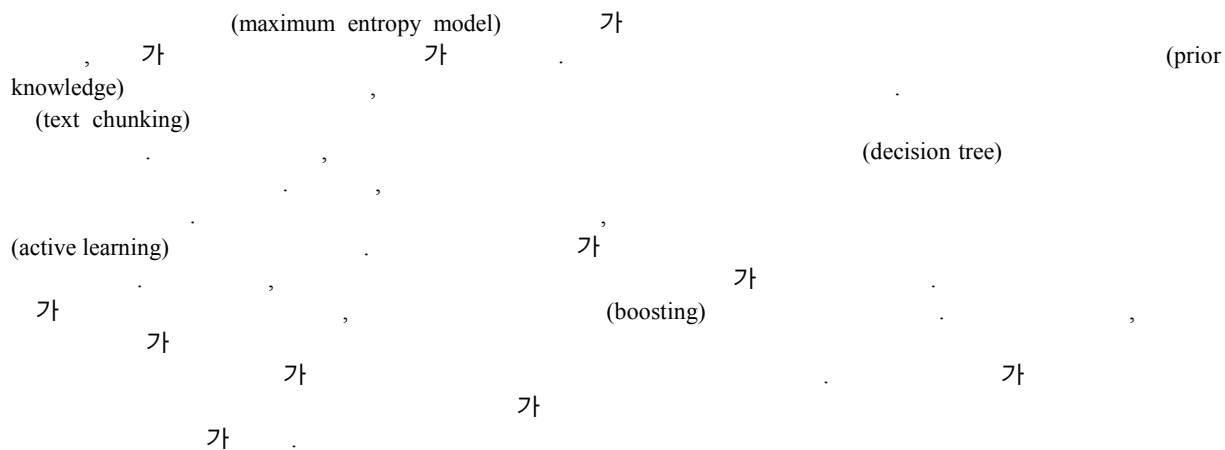


Learning Text Chunking Using Maximum Entropy Models

Seong-Bae Park^o and Byoung-Tak Zhang
School of Computer Science and Engineering
Seoul National University

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

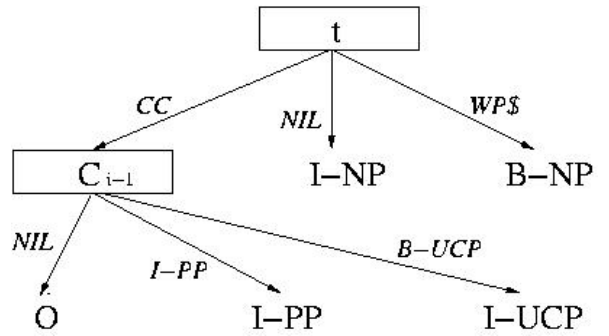


1.

Koeling
 (text chunking) [9].
 Abney 가 (prior knowledge) (feature)
 [1], 가
 n , n
 m
 가 $O(n^3)$
 Ramshaw Marcus 가
 [14],
 가 [2,5,7,10].
 (maximum entropy) 가 then if-
 [3]

if-then

(active learning)



1.

2

가

$p(x, y)$ Generalized Iterative Scaling (GIS) μ_i

, 3

[6]

4

\tilde{p}

f_i

μ_i
 K_i , μ_i

, 5

. 6

$$\mu_i = \mu_i \cdot \frac{K_i}{E_{\tilde{p}}[f_i]}$$

\tilde{p} μ_i 가

7

2.

$$\tilde{p}(x, y) = \frac{\prod_i \mu_i^{f_i(x, y)}}{Z}$$

f_i 가 가

(solution)가

[2].

가

가 가

f_i 가

f_i $E_{\tilde{p}}[f_i]$

$$p(x, y) = \frac{\prod_i \mu_i^{f_i(x, y)}}{Z}$$

$E_{\tilde{p}}[f_i]$

가

, μ_i

$E_{\tilde{p}}[f_i]$

$$\mu_i = \exp(\lambda_i)$$

if-then

if-then

, λ_i
 $f(x, y)$

(x, y)

$S = \{(x_1, y_1), \dots, (x_N, y_N)\}$
 K_i

n -gram

$y_N\}$

$f_i(x, y)$
 p

$$E_p[f_i] = \sum_{j=1}^N p(x_j, y_j) f_i(x_j, y_j) = K_i$$

K_i

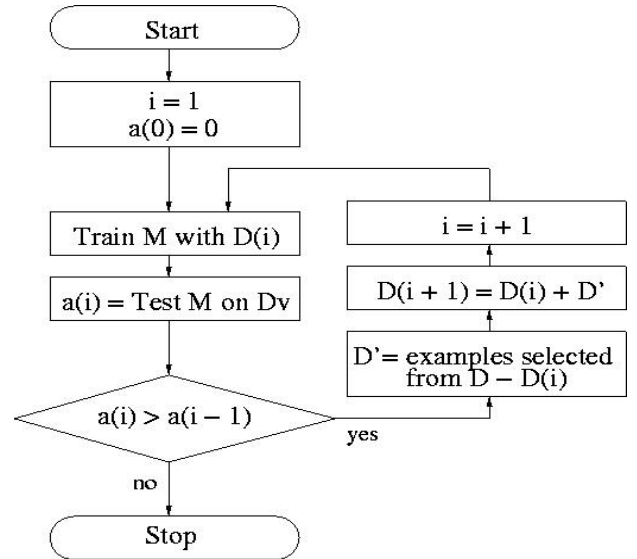
S

\tilde{p}

$$E_{\tilde{p}}[f_i] = \sum_{j=1}^N \tilde{p}(x_j, y_j) f_i(x_j, y_j) = E_p[f_i]$$

3.

