

Information Retrieval & Social Web

CS 525/DS 595

Worcester Polytechnic Institute

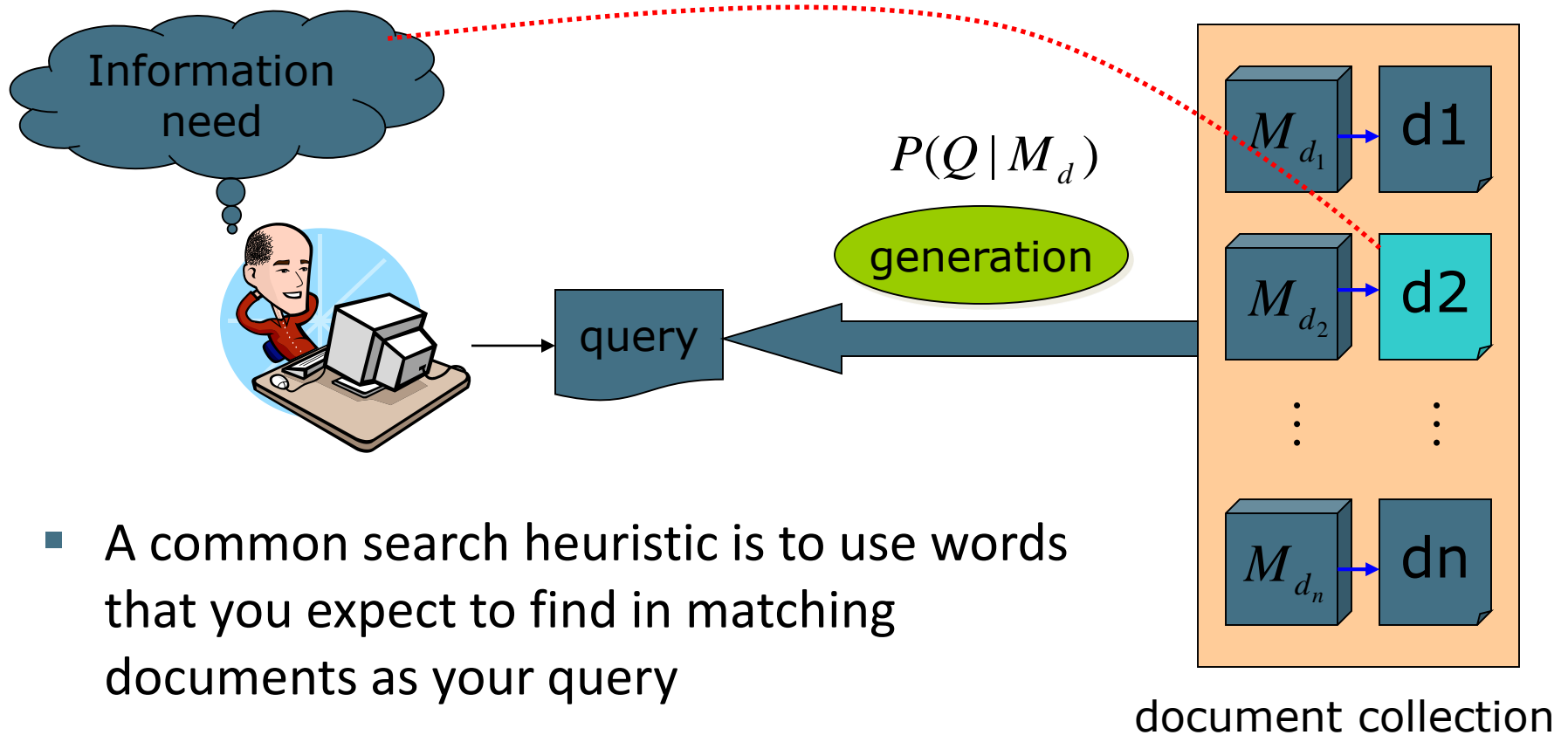
Department of Computer Science

Instructor: Prof. Kyumin Lee

Previous Class...

Statistical Language
Models

IR based on Language Model (LM)



- 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
 - Many bitly's or diggs
 - (Pagerank)



