

Large-Scale Programming and Testing

Search Engines by Croft, Metzler, and Strohman, 2009

Fall 2017 – CSCI 4963/6963 – Week 05

David Goldschmidt – goldschmidt@gmail.com

Office: Amos Eaton 115

Office hours: Mon/Thu 1:00-1:50PM; Wed 1:00-2:50PM

Politeness policies

- Web crawlers adhere to politeness policies
 - In general, a crawler sends a **GET** request to a specific webserver every few seconds (or minutes)
 - A website might have a **robots.txt** file, which specifies what crawlers can and cannot crawl
 - Also may specify one or more sitemaps...

Example robots.txt file

```
User-agent: *
Disallow: /private/
Disallow: /confidential/
Disallow: /other/
Allow: /other/public

User-agent: FavoredCrawler
Disallow:

Sitemap: http://xyz.com/sitemap.xml.gz
```

Sitemaps

- Default priority is set at 0.5
- Some URLs might not be easily discovered by crawlers...

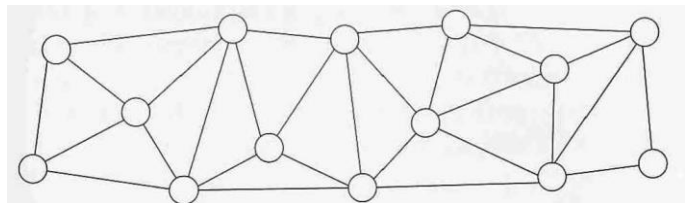
Example sitemap.xml file

```
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
  <url>
    <loc>http://www.company.com/</loc>
    <lastmod>2008-01-15</lastmod>
    <changefreq>monthly</changefreq>
    <priority>0.7</priority>
  </url>
  <url>
    <loc>http://www.company.com/inventory?item=iphone</loc>
    <changefreq>weekly</changefreq>
  </url>
  <url>
    <loc>http://www.company.com/news/</loc>
    <changefreq>daily</changefreq>
  </url>
  ...
</urlset>
```



Are you connected?

- The Internet (1969) is a network that is:
 - Global
 - Decentralized
 - Redundant
 - Heterarchical
 - Ever-changing
- Made up of many smaller networks, which in turn are made up of varying numbers of machines
- Made up of many different types of machines



Weaving the Web

- The World Wide Web (1989) is:
 - Global
 - Decentralized
 - Redundant (sometimes)
 - Heterarchical
 - Ever-changing
 - Made up of many websites, which in turn are composed of varying numbers of webpages
 - Made up of many different types of webpages (static versus dynamic)

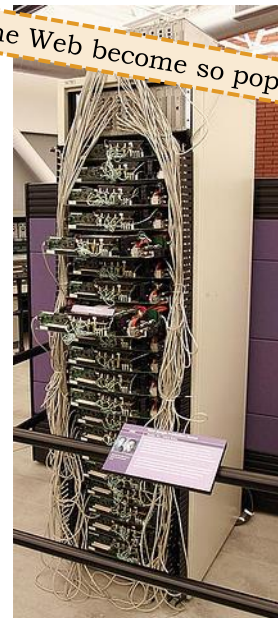


Who invented the Web?

- Sir Tim Berners-Lee at CERN



Why did the Web become so popular?



Why are links so important?

- Links (i.e., URLs) are useful to us humans for navigating websites and finding things
- Links are also extremely useful to search engines



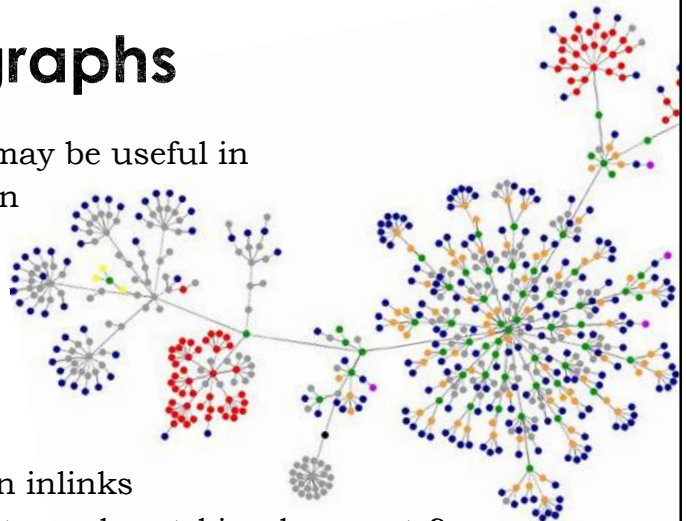
```
<a href="http://nytimes.com">the latest news</a>
```

- Anchor (or link) text helps with ranking
 - Link text summarizes the content of the destination page
 - Link text is succinct, descriptive, and often coincides with query text
 - Link text is often written by a non-biased third party (less spam?)



