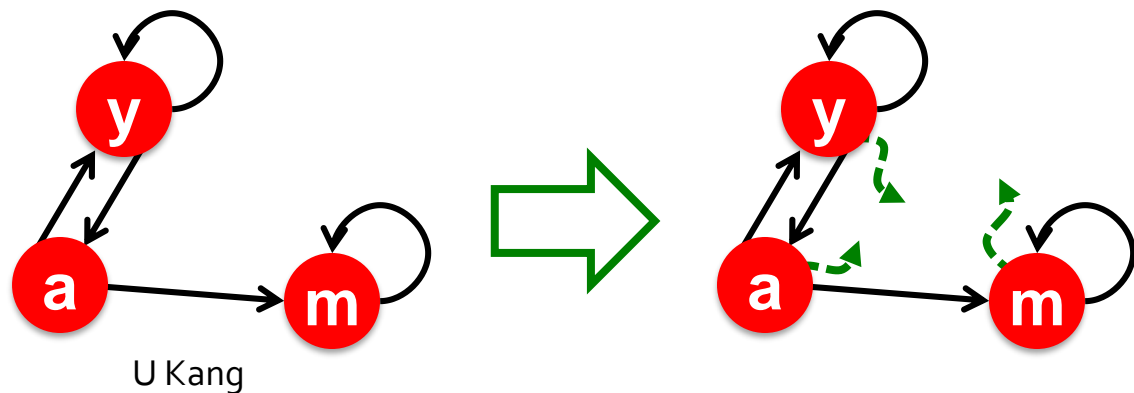
$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} =$	1/3	2/6	3/12	5/24		0
	1/3	1/6	2/12	3/24	...	0
	1/3	3/6	7/12	16/24		1

All the PageRank score gets “trapped” in node m.



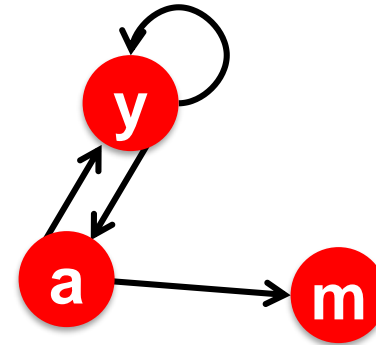
Solution: Teleports!

- **The Google solution for spider traps: At each time step, the random surfer has two options**
 - With prob. β , follow a link at random
 - With prob. $1-\beta$, jump to some random page
 - Common values for β are in the range 0.8 to 0.9
- **Surfer will teleport out of spider trap within a few time steps**





Problem: Dead Ends



	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	0

■ Power Iteration:

- Set $r_j = 1$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
 - And iterate

$$r_y = r_y/2 + r_a/2$$

$$r_a = r_y/2$$

$$r_m = r_a/2$$

■ Example:

Iteration 0, 1, 2, ...

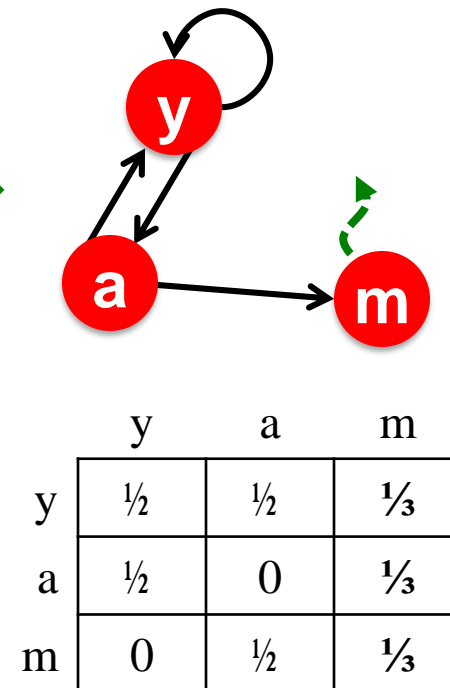
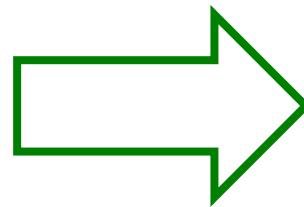
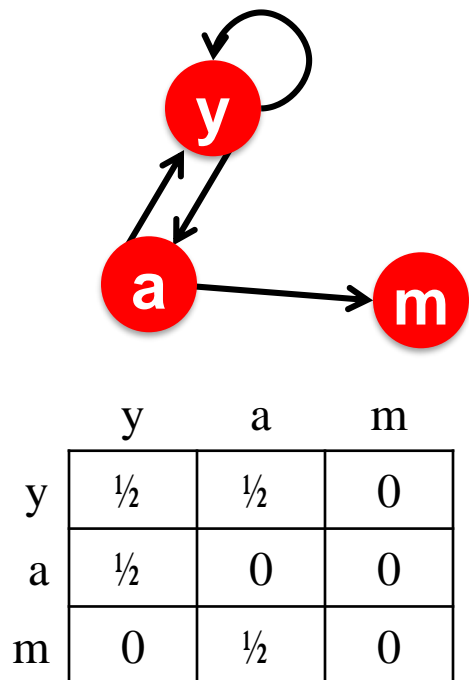
$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{matrix}$$

Here the PageRank “leaks” out since the matrix is not column stochastic.



Solution: Always Teleport!

- **Teleports:** Follow random teleport links with probability 1.0 from dead-ends
 - Adjust matrix accordingly





Why Teleports Solve the Problem?

Why are dead-ends and spider traps a problem and **why do teleports solve the problem?**

- **Spider-traps** are not a problem, but with traps PageRank scores are **not** what we want
 - **Solution:** Never get stuck in a spider trap by teleporting out of it in a finite number of steps
- **Dead-ends** are a problem
 - The matrix is not column stochastic so our initial assumptions are not met
 - **Solution:** Make matrix column stochastic by always teleporting when there is nowhere else to go



Solution: Random Teleports

■ Google's solution:

At each step, random surfer has two options:

- With probability β , follow a link at random
- With probability $1-\beta$, jump to some random page

■ PageRank equation [Brin-Page, 98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{N}$$

d_i ... out-degree
of node i

This formulation assumes that M has no dead ends. We can either preprocess matrix M to remove all dead ends or explicitly follow random teleport links with probability 1.0 from dead-ends.



