

Web-Document Retrieval Using Genetic Algorithms

2001 2

Web-Document Retrieval Using Genetic Algorithms

2000 10

2000 12

印

印

印

HTML

가

가

. HTML

HTML

가

가

가 TREC

가

가

1.	1
1.1	1
1.2	3
2.	4
2.1	4
2.2	7
3.	8
3.1	8
3.1.1	8
3.1.2	10
3.1.3	11
3.2	(SCAIR)	14
4.	17
4.1	17
4.2	가	20
5.	25
5.1	I :	26
5.1.1	26
5.1.2	27
5.2	II : TREC	30
5.2.1	TREC	30
5.2.2	33
6.	39
	40

1.

1.1

1969 .
HTML(HyperText Markup Language) WWW(World
Wide Web) 가
3 가 [Lawrence and Giles, 1998].
가
가
가
가
가
가
HTML HTML
HTML 가
가 [Boyan et al., 1996].
HTML
HTML

(anchor),

. HTML

HTML

가

가 TREC(Text

REtrieval Conference)

가

1.2

. 2

. 3

4

. 5

가

TREC

6

.

2.

1

, 가

가 가

,

2.1

가

,
1998][Weiss et al., 1996].

[Bharat and Henzinger, 1998][Picard,

가

.

. a b 가 , a

b

.

Spertus

가

[Spertus, 1997].

a 가

b

c 가

b c

가 .
[Chakrabarti et al., 1997].
가 , 가
가 가 .
(Google) [Brin and Page, 1998].
(PageRank)
, 가
가
가 , 가
가 가 .
, .
LASER
[Boyan et al., 1996].
HTML
,
- .
가
(simulated annealing) .
, (AltaVista) (Yahoo)

, (Lycos)[Mauldin, 1997],
가
. HTML 가
가
.
Cutler [Cutler et al., 1999].

2.2

(description)

[Gordon, 1988].

(population)

가

가

[Yang et al., 1993][Yang and Honavar, 1998]. Yang

(feature

selection)

(classifier)

3.

(vector space) , (probabilistic) (boolean) ,
AND, OR ,
SCAIR .

3.1

3.1.1

AND, OR, NOT

(similarity)

MIN MAX 가
 MIN MAX 가
 [Salton et al., 1983][Turtle and Croft, 1991][Callan et al., 1992]. AND
 OR

P-norm 가
 P-norm A_1, A_2, \dots, A_n 가
 $d_{A_1}, d_{A_2}, \dots, d_{A_n}$ D_n $(d_{A_1}, d_{A_2}, \dots, d_{A_n})$
 $d_{A_1} \text{ OR } d_{A_2} \text{ OR } \dots \text{ OR } d_{A_n}$ OR
 가 가 n 가 0
 $d_{A_1} \text{ AND } d_{A_2} \text{ AND } \dots \text{ AND } d_{A_n}$ AND
 n 가 1 , 가 가 1 가
 가 OR $(0, 0, \dots, 0)$
 , AND $(1, 1, \dots, 1)$
 가 .

P-norm (1) .
 p .

$$Q_{OR_p} = (A_1, a_1) OR_p (A_2, a_2) OR_p \dots OR_p (A_n, a_n) \quad (1)$$

$$Q_{AND_p} = (A_1, a_1) \text{ AND}_p (A_2, a_2) \text{ AND}_p \dots \text{ AND}_p (A_n, a_n)$$

가 a_i . P-norm
(2)

$$\text{SIM}(Q_{or_p}, D) = \sqrt[p]{\frac{a_1^p d_{A_1}^p + a_2^p d_{A_2}^p + \dots + a_n^p d_{A_n}^p}{a_1^p + a_2^p + \dots + a_n^p}}$$

$$\text{SIM}(Q_{and_p}, D) = 1 - \sqrt[p]{\frac{a_1^p (1 - d_{A_1})^p + a_2^p (1 - d_{A_2})^p + \dots + a_n^p (1 - d_{A_n})^p}{a_1^p + a_2^p + \dots + a_n^p}} \quad (2)$$

$$\text{SIM}(Q_{not}, D) = 1 - \text{SIM}(Q, D)$$

3.1.2

(binary) 가

가 . 가
, .
가 (3)

$$d_i = (w_{i1}, w_{i2}, \dots, w_{in}) \quad (3)$$

d_i , w_{ik} d_i t_k
가 . i 가 0 .

