

Document Ranking of Web Document Retrieval Systems

* , **

Dong-Un An · In-Ho Kang

1.	4. 가
2.	5.
3.	•

가 .
, 가
, 가
, 가
가
가
URL

*
(Associate Professor, School of Information and Electronics Engineering, Chonbuk National University, duan@chonbuk.ac.kr)

**
(Department of Computer Science, KAIST, ihkang@csone.kaist.ac.kr)
· : 2003 6 18
· : 2003 6 23

가 URL
URL
, URL

ABSTRACT

The Web is rich with various sources of information. It contains the contents of documents, multimedia data, shopping materials and so on. Due to the massive and heterogeneous web document collections, users want to find various types of target pages. We can classify user queries as three categories according to users' intent, content search, the site search, and the service search. In this paper, we present that different strategies are needed to meet the need of a user. Also we show the properties of content information, link information and URL information according to the class of a user query. In the content search, content information showed the good result. However, we lost the performance by combining link information and URL information. In the site search, we could increase the performance by combining link information and URL information.

KEYWORDS

Content Search, Site Search, Service Search, Content Information, Link Information, URL Information

1.

(Salton et al. 1998).

가 가
가 . 가

가

(Page et al. 1998).

(pre-

가

cision)

, (recall)

(PageRank)
(Brin et al. 1998). (www.google.com)

- (informational)
 - (navigational)
 - (transactional)
- 가

“ What is a prime factor? ”
 “ prime factor ” 가
 ‘ prime factor ’ 가

가 (Croft 2000).

“ Where is the site of John Hopkins Medical Institutions? ” “ John Hopkins Medical Institutions ” 가 John Hopkins Medical Institution

가

가 (Westerveld et al. 2001; Yang 2001).

가

“ Where can I buy concert tickets? ” “ buy concert tickets ”
가

‘ Mutual Information ’
 , ‘ mutual funds ’ ‘ information ’
 가 가

가

가

가

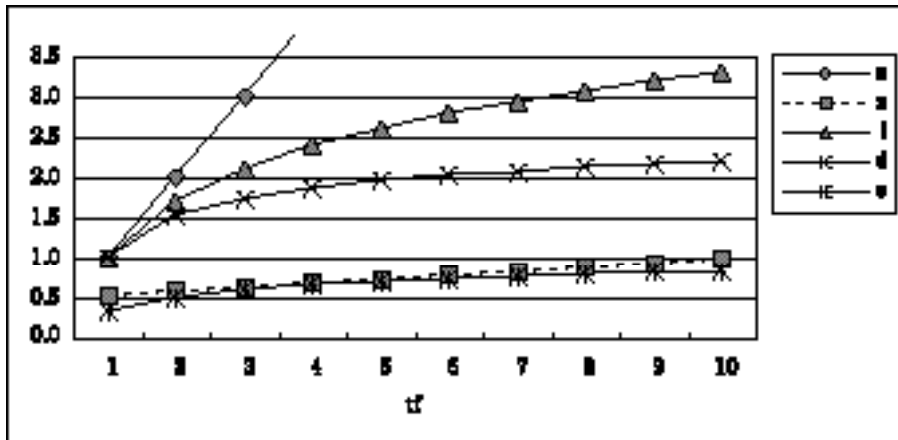
‘ Britney’s Fan Club ’
 (www.yahoo.com) (www.lycos.com)
 가

2.

가

(Broder 2002).

URL	가	SMART	tf
가	URL	가	(Salton 1989).
2.1		$b: 0.01$	
		$n: tf$	
		$a = 0.5 + \frac{tf}{\max tf}$	
		$l: 1.0 + \log tf$	
text	anchor	max tf	
text	(Croft 2000). Anchor	tf	
	가	OKAPI	가
		(Robertson et al. 1994).	
가	Anchor text	$d: 1 + \log(1 + \log(tf))$	
	가	$o: \frac{tf}{2 + tf}$	
	가		
	tf(term frequency)	df	가
df(document frequency)	tf		
	가	$n: 1.0$	
	df	$t: \log \frac{N}{df}$	
가	가		
	df가	N	
	df가	tf	가
		<	1>
	(Salton et al. 1983). tf		
df		2.2	
	가	TFIDF	
	, tf	df	



< 1 > 가 tf

가 .

