


Note to other teachers and users of these slides: We would be delighted if you found this our material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <http://www.mmds.org>

## Analysis of Large Graphs: Link Analysis, PageRank

Mining of Massive Datasets  
Jure Leskovec, Anand Rajaraman, Jeff Ullman  
Stanford University  
<http://www.mmds.org>



## Web Search: 2 Challenges

### 2 challenges of web search:

#### ■ (1) Web contains many sources of information Who to “trust”?

- Web is huge, full of untrusted documents, random things, web spam, etc.
- **Trick:** Trustworthy pages may point to each other!

#### ■ (2) What is the “best” answer to query “newspaper”?

- No single right answer
- **Trick:** Pages that actually know about newspapers might all be pointing to many newspapers

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

2

## Google's PageRank (Brin/Page 98)



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

3

## Google's PageRank (Brin/Page 98)

- A technique for estimating page quality
  - Based on web link graph
- Results are combined with IR score
  - Think of it as:  $\text{TotalScore} = \text{IR score} * \text{PageRank}$
  - In practice, search engines use many other factors
  - (for example, Google says it uses more than 200 features)

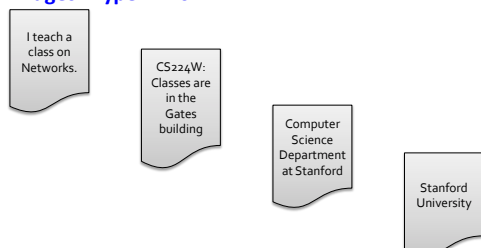
J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

4

## Web as a Graph

### ■ Web as a directed graph:

- **Nodes:** Webpages
- **Edges:** Hyperlinks



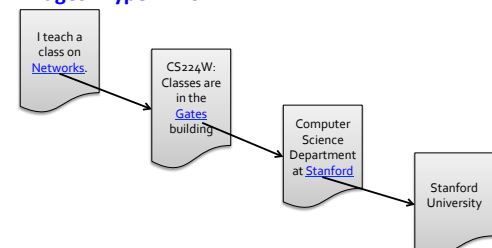
J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

5

## Web as a Graph

### ■ Web as a directed graph:

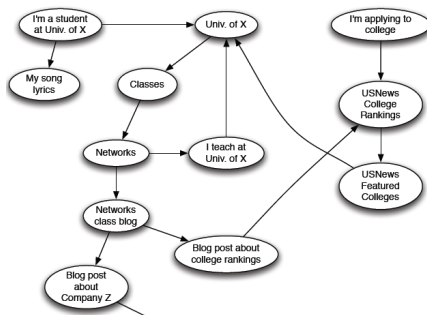
- **Nodes:** Webpages
- **Edges:** Hyperlinks



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

6

## Web as a Directed Graph



J. Leskovec, A. Rajaraman, I. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

7

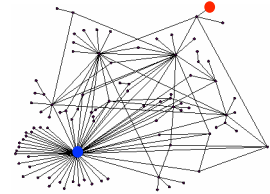
## Ranking Nodes on the Graph

- All web pages are not equally "important"

[www.joe-schmoe.com](http://www.joe-schmoe.com) vs. [www.stanford.edu](http://www.stanford.edu)

- There is large diversity in the web-graph node connectivity.

Let's rank the pages by the link structure!



J. Leskovec, A. Rajaraman, I. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

8

## PageRank: The "Flow" Formulation

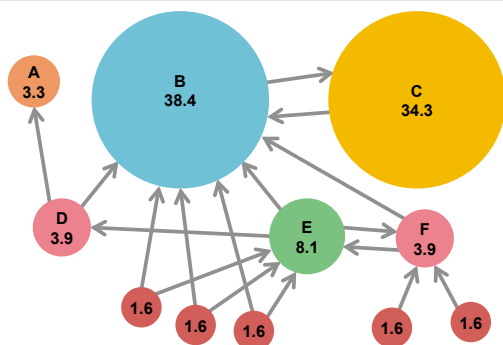
## Links as Votes

- 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:
  - [www.stanford.edu](http://www.stanford.edu) has 23,400 in-links
  - [www.joe-schmoe.com](http://www.joe-schmoe.com) has 1 in-link
- Are all in-links equal?
  - Links from important pages count more
  - Recursive question!

