

가 *

A Study on Information Resource Evaluation for Text C ategorization

(EunKyung Chung)**

가가 (full text) 가 가

가가

ABSTRACT

The purpose of this study is to examine whether the information resources referenced by human indexers during indexing process are effective on Text Categorization. More specifically, information resources from bibliographic information as well as full text information were explored in the context of a typical scientific journal article data set. The experiment results pointed out that information resources such as citation, source title, and title were not significantly different with full text. Whereas keyword was found to be significantly different with full text. The findings of this study identify that information resources referenced by human indexers can be considered good candidates for text categorization for automatic subject term assignment.

: text Categorization, automatic indexing, information resources, subject indexing procedure

*

**

(echung@ewha.ac.kr)

:2007 12 1

:2007 12 1

1.

(full text)

(feature selection),

가

가

(Cunningham, Witten, & Litten, 1999; Sebastiani, 2002; 2005).

가

가가

(Moens, 2000).

(Collocate)

가

(O'Connor, 1996).

가

가

, 가가

가

가

가

가

2.

가

(machine learning)

2.1

(SVM)

(text categorization)

(machine learning)

(classifier)

(Naive Bayes), (Support Vector Machine), (Neural Network), kNN (k - nearest neighbors) (Sebastini, 2002).

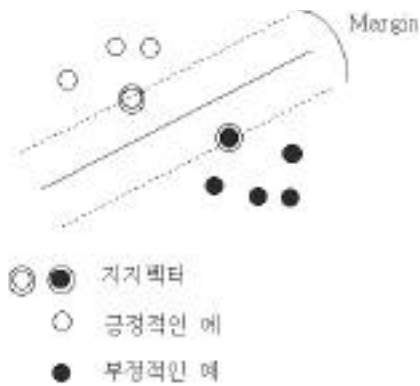
(SVM) Joachims (1998)

가 (full text) (Cunningham, Witten, & Litten, 1999; Sebastiani, 2002; 2005).

(Sebastiani, 2005). < 1>

(margin)

(Joachims, 1998).
2.2



< 1> (SVM)

Larkey(1999) kNN (20) 가 Efron, Elsas, Marchionini, & Zhang (2004)

. , Zhang, et al. (2004)

. Zhang,
et al. (2004) , Slattery
(2002)

가

.
Diaz, Ranilla, Montanes,
Fernandez, & Combarro (2004)
(contextual information) (i.e.
)

가
(full text)

3.

가
,
가가

가 .

가가

가가

(Information Resources)

가 .
가가
가

. ,
가

. Mai (2005)가 ,
(Mitchell, et al., 2003)

, , , ,

(International Standard Organization,
ISO 5963:1985) , , ,

. Chan
(1987) Sauperl (2002)

. Chan
Sauperl Chan (1981),
Foskett (1996), Taylor (2003) ,

, , , ,

. Foskett ,
, Taylor /

, ,
. , Chan

(document attributes)

. Chan , , , ,

,
, 가

가 (Chu & O'Brien, 1993; Jeng, 1996; Sauperl, 2002; 2004).
Jeng 가 가
. Sauperl

가
Chu & O'Brien

4.

4.1

가가
가
가

(full text)

4.2

가
(full-text) 가
가 가
(INSPEC)
가

8
3 5 1 5
(Engineering Village 2, n.d.).

1,000 가
20
50

(homogeneous data set)
(heterogeneous data set)

(computer science)

(information technology) 2006 .
 10
 . 1969
 (software) 10
 .
 . 2000 2006
 . < 1> .
 10 가
 , 4.3
 10
 . (full - text) , ,
 5 , , ,
 (engineering village 2) .
 1) , 2) 8
 20 , 3) ,
 , 4) PDF
 , 5) 2000 . ,
 < 1>

software architecture	computer architecture
software development management	computer graphics
software libraries	computer interface
software maintenance	discrete systems
software metrics	information management
software portability	knowledge based systems
software prototyping	pattern recognition
software quality	reliability
software reliability	software engineering
software reusability	user interfaces

PDF ASCII (feature selection)

Adobe Acrobat Capture 8

가 , 8가
(full-text)

4가 < 2>

(full text)

()

8가

(Information Gain), (DF),
(X^2 statistics)

가

가

8가

Porter (1980)

< 2> 8

	121,699	61,023	60,676
	472,985	244,502	228,483
	275,253	149,797	125,456
	6,602,440	3,507,097	3,095,343
	517,197	275,728	241,469
	32,861	15,625	17,237
	4,074	2,026	2,048
	7,553	3,690	3,863
	8,034,062	4,259,488	3,774,575

WEKA (Witten & Frank, 2000)

(Micro-averaging)

WEKA

ARFF (.arff)

(Diaz, et al., 2004).