Quantitative

Today: Link Analysis

- Anchor text
- PageRank

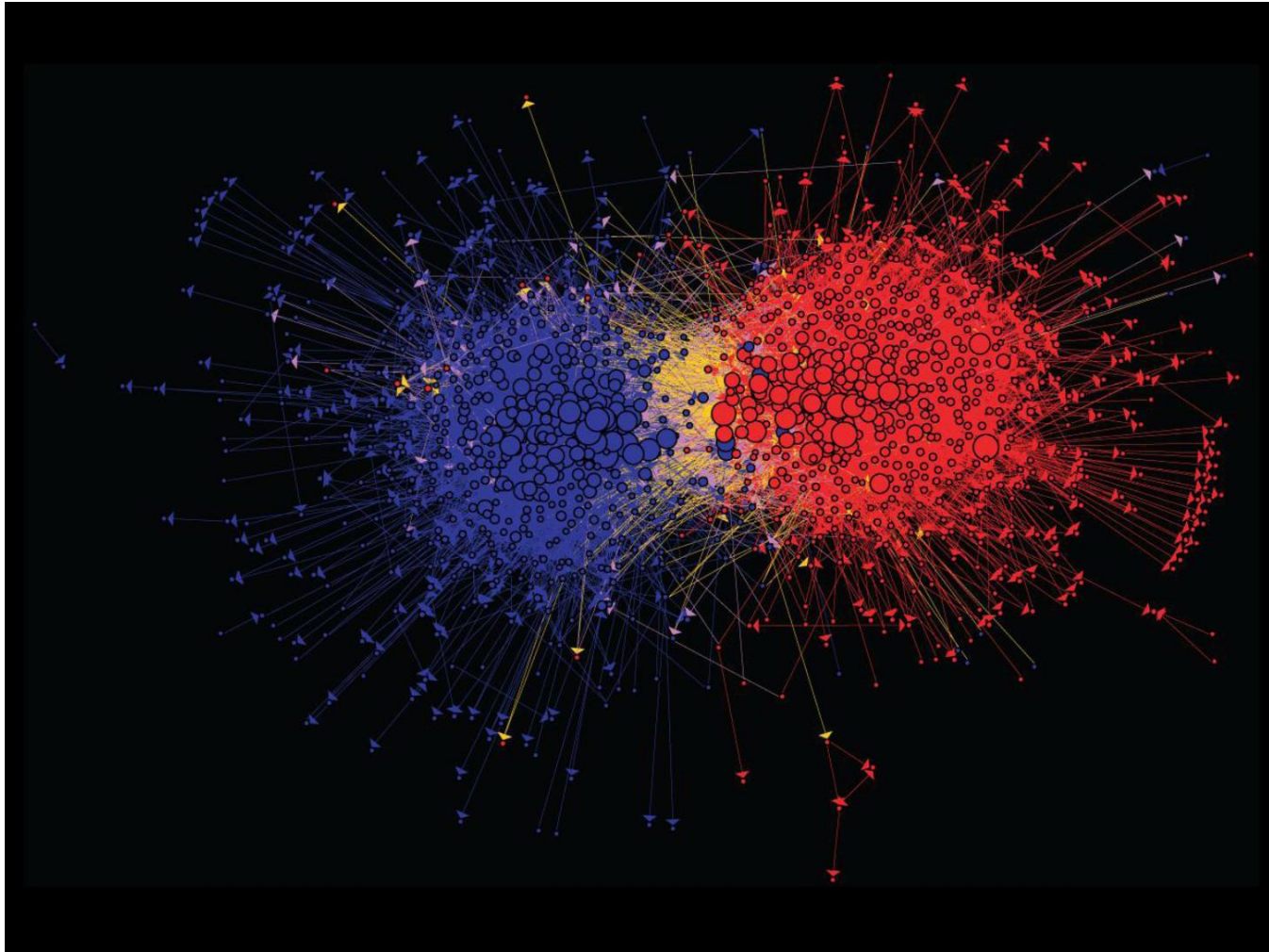
Graph Data: Social Networks



Facebook social graph

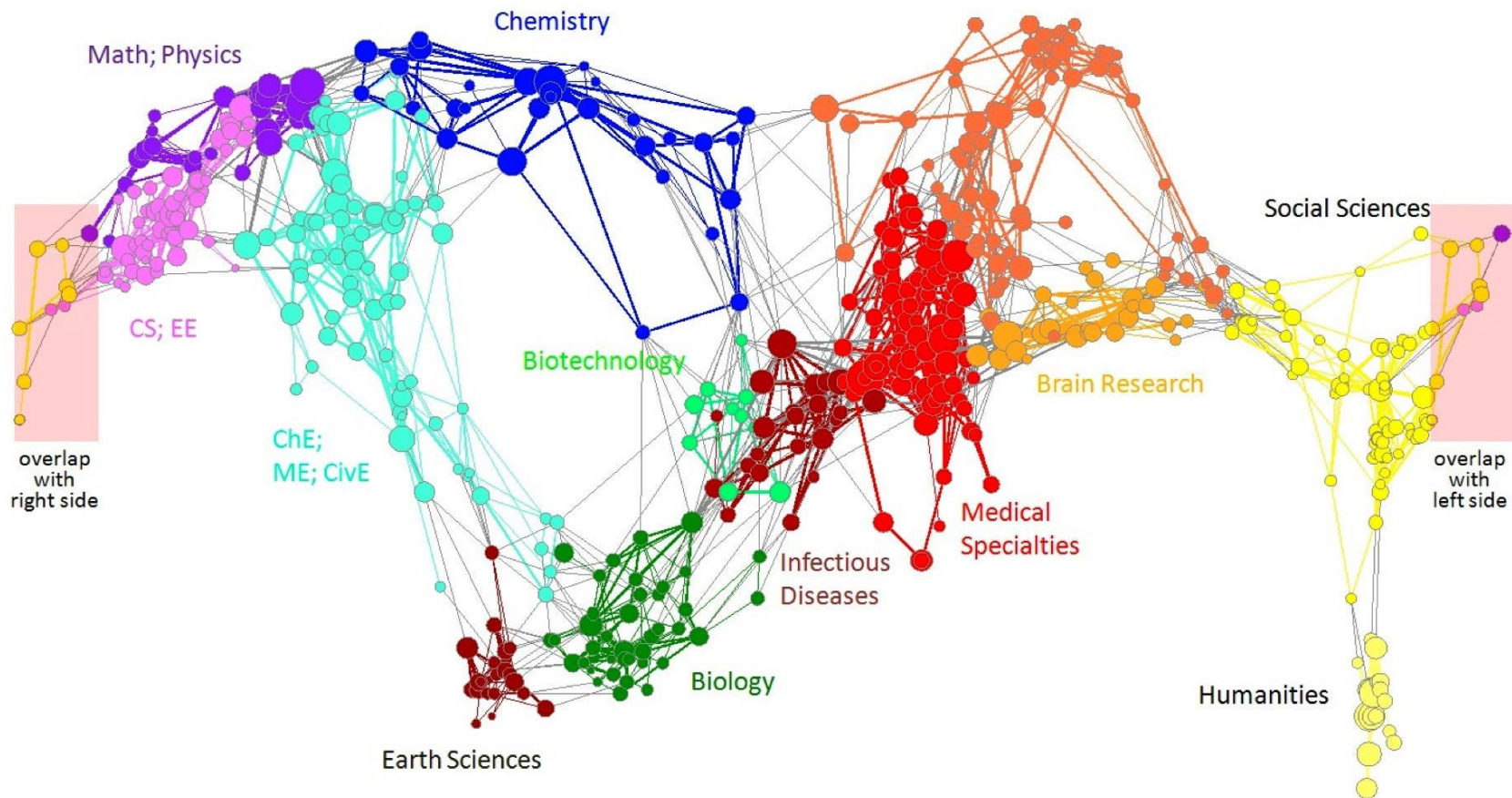
4-degrees of separation [Backstrom-Boldi-Rosa-Ugander-Vigna, 2011]

Graph Data: Media Networks



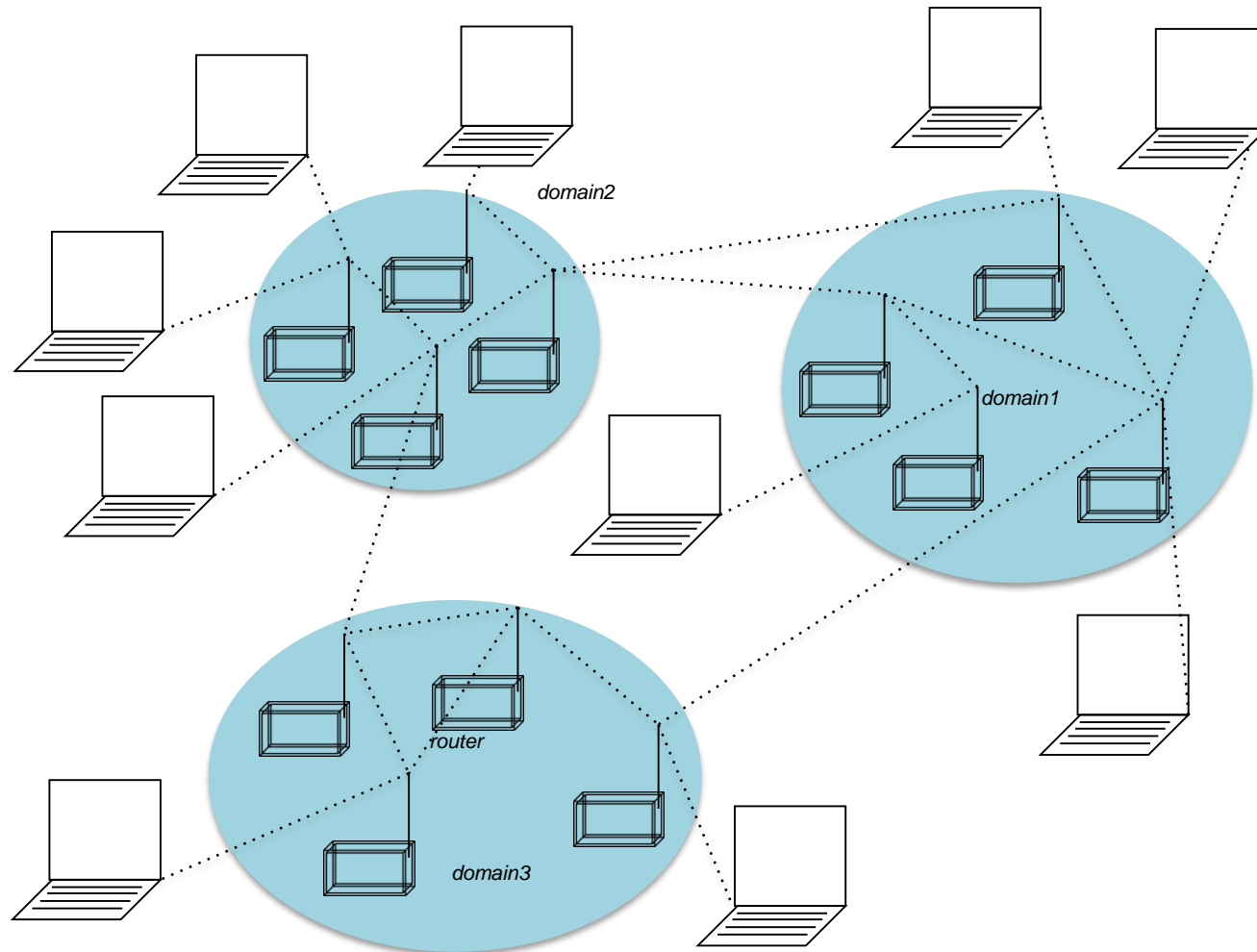
Connections between political blogs
Polarization of the network [Adamic-Glance, 2005]