The Google Matrix

- **PageRank equation** [Brin-Page, '98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{N}$$

- **In matrix form:**

$$\begin{aligned} \square r &= \beta M r + (1 - \beta) \left[\frac{1}{N} \right]_{N \times N} r \\ &= \underbrace{\{ \beta M + (1 - \beta) \left[\frac{1}{N} \right]_{N \times N} \}} r \end{aligned}$$

This is called the “Google Matrix”

$\left[\frac{1}{N} \right]_{N \times N}$...N by N matrix
where all entries are 1/N

1/3	1/3	1/3
1/3	1/3	1/3
1/3	1/3	1/3

E.g., for N=3



The Google Matrix

- **PageRank equation** [Brin-Page, '98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{N}$$

- **The Google Matrix A:**

$$A = \beta M + (1 - \beta) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \frac{1}{N}$$

[1/N]_{N×N}...N by N matrix
where all entries are 1/N

- **We have a recursive problem: $\mathbf{r} = \mathbf{A} \cdot \mathbf{r}$**

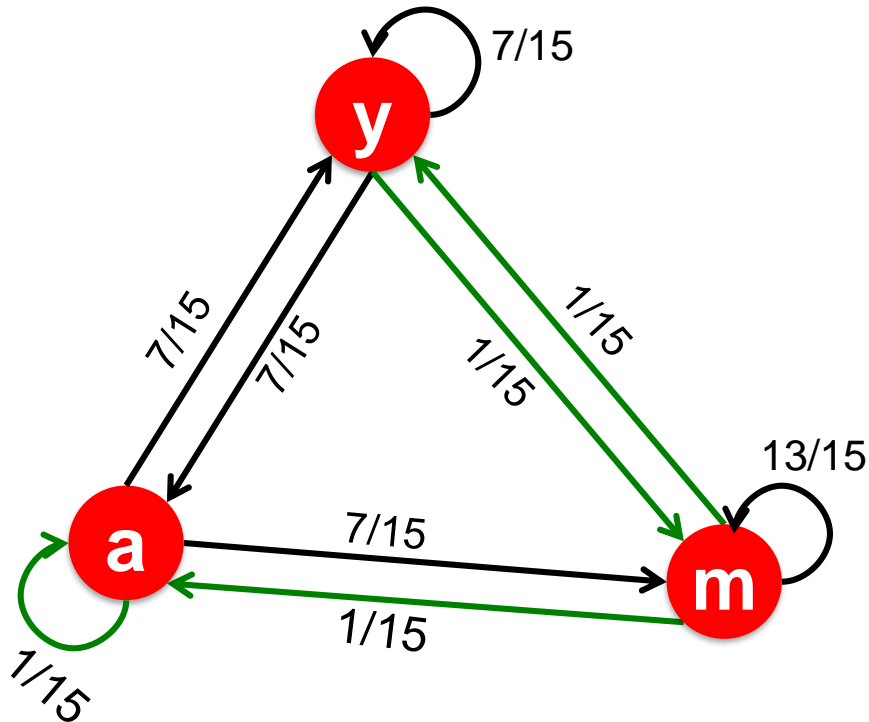
And the Power method still works!

- **What is β ?**

□ In practice $\beta = 0.8, 0.9$ (make ~ 5 steps on avg., jump)



Random Teleports ($\beta = 0.8$)



$$0.8 \begin{matrix} \mathbf{M} \\ \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} \end{matrix} + 0.2 \begin{matrix} [1/N]_{N \times N} \\ \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix} \end{matrix}$$

$$\begin{matrix} \mathbf{A} \\ \begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix} \end{matrix}$$

y	=	1/3	0.33	0.24	0.26	7/33
a		1/3	0.20	0.20	0.18	5/33
m		1/3	0.46	0.52	0.56	21/33



Outline

PageRank: Google Formulation

 **Computing PageRank**

Topic-Specific PageRank

Measuring Proximity in Graphs



Computing Page Rank

■ Key step is matrix-vector multiplication

□ $r^{new} = A \cdot r^{old}$

■ Easy if we have enough main memory to hold A , r^{old} , r^{new}

■ Say $N = 1$ billion pages

$$A = \beta \cdot M + (1-\beta) [1/N]_{N \times N}$$

□ We need 4 bytes for each entry (say)

□ Total 2 billion entries for 2 vectors (r^{old} , r^{new}): ~ 8 GB

$$A = 0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$$

□ Matrix A has N^2 entries

■ $N^2 = 10^{18}$ (1000 Peta) is a large number!

■ We need to exploit sparsity of M

$$= \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}$$



Sparse Matrix Formulation

- $\mathbf{r} = \mathbf{A} \cdot \mathbf{r}$, where $\mathbf{A} = \beta \mathbf{M} + (1 - \beta) \left[\frac{1}{N} \right]_{N \times N}$

- Main idea: do not construct \mathbf{A} explicitly

- Specifically:

- $\mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + (1 - \beta) \left[\frac{1}{N} \right]_{N \times N} \mathbf{r}$

- $= \beta \mathbf{M} \cdot \mathbf{r} + (1 - \beta) \left[\frac{1}{N} \right]_N$

Note: Here we assume \mathbf{M} has no dead-ends.

$[x]_N$... a vector of length N with all entries x



Sparse Matrix Formulation

- The **PageRank equation**

$$\mathbf{r} = \beta \mathbf{M} \cdot \mathbf{r} + \left[\frac{\mathbf{1} - \beta}{N} \right]_N$$

- where $[(\mathbf{1}-\beta)/N]_N$ is a vector with all N entries $(\mathbf{1}-\beta)/N$

- \mathbf{M} is a **sparse matrix!**

- 10 links per node, approx $10N$ entries

- So in each iteration, we need to:

- Compute $\mathbf{r}^{\text{new}} = \beta \mathbf{M} \cdot \mathbf{r}^{\text{old}}$

- Add a constant value $(\mathbf{1}-\beta)/N$ to each entry in \mathbf{r}^{new}

- **Note if \mathbf{M} contains dead-ends then $\sum_j r_j^{\text{new}} < 1$ and we also have to renormalize \mathbf{r}^{new} so that it sums to 1**



PageRank: The Complete Algorithm

■ Input: Graph G and parameter β

- Directed graph G (can have **spider traps** and **dead ends**)
- Parameter β

■ Output: PageRank vector r^{new}

- **Set:** $r_j^{old} = \frac{1}{N}$

- **repeat until convergence:** $\sum_j |r_j^{new} - r_j^{old}| < \varepsilon$

- $\forall j: r_j^{\prime new} = \sum_{i \rightarrow j} \beta \frac{r_i^{old}}{d_i}$

- **Now re-insert the leaked PageRank:**

- $\forall j: r_j^{new} = r_j^{\prime new} + \frac{1-S}{N}$ **where:** $S = \sum_j r_j^{\prime new}$

- $r^{old} = r^{new}$

If the graph has no dead-ends then the amount of leaked PageRank is $1-\beta$. But since we have dead-ends the amount of leaked PageRank may be larger. We have to explicitly account for it by computing S .



