

1.2.1 Solving the Flow Equations

□ 3 equations, 3 unknowns, no constants

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

Flow equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

□ Additional constraint forces uniqueness:

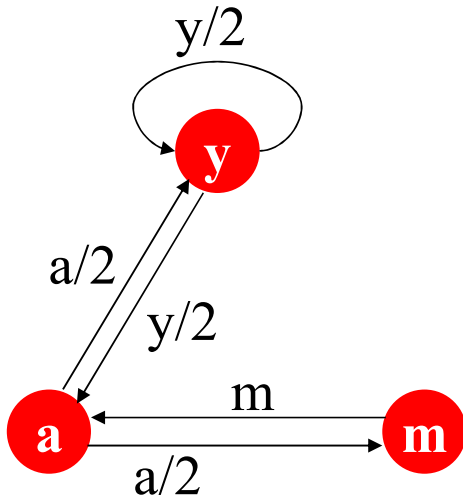
- $r_y + r_a + r_m = 1$
- **Solution:** $r_y = \frac{2}{5}, r_a = \frac{2}{5}, r_m = \frac{1}{5}$

□ Gaussian elimination method (高斯消元法/高斯消去法) works for small examples, but we need a better method for large web-size graphs ➡ **We need a new formulation!**

1.2.2 PageRank: Matrix Formulation

1、Stochastic adjacency matrix(邻接矩阵) 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
- **Note:** M also is called **Web transition matrix(Web转移矩阵)** in some literature



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

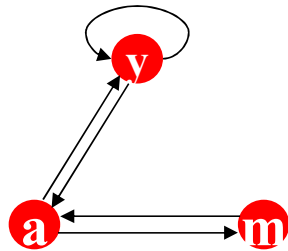
【备注】：

随机矩阵是这样一类方阵，其元素为非负实数，且行和或列和为1。
如果列和为1，则称为**列随机矩阵**；如果行和为1，则称为**行随机矩阵**。
如果行和和列和都为1，则称为**双随机矩阵**

1.2.2 PageRank: Matrix Formulation

□ **2、Rank vector(秩向量) r :** vector with an entry per page

- r_i is the importance score of page i
- Initial, each page has $1/n$ importance score, when total n pages.
- $\sum_i r_i = 1$