d 가 $(w_{i1}, w_{i2}, \dots, w_{in})$ q 가 $(w_{i1}, w_{i2}, \dots, w_{in})$
 d q (4)
 (cosine coefficient similarity)

$$\text{sim}(d_i, q) = \frac{\vec{d}_i \cdot \vec{q}}{|\vec{d}_i| |\vec{q}|} \quad (4)$$

가
 가
 (term frequency), (document frequency), (document
 length normalization) 가 .
 가 .
 가 .
 가 .

3.1.3

1960 [Maron and Kuhns,
 1960], Robertson
 [Robertson and Sparck, 1976].

가

가

가

X 가 (5)

, x_i

i

$x_i = 0$

1

$$X = (x_1, x_2, \dots, x_n)$$

(5)

가

$g(X)$

(6)

$$g(X) = \log \frac{\Pr(X | rel)}{\Pr(X | nonrel)}$$

(6)

$\Pr(X | rel)$

가

X

$\Pr(X | nonrel)$

가

X

X 가

가

$g(X)$

가

$g(X)$

$\Pr(X | rel)$

$\Pr(X | nonrel)$

(7) $X = \sum_{i=1}^n x_i$ 가 , $g(X)$

· X
· 가 0 1 가 .

$$g'(X) = \sum_{i=1}^n x_i \log \frac{p_i(1 - q_i)}{q_i(1 - p_i)} + C$$

$$p_i = \Pr(x_i = 1 | rel) \tag{7}$$

$$q_i = \Pr(x_i = 1 | nonrel)$$

C
· $g'(X)$ x_i p_i q_i
· $g'(X)$
·
· 가

· p_i q_i
Croft $p_i = 0.5$ $q_i = n_i/N$ [Croft and
Harper, 1979].

3.2 (SCAIR)

TREC

SCAIR(SCAI Information Retrieval Engine)

[Shin and Zhang, 1998]. SCAIR

[Salton et al., 1975].

(term)

가

SCAIR

(8)

$tf \cdot idf$

frequency)

[Salton, 1989].

$tf \cdot idf$

tf

(term

idf

(inverse document frequency)

가

$f \cdot idf$

$$w_{di} = tf_{di} \cdot \log\left(\frac{N}{df_i}\right)$$

(8)

w_{di} : d i 가

tf_{di} : d i

N :

df_i : i

SCAIR

가

$tf \cdot idf$

[Broglia et al., 1995] 2-

[Robertson et al., 1995]

2-

가 (9) (10) ,

가 f .

$$w_{di} = \left[0.4 \times H + 0.6 \times \frac{\log(f_{di} + 0.5)}{\log(\max f_{di} + 1.0)} \right] \times \frac{\log \frac{N}{n}}{\log N} \quad (9)$$

w_{di} : d i 가

f_{di} : d i

N :

n : i 가

H : 1

$$w_{di} = \frac{f_{di}}{k_1 \left((1-b) + b \frac{\text{document length}}{\text{average document length}} \right) + f_{di}} \times \log \frac{N - n + 0.5}{n + 0.5} \quad (10)$$

w_{di} : d i 가

f_{di} : d i

N :

n : i 가

k_1 : 2.0

b : 0.75

(11) .

$$sim(d, q) = \sum_{k=1}^n (\alpha_{dk} \times w_{dk}) \times w_{qk} \quad (11)$$

w_{dk} : d k 가

w_{qk} : q k 가

α_{dk} : d k

가 a_{dk} k 가 HTML
 . , k 가 a_{dk}
 1.0 . 가

HTML HTML
 HTML . 가
 가 . 가
 . 가
 (11) .

