

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>



New Topic: Graph Data!

High dim. data

Locality
sensitive
hashing

Clustering

Dimensional
ity
reduction

Graph data

PageRank,
SimRank

Community
Detection

Spam
Detection

Infinite data

Filtering
data
streams

Web
advertising

Queries on
streams

Machine learning

SVM

Decision
Trees

Perceptron,
kNN

Apps

Recommen
der systems

Association
Rules

Duplicate
document
detection

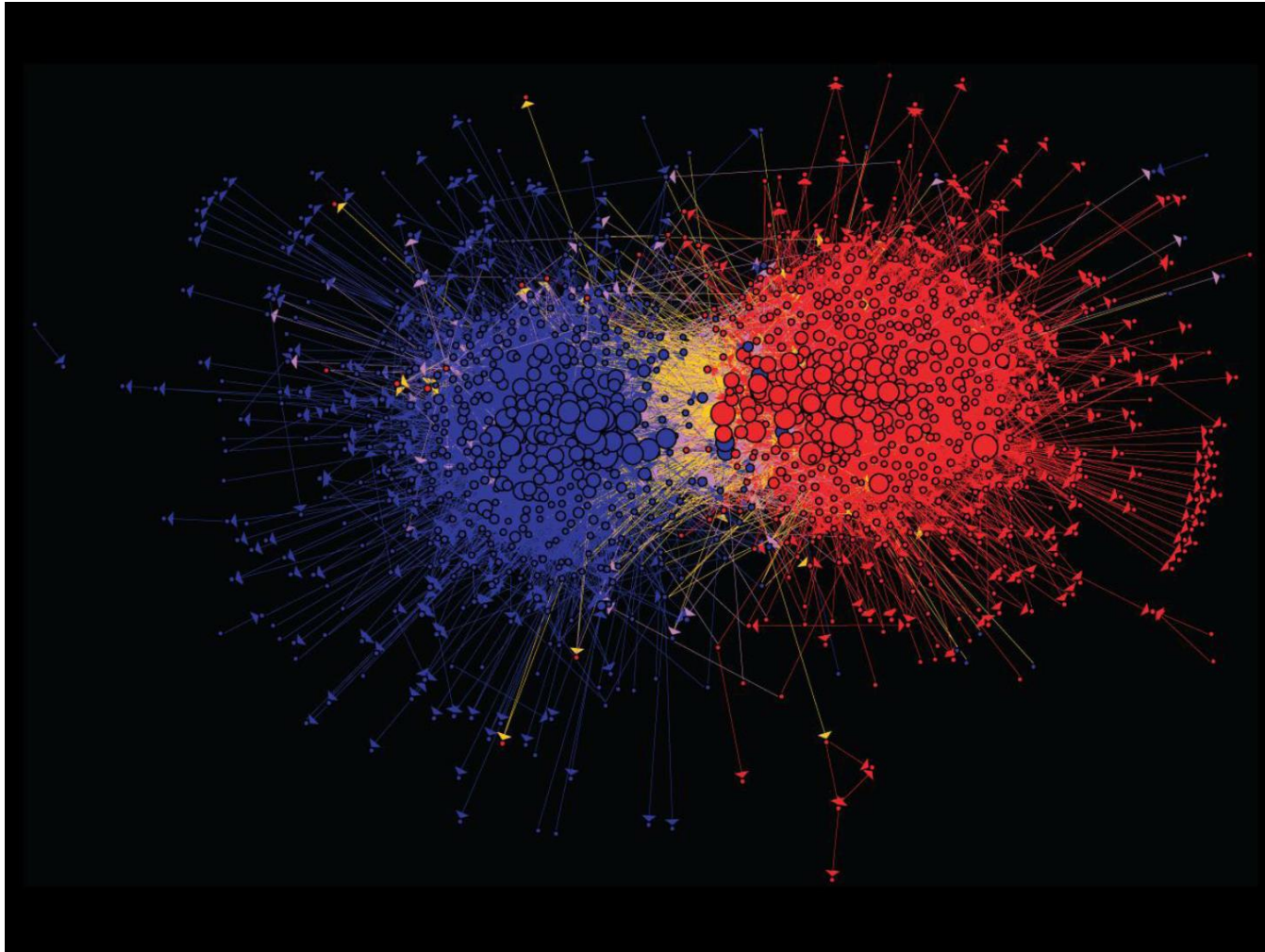
Graph Data: Social Networks



Facebook social graph

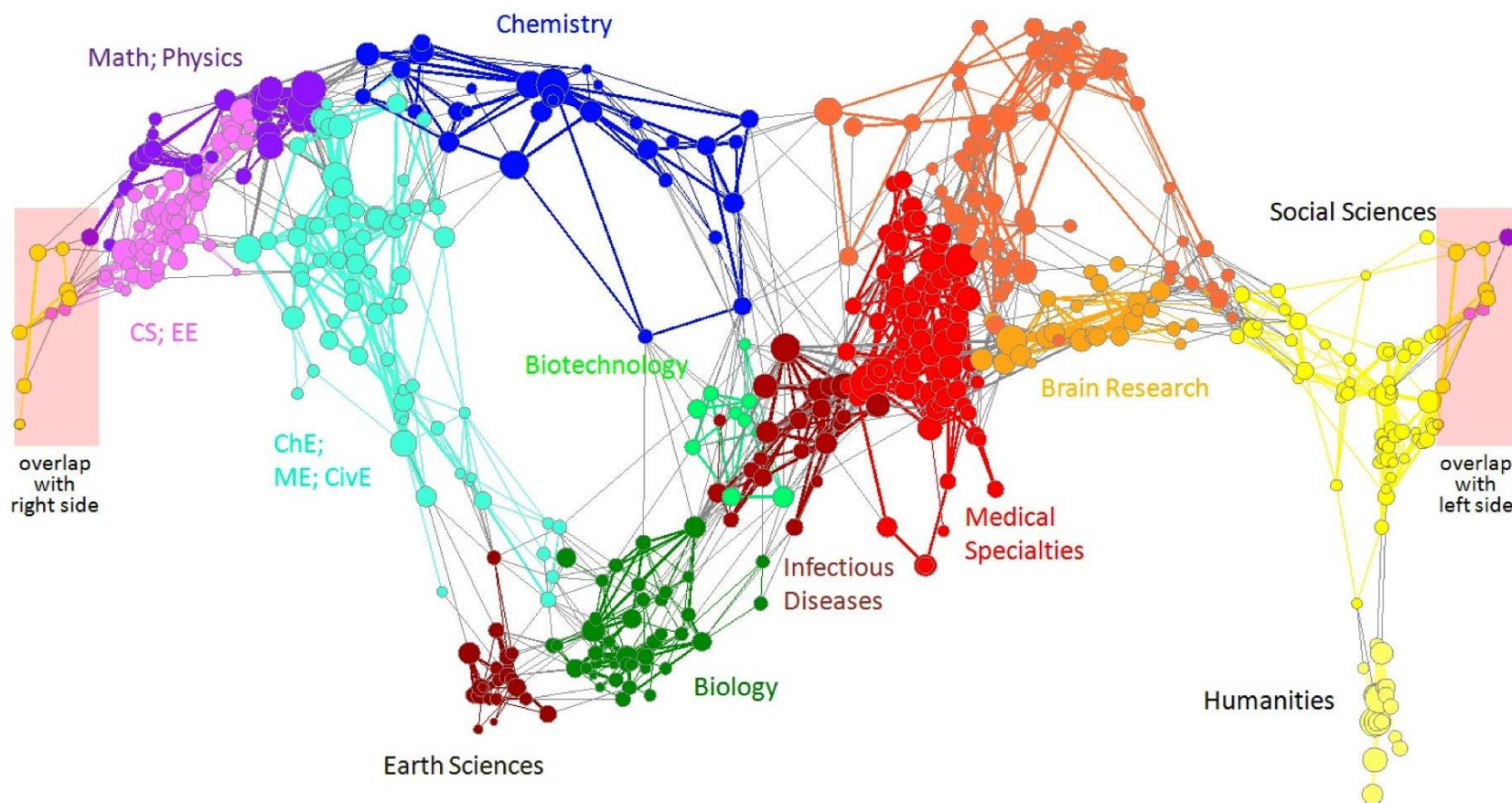
4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

Graph Data: Media Networks



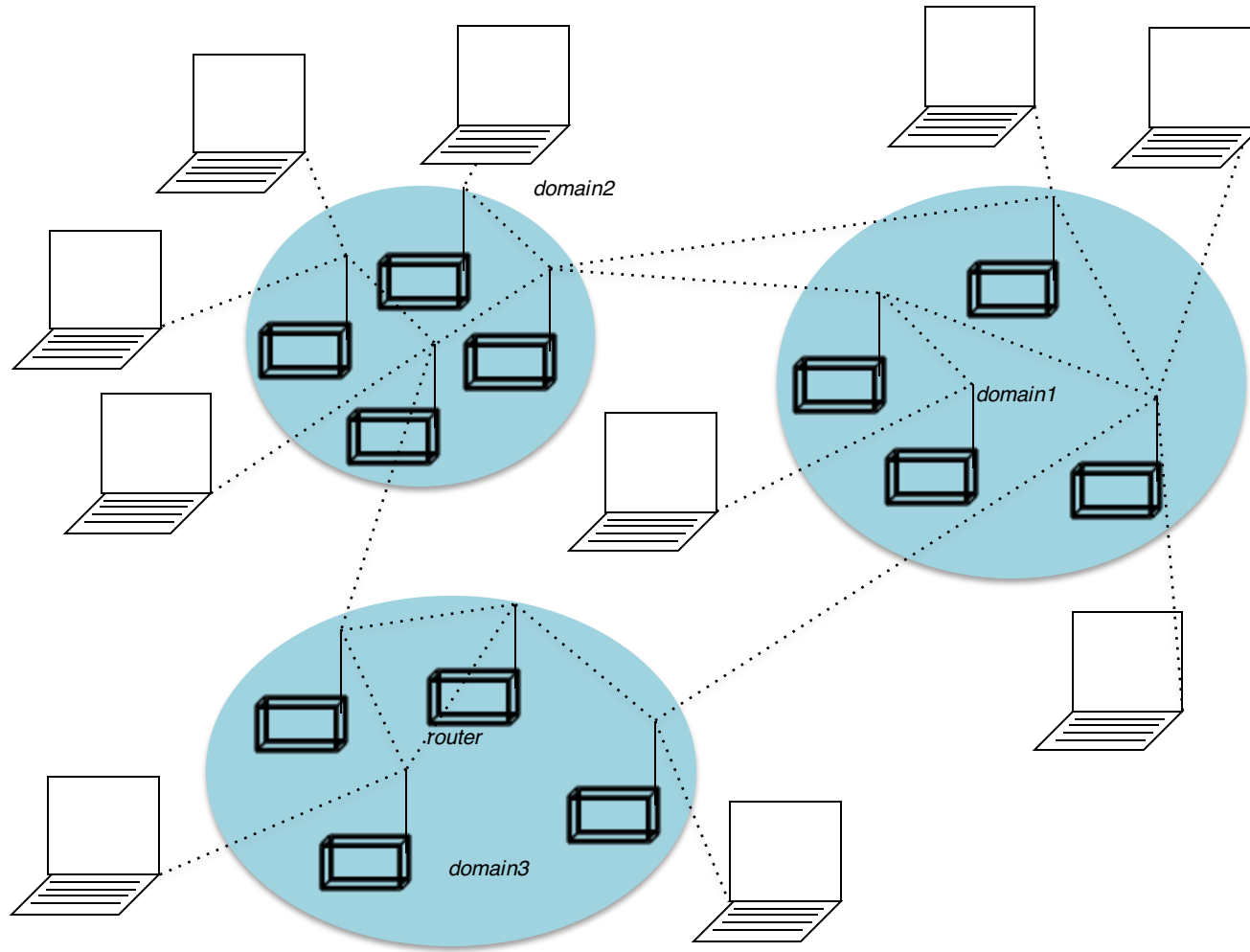
Connections between political blogs
Polarization of the network [Adamic-Glance, 2005]

