



Web Crawler

Lecture 3

CS 4250 – Web Search and Recommender Systems

1

Web Crawler

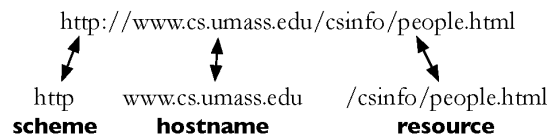
- **Finds** and **downloads** web pages automatically
 - provides the collection for searching
- Web is huge and constantly growing
- Web is not under the control of search engine providers
- Web pages are constantly changing
- Crawlers also used for other types of data

CS 4250 – Web Search and Recommender Systems

2

Retrieving Web Pages

- Every page has a unique *uniform resource locator* (URL)
- Web pages are stored on web servers that use HTTP to exchange information with client software
- e.g.,

`http://www.cs.umass.edu/csinfo/people.html`

scheme **hostname** **resource**

CS 4250 – Web Search and Recommender Systems

3

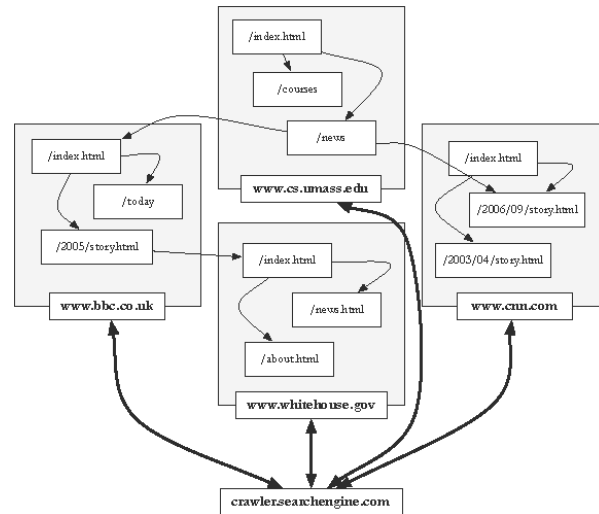
Retrieving Web Pages

- Web crawler client program connects to a *domain name system* (DNS) server
- DNS server translates the hostname into an *internet protocol* (IP) address
- Crawler then attempts to connect to server host using specific *port*
- After connection, crawler sends an HTTP request to the web server to request a page
 - usually a GET request

CS 4250 – Web Search and Recommender Systems

4

Crawling the Web



CS 4250 – Web Search and Recommender Systems

5

Web Crawler

- Starts with a set of *seeds*, which are a set of URLs given to it as parameters
- Seeds are added to a URL request queue
- Crawler starts fetching pages from the request queue
- Downloaded pages are parsed to find link tags that might contain other useful URLs to fetch
- New URLs added to the crawler's request queue, or *frontier*
- Continue until no more new URLs or disk full

CS 4250 – Web Search and Recommender Systems

6

Web Crawling

- Web crawlers spend a lot of time waiting for responses to requests
- To reduce this inefficiency, web crawlers use threads and fetch hundreds of pages at once
- Crawlers could potentially flood sites with requests for pages
- To avoid this problem, web crawlers use *politeness policies*
 - e.g., delay between requests to same web server

CS 4250 – Web Search and Recommender Systems

7

Politeness policies

- Reasonable web crawlers
 - do not fetch more than one page at a time from a particular web server
 - wait at least a few seconds between requests to the same web server.
- Single queue per web server
- Question: *Suppose a web crawler can fetch 100 pages each second, and that its politeness policy dictates that it cannot fetch more than one page each 30 seconds from a particular web server. URLs from how many different Web servers should be in the request queue at any given time to achieve high throughput?*
 - **> 3000**

CS 4250 – Web Search and Recommender Systems

8

Controlling Crawling

- Even crawling a site slowly will anger some web server administrators, who object to any copying of their data
- Robots.txt file can be used to control crawlers

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

User-agent: FavoredCrawler
Disallow:

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

CS 4250 – Web Search and Recommender Systems

9

Simple Crawler Thread

```
procedure CRAWLERTHREAD(frontier)
  while not frontier.done() do
    website ← frontier.nextSite()
    url ← website.nextURL()
    if website.permitsCrawl(url) then
      text ← retrieveURL(url)
      storeDocument(url, text)
      for each url in parse(text) do
        frontier.addURL(url)
      end for
    end if
    frontier.releaseSite(website)
  end while
end procedure
```

CS 4250 – Web Search and Recommender Systems

10

Controlling Crawling

1. The crawling thread first retrieves a website from the frontier.
2. The crawler then identifies the next URL in the website's queue.
3. In `permitsCrawl`, the crawler checks to see if the URL is okay to crawl according to the website's `robots.txt` file.
4. If it can be crawled, the crawler uses `retrieveURL` to fetch the document contents.
5. This is the most expensive part of the loop, and the crawler thread may block here for many seconds.
6. Once the text has been retrieved, `storeDocument` stores the document text in a document database (discussed later in this chapter).
7. The document text is then parsed so that other URLs can be found.
8. These URLs are added to the frontier, which adds them to the appropriate website queues.
9. When all this is finished, the website object is returned to the frontier, which takes care to enforce its politeness policy by not giving the website to another crawler thread until an appropriate amount of time has passed.
10. In a real crawler, the timer would start immediately after the document was retrieved, since parsing and storing the document could take a long time.

CS 4250 – Web Search and Recommender Systems

11

Freshness

- Web pages are constantly being added, deleted, and modified
- Web crawler must continually revisit pages it has already crawled to see if they have changed in order to maintain the *freshness* of the document collection
 - *stale* copies no longer reflect the real contents of the web pages

CS 4250 – Web Search and Recommender Systems

12

Freshness

- HTTP protocol has a special request type called HEAD that makes it easy to check for page changes

- returns information about page, not page itself

Client request: HEAD /csinfo/people.html HTTP/1.1
Host: www.cs.umass.edu

HTTP/1.1 200 OK
Date: Thu, 03 Apr 2008 05:17:54 GMT
Server: Apache/2.0.52 (CentOS)
Last-Modified: Fri, 04 Jan 2008 15:28:39 GMT
Server response: ETag: "239c33-2576-2a2837c0"
Accept-Ranges: bytes
Content-Length: 9590
Connection: close
Content-Type: text/html; charset=ISO-8859-1

CS 4250 – Web Search and Recommender Systems

13

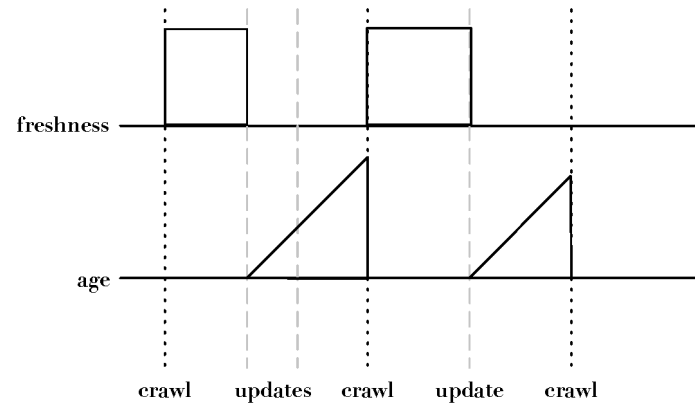
Freshness

- Not possible to constantly check all pages
 - must check important pages and pages that change frequently
- Freshness is the proportion of pages that are fresh
- Optimizing for this metric can lead to bad decisions, such as not crawling popular sites
- Age is a better metric

CS 4250 – Web Search and Recommender Systems

14

Freshness vs. Age



CS 4250 – Web Search and Recommender Systems

15

Freshness vs. Age

Freshness

$$F_p(t) = \begin{cases} 1 & \text{if } p \text{ is equal to the local copy at time } t \\ 0 & \text{otherwise} \end{cases}$$

Age

$$A_p(t) = \begin{cases} 0 & \text{if } p \text{ is not modified at time } t \\ t - \text{modification time of } p & \text{otherwise} \end{cases}$$

