# Information Retrieval & Social Web
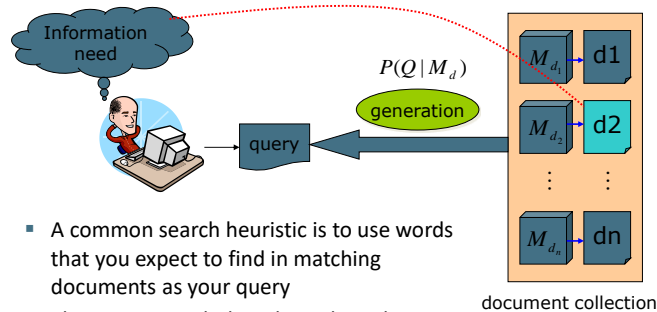
CS 525/DS 595
Worcester Polytechnic Institute
Department of Computer Science
Instructor: Prof. Kyumin Lee

---

# Previous Class…

Statistical Language Models

---

# IR based on Language Model (LM)



$P(Q|M_d)$

generation

query

Information need

document collection

- A common search heuristic is to use words that you expect to find in matching documents as your query
- The LM approach directly exploits that idea!

---

# Basic mixture model summary

- General formulation of the LM for IR

$$P(q|d) \propto \prod_{t \in q}((1-\lambda)P(t|M_c) + \lambda P(t|M_d))$$

general language model

individual-document model

- The user has a document in mind, and generates the query from this document.
- The equation represents the probability that the document that the user had in mind was in fact this one.

## Example

- Document collection (2 documents)
  - $d_1$: Xerox reports a profit but revenue is down
  - $d_2$: Lucent narrows quarter loss but revenue decreases further
- Model: MLE unigram from documents; $\lambda = \frac{1}{2}$
- Query: *revenue down*
  - $P(Q|d_1) = [(1/8 + 2/16)/2] \times [(1/8 + 1/16)/2]$
    $= 1/8 \times 3/32 = 3/256$
  - $P(Q|d_2) = [(1/8 + 2/16)/2] \times [(0 + 1/16)/2]$
    $= 1/8 \times 1/32 = 1/256$
- Ranking: $d_1 > d_2$

## Previous Class…

Statistical Language Models

Crawler

## Previous Class…

Statistical Language Models

Crawler

Web APIs

## Available Web APIs

- Twitter: https://dev.twitter.com/
- Flickr: http://www.flickr.com/services/api/
- Google Maps: https://developers.google.com/maps/
- Facebook: http://developers.facebook.com/
- Foursquare: https://developer.foursquare.com/
- Yahoo Boss API: http://developer.yahoo.com/search/boss/
- Wikipedia API: http://www.mediawiki.org/wiki/API:Main_page
- Youtube API: http://code.google.com/apis/youtube/overview.html
- Openstreetmap API: http://wiki.openstreetmap.org/wiki/API
- Halo API: https://developer.haloapi.com/

- List of APIs: https://www.reddit.com/r/webdev/comments/3wrswc/what_are_some_fun_apis_to_play_with/

## Static quality scores

- We want top-ranking documents to be both *relevant* and *authoritative*
- *Relevance* is being modeled by cosine scores
- *Authority* is typically a query-independent property of a document
- Examples of authority signals
  - Wikipedia among websites
  - Articles in certain newspapers
  - A paper with many citations ← Quantitative
  - Many bitly's or diggs
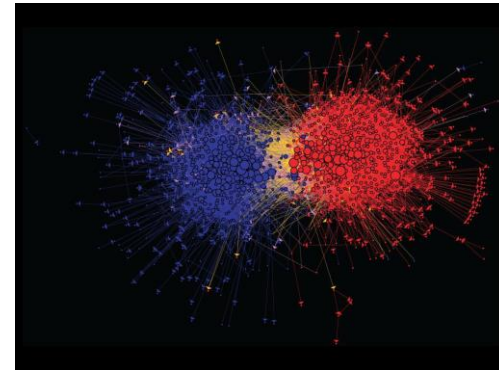  - (Pagerank) ←

