

COMP9313: Big Data Management

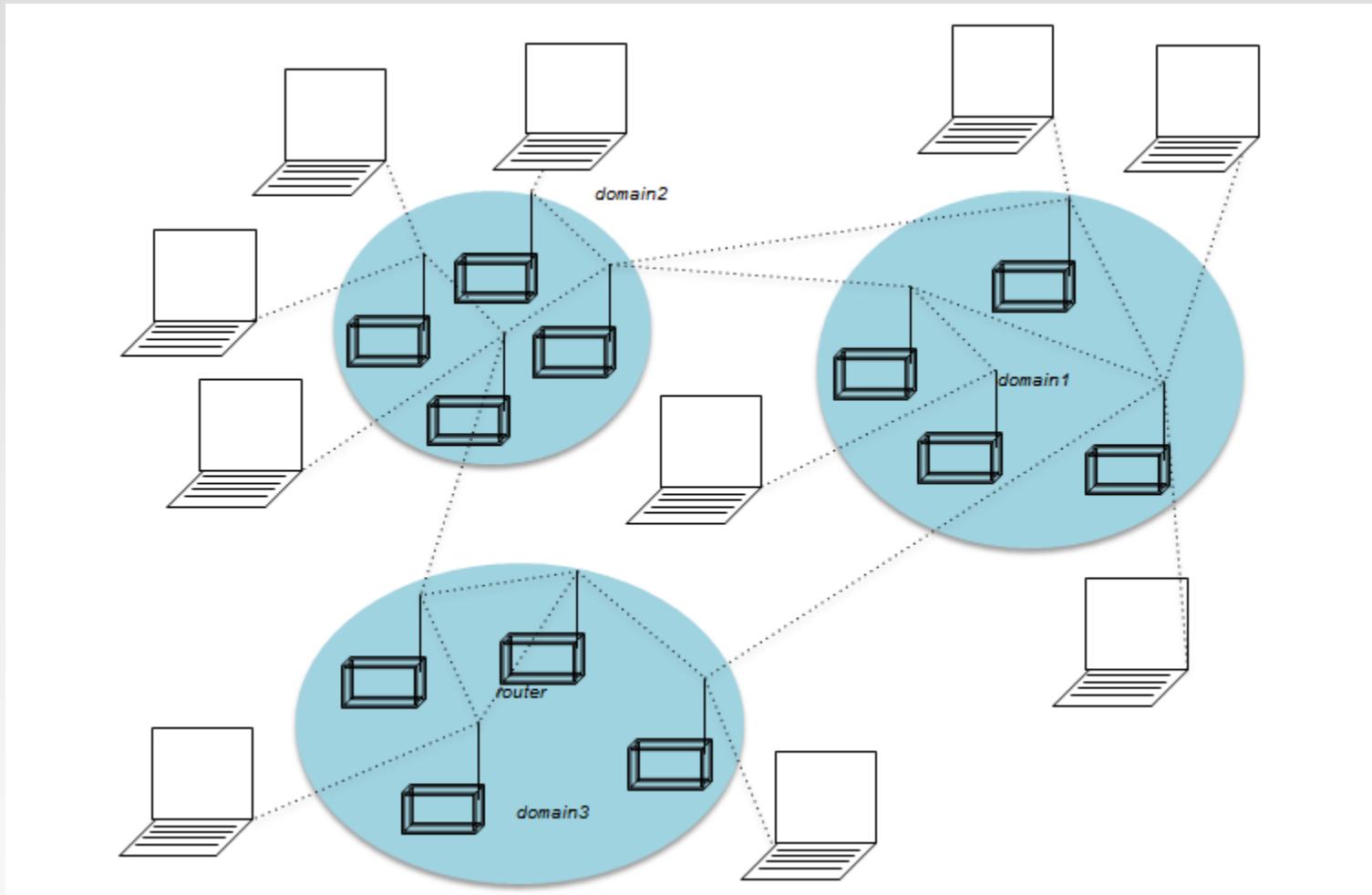


Lecturer: Xin Cao

Course web site: <http://www.cse.unsw.edu.au/~cs9313/>

Chapter 8.2: Link Analysis - PageRank

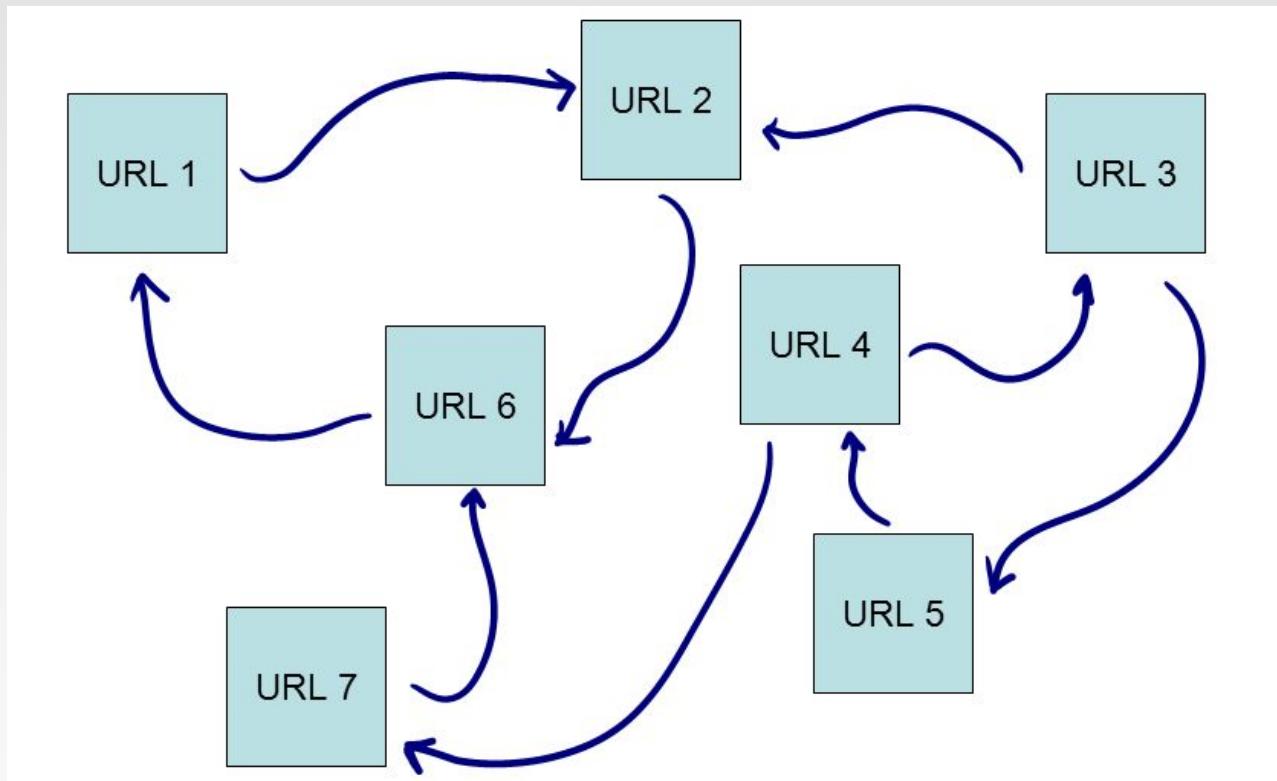
Graph Data: Communication Nets



Internet

Web as a Directed Graph

- ❖ Web as a directed graph:
 - Nodes: Webpages
 - Edges: Hyperlinks



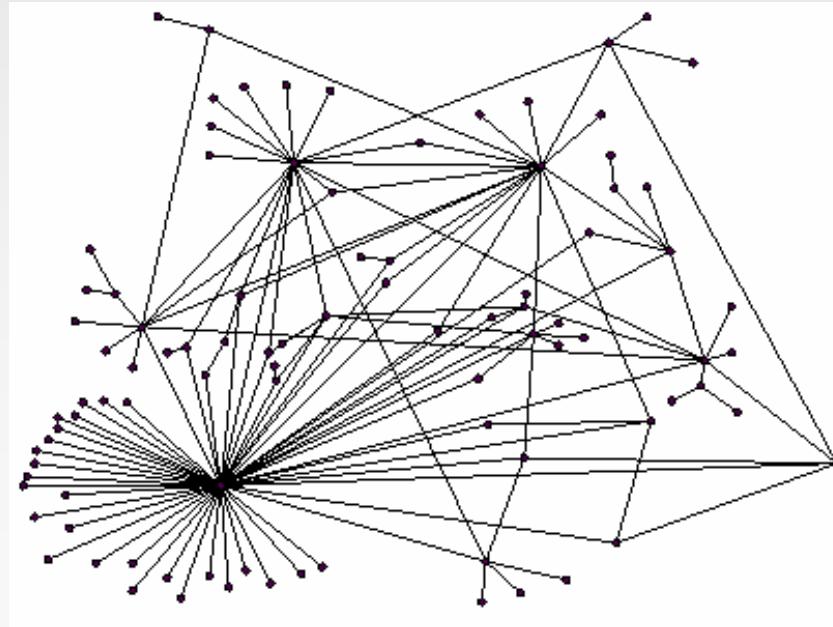
Broad Question

- ❖ How to organize the Web?
- ❖ First try: Human curated Web directories
 - Yahoo, LookSmart, etc.
- ❖ Second try: Web Search
 - Information Retrieval investigates:
Find relevant docs in a small
and trusted set
 - ▶ Newspaper articles, Patents, etc.
 - But: Web is huge, full of untrusted documents, random things, web spam, etc.
- ❖ What is the “best” answer to query “newspaper”?
 - No single right answer