Graph Data: Information Nets



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

Graph Data: Communication Nets



Web as a Graph

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

I teach a
class on
Networks.

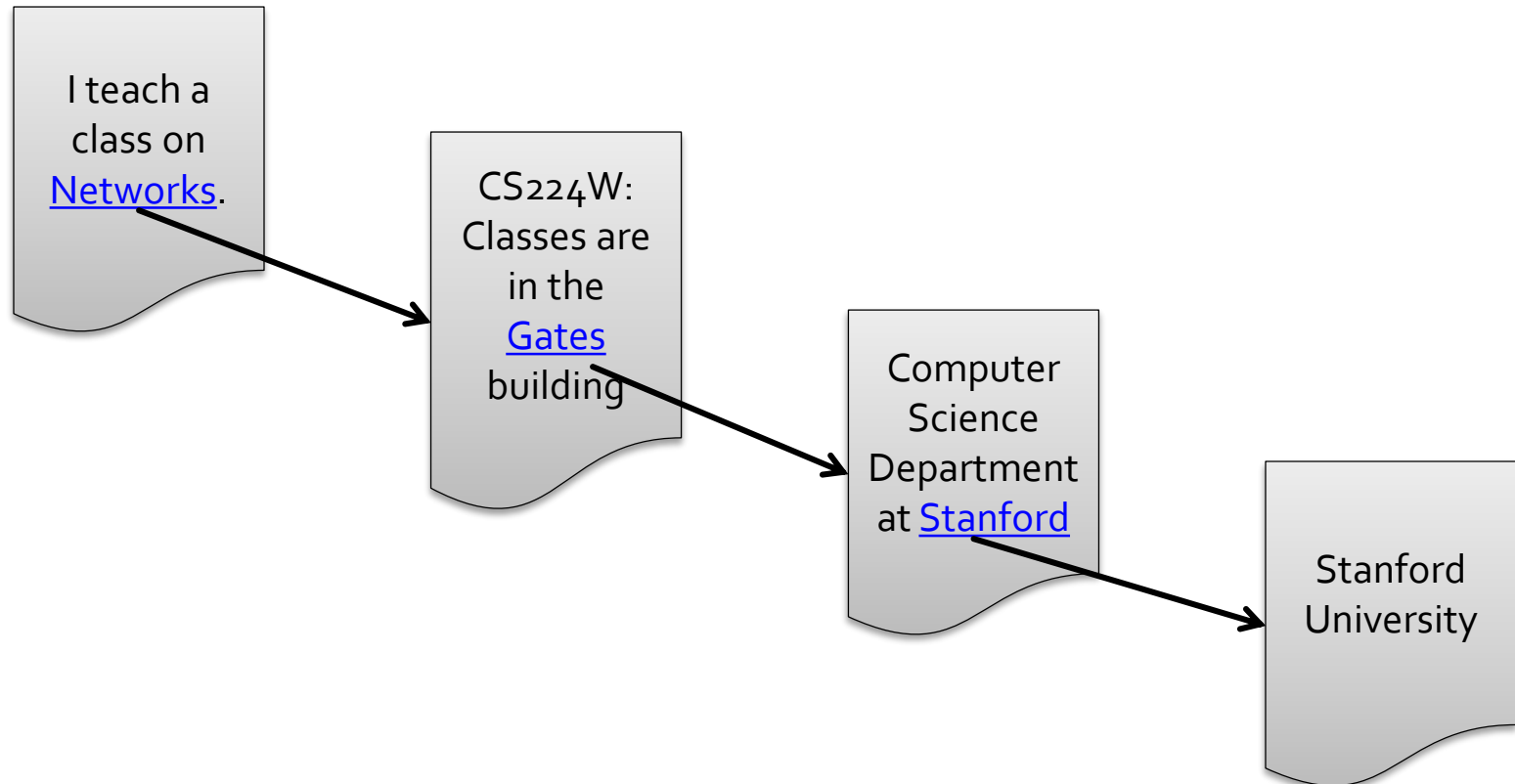
CS224W:
Classes are
in the
Gates
building