CS 4250 – Web Search and Recommender Systems

16

Focused Crawling

- Attempts to download only those pages that are about a particular topic
 - used by *vertical search* applications
- Rely on the fact that pages about a topic tend to have links to other pages on the same topic
 - popular pages for a topic are typically used as seeds
- Crawler uses *text classifier* to decide whether a page is on topic

CS 4250 – Web Search and Recommender Systems

17

Deep Web

- Sites that are difficult for a crawler to find are collectively referred to as the *deep* (or *hidden*) *Web*
 - much larger than conventional Web
- Three broad categories:
 - private sites
 - no incoming links, or may require log in with a valid account
 - form results
 - sites that can be reached only after entering some data into a form
 - scripted pages
 - pages that use JavaScript, Flash, or another client-side language to generate links

CS 4250 – Web Search and Recommender Systems

18

Sitemaps

- Sitemaps contain lists of URLs and data about those URLs, such as modification time and modification frequency
- Generated by web server administrators
- Tells crawler about pages it might not otherwise find
- Gives crawler a hint about when to check a page for changes

CS 4250 – Web Search and Recommender Systems

19

Sitemap Example

```
<?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/items?item=truck</loc>
    <changefreq>weekly</changefreq>
  </url>
  <url>
    <loc>http://www.company.com/items?item=bicycle</loc>
    <changefreq>daily</changefreq>
  </url>
</urlset>
```

CS 4250 – Web Search and Recommender Systems

20

Desktop Crawls

- Used for desktop search and enterprise search
- Differences to web crawling:
 - Much easier to find the data
 - Responding quickly to updates is more important
 - Must be conservative in terms of disk and CPU usage
 - Many different document formats
 - Data privacy very important

CS 4250 – Web Search and Recommender Systems

21

Document Feeds

- Many documents are *published*
 - created at a fixed time and rarely updated again
 - e.g., news articles, blog posts, press releases, email
- Published documents from a single source can be ordered in a sequence called a *document feed*
 - new documents found by examining the end of the feed

CS 4250 – Web Search and Recommender Systems

22

Document Feeds

- Two types:
 - A *push feed* alerts the subscriber to new documents
 - A *pull feed* requires the subscriber to check periodically for new documents
- Most common format for pull feeds is called *RSS*
 - Really Simple Syndication, RDF Site Summary, Rich Site Summary, ...

CS 4250 – Web Search and Recommender Systems

23

RSS Example

```
<?xml version="1.0"?>
<rss version="2.0">
  <channel>
    <title>Search Engine News</title>
    <link>http://www.search-engine-news.org/</link>
    <description>News about search engines.</description>
    <language>en-us</language>
    <pubDate>Tue, 19 Jun 2008 05:17:00 GMT</pubDate>
    <ttl>60</ttl>

    <item>
      <title>Upcoming SIGIR Conference</title>
      <link>http://www.sigir.org/conference</link>
      <description>The annual SIGIR conference is coming!
        Mark your calendars and check for cheap
        flights.</description>
      <pubDate>Tue, 05 Jun 2008 09:50:11 GMT</pubDate>
      <guid>http://search-engine-news.org#500</guid>
    </item>
```

CS 4250 – Web Search and Recommender Systems

24

RSS Example

```
...  
<item>  
  <title>New Search Engine Textbook</title>  
  <link>http://www.cs.umass.edu/search-book</link>  
  <description>A new textbook about search engines  
    will be published soon.</description>  
  <pubDate>Tue, 05 Jun 2008 09:33:01 GMT</pubDate>  
  <guid>http://search-engine-news.org#499</guid>  
</item>  
</channel>  
</rss>
```

CS 4250 – Web Search and Recommender Systems

25

RSS

- RSS feeds are accessed like web pages
 - using HTTP GET requests to web servers that host them
- Easy for crawlers to parse
- Easy to find new information

CS 4250 – Web Search and Recommender Systems

26

Conversion

- Text is stored in hundreds of incompatible file formats
 - e.g., raw text, RTF, HTML, XML, Microsoft Word, ODF, PDF
- Other types of files also important
 - e.g., PowerPoint, Excel
- Typically use a conversion tool
 - converts the document content into a tagged text format such as HTML or XML
 - retains some of the important formatting information

