

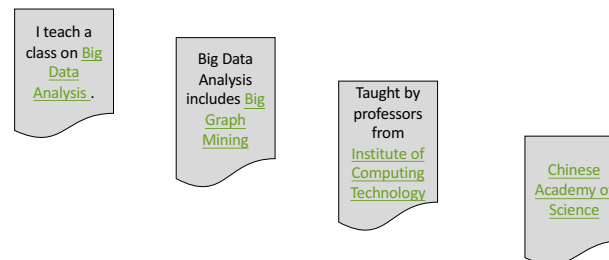
# 大数据分析

PageRank

刘盛华

## Web as a Graph

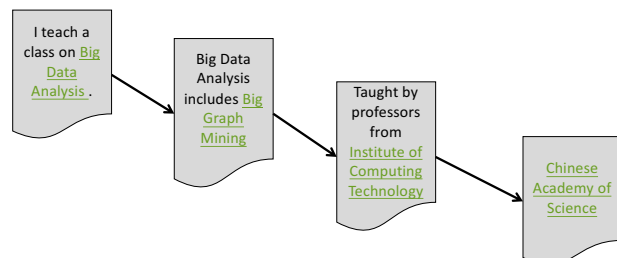
- Web as a directed graph: 有向图  
 结点: 网页  
 边: 超链接
- Nodes: Webpages
- Edges: Hyperlinks



2

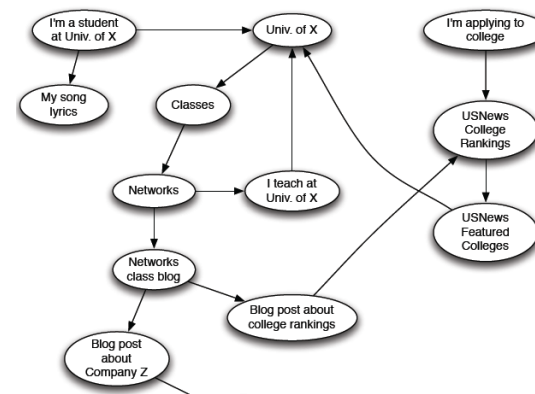
## Web as a Graph

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



3

## Web as a Directed Graph



4

## Broad Question

### How to organize the Web?

#### First try: **Human curated Web directories**

- Yahoo, DMOZ, LookSmart

#### 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.

信息检索调查:在一个小而可信的集合中查找相关文档

网络是巨大的, 充满了不可信的文件, 随机的东西, 网络垃圾邮件, 等等



## Web Search: 2 Challenges

### 2 challenges of web search:

#### (1) Web contains many sources of information Who to "trust"?

网络上有很多信息源, 谁来“信任”? 问  
窍门:值得信任的页面可能会相互指向对方

- Trick:** Trustworthy pages may point to each other!

#### (2) What is the "best" answer to query "newspaper"?

问“报纸”的“最佳”答案是什么?  
没有一个正确答案

- No single right answer 诀窍:真正了解报纸的页面可能都指向许多报纸

- Trick:** Pages that actually know about newspapers might all be pointing to many newspapers

## Ranking Nodes on the Graph

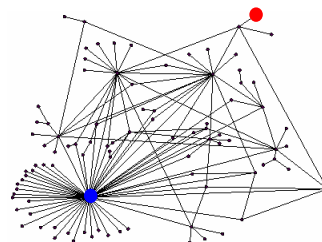
### All web pages are not equally "important"

[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!

网络图节点的连接具有很大的多样性。让我们按照链接结构对页面进行排序!



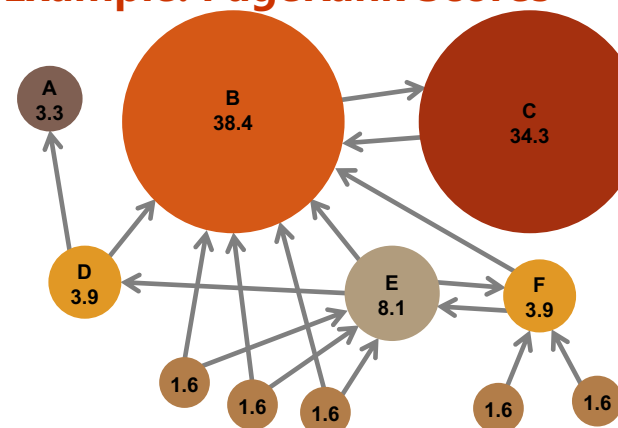
## PageRank: The "Flow" Formulation

## Links as Votes

- **Idea: Links as votes** 连接当作投票
  - 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 are equal?**
  - Links from important pages count more 重要页面的链接值更大
  - Recursive question! 递归问题