$$r = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

□ Now, the flow equations can be written

$$r = M \cdot r$$

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

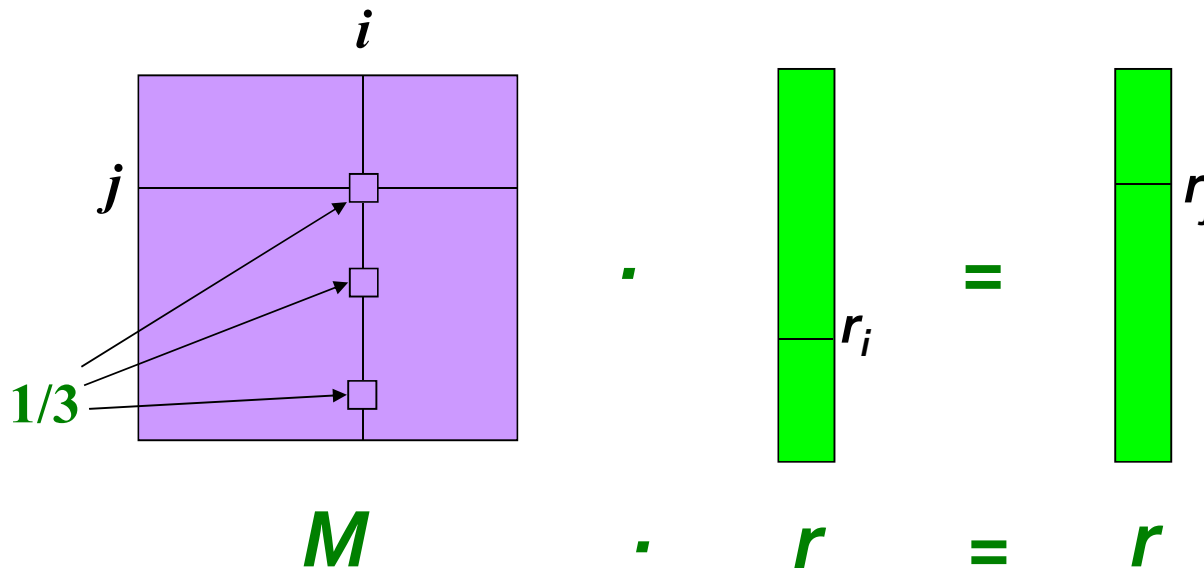
□ **M fixed. How to calculate r ?**

1.2.2 Example

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

Flow equation in the matrix form: $M \cdot r = r$

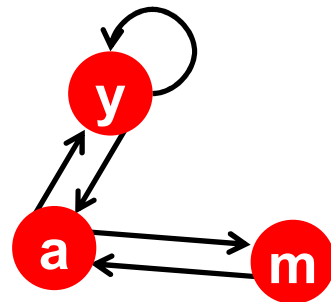
➤ Suppose page i links to 3 pages, including j



Both two forms mean the same thing!

1.2.2 Example

□ With the given example, we can give the “flow” model, and the matrix formulation:



(2) the matrix formulation

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

(1) the “flow” model

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

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

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

1.2.2 Eigenvector Formulation

□ The flow equations can be written

$$r = M \cdot r$$

【备注线性代数相关知识】： x is an eigenvector (特征向量) with the corresponding eigenvalue (特征值)

λ if: $Ax = \lambda x$

□ So the rank vector r is an eigenvector (特征向量) of the stochastic web matrix M

➤ In fact, its first or principal eigenvector (主特征向量, 最大特征值对应的特征向量) with corresponding eigenvalue 1 (最大特征值为1)

- Largest eigenvalue of M is 1, since M is column stochastic (columns sum to 1)

□ We can now efficiently solve for r ! The method is called **Power iteration** (幂迭代法)

1.2.2 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 $\|\mathbf{r}^{(t+1)} - \mathbf{r}^{(t)}\|_1 < \varepsilon$

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

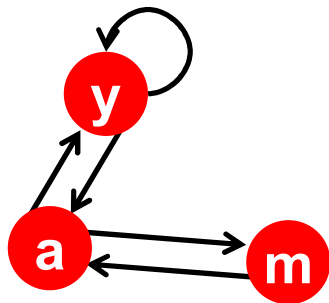
Can use any other vector norm, e.g., Euclidean

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

d_i out-degree of node i

1.2.2 PageRank: How to solve?

Example:



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

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

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

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

$$\begin{bmatrix} \mathbf{r}_y \\ \mathbf{r}_a \\ \mathbf{r}_m \end{bmatrix} = \begin{bmatrix} 1/3 \\ 1/3 \\ 1/3 \end{bmatrix}$$

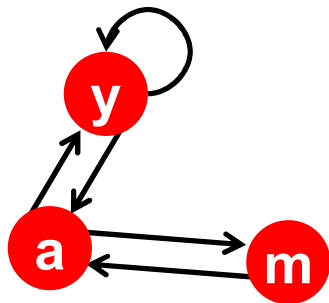
Iteration 0, 1, 2, ...

Power Iteration:

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

1.2.2 PageRank: How to solve?

Example:



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

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

Power Iteration:

➤ Set $r_j = 1/N$

➤ **1:** $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$

➤ **2:** $r = r'$

➤ Go to 1

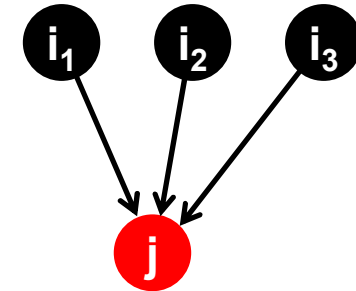
$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{matrix} 1/3 & 1/3 & 5/12 & 9/24 & & 6/15 \\ 1/3 & 3/6 & 1/3 & 11/24 & \dots & 6/15 \\ 1/3 & 1/6 & 3/12 & 1/6 & & 3/15 \end{matrix}$$

Iteration 0, 1, 2, ...