Links and webgraphs

- The actual links themselves may be useful in describing a target webpage in terms of the target's:
 - Popularity
 - Importance
 - Authority
 - Incoming link (inlink) count
- Link analysis often focuses on inlinks
 - How can we use inlinks to better rank matching documents?
 - The challenge here is how we actually obtain the inlinks



PageRank

- PageRank is an iterative link analysis algorithm
 - Use “the link structure...to calculate a quality ranking for each web page”

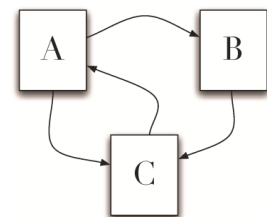


- Original Google paper by the original Google guys:
<http://infolab.stanford.edu/~backrub/google.html>

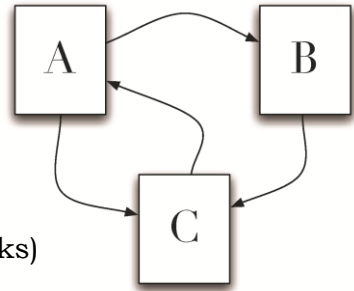


How does PageRank work?

- The PageRank of page A (i.e., $PR(A)$) is the probability that a “random surfer” visits that page
- Simulate browsing the Web as a random surfer:
 - Set constant λ
 - Choose a random number r between 0 and 1
 - If $r < \lambda$ then go to a random page
 - Otherwise, follow a random link from the current page
 - Go back to Step 2
- Navigating to a random page avoids getting stuck in pages with no links, with broken links, or that form cycles



PageRank



- PageRank of page C is the probability that a random surfer is viewing page C:
 - Based on both inlinks and outgoing links (outlinks)
 - $C(A)$ is the number of outlinks from page A
 - $PR(C) = PR(A)/C(A) + PR(B)/C(B)$
- We assume PR forms an even distribution across all pages
 - Initially, $PR(A) = PR(B) = PR(C) = 0.333$
 - Then, $PR(C) = 0.333/2 + 0.333/1 = 0.500$
 - And $PR(B) = PR(A)/C(A) = 0.333/2 = 0.166$
 - And $PR(A) = PR(C)/C(C) = 0.500/1 = 0.500$

Keep going until convergence!



PageRank

- For each page u :
$$PR(u) = \sum_{v \in B_u} \frac{PR(v)}{L_v}$$
 - Here, B_u is the set of pages that point to page u
 - And, L_v is the number of deduplicated outlinks from page v
- We can account for the “random jumps” by incorporating constant λ into the equation (and N is number of pages):