Sparse Matrix Encoding

- **Encode sparse matrix using only nonzero entries**
 - Space proportional roughly to number of links
 - Assuming $N = 1$ billion,
10N edges would require $4 * 10 * 1$ billion = 40GB
 - **Still won't fit in memory, but will fit on disk**

source node	degree	destination nodes
0	3	1, 5, 7
1	5	17, 64, 113, 117, 245
2	2	13, 23



Basic Algorithm: Update Step

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix \mathbf{M} on disk
- 1 step of power-iteration is:

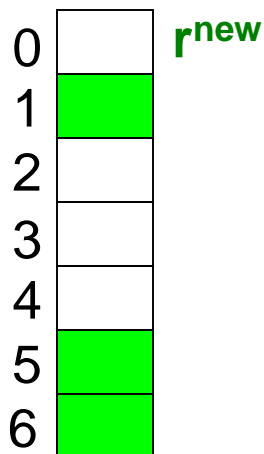
Initialize all entries of $r^{new} = (1-\beta) / N$

For each page i (of out-degree d_i):

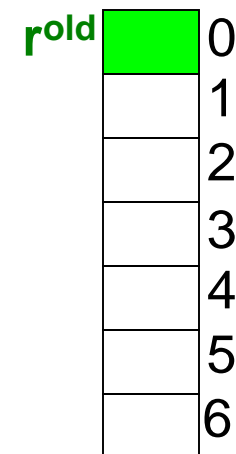
Read into memory: $i, d_i, dest_1, \dots, dest_{d_i}, r^{old}(i)$