## Today: Link Analysis

- Anchor text
- PageRank

## Graph Data: Social Networks



**Facebook social graph**
4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

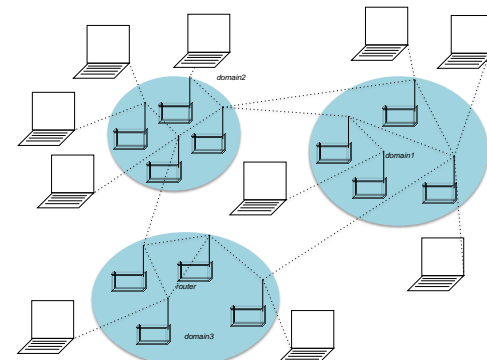## Graph Data: Media Networks



**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]
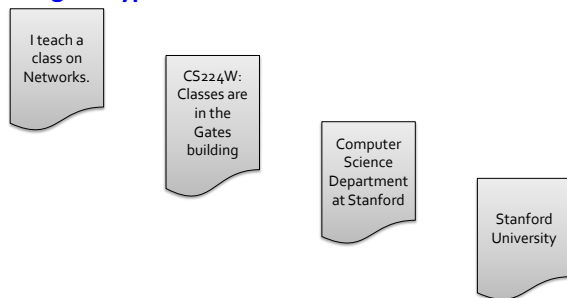
## Graph Data: Communication Nets



# Internet

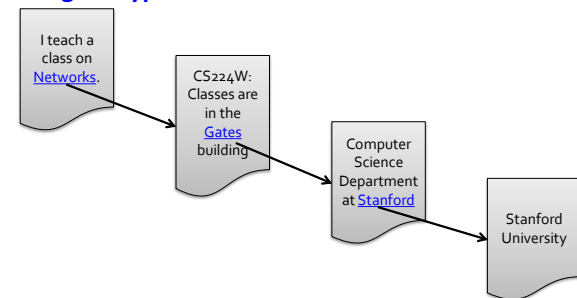## Web as a Graph

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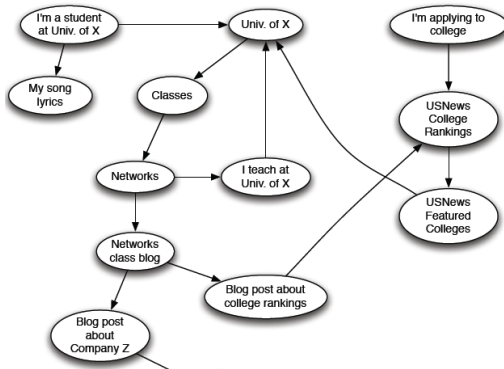  - **Edges: Hyperlinks**



## Web as a Graph

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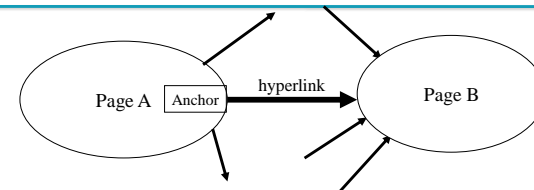
## Web as a Directed Graph



## 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.



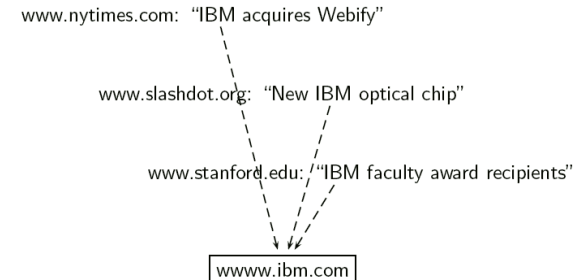# Anchor Text

## The Web as a Directed Graph



- **Assumption 1:** a hyperlink is a quality signal
  - A hyperlink between pages denotes author perceived relevance
- **Assumption 2:** The anchor text describes the target page
  - we use anchor text somewhat loosely here: the text surrounding the hyperlink. Example: "You can find cheap cars <a href= …>here</a>"