Ranking Nodes on the Graph

- ❖ All web pages are not equally “important”
 - <http://xxx.github.io/> vs. <http://www.unsw.edu.au/>
- ❖ There is large diversity in the web-graph node connectivity. Let's rank the pages by the link structure!



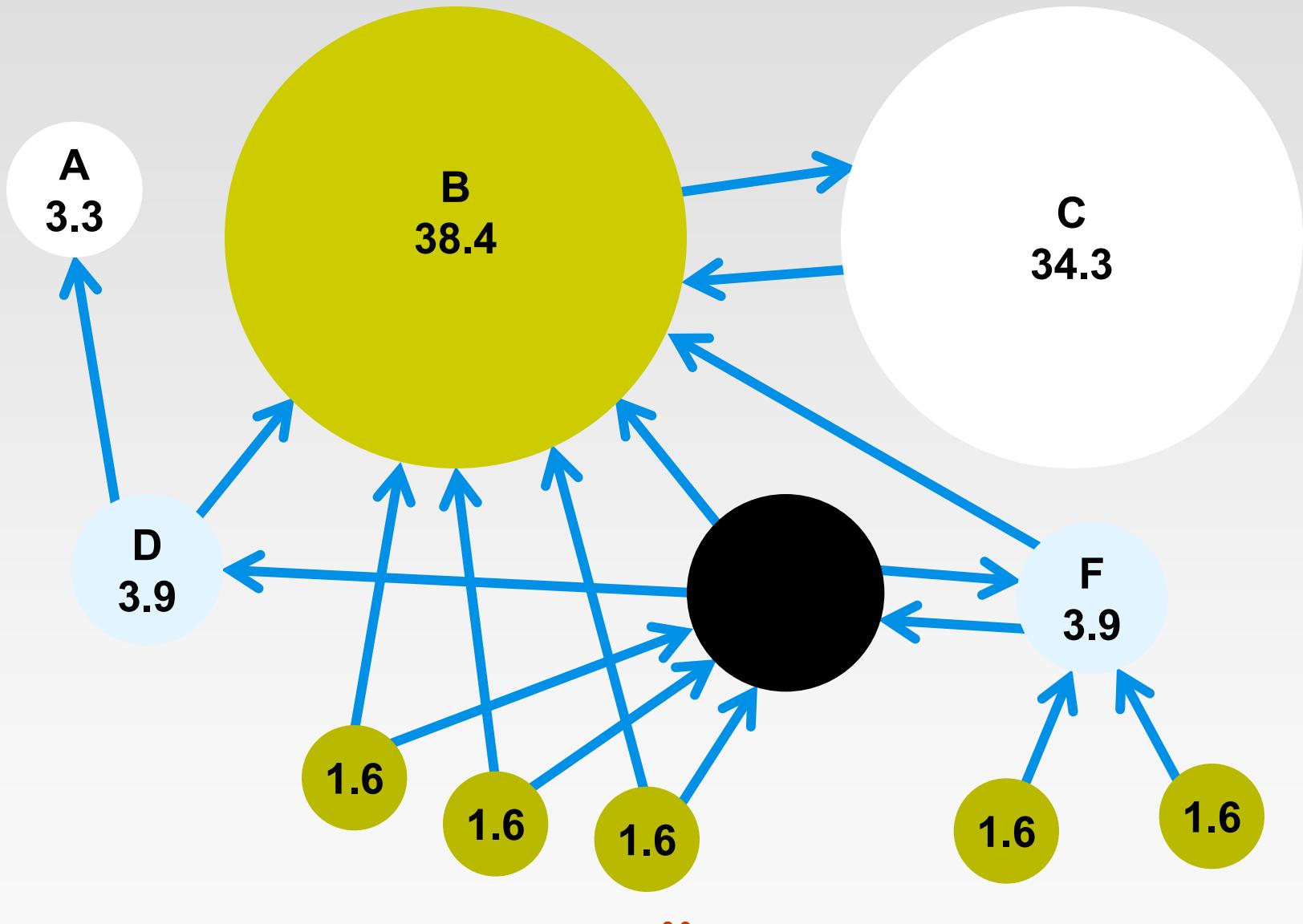
Link Analysis Algorithms

- ❖ We will cover the following Link Analysis approaches for computing importance of nodes in a graph:
 - Page Rank
 - Topic-Specific (Personalized) Page Rank
 - HITS
 -

Links as Votes

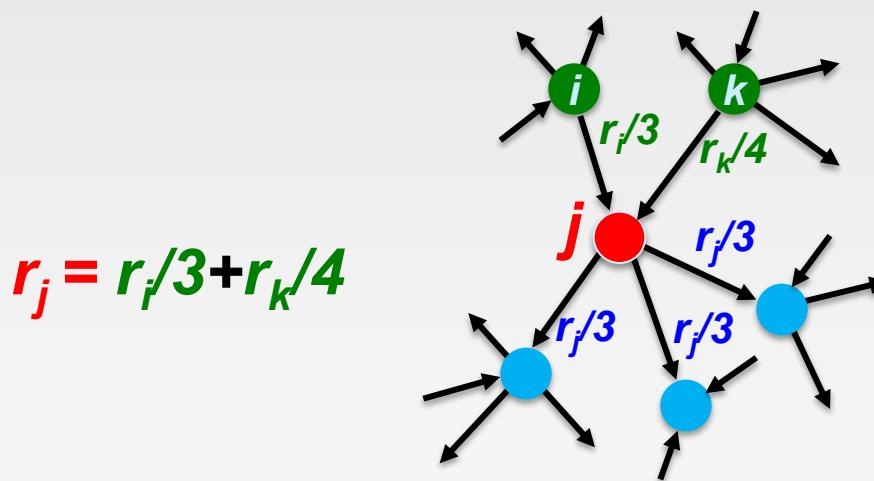
- ❖ Idea: Links as votes
 - Page is more important if it has more links
 - In-coming links? Out-going links?
- ❖ Think of in-links as votes:
 - <http://www.unsw.edu.au/> has 23,400 in-links
 - <http://xxx.github.io/> has 1 in-link
- ❖ Are all in-links equal?
 - Links from important pages count more
 - Recursive question!