$$PR(u) = \frac{\lambda}{N} + (1 - \lambda) \sum_{v \in B_u} \frac{PR(v)}{L_v}$$

Typically, λ is low, e.g., 0.15



```

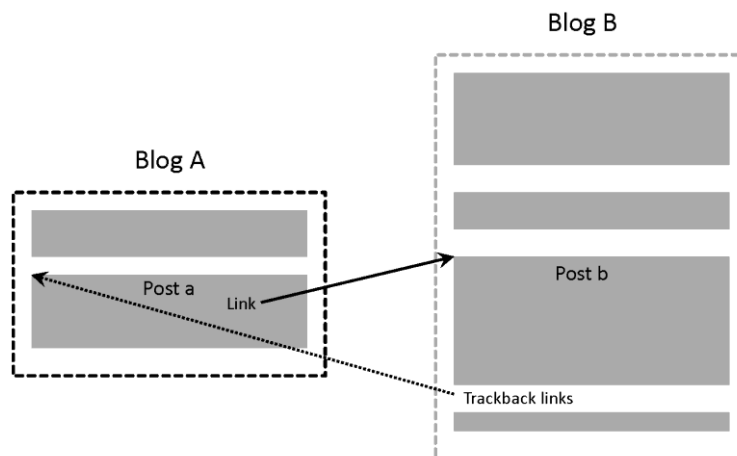
1: procedure PAGERANK( $G$ )
2:    $\triangleright G$  is the web graph, consisting of vertices (pages) and edges (links).
3:    $(P, L) \leftarrow G$   $\triangleright$  Split graph into pages and links
4:    $I \leftarrow$  a vector of length  $|P|$   $\triangleright$  The current PageRank estimate
5:    $R \leftarrow$  a vector of length  $|P|$   $\triangleright$  The resulting better PageRank estimate
6:   for all entries  $I_i \in I$  do
7:      $I_i \leftarrow 1/|P|$   $\triangleright$  Start with each page being equally likely
8:   end for
9:   while  $R$  has not converged do
10:    for all entries  $R_i \in R$  do
11:       $R_i \leftarrow \lambda/|P|$   $\triangleright$  Each page has a  $\lambda/|P|$  chance of random selection
12:    end for
13:    for all pages  $p \in P$  do
14:       $Q \leftarrow$  the set of pages  $p$  such that  $(p, q) \in L$  and  $q \in P$ 
15:      if  $|Q| > 0$  then
16:        for all pages  $q \in Q$  do
17:           $R_q \leftarrow R_q + (1 - \lambda)I_p/|Q|$   $\triangleright$  Probability  $I_p$  of being a page  $p$ 
18:        end for
19:      else
20:        for all pages  $q \in P$  do
21:           $R_p \leftarrow R_p + (1 - \lambda)I_p/|P|$ 
22:        end for
23:      end if
24:       $I \leftarrow R$   $\triangleright$  Update our current PageRank estimate
25:    end for
26:  end while
27:  return  $R$ 
28: end procedure

