



## Collection Building II

---

CS 510 Spring 2010

1



## Scoped Crawling

---

Limit content being sought

- Topic
- Location
- Format
- Genre
- Language

Weight function  $w(p)$  will reflect degree to which  $p$  is in scope

Evaluate recall

- % of hand-labeled pages in result
- overlap starting from different seeds

CS 510 Spring 2010

2



## Topical Crawlers

---

- 1<sup>st</sup> generation: explore depth  $d$  from each topical page
- 2<sup>nd</sup> generation: differential exploration  
rank potential of link  $p \rightarrow q$

CS 510 Spring 2010

3



## Topical Crawler II

---

- 3<sup>rd</sup> generation: train topic classifier from directory
  - User gives topics
  - Follow links from pages the classifier finds most relevant
  - Also use parent topics
  - Identify authoritative hub pages, favor their links

CS 510 Spring 2010

4



## Incremental Crawling Weighted Freshness (WF)

---

$$WF(t) = \sum_{p \in C(t)} w(p) * f(p, t)$$

Where

- $C(t)$  = pages crawled up to time  $t$
- $w(p)$  = weight of page  $p$
- $f(p, t)$  = freshness of local copy of  $p$  at time  $t$

Want to maximize steady-state average of  
 $WF(t)$

CS 510 Spring 2010

5



## Choice

---

- Download new page: improve coverage
  - Revisit old page: improve currency
- Both can improve link info

CS 510 Spring 2010

6



## Maximizing Freshness

---

Optimal strategy depends on freshness function

General process

- Estimate change rate of each page  $p$
- Assign portion of crawl rate per page:  $r(p)$
- Scheduling: Find a crawl order that matches the page rates

CS 510 Spring 2010

7



## Binary Freshness

---

$f(p, t) = 1$  if stored copy = live, else 0

- Revisiting proportional to change rates is sub-optimal
- Studied by Cho and García-Molina: Uniform does better

CS 510 Spring 2010

8

## Binary Freshness Example

### Two pages

- e1 changes 9 times/day
  - e2 changes 1 time/day
- uniformly random in interval  
fetch at mid-interval

Look at effect on freshness of different #  
of fetches per day

(duration of freshness)\*  
(prob. changes before mid-interval)

CS 510 Spring 2010

9

## Refresh Rates

row	$f_1 + f_2$	$f_1$	$f_2$	benefit	best
(a)	1	1	0	$\frac{1}{2} \times \frac{1}{18} = \frac{1}{36}$	0 1
(b)		0	1	$\frac{1}{2} \times \frac{1}{2} = \frac{9}{36}$	
(c)	2	2	0	$\frac{1}{2} \times \frac{1}{18} + \frac{1}{2} \times \frac{1}{18} = \frac{2}{36}$	0 2
(d)		1	1	$\frac{1}{2} \times \frac{1}{18} + \frac{1}{2} \times \frac{1}{2} = \frac{10}{36}$	
(e)		0	2	$\frac{1}{3} \times \frac{2}{3} + \frac{1}{3} \times \frac{1}{3} = \frac{12}{36}$	
(f)	5	3	2	$\frac{3}{36} + \frac{12}{36} = \frac{30}{72}$	2 3
(g)		2	3	$\frac{2}{36} + \frac{6}{36} = \frac{31}{72}$	
(h)	10	9	1	$\frac{9}{36} + \frac{1}{4} = \frac{36}{72}$	7 3
(i)		7	3	$\frac{7}{36} + \frac{6}{36} = \frac{41}{72}$	
(j)		5	5	$\frac{5}{36} + \frac{15}{36} = \frac{40}{72}$	

CS 510 Winter 2007

10

## Observations

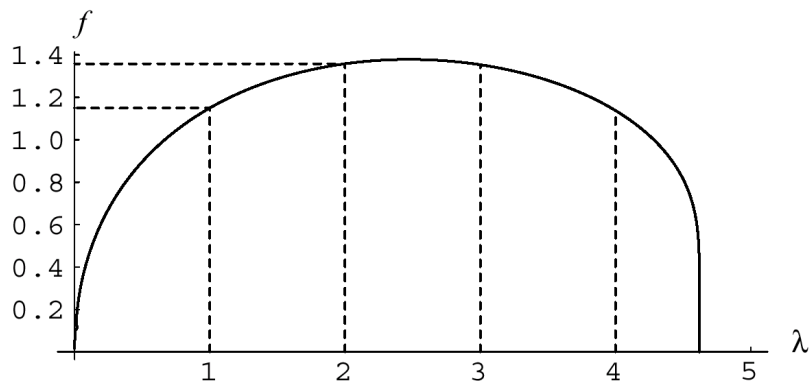
- Uniform method is ALWAYS better than proportional method for freshness
- Optimal refresh rate graph ALWAYS has the same shape

CS 510 Winter 2007

11

## Freshness: Optimal Rate

5 pages that change 1, 2, 3, 4, and 5 times a day



CS 510 Winter 2007

12

## Temporal Freshness

$\text{age}(p, t)$  = length of time cached  $p$  has been different from on-line  $p$  at time  $t$