J. Leskovec, A. Rajaraman, I. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

10

## Example: PageRank Scores



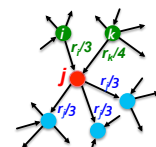
J. Leskovec, A. Rajaraman, I. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

11

## 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

$$r_j = r_j/3 + r_k/4$$



J. Leskovec, A. Rajaraman, I. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

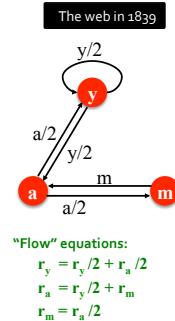
12

## 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



## Solving the Flow Equations

- 3 equations, 3 unknowns, no constants

$$\begin{aligned} r_y &= r_y/2 + r_a/2 \\ r_a &= r_y/2 + r_m \\ r_m &= r_a/2 \end{aligned}$$

- No unique solution
- All solutions equivalent modulo the scale factor

- Additional constraint forces uniqueness:

- $r_y + r_a + r_m = 1$
- Solution:  $r_y = 2/5, r_a = 2/5, r_m = 1/5$

- 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,  $M_{ji} = 1/d_i$  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

$$r = M \cdot r$$

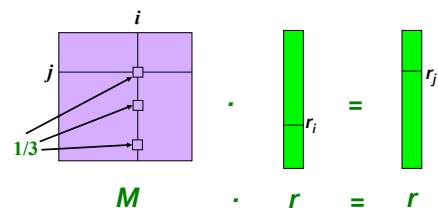
$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

## Example

- Remember the flow equation:  $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- Flow equation in the matrix form

$$M \cdot r = r$$

- Suppose page  $i$  links to 3 pages, including  $j$

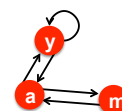


## Eigenvector Formulation

- The flow equations can be written
 
$$r = M \cdot r$$
- So the rank vector  $r$  is an eigenvector of the stochastic web matrix  $M$ 
  - In fact, its first or principal eigenvector, with corresponding eigenvalue 1
    - Largest eigenvalue of  $M$  is 1 since  $M$  is column stochastic
      - Why? We know  $r$  is unit length and each column of  $M$  sums to one, so  $Mr \leq r$
- We can now efficiently solve for  $r$ !
  - The method is called Power iteration

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

## Example: Flow Equations & $M$



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

$$r = M \cdot r$$

$$\begin{aligned} r_y &= r_y/2 + r_a/2 \\ r_a &= r_y/2 + r_m \\ r_m &= r_a/2 \end{aligned}$$

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

## 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:  $\mathbf{r}^{(t+1)} = \mathbf{M} \cdot \mathbf{r}^{(t)}$ 

$$r_j^{(t+1)} = \sum_{i \rightarrow j} \frac{r_i^{(t)}}{d_i}$$
 $d_i \dots \text{out-degree of node } i$
  - Stop when  $\|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}\|_1 < \epsilon$ 
    - $\|\mathbf{x}\|_1 = \sum_{1 \leq i \leq N} |x_i|$  is the  $L_1$  norm
    - Can use any other vector norm, e.g., Euclidean

J. Leskovec, A. Rajaraman, J. Ullman, Mining of Massive Datasets, <http://www.mmds.org>

19

## PageRank: How to solve?

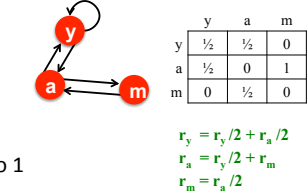
### Power Iteration:

- Set  $r_i = 1/N$
- 1:  $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- 2:  $\mathbf{r} = \mathbf{r}'$
- If not converged: goto 1