Graph Data: Information Nets



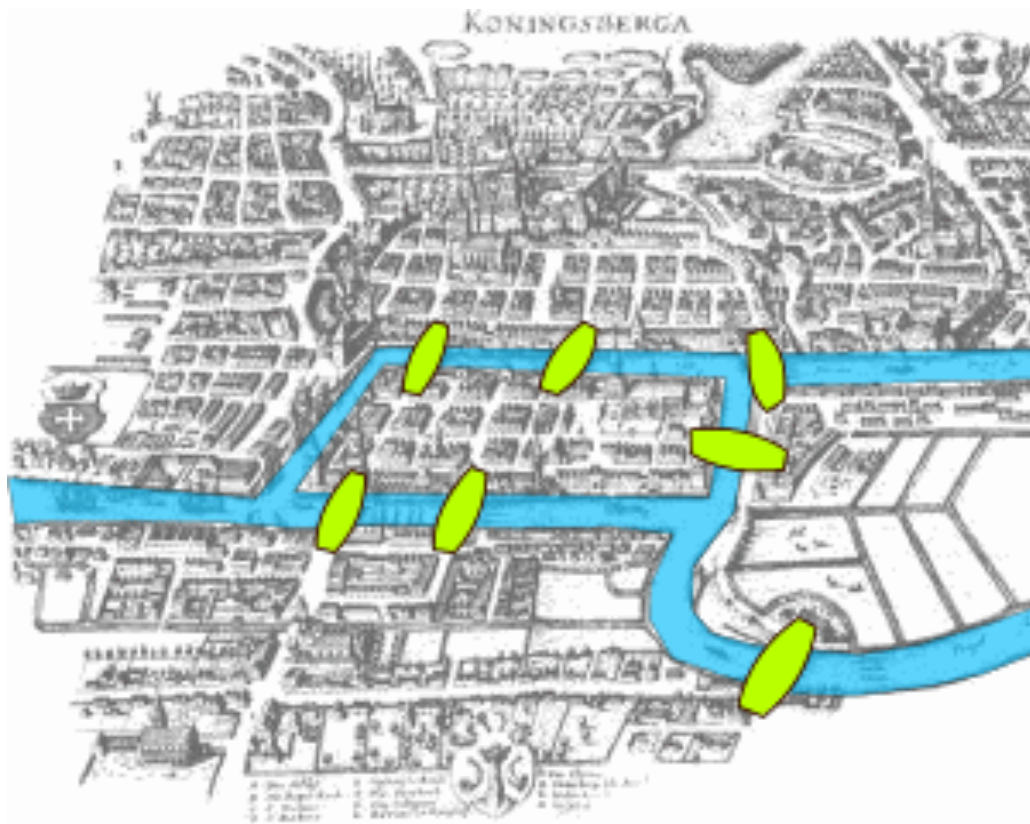
Citation networks and Maps of science
[Börner et al., 2012]

Graph Data: Communication Nets



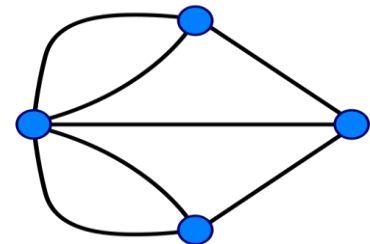
Internet

Graph Data: Technological Networks



Seven Bridges of Königsberg [Euler, 1735]

Return to the starting point by traveling each link of the graph once and only once.



Web as a Graph

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

I teach a
class on
Networks.

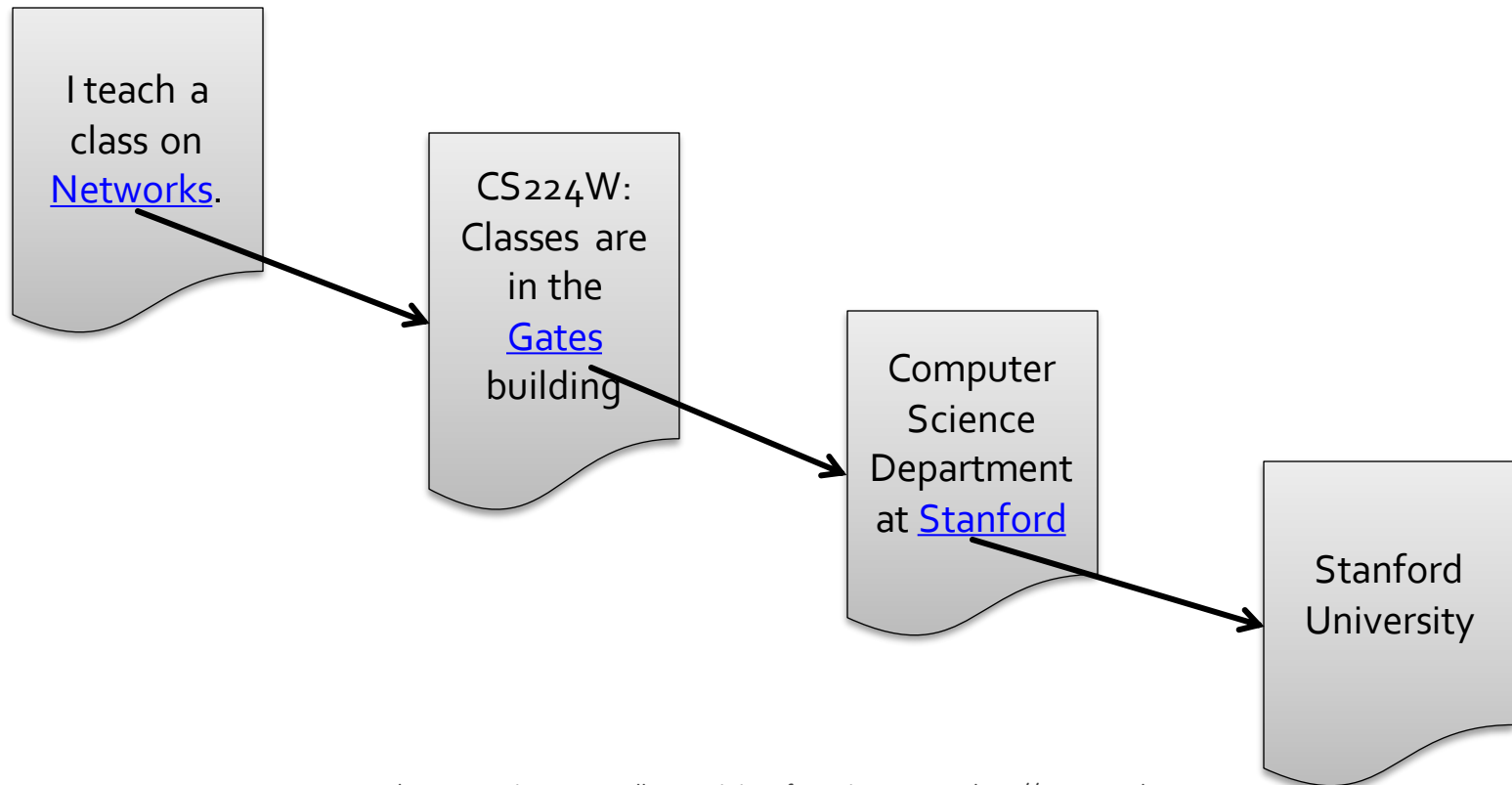
CS224W:
Classes are
in the
Gates
building

