

{sbpark,btzhang}@scai.snu.ac.kr

# Text Categorization Using Both Lexical Information and Syntactic Information

Seong-Bae Park<sup>o</sup> Byoung-Tak Zhang

School of Computer Science and Engineering, Seoul National University

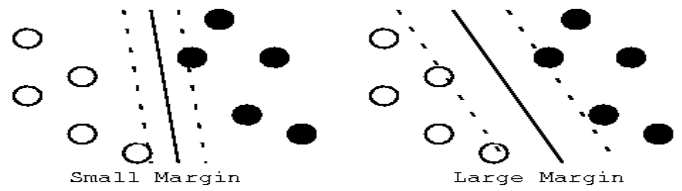
가

Reuters-21578

Support Vector Machine 97%

0.63% 가

1. (text categorization)



1. SVMs

가

(machine learning) TF•IDF

(syntactic information)

21578, SVM 96.89%

, 0.63%

가 [6] (formal definition)가

## 2. SVM

Support Vector Machines

SVMs

가

[5]

(full parsing)

(text chunking)[2]

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \quad \mathbf{x}_i \in R^n, \quad y_i \in \{+1, -1\}$$

$$(\mathbf{w} \cdot \mathbf{x}) + b = 0$$

(classifier)

margin

Support

Margin 1

SVMs margin

margin(d)

Vector Machines

$$(\mathbf{w} \cdot \mathbf{x}) + b = \pm 1, \quad d = 2 / \|\mathbf{w}\|$$

Reuters-21578

SVM

$$y_i [(\mathbf{w} \cdot \mathbf{x}) + b] \geq 1$$

$\|\mathbf{w}\|$

0.63% Reuters-

$\mathbf{w}$   $b$

3. 3.1

$TF \cdot IDF$  가  $d(i)$  ,  $TF$  ,  $IDF$

$$d(i) = TF(w_i, d) \cdot IDF(w_i)$$

$$IDF(w_i) = \log\left(\frac{|D|}{DF(w_i)}\right)$$

document frequency DF  $w_i$  가

Feature	Description
SF1	Detected NPs / total detected chunks
SF2	Detected VPs / total detected chunks
SF3	Detected PPs / total detected chunks
SF4	Detected Os / total detected chunks
SF5	Words included in NPs / detected NPs
SF6	Words included in VPs / detected VPs
SF7	Words included in PPs / detected PPs
SF8	Words included in Os / detected Os

1. (feature).

$TF \cdot IDF$  algorithm bow [7] Porter's stemming stoplist

3.3 SVM

3.2

가 margin SVMs (chunking)

SVMs 가 margin SVMs margin 가

90 가

4.

CoNLL-2000 shared task<sup>1</sup> SVMs CoNLL-2000 1987 23 가 NP, VP, PP, O O B-X I-X I-X X B-X X SVMs 가 SVM

4.1

Carnegie Reuters-21578 가 135 10 (“ModLewis”, “ModApte”, “ModHayes”가 “ModApte” 9,603 3,299 SVMs bow  $TF \cdot IDF$  SVMs 2000 shared task Reuters-21578 가 가 Brill's Tagger[3]

(pairwise classification)

1

Stamatatos (style markers)

(author)

[1].

가 1

4

(grammatical phrase)

4 가

가

, Brill's Tagger

,

,

,

,

, CoNLL-2000

SVMs

<sup>1</sup> <http://lcg-www.uia.ac.be/conll2000/chunking>

4.2

2 TF-IDF 가 8  
 Support Vector Machine  
 'Lexical' TF-IDF SVM  
 'Syntactic' Reuters-21578  
 10 가  
 가  
 'Grain' 'Earn', 'Corn'  
 'Acq'  
 Reuters-21578  
 가

Class	Accuracy		
	Syntactic	Lexical	Both
Earn	91.42%	95.09%	95.30%
Acq	78.21%	93.76%	93.73%
Money-fx	94.75%	96.15%	96.15%
Grain	95.48%	95.48%	95.51%
Crude	94.27%	97.00%	97.00%
Trade	96.45%	97.79%	97.79%
Interest	96.03%	97.18%	97.18%
Ship	97.30%	98.15%	98.15%
Wheat	97.85%	98.91%	98.91%
Corn	98.30%	99.12%	99.18%

2.

margin , 3.3  
 3 Increase  
 0.63%  
 'Acq' 가  
 1.48% 'Wheat' 'Corn'  
 98.91% 가  
 가가

6.

Support Vector  
 Machine  
 Support Vector Machine , Reuters-  
 21578 97.52%  
 0.63%

Class	Accuracy	Increase
Earn	96.61%	1.31%
Acq	95.21%	1.48%
Money-fx	97.12%	0.97%
Grain	95.51%	0.00%
Crude	97.67%	0.67%
Trade	98.42%	0.63%
Interest	97.67%	0.49%
Ship	98.58%	0.43%
Wheat	99.15%	0.24%
Corn	99.27%	0.09%
Average	97.52%	0.63%

3.

Co-  
 Training[8] Co-  
 Training 가  
 TREC filtering  
 (unlabeled  
 data)  
 (KOSEF)  
 (AITre) BK 21

[1] E. Stamatatos, N. Fakotakis, and G. Kokkinakis, "Automatic Text Categorization in Terms of Genre and Author," *Computational Linguistics*, Vol. 26, No. 4, pp. 471-495, 2000.  
 [2] T. Kodoh and Y. Matsumoto, "Use of Support Vector Learning for Chunk Identification," In *Proceedings of CoNLL-2000 and LLL-2000*, pp. 142-144, 2000.  
 [3] E. Brill. Rule based tagger. <http://www.cs.jhu.edu/~brill/>.  
 [4] T. Joachims. *SVM<sup>light</sup>* version 3.02. [http://www-ai.cs.uni-dortmund.de/SOFTWARE/SVM\\_LIGHT/svm\\_light.eng.html](http://www-ai.cs.uni-dortmund.de/SOFTWARE/SVM_LIGHT/svm_light.eng.html).  
 [5] T. Joachims, "Text Categorization with Support Vector Machines: Learning with Many Relevant Features," In *Proceedings of ECML 98*, pp. 137-142, 1998.  
 [6] D. Biber, "Dimensions of Register Variation: A Cross-Linguistic Comparison," *Cambridge University Press*, 1995.  
 [7] A. McCallum and Andrew Kachites. "Bow: A toolkit for statistical language modeling, text retrieval, classification and clustering." <http://www.cs.cmu.edu/~mccallum/bow>. 1996.  
 [8] A. Blum and T. Mitchell, "Combining labeled and unlabeled data with co-training," In *Proceedings of 11<sup>th</sup> Annual Conference of Computational Learning Theory*, pp. 92-100, 1998.