9

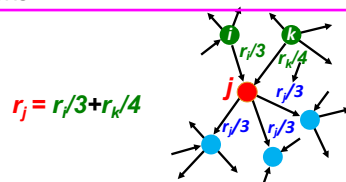
## Example: PageRank Scores



10

## 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



11

## PageRank: The "Flow" Model

来自重要网页的投票更值钱,

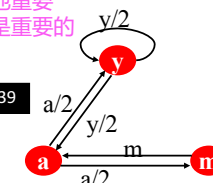
- A "vote" from an important page is worth more
- A page is important if it is pointed to by other important pages 如果一个网页被其他重要网页链接, 那他也是重要的
- Define a "rank"  $r_j$  for page  $j$

可能考矩阵形式

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

 $d_i$  ... out-degree of node  $i$ 

The web in 1839



"Flow" equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

12

## Solving the Flow Equations

- 3 equations, 3 unknowns, no constants  
3方程, 3未知数, 无常数, 无唯一解  
Flow equations:  

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$
- No unique solution
- All solutions equivalent modulo the scale factor  
所有解都等于尺度因子取模
- Additional constraint forces uniqueness: 附加约束强制得到唯一性  

$$r_y + r_a + r_m = 1$$
- 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  
高斯消去法适用于小例子, 但是我们需要一个更好的方法来处理大型网络图。
- We need a new formulation!  
我们需要一个新的公式

13

## PageRank: Matrix Formulation 矩阵公式

- Stochastic adjacency matrix  $M$   
随机邻接矩阵  
  - Let page  $i$  has  $d_i$  out-links  
外链这里列指向行
  - If  $i \rightarrow j$ , then  $M_{ji} = \frac{1}{d_i}$  else  $M_{ji} = 0$ 
    - $M$  is a column stochastic matrix  
列随机矩阵
    - Columns sum to 1  
列和为1
- Rank vector  $r$ : vector with an entry per page  
每页一个元素  
  - $r_i$  is the importance score of page  $i$
  - $\sum_i r_i = 1$  这里一定考
- The flow equations can be written  

$$r = M \cdot r$$

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

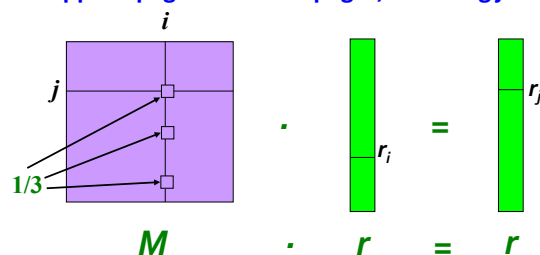
新的链接等于M乘旧的链接

14

## 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  $j$



15

## Eigenvector Formulation

- The flow equations can be written  

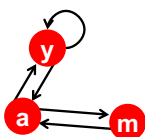
$$r = M \cdot r$$

特征值是1的特征向量
- 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  
相对于特征值1的主特征向量
- Largest eigenvalue of  $M$  is 1 since  $M$  is column stochastic (with non-negative entries)  
M最大的特征值就是1, 因为M时列随机的(非负元素)
- We know  $r$  is unit length and each column of  $M$  sums to one, so  $Mr \leq 1$
- We can now efficiently solve for  $r$ !  
The method is called Power iteration  
我们现在可以有效地解出r! 这种方法称为幂次迭代

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

16

## Example: Flow Equations & M



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

$$r = M \cdot r$$

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

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

17

## 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:  $r^{(0)} = [1/N, \dots, 1/N]^T$

- Iterate:  $r^{(t+1)} = M \cdot r^{(t)}$

- Stop when  $|r^{(t+1)} - r^{(t)}|_1 < \epsilon$

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

$d_i \dots$  out-degree of node  $i$

$|x|_1 = \sum_{1 \leq i \leq N} |x_i|$  is the  $L_1$  norm  
Can use any other vector norm, e.g., Euclidean

18

## PageRank: How to solve?

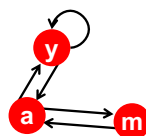
- Power Iteration:**

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

- Example:**

$$\begin{pmatrix} r_y \\ r_a \\ r_m \end{pmatrix} = \begin{pmatrix} 1/3 \\ 1/3 \\ 1/3 \end{pmatrix}$$

Iteration 0, 1, 2, ...