## [document text only] vs. [document text + anchor text]

- Searching on [document text + anchor text] is often more effective than searching on [document text only].
- Example: Query *IBM*
  - Matches IBM's copyright page
  - Matches many spam pages
  - Matches IBM wikipedia article
  - May not match IBM home page! (if IBM home page is mostly graphical)
- Searching on anchor text is better for the query IBM.
- **Represent each page by all the anchor text pointing to it.**
- In this representation, the page with the most occurrences of IBM is www.ibm.com.
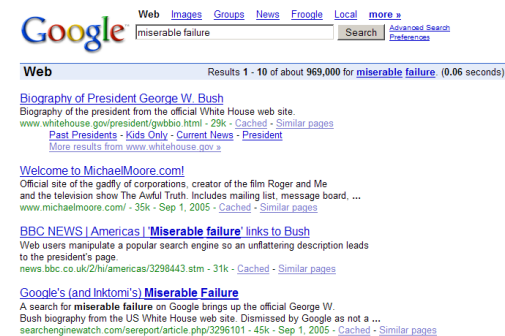
## Anchor text containing *IBM* pointing to www.ibm.com

www.nytimes.com: "IBM acquires Webify"

www.slashdot.org: "New IBM optical chip"

www.stanford.edu: "IBM faculty award recipients"

wwww.ibm.com

## Indexing anchor text

- Thus: anchor text is often a better description of a page's content than the page itself
- Anchor text can be weighted more highly than document text (based on Assumptions 1 & 2)
- Indexing anchor text can have unexpected side effects - Google bombs.
- A Google bomb is a search with "bad" results due to maliciously manipulated anchor text
- Google introduced a new weighting function in January 2007 that fixed many Google bombs

## Google bomb example

# Web Search: Pre-History

## Brief (non-technical) history of Web Search

- Early keyword-based engines ca. 1995-1997
  - Altavista, Excite, Infoseek, Inktomi, Lycos,
- Paid placement ranking: Goto.com (morphed into Overture.com → Yahoo!)
  - Your search ranking depended on how much you paid
  - Auction for keywords: *casino* was expensive!

## Brief (non-technical) history

- 1998+: Link-based ranking pioneered by Google
  - Blew away all early engines
  - Great user experience in search of a business model
  - Meanwhile Goto/Overture's annual revenues were nearing $1 billion
- Result: Google added paid-placement "ads" to the side, independent of search results
  - Yahoo follows suit, acquiring Overture (for paid placement) and Inktomi (for search)
- 2005+: Google gains search share, dominating in Europe and very strong in North America
  - 2009: Yahoo! and Microsoft propose combined paid search offering

## Web search basics



# PageRank

## Link-based ranking

- Query processing with link-based ranking:
  - First retrieve all pages meeting the query (say **venture capital**)
  - Order these by their link popularity (= citation frequency, first generation)
  - . . . or by Pagerank (second generation)

• Simple link popularity (= number of inlinks of a page) is easy to spam.
• Why?



## PageRank:
## Recursive formulation

• Each link's vote is proportional to the **importance of its source page**
• If page P with importance x has n outlines, each link gets x/n votes
• Page P's own importance is the sum of the vote on its inlinks



$$y = y/2 + a/2$$
$$a = y/2 + m$$
$$m = a/2$$

# PageRank basics

- Imagine a web surfer doing a random walk on the web
  - Start at a random page
  - At each step, go out of the current page along one of the links on that page, equiprobably

  1/3
  1/3
  1/3

- "In the steady state" each page has a long-term visit rate - use this as the page's score.
- **PageRank = steady state probability**
  **= long-term visit rate**

# Markov chains

- A Markov chain consists of n states, plus an n×n <u>transition probability matrix</u> **P**.
- **state = page**
- At each step, we are on exactly one of the states.
- For $1 \le i, j \le n$, the matrix entry $P_{ij}$ tells us the probability of $j$ being the next state (page), given we are currently on page (state) $i$.

$$d_i \xrightarrow{P_{ij}} d_j$$

# Markov chains

- Clearly, for all i, $\sum_{j=1}^{N} P_{ij} = 1$
- Markov chains are abstractions of random walks.