1.2.3 Random Walk Interpretation

Imagine a random web surfer(随机冲浪者):

- 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



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

Let:

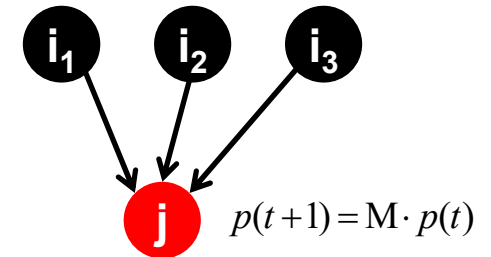
- $\mathbf{p}(t)$... vector whose i^{th} coordinate is the prob. that the surfer is at page i at time t
- So, $\mathbf{p}(t)$ is a probability distribution over pages(概率分布向量)

1.2.3 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 (稳态分布)** of a random walk

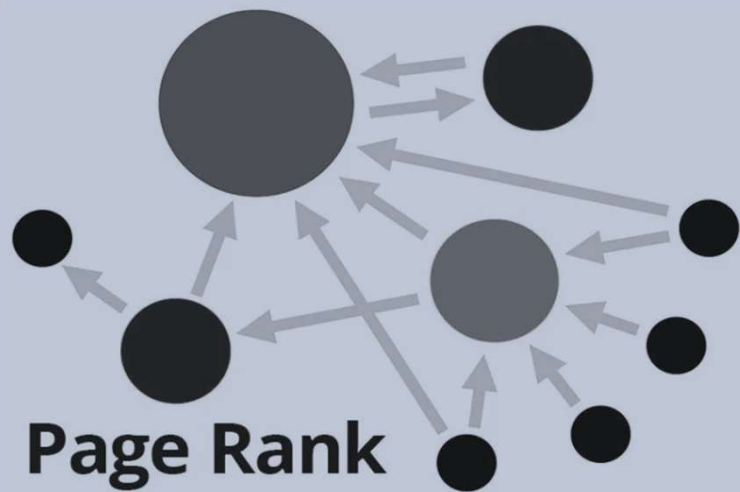
- ### □ Our original rank vector r satisfies $r = M \cdot r$

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

1.2.3 Existence and Uniqueness

- A central result from the theory of random walks (a.k.a. **Markov processes**, 马尔科夫过程):

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 $t = 0$



Section 1.3: PageRank, The Google Formulation

Content

1

PageRank Problems

2

Spider Traps → Teleports

3

Dead Ends → Teleports

4

Google Matrix

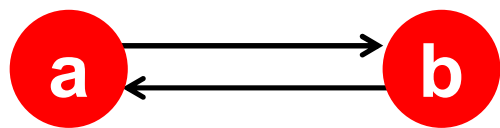
1.3.1 PageRank: Three Questions

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

- ❑ Does this converge(收敛)?
- ❑ Does it converge to what we want?
- ❑ Are results reasonable?

1.3.1 Does this converge?

□ Example:



	a	b
a	0	1
b	1	0

$$r = Mr$$

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

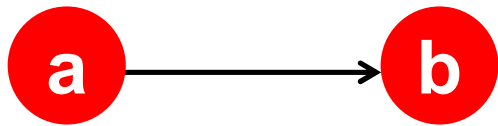
$$\begin{array}{rcl} r_a & = & 1 \quad 0 \quad 1 \quad 0 \\ r_b & = & 0 \quad 1 \quad 0 \quad 1 \end{array}$$

iteration 0, 1, 2, ...

无法收敛！ 这是一个典型的**蜘蛛陷阱**问题(后续马上会讲)

1.3.1 Does it converge to what we want?

□ Example:



	a	b
a	0	0
b	1	0

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

$$\begin{array}{lcl} r_a & = & 1 \quad 0 \quad 0 \quad 0 \\ r_b & & 0 \quad 1 \quad 0 \quad 0 \end{array}$$

iteration 0, 1, 2, ...