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

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

19

## PageRank: How to solve?

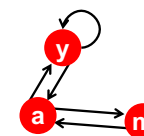
- Power Iteration:**

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

- Example:**

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

Iteration 0, 1, 2, ...



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

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

20

## Why Power Iteration works? (1) Details!

- **Power iteration:** 求主特征向量(对应最大特征值的向量)的方法  
A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

$$\begin{aligned} \square \mathbf{r}^{(1)} &= M \cdot \mathbf{r}^{(0)} \\ \square \mathbf{r}^{(2)} &= M \cdot \mathbf{r}^{(1)} = M(M\mathbf{r}^{(0)}) = M^2 \cdot \mathbf{r}^{(0)} \\ \square \mathbf{r}^{(3)} &= M \cdot \mathbf{r}^{(2)} = M(M^2\mathbf{r}^{(0)}) = M^3 \cdot \mathbf{r}^{(0)} \end{aligned}$$

- **Claim:**  
Sequence  $M \cdot \mathbf{r}^{(0)}, M^2 \cdot \mathbf{r}^{(0)}, \dots, M^k \cdot \mathbf{r}^{(0)}, \dots$  approaches the dominant eigenvector of  $M$

21

## Why Power Iteration works? (2) Details!

- **Claim:** Sequence  $M \cdot \mathbf{r}^{(0)}, M^2 \cdot \mathbf{r}^{(0)}, \dots, M^k \cdot \mathbf{r}^{(0)}, \dots$  approaches the dominant eigenvector of  $M$

- **Proof:** 接近M的主导特征向量

- Assume  $M$  has  $n$  linearly independent eigenvectors,  $x_1, x_2, \dots, x_n$  with corresponding eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$ , where  $\lambda_1 > \lambda_2 > \dots > \lambda_n$
- Vectors  $x_1, x_2, \dots, x_n$  form a basis and thus we can write:  
 $\mathbf{r}^{(0)} = c_1 x_1 + c_2 x_2 + \dots + c_n x_n$
- $M\mathbf{r}^{(0)} = M(c_1 x_1 + c_2 x_2 + \dots + c_n x_n)$
- $= c_1(Mx_1) + c_2(Mx_2) + \dots + c_n(Mx_n)$
- $= c_1(\lambda_1 x_1) + c_2(\lambda_2 x_2) + \dots + c_n(\lambda_n x_n)$
- Repeated multiplication on both sides produces
- $M^k \mathbf{r}^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \dots + c_n(\lambda_n^k x_n)$

22

## Why Power Iteration works? (3) Details!

- **Claim:** Sequence  $M \cdot \mathbf{r}^{(0)}, M^2 \cdot \mathbf{r}^{(0)}, \dots, M^k \cdot \mathbf{r}^{(0)}, \dots$  approaches the dominant eigenvector of  $M$

- **Proof (continued):**

- Repeated multiplication on both sides produces  
 $M^k \mathbf{r}^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \dots + c_n(\lambda_n^k x_n)$

$$\square M^k \mathbf{r}^{(0)} = \lambda_1^k \left[ c_1 x_1 + c_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k x_2 + \dots + c_n \left( \frac{\lambda_n}{\lambda_1} \right)^k x_n \right]$$

- Since  $\lambda_1 > \lambda_2$  then fractions  $\frac{\lambda_2}{\lambda_1}, \frac{\lambda_3}{\lambda_1}, \dots < 1$   
and so  $\left( \frac{\lambda_i}{\lambda_1} \right)^k = 0$  as  $k \rightarrow \infty$  (for all  $i = 2 \dots n$ ).

- Thus:  $M^k \mathbf{r}^{(0)} \approx c_1(\lambda_1^k x_1)$

- Note if  $c_1 = 0$  then the method won't converge 初始的 $r_0$ , 第一个向量不能为0??? ? ? ?