2.2.1

2.2.2

가 (authoritative)

(Page-

Rank)

(hub)가 가

(authority)

가 .

(hub)

p

(Kleinberg 1999).

$$authority(p) = \sum_{q, q \text{ Points } p} hub(q)$$

$$hub(p) = \sum_{q, p \text{ Points } q} authority(q)$$

가

가

(Brin et al. 1998).

$$PR(A) = (1 - d) + d \times (PR(T_1) / C(T_1) + K + PR(T_n) / C(T_n))$$

100-200

PR(A) A , PR
 (Ti) A 가 Ti
 , C(Ti) Ti가 가
 , d 가

papers)
 · file: index.html
 (. http://trec.nist.gov/pubs/trec9/
 t9proceedings.html)

Westerveld(2001) URL 가
 (URLprior)
 . Westerveld(2001)
 , P(entrypage|URL type = t)

(Page et al. 1998).

WT10g <
 1> .

2.3 URL

< 1> WT10g

가 .
 가 가

URL 가
 URL 가
 URL
 URL

URL type	# entry pages	# WT10g
root	38 (71.7%)	11,680 (0.6%)
subroot	7 (13.2%)	37,959 (2.2%)
path	3 (5.7%)	83,734 (4.9%)
file	3 (5.7%)	1,557,719 (92.1%)

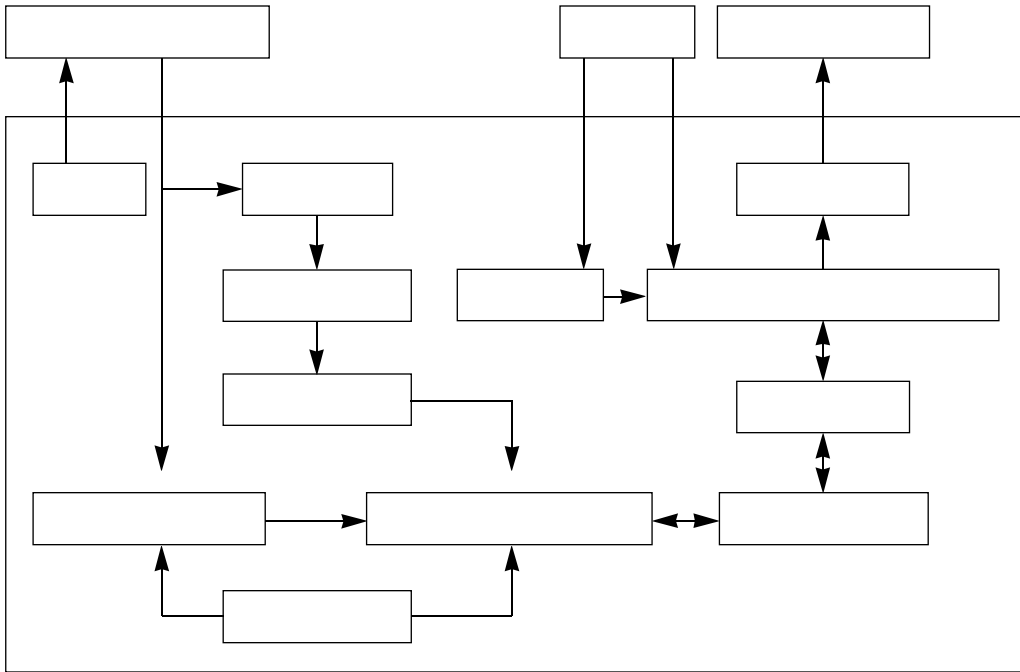
URL URL
 ‘ / ’ . Westerveld
 (2001) URL 4가

URL

- root:
 (. http://trec.nist.gov)
- subroot:
 가
 (. http://trec.nist.gov/pubs)
- path:
 가
 (. http://trec.nist.gov/pubs/trec9/

3.