### Example:

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{bmatrix} 1/3 & 1/3 & 5/12 & 9/24 & \dots & 6/15 \\ 1/3 & 3/6 & 1/3 & 11/24 & \dots & 6/15 \\ 1/3 & 1/6 & 3/12 & 1/6 & \dots & 3/15 \end{bmatrix}$$

Iteration 0, 1, 2,



J. Leskovec, A. Rajaraman, J. Ullman, Mining of Massive Datasets, <http://www.mmds.org>

20

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

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

J. Leskovec, A. Rajaraman, J. Ullman, Mining of Massive Datasets, <http://www.mmds.org>

21

## The Stationary Distribution

- Where is the surfer at time  $t+1$ ?**
  - Follows a link uniformly at random
- Suppose the random walk reaches a state
  - then  $\mathbf{p}(t)$  is **stationary distribution** of a random walk
- Our original rank vector  $\mathbf{r}$  satisfies  $\mathbf{r} = \mathbf{M} \cdot \mathbf{r}$** 
  - So,  $\mathbf{r}$  is a **stationary distribution** for the random walk

$$\mathbf{p}(t+1) = \mathbf{M} \cdot \mathbf{p}(t)$$

J. Leskovec, A. Rajaraman, J. Ullman, Mining of Massive Datasets, <http://www.mmds.org>

22

## 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  $t = 0$

J. Leskovec, A. Rajaraman, J. Ullman, Mining of Massive Datasets, <http://www.mmds.org>

23

## PageRank: The Google Formulation

## 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?

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

25

## Does this converge?



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

### Example:

$$\begin{matrix} r_a \\ r_b \end{matrix} = \begin{matrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{matrix}$$

Iteration 0, 1, 2, ...

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

26

## Does it converge to what we want?



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

### Example:

$$\begin{matrix} r_a \\ r_b \end{matrix} = \begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{matrix}$$

Iteration 0, 1, 2, ...

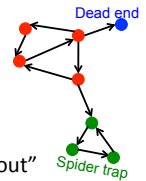
J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

27

## 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



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

28

## Problem: Spider Traps

### Power Iteration:

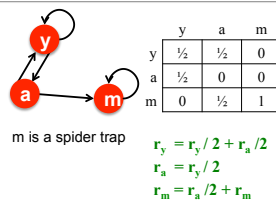
- Set  $r_i = 1/N$
- 1:  $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- 2:  $r = r'$
- If not converged: goto 1

### Example:

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

Iteration 0, 1, 2, ...

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

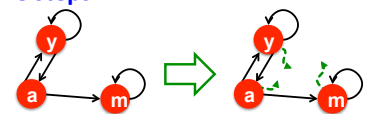


J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

29

## Solution: Random Teleports!

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



J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

30

## Problem: Dead Ends

### Power Iteration:

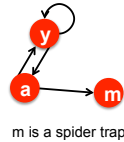
- Set  $r_i = 1/N$
- 1:  $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
- 2:  $r = r'$
- If not converged: goto 1

### Example:

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

Iteration 0, 1, 2.

Here the PageRank "leaks" out since the matrix is not stochastic.

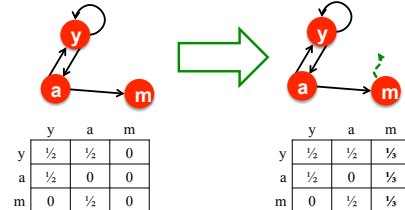


	y	a	m
y	1/2	1/2	0
a	1/2	0	0
m	0	1/2	0

$$\begin{aligned} r_y &= r_y / 2 + r_a / 2 \\ r_a &= r_y / 2 \\ r_m &= r_a / 2 \end{aligned}$$

## Solution: Always Teleport!

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



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

	y	a	m
y	1/2	1/2	1/3
a	1/2	0	1/3
m	0	1/2	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

## Solution: Random Teleports

- Google's solution that does it all:**
  - At each step, random surfer has two options:
    - With probability  $\beta$ , follow a link at random
    - With probability  $1-\beta$ , jump to some random page
- PageRank equation [Brin-Page, 98]