# Example web graph

And the corresponding link matrix

|  | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|---|---|---|---|---|---|---|---|
| $d_0$ | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| $d_1$ | 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| $d_2$ | 1 | 0 | 1 | 1 | 0 | 0 | 0 |
| $d_3$ | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| $d_4$ | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| $d_5$ | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| $d_6$ | 0 | 0 | 0 | 1 | 1 | 0 | 1 |

## Transition probability matrix P

|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 0     | 0     | 1     | 0     | 0     | 0     | 0     |
| $d_1$ | 0     | 1     | 1     | 0     | 0     | 0     | 0     |
| $d_2$ | 1     | 0     | 1     | 1     | 0     | 0     | 0     |
| $d_3$ | 0     | 0     | 0     | 1     | 1     | 0     | 0     |
| $d_4$ | 0     | 0     | 0     | 0     | 0     | 0     | 1     |
| $d_5$ | 0     | 0     | 0     | 0     | 0     | 1     | 1     |
| $d_6$ | 0     | 0     | 0     | 1     | 1     | 0     | 1     |

Transition probability matrix

|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 0.00  | 0.00  | 1.00  | 0.00  | 0.00  | 0.00  | 0.00  |
| $d_1$ | 0.00  | 0.50  | 0.50  | 0.00  | 0.00  | 0.00  | 0.00  |
| $d_2$ | 0.33  | 0.00  | 0.33  | 0.33  | 0.00  | 0.00  | 0.00  |
| $d_3$ | 0.00  | 0.00  | 0.00  | 0.50  | 0.50  | 0.00  | 0.00  |
| $d_4$ | 0.00  | 0.00  | 0.00  | 0.00  | 0.00  | 0.00  | 1.00  |
| $d_5$ | 0.00  | 0.00  | 0.00  | 0.00  | 0.00  | 0.50  | 0.50  |
| $d_6$ | 0.00  | 0.00  | 0.00  | 0.33  | 0.33  | 0.00  | 0.33  |

## Long-term visit rate

- Recall: PageRank = long-term visit rate

- Long-term visit rate of page *d* is the probability that a web surfer is at page *d* at a given point in time.

- Next: what properties must hold of the web graph for the long-term visit rate to be well defined?

## Not quite enough

- The web is full of dead-ends.
  - Random walk can get stuck in dead-ends.
  - Makes no sense to talk about long-term visit rates.



## Teleporting

- At a dead end, jump to a random web page.
- At any non-dead end, with probability 10%, jump to a random web page.
  - With remaining probability (90%), go out on a random link.
  - 10% - a parameter.

## Teleporting Matrix

- Recall: At a dead end, jump to a random web page

|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_1$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_2$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_3$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_4$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_5$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |
| $d_6$ | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   | 1/7   |

## Result of teleporting

- With teleporting, we cannot get stuck in a dead end

- There is a long-term rate at which any page is visited (not obvious, will show this).

- How do we compute this visit rate?

## Formalization of "visit": Probability vectors

- A probability (row) vector $\mathbf{x} = (x_1, \dots x_n)$ tells us where the walk is at any point.
- E.g., $(000\dots1\dots000)$ means we're in state $i$.
  $\quad\quad 1 \quad\quad i \quad\quad n$

- More generally, the vector $\mathbf{x} = (x_1, \dots x_n)$ means the walk is in state $i$ with probability $x_i$.

$$\sum_{i=1}^{n} x_i = 1.$$

## Change in probability vector

- If the probability vector is $\mathbf{x} = (x_1, \dots x_n)$ at this step, what is it at the next step?
- Recall that row $i$ of the transition prob. Matrix **P** tells us where we go next from state $i$.
- So from **x**, our next state is distributed as **xP**.

## Steady state example

- The steady state looks like a vector of probabilities $\mathbf{a} = (a_1, \dots a_n)$:
- $a_i$ is the probability that we are in state $i$.



What is the steady state in this example?

## Steady state example

- The steady state looks like a vector of probabilities $\mathbf{a} = (a_1, \dots a_n)$:
- $a_i$ is the probability that we are in state $i$.



For this example, $a_1=1/4$ and $a_2=3/4$.