CS 4250 – Web Search and Recommender Systems

27

Character Encoding

- A character encoding is a mapping between bits and glyphs
 - i.e., getting from bits in a file to characters on a screen
 - Can be a major source of incompatibility
- ASCII is basic character encoding scheme for English
 - encodes 128 letters, numbers, special characters, and control characters in 7 bits, extended with an extra bit for storage in bytes

CS 4250 – Web Search and Recommender Systems

28

Character Encoding

- Other languages can have many more glyphs
 - e.g., Chinese has more than 40,000 characters, with over 3,000 in common use
- Many languages have multiple encoding schemes
 - e.g., CJK (Chinese-Japanese-Korean) family of East Asian languages, Hindi, Arabic
 - must specify encoding
- Unicode developed to address encoding problems

CS 4250 – Web Search and Recommender Systems

29

Unicode

- Single mapping from numbers to glyphs that attempts to include all glyphs in common use in all known languages
- Unicode is a mapping between numbers and glyphs
 - Most common: **UTF-8**, UTF-16, UTF-32

CS 4250 – Web Search and Recommender Systems

30

Unicode

- Proliferation of encodings comes from a need for compatibility and to save space
 - UTF-8 uses one byte for English (ASCII), as many as 4 bytes for some traditional Chinese characters
 - variable length encoding
 - UTF-32 uses 4 bytes for every character

CS 4250 – Web Search and Recommender Systems

31

Unicode

Decimal	Hexadecimal	Encoding			
0–127	0–7F	0xxxxxxx			
128–2047	80–7FF	110xxxxx	10xxxxxx		
2048–55295	800–D7FF	1110xxxx	10xxxxxx	10xxxxxx	
55296–57343	D800–DFFF	Undefined			
57344–65535	E000–FFFF	1110xxxx	10xxxxxx	10xxxxxx	
65536–1114111	10000–10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

- e.g., Greek letter pi (π) is Unicode symbol number 960
- In binary, 00000011 11000000 (3C0 in hexadecimal)
- Final encoding is 11001111 10000000 (CF80 in hexadecimal)

CS 4250 – Web Search and Recommender Systems

32

Storing the Documents

- Many reasons to store converted document text
 - provides efficient access to text for snippet generation, information extraction, etc.
 - keep copies of the documents around instead of trying to fetch them again the next time you want to build an index
- Database systems can provide document storage for some applications
 - web search engines use customized document storage systems

CS 4250 – Web Search and Recommender Systems

33

Storing the Documents

- Requirements for document storage system:
 - Random access
 - request the content of a document based on its URL
 - hash function based on URL is typical
 - Compression and large files
 - reducing storage requirements and efficient access
 - Update
 - handling large volumes of new and modified documents

CS 4250 – Web Search and Recommender Systems

34

Detecting Duplicates

- Duplicate and near-duplicate documents occur in many situations
 - Copies, versions, plagiarism, spam, mirror sites
 - 30% of the web pages in a large crawl are exact or near duplicates of pages in the other 70%
- Duplicates consume significant resources during crawling, indexing, and search
 - Little value to most users

CS 4250 – Web Search and Recommender Systems

35

Duplicate Detection

- *Exact* duplicate detection is relatively easy
 - *Checksum* techniques
 - A checksum is a value that is computed based on the content of the document
 - e.g., sum of the ASCII codes in the document file
- | | | | | | | | | | | | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|-----|
| T | r | o | p | i | c | a | l | | f | i | s | h | Sum |
| 54 | 72 | 6F | 70 | 69 | 63 | 61 | 6C | 20 | 66 | 69 | 73 | 68 | 508 |
- Possible for files with different text to have same checksum
 - Functions such as a *cyclic redundancy check* (CRC), have been developed that consider the positions of the bytes

CS 4250 – Web Search and Recommender Systems

36

Near-Duplicate Detection

- More challenging task
 - Are web pages with same text context but different advertising or format near-duplicates?
- A near-duplicate document is defined using a threshold value for some similarity measure between pairs of documents
 - e.g., document $D1$ is a near-duplicate of document $D2$ if more than 90% of the words in the documents are the same