Computer
Science
Department
at Stanford

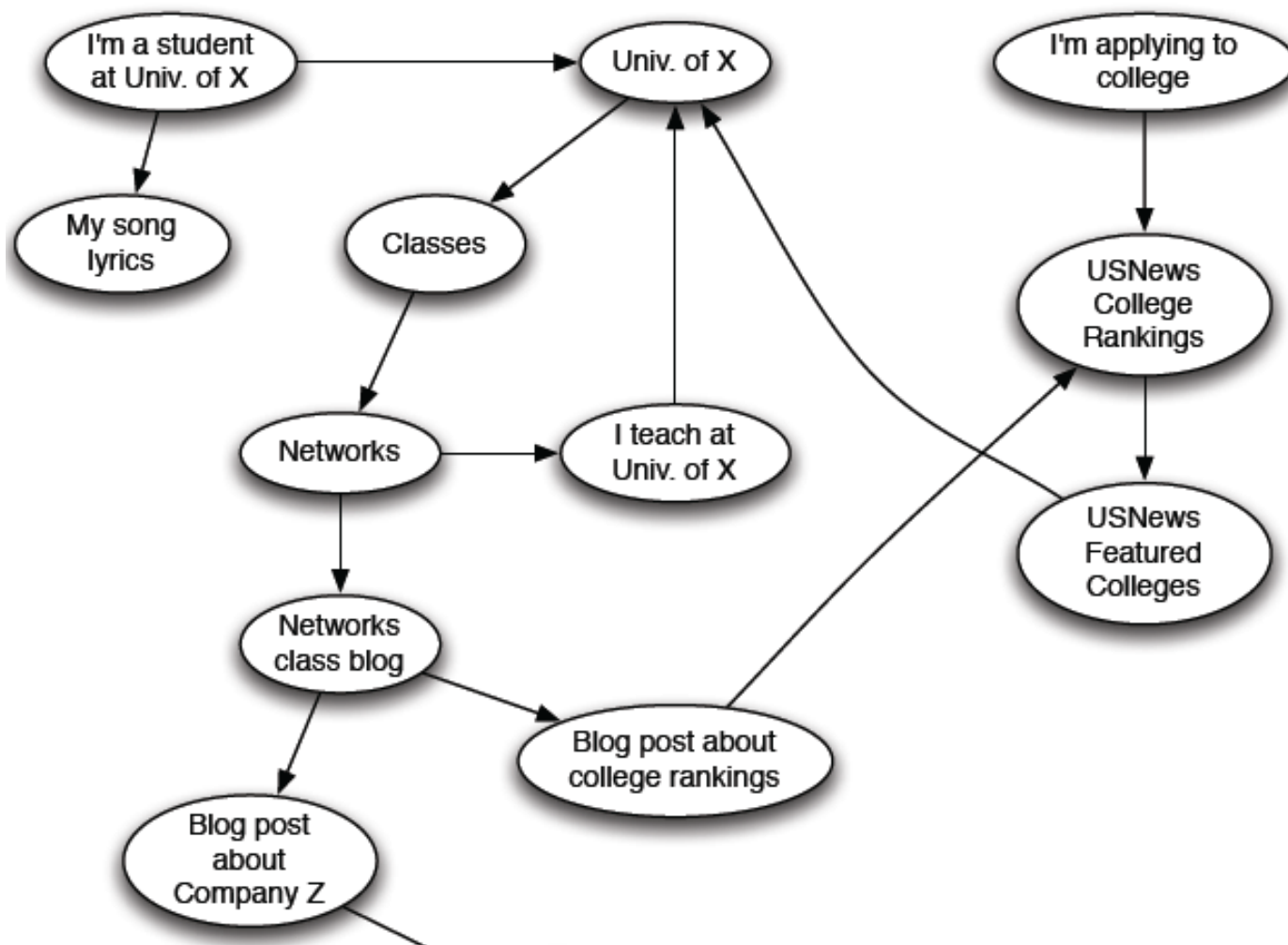
Stanford
University

Web as a Graph

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



Web as a Directed Graph



SHORT HISTORY of INFORMATION RETRIEVAL

- **Classic Document Collections:** Hierarchical Indexing, Libraries
- **1° Digital Revolution:** Centralized Database Systems and Search Systems (Programming Languages for Queries, SQL)
- **2° Digital Revolution (1989): The World Wide Web:**
 - HTML Pages, Hyperlinks. Decentralized Information System.
 - IR on the WEB: WEB Search Engines and the Link Analysis

Information Retrieval on WWW

Main Differences and New Challenges in WWW IR:

- Huge Size
- Evolving
- Self-Organized & Distributed (no standard rules)
- Hyperlinked

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**, dynamics, full of untrusted documents, random things, web spam, etc.

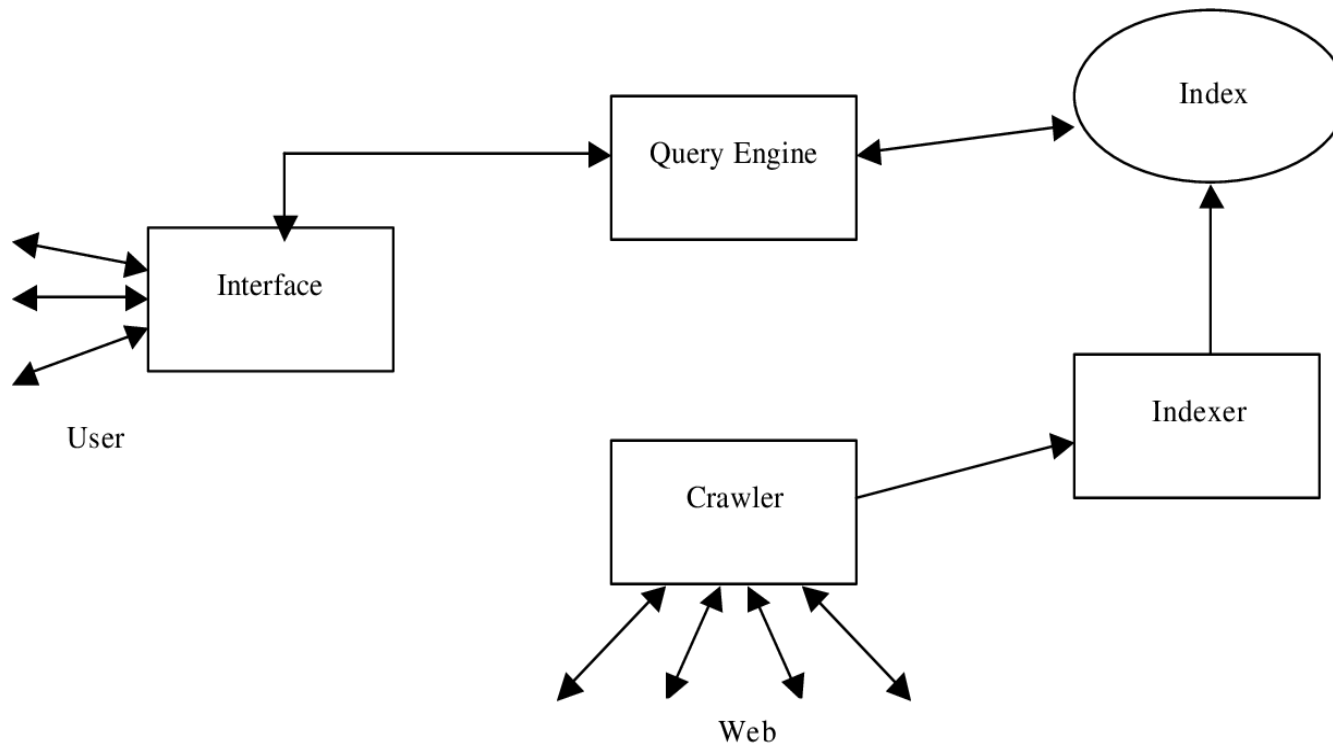


WWW Search Engines

- **MAIN TASK:** Given a set S of words, find the *most relevant* Web Pages for S
- The term «*most relevant*» *hides the most challenging aspect of WWW IR: How to select the first few pages which are more relevant for S among millions of them?*

WWW Search Engines

- **WWW Search Engines** are complex systems formed by several SW modules (see for instance [Langville_Meyer_06])



Web Search: 2 Challenges

2 challenges of web search:

- (1) Web contains many sources of information
Who to “trust”? How to recover them?
 - **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

WWW Search Engines: The Query Module

- **Query Module (QM):** it converts a user's natural language **Query** into a language the *WWW-Search Engine System* can understand (usually numbers), and consults the **Index Module**
- **QM** consults the **content index** and its **inverted file** to select a set **P** of pages that contain the **Query terms T**, i.e.,

$$P := \{ \text{Relevant Pages for } T \}$$

- Then, **QM** passes **P** to the **Ranking Module**

WWW Search Engines: The Ranking Module

- **Ranking Module:**
 - (i) Compute and Assign an **Overall Score** to every **page p** in **P**,
 - (ii) Set **P** is then returned to the *User* in order of the **Overall Score**.
- **Overall Score** is based on **two scores** (computed by the Module):
- **Content Score:** it depends on several parameters: # **query-words**'s occurrences in **p**; **query-word**'s presence in the *title* or in *bold* in **p**.
- **Popularity Score** (**focus** of our course): it is determined from a **Link Analysis** of the **Web's hyperlink structure** (i.e. a Directed Graph).
- **Note:** The **Ranking Module** is perhaps the **most important component** of the Search Process because the output of the query module often results in too many (thousands of) **Relevant Pages** that the User must sort through.

Popularity Score & Link Analysis: The WWW IR Revolution

- «*Nobody wants to be picked last for teams in gym class. Likewise, nobody wants their webpage to appear last in the list of relevant pages for a search query. As a result, many grown-ups transfer their high school wishes to be the “Most Popular” to their webpages*»
- By 1998, the traditional **Content Score** was buckling under the **Web’s massive size** and the **death grip of spammers**. In 1998, the **Popularity Score** came to the rescue of the **Content Score**. The **Popularity Score** became a crucial complement to the **Content Score** and provided *Search Engines* with **impressively accurate results** for **all types of queries**. The **Popularity Score**, also known as the *importance score*, harnesses the information in the immense **Graph** created by the **Web’s hyperlink structure**. Thus, models exploiting the **Web’s hyperlink structure** are called **Link-Analysis** models.
- **Note:** The impact that these link-analysis models have had is truly awesome. Since 1998, the use of *Web Search Engines* has increased dramatically. In fact, an April 2004 survey by *Websense, Inc.*, reported that half of the respondents would rather forfeit their habitual morning cup of coffee than their connectivity. That’s because today’s *Search Tools* allow **Users** to answer in seconds **queries** that were impossible just a decade ago (from fun searches for pictures, quotes, and snooping amateur detective work to more serious searches for academic research papers and patented inventions).

The WWW IR Revolution: The Page-Rank Algorithm

Before 1998, the **Web Graph** was largely an untapped source of information. *While **Computer Science Researchers** like Kleinberg and Brin and Page recognized this **Graph's Potential**, most people wondered just what the **Web Graph** had to do with *Search_Engine* results ☺☺.*

Key Idea of the PR Algorithm: The connection is understood by viewing a **hyperlink** as a **recommendation**. A **hyperlink** from my homepage to your page is my **endorsement** of your page. Thus, a page with more **recommendations** (which are realized through **in-links**) must be more important than a page with a few **in-links**.

The Page-Rank Algorithm

Key Issue:

- Similar to other Recommendation Systems (biblio citations, letters of references), the **status** of the recommender is also important: one personal endorsement from B. Obhama probably does more to strengthen a job app than 20 endorsements from 20 unknown teachers and colleagues. On the other hand, if the job interviewer learns that B. Obhama is very generous with his praises of employees, and he (or his secretary) has written over 40,000 recommendations in his life, then his recommendation suddenly drops in weight. **Thus, weights signifying the status of a recommender must be lowered for recommenders with little discrimination. In fact, the weight of each endorsement should be tempered by the total number of recommendations made by the recommender.**

Actually, this is exactly how Google's PageRank popularity score works.