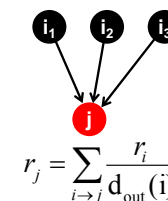
23

## Random Walk Interpretation 随机游走解释

想象一个随机的网络冲浪者

- **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 从开始在某页结束, 这个过程无限重复



- **Let:**

- $p(t) \dots$  vector whose  $i^{\text{th}}$  coordinate is the prob. that the surfer is at page  $i$  at time  $t$
- So,  $p(t)$  is a probability distribution over pages  
P是页面的概率分布

24

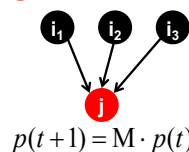
## The Stationary Distribution

什么时候达到稳态

### 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

25

## Existence and Uniqueness

存在性和唯一性

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

随机游动理论(又称马尔可夫过程)的一个中心结果是:

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  $t = 0$

对于满足一定条件的图, 其平稳分布是唯一的, 且无论  $t = 0$  时的初始概率分布如何, 最终都会达到平稳分布

ref to:

Perron–Frobenius theorem [Nonnegative Matrix, irreducible (connected), primitivity (**k-connected**)]

不管初始是什么值, 最终都会达到稳定解

A自乘k次, 那么所有的点都是连接起来了, 这样第一大特征值全是正的

26

## 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?

这是收敛?  
是否收敛于想要的结果?  
结果合理吗?

28

## Does this converge?



### ■ Example:

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

Iteration 0, 1, 2, ...

29

## Does it converge to what we want?



### ■ Example:

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

Iteration 0, 1, 2, ...

30

## PageRank: Problems

### 2 problems:

#### ■ (1) Some pages are **dead ends** (have no out-links)

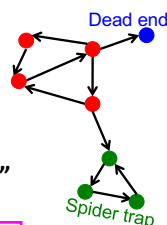
- Random walk has “nowhere” to go to
- Such pages cause importance to “leak out”

重要性消失

#### ■ (2) **Spider traps**: (all out-links are within the group)

- Random walk gets “stuck” in a trap
- And eventually spider traps absorb all importance

group里的点重要性都很大



31

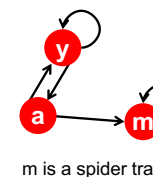
## Problem: Spider Traps

### ■ Power Iteration:

- Set  $r_j = 1$

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

- And iterate



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

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

$$r_a = r_y/2$$

$$r_m = r_a/2 + r_m$$

### ■ Example:

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

Iteration 0, 1, 2, ...

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

32

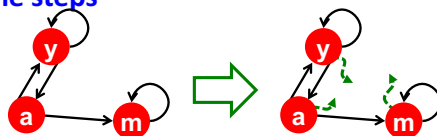


## Solution: Teleports! 瞬移

- 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



33

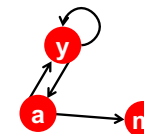
## Problem: Dead Ends

- **Power Iteration:**

- Set  $r_j = 1$

- $r_j = \sum_i \frac{r_i}{d_i}$

- And iterate



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

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

$$r_a = r_y/2$$

$$r_m = r_a/2$$

- **Example:**

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

Iteration 0, 1, 2, ...

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

34

## 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

把一列全为0的换成均值1/d

35

## Why Teleports Solve the Problem?

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

蜘蛛陷阱，不是问题，只不过不是我们想要的结果

- Spider-traps are not a problem (**converge**), 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 不收敛

- The matrix is not column stochastic (**zero column**) so our initial assumptions are not met

- **Solution:** Make matrix column stochastic by **always teleporting** when there is nowhere else to go

36

## 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]

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

$d_i \dots$  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 (add  $1/N$  in  $M$ ) or explicitly follow random teleport links with probability  $1.0$  from dead-ends ( $\beta=0$ ).

这个公式假定M没有dead ends。我们可以预处理矩阵M删除所有dead ends(加1/N)或显式地遵循概率1.0随机瞬移( $\beta=0$ )。

37

## 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 $A$ :

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

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

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

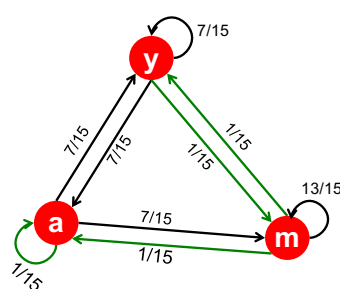
And the Power method still works!