4.4 가

가

가 , F- 가
 . (Lewis, 1995; Sebastiani, 2002; Yang, 1999)

4.5

WEKA (Witten & Frank, 2000) WEKA

F- (van Rijsbergen, 1979)

$$(R) = a / (a + c)$$

$$(P) = a / (a + b)$$

$$F = 2PR / (P + R)$$

< 3> 가

	a	b
	c	d

WEKA (SVM)
 "weka.classifiers.functions.SMO"

8

3가

(full text)
 (baseline)

7

가

가 (Macro-averaging)

(Micro-averaging)

(Yang, 1999).

(Macro-averaging)

5.

5.1 가
 가 ()
 , ,
 가 , F-

< 4>

가

< 4>

			F -
	0.277	0.234	0.230
	0.344	0.290	0.308
	0.206	0.193	0.193
	0.349	0.283	0.300
	0.283	0.213	0.231
	0.387	0.366	0.368
	0.323	0.304	0.299
	0.319	0.293	0.296

< 5>

가

< 5>

			F -
	0.283	0.275	0.262
	0.345	0.310	0.322
	0.249	0.244	0.243
	0.287	0.258	0.261
	0.286	0.234	0.247
	0.402	0.382	0.384
	0.267	0.269	0.262
	0.310	0.292	0.295

가< 6>

, F- 가

가

가

.< 6>

가

< 6>

			F -
	0.425	0.393	0.389
	0.493	0.459	0.469
	0.371	0.349	0.348
	0.564	0.505	0.520
	0.444	0.426	0.427
	0.587	0.568	0.569
	0.461	0.420	0.432
	0.487	0.422	0.437

가

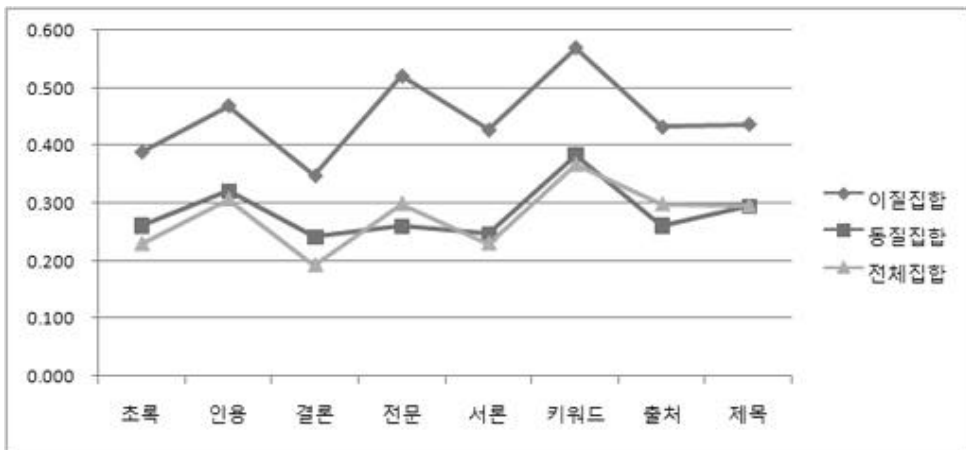
(, 가)
 < 2> 5.2 가

F- 8 F- t- , 8

, < 7> (full text)

7 t- 95% 가 , < 2>

가



< 2> (F-)

< 9> F- t- ()

		0.148	0.335	0.270	0.006	0.441	0.346
			0.001	0.170	0.069	0.526	0.544
				0.002	0.001	0.120	0.047
					0.001	0.914	0.732
						0.004	0.001
							0.879

가 가
가 , , , ,
(full-text)
가 가

5.3

가 가 . < 10>

< 10> F- t-

	t	p	t	p	t	p
	3.446	0.003	-0.025	0.981	3.273	0.010
	-0.249	0.806	-0.857	0.414	1.252	0.242
	4.481	0.000	0.500	0.629	4.333	0.002
	2.587	0.018	0.310	0.763	2.275	0.051
	-2.231	0.038	-5.212	0.001	-0.753	0.471
	0.043	0.966	-0.015	0.988	1.402	0.194
	0.110	0.913	-1.267	0.237	1.321	0.219

*p<0.05

< 11> t-

	t	p	t	p	t	p
	1.443	0.165	-0.441	0.669	1.763	0.112
	-0.160	0.875	-0.601	0.563	1.398	0.196
	3.432	0.003	0.300	0.771	4.176	0.002
	1.787	0.090	0.339	0.742	1.619	0.140
	-2.935	0.008	-3.945	0.003	-1.057	0.318
	-0.852	0.405	-0.852	0.405	1.667	0.128
	-0.354	0.727	-0.354	0.727	1.733	0.117

*p<0.05

< 12> t-

	t	p	t	p	t	p
	2.356	0.029	0.120	0.907	2.306	0.047
	0.131	0.897	-0.925	0.379	1.096	0.302
	4.523	0.000	1.086	0.306	2.611	0.028
	2.389	0.027	0.010	0.992	2.034	0.073
	-0.924	0.367	-4.757	0.001	-0.286	0.781
	0.704	0.490	0.476	0.645	1.235	0.248
	0.644	0.527	-.510	0.623	0.938	0.373

*p<0.05

95% F- , < 11> F-

F-

t- 가 < 12>

(full text)

6.