4.

가
가 .
NP
(evolutionary computation)
[Zhang, 1995].
, 가 .

4.1

(genetic algorithm),
(genetic programming), (evolutionary programming),
(evolutionary strategy) , , ,

DNA,
(chromosome) . (individual) (crossover)
(mutation)
가 ,

가

,
 (population)
 .
 (gene) 가 가 .
 (genotype) , (phenotype)
 .
 , .
 1 [Ballard, 1997].
 .
 n , n . n
 k , , ,
 . k k
 . ,
 가 가 . k
 k/n (generation gap)
 , 가 1 가 , 가
 (generational GA) . , 가 $1/n$
 가 , , 가
 (steady-state GA) .
 가

Choose a population size.
 Choose the number of generations N_g .
 Initialize the population.
 Repeat the following for N_g generations:

1. Select a given number of pairs of individuals from the population probabilistically after assigning each structure a probability proportional to observed performance.
2. Copy the selected individuals(s), then apply *operators* to them to produce new individual(s).
3. Select other individuals at random and replace them with the new individuals.
4. Observe and record the fitness of the new individuals.

Output the fittest individual as the answer.

1:

(parent)가 .
 (offspring)

가

가 가

(loop)

가

4.2 가

(hyperlink)

가

가

가

(deterministic)

가

가

HTML

, HTML

가

가

. 가

2

3


```
for g = 1 to gmax
```

가

```
for i = 1 to M
```

```
    p1, p2
```

```
    offspringi = crossover(p1, p2)
```

```
    offspringi = mutation(offspringi)
```

```
end for
```

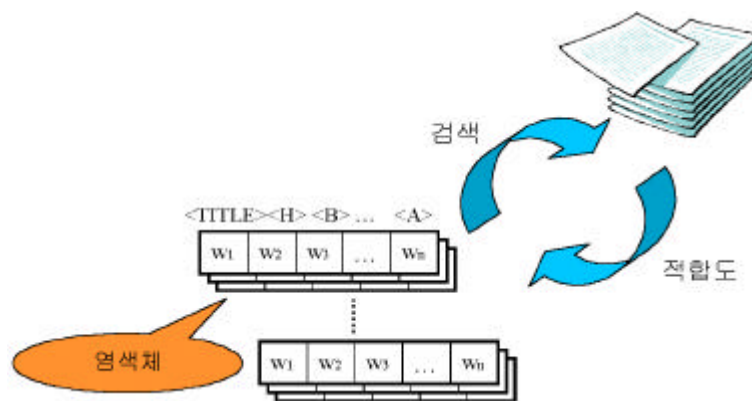
```
M
```

```
end for
```

```
N
```

```
return
```

2: 가



3:

(fitness function) . 가

TREC

11 (11-point average precision)

[Voorhees and Harman, 1999]. 11 (recall) 0.0

0.1 11 (precision)

P (relevant

document) , R

(12),

(13) , (14) .

11 가

$$P = \frac{\text{Number of retrieved relevant documents}}{\text{Total number of retrieved documents}} \quad (12)$$

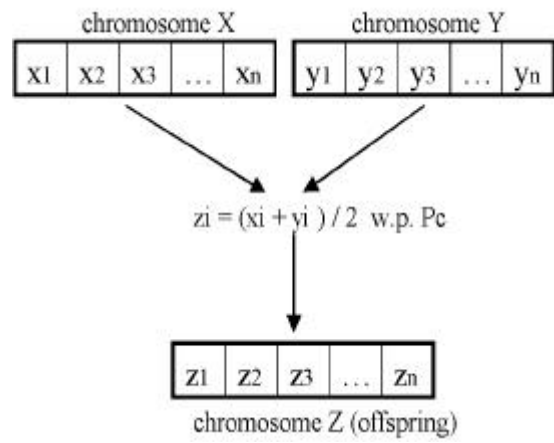
$$R = \frac{\text{Number of retrieved relevant documents}}{\text{Total number of relevant documents in collection}} \quad (13)$$

$$fitness = \frac{1}{\sum_{i=1}^N r(d_i)} \sum_{i=1}^N \frac{1}{i} \sum_{j=1}^i r(d_j) \quad (14)$$