가
 (Condor)



< 2 >

가

Inference Net-

가

가

work

가

가

< 2 >

$$w_t = TF \times IDF \times P_t$$

(Baeza-Yates et al. 1999;

Salton et al. 1983).

where

$$TF = \frac{tf}{2 \times (0.25 + 0.75 \times \frac{doc\ len}{avg\ doc\ len}) + tf}$$

$$IDF = \log\left(\frac{N - df + 0.5}{df + 0.5}\right)$$

2.0: high

$p_t = 1.5$: important

1: others

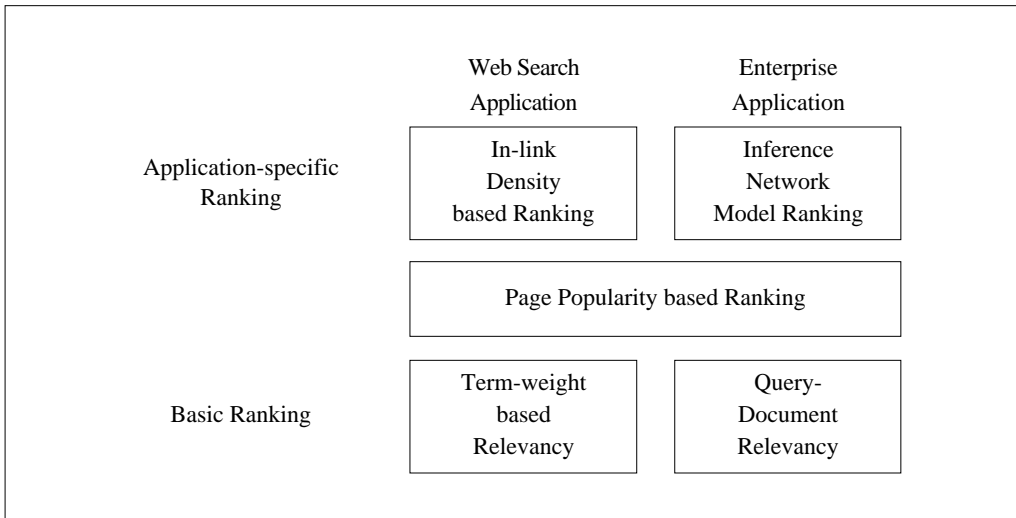
(, , ,)

가 가 . 가

, , , tf df

OKAPI

(P_t)



< 3>

4.1 가

(< 3>).

WT10g

(Bailey et al. to

가 appear). WT10g 10G

가 가 CSIRO (CSIRO 2001).

TREC-2001 Web Track ad-hoc task (501-550) home-

page finding task (1-145)

(Hawking et al. 2001). Ad-hoc

task

, homepage finding task

4. 가

TOPIC HOME

(average precision)

(Hawking et al. 2001).

가

R d가 , df , tf
 r(d) , URLprior
 (CMB) 가 .
 URLprior

$$P_{avg} = \frac{1}{|R|_d} \frac{|R_{r(d)}|}{r(d)}$$

$$rel(d) = 0.65 \times \text{Content Information} + 0.25 \times \text{URL Information} + 0.1 \times \text{Link Information}$$

가 MRR (Mean Reciprocal Rank) (Hawking et al. 2001). MRR

‘ and ’

가 ‘ sum ’ . ‘ and ’

. n MRR 가

r ‘ sum ’

가 가

$$MRR = \frac{1}{n} \sum_i \frac{1}{r_i}$$

< 2 >

Anchor and CMB

anchor text

, ‘ and ’