## How to compute the steady-state?

- Recall, regardless of where we start, we eventually reach the steady state $\mathbf{a}$.
- Start with any distribution (say $\mathbf{x}=(10\dots0)$).
- After one step, we're at $\mathbf{xP}$;
- after two steps at $\mathbf{xP}^2$, then $\mathbf{xP}^3$ and so on.
- "Eventually" means for "large" $k$, $\mathbf{xP}^k = \mathbf{a}$.
- Algorithm: multiply $\mathbf{x}$ by increasing powers of $\mathbf{P}$ until the product looks stable.
- This is called the power method

## Power method: example

Two-node example: $\vec{x} = (0.5, 0.5)$, $P = \begin{pmatrix} 0.25 & 0.75 \\ 0.25 & 0.75 \end{pmatrix}$

$\vec{x}P = (0.25, 0.75) = \vec{x}_2$

$\vec{x}_2 P = (0.25, 0.75)$

Convergence in one iteration!

## Exercise on PageRank

Transition probability matrix of a surfer's walk with teleportation:

P = (1- α) * transition matrix + α * teleporting matrix

- Consider a Web graph with three nodes 1, 2, and 3. The links are as follows: 1->2, 3->2, 2->1, 2->3. Write down the transition probability matrices P and pagerank scores for the surfer's walk with teleporting, with the value of teleport probability α=0.5.

P =

| 0 | 1 | 0 |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |

Each 1 divied by the number of ones in this row

(1- α)*

| 0 | 1 | 0 |
|---|---|---|
| ½ | 0 | ½ |
| 0 | 1 | 0 |

+

α*

| 1/3 | 1/3 | 1/3 |
|---|---|---|
| 1/3 | 1/3 | 1/3 |
| 1/3 | 1/3 | 1/3 |

=

| 1/6 | 2/3 | 1/6 |
|---|---|---|
| 5/12 | 1/6 | 5/12 |
| 1/6 | 2/3 | 1/6 |

## Exercise on PageRank (Cont'd)

Remember

$\vec{x}_1 = \vec{x}_0 P$
$\vec{x}_2 = \vec{x}_1 P$
$\vec{x}_2 = \vec{x}_1 P$
...
...
...
Until converged

$\vec{x}_0 = $

| 1 | 0 | 0 |
|---|---|---|

P=

| 1/6 | 2/3 | 1/6 |
|---|---|---|
| 5/12 | 1/6 | 5/12 |
| 1/6 | 2/3 | 1/6 |

$\vec{x}_1 = $

| 1/6 | 2/3 | 1/6 |
|---|---|---|

$\vec{x}_2 = $

| 1/3 | 1/3 | 1/3 |
|---|---|---|

$\vec{x}_3 = $

| 1/4 | 1/2 | 1/4 |
|---|---|---|

...
...

$\vec{x}_k = $

| 5/18 | 4/9 | 5/18 |
|---|---|---|

converged

# Example web graph



And the corresponding link matrix
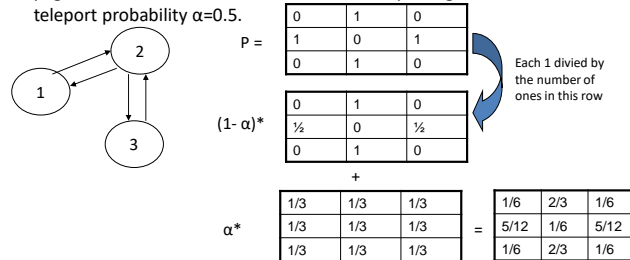
|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| $d_1$ | 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| $d_2$ | 1 | 0 | 1 | 1 | 0 | 0 | 0 |
| $d_3$ | 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| $d_4$ | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| $d_5$ | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| $d_6$ | 0 | 0 | 0 | 1 | 1 | 0 | 1 |