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{n} \quad \dots \text{out-degree of node } i$$

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]

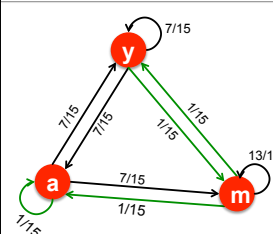
$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1 - \beta) \frac{1}{n}$$

- The Google Matrix A:**

$$A = \beta M + (1 - \beta) \frac{1}{n} \mathbf{e} \cdot \mathbf{e}^T \quad \mathbf{e} \text{ vector of all 1s}$$

- We have a recursive problem:  $r = A \cdot r$ 
  - Power iteration still works!
- What is  $\beta$ ?**
  - In practice  $\beta=0.8, 0.9$  (make 5 steps and jump)

## Random Teleports ( $\beta = 0.8$ )



$$0.8 \begin{pmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{pmatrix} + 0.2 \begin{pmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{pmatrix}$$

	y	a	m
y	7/15	7/15	1/15
a	7/15	1/15	1/15
m	1/15	7/15	13/15

**A**

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

## How do we actually compute the PageRank?

## 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  $A$ ,  $r^{\text{old}}$ ,  $r^{\text{new}}$

- 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!

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

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

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

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

28

## 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  $\beta/d_i$
  - Equivalent:** Tax each page a fraction  $(1-\beta)$  of its score and redistribute evenly

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

29

## Rearranging the Equation

- $r = A \cdot r$ , where  $A_{ij} = \beta M_{ij} + \frac{1-\beta}{N}$
- $r_i = \sum_{j=1}^N A_{ij} \cdot r_j$
- $$r_i = \sum_{j=1}^N \left[ \beta M_{ij} + \frac{1-\beta}{N} \right] \cdot r_j$$

$$= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \sum_{j=1}^N \frac{1-\beta}{N} r_j$$

$$= \sum_{j=1}^N \beta M_{ij} \cdot r_j + \frac{1-\beta}{N} \text{ since } \sum r_j = 1$$
- So we get:  $r = \beta M \cdot r + \left[ \frac{1-\beta}{N} \right]_N$

Note: Here we assumed  $M$  has no dead-ends.

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

J. Leskovec, A. Rajaraman, J. Ullman (Stanford University) Mining of Massive Datasets

30

## Sparse Matrix Formulation

- We just rearranged the **PageRank equation**
  - where  $[(1-\beta)/N]_N$  is a vector with all  $N$  entries  $(1-\beta)/N$
- $M$  is a **sparse matrix!** (with no dead-ends)
  - 10 links per node, approx 10N entries
- So in each iteration, we need to:
  - Compute  $r^{\text{new}} = \beta M \cdot r^{\text{old}}$
  - Add a constant value  $(1-\beta)/N$  to each entry in  $r^{\text{new}}$ 
    - Note if  $M$  contains dead-ends then and we also have to renormalize  $r^{\text{new}}$  so that it sums to 1

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, <http://www.mmds.org>

31

## PageRank: The Complete Algorithm

- Input:** Graph  $G$  and parameter  $\beta$ 
  - Directed graph  $G$  with spider traps and dead ends
  - Parameter  $\beta$
- Output:** PageRank vector  $r$ 
  - Set:  $r_j^{(0)} = \frac{1}{N}$ ,  $t = 1$
  - do:
    - $\forall j: r_j^{(t)} = \sum_{i \rightarrow j} \beta \frac{r_i^{(t-1)}}{d_i}$
    - $r_j^{(t)} = 0$  if in-deg. of  $j$  is 0
    - Now re-insert the leaked PageRank:
      - $\forall j: r_j^{(t)} = r_j^{(t)} + \frac{1-S}{N}$  where:  $S = \sum_j r_j^{(t)}$
    - $t = t + 1$
  - while  $\sum_j |r_j^{(t)} - r_j^{(t-1)}| > \epsilon$

J. Leskovec, A. Rajaraman, J. Ullman (Stanford University) Mining of Massive Datasets

32