For $j = 1 \dots d_i$

$$r^{new}(dest_j) += \beta r^{old}(i) / d_i$$

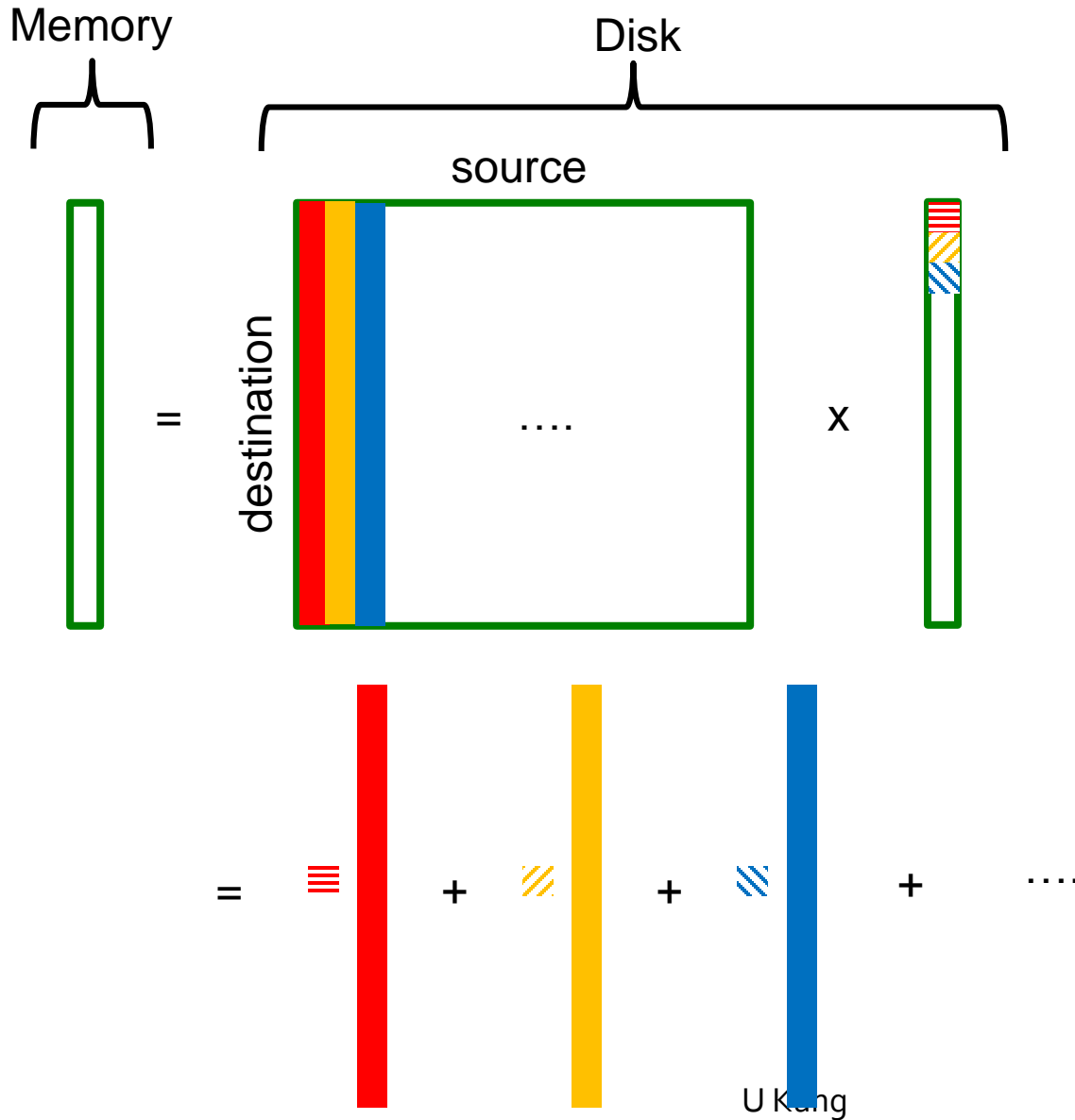


source	degree	destination
0	3	1, 5, 6
1	4	17, 64, 113, 117
2	2	13, 23





Basic Algorithm: Update Step



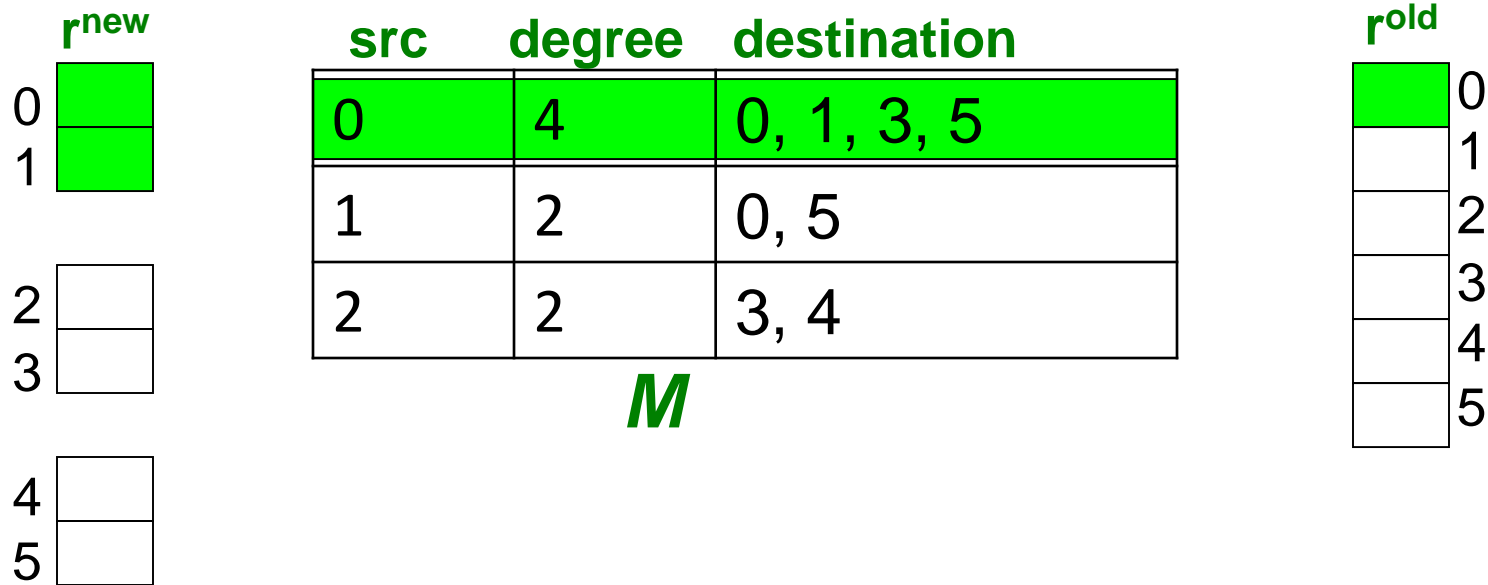


Analysis

- Assume enough RAM to fit r^{new} into memory
 - Store r^{old} and matrix M on disk
- In each iteration, we have to:
 - Read r^{old} and M
 - Write r^{new} back to disk
 - Cost (disk I/O) per iteration of Power method:
 $= 2|r| + |M|$
- Question:
 - What if we could not even fit r^{new} in memory?



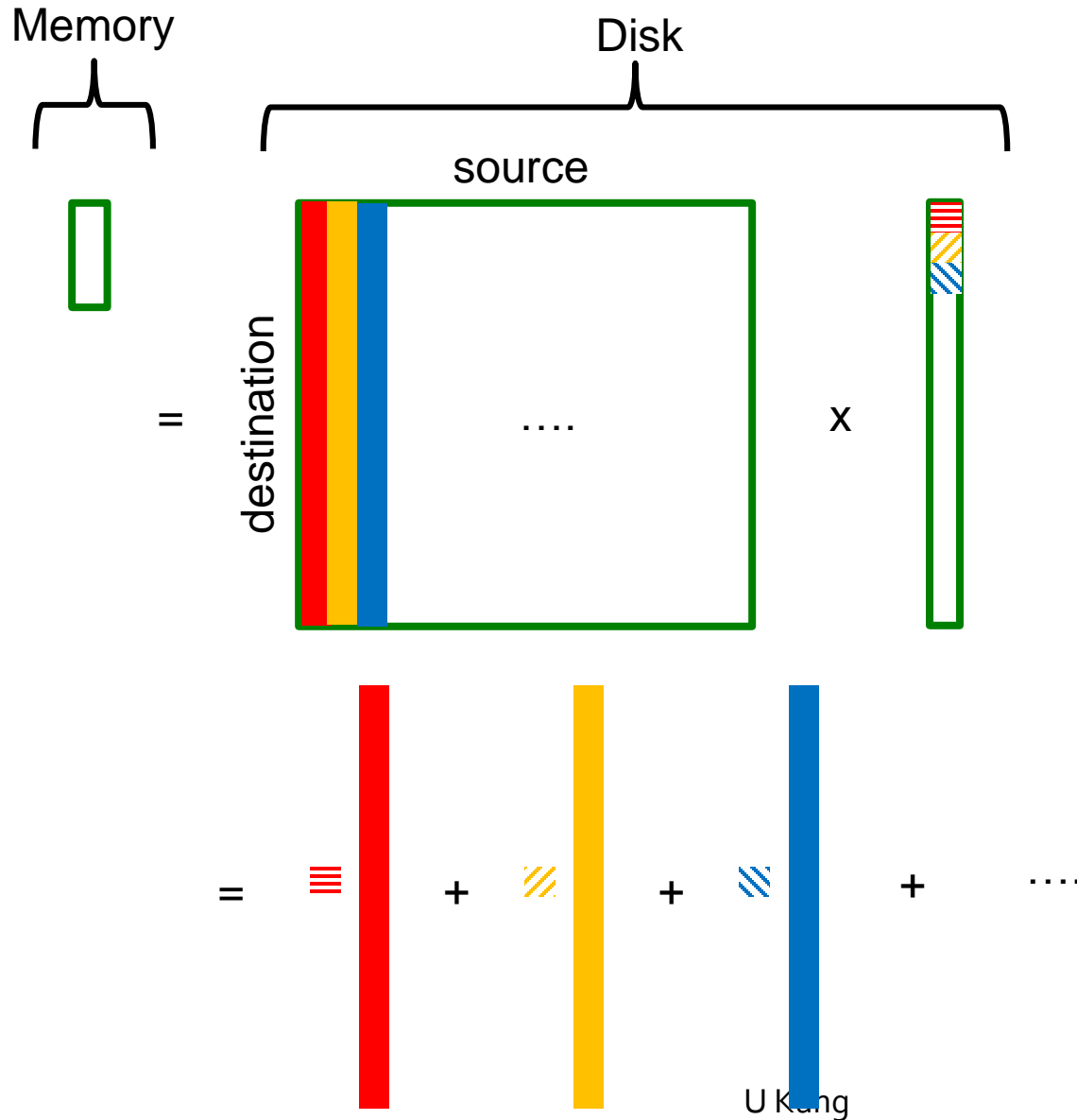
Block-based Update Algorithm



- Break r^{new} into k blocks that fit in memory
- Scan M and r^{old} once for each block