Example: PageRank Scores



Simple Recursive Formulation

- ❖ Each link's vote is proportional to the **importance** of its source page
- ❖ If page j with importance r_j has n out-links, each link gets r_j / n votes
- ❖ Page j 's own importance is the sum of the votes on its in-links

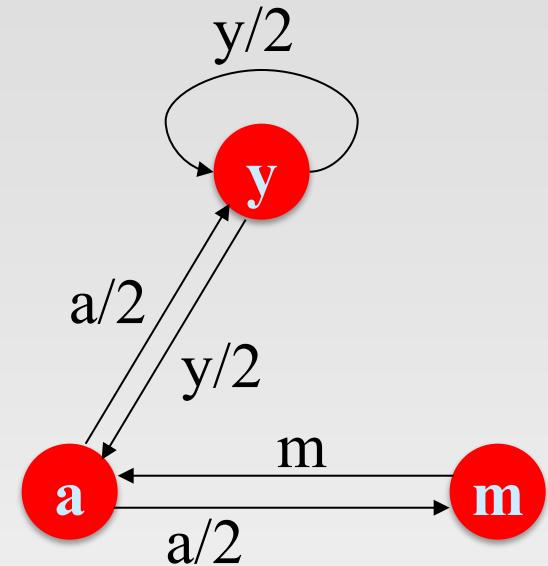


PageRank: The “Flow” Model

- ❖ A “vote” from an important page is worth more
- ❖ A page is important if it is pointed to by other important pages
- ❖ Define a “rank” r_j for page j

$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

d_i ... out-degree of node i



“Flow” equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

Solving the Flow Equations

- Flow equations:**
- ❖ 3 equations, 3 unknowns, no constants
 - No unique solution
 - All solutions equivalent modulo the scale factor
 - ❖ Additional constraint forces uniqueness:
 - $r_y + r_a + r_m = 1$
 - **Solution:** $r_y = \frac{2}{5}$, $r_a = \frac{2}{5}$, $r_m = \frac{1}{5}$
 - ❖ Gaussian elimination method works for small examples, but we need a better method for large web-size graphs
 - ❖ We need a new formulation!

PageRank: Matrix Formulation

- ❖ Stochastic adjacency matrix M

- Let page i has d_i out-links

- If $i \rightarrow j$, then $M_{ji} = \frac{1}{d_i}$ else $M_{ji} = 0$

- ▶ M is a **column stochastic matrix**

- Columns sum to 1

- ❖ **Rank vector r :** vector with an entry per page

- r_i is the importance score of page i

- $\sum_i r_i = 1$

- ❖ The flow equations can be written

$$\mathbf{r} = M \cdot \mathbf{r}$$

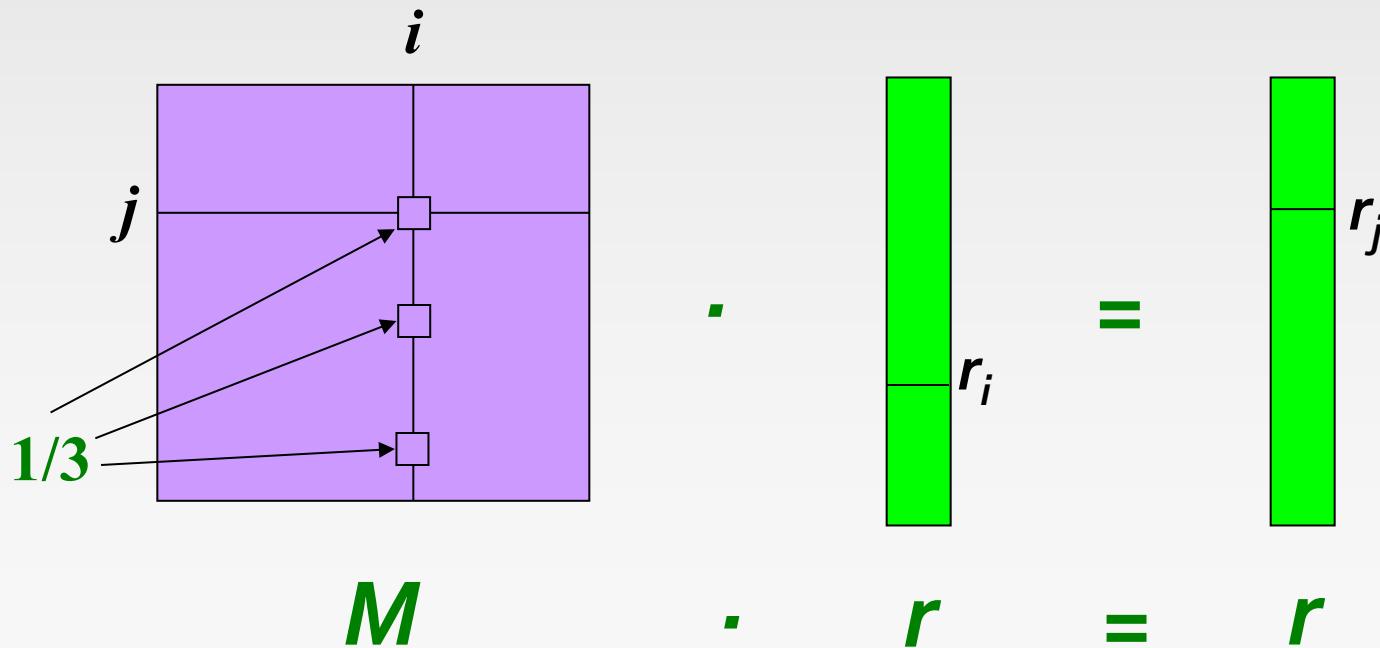
Example

- ❖ Remember the flow equation: $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$

- ❖ Flow equation in the matrix form

$$\mathbf{M} \cdot \mathbf{r} = \mathbf{r}$$

- Suppose page i links to 3 pages, including j



Eigenvector Formulation

- ❖ The flow equations can be written

$$\mathbf{r} = \mathbf{M} \cdot \mathbf{r}$$

- ❖ So the **rank vector r** is an **eigenvector** of the stochastic web matrix \mathbf{M}

- In fact, its first or principal eigenvector, with corresponding eigenvalue **1**

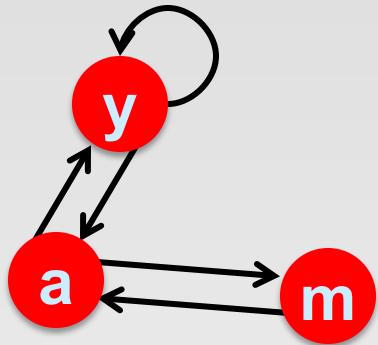
- Largest eigenvalue of \mathbf{M} is **1** since \mathbf{M} is column stochastic (with non-negative entries)