收敛了但不是我们想要的结果！ 这是一个典型的**死角**问题(后续马上会讲)

1.3.1 Are results reasonable?

□ 以下为同一个网页的部分内容截图示例:

美国绝密战争计划，居然这样泄露了

牛弹琴 03-25 08:22

摘要：

美国国家安全顾问沃尔兹在加密应用上讨论轰炸也门行动时，错将《大西洋月刊》主编戈德堡加入群聊，导致绝密战争计划泄露。更有趣的是，群聊还曝光了副总统万斯反对轰炸行动、高层对欧洲的厌恶以及美国要经济回报等不为人知的内幕。



最热文章

1. 特朗普关税立场软化迹象力挺美股、打压美债，纳指涨超2%，特斯拉涨近12%
2. 华尔街见闻早餐FM-Radio | 2025年3月25日
3. 特朗普：未来几天宣布汽车、木材和芯片关税，可能对多国减免关税征收
4. 美国绝密战争计划，居然这样泄露了
5. 报道：小米集团拟以先旧后新方式配售股份，融资最多53亿美元

华尔街见闻 关于我们 广告投放 版权与商务合作 联系方式 意见反馈	法律信息 版权声明 用户协议 付费内容订阅协议 商品购买协议 隐私政策
声明 未经许可，任何人不得复制、转载、或以其他方式使用本网站的内容。 评论前请阅读网站“ 跟帖评论自律管理承诺书 ”	违法和不良信息举报 举报电话：021-60675200（周一到周五9:30-11:30，13:00-18:30） 举报邮箱：contact@wallstreetcn.com 网站举报： 点击这里 违法和不良信息举报受理和处置管理办法 清朗·财经违规内容专项整治公告 涉企虚假不实信息举报专区 提交举报

每个链接被跳转的概率是一样么？

NO!

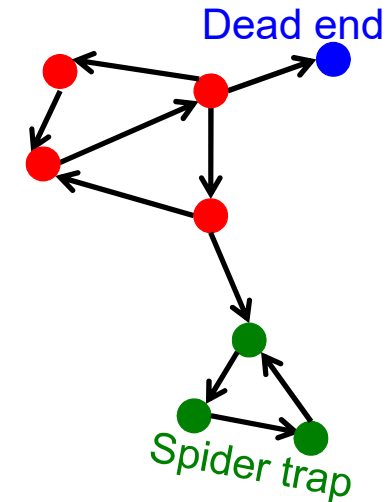
Random Surfer Model vs Reasonable Surfer Model

【备注】在我们后续的学习中还是采用基础的Random Surfer Model

1.3.1 PageRank: Problems

There maybe 2 problems:

- (1) Some pages are **dead ends** (死角问题, 死胡同问题, 或称终止点 have no out-links)
 - Random walk has “nowhere” to go to
 - Such pages cause importance to “leak out”

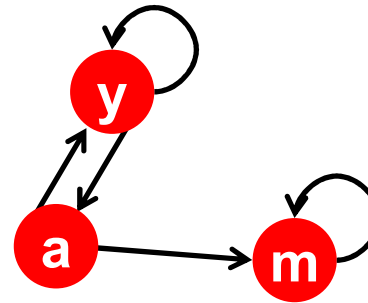


- (2) **Spider traps:** (蜘蛛陷阱问题, 或称采集器陷阱, all out-links are within the group)
 - Random walked gets “stuck” in a trap
 - And eventually spider traps absorb all importance

1.3.2 Problem: Spider Traps

Power Iteration:

- Set $r_j = 1/N$
- $r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$
 - 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:

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

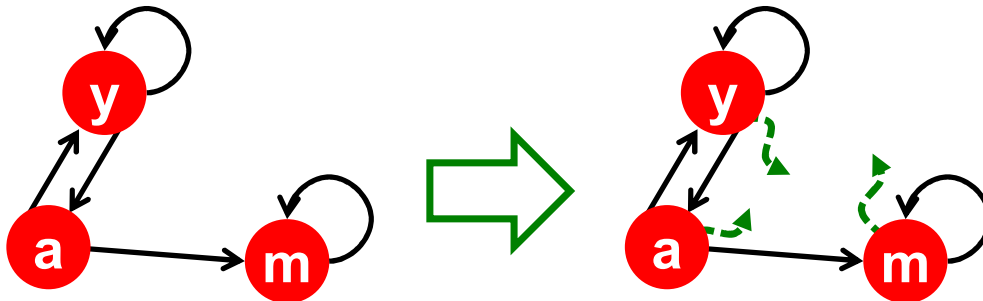
Iteration 0, 1, 2, ...

All the PageRank score gets "trapped" in node m.

1.3.2 Solution: Teleports(随机跳转)!

❑ The Google solution (called **taxation method**, **抽税法**) for spider traps: **At each time step, the random surfer has two options**

- With prob. β (也叫**阻尼系数**), **follow a link** at random. Note, common values for β are in the range 0.8 to 0.9
- With prob. $1-\beta$, jump to **some random page**

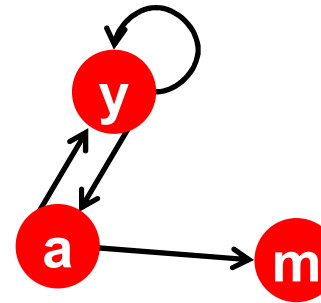


❑ Surfer will teleport out of spider trap within a few time steps

1.3.3 Problem: Dead Ends

Power Iteration:

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



m is a dead end

	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:

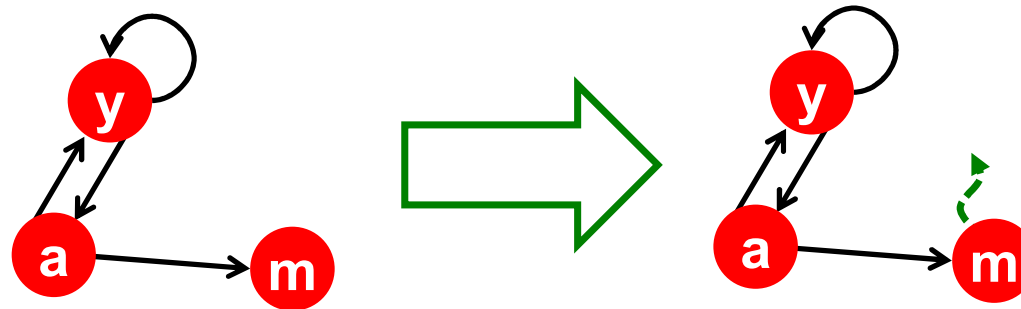
$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{bmatrix} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{bmatrix}$$

Iteration 0, 1, 2, ...

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

1.3.3 Solution: Always Teleport!

□ **Taxation (抽税法, or 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	$\frac{1}{2}$	$\frac{1}{3}$

1.3.3 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

1.3.4 Solution: Random Teleports

- Google's solution (Taxation, 抽税法) 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

【备注】 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.

1.3.4 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 :