Computer
Science
Department
at Stanford

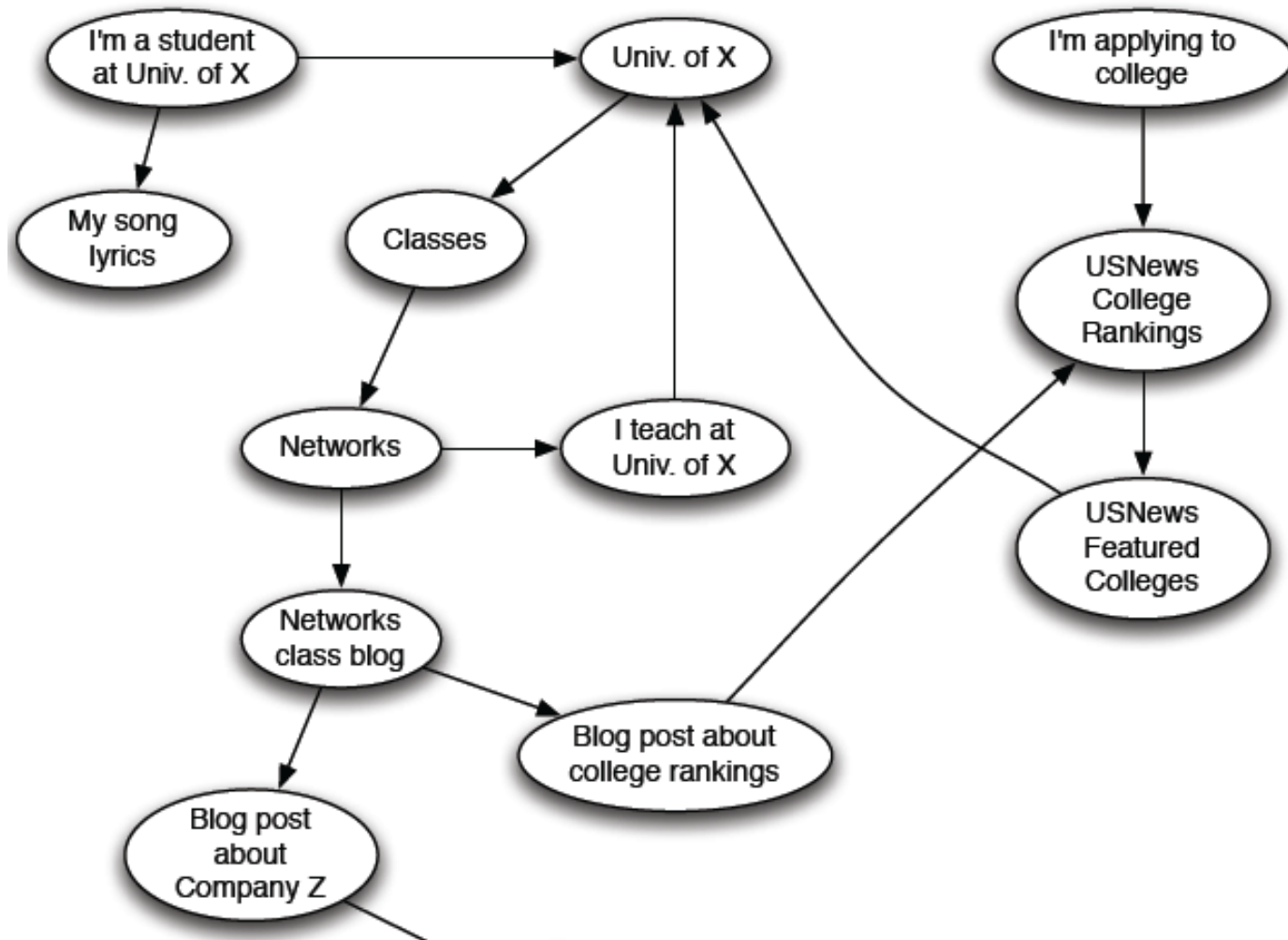
Stanford
University

Web as a Graph

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



Web as a Directed Graph



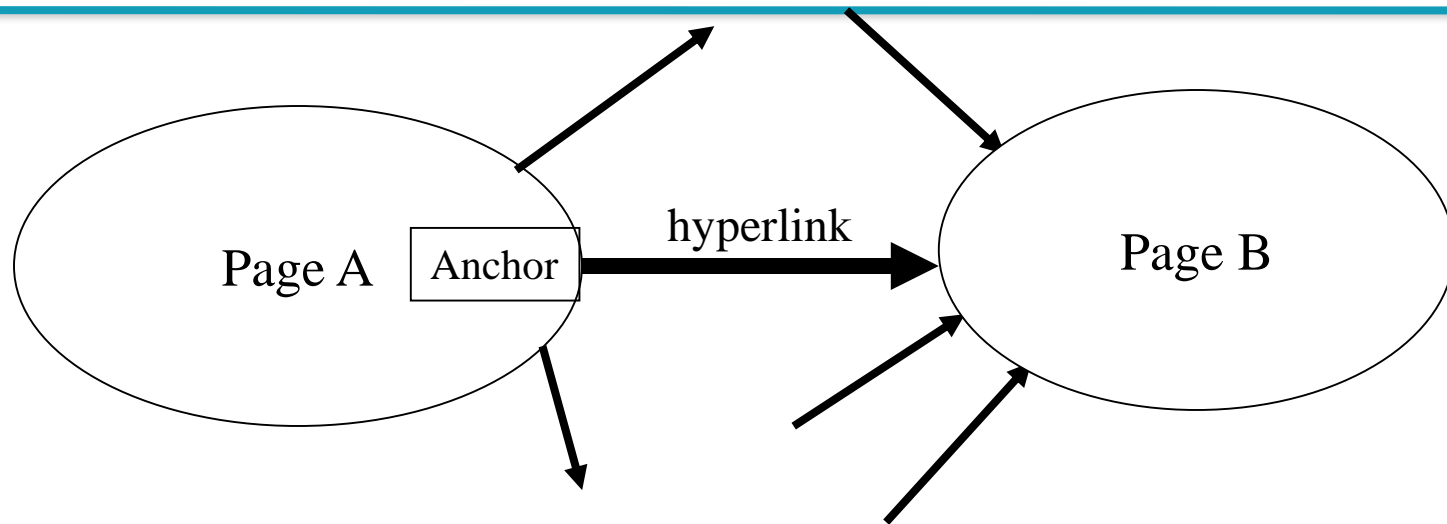
Broad Question

- **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 here”

[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”



www.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

Results 1 - 10 of about 969,000 for [miserable failure](#). (0.06 seconds)

[Biography of President George W. Bush](#)

Biography of the president from the official White House web site.

www.whitehouse.gov/president/gwbbio.html - 29k - [Cached](#) - [Similar pages](#)

[Past Presidents](#) - [Kids Only](#) - [Current News](#) - [President](#)

[More results from www.whitehouse.gov »](#)

[Welcome to MichaelMoore.com!](#)

Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ...

www.michaelmoore.com/ - 35k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

[BBC NEWS | Americas | 'Miserable failure' links to Bush](#)

Web users manipulate a popular search engine so an unflattering description leads to the president's page.

news.bbc.co.uk/2/hi/americas/3298443.stm - 31k - [Cached](#) - [Similar pages](#)

[Google's \(and Inktomi's\) Miserable Failure](#)

A search for **miserable failure** on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ...

searchenginewatch.com/sereport/article.php/3296101 - 45k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)



News Front Page



Africa

Americas

Asia-Pacific

Europe

Middle East

South Asia

UK

Business

Health

Science & Environment

Technology

Entertainment

Also in the news

Video and Audio

Programmes

Have Your Say

In Pictures

Country Profiles

Special Reports


RELATED BBC SITES


SPORT

WEATHER

ON THIS DAY

Last Updated: Sunday, 7 December, 2003, 15:04 GMT

 E-mail this to a friend

 Printable version

'Miserable failure' links to Bush

George W Bush has been Google bombed.

Web users entering the words "miserable failure" into the popular search engine are directed to the biography of the president on the White House website.

The trick is possible because Google searches more than just the contents of web pages - it also counts how often a site is linked to, and with what words.

Thus, members of an online community can affect the results of Google searches - called "Google bombing" - by linking their sites to a chosen one.

Weblogger Adam Mathes is credited with inventing the practice in 2001, when he used it to link the phrase "talentless hack" to a friend's website.

The search engine can be manipulated by a fairly small group of users, one report suggested.

Newsday newspaper says as few as 32 web pages with the words "miserable failure" link to the Bush biography.

The Bush administration has



Bush has been the target of similar pranks before

SEE ALSO:

- ▶ WMD spoof is internet hit
04 Jul 03 | West Midlands
- ▶ Google hit by link bombers
13 Mar 02 | Science/Nature

RELATED INTERNET LINKS:

- ▶ White House
- ▶ Google bombing

The BBC is not responsible for the content of external internet sites

TOP AMERICAS STORIES

- ▶ US lifts lid on WikiLeaks probe
- ▶ Iran scientist heads home
- ▶ Argentina legalises gay marriage

 | News feeds

“ If you are George Bush and typed the country's name in the address bar, make sure that it is spelled correctly (IRAQ) ”

Prank website

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!