- We know r is unit length and each column of \mathbf{M} sums to one, so $\mathbf{M}r \leq 1$

NOTE: x is an eigenvector with the corresponding eigenvalue λ if:
 $Ax = \lambda x$

- ❖ We can now efficiently solve for r !
 - The method is called Power iteration

Example: Flow Equations & M



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

$$r = M \cdot r$$

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

$$\begin{bmatrix} y \\ a \\ m \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1 \\ 0 & 1/2 & 0 \end{bmatrix} \begin{bmatrix} y \\ a \\ m \end{bmatrix}$$

Power Iteration Method

- ❖ Given a web graph with n nodes, where the nodes are pages and edges are hyperlinks
- ❖ **Power iteration:** a simple iterative scheme
 - Suppose there are N web pages
 - Initialize: $\mathbf{r}^{(0)} = [1/N, \dots, 1/N]^T$
 - Iterate: $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$
 - Stop when $|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}|_1 < \varepsilon$

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

d_i out-degree of node i

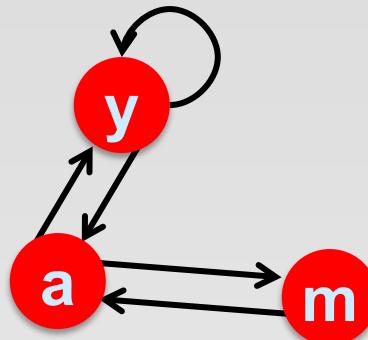
$|\mathbf{x}|_1 = \sum_{1 \leq i \leq N} |x_i|$ is the L_1 norm

Can use any other vector norm, e.g., Euclidean

PageRank: How to solve?

❖ Power Iteration:

- Set $r_j = 1/N$
- 1: $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- 2: $r = r'$
- Goto 1



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	1
m	0	$\frac{1}{2}$	0

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

❖ Example:

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 \\ 1/3 \\ 1/3 \end{matrix} \quad \dots \quad \begin{matrix} 6/15 \\ 6/15 \\ 3/15 \end{matrix}$$

Iteration 0, 1, 2, ...

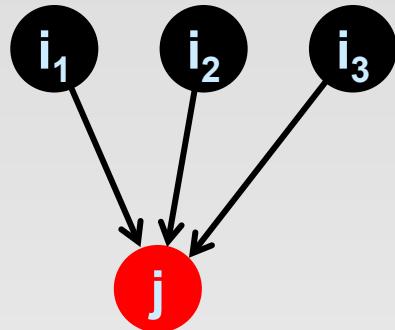
Random Walk Interpretation

- ❖ **Imagine a random web surfer:**

- At any time t , surfer is on some page i
- At time $t + 1$, the surfer follows an out-link from i uniformly at random
- Ends up on some page j linked from i
- Process repeats indefinitely

- ❖ **Let:**

- $p(t)$... vector whose i^{th} coordinate is the prob. that the surfer is at page i at time t
- So, $p(t)$ is a probability distribution over pages



$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_{\text{out}}(i)}$$

The Stationary Distribution

- ❖ Where is the surfer at time $t+1$?

- Follows a link uniformly at random

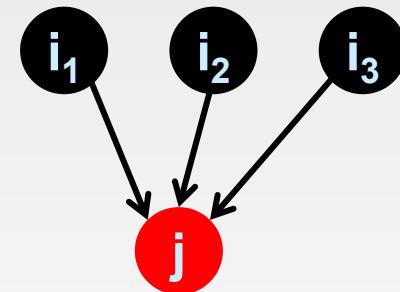
$$p(t+1) = M \cdot p(t)$$

- ❖ Suppose the random walk reaches a state $p(t+1) = M \cdot p(t) = p(t)$

then $p(t)$ is **stationary distribution** of a random walk

- ❖ Our original rank vector r satisfies $r = M \cdot r$

- So, r is a stationary distribution for the random walk



$$p(t+1) = M \cdot p(t)$$

Existence and Uniqueness

- ❖ A central result from the theory of random walks (a.k.a. Markov processes):

For graphs that satisfy **certain conditions**,
the **stationary distribution is unique** and
eventually will be reached no matter what the
initial probability distribution at time $t = 0$

PageRank: Two Questions

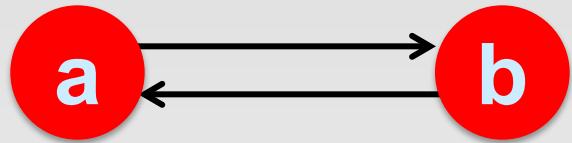
$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

or
equivalently

$$r = Mr$$

- ❖ Does this converge?
- ❖ Does it converge to what we want?

Does this converge?



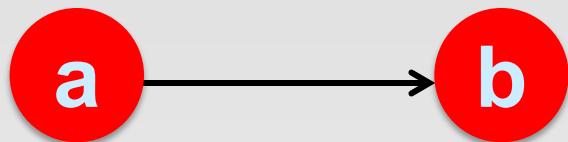
$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$

❖ Example:

r_a	1	0	1	0	...
r_b	0	1	0	1	...

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:

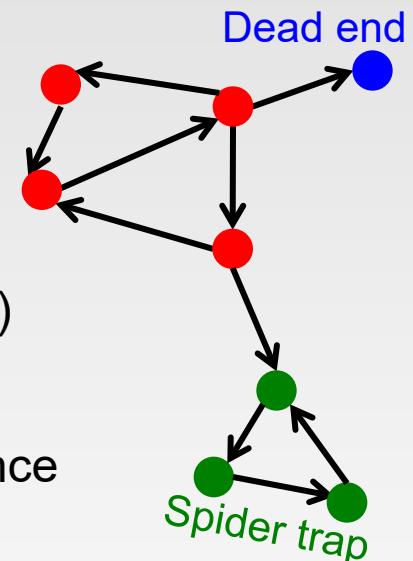
r_a	1	0	0	0
r_b	0	1	0	0

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”

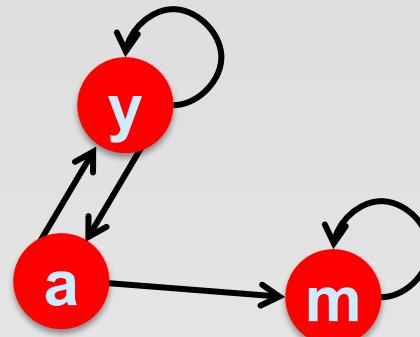


- ❖ (2) **Spider traps:** (all out-links are within the group)
 - Random walked gets “stuck” in a trap
 - And eventually spider traps absorb all importance

Problem: Spider Traps

❖ Power Iteration:

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



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	0
m	0	$\frac{1}{2}$	1

m is a spider trap

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

$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

❖ Example:

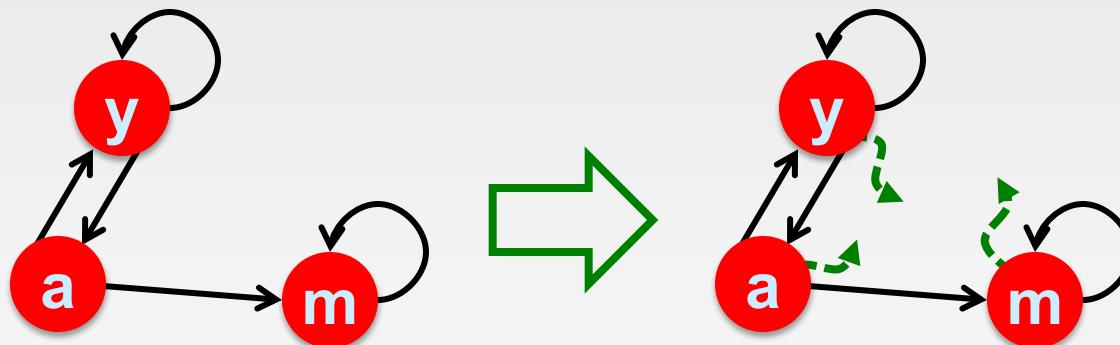
$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 3/6 & 7/12 & 16/24 & 1 \end{matrix}$$

Iteration 0, 1, 2, ...