### What is $\beta$ ?

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

38

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



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

$M$                        $[1/N]_{N \times N}$

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

$A$

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

39

# 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!

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

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

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

41

## 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

随机传送等于:

添加一个从  $i$  到其他页面的传送链接, 并设置转换概率为  $(1-\beta)/N$

降低每个外链的概率, 从  $1/d_i$  到  $\beta/d_i$

相当于: 对每一页征收得分除以  $(1-\beta)$  的税, 然后重新平均分配

42

## Rearranging the Equation

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

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

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

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

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

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

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

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

43

## Sparse Matrix Formulation

- We just rearranged the **PageRank equation**

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

- 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  $\sum_j r_j^{\text{new}} < 1$  and we also have to renormalize  $r^{\text{new}}$  so that it sums to 1

44

## PageRank: The Complete Algorithm

- **Input:** Graph  $G$  and parameter  $\beta$ 
  - Directed graph  $G$  (can have spider traps and dead ends)
  - Parameter  $\beta$

- **Output:** PageRank vector  $r^{new}$

- Set:  $r_j^{old} = \frac{1}{N}$
- repeat until convergence:  $\sum_j |r_j^{new} - r_j^{old}| > \epsilon$ 
  - $\forall j: r_j^{new} = \sum_{i \rightarrow j} \beta \frac{r_i^{old}}{d_i}$   
 $r_j^{new} = 0$  if in-degree of  $j$  is 0
  - Now re-insert the leaked PageRank:  
 $\forall j: r_j^{new} = r_j^{new} + \frac{1-S}{N}$
  - $r^{old} = r^{new}$  where:  $S = \sum_j r_j^{new}$

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  $S$ .

如果图没有死角则1-β泄露的PageRank的数量。但由于我们有死胡同，泄露的PageRank的数量可能会更大。我们必须通过计算S来明确地解释它

## Sparse Matrix Encoding

- Encode sparse matrix using only nonzero entries
  - Space proportional roughly to number of links
  - Say 10N, or  $4 \times 10^1$  billion = 40GB
  - Still won't fit in memory, but will fit on disk

source node	degree	目的节点 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  $M$  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  $j = 1 \dots d_i$   
 $r^{new}(dest_j) += \beta r^{old}(i) / d_i$

	$r^{new}$	source	degree	destination	$r^{old}$
0		0	3	1, 5, 6	0
1		1	4	17, 64, 113, 117	1
2		2	2	13, 23	2
3					3
4					4
5					5
6					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?

如果内存装不下  $r_{new}$  怎么办？分块

## Block-based Update Algorithm

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

$r^{new}$	src	degree	destination	$r^{old}$
0	0	4	0, 1, 3, 5	0
1	1	2	0, 5	1
2	2	2	3, 4	2
3				3
4				4
5				5

$M$

49

## Analysis of Block Update

- Similar to nested-loop join in databases 类似于数据库中的嵌套循环连接
  - 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  
 $M$ 很大, 应该避免读 $M$ 次数为 $k$

50

## Block-Stripe Update Algorithm

先根据 $R_{new}$ 把 $M$ 分好, 会重复存很多, 但是空间复杂度小了

$r^{new}$	src	degree	destination	$r^{old}$
0	0	4	0, 1	0
1	1	3	0	1
2	2	2	1	2
3	0	4	3	3
4	2	2	3	4
5	0	4	5	5
	1	3	5	
	2	2	4	

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

51

## Block-Stripe Analysis

- Break  $M$  into stripes  $M$ 分成条! 每个条带只包含 $r^{new}$ 相应块中的目标节点
  - Each stripe contains only destination nodes in the corresponding block of  $r^{new}$
- Some additional overhead per stripe 每个条纹 $q$ 有一些额外的开销, 但这通常是值得的
  - But it is usually worth it
- Cost per iteration of Power method:  
 $= |M|(1+\epsilon) + (k+1)|r|$