The Page-Rank Algorithm

Main Idea of PageRank Algorithm:

«a Web-Page is *important* if it is pointed to by other *important* Web-Pages»

Sounds *circular*, doesn't it? 😊😊

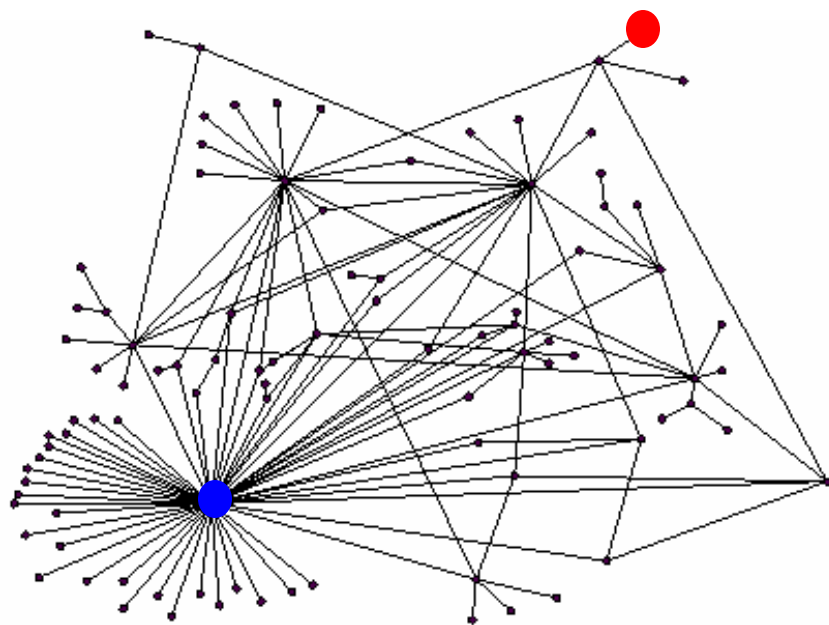
We will see this can be *efficiently* implemented thanks to a beautifully simple *mathematical algorithm*.

Ranking Nodes on the Graph

- All web pages are not equally “important”

www.joe-schmoe.com vs. www.stanford.edu

- There is large diversity in the web-graph node connectivity.
Let's rank the pages by the link structure!



Link Analysis Algorithms

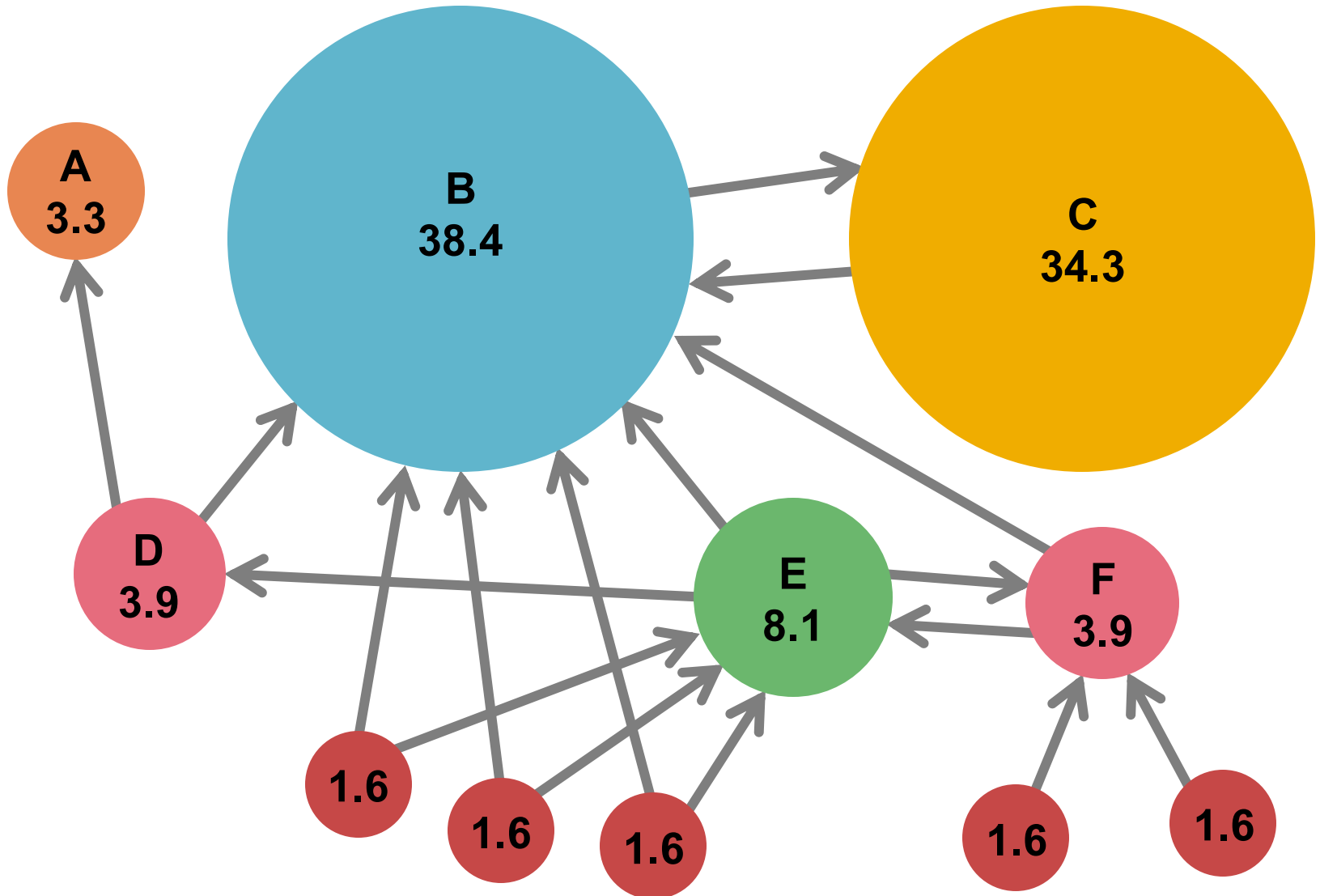
- We will cover the following **Link Analysis approaches** for computing **importances** of nodes in a graph:
 - Page Rank
 - Topic-Specific (Personalized) Page Rank
 - Web Spam Detection Algorithms

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 has 23,400 in-links
 - www.joe-schmoe.com has 1 in-link
- **Are all in-links are equal?**
 - Links from important pages count more
 - Recursive question! How to solve it?

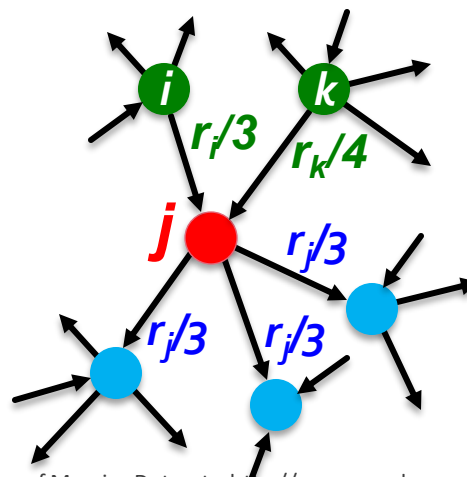
Example: PageRank Scores



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

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



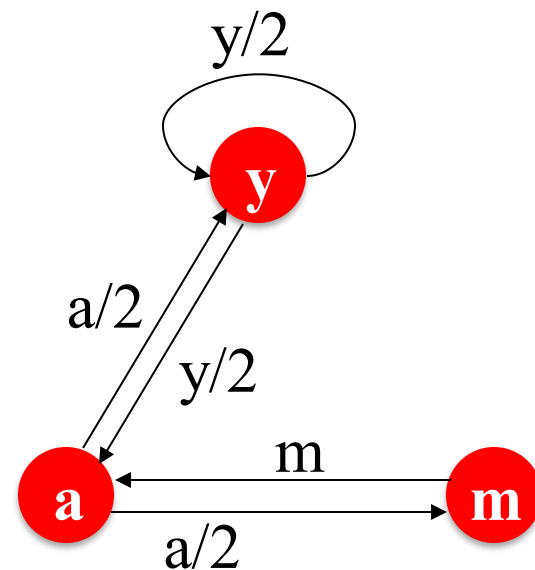
PageRank: The “Flow” Model

- 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

The web in 1839



“Flow” equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

Solving the 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!**

Flow equations:

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

NOTA

- LA PARTE DI ALGEBRA LINEARE E DI MARKOV CHAIN E' FATTA BENE NELLE NOTE:
- [1B-Google_slides.pdf](#)
- disponibili su Teams (folder: *materiale mining massive data*)

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 $\rightarrow \forall i : \sum_j M_{ji} = 1$
 - Columns sum to 1

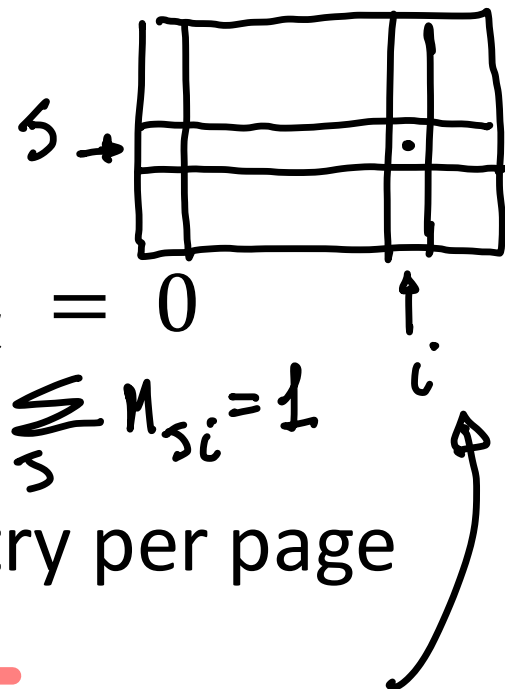
- **Rank vector r** : vector with an entry per page