```

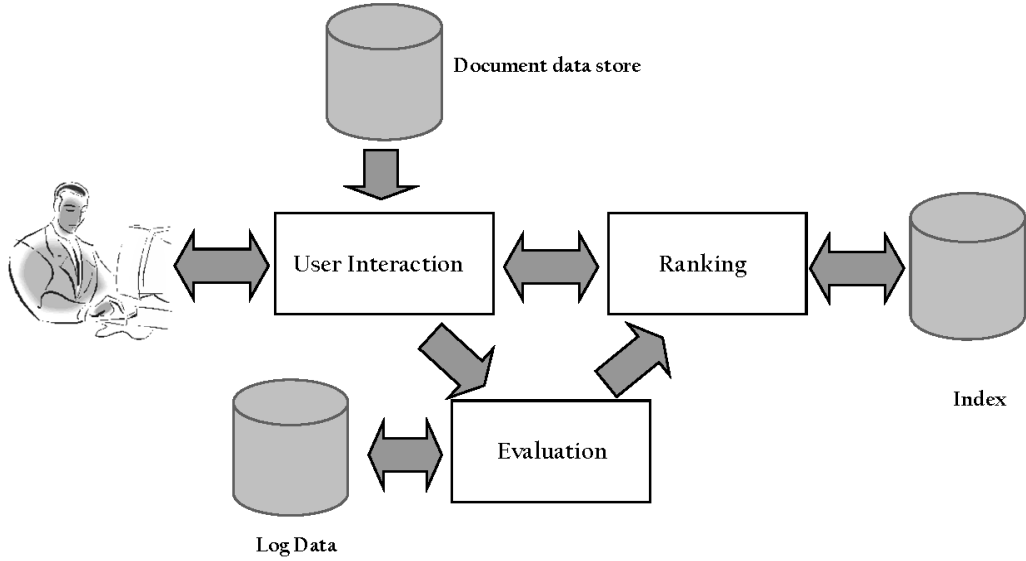


Link quality and avoiding spam

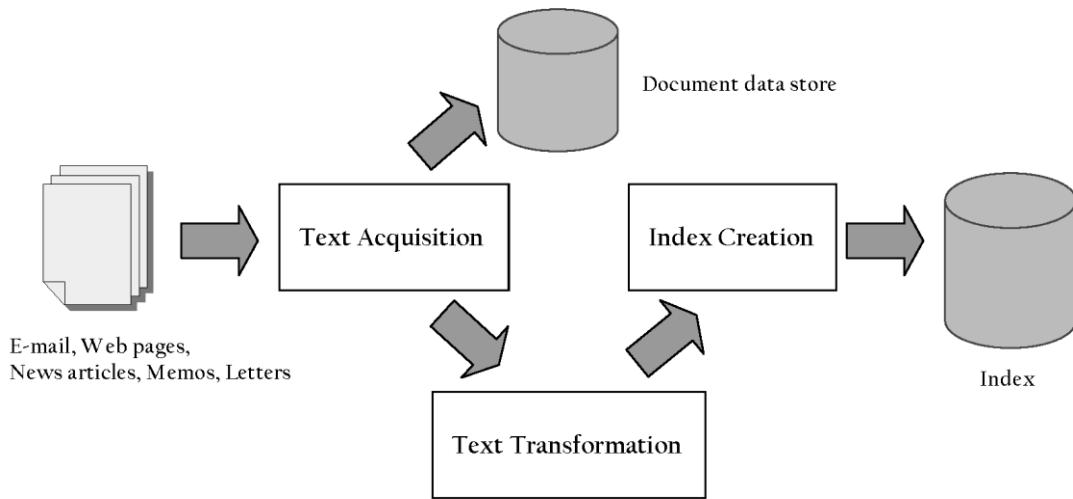
- A cycle tends to negate the effectiveness of the PR algorithm



Ranking and querying



Acquisition and indexing



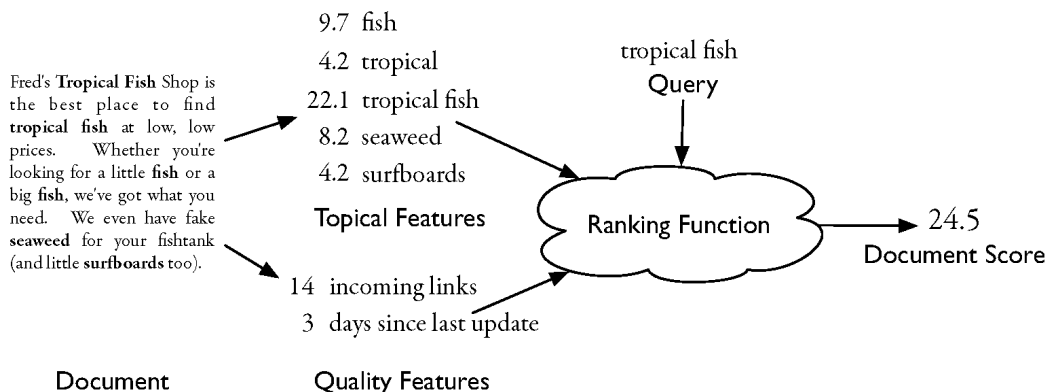
Document features

- Document features are obtained during text transformation
 - A document feature is some extractable characteristic of the document expressed numerically
 - For example, a topical feature estimates the degree to which the document is about a particular topic
 - As another example, quality features include inlink counts, the number of days since a page was last updated, etc.
 - Document features can translate into index terms



Abstract model of ranking

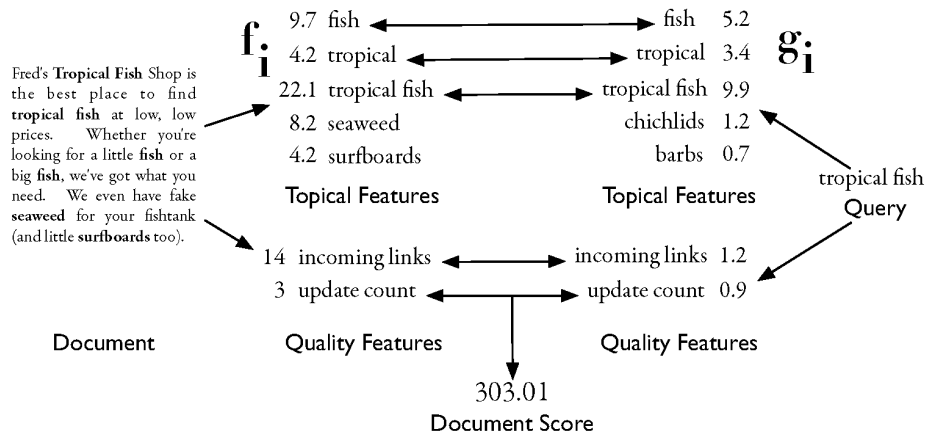
- Regardless of the ranking function, the abstract model below provides an overview of implementation