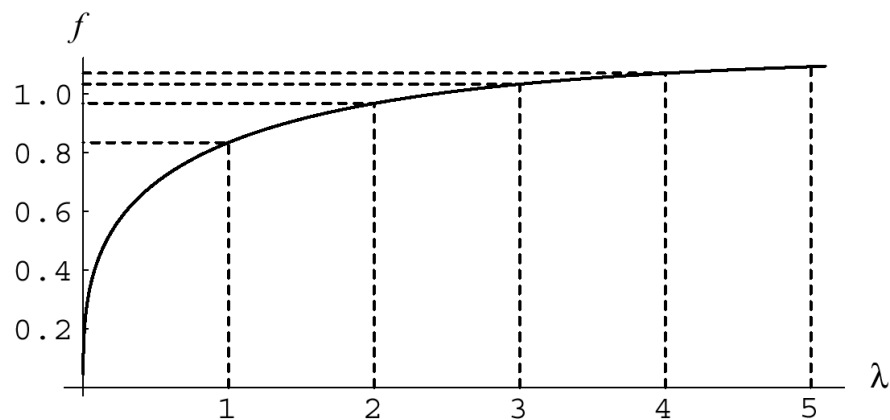
Rate freshness inversely to age

Optimal strategy crawls faster-changing pages more often

CS 510 Spring 2010

13

## Age: Optimal Rate



CS 510 Winter 2007

14



## Other Metrics

---

- Percent of content in common with current page
  - Need to model longevity of new content
- Embarrassment level: chance that the user will click on a stale page

CS 510 Spring 2010

15



## Undesirable Content: Duplicates

---

- Reduce caching requirements, or prune search results
- Real issue is near-duplicates
  - Date stamp, formatting, home link

CS 510 Spring 2010

16





## Shingling

---


Use overlapping “shingles”, hashing and permutations to produce a *sketch*, from which similarity can be derived

w-shingle: w consecutive terms from a document

Overlaps with next shingle by w-1 terms

CS 510 Spring 2010

17



## Example

---

D1: A new ESP device was announced by Apple Computer at their annual meeting on Tuesday

D2: A new ESP device was announced by Apple Corporation at their annual meeting on Tuesday

### 3-shingles of D1:

A new ESP, new ESP device, ESP device was, device was announced, was announced by, ...

$S(D)$  = set of (unique) shingles of document D

Document similarity = Jaccard Coefficient

$$\text{SIM}(D1, D2) = \frac{|S(D1) \cap S(D2)|}{|S(D1) \cup S(D2)|}$$

CS 510 Spring 2010

18

## Fast Estimate of SIM(D1,D2)

First, apply hash function  $h$  to all shingles

$$H(D) = \{h(s) | s \in S(D)\}$$

Assume  $h \rightarrow 1..1000$  [ $2^{64}$  more usual]

$$H(D1) =$$

{150, 460, 26, 821, 337, 698, 14, 875, 512, 504, 219, 417, 733}

$$H(D2) =$$

{150, 460, 26, 821, 337, 698, 711, 43, 208, 504, 219, 417, 733}

## Fast Estimate, cont.

$$\text{sorted } H(D1) =$$

{14, 26, 150, 219, 337, 417, 460, 504, 512, 698, 733, 821, 875}

$$\text{sorted } H(D2) =$$

{26, 43, 150, 208, 219, 337, 417, 460, 504, 698, 711, 733, 821}

- Consider the smallest in each

$$\min(H(D1)) = 14$$

$$\min(H(D2)) = 26$$

- Take  $(n + 200) \bmod 1000$ , then take smallest

$$+200 \quad +400 \quad +600 \quad +800$$

$$H(D1)$$

$$H(D2)$$



## Sketch of Document

---

Collect the list of mins

For  $H(D1)$  [14, 21, 98, 17, 19]

For  $H(D2)$  [26, 21, 98, 17, 8]

% similarity of sketch estimates  
 $SIM(D1, D2)$



## How Contrived Was That?

---

	<u>H1</u>	<u>H2</u>
Was trying to sample		
the hash sets of both		
	14	0
	26	1
	43	1
	150	1
	208	1
(But I didn't want to	219	1
sample over 1..1000)	337	1
	417	1
	460	1
	512	0
	698	1
	711	1
	733	1
	821	1
	875	0



## Getting a Random Sample

My sampling wasn't random

Take a bunch of random permutations

$\Pi_1, \Pi_2, \dots, \Pi_{200}$

$m_i(D) = \min(\{\Pi_i(v) \mid v \in H(D)\})$

Sketch of  $D =$

$[m_1(D), m_2(D), \dots, m_{200}(D)]$

CS 510 Spring 2010

23



## Estimate Pairwise Similarities

If  $m_j(D1) = m_j(D2)$ , generate  $(D1, D2, 1)$

Notes:

- Random permutation: only needs to be "min-wise" independent – any element has equal chance of being mapped to smallest
- Pre-filter: Sketch of sketches

CS 510 Spring 2010

24



## Crawler Traps

---

Populate a large URL space with generated content

- calendar pages
- malicious

BEAST: Budget Enforcement with Anti-Spam Techniques

Give each site a “crawl budget” & prioritize based on remaining budget

CS 510 Spring 2010

25



## Web Spam

---

1. Keyword stuffing
2. Link spam: trying to influence rank
3. Cloaking: different content to crawlers vs. people

CS 510 Spring 2010

26



## Handling 1 & 2

---

- Classification
- Which features?
  - Hyperlink features
  - Term & phrase frequency
  - DNS lookup stats
  - HTML markup structure

CS 510 Spring 2010

27



## Handling Cloaking

---

- Giving different content based on `user-agent`
    - Need to come in with different user agent or from different IP and compare results
  - Redirection Spam
    - Script will forward user to another URL
- Use some kind of time-bound evaluation of scripts?

CS 510 Spring 2010

28