Block-based Update Algorithm



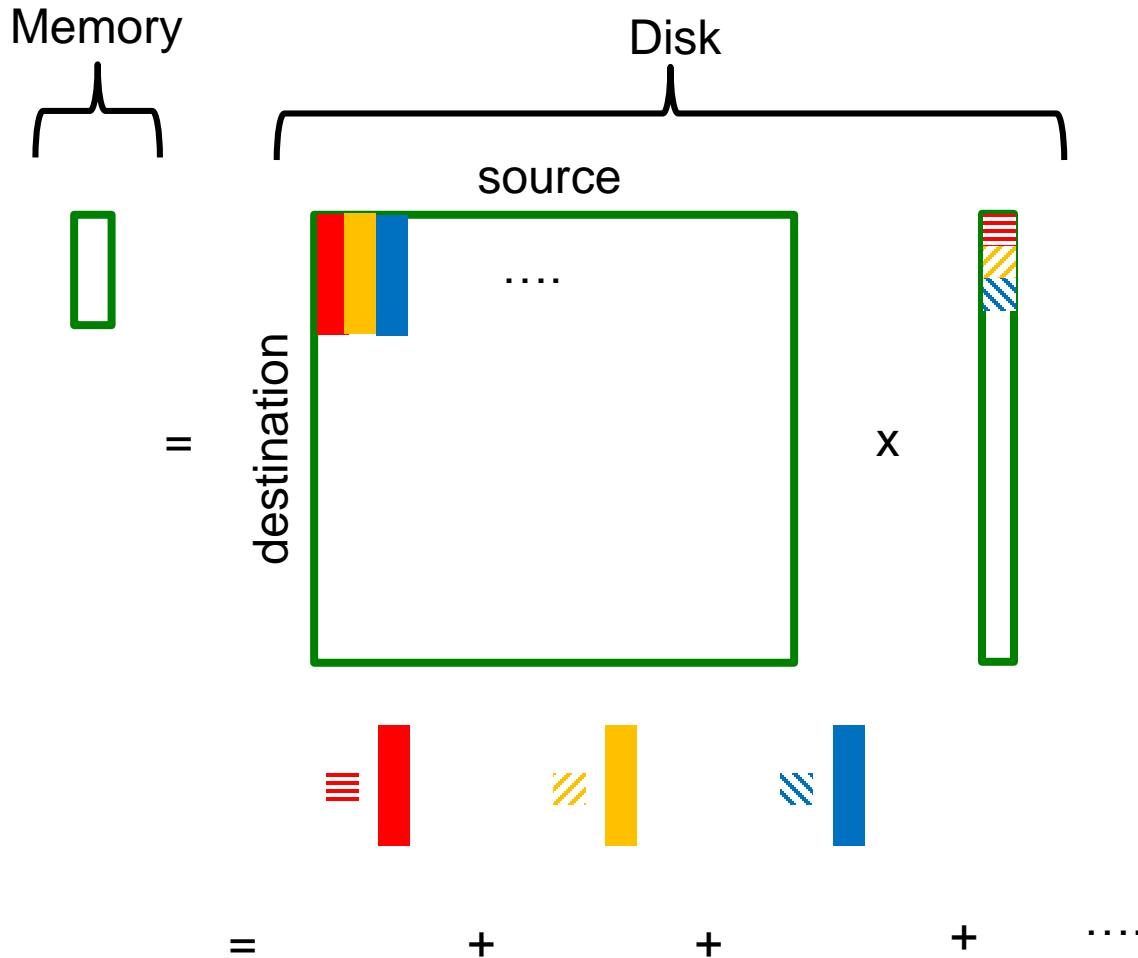


Analysis of Block Update

- **Similar to nested-loop join in databases**
 - Break r^{new} into k blocks that fit in memory
 - Scan M and r^{old} once for each block
- **Total cost:**
 - k scans of M and r^{old}
 - **Cost per iteration of Power method:**
 $k(|M| + |r|) + |r| = k|M| + (k+1)|r|$
- **Can we do better?**
 - **Hint:** M is much bigger than r (approx 10-20x), so we must avoid reading it k times per iteration

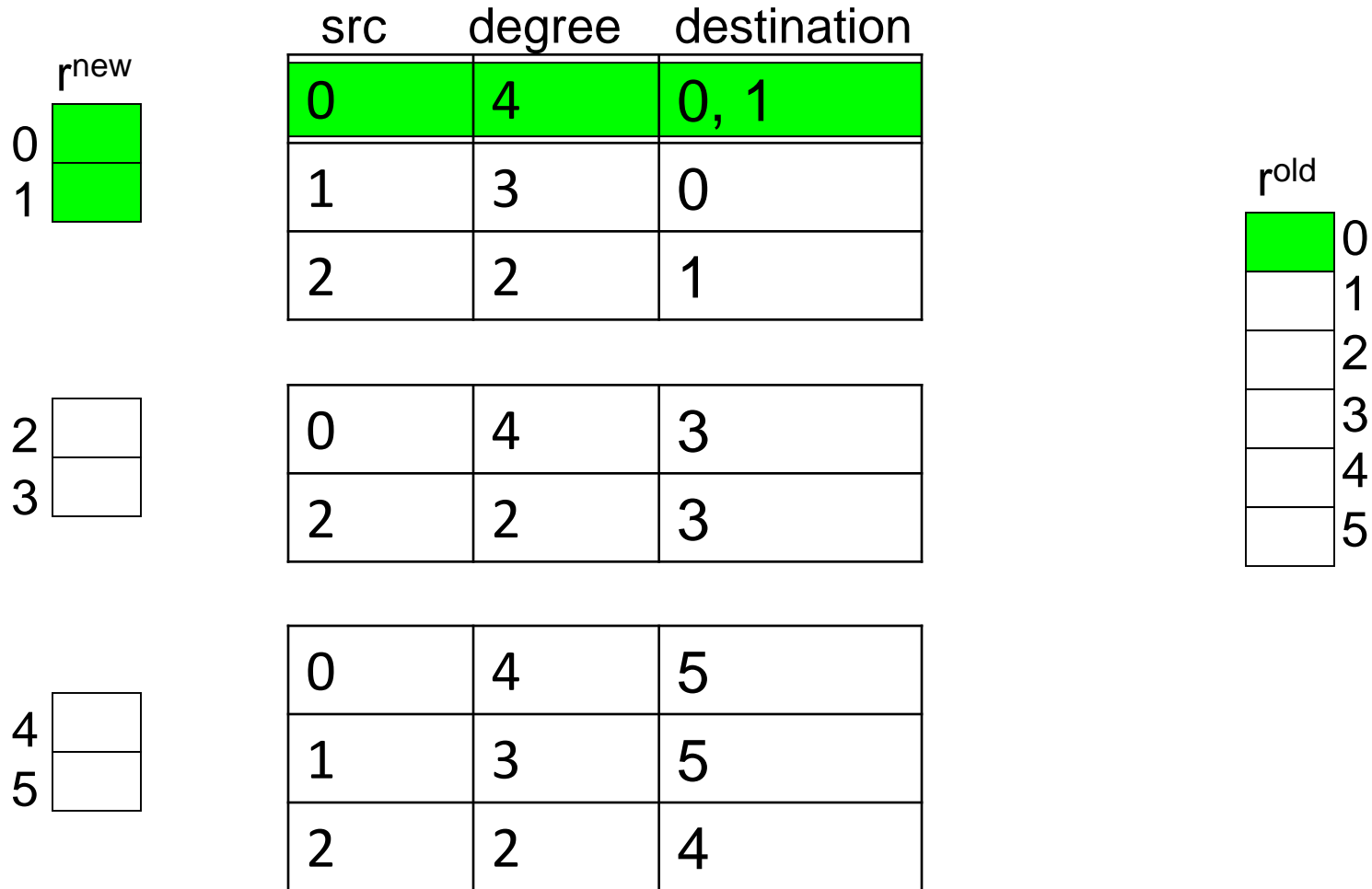


Block-Stripe Update Algorithm





Block-Stripe Update Algorithm



Break M into stripes! Each stripe contains only destination nodes in the corresponding block of r^{new}



Block-Stripe Analysis

- Break M into stripes
 - Each stripe contains only destination nodes in the corresponding block of r^{new}
- Some additional overhead per stripe
 - But it is usually worth it
- **Cost per iteration of Power method:**
 $= |M|(1+\varepsilon) + (k+1)|r|$