4.2

< 2 >

text

가 anchor
 Anchor
 anchor text 가
 Common
 가

	TOPIC	HOME
model	P _{avg}	MRR
Anchor and	0.031	0.297
Anchor and CMB	0.031	0.431
Anchor sum	0.034	0.351
Anchor sum CMB	0.034	0.583
Common and	0.131	0.294
Common and CMB	0.122	0.580
Common sum	0.182	0.355
Common sum CMB	0.169	0.673
MAX	0.226	0.774
AVG	0.145	0.432

Common sum CMB
 가
 Common sum
 0.173
 URLprior
 . < 2>
 MAX AVG TREC-2001 Anchor and
 'and' 가
 < 2> Com-
 mon Anchor
 Common Anchor
 'sum'
 anchor text
 anchor text 'and'
 URLprior
 URLprior
 URL 가
 tf df 가
 가
 a df
 가
 TFIDF
 OKAPI <
 3> . < 3> Content Info.
 'and'
 'and' 가 'sum'
 'and'
 가
 가
 OKAPI TFIDF
 URLprior URL
 TFIDF
 OKAPI 가
 Common sum CMB Anchor and CMB
 0.730 TFIDF

< 3> 가

model	TFIDF		OKAPI	
	TOPIC	HOME	TOPIC	HOME
Content Info.	0.170	0.340	0.182	0.355
CMB	0.159	0.640	0.169	0.673

‘ sum ’
 URL
 anchor text
 가 ‘ and ’
 가 가
 URL
 가
 TFIDF OKAPI

5.

가

가

가

가

tf df

URL

URL

Baeza-Yates, R., & Ribeiro-Neto, B. 1999. Modern Information Retrieval. ACM Press.

Bailey, P., Craswell, N., & Hawking, D. (to appear). “ Engineering a multipurpose test collection for web retrieval experiments”. Information Processing and Management

Brin, S., & Page, L. 1998. “ The Anatomy of a Large-Scale Hypertextual Web Search Engine ”. Computer Networks and ISDN Systems, 30(1-7): 107-117.

-
- Broder, A. 2002. "A Taxonomy of Web Search". SIGIR Forum, 36(2).
- Croft, W. B. 2000. Combining Approaches to Information Retrieval. In W. B. Croft (Ed.), *Advances in Information Retrieval: Recent Research from the Center for Intelligent Information Retrieval*, Kluwer Academic Publishers: 1-36.
- CSIRO. (2001). "Web Research Collections - Trec Web Track." <<http://www.ted.cmis.csiro.au/TRECWeb/>>.
- Hawking, D., & Craswell, N. 2001. "Overview of the Trec-2001 Web Track". In *Text REtrieval Conference (trec-10)*: 61-67.
- Kleinberg, J. M. 1999. "Authoritative Sources in a Hyperlinked Environment". *Journal of the ACM*, 46(5): 604-632.
- Page, L., Brin, S., Motwani, R. & Winograd, T. 1998. *The Pagerank citation ranking: Brining Order to the Web* (Tech. Rep.). Stanford Digital Library Technologies Project.
- Robertson, S. E., Walker, S., Jones, S., Hancock-Beaulieu, M., & Gatford, M. 1994. "Okapi at trec-3". In *Text REtrieval Conference (trec-2)*: 109-126.
- Salton, G. 1989. *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer Reading*, MA: Addison-Wesley.
- Salton, G. & McGill, M. 1983. *Introduction to Modern Information Retrieval*, New York: McGraw-Hill Press.
- Westerveld, T., Kraaij, W., & Hiemstra, D. 2001. "Retrieving Web Pages using Content, Links, Urls and Anchors". In *Text REtrieval Conference (trec-10)*: 663-672.
- Yang, K. 2001. "Combining Text and Link-Based Retrieval Methods for Web IR". In *Text REtrieval Conference (trec-10)*: 609-618.