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

Solution: Teleport!

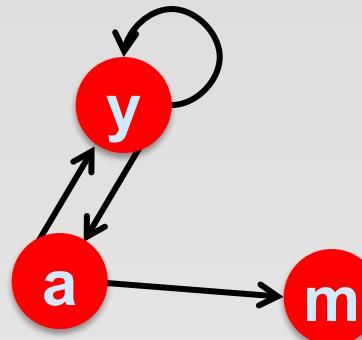
- ❖ The Google solution for spider traps: **At each time step, the random surfer has two options**
 - With prob. β , follow a link at random
 - With prob. $1-\beta$, jump to some random page
 - Common values for β are in the range 0.8 to 0.9
- ❖ Surfer will teleport out of spider trap within a few time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



Problem: Dead Ends

❖ Power Iteration:

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



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

$$\mathbf{r}_y = \mathbf{r}_y/2 + \mathbf{r}_a/2$$

$$\mathbf{r}_a = \mathbf{r}_y/2$$

$$\mathbf{r}_m = \mathbf{r}_a/2$$

❖ Example:

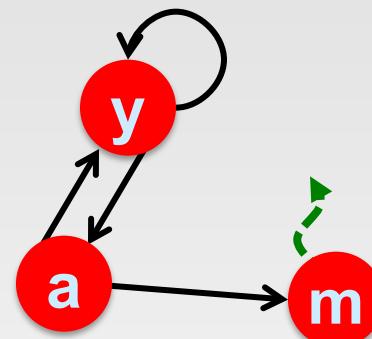
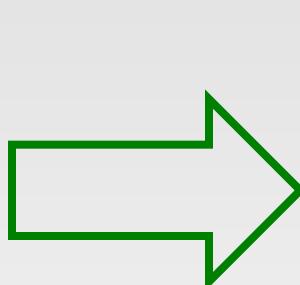
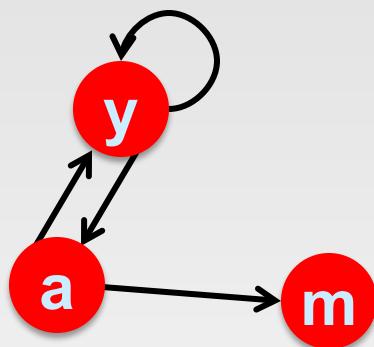
$$\begin{bmatrix} \mathbf{r}_y \\ \mathbf{r}_a \\ \mathbf{r}_m \end{bmatrix} = \begin{matrix} 1/3 & 2/6 & 3/12 & 5/24 & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & 0 \end{matrix}$$

Iteration 0, 1, 2, ...

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

Solution: Always Teleport!

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



	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	0
a	$\frac{1}{2}$	0	0
m	0	$\frac{1}{2}$	0

	y	a	m
y	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$
a	$\frac{1}{2}$	0	$\frac{1}{3}$
m	0	$\frac{1}{2}$	$\frac{1}{3}$

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

Google's Solution: Random Teleports

- ❖ Google's solution that does it all:

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]

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

- ❖ The Google Matrix A :

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

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

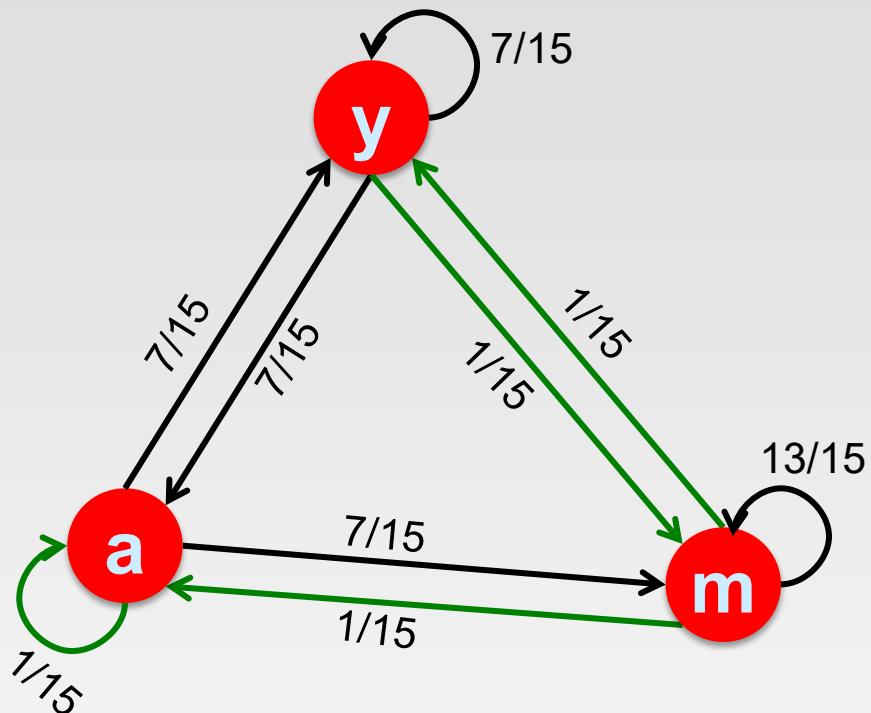
- ❖ We have a recursive problem: $r = A \cdot r$

And the Power method still works!

- ❖ What is β ?

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

Random Teleports ($\beta = 0.8$)



$$\begin{array}{c}
 M \\
 \boxed{\begin{matrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{matrix}} \\
 + 0.2 \\
 \boxed{\begin{matrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{matrix}} \\
 [1/N]_{NxN} \\
 A \\
 \boxed{\begin{matrix} y & 7/15 & 7/15 & 1/15 \\ a & 7/15 & 1/15 & 1/15 \\ m & 1/15 & 7/15 & 13/15 \end{matrix}}
 \end{array}$$

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

Computing Page Rank

- ❖ Key step is matrix-vector multiplication

- $r^{\text{new}} = A \cdot r^{\text{old}}$

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

- ❖ Say $N = 1$ billion pages

- We need 4 bytes for each entry (say)
 - 2 billion entries for vectors, approx 8GB
 - Matrix A has N^2 entries
 - ▶ 10^{18} is a large number!