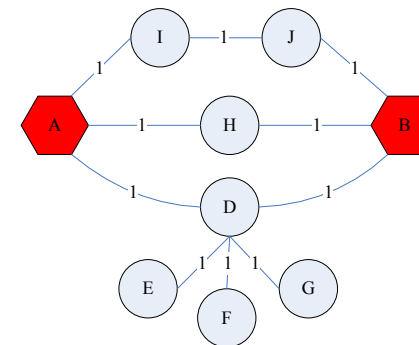
52

## Application to Measuring Proximity in Graphs

应用于测量接近度的图表  
带重启的随机游走: S是单个元素

**Random Walk with Restarts: S is a single element**

## Proximity on Graphs

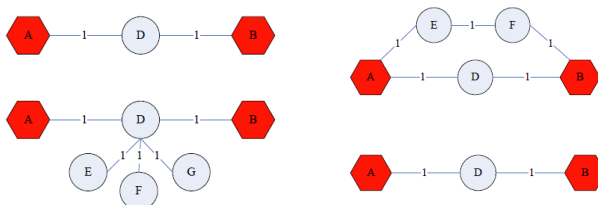


**a.k.a.: Relevance, Closeness, 'Similarity'...**

54

## Good proximity measure?

- Shortest path is not good:

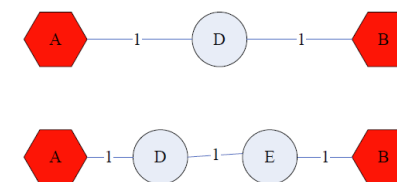


- No effect of degree-1 nodes (E, F, G)!
- Multi-faceted relationships

55

## Good proximity measure?

- Network flow is not good:

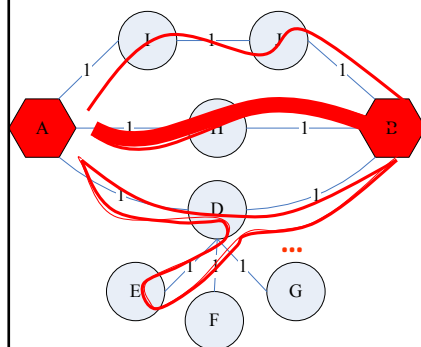


- Does not punish long paths 不惩罚长路吗

56

什么是好的相似的定义?

## What is good notion of proximity?



- Multiple connections
- Quality of connection
  - Direct & Indirect connections
  - Length, Degree, Weight...

连接质量  
直接连接和间接连接  
长度、度、重量...

57

## SimRank: Idea

- **SimRank**: Random walks from a fixed node on  $k$ -partite graphs

- **Setting**:  $k$ -partite graph with  $k$  types of nodes

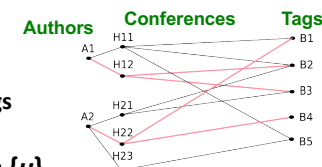
- E.g.: Authors, Conferences, Tags

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

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

- **Problem**:

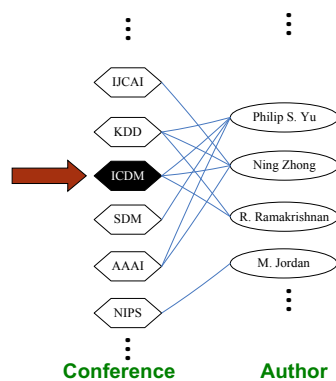
- Must be done once for each node  $u$  每个结点 $u$ 都得做一次
- Suitable for **sub-Web-scale** applications 不适应大的图



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

58

## SimRank: Example



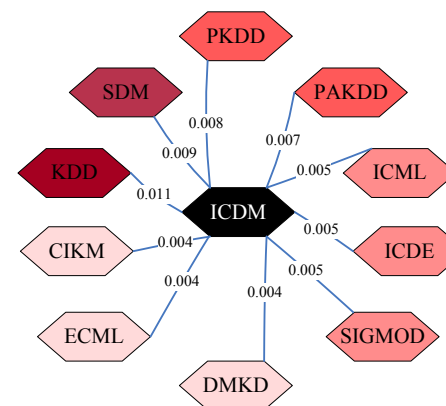
问:与ICDM最相关的会议是什么?  
A:具有传送集 $S=\{ICDM\}$ 的特定主题PageRank

**Q:** What is most related conference to ICDM?

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

59

## SimRank: Example



60

## PageRank: Summary

- **"Normal" PageRank:** 标准页面排序
  - Teleports uniformly at random to any node 均匀随机瞬移到任意结点
  - 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:** 特定主题的页面排序也称为个性化页面排序
  - Teleports to a topic specific set of pages 跳到特定主题页面的集合
  - 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:** 带重启的随机游走
  - Topic-Specific PageRank where teleport is always to the same node.  $S = [0, 0, 0, 0, 1, 0, 0, 0, 0, 0]$

61

# Questions?