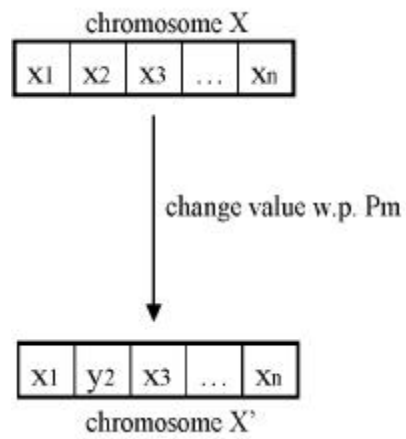
$r(d)$: d
 (d가 1 , 0)

N :

1992]. (arithmetical crossover) [Michalewicz,
가
(mutation) . 4
5 가



4:



5:

5.

가 가

가 가

HTML

가 (title), (header), (bold),
(italic), (anchor) . 가
<TITLE>, <Hx>, , <I>, <A> .

가
2

가

가

가

가

TREC

SCAIR

가

f · idf

5.1 I :

5.1.1

I

가

(call-for-papers homepage)

(conference homepage)

100

가 가

A

B

가 , 가

6

(title)

(description)

II

TREC

가

가

6

‘genetic’, ‘algorithms’, ‘conference’, ‘especially’, ‘information’, ‘retrieval’

가 . <title>, <desc>, <narr>
 가 가 , 100

10 가 .

<title> genetic algorithms
 <desc> Description:
 Is there a conference on genetic algorithms,
 especially containing call for papers on information retrieval?
 <narr> Narrative:
 none.

6: I

5.1.2

I ,
 A B .

- (population size) : 100
- (number of generation) : 30
- (probability of mutation) : 0.04
- 가 : 0.0 4.0