Limitations in Page Rank

- **Measures generic popularity of a page**
 - Biased against topic-specific authorities
 - **Solution:** Topic-Specific PageRank (**next**)
- **Uses a single measure of importance**
 - Other models of importance
 - **Solution:** Hubs-and-Authorities
- **Susceptible to Link spam**
 - Artificial link topographies created in order to boost page rank
 - **Solution:** TrustRank



Outline

- PageRank: Google Formulation
- Computing PageRank
-  **Topic-Specific PageRank**
- Measuring Proximity in Graphs



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:** Answer the query “Jaguar” differently



U Kang





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 picks a page from a set S
 - S contains only pages that are relevant to the topic
 - E.g., Yahoo or DMOZ pages for a given topic/query
 - For each teleport set S , we get a different vector r_S



Matrix Formulation

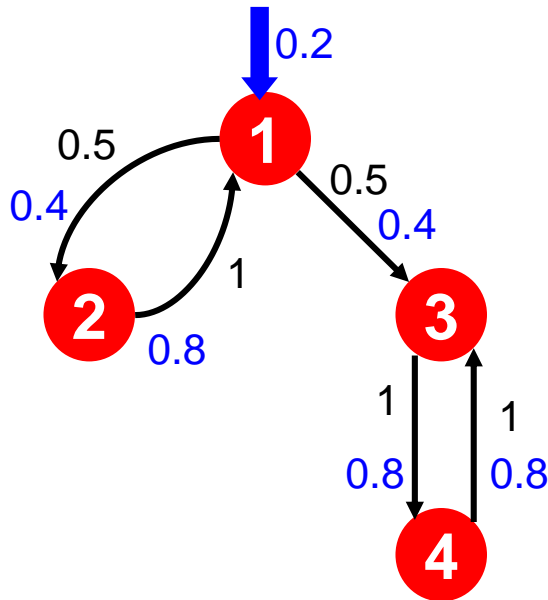
- To make this work all we need is to update the teleportation part of the PageRank formulation:

$$A_{ij} = \begin{cases} \beta M_{ij} + (1 - \beta)/|S| & \text{if } i \in S \\ \beta M_{ij} + 0 & \text{otherwise} \end{cases}$$

- A is column stochastic!
- We weighted all pages in the teleport set S equally
 - Could also assign different weights to pages!
- 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, ...

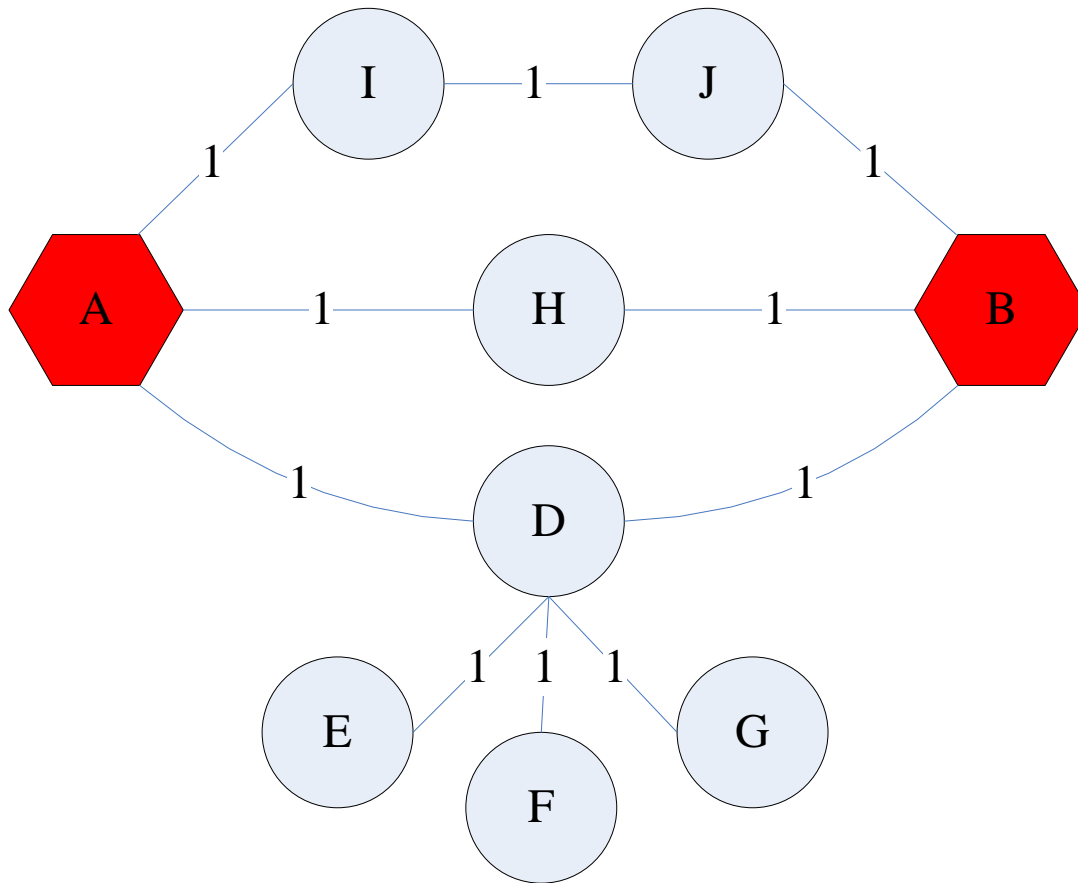


Outline

- PageRank: Google Formulation
- Computing PageRank
- Topic-Specific PageRank
-  **Measuring Proximity in Graphs**



Proximity on Graphs

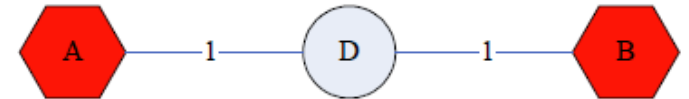
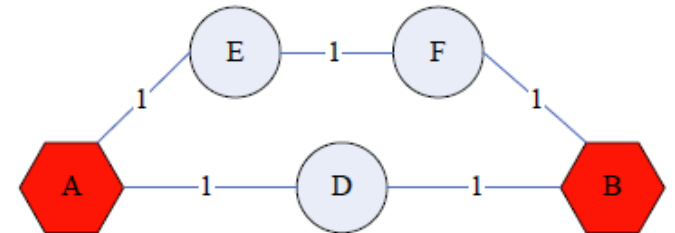
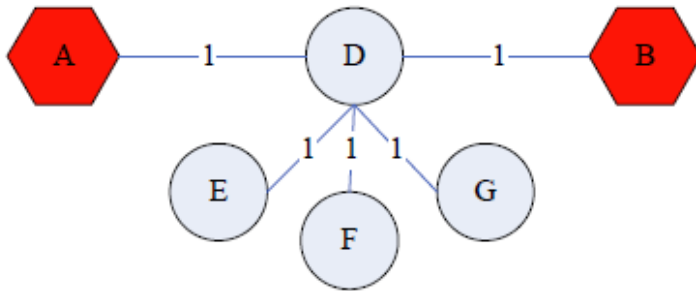
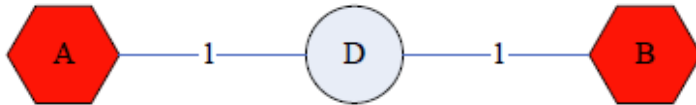


a.k.a.: Relevance, Closeness, 'Similarity'...