- r_i is the importance score of page i

- $\sum_i r_i = 1 \rightarrow$ NORMALIZATION

- **The flow equations can be written**

$$r = M \cdot r$$



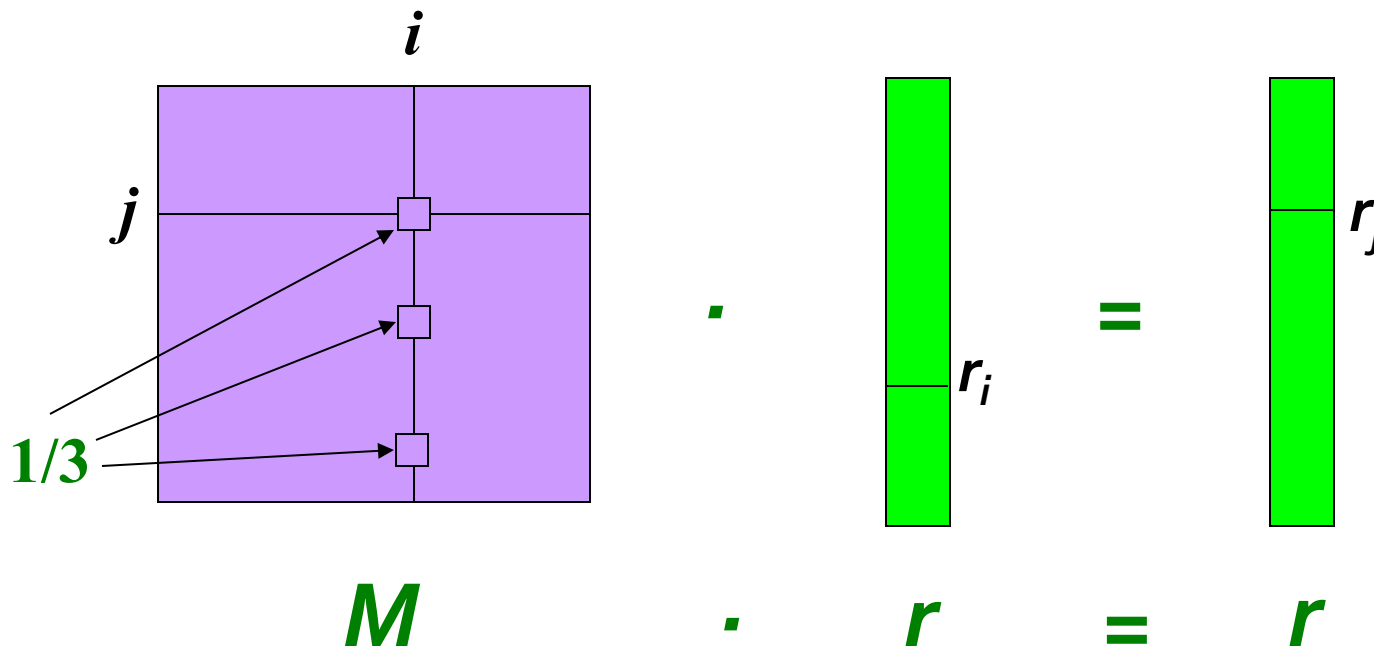
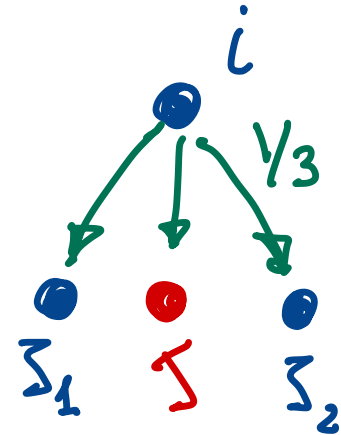
$$r_j = \sum_{i \rightarrow j} \frac{r_i}{d_i}$$

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

$$r = M \cdot r \quad (\lambda=1)$$

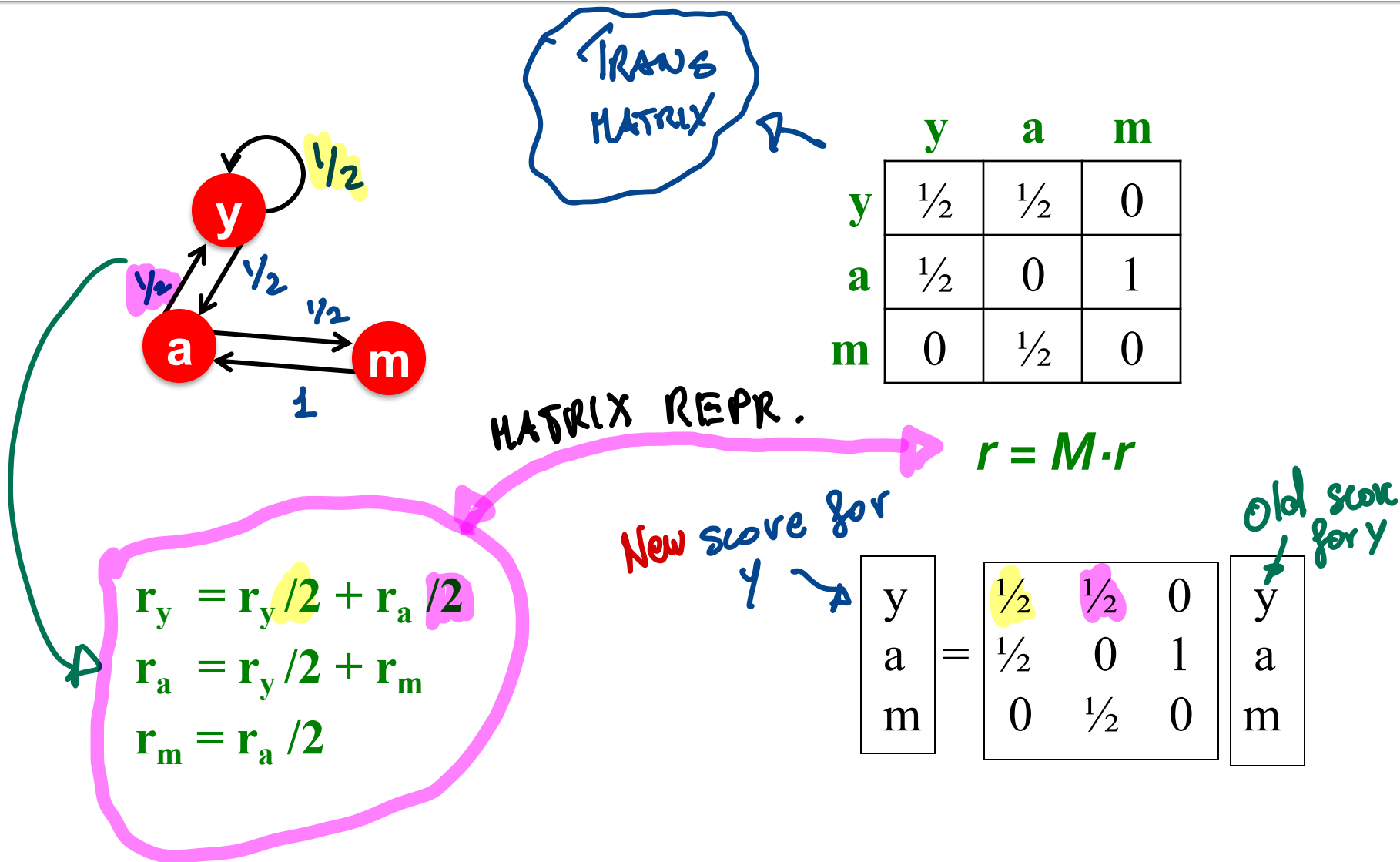
- So the rank vector r is an eigenvector of the stochastic web matrix M

- In fact, its first or principal eigenvector, with corresponding eigenvalue 1 ~~$\lambda=1$~~
 - Largest eigenvalue of M is 1 since M is column stochastic (with non-negative entries)
 - We know r is unit length and each column of M sums to one, so $Mr \leq 1$ in L_1 norm

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



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: $\mathbf{r}^{(0)} = [1/N, \dots, 1/N]^T \rightarrow$ UNIF. START

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

- 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

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

$d_i \equiv$ out-degree of node i

“ALMOST” a FIXED POINT

|||

EIGEN VECTOR
of \mathbf{M} with
 $\lambda = 1$

PageRank: How to solve?

■ Power Iteration:

- Set $r_j = 1/N$

- 1: $r'_j = \sum_{i \rightarrow j} \frac{r_i}{d_i} \quad \forall j=1 \dots N$

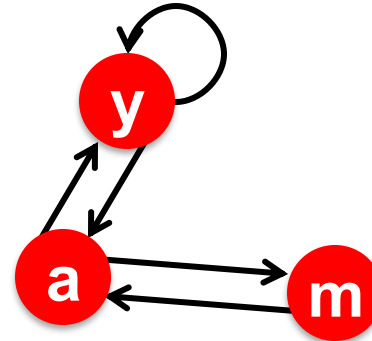
- 2: $r = r'$

- Goto 1

■ Example: $t=0$

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

Iteration 0, 1, 2, ...



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

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

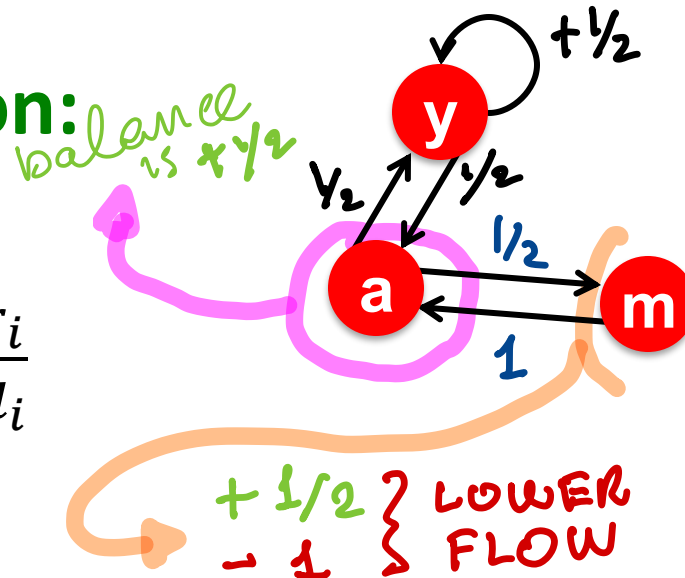
$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

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

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

$$r_a = r_y/2 + r_m$$

$$r_m = r_a/2$$

■ Example:

	$t=0$	$t=1$	$t=2$	$t=3$
r_y	1/3	1/3	5/12	9/24	
r_a	1/3	3/6	1/3	11/24	...
r_m	1/3	1/6	3/12	1/6	

Iteration 0, 1, 2, ...

6/15
6/15
3/15

fixed point
in
Eigen Vector!

Why Power Iteration works? (1)

Details!