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

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

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

Matrix Formulation

- ❖ Suppose there are N pages
- ❖ Consider page i , with d_i out-links
- ❖ We have $M_{ji} = 1/d_i$ when $i \rightarrow j$
and $M_{ji} = 0$ otherwise
- ❖ The random teleport is equivalent to:
 - Adding a **teleport link** from i to every other page and setting transition probability to $(1-\beta)/N$
 - Reducing the probability of following each out-link from $1/d_i$ to β/d_i
 - **Equivalent:** Tax each page a fraction $(1-\beta)$ of its score and redistribute evenly

Rearranging the Equation

❖ $r = A \cdot r$, where $A_{ji} = \beta M_{ji} + \frac{1-\beta}{N}$

❖ $r_j = \sum_{i=1}^N A_{ji} \cdot r_i$

❖ $r_j = \sum_{i=1}^N \left[\beta M_{ji} + \frac{1-\beta}{N} \right] \cdot r_i$
= $\sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N} \sum_{i=1}^N r_i$

= $\sum_{i=1}^N \beta M_{ji} \cdot r_i + \frac{1-\beta}{N}$ since $\sum r_i = 1$

❖ So we get: $r = \beta M \cdot r + \left[\frac{1-\beta}{N} \right]_N$

Note: Here we assumed **M** has no dead-ends

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

Sparse Matrix Formulation

- ❖ We just rearranged the **PageRank equation**

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

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

- ❖ \mathbf{M} is a **sparse matrix!** (with no dead-ends)

- 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 $(1-\beta)/N$ 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:** $\sum_j |r_j^{new} - r_j^{old}| > \varepsilon$

- ▶ $\forall j: r_j'^{new} = \sum_{i \rightarrow j} \beta \frac{r_i^{old}}{d_i}$
- $r_j'^{new} = 0$ if in-degree of j is 0

- ▶ Now re-insert the leaked PageRank:

$$\forall j: r_j^{new} = r_j'^{new} + \frac{1-S}{N} \quad \text{where: } S = \sum_j r_j'^{new}$$

- ▶ $r^{old} = r^{new}$

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
 - Say $10N$, or 4×10^9 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 $j = 1 \dots d_i$

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

0	r^{new}
1	
2	
3	
4	
5	
6	

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

r^{old}	0
1	
2	
3	
4	
5	
6	

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 per iteration of Power method:**
 $= 2|r| + |M|$
- ❖ Question:
 - What if we could not even fit r^{new} in memory?
 - Split r^{new} into blocks. Details ignored

Some Problems with Page Rank

- ❖ **Measures generic popularity of a page**
 - Biased against topic-specific authorities
 - **Solution:** Topic-Specific (Personalized) PageRank (**next**)
- ❖ **Uses a single measure of importance**
 - Other models of importance
 - **Solution:** Hubs-and-Authorities

PageRank in MapReduce

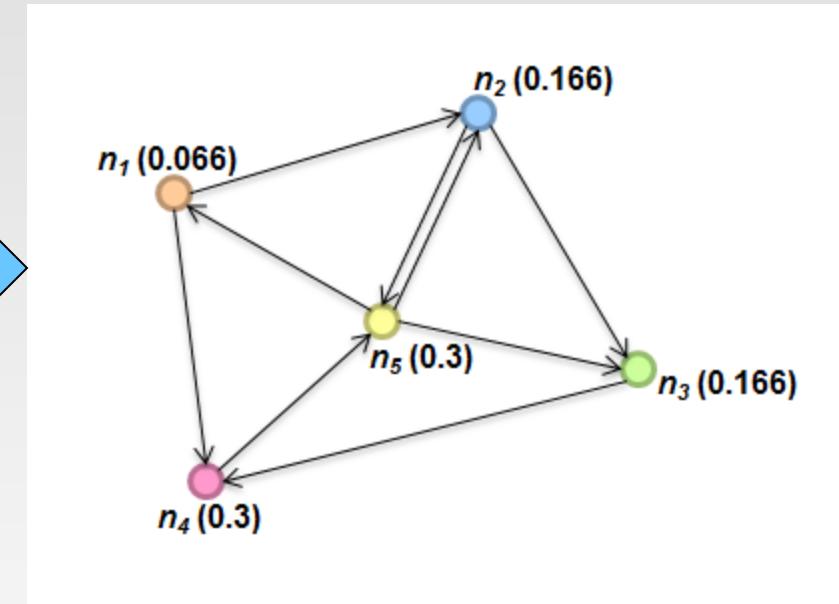
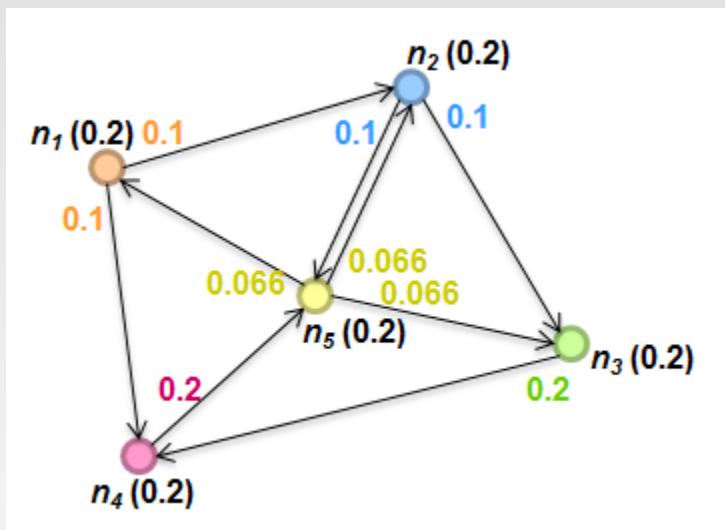
PageRank Computation Review

- ❖ Properties of PageRank
 - Can be computed iteratively
 - Effects at each iteration are local
- ❖ Sketch of algorithm:
 - Start with seed r_i values
 - Each page distributes r_i “credit” to all pages it links to
 - Each target page t_j adds up “credit” from multiple in-bound links to compute r_j
 - Iterate until values converge

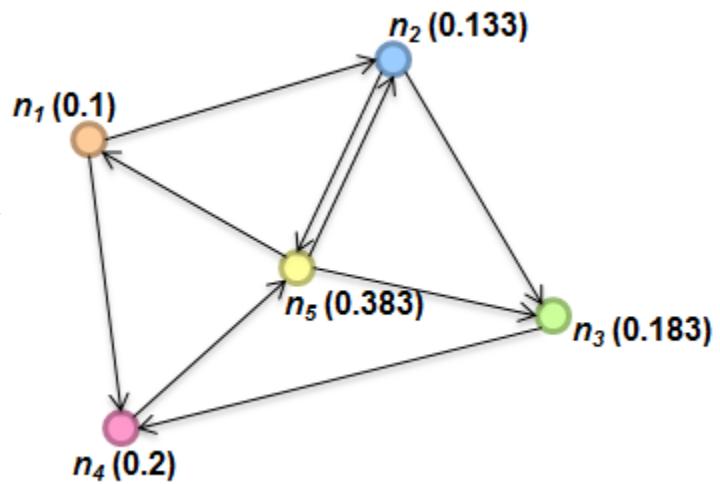
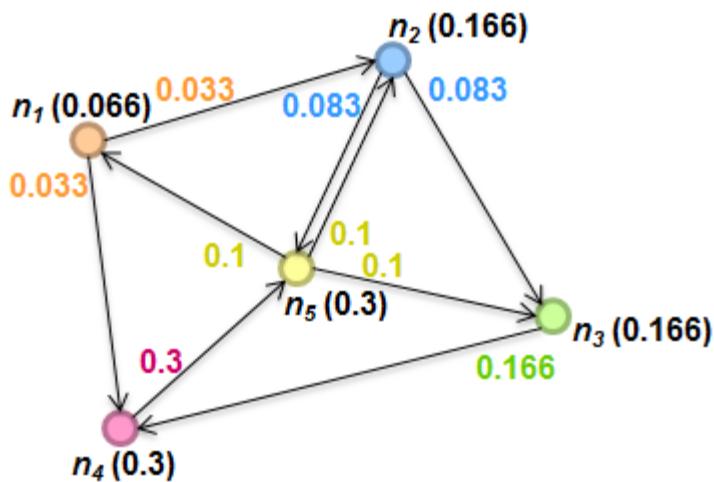
Simplified PageRank

- ❖ First, tackle the simple case:
 - No teleport
 - No dead ends
- ❖ Then, factor in these complexities...
 - How to deal with the teleport probability?
 - How to deal with dead ends?