Good proximity measure?

- **Shortest path is not good:**

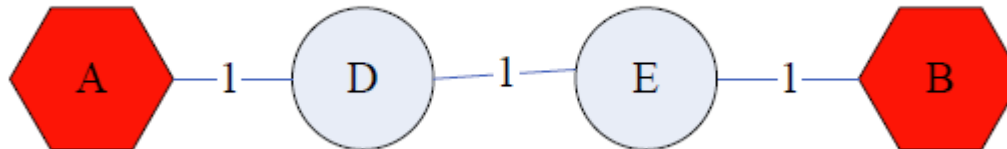
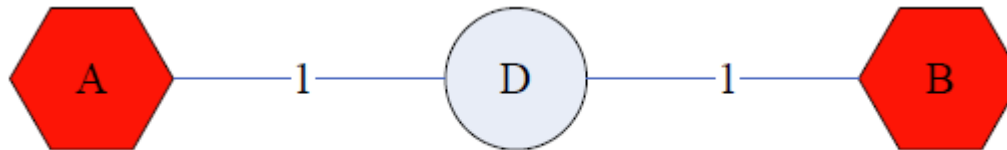


- **No effect of degree-1 nodes (E, F, G)!**
- Multi-faceted relationships



Good proximity measure?

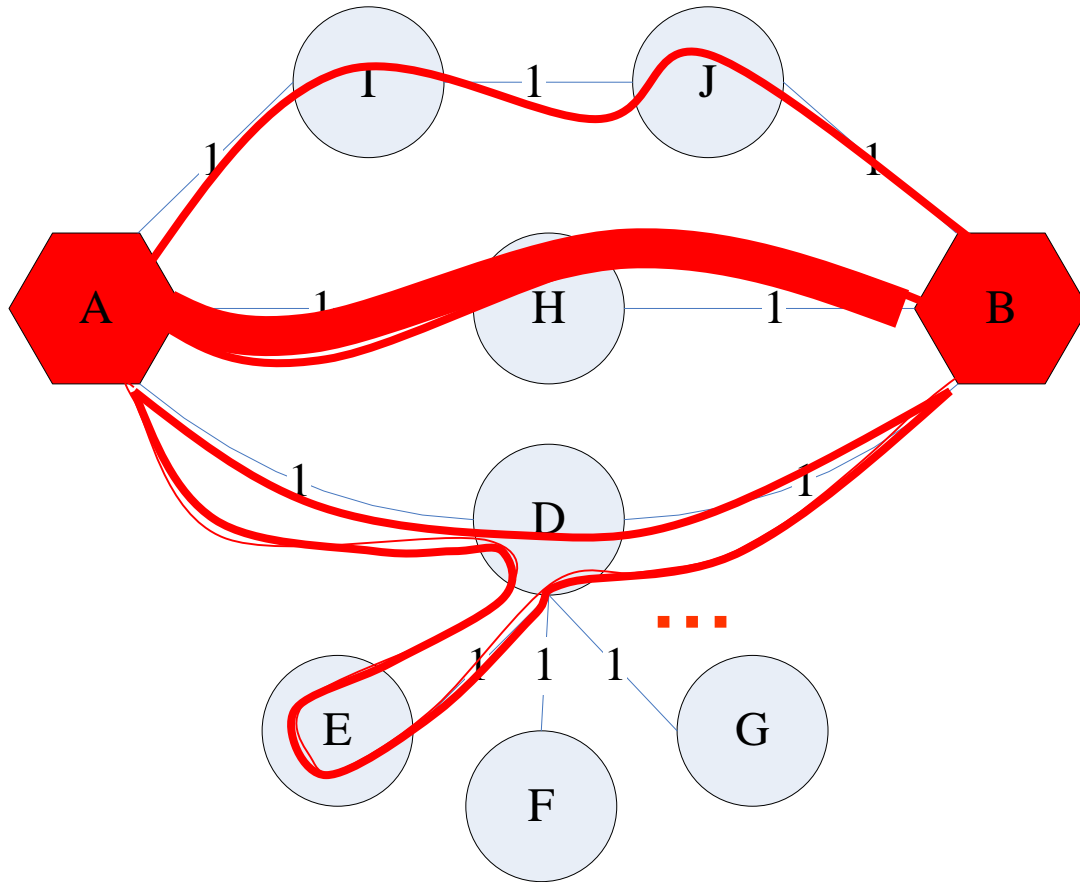
- Network flow is not good:



- Does not punish long paths



What is good notion of proximity?



- Multiple connections
- Quality of connection
 - Length, Degree, Weight...