가

F-

가

(text categorization) 가
가

가
Efron, Elsas, Marchionini, &
Zhang (2004), Zhang, et al. (2004),
Slattery (2002)

(full text)

가

t-

가가

가가

가

(full text)

t-

가

가

가

가

가

Chan, L.M. (1981). *Cataloging and classification: An introduction*. New York City, NY: McGraw-Hill.

Chan, L.M. (1987). Instructional materials used in teaching cataloging and classification. *Cataloging and Classification*, (7) : 131-144.

Chu, C.M. & O'Brien, A. (1993). Subject analysis: The critical first stage in indexing. *Journal of Information Science* (19) : 439-454.

Cunningham, S.J., Witten, I.H., & Litten, J. (1999). Applications of machine learning in information retrieval. *Annual Review of Information Science*

and Technology, (34) : 341-384.

Diaz, I., Ranilla, J., Montanes, E., Fernandez, J., & Combarro, E. (2004). Improving performance of text categorization by combining filtering and support vector machines, *Journal of the American Society for Information Science and Technology*, 55(7) : 579-592.

Efron, M., Marchionini, G., Elsas, J., & Zhang, J. (2004). Machine learning for information architecture in a large governmental website. *Proceedings of the 2004 Joint ACM/IEEE Conference on Digital*

- Libraries 151-159.
- Engineering Village 2. (n.d.). Retrieved November 11, 2006, from <http://www.ingvillage2.org/controller/servlet/Controller>.
- Foskett, A.C. (1996). *The Subject Approach to Information*. London: Library Association Publishing.
- ISO 5963: 1985. (1985). *Documentation - methods for examining documents: Determining their subjects and selecting indexing terms* International Standards Organization.
- Jeng, L.H. (1996). Using verbal reports to understand cataloging expertise: Two cases, *Library Resources and Technical Services* 40(4) : 343-358.
- Joachims, T. (1998). Text categorization with support vector machine: Learning with many relevant features, *Proceedings of the 10th European Conference on Machine Learning*, 137-142.
- Larkey, L.S. (1999). A patent search and classification system. *Proceedings of the 4th ACM Conference on Digital Libraries* 179-187.
- Lewis, D.D. (1995). Evaluating and optimizing autonomous text categorization systems. Unpublished Doctoral Dissertation, University of Massachusetts, Massachusetts.
- Mai, J.E. (2005). *Analysis in indexing: document and domain centered approaches*, *Information Processing and Management*, (41) : 599-611.
- Mitchell, J.S. et al. (Eds.). (2003). *Dewey Decimal Classification and Relative Index*. Dublin, OH: OCLC Online Library Computer, Inc.
- Moens, M.F. (2000). *Automatic Indexing and Abstracting of Document Texts*. Norwell, MS: Kluwer Academic Publishers.
- O'Connor, B.C. (1996). *Explorations in Indexing and Abstracting: pointing, virtue, and power*. CO: Libraries Unlimited.
- Porter, M.F. (1980). An algorithm for suffix stripping, *Program*, (14) : 130-137.
- Sauperl, A. (2002). *Subject determination during the cataloging process*. Lanham, MD; Scarecrow Press.
- Sauperl, A. (2004). *Catalogers' common ground and shared knowledge*. *Journal of the American Society for Information Science and Technology* 55(1) : 55-63.

- Sebastiani, F. (2002). Hypertext categorization. In A. Zanasi (Eds.), Text Mining and Its Applications (pp. 109-129), Southampton, U.K.: WIT Press.
- Sebastiani, F. (2005). Text categorization. In A. Zanasi (Eds.), Text mining and its applications (pp. 109-129), Southampton, U.K.: WIT Press.
- Slattery, S. (2002). Hypertext categorization Unpublished Doctoral Dissertation. School of Computer Science. Carnegie Mellon University.
- Taylor, A.G. (2003). The organization of information (2nd ed.). Englewood, CO; Libraries Unlimited.
- van Rijsbergen, C.J. (1979). Information Retrieval Butterworths, London.
- Witten, I.H. & Frank, E. (2000). Data Mining: Practical Machine Learning Tools and Techniques with JAVA Implementations. CA: San Diego, Academic Press.
- Yang, Y. 1999. An evaluation of statistical approaches to text categorization. Information Retrieval (1) : 69-90.
- Zhang, B., Goncalves, M.A., Fan, W., Chen, Y., Fox, E.A., Calado, P. & Cristo, M. (2004). Combining structural and citation-based evidence for text categorization, Proceedings of the 13th ACM Conference on Information and Knowledge Management, 162-163.