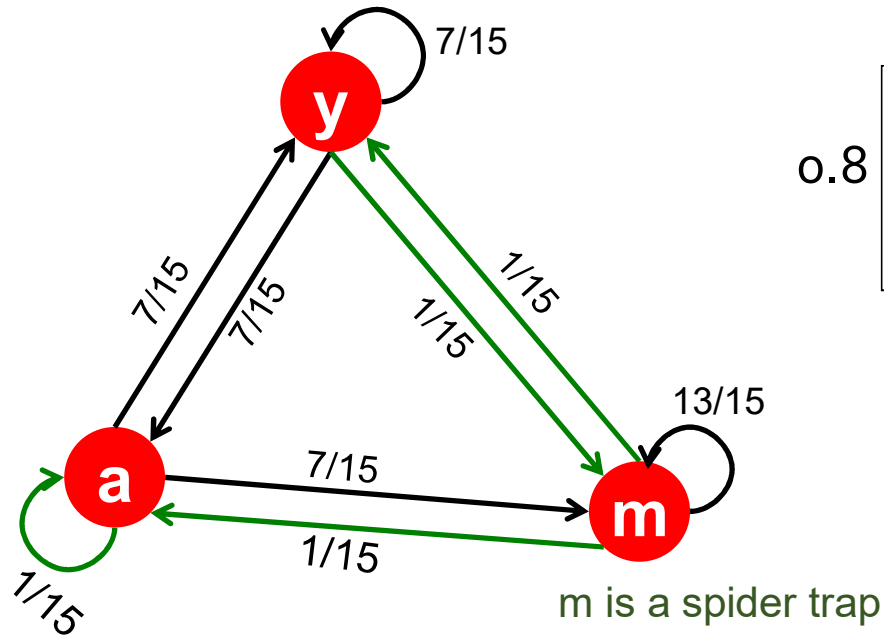
$$A = \beta M + (1 - \beta) \left[\frac{1}{N} \right]_{N \times N}$$

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

□ What is β ? In practice $\beta = 0.8, 0.9$ (make 5 steps on avg., jump)

□ We have a recursive problem: $\mathbf{r} = A \cdot \mathbf{r}$, and the Power iteration method still works!

1.3.4 Random Teleports ($\beta = 0.8$)

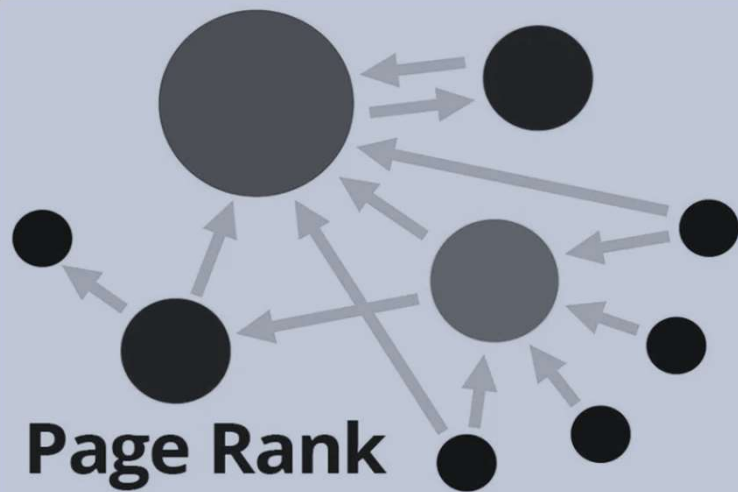


$$\begin{aligned}
 & \mathbf{M} \\
 & 0.8 \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 \\ 0 & 1/2 & 1 \end{bmatrix} + 0.2 \begin{bmatrix} 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/3 \end{bmatrix} \\
 & \mathbf{A} \\
 & \begin{matrix} y \\ a \\ m \end{matrix} \begin{bmatrix} 7/15 & 7/15 & 1/15 \\ 7/15 & 1/15 & 1/15 \\ 1/15 & 7/15 & 13/15 \end{bmatrix}
 \end{aligned}$$

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

$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{bmatrix} 1/3 & 0.33 & 0.24 & 0.26 & & 7/33 \\ 1/3 & 0.20 & 0.20 & 0.18 & \dots & 5/33 \\ 1/3 & 0.46 & 0.52 & 0.56 & & 21/33 \end{bmatrix}$$

Previously, all the PageRank score gets “trapped” in node m.
But now taxation solves it!



Section 1.4: Computing PageRank

Content

- 1 The Complete PageRank Algorithm
- 2 Basic Update Step
- 3 Block-based Update Step
- 4 Block-Stripe Update Step

1.4.1 Computing PageRank

□ Key step is matrix-vector multiplication

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

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

□ Say $N = 1$ billion (十亿) pages

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

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

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

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

【备注】观察矩阵 \mathbf{A} 和 \mathbf{M} 可知,
Matrix **A** dense matrix!
Matrix **M** sparse matrix