ALTAVISTA

Technology, Inc.

View Multimedia From Our Vantage Point



Buy and insure new cars & trucks online

**Car Buying & Car Insurance
Pain Relief**



[Click here for advertising information - reach millions every month!](#)

Search and Display the Results

Search with Digital's Alta Vista [[Advanced Search](#)]



Add your URL to the most
popular Search Engines



**Free
Software**

Download Now...



Contests

Make Me Laugh...



Creative Web

Create a Site...

**FREE
WEB
SITES !**

[Create Your Personal Web Page For Free With Howdy!](#)

**FREE
WEB
SITES !**

ALTAVISTA
Technology, Inc.

[\[Creative\]](#)[\[Search\]](#)[\[Humor\]](#)

Search for information about:

in

Infoseek Guide is best viewed with:



Want personalized news? [Get Personal now!](#)

Basic Search Tips:

- Click in the box above and type a few words that describe what you want to find. For example, typing **growing orchids indoors** will find sites about caring for orchids.
- If you are looking for a person or place, type the name, starting with capital letters. For example, typing **Florence Italy** will find sites about this famous city.
- These detailed [search tips](#) describe how to use the features of Infoseek Guide to find what you are looking for.
- For the broadest results, you can search the entire **World Wide Web**.
- To restrict your search to hand-picked and categorized sites, choose **Infoseek Select Sites**.
- Or just search for a category within Infoseek Select by choosing **Categories of Sites**.
- To search through Internet discussion forums (similar to bulletin boards), choose **Usenet Newsgroups**.
- To search for someone's e-mail address, choose **E-mail Addresses**.
- To search through news stories within the past month, choose **Reuters News**.
- To search through answers to Frequently Asked Questions, choose **Web FAQs**.

Explore these popular Infoseek Select topics:

- ▶ [Arts & Entertainment](#)
- ▶ [Business & Finance](#)
- ▶ [Computers & Internet](#)
- ▶ [Education](#)
- ▶ [Government & Politics](#)
- ▶ [Health & Medicine](#)
- ▶ [Living](#)
- ▶ [News](#)
- ▶ [Reference](#)
- ▶ [Science & Technology](#)
- ▶ [Sports](#)
- ▶ [Travel](#)

Try [Infoseek Personal](#), your personalized news service

[Customize](#) | [Add Site](#) | [Help](#) | [Feedback](#)

[Download iSeek](#) | [About Infoseek](#)



[Click here to try Microsoft Money 97 FREE](#)



**It's amazing where
Go Get It will get you.**

Find:

Go Get It

[Enhance your search.](#)



[New Search](#) • [TopNews](#) • [Sites by Subject](#) • [Top 5% Sites](#) • [City Guide](#) • [Pictures & Sounds](#)
[PeopleFind](#) • [Point Review](#) • [Road Maps](#) • [Software](#) • [About Lycos](#) • [Club Lycos](#) • [Help](#)

[Add Your Site to Lycos](#)

Copyright © 1996 Lycos[™], Inc. All Rights Reserved.
Lycos is a trademark of Carnegie Mellon University.

[Questions & Comments](#)



search



reviews



city.net



live!



tours



news

people finder

maps

yellow pages



"Turbo Search!"

[Download
Excite Direct](#)[Take an
ExciteSeeing Tour](#)[Excite on TV](#)[Make your website
searchable, FREE!](#)**Excite Search:** twice the power of the competition.

What:

search

Where:

World Wide Web

[\[Help\]](#)[\[Advanced Search\]](#)INTEGRATED BROWSING, EMAIL,
NEWSGROUPS AND PAGE CREATION.**Excite Reviews:** site reviews by the web's [best editorial team](#).

- [Arts](#)
- [Business](#)
- [Computing](#)
- [Education](#)
- [Entertainment](#)
- [Health](#)
- [Hobbies](#)
- [Life & Style](#)
- [Money](#)
- [News & Reference](#)
- [Personal Pages](#)
- [Politics & Law](#)
- [Regional](#)
- [Science](#)
- [Shopping](#)
- [Sports](#)

Excite City.NetPlan your weekend, your
travels.

Find-A-Destination

[Take me there!](#)
[Maps](#) ◦ [Top Cities](#) ◦
[Concierge](#)
ExciteSeeing Tours

Choose from hundreds.

- [X-Files: The truth is out there!](#)
- [Dr. Ruth's guide to safer sex](#)
- [Windows 95 shareware and freeware](#)
- [Celebrating Thanksgiving](#)
- [Investing in high-tech stocks](#)
- [New to the Net?](#)

Excite Live!

Your news, your way.

- [Latest news](#)
- [Sports scores](#)
- [Local weather](#)
- [Movie reviews](#)
- [Stock quotes](#)
- [TV listings](#)
- [Horoscopes](#)
- [Site reviews](#)

Excite Reference

Just the facts, ma'am.

- [Yellow Pages](#)
- [People Finder](#)
- [Email Lookup](#)
- [Maps](#)
- [Shareware](#)
- [Dictionary](#)



Search the web using Google!

10 results ▾

Google Search

I'm feeling lucky

Index contains ~25 million pages (soon to be much bigger)

About Google!

[Stanford Search](#) [Linux Search](#)

Get Google! updates monthly!

your e-mail

Subscribe

[Archive](#)

Copyright ©1997-8 Stanford University



Search the web using Google!

Google Search

I'm feeling lucky

Special Searches

[Stanford Search](#)

[Linux Search](#)

[Help!](#)

[About Google!](#)

[Company Info](#)

[Google! Logos](#)

Get Google!

updates monthly:

your e-mail


Subscribe

[Archive](#)

Copyright ©1998 Google Inc.



 [Jobs@Google](#)

 [About Google](#)

Search the web using Google

[Google Launches! Read the press release.](#)

©1999 Google Inc.

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

All

News

Shopping

Images

Apps

More ▾

Search tools

About 1,460,000,000 results (0.33 seconds)

Trade up to a new iPhone

Ad www.apple.com/ ▾

Trade in your current smartphone and get up to \$350 in credit.

Get instant credit · Get a gift card

iPhone 6s

The only thing that's changed is everything. [Learn more.](#)

Buy now

Order now and get free shipping.
Or choose free in-store pickup.