CS 4250 – Web Search and Recommender Systems

37

Near-Duplicate Detection

- *Search*:
 - find near-duplicates of a document D
 - $O(N)$ comparisons required
- *Discovery*:
 - find all pairs of near-duplicate documents in the collection
 - $O(N^2)$ comparisons
- IR techniques are effective for search scenario
- For discovery, other techniques used to generate compact representations

CS 4250 – Web Search and Recommender Systems

38

Fingerprints

1. The document is parsed into words. Non-word content, such as punctuation, HTML tags, and additional whitespace, is removed.
2. The words are grouped into contiguous *n-grams* for some *n*. These are usually overlapping sequences of words, although some techniques use non-overlapping sequences.
3. Some of the n-grams are selected to represent the document.
4. The selected n-grams are hashed to improve retrieval efficiency and further reduce the size of the representation.
5. The hash values are stored, typically in an inverted index.
6. Documents are compared using overlap of fingerprints

CS 4250 – Web Search and Recommender Systems

39

Fingerprint Example

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

(a) Original text

tropical fish include, fish include fish, include fish found, fish found in, found in tropical, in tropical environments, tropical environments around, environments around the, around the world, the world including, world including both, including both freshwater, both freshwater and, freshwater and salt, and salt water, salt water species

(b) 3-grams

938 664 463 822 492 798 78 969 143 236 913 908 694 553 870 779

(c) Hash values

664 492 236 908

(d) Selected hash values using $0 \bmod 4$

CS 4250 – Web Search and Recommender Systems

40

Simhash

- Similarity comparisons using word-based representations more effective at finding near-duplicates
 - Problem is efficiency
- Simhash combines the advantages of the word-based similarity measures with the efficiency of fingerprints based on hashing
- Similarity of two pages as measured by the cosine correlation measure is proportional to the number of bits that are the same in the simhash fingerprints

CS 4250 – Web Search and Recommender Systems

41

Simhash

1. Process the document into a set of features with associated weights. We will assume the simple case where the features are words weighted by their frequency.
2. Generate a hash value with b bits (the desired size of the fingerprint) for each word. The hash value should be unique for each word.
3. In b -dimensional vector V , update the components of the vector by adding the weight for a word to every component for which the corresponding bit in the word's hash value is 1, and subtracting the weight if the value is 0.
4. After all words have been processed, generate a b -bit fingerprint by setting the i th bit to 1 if the i th component of V is positive, or 0 otherwise.

CS 4250 – Web Search and Recommender Systems

42

Simhash Example

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

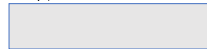
(a) Original text

tropical 2 fish 2 include 1 found 1 environments 1 around 1 world 1
including 1 both 1 freshwater 1 salt 1 water 1 species 1

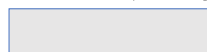
(b) Words with weights

tropical	01100001	fish	10101011	include	11100110
found	00011110	environments	00101101	around	10001011
world	00101010	including	11000000	both	10101110
freshwater	00111111	salt	10110101	water	00100101
species	11101110				

(c) 8 bit hash values



(d) Vector V formed by summing weights



(e) 8-bit fingerprint formed from V

CS 4250 – Web Search and Recommender Systems

43

Simhash Example

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

(a) Original text

tropical 2 fish 2 include 1 found 1 environments 1 around 1 world 1
including 1 both 1 freshwater 1 salt 1 water 1 species 1

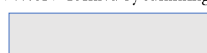
(b) Words with weights

tropical	01100001	fish	10101011	include	11100110
found	00011110	environments	00101101	around	10001011
world	00101010	including	11000000	both	10101110
freshwater	00111111	salt	10110101	water	00100101
species	11101110				

(c) 8 bit hash values

1 -5 9 -9 3 1 3 3

(d) Vector V formed by summing weights



(e) 8-bit fingerprint formed from V

CS 4250 – Web Search and Recommender Systems

44

Simhash Example

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

(a) Original text

tropical 2 fish 2 include 1 found 1 environments 1 around 1 world 1
including 1 both 1 freshwater 1 salt 1 water 1 species 1