A, B 10
 가 가 ,
 가 가 가

7 A B 가

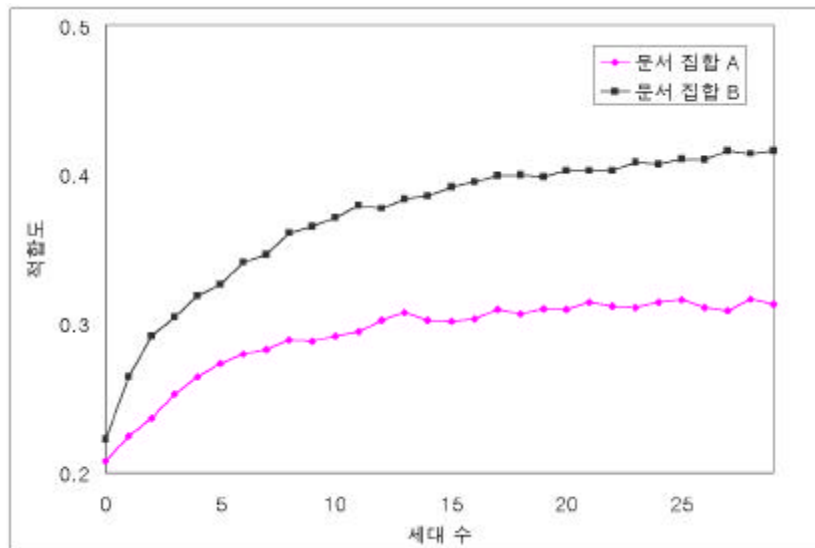
가

A B 가

A B

B , 가 A

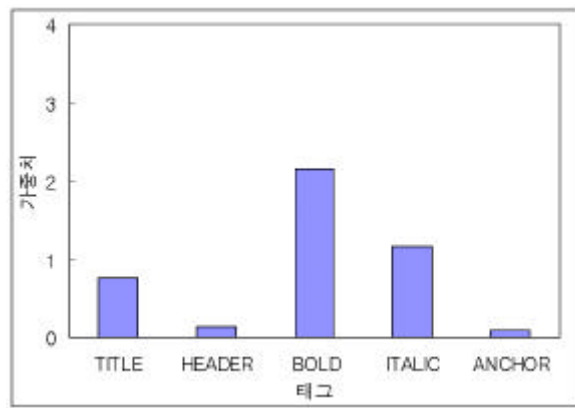
가 가



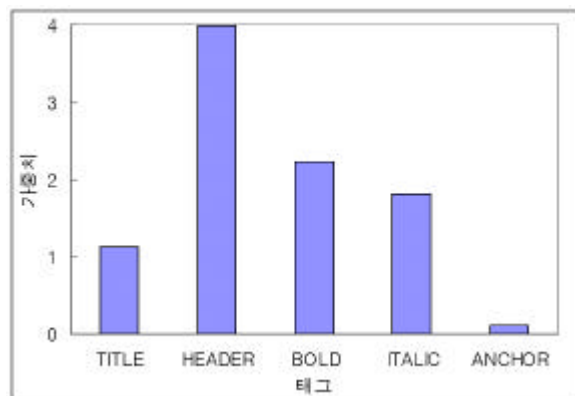
7:

30 가

, , , , A
 0.77, 0.13, 2.15, 1.16, 0.08 , B 1.12, 3.98, 2.23,
 1.80, 0.10 . 8, 9
 . A 가 가
 , B 가 가
 . B



8: A



9: B

5.2 II : TREC

5.2.1 TREC

II ,
가

TREC

TREC(Text REtrieval Conference)

NIST

가

(competition)

[NIST]. TREC

(web track)

(filtering track),

(question

answering track),

(cross-language track)

II

TREC

TREC

1999

TREC-8

(WT2g)

Internet Archive [Archive]

2 가

, 247,491

<title> foreign minorities, Germany
 <desc> Description:
 What language and cultural differences impede the integration
 of foreign minorities in Germany?
 <narr> Narrative:
 A relevant document will focus on the causes of
 the lack of integration in a significant way;
 that is, the mere mention of immigration difficulties is not
 relevant. Documents that discuss immigration problems
 unrelated to Germany are also not relevant.

10: TREC

(description), (topic) , (title),
 (narrative) . 10
 가 가
 가 가
 가

(pooling method)

[Voorhees and Harman, 1999][Zobel, 1998].

가

(pool)

,

가

, 가

100

가

가

가

가

5.2.2

WT2g , TREC-8 401 420
.
20
401 410
가 411 420 .
가 가
200
11 .
200
II
· (population size) : 100
· (number of generation) : 25
· (probability of mutation) : 0.04
· 가 : 0.0 4.0

20

가 I 20 가
가 가 . 가

1.
가 .

2. 1 가 가
가 , 가 가

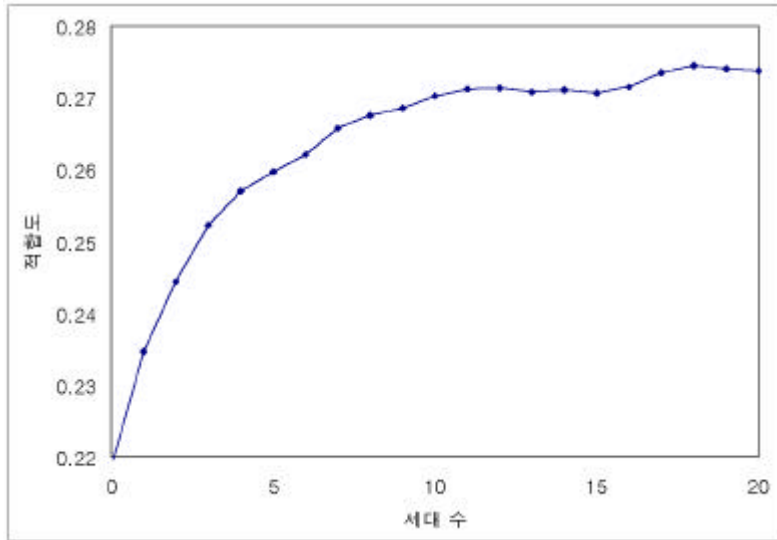
3. 가 1.0 가
가 .

