

Resolving Part-of-Speech Tagging Ambiguities by a Maximum Entropy Boosting Model

Seong-Bae Park^o and Yung Taek Kim
School of Computer Science and Engineering, Seoul National University
^o{sbpark, btzhang}@bi.snu.ac.kr

가 ,
(classification problem)
(maximum entropy boosting model)
96.78%

1. (maximum entropy boosting model)

가 [2]. (text chunking)
96% [1].

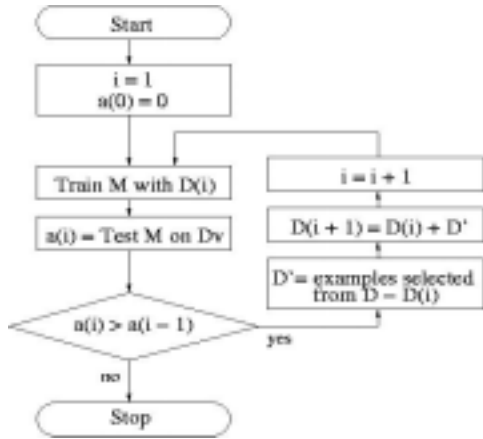
2.

Weischedel et al.

[1]. ,
가 (machine learning) (classifier)
(classification problem) (modeler) 가
가 GIS

, Park Zhang

가



1.

		가	
w_{i-2}	$i-2$	w_{i-2}	$i-2$
w_{i-1}	$i-1$	w_{i-1}	$i-1$
w_{i+1}	$i+1$	w_i	i
w_{i+2}	$i+2$	w_{i+1}	$i+1$
t_{i-2}	$t-2$	w_{i+2}	$i+2$
t_{i-1}	$i-1$	t_{i-2}	$t-2$
prefix(w_i, j)	j	t_{i-1}	$i-1$
suffix(w_i, j)	j	w_i	i
hasnumber	가	가?	
hasupper	가	가?	
hashyper		가?	

$$p(t_i | h_i) = \frac{1}{Z} \exp\left(\sum_i \lambda_i f_i(h_i, t_i)\right)$$

1. D , M , D_v
 가
 Park Zhang[1]

$f_i(h_i, t_i)$, λ_i , f_i , 가

$$h_i = \{w_{i-2}, w_{i-1}, w_i, w_{i+1}, w_{i+2}, t_{i-2}, t_{i-1}\}$$

(decision tree)

if-then

if-then

n -gram

(active learning)

1

AdaBoost

3.

w_1, \dots, w_N

POS

$p(t_1, \dots, t_N | w_1, \dots, w_N)$

t_1, \dots, t_N

가

가

$$p(t_1, \dots, t_N | w_1, \dots, w_N) = \prod_{i=1}^N p(t_i | h_i)$$

h_i , w_i

$p(t_i | h_i)$

4.

4.1

Penn Treebank II Wall Street Journal

1,173,765, 49,206

, 45

가

. 60% 704,251

20% 234,819

, 20%

2.

AdaBoost.MI	96.72%
	96.89%
	96.36%
	96.78%

3.

	96.78%
	92.19%

234,695

23,237

3,713,

3,199

4,578

4.2

2

AdaBoost.MI [3]

AdaBoost

“attribute=value”

가

[4].

96.36%

‘AdaBoost.MI’

96.78%

가

‘AdaBoost.MI’

가

가

가

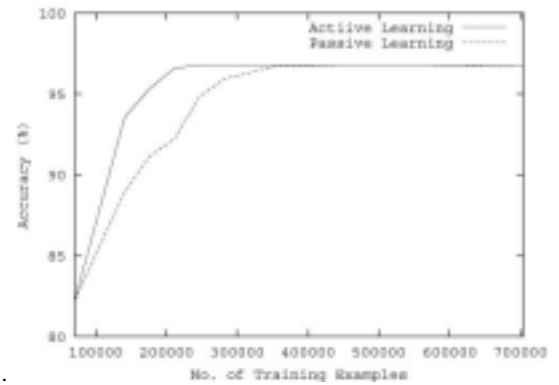
3

가 92.19%

2

40%

25%



2.

5.

Wall Street Journal

96.78%

가

가

BrainTech

BK 21

[1] R. Weischedel, M. Meteor, R. Schwartz, L. Ramshaw, and J. Palmucci, “Coping with ambiguity and unknown words through probabilistic models,” *Computational Linguistics*, 19(2), pp. 359-382, 1994.

[3] S.-B. Park and B.-T. Zhang, “A boosted maximum entropy model for learning text chunking,” In *Proceedings of the 19th International Conference on Machine Learning*, pp. 482-489, 2002.

[3] S. Abney, R. Schapire, and Y. Singer, “Boosting Applied to Tagging and PP-attachment,” In *Proceedings of the Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, pp. 38-45, 1999.

[4] A. Ratnaparkhi, “A Maximum Entropy Model for Part-of-speech Tagging,” In *Proceedings of the Empirical Methods in Natural Language Processing*, pp. 133-142, 1996.