Sample PageRank Iteration (1)



Sample PageRank Iteration (2)



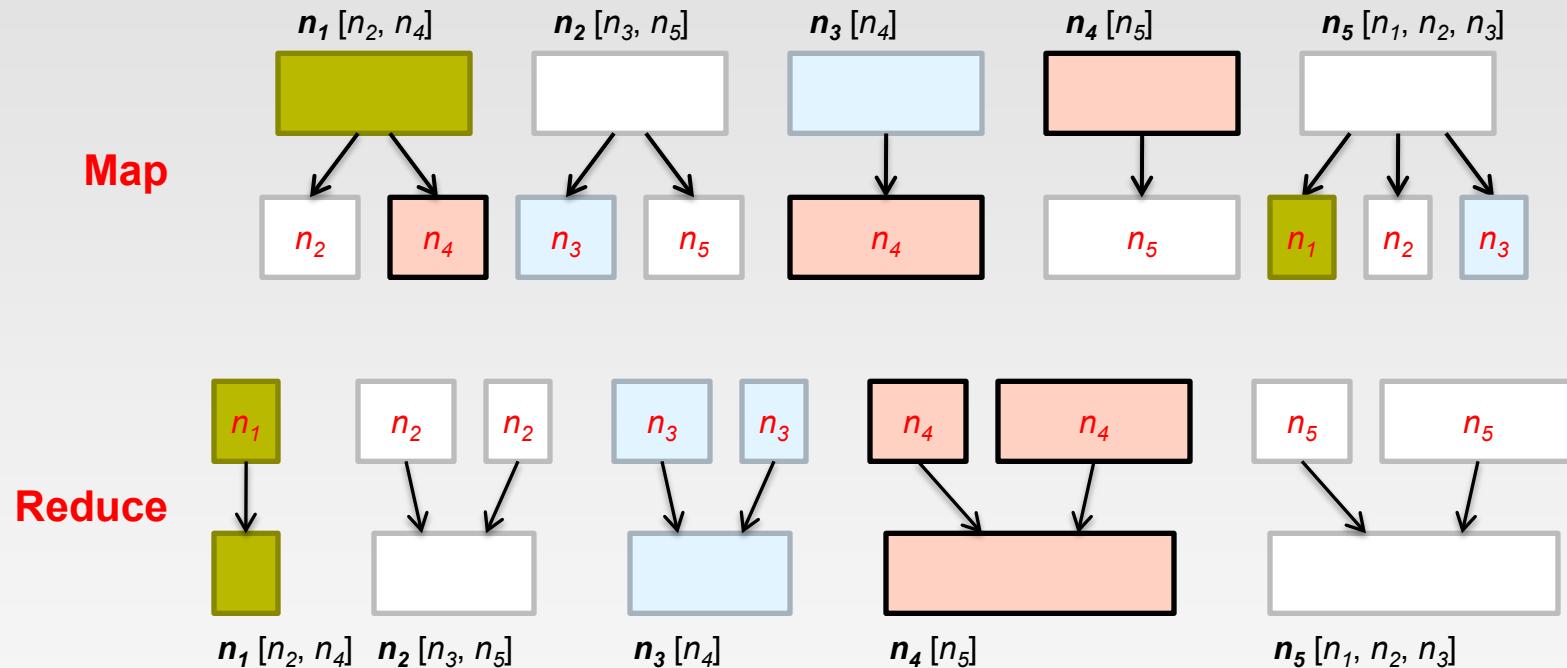
PageRank in MapReduce

- ❖ One iteration of the PageRank algorithm involves taking an estimated PageRank vector r and computing the next estimate r' by

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

- ❖ Mapper: input – a line containing node u , r_u , a list of out-going neighbors of u
 - For each neighbor v , emit(v , $r_u/\deg(u)$)
 - Emit (u , a list of out-going neighbors of u)
- ❖ Reducer: input – (node v , a list of values $\langle r_u/\deg(u), \dots \rangle$)
 - Aggregate the results according to the equation to compute r'_v
 - Emit node v , r'_v , a list of out-going neighbors of v

PageRank in MapReduce (One Iteration)



PageRank Pseudo-Code

```
1: class MAPPER
2:   method MAP(nid  $n$ , node  $N$ )
3:      $p \leftarrow N.\text{PAGERANK}/|N.\text{ADJACENCYLIST}|$ 
4:     EMIT(nid  $n, N$ )                                ▷ Pass along graph structure
5:     for all nodeid  $m \in N.\text{ADJACENCYLIST}$  do
6:       EMIT(nid  $m, p$ )                            ▷ Pass PageRank mass to neighbors
7:
8: class REDUCER
9:   method REDUCE(nid  $m, [p_1, p_2, \dots]$ )
10:     $M \leftarrow \emptyset$ 
11:    for all  $p \in \text{counts } [p_1, p_2, \dots]$  do
12:      if IsNODE( $p$ ) then
13:         $M \leftarrow p$                                 ▷ Recover graph structure
14:      else
15:         $s \leftarrow s + p$                           ▷ Sums incoming PageRank contributions
16:     $M.\text{PAGERANK} \leftarrow s$ 
17:    EMIT(nid  $m, \text{node } M$ )
```

Complete PageRank

- ❖ Two additional complexities
 - What is the proper treatment of dangling nodes?
 - How do we factor in the random jump factor?
- ❖ Solution:
 - If a node's adjacency list is empty, distribute its value to all nodes evenly.
 - ▶ In mapper, for such a node i , emit $(\text{nid } m, r/N)$ for each node m in the graph
 - Add the teleport value
 - ▶ In reducer, $\text{M.PageRank} = \beta * s + (1 - \beta) / N$