■ Power iteration:

[A method for finding dominant eigenvector (the vector corresponding to the largest eigenvalue)

- $\mathbf{r}^{(1)} = \mathbf{M} \cdot \mathbf{r}^{(0)}$

- $\mathbf{r}^{(2)} = \mathbf{M} \cdot \mathbf{r}^{(1)} = \mathbf{M}(\mathbf{M}\mathbf{r}^{(1)}) = \mathbf{M}^2 \cdot \mathbf{r}^{(0)}$

- $\mathbf{r}^{(3)} = \mathbf{M} \cdot \mathbf{r}^{(2)} = \mathbf{M}(\mathbf{M}^2\mathbf{r}^{(0)}) = \mathbf{M}^3 \cdot \mathbf{r}^{(0)}$

■ Claim:

[Sequence $\mathbf{M} \cdot \mathbf{r}^{(0)}, \mathbf{M}^2 \cdot \mathbf{r}^{(0)}, \dots \mathbf{M}^k \cdot \mathbf{r}^{(0)}, \dots$ approaches the dominant eigenvector of \mathbf{M}]

PROOF \rightarrow next slides

Why Power Iteration works? (2)

Details!

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

- **Proof:**

- 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 $\lambda_1 > \lambda_2 > \dots > \lambda_n$ $\det(M) \neq 0$

$\lambda_1 = 1$

- Vectors x_1, x_2, \dots, x_n form a basis and thus we can write:
 $\forall \text{ any } r^{(0)}: r^{(0)} = c_1 x_1 + c_2 x_2 + \dots + c_n x_n$, for some $\langle c_1, \dots, c_n \rangle$

- $M r^{(0)} = M(c_1 x_1 + c_2 x_2 + \dots + c_n x_n)$
 $= c_1 (M x_1) + c_2 (M x_2) + \dots + c_n (M x_n)$
 $= c_1 (\lambda_1 x_1) + c_2 (\lambda_2 x_2) + \dots + c_n (\lambda_n x_n)$

\rightarrow MUST BE $\neq \emptyset$!

- **Repeated multiplication on both sides produces**

$$M^k r^{(0)} = c_1 (\lambda_1^k x_1) + c_2 (\lambda_2^k x_2) + \dots + c_n (\lambda_n^k x_n)$$

\rightarrow largest = 1

$\wedge 1$

$\wedge 1$

Why Power Iteration works? (3)

Details!

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

- Repeated multiplication on both sides produces

$$M^k r^{(0)} = c_1(\lambda_1^k x_1) + c_2(\lambda_2^k x_2) + \dots + c_n(\lambda_n^k x_n)$$

- $M^k r^{(0)} = \lambda_1^k \left[\underbrace{c_1 x_1}_{=1} + c_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k x_2 + \dots + c_n \left(\frac{\lambda_n}{\lambda_1} \right)^k x_n \right]$

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

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

converge

- Note if $c_1 = 0$ then the method won't converge

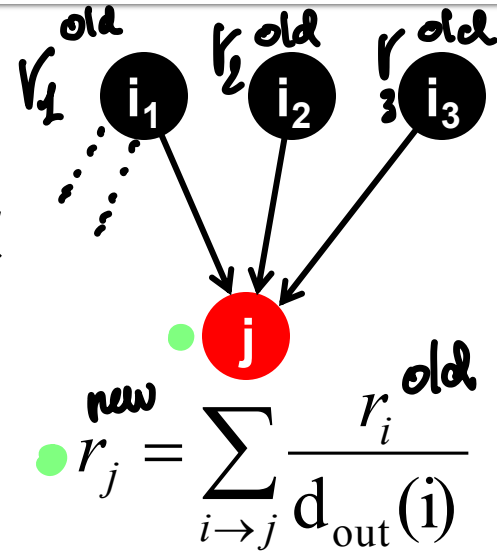
in NORM L_1

the initial vector must have a > 0 comp. on x_1

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

HARKOV
CHAIN

VERY
IMPORT

$$\sum_i p_i(t) = 1$$

Markov Chain

$n = \text{time slots } t$

- A Markov chain is a sequence X_1, X_2, X_3, \dots of random variables ($\sum_v \text{all possible values of } X P(X=v) = 1$) with the property:

- **Markov property:** the conditional probability distribution of the next future state X_{n+1} given the present and past states is a function of the present state X_n alone $\forall x, x_0, \dots, x_n$ $t=n$

$$\Pr(X_{n+1} = x | X_0 = x_0, X_1 = x_1, \dots, X_n = x_n) = \Pr(X_{n+1} = x | X_n = x_n).$$

- If the state space is finite then the transition probabilities can be described with a matrix $P_{ij} = P(X_{n+1} = j | X_n = i)$, $i, j = 1, \dots, m=N$ $P_n[1 \rightarrow 2]$

* $\begin{pmatrix} .85 & .15 \\ .38 & .62 \end{pmatrix} = \begin{pmatrix} P(X_{n+1} = 1 | X_n = 1) & P(X_{n+1} = 2 | X_n = 1) \\ P(X_{n+1} = 1 | X_n = 2) & P(X_{n+1} = 2 | X_n = 2) \end{pmatrix}$

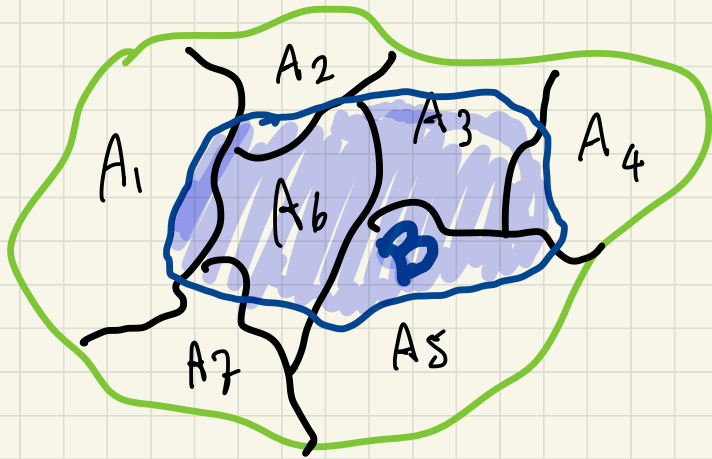
Note: $P = M^T$

column $j=1$ $i=2$ \hookrightarrow link $(2,1)$

* [rows sum to 1]

PROBABILITY SPACE $\langle \Omega, \Pr \rangle$ $|\Omega| = N$

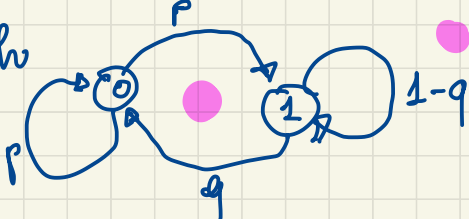
Sets $\{A_1, \dots, A_k\}$ a partition of Ω then, $\forall B \subseteq \Omega$,
it holds:
$$\Pr[B] = \sum_{j=1}^k \Pr[B|A_j] \cdot \Pr[A_j]$$



We use this law
for computing
MARKOV CHAIN'S
DISTRIBUTIONS!

$$\bar{\pi}(t_{i+1}) = F(\bar{\pi}(t_i))$$

~~6/4~~ H.C. $\left\{ \begin{array}{l} \mathcal{S} = \{0, 1\} ; \text{ with} \\ \{X_0, X_1, \dots, X_t, \dots\} \end{array} \right.$ 1-p



$\left\{ \begin{array}{l} P_r[X_t = 1 | X_{t-1} = 0] = p \\ P_r[X_t = 0 | X_{t-1} = 1] = q \end{array} \right.$

TRANSITION MATRIX $P = \begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix}$

Consider $\bar{\pi}(t) = \begin{pmatrix} \pi_0(t) \\ \pi_1(t) \end{pmatrix} \equiv \text{distr. of } X_t$ (it is a col. vector!)

$\bar{\pi} = \langle P[i, j] \rangle \equiv P_r[X_t = j | X_{t-1} = i] \forall t!$

COL. ROW

How to compute $\bar{\pi}(t+1)$? We need to know $\bar{\pi}(t)$:

$\forall j = 0, 1$, we have (i) $P_r[X_t = j] \equiv \sum_{i=0,1} P_r[X_t = j | X_{t-1} = i] \cdot P_r[X_{t-1} = i]$

In MATRIX ALGEBRA, (i) is equivalent to: *

$$\bar{\pi}(t+1) = P \cdot \bar{\pi}(t) \Rightarrow \bar{\pi}(t+1) = P \cdot (P \cdot \bar{\pi}(t-1)) = \dots = P^t \cdot \bar{\pi}(0) \rightarrow \text{STARTING DISTRIB.}$$

(**) To define a H.C., we need:

$X \equiv \langle \mathcal{S}, P \rangle$, i.e., SET of STATES \wedge TRANS. MATRIX

MAIN QUESTIONS: \rightarrow arbitrary distrib.

i) starting from $\pi(0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$: $\underbrace{\text{PR}[\text{Agent is INFECTED at round } t=3]}_{(a)} = ?$

Answer
POWER METHOD

$$P^3 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix}^3 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} (a) \\ 1-(a) \end{pmatrix}$$

$\left[\begin{array}{l} \text{State 0} \equiv \text{SUSCEPTIBLE} \\ \text{State 1} \equiv \text{INFECTED} \end{array} \right]$

ii) starting from **any** $\pi(0)$, after a large $\tilde{t} \rightarrow t \rightarrow \infty$, can we find a **stationary distribution** π s.t. $\forall t \geq \tilde{t}$, we have

$$\pi(t) = P \cdot \pi(t-1) = \pi \quad ?$$

Answer

We need to find a vector $\tilde{\pi}$ s.t. $P \cdot \tilde{\pi} = \tilde{\pi}$, i.e.,
an **eigenvector** $\tilde{\pi}$ with **eigenvalue** $\lambda = 1$

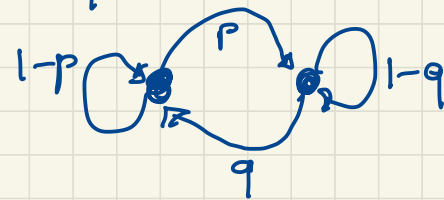
MAIN QUESTIONS ON MARKOV CHAIN: $X \equiv \langle \Omega, P \rangle$

- (1) When $P^t \pi(0)$ converge?
- (2) Where? which is the meaning?
- (3) CONVERGENCE TIME?