4. 가 , 200

5. 가 가

6. .

가 (1, 2)
, (feature selection) (4, 5, 6)



11:

11

가

가 가

8

가

가 가

가

가

가

가

가

가

가가

HTML	가
<TITLE>	0.6
<Hx>	1.6
	0.7
<I>	0.6
<A>	1.6

1: 가

가 1 . , , , ,
0.6, 1.6, 0.7, 0.6, 1.6 가 .
, > > , 가 .
,
가 ,
가 .
가 , 가
가 가 가
가 가

	11
	0.2383
가	0.2503

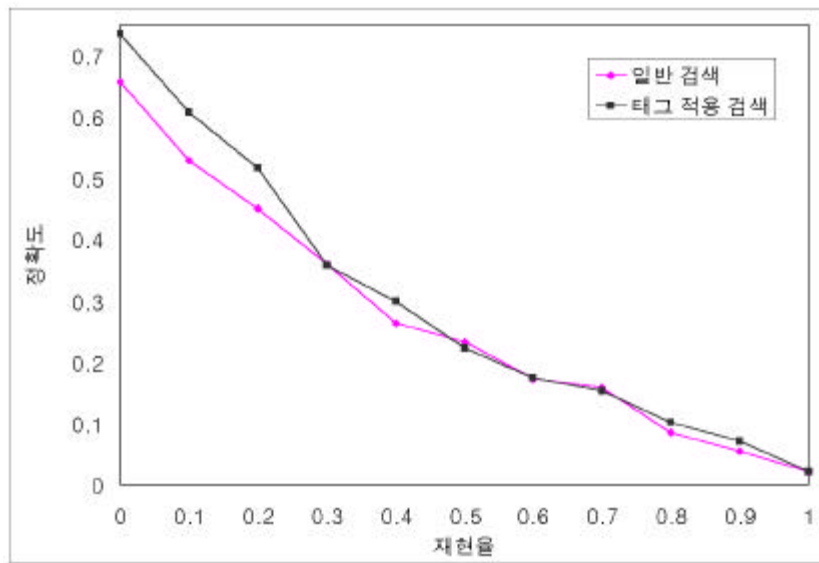
2:

2 가

11

0.2383, 가

0.2503



12:

12 가

- (Average Precision-Recall)

0.3

가

가

가

3

가

		가
0.0	0.6574	0.7350
0.1	0.5284	0.6081
0.2	0.4514	0.5158
0.3	0.3601	0.3580
0.4	0.2645	0.2990
0.5	0.2339	0.2227
0.6	0.1735	0.1745
0.7	0.1583	0.1525
0.8	0.0862	0.1007
0.9	0.0561	0.0705
1.0	0.0219	0.0219

3:

6.

가 . 가
가 가
가 . 가
가 , 가 . 가
가 . 가
가 . 가
가 . 가
가 . HTML
가 (semi-structured) . 가
가 XML
가 가 가 .
(personalized) 가
가