(b) Words with weights

tropical	01100001	fish	10101011	include	11100110
found	00011110	environments	00101101	around	10001011
world	00101010	including	11000000	both	10101110
freshwater	00111111	salt	10110101	water	00100101
species	11101110				

(c) 8 bit hash values

1 -5 9 -9 3 1 3 3

(d) Vector V formed by summing weights

1 0 1 0 1 1 1 1

(e) 8-bit fingerprint formed from V

CS 4250 – Web Search and Recommender Systems

45


Removing Noise

- Many web pages contain text, links, and pictures that are not directly related to the main content of the page
- This additional material is mostly *noise* that could negatively affect the ranking of the page
- Techniques have been developed to detect the content blocks in a web page
 - Non-content material is either ignored or reduced in importance in the indexing process

CS 4250 – Web Search and Recommender Systems

46

Noise Example

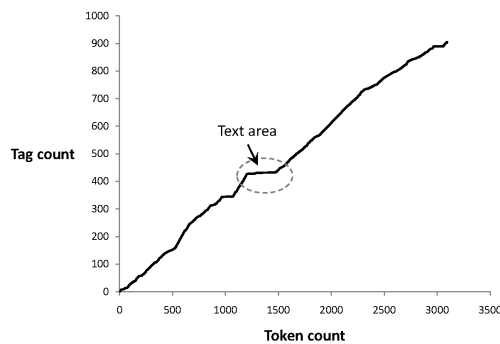


The screenshot shows a CNN.com article from November 2007. The main headline is "Aquarium plays whale shark matchmaker". A red rectangular box highlights a paragraph in the main text area, which is pointed to by an arrow labeled "Content block". The highlighted text discusses the Georgia Aquarium's efforts to find a mate for a whale shark named Alice. The page includes a sidebar with navigation links, a "SCIENCE & SPACE" section header, and various advertisements on the right side.

47

Finding Content Blocks

- Cumulative distribution of tags in the example web page



- Main text content of the page corresponds to the "plateau" in the middle of the distribution

CS 4250 – Web Search and Recommender Systems

48

Finding Content Blocks

- Represent a web page as a sequence of bits, where $b_n = 1$ indicates that the n th token is a tag
- Optimization problem where we find values of i and j to maximize both the number of tags below i and above j and the number of non-tag tokens between i and j
- i.e., maximize

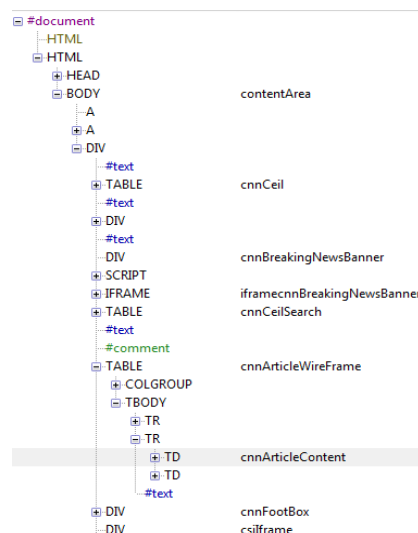
$$\sum_{n=0}^{i-1} b_n + \sum_{n=i}^j (1 - b_n) + \sum_{n=j+1}^{N-1} b_n$$

CS 4250 – Web Search and Recommender Systems

49

Finding Content Blocks

- Other approaches use DOM structure and visual (layout) features

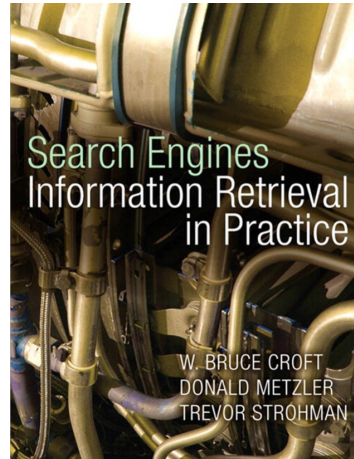


CS 4250 – Web Search and Recommender Systems

50

Reading

- Chapter 3



CS 4250 – Web Search and Recommender Systems