In the news



Used iPhone 6 could be the bargain you're looking for

CNET · 15 hours agoThis is especially true of Apple's **iPhone**. But when is the best time to get the great deal?[iPhone 7 Plus to Boast Dual Rear Cameras: Report](#)**PetaPixel** · 10 hours ago[Apple Loop: New iPhone Leaks, iPad Air 3's Launch Date, iOS 9.2.1 Reveals Secret iPhone Powers](#)**Forbes** · 1 day ago[More news for iphone](#)

iPhone - Apple

www.apple.com/iphone/ ▾ **Apple** ▾**iPhone 6s**. With the most powerful technology and most intuitive operating system ever. It's here, and yours to explore.[Compare iPhone Models](#) · [Where to buy iPhone](#) · [Accessories](#) · [iPhone in Business](#)

Apple iPhone 6s - 64 GB - Rose Gold - Verizon - CDMA/GSM



4.8 ★★★★★ 4,312 user reviews

Shop now

Sponsored ⓘ

Rose Gold ▾ **64 GB** ▾ **Verizon - CDMA/GSM** ▾**\$299.00** · [Apple Store](#)
With contract

Free shipping

\$299.99 · [Best Buy](#)

Free shipping

[View all sellers and prices](#)

worcester polytechnic institute academic calendar



All

Shopping

News

Images

Maps

More

Settings

Tools

About 119,000 results (0.60 seconds)

Academic Calendar & Catalogs - Worcester Polytechnic Institute

<https://www.wpi.edu/academics/calendar-catalogs> ▼

The information on this page is accurate as of the date of publication. However, all future **academic** calendars are reviewed annually, published for planning purposes, and are subject to change. Important Dates. Feb15. **Academic** Advising Day. 8:00 am to 11:00 pm. Feb23. Reading Day. 8:00 am to 11:00 pm. Mar2.

[PDF] Undergraduate (PDF)

https://www.wpi.edu/sites/default/files/UG_17-18_20170612.pdf ▼

Jun 12, 2017 - **CALENDAR**. 2017-2018. S M T W R F S S M T W R F S. JUL 16 17 18 19 20 21 22. 4. 5 6 7 8 9 10. 23 24 25 26 27 28 29 FEB 11 12 13 14 15 16 17 FEBRUARY 15. ACAD. ADV. DAY. (PROJ. OPPORTUNITIES). 30 31 1 2 3 4 5. 18 19 20 21 22 23 24 FEBRUARY 23. READING/MAKEUP DAY. 6 7 8 9 10 11 ...

University Calendar - Worcester Polytechnic Institute

<https://www.wpi.edu/news/calendar> ▼

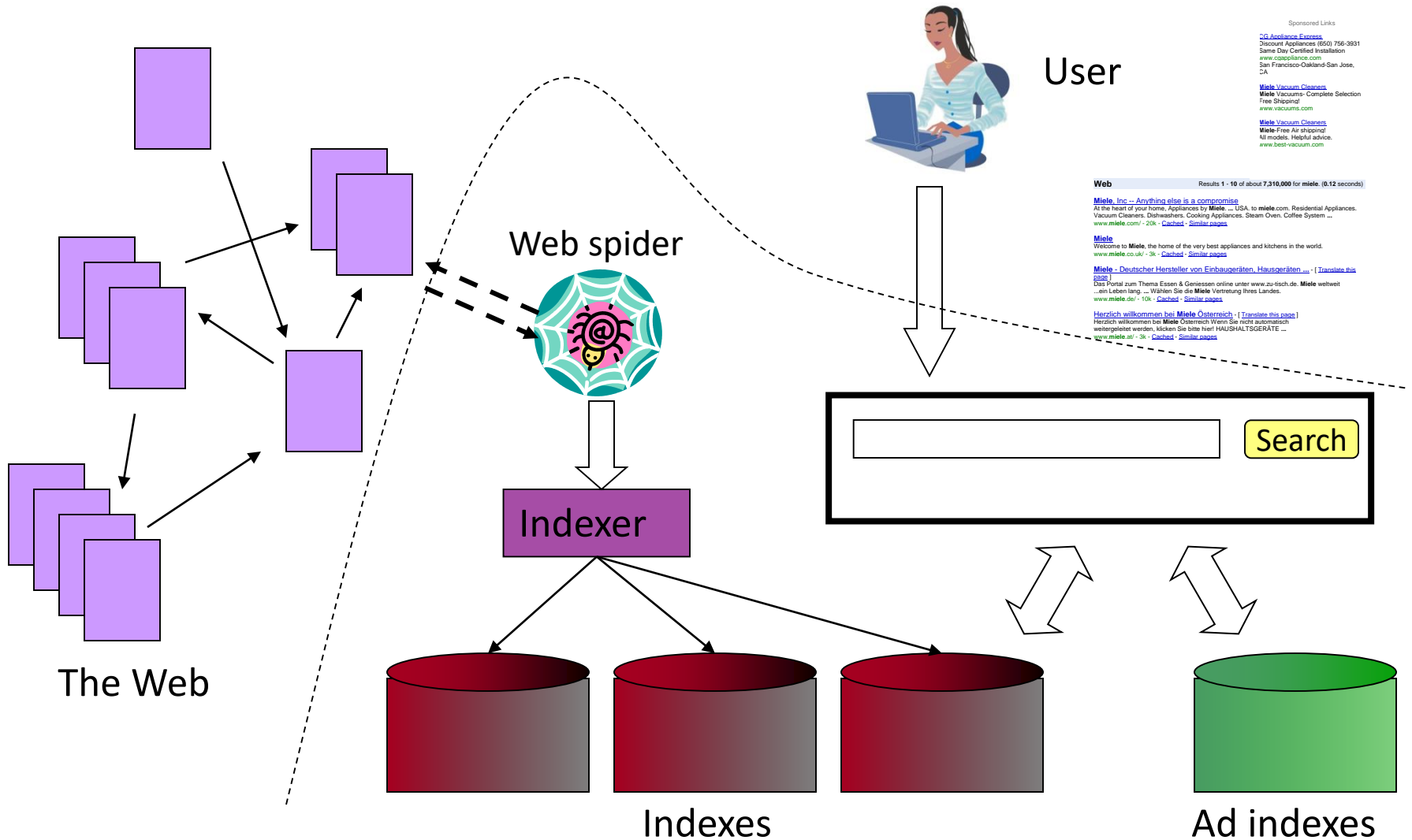
Academic Calendars · Varsity Athletics **Calendar** · Annual Events · Campus Dining · Residence Halls · Add Your Event. WPI in the World. Global Impact Program. Global Projects Program · About WPI · Bookstore · Canvas · Careers · Directories · Library · Offices · Worcester. **WORCESTER POLYTECHNIC INSTITUTE**

[PDF] undergraduate calendar 2018-2019

<https://www.wpi.edu/.../Academic.../Academic-Calendar/Future%20Calendar%20-%20...> ▼

UNDERGRADUATE. **CALENDAR**. 2017-2018. S M T W R F S. S M T W R F S. JUL 16 17 18 19 20 21 22. 4. 5 6 7 8 9 10. 23 24 25 26 27 28 29. FEB 11 12 13 14 15 16 17. FEBRUARY 15. ACAD. ADV. DAY. (PROJ. OPPORTUNITIES). 30 31 1 2 3 4 5. 18 19 20 21 22 23 24. FEBRUARY 23. READING DAY. 6 7 8 9 10 11 12.

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?

Mechanical Turk is a marketplace for work.

We give businesses and developers access to an on-demand, scalable workforce.
Workers select from thousands of tasks and work whenever it's convenient.