- [Archive] Internet Archive, *Building an Internet Library*,
<http://www.archive.org>.
- [Ballard, 1997] Ballard, D. H., *An Introduction to Natural Computation*,
MIT Press, pp. 263-275, 1997.
- [Bharat and Henzinger, 1998] Bharat, K. and Henzinger, M. R., Improved
Algorithms for Topic Distillation in a Hyperlinked Environment,
Proceedings of the ACM SIGIR '98 Conference, pp. 104-111, 1998.
- [Boyan et al., 1996] Boyan, J., Freitag, D., and Joachims, T., A Machine
Learning Architecture for Optimizing Web Search Engines,
*Proceedings of the AAAI Workshop on Internet-Based Information
Systems*, pp. 1-8, 1996.
- [Brin and Page, 1998] Brin, S. and Page, L., The Anatomy of a Large-scale
Hypertextual Web Search Engine, *The Seventh International World
Wide Web Conference (WWW7)*, pp. 107-117, 1998.
- [Broglia et al., 1995] Broglia, J., Callan, J. P., Croft, W. B., and Nachbar, D.
W., Document Retrieval and Routing Using The INQUERY System,
The Third Text REtrieval Conference (TREC-3), pp. 29-38, 1995.
- [Callan et al., 1992] Callan, J. P., Croft, W. B. and Harding, S. M., The
INQUERY Retrieval System, *Proceedings of the Third International
Conference on Database and Expert Systems Applications*, Springer,
pp. 78-83, 1992.

- [Chakrabarti et al., 1997] Chakrabarti, S., Dom, B., Gibson, D., Kumar, S. R., Raghavan, P., Rajagopalan, S., and Tomkins, A., Experiments in Topic Distillation, *A CM - SIGIR '98 Post-Conference Workshop on Hypertext Information Retrieval for the Web*, 1998.
- [Croft and Harper, 1979] Croft, W. B. and Harper, D. J., "Using Probabilistic Models of Document Retrieval without Relevance Information, *Journal of Documentation*, 35(4), pp. 285-295, 1979.
- [Cutler et al., 1999] Cutler, M., Deng, H., Maniccam, S and Meng, W., A New Study on Using HTML Structures to Improve Retrieval, *The Eleventh IEEE Conference on Tools with AI*, pp. 406-409, 1999.
- [Goldberg, 1989] Goldberg, D. E., *Genetic Algorithms in Search, Optimization, and Machine Learning*, Addison-Wesley, 1989.
- [Gordon, 1988] Gordon, M., Probabilistic and Genetic Algorithms for Document Retrieval, *Communications of the ACM* 31, pp. 1208-1218, 1988.
- [Holland, 1975] Holland, J. H., *Adaptation in Natural and Artificial Systems*, University of Michigan Press, 1975.
- [Lawrence and Giles, 1998] Lawrence, S. and Giles, C. L., Searching the World Wide Web, *Science*, Vol. 280, pp. 98-100, 1998.
- [Maron and Kuhns, 1960] Maron, M. E. and Kuhns, J. L., On Relevance, Probabilistic Indexing and Information Retrieval, *Association for Computing Machinery*, 7(3), pp. 216-244, 1960.
- [Mauldin, 1997] Mauldin, M. L., Lycos: Design Choices in an Internet Search Service, *IEEE Expert*, 12(1), pp. 8-11, 1997.

- [Michalewicz, 1992] Michalewicz, Z., *Genetic Algorithms + Data Structures = Evolutionary Programs*, Springer, pp. 104-105, 1992.
- [NIST] NIST, *Text REtrieval Conference homepage*, <http://trec.nist.gov>.
- [Picard, 1998] Picard, J., Modeling and Combining Evidence Provided by Document Relationships Using Probabilistic Argumentation Systems, *Proceedings of the ACM SIGIR '98 Conference*, pp. 182-189, 1998.
- [Robertson and Sparck, 1976] Robertson, S. E. and Sparck Jones, K., Relevance Weighting of Search Terms, *Journal of the American Society for Information Science*, 27, pp. 129-146, 1976.
- [Robertson et al., 1995] Robertson, S. E. et al., Okapi at TREC-3, *The Third Text REtrieval Conference (TREC-3)*, pp. 109-126, 1995.
- [Salton et al., 1975] Salton, G., Wong, A. and Yang, C. S., A Vector Space Model for Automatic Indexing, *Communications of the ACM* 18, pp. 613-620, 1975.
- [Salton et al., 1983] Salton, G., Fox, E. A., and Wu, H., Extended Boolean Information Retrieval, *Communications of the ACM*, Vol. 26, No. 11, pp. 1022-1036, 1983.
- [Salton, 1989] Salton, G., *Automatic Text Processing*, Addison-Wesley, pp. 279-281, 1989.
- [Shin and Zhang, 1998] Shin, D. H. and Zhang, B. T., A Two-Stage Retrieval Model for the TREC-7 Ad Hoc Task, *The Seventh Text REtrieval Conference (TREC-7)*, pp. 501-507, 1998.
- [Spertus, 1997] Spertus, E., ParaSite: Mining Structural Information on the Web, *The Sixth International World Wide Web Conference (WWW6)*, pp. 1205-1215, 1997.