Consider birth-death H.C. (●): $\Omega = \{0, 1\}$; $p \equiv$ birth prob
 $q \equiv$ death prob.

$$\forall \pi = \langle \pi_0, \pi_1 \rangle, \pi(0) = \pi:$$

$$\begin{pmatrix} \pi_0(1) = \pi_0 \cdot (1-p) + \pi_1 \cdot q \\ \pi_1(1) = \pi_0 \cdot p + \pi_1 \cdot (1-q) \end{pmatrix} \equiv P \cdot \pi(0),$$



In general: $\pi(t) = P^t \cdot \pi(0)$, when stops? CNES to stop is

$\exists T \in \mathbb{N}$ s.t. (2) $\bar{\pi}(T) = P \cdot \bar{\pi}(T-1) \equiv$ a fixed point of the process $\{\pi(0), \dots, \bar{\pi}(t), \dots$

Indeed, if (2) holds Then:

$$\forall t \geq T \Rightarrow \bar{\pi}(t) = P \cdot \bar{\pi}(t-1) = \dots = P \cdot \overbrace{P \cdot \dots \cdot P}^{\pi(T)} \cdot \bar{\pi}(T-1) = \bar{\pi}(T-1)$$

So the H.C.

reaches a STATIONARY DISTRIBUTION $\bar{\pi}$ iff

$\bar{\pi} = P \bar{\pi} \iff \bar{\pi}$ is an eigenvector of P with (MAX) eigenvalue $\lambda = 1$ $\bar{\pi} = P \cdot (\lambda \cdot \pi)$

EX (.) $\begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix} \Rightarrow$ consider $\tilde{\pi} = \begin{pmatrix} q/p+q \\ p/p+q \end{pmatrix}$

$$\begin{pmatrix} 0 & 1 \\ 1 & p \end{pmatrix} \begin{pmatrix} q/p+q \\ p/p+q \end{pmatrix} = \begin{pmatrix} (1-p) \cdot \frac{q}{p+q} + p \frac{p}{p+q} = \frac{q - pq + pq}{p+q} \\ \text{SIMILAR } \uparrow \text{ CALCULUS} \end{pmatrix} = \begin{pmatrix} q/p+q \\ p/p+q \end{pmatrix}$$

so, $\tilde{\pi} = \begin{pmatrix} q/p+q \\ p/p+q \end{pmatrix}$ is an eigenvector of P and is

A STATIONARY DISTRIBUTION for H.C. $P \equiv (.)$

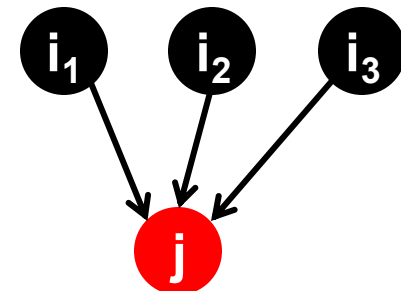
Q. is it UNIQUE? Q. starting π , $P^t \pi \rightarrow \tilde{\pi}$??

The Stationary Distribution

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

- Follows a link uniformly at random

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



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

- Suppose the random walk reaches a state

$$p(t+1) = M \cdot p(t) = p(t) \quad \text{fixed point}$$

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

→ column vector

$$M = (M(i,j)) = \underset{\substack{\downarrow \\ \text{SURFER}}}{\Pr[i \rightarrow j]} = \frac{r_i}{d_i}$$

Existence and Uniqueness

stationary pag-rank \vec{v} exists? is it UNIQUE?

- 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 $t = 0$

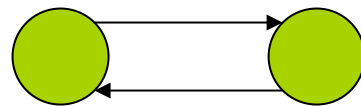
Ergodic Markov chains

□ A Markov chain is **ergodic** if:

■ Informally: *there is a path from any state to any other; and the states are not partitioned into sets such that all state transitions occur cyclically from one set to another.*

■ Formally: for any start state, after a finite transient time T_0 , the probability of being in any state at any fixed time $T > T_0$ is nonzero.

BIPARTITE



Not ergodic (even/odd).

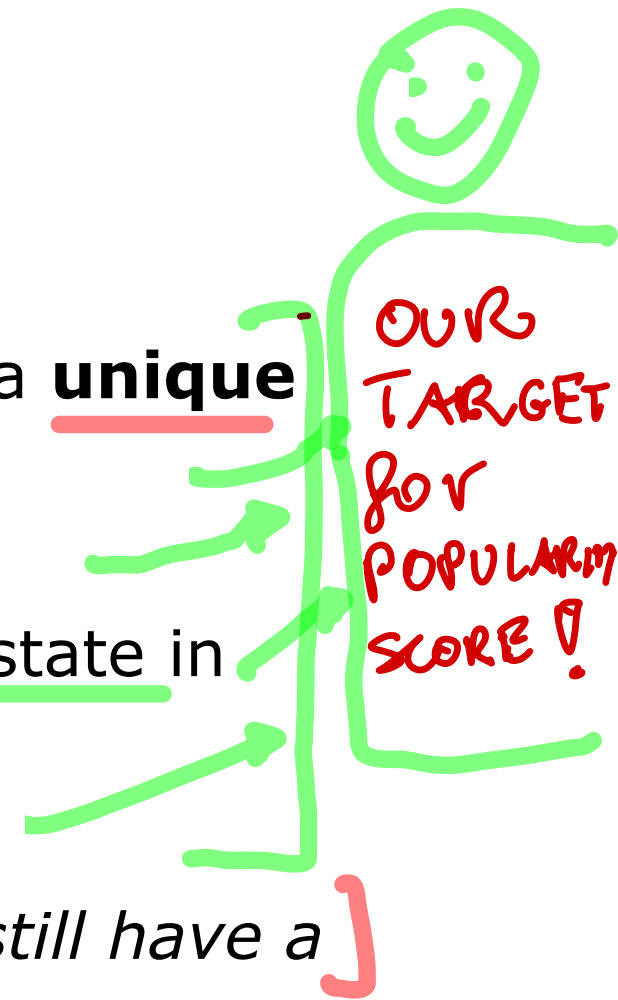
Not ergodic: the probability to be in a state, at a fixed time, e.g., after 500 transitions, is always either 0 or 1 according to the initial state.

GRAPH Theory!
→ strong CONNECT.
+
if length!

→ $\forall i \in [N], \text{vector}$
 $P^{(t)}(i) > 0, \forall t \geq T_0$

Ergodic Markov chains

- For any ergodic Markov chain, there is a **unique long-term visit rate** for each state
 - *Steady-state probability distribution*
- Over a long time-period, we visit each state in proportion to this rate
- It doesn't matter where we start.
- *Note: non ergodic Markov chains may still have a steady state.*



PageRank: The Google Formulation

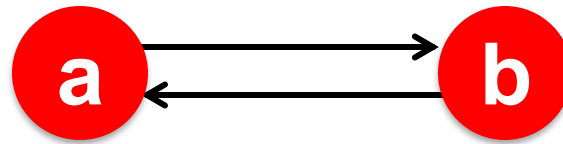
PageRank: Three Questions

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

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

Does this converge?

MAIN ISSUE: WEB GRAPH IS NOT ERGODIC ☹️



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

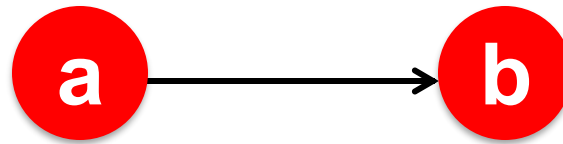
■ Example:

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

Iteration 0, 1, 2, ...

Does it converge to what we want?

it converges! But....



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

■ Example:

$$\begin{array}{c} r_a \\ r_b \end{array} = \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \end{array}$$

Iteration 0, 1, 2, ...

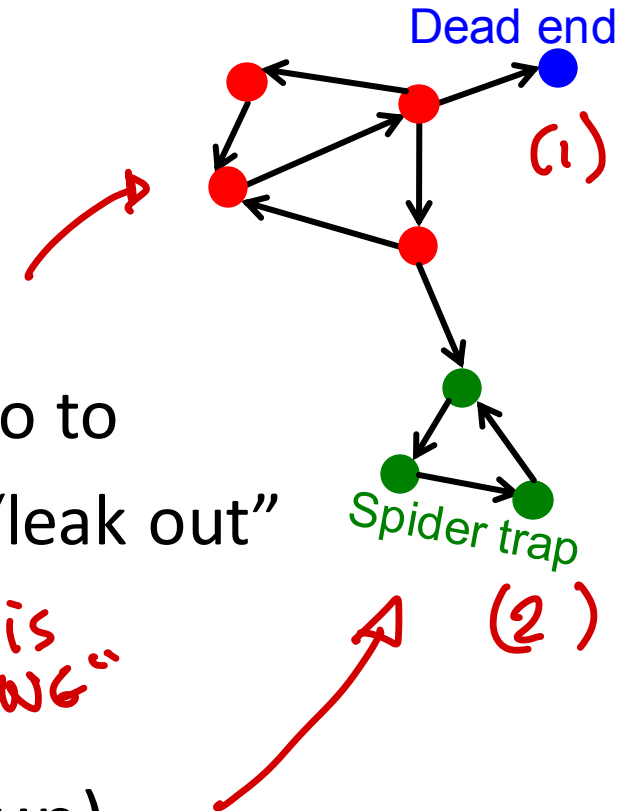


↓ Fixed Point no meaning!

PageRank: Problems

2 problems:

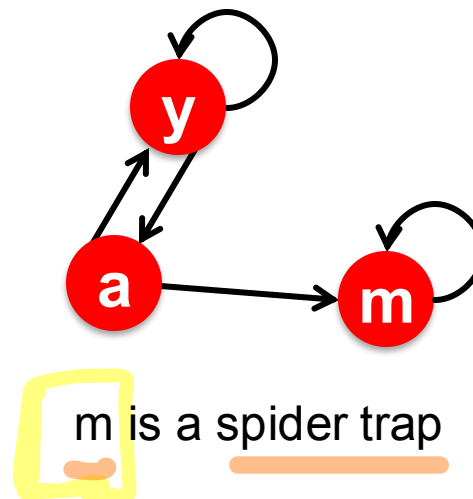
- (1) Some pages are **dead ends** (have no out-links)
 - Random walk has “nowhere” to go to
 - Such pages cause importance to “leak out”
- (2) **Spider traps:** (all out-links are within the group)
 - Random walked gets “stuck” in a trap
 - And eventually spider traps absorb all importance



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	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:

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

Iteration 0, 1, 2, ...

0
0
1



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

Solution: Teleports!