Random Walk with Restart: Idea

- **RWR:** Random walks from a **fixed node**

- **E.g.,** k -partite graph with k types of nodes

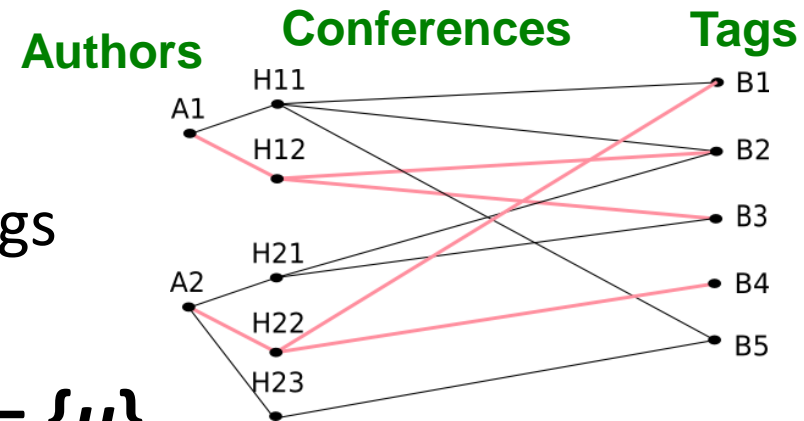
- E.g.: Authors, Conferences, Tags

- **Topic Specific PageRank** from node u : **teleport set** $S = \{u\}$

- **Resulting scores** measures similarity to node u

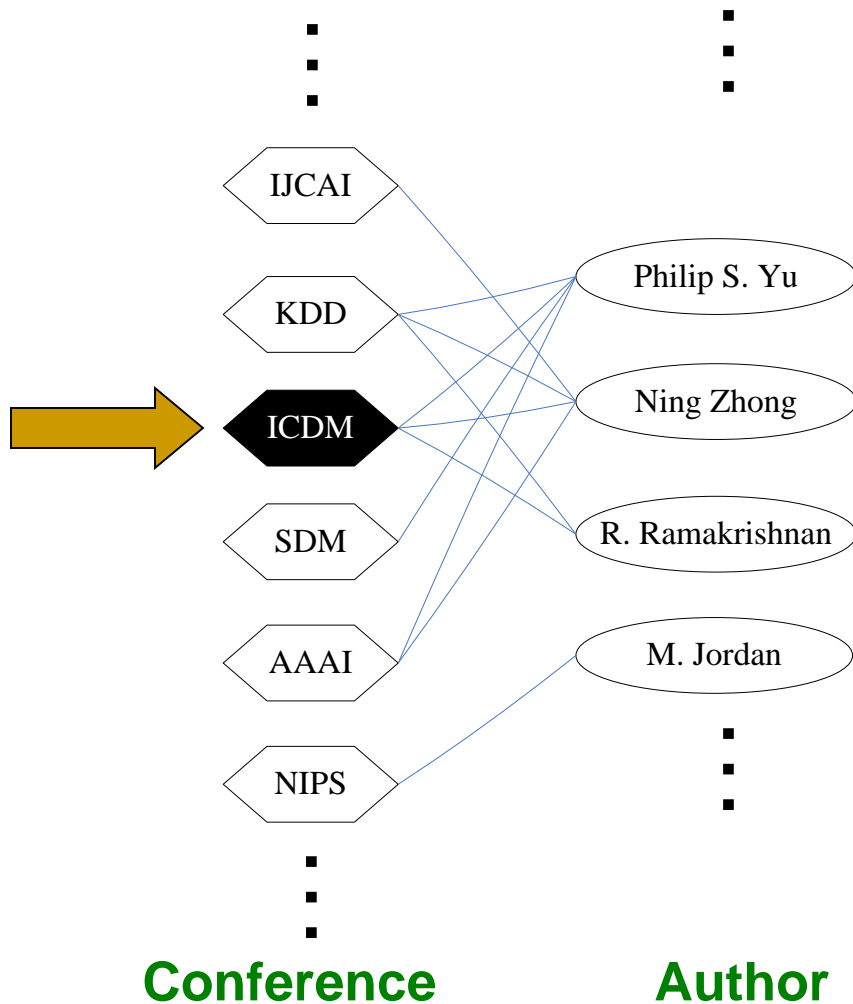
- **Problem:**

- Must be done once for each node u
- Suitable for sub-Web-scale applications





RWR: Example

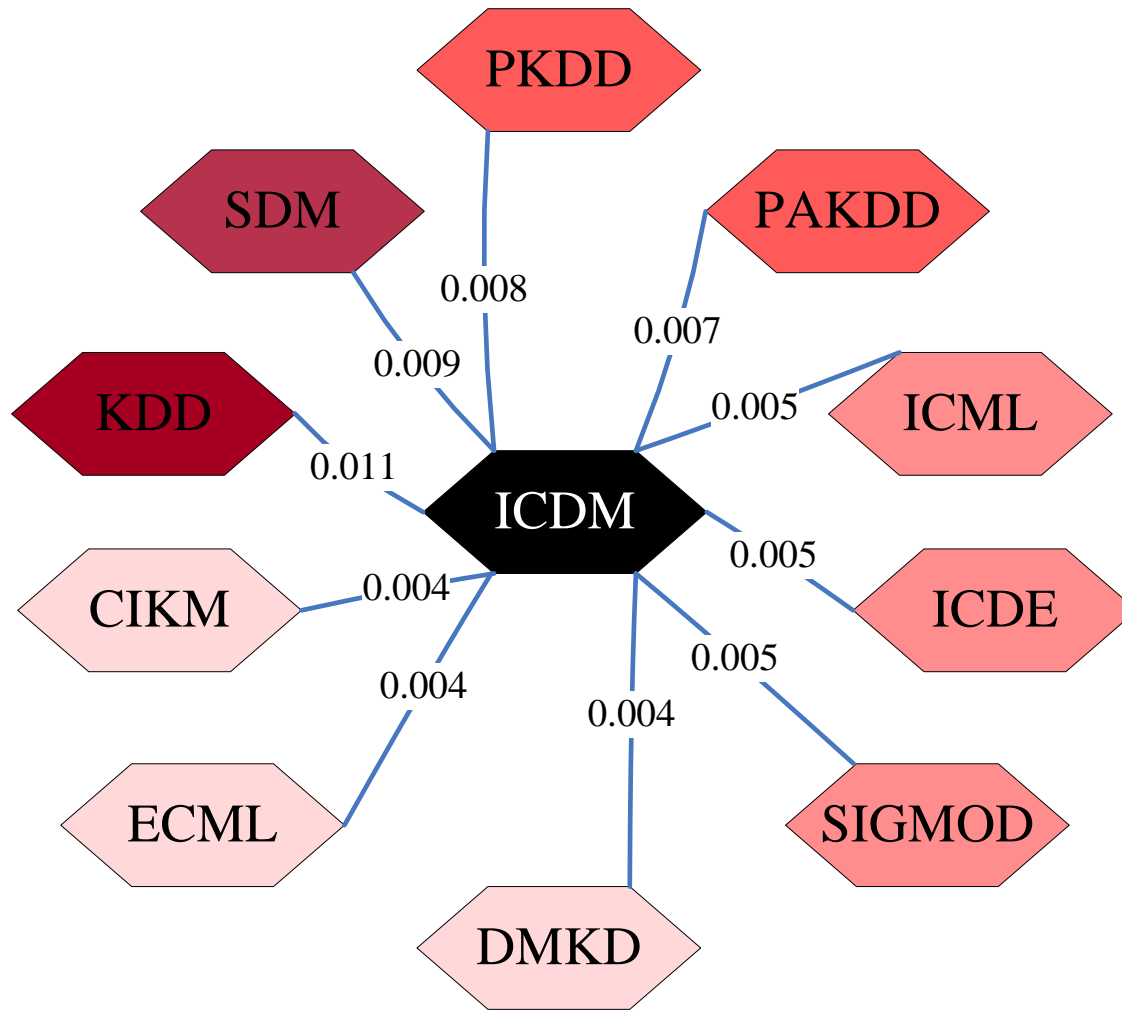


Q: What is the most related conference to ICDM?

A: Topic-Specific PageRank with teleport set $S=\{\text{ICDM}\}$



RWR: Example





What You Need to Know

■ “Normal” PageRank:

- Teleports uniformly at random to any node
- All nodes have the same probability of surfer landing there:
 $S = [0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1]$

■ Topic-Specific PageRank also known as Personalized PageRank:

- Teleports to a topic specific set of pages
- Nodes can have different probabilities of surfer landing there: $S = [0.1, 0, 0, 0.2, 0, 0, 0.5, 0, 0, 0.2]$

■ Random Walk with Restarts:

- Topic-Specific PageRank where teleport is always to the same node. $S = [0, 0, 0, 0, \mathbf{1}, 0, 0, 0, 0, 0]$



Questions?