- [Turtle and Croft, 1991] Turtle, H. and Croft, W. B., Evaluation of an Inference Network-based Retrieval Model, *ACM Transactions on Information Systems*, Vol. 9, No. 3, pp. 187-222, 1991.
- [Voorhees and Harman, 1999] Voorhees, E. M. and Harman, D., Overview of the Eighth Text Retrieval Conference, *The Eighth Text REtrieval Conference (TREC-8)*, pp. 1-27, 1999.
- [Weiss et al., 1996] Weiss, Ron., Véléz, B., and Sheldon, M. A., HyPursuit: A Hierarchical Network Search Engine that Exploits Content-Link Hypertext Clustering, *Proceedings of the Seventh ACM Conference on Hypertext*, pp. 180-193, 1996.
- [Yang et al., 1993] Yang, J., Korfhage, R. R., and Rasmussen, E., Query Improvement in Information Retrieval using Genetic Algorithms: A Report on the Experiments of the TREC Project, *The First Text REtrieval Conference (TREC-1)*, pp. 31-58, 1993.
- [Yang and Honavar, 1998] Yang, J. and Honavar, V., *Feature Extraction, Construction and Selection - A Data Mining Perspective*, Kluwer Academic Publishes, pp. 117-136, 1998.
- [Zhang, 1995] Zhang, B.-T., Learning and Optimization by Artificial Evolution (in Korean), *The Institute of Control, Automation and Systems Engineers Magazine*, Vol. 1, No. 3, pp. 52-61, 1995.
- [Zobel, 1998] Zobel, J., How Reliable are the Results of Large-Scale Information Retrieval Experiments?, *Proceedings of the ACM SIGIR '98 Conference*, pp. 307-314, 1998.

Abstract

This paper presents a method for web-document retrieval by learning importance factors for tags which are used for document structuring. Web documents are usually written in Hypertext Markup Language (HTML). HTML consists of tags which make a document into a specific form and a homepage is designed using the tags according to its object. In this paper, we propose a method for improving the retrieval performance using the information of HTML tags. The importance factors for the tags are learned using a genetic algorithm. A tag is mapped into a gene and a set of tags represented as a chromosome. The results obtained by genetic learning are the weights for tag importance, and provided the retrieval engine as the weights of documents.

Experiments have been performed on an artificial dataset and a large collection of TREC (Text REtrieval Conference) documents. Our empirical results show that this algorithm learns the weights by tag importance factors, and can improve the retrieval performance on top-ranked documents.

Keywords: Information Retrieval, Web-Document, Tag Weight, Genetic Algorithm

가

가
가

가

가

가

2

?

NLP/AI 가

가

Also, Thanks To:

My brother his wife... , ...
 , , , , ... , ...
 , , , ... , , , , , ,
 , , , ? ...
 , , , , , , , , , ,

 , , , , , , , ...
 ... 가 , , , , , , ,
 ... , ...
 , , , , , , ,
 ... CCMGer , , , , , , ,
 , ... , ... Stryper, Michael W.
Smith, Impellitteri... dog Roxi...
Gromit...
...