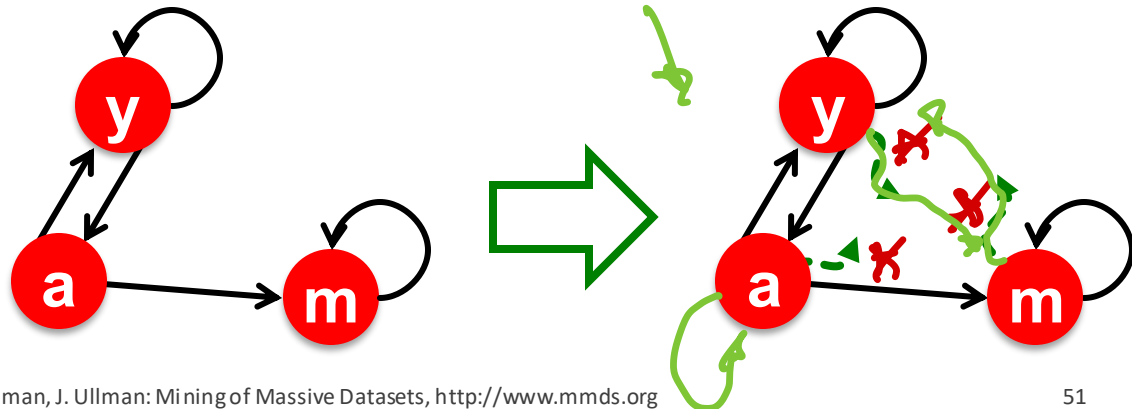
- The Google solution for spider traps: **At each time step, the random surfer has two options**

- 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**

TELEPORTS BRIDGES

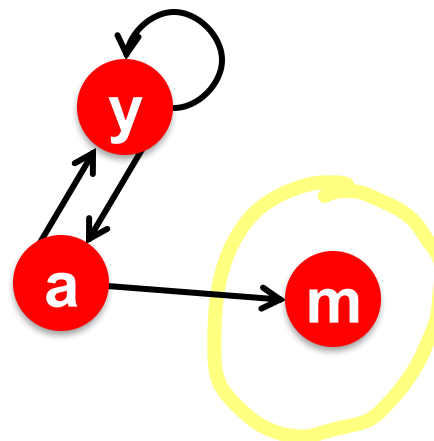
* the walker does not follow the link!
it jumps to any page of the WEB u.a.v.



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

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

$$r_a = r_y/2$$

$$r_m = r_a/2$$

■ Example:

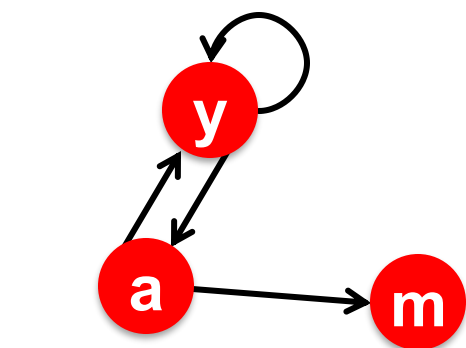
$$\begin{bmatrix} r_y \\ r_a \\ r_m \end{bmatrix} = \begin{array}{cccccc} 1/3 & 2/6 & 3/12 & 5/24 & & 0 \\ 1/3 & 1/6 & 2/12 & 3/24 & \dots & 0 \\ 1/3 & 1/6 & 1/12 & 2/24 & & 0 \end{array}$$

Iteration 0, 1, 2, ...

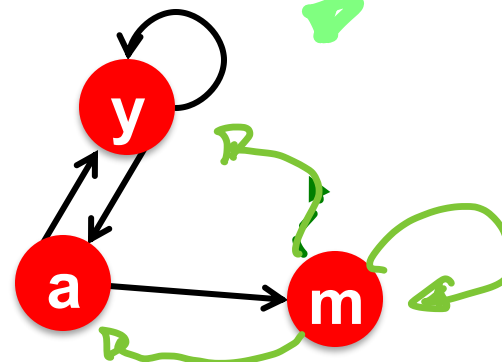
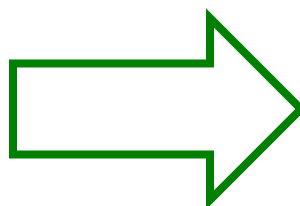
Here the PageRank “leaks” out since the matrix is not stochastic.

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

$\beta = 0!$

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]

OLD
FLUID
EQ.s

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

TELEPORT
TERM

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:**

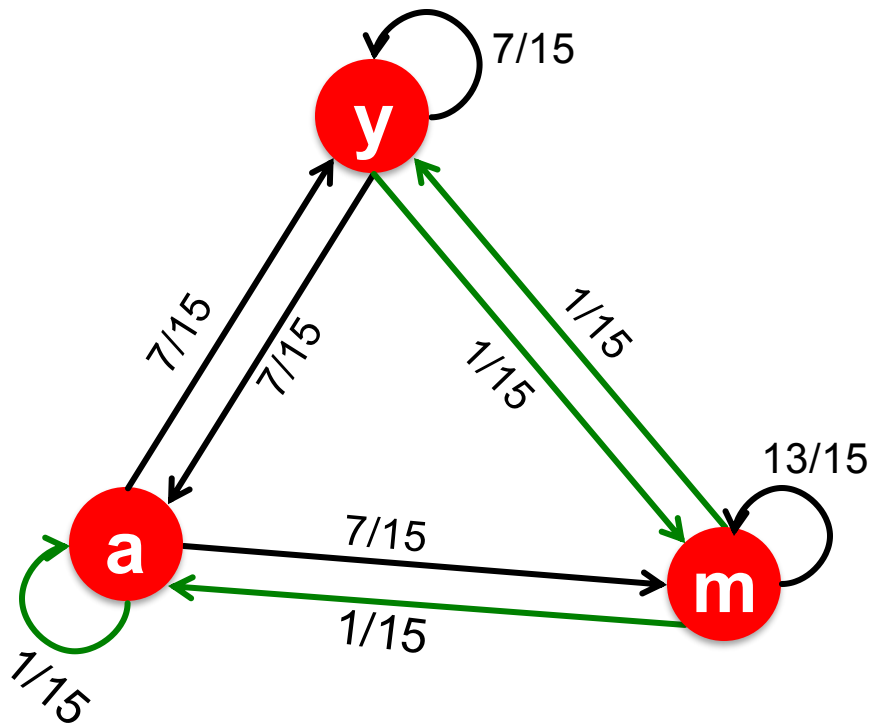
$[1/N]_{N \times N} \dots N$ by N matrix
where all entries are $1/N$

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

- **We have a recursive problem: $r = A \cdot r$**
And the Power method still works! : $v^{(t)} = A^{(t)} \cdot v^{(0)}$
- **What is β ?**

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

Random Teleports ($\beta = 0.8$)



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

\downarrow
Teleport

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

y	=	1/3	0.33	0.24	0.26
a		1/3	0.20	0.20	0.18
m		1/3	0.46	0.52	0.56

A

7/33
5/33
21/33

MAX
Eigen
Vector
of A
popularity
SCORES

**How do we actually compute
the PageRank?**

Computing Page Rank

- **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 pages**

- We need 4 bytes for each entry (say)

- 2 billion entries for vectors, approx 8GB

- **Matrix \mathbf{A} has N^2 entries**

- 10^{18} is a large number!

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

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

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

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

Rearranging the Equation

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

Note: Here we assumed \mathbf{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:**
 $\forall j: r_j^{new} = r_j'^{new} + \frac{1-S}{N}$ **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 $4 \times 10 \times 1$ 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 $j = 1 \dots d_i$

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

0	
1	
2	
3	
4	
5	
6	

r^{new}

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

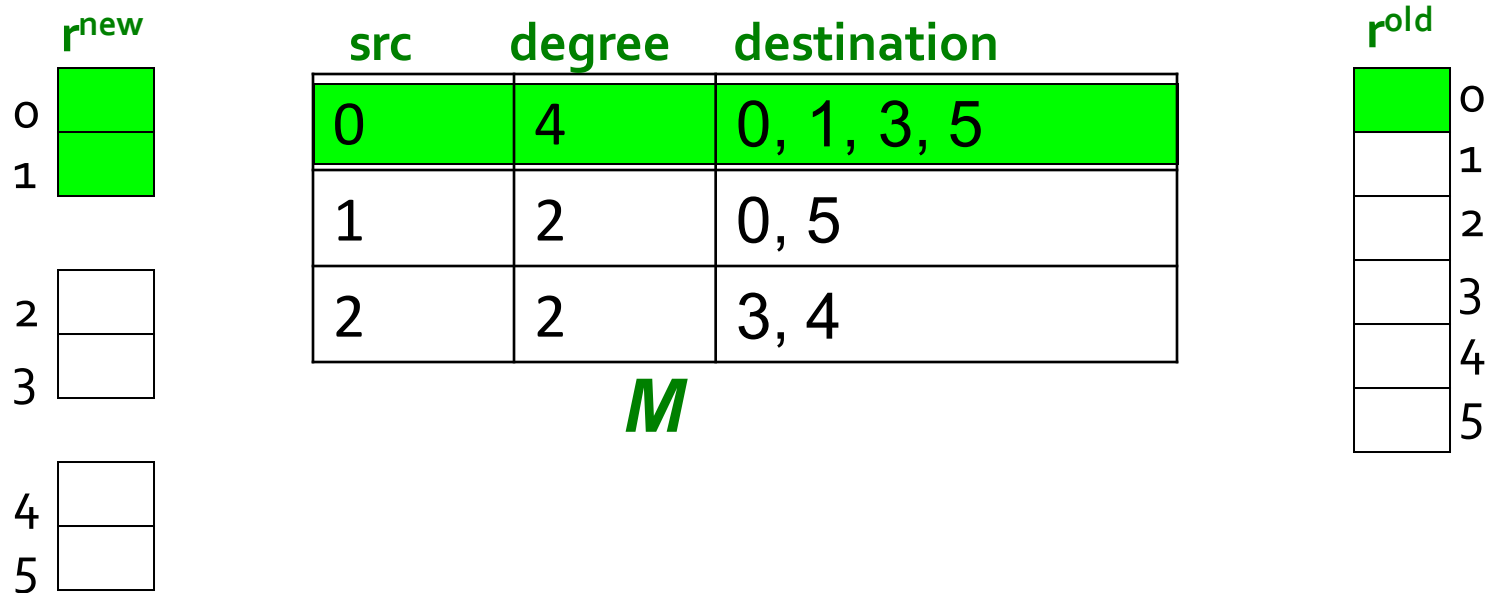
	0
	1
	2
	3
	4
	5
	6

r^{old}

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?

Block-based Update Algorithm



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

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

r^{new}

0	
1	

src	degree	destination
0	4	0, 1
1	3	0
2	2	1

2	
3	

0	4	3
2	2	3

4	
5	

0	4	5
1	3	5
2	2	4

r^{old}

	0
	1
	2
	3
	4
	5

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:
 $= |M|(1+\varepsilon) + (k+1)|r|$

Some Problems with 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