Graphs and MapReduce

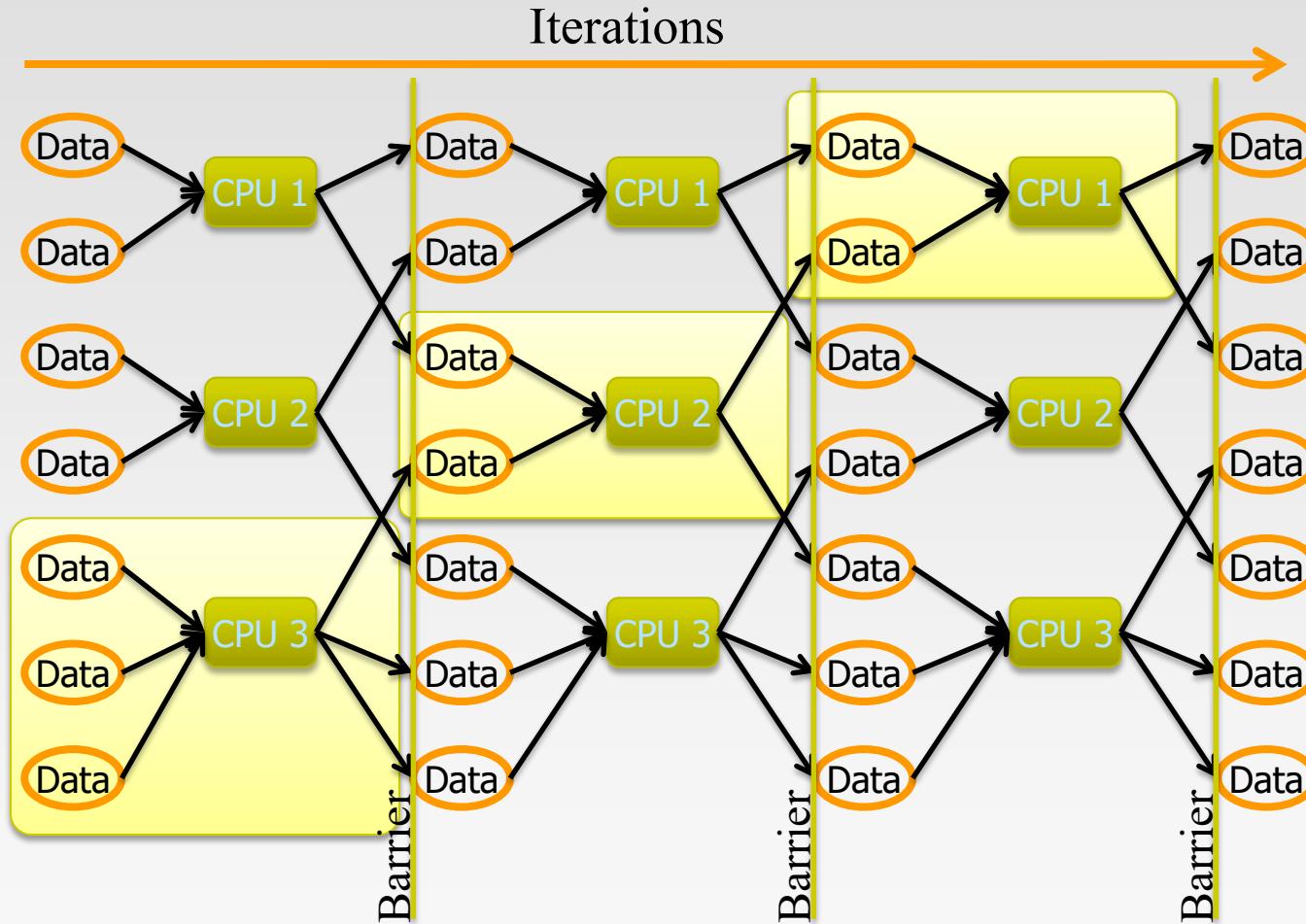
- ❖ Graph algorithms typically involve:
 - Performing computations at each node: based on node features, edge features, and local link structure
 - Propagating computations: “traversing” the graph
- ❖ Generic recipe:
 - Represent graphs as adjacency lists
 - Perform local computations in mapper
 - Pass along partial results via outlinks, keyed by destination node
 - Perform aggregation in reducer on inlinks to a node
 - Iterate until convergence: controlled by external “driver”
 - Don’t forget to pass the graph structure between iterations

Issues with MapReduce on Graph Processing

- ❖ MapReduce Does not support iterative graph computations:
 - External driver. Huge I/O incurs
 - No mechanism to support global data structures that can be accessed and updated by all mappers and reducers
 - ▶ Passing information is only possible within the local graph structure – through adjacency list
 - ▶ Some “wasted” computations are needed
- ❖ MapReduce algorithms are often impractical on large, dense graphs.
 - The amount of intermediate data generated is on the order of the number of edges.
 - For dense graphs, MapReduce running time would be dominated by copying intermediate data across the network.

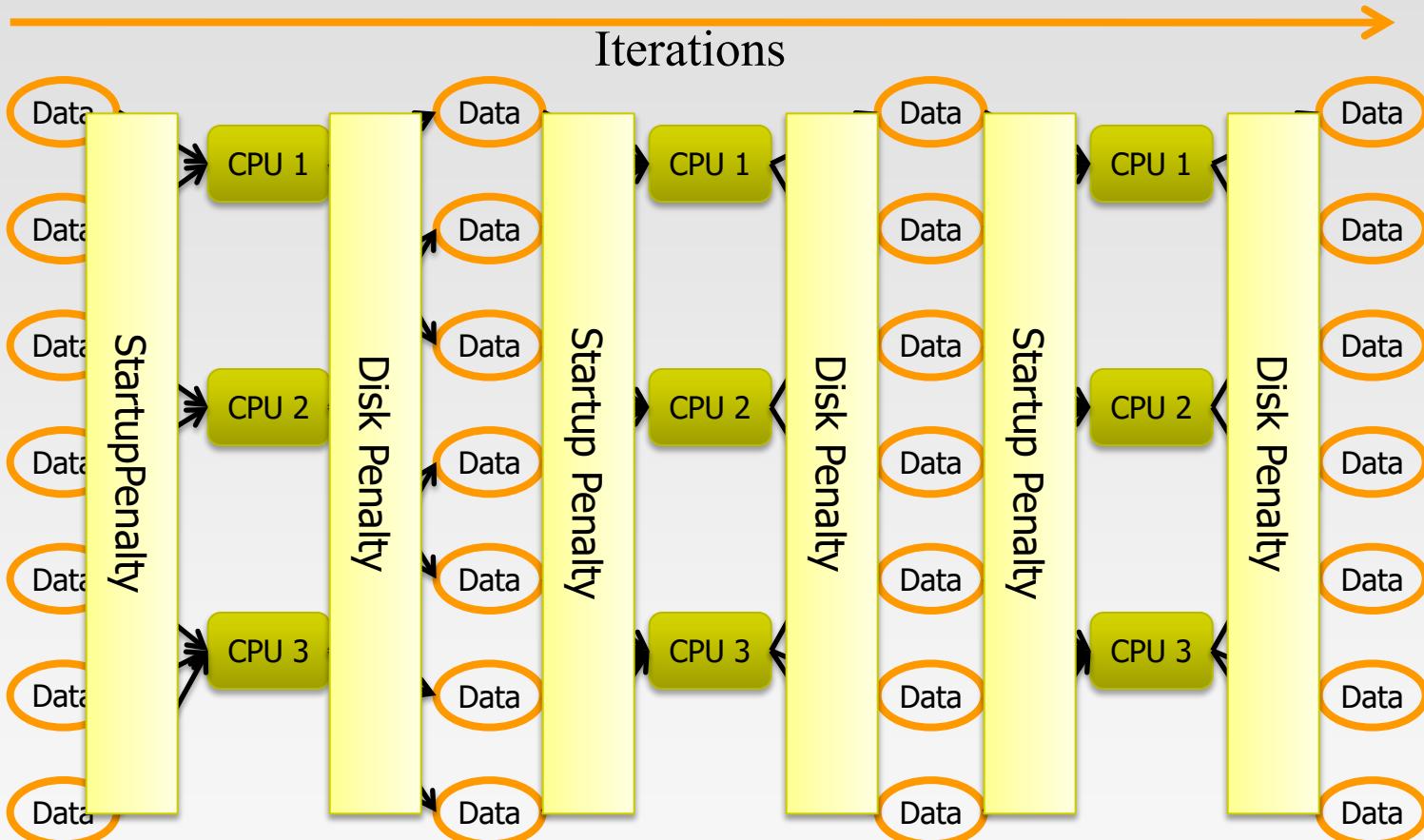
Iterative MapReduce

- ❖ Only a subset of data needs computation:



Iterative MapReduce

- ❖ System is not optimized for iteration:



References

- ❖ Chapter 5. Mining of Massive Datasets.
- ❖ Chapter 5. Data-Intensive Text Processing with MapReduce

End of Chapter 8.2