162,119 HITs available. [View them now.](#)

Make Money by working on HITs

HITs - *Human Intelligence Tasks* - are individual tasks that you work on. [Find HITs now.](#)

As a Mechanical Turk Worker you:

- Can work from home
- Choose your own work hours
- Get paid for doing good work

Find an
interesting task

Work

Earn
money



[Find HITs Now](#)

or [learn more about being a Worker](#)

Get Results from Mechanical Turk Workers

Ask workers to complete HITs - *Human Intelligence Tasks* - and get results using Mechanical Turk. [Register Now](#)

As a Mechanical Turk Requester you:

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results

Fund your
account

Load your
tasks

Get
results

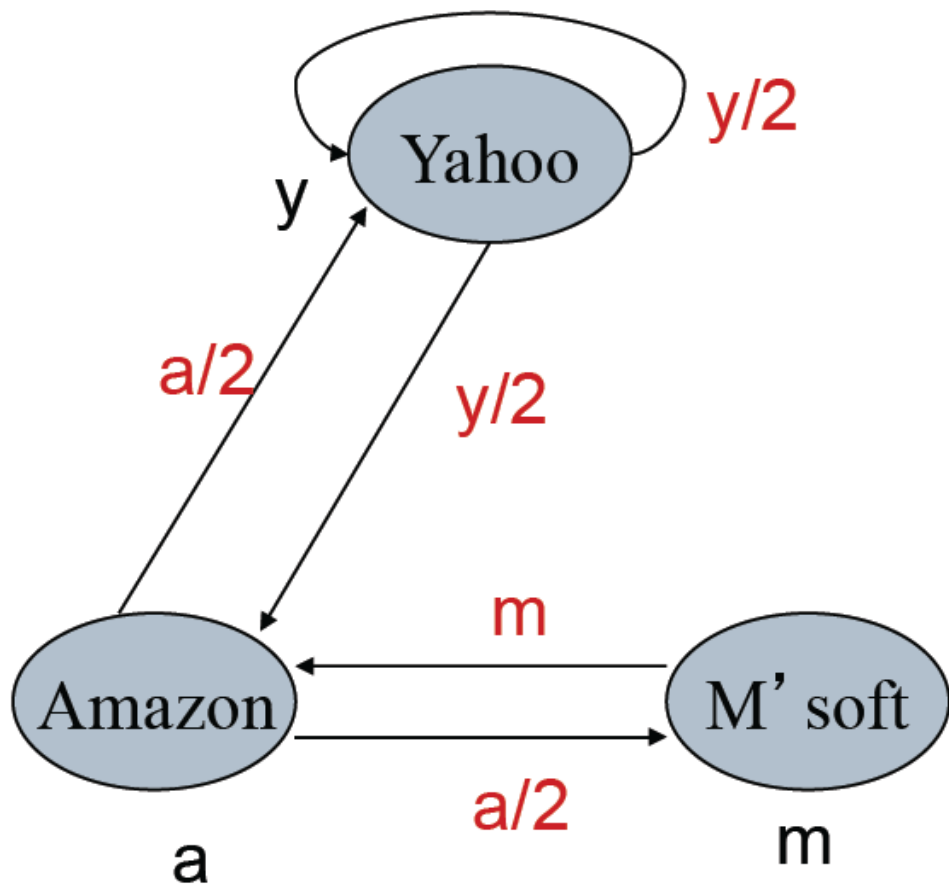


[Get Started](#)

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



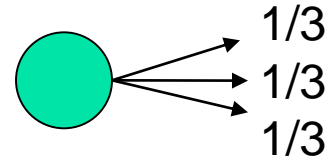
$$y = y/2 + a/2$$

$$a = y/2 + m$$

$$m = a/2$$

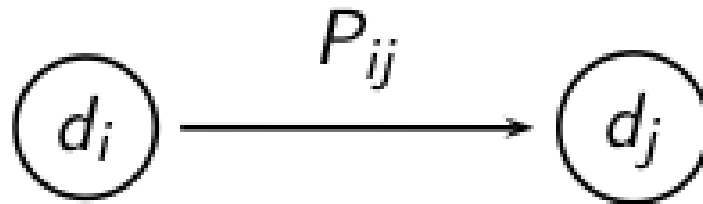
PageRank basics

- 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
- “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

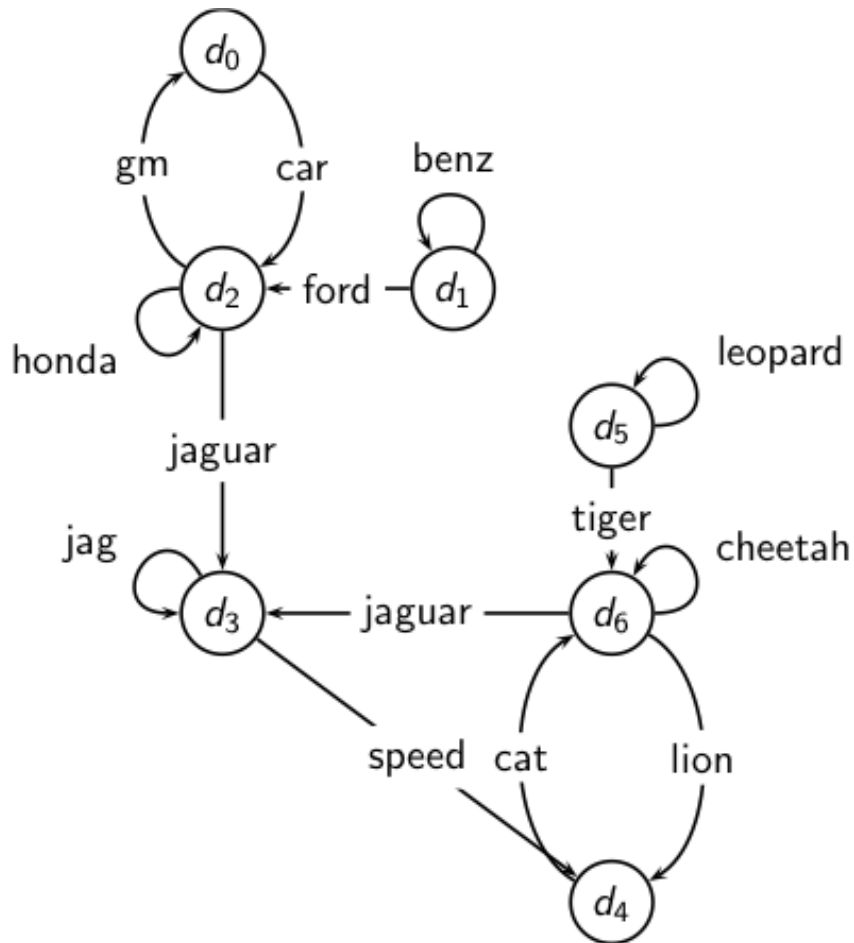
- A Markov chain consists of n states, plus an $n \times n$ transition probability matrix \mathbf{P} .
- **state = page**
- At each step, we are on exactly one of the states.
- For $1 \leq i, j \leq n$, the matrix entry P_{ij} tells us the probability of j being the next state (page), given we are currently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



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



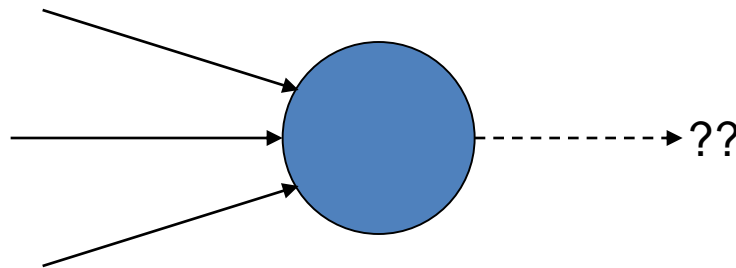
Transition probability matrix							
	d_0	d_1	d_2	d_3	d_4	d_5	d_6
d_0	0.00	0.00	1.00	0.00	0.00	0.00	0.00
d_1	0.00	0.50	0.50	0.00	0.00	0.00	0.00
d_2	0.33	0.00	0.33	0.33	0.00	0.00	0.00
d_3	0.00	0.00	0.00	0.50	0.50	0.00	0.00
d_4	0.00	0.00	0.00	0.00	0.00	0.00	1.00
d_5	0.00	0.00	0.00	0.00	0.00	0.50	0.50
d_6	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., $(\underset{1}{000}\dots\underset{i}{1}\dots\underset{n}{000})$ means we're in state i .
- More generally, the vector $\mathbf{x} = (x_1, \dots, x_n)$ means the walk is in state i with probability x_i .