A more concrete ranking model

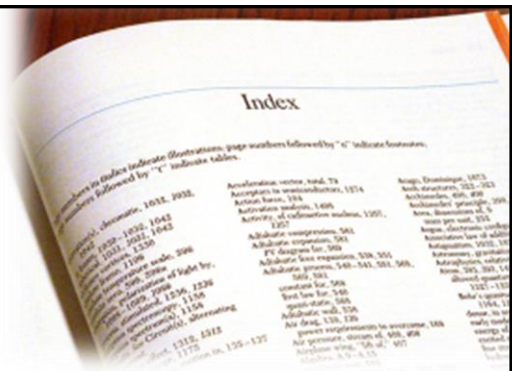
$$R(Q, D) = \sum_i g_i(Q) f_i(D)$$

f_i is a document feature function
 g_i is a query feature function



Indexing

- An index is a data structure designed to make search very fast and efficient
- Text-based search typically requires an inverted index
 - The index is inverted because we associate documents with words rather than identifying words within or as part of a document
 - Each index term is associated with an inverted list that contains:
 - A list of documents
 - A list of word occurrences within each document
 - Word counts and positional information
 - Metadata identifying semantics



Inverted indexes

- Each entry in an inverted index is called a posting
 - The part of the posting that refers to a specific document or location is called a pointer
 - Further, each document within the collection is given a unique document number
 - Lists are usually document-ordered
 - Sorted by document number



Example document collection

- Assume each sentence below is a separate document
 - S_1 Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.
 - S_2 Fishkeepers often use the term tropical fish to refer only those requiring fresh water, with saltwater tropical fish referred to as marine fish.
 - S_3 Tropical fish are popular aquarium fish, due to their often bright coloration.
 - S_4 In freshwater fish, this coloration typically derives from iridescence, while salt water fish are generally pigmented.



Example

- Inverted index for given documents $S_1, S_2, S_3,$ and S_4
- Individual word occurrences have been deduplicated and ignored
- What does this structure tell us?
- What is missing?

and	1	only	2
aquarium	3	pigmented	4
are	3 4	popular	3
around	1	refer	2
as	2	referred	2
both	1	requiring	2
bright	3	salt	1 4
coloration	3 4	saltwater	2
derives	4	species	1
due	3	term	2
environments	1	the	1 2
fish	1 2 3 4	their	3
fishkeepers	2	this	4
found	1	those	2
fresh	2	to	2 3
freshwater	1 4	tropical	1 2 3
from	4	typically	4
generally	4	use	2
in	1 4	water	1 2 4
include	1	while	4
including	1	with	2
iridescence	4	world	1
marine	2		
often	2 3		



Example

- Inverted index for given documents $S_1, S_2, S_3,$ and $S_4,$ with word counts
- What does this structure tell us?
- What is missing?

and	1:1	only	2:1
aquarium	3:1	pigmented	4:1
are	3:1 4:1	popular	3:1
around	1:1	refer	2:1
as	2:1	referred	2:1
both	1:1	requiring	2:1
bright	3:1	salt	1:1 4:1
coloration	3:1 4:1	saltwater	2:1
derives	4:1	species	1:1
due	3:1	term	2:1
environments	1:1	the	1:1 2:1
fish	1:2 2:3 3:2 4:2	their	3:1
fishkeepers	2:1	this	4:1
found	1:1	those	2:1
fresh	2:1	to	2:2 3:1
freshwater	1:1 4:1	tropical	1:2 2:2 3:1
from	4:1	typically	4:1
generally	4:1	use	2:1
in	1:1 4:1	water	1:1 2:1 4:1
include	1:1	while	4:1
including	1:1	with	2:1
iridescence	4:1	world	1:1
marine	2:1		
often	2:1 3:1		



Document fields

- A document field is a section of a document that has identifiable additional semantic meaning
 - e.g., date, from:, to:, etc.
 - e.g., title, author, copyright, publisher, isbn, etc.
- Implementation options:
 - Use separate inverted lists for each field
 - Add extra information about fields to postings
 - Use extent lists....



Extent lists

- An extent is defined as a contiguous region of a document (typically with special meaning)
 - We can represent extents using word position ranges
 - The inverted list records all extents for a given field