# Transition matrix with teleporting

|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 0.00 | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| $d_1$ | 0.00 | 0.50 | 0.50 | 0.00 | 0.00 | 0.00 | 0.00 |
| $d_2$ | 0.33 | 0.00 | 0.33 | 0.33 | 0.00 | 0.00 | 0.00 |
| $d_3$ | 0.00 | 0.00 | 0.00 | 0.50 | 0.50 | 0.00 | 0.00 |
| $d_4$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.00 |
| $d_5$ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.50 | 0.50 |
| $d_6$ | 0.00 | 0.00 | 0.00 | 0.33 | 0.33 | 0.00 | 0.33 |

α = 0.14

P =

|       | $d_0$ | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $d_0$ | 0.02 | 0.02 | 0.88 | 0.02 | 0.02 | 0.02 | 0.02 |
| $d_1$ | 0.02 | 0.45 | 0.45 | 0.02 | 0.02 | 0.02 | 0.02 |
| $d_2$ | 0.31 | 0.02 | 0.31 | 0.31 | 0.02 | 0.02 | 0.02 |
| $d_3$ | 0.02 | 0.02 | 0.02 | 0.45 | 0.45 | 0.02 | 0.02 |
| $d_4$ | 0.02 | 0.02 | 0.02 | 0.02 | 0.02 | 0.02 | 0.88 |
| $d_5$ | 0.02 | 0.02 | 0.02 | 0.02 | 0.02 | 0.45 | 0.45 |
| $d_6$ | 0.02 | 0.02 | 0.02 | 0.31 | 0.31 | 0.02 | 0.31 |

## Power method convergence

| | x | xP¹ | xP² | xP³ | xP⁴ | xP⁵ | xP⁶ | xP⁷ | xP⁸ | xP⁹ | xP¹⁰ | xP¹¹ | xP¹² | xP¹³ |
|-------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| $d_0$ | 0.14 | 0.06 | 0.09 | 0.07 | 0.07 | 0.06 | 0.06 | 0.06 | 0.06 | 0.05 | 0.05 | 0.05 | 0.05 | 0.05 |
| $d_1$ | 0.14 | 0.08 | 0.06 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| $d_2$ | 0.14 | 0.25 | 0.18 | 0.17 | 0.15 | 0.14 | 0.13 | 0.12 | 0.12 | 0.12 | 0.12 | 0.11 | 0.11 | 0.11 |
| $d_3$ | 0.14 | 0.16 | 0.23 | 0.24 | 0.24 | 0.24 | 0.24 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 |
| $d_4$ | 0.14 | 0.12 | 0.16 | 0.19 | 0.19 | 0.20 | 0.21 | 0.21 | 0.21 | 0.21 | 0.21 | 0.21 | 0.21 | 0.21 |
| $d_5$ | 0.14 | 0.08 | 0.06 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| $d_6$ | 0.14 | 0.25 | 0.23 | 0.25 | 0.27 | 0.28 | 0.29 | 0.29 | 0.30 | 0.30 | 0.30 | 0.30 | 0.31 | 0.31 |

## Pagerank summary

- Preprocessing:
  - Given graph of links, build matrix **P**.
  - From it compute **a**.
  - The entry $a_i$ is a number between 0 and 1: the pagerank of page $i$.
- Query processing:
  - Retrieve pages meeting query.
  - Rank them by their pagerank.
  - Order is **query-*independent***.

## PageRank issues

- Real surfers are not random surfers – Markov model is not a good model of surfing.
  - Issues: back button, short vs. long paths, bookmarks, directories – and search!
- Simple PageRank ranking (as described on previous slide) produces bad results for many pages.
  - Consider the query *video service*
  - The Yahoo home page (i) has a very high PageRank and (ii) contains both words.
  - If we rank all Boolean hits according to PageRank, then the Yahoo home page would be top-ranked.
  - Clearly not desirable
- In practice: rank according to weighted combination of many factors, including raw text match, anchor text match, PageRank and many other factors

## How important is PageRank?

- Frequent claim: PageRank is the most important component of web ranking.
- The reality:
  - There are several components that are at least as important: e.g., anchor text, indexing , zone weighting, phrases ...
- Rumor has it that PageRank in his original form (as presented here) now has a negligible impact on ranking!
- However, variants of a page's PageRank are still an essential part of ranking.
- Addressing link spam is difficult and crucial.

What is PageRank?