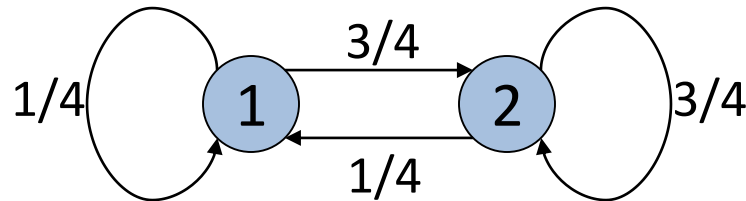
$$\sum_{i=1}^n x_i = 1.$$

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 \mathbf{P} tells us where we go next from state i .
- So from \mathbf{x} , our next state is distributed as \mathbf{xP} .

Steady state example

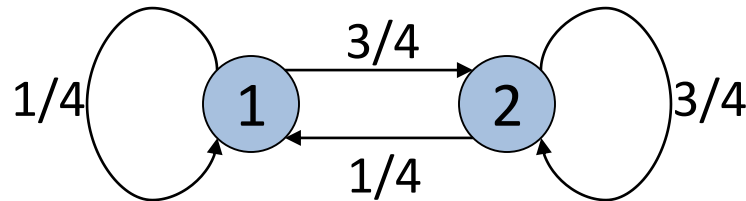
- The steady state looks like a vector of probabilities $\mathbf{a} = (a_1, \dots, a_n)$:
- a_i is the probability that we are in state i .



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

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

$$\vec{x}_2P = (0.25, 0.75)$$

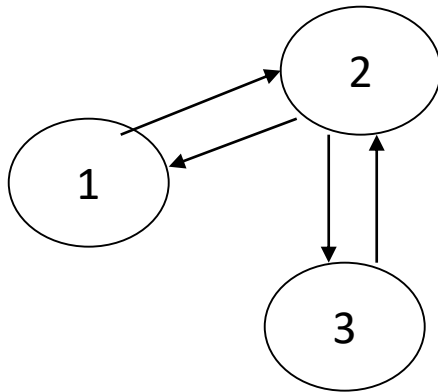
Convergence in one iteration!

Exercise on PageRank

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

$$P = (1 - \alpha) * \text{transition matrix} + \alpha * \text{teleporting matrix}$$

- 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 $\alpha=0.5$.



$P =$

0	1	0
1	0	1
0	1	0

$(1 - \alpha)^*$

0	1	0
$\frac{1}{2}$	0	$\frac{1}{2}$
0	1	0

+

α^*

$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

=

$\frac{1}{6}$	$\frac{2}{3}$	$\frac{1}{6}$
$\frac{5}{12}$	$\frac{1}{6}$	$\frac{5}{12}$
$\frac{1}{6}$	$\frac{2}{3}$	$\frac{1}{6}$

Each 1 divided by the number of ones in this row

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

$$\vec{x}_0 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

P=

1/6	2/3	1/6
5/12	1/6	5/12
1/6	2/3	1/6

$$\vec{x}_1 = \begin{bmatrix} 1/6 & 2/3 & 1/6 \end{bmatrix}$$

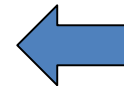
$$\vec{x}_2 = \begin{bmatrix} 1/3 & 1/3 & 1/3 \end{bmatrix}$$

$$\vec{x}_3 = \begin{bmatrix} 1/4 & 1/2 & 1/4 \end{bmatrix}$$

...

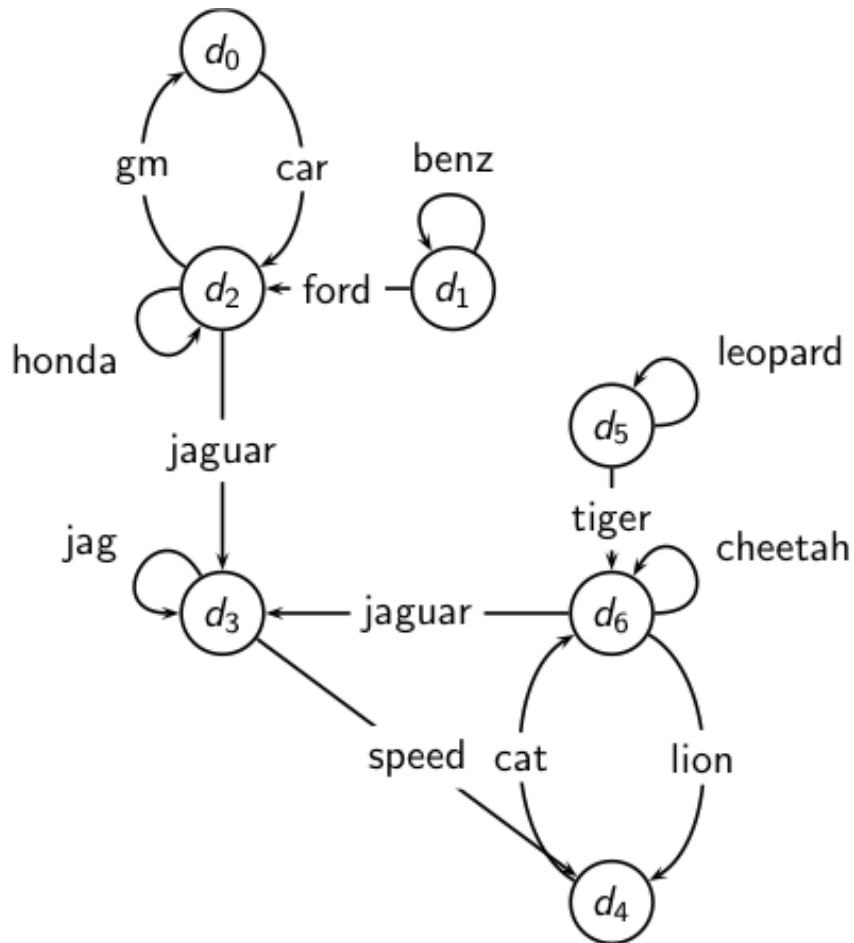
...

$$\vec{x}_k = \begin{bmatrix} 5/18 & 4/9 & 5/18 \end{bmatrix}$$



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



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

$\alpha = 0.14$

Power method convergence

	x	xP^1	xP^2	xP^3	xP^4	xP^5	xP^6	xP^7	xP^8	xP^9	xP^{10}	xP^{11}	xP^{12}	xP^{13}
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 \mathbf{P} .
 - From it compute \mathbf{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?