- if-then
1. $f_1(t, C_{i-1}, C) = \begin{cases} 1 & \text{if } t = WPS \text{ and } C = B - NP \\ 0 & \text{otherwise} \end{cases}$
 2. $f_2(t, C_{i-1}, C) = \begin{cases} 1 & \text{if } t = NIL \text{ and } C = I - NP \\ 0 & \text{otherwise} \end{cases}$
 3. $f_3(t, C_{i-1}, C) = \begin{cases} 1 & \text{if } C_{i-1} = B - UCP \text{ and } t = CC \text{ and } C = I - UCP \\ 0 & \text{otherwise} \end{cases}$
 4. $f_4(t, C_{i-1}, C) = \begin{cases} 1 & \text{if } C_{i-1} = I - PP \text{ and } t = CC \text{ and } C = I - PP \\ 0 & \text{otherwise} \end{cases}$
 5. $f_5(t, C_{i-1}, C) = \begin{cases} 1 & \text{if } C_{i-1} = NIL \text{ and } t = CC \text{ and } C = O \\ 0 & \text{otherwise} \end{cases}$



2. , D . M , Dv

μ_i

, μ_i

GIS

(first-order feature)

(high-order feature)

100 가 t, C_{i-1}, C_i 가 10⁶ 가 5 가 1 가 1
 Quinlan C4.5 release 8[13] n-gram

$$p^* = \arg \max_{p \in C} H(p)$$

$$E_p[f_i] = E_{\tilde{p}}[f_i]$$

\tilde{p} 가 (log-likelihood)

$$L_{\tilde{p}}(p) = \sum_{(x,y)} \tilde{p}(x,y) \log p(y|x) = \log \prod_{(x,y)} p(y|x)^{\tilde{p}(x,y)}$$

$$p^* = \arg \max_{p \in C} H(p) = \arg \max_{p \in Q} L_{\tilde{p}}(p)$$

$$Q = \{p \mid p(x) = \pi \prod_i \alpha_i^{f_i(x)}, 0 < \alpha < \infty\}$$

가 가 (variance)

4.

GIS

$$E_{\tilde{p}}[f_i]$$

M

N

$$O(M \cdot N)$$

M

$$E_{\tilde{p}}[f_i]$$

M

가

[3,12], N

가 (uncertainty)

$E_{\tilde{p}}[f_i]$ M , N , e (posterior distribution)
(uniform distribution)

$E_{\tilde{p}}[f_i]$ (relative entropy)

$E_{\tilde{p}}[f_i] \approx \frac{1}{n} \sum_{j=1}^n f_i(x_j, y_j)$ *Kullback-Leibler divergence*(KL-divergence)

$(x_1, y_1), \dots, (x_n, y_n)$
(random sample)

$$KL(p \parallel q) = \sum_{c_i \in C} p(c_i) \ln \frac{p(c_i)}{q(c_i)}$$

p , $E_{\tilde{p}}[f_i]$ 가, C 가 e^*

$$e^* = \arg \max_{e_j \in D-D(i)} KL(U \parallel p(C | e_j))$$

Monte Carlo Markov Chain

(MCMC) 가 Gibbs sampling, Metropolis sampling,

perfect sampling

가 p

5. AdaBoost

CoNLL-2000

CoNLL-2000

12 (NP, PP, VP, O)가

95.10%
(recall)

O Zhang AdaBoost
AdaBoost 가
[18].

AdaBoost (weak classifier)

2

가, AdaBoost 가

D , $D(0)$ 가 $D(0) \subset D$

(loop) 가

D_v , $D(i)$, i

AdaBoost 가 (committee model)

λ , $D(i)$, $D - D(i)$, $D(i+1)$

6.

CoNLL-2000 Shared Task [11]
[5].

‘Query by Committee’(QBC)
[15]. Freund et al.

CoNLL-2000 Wall Street Journal(WSJ)

[8]. 가 n 가

WSJ section 15-18

$O(1/n)$

211,727 WSJ

section 20 47,377

QBC

version space

가 (3).

가

Brill tagger

He PRP B-NP
 reckons VBZ B-VP
 the DT B-NP
 current JJ I-NP
 deficit NN I-NP
 will MD B-VP
 narrow VB I-VP
 to TO B-PP
 only RB B-NP
 # # I-NP
 1.8 CD I-NP
 billion CD I-NP
 in IN B-PP
 September NNP B-NP
 O

3. CoNLL-2000

가 가 , B-NP
 , I-NP , 11
 가 O 가
 , O
 23 O
 CoNLL-
 2000
 148,181 Dv 63,546 D

[11]

1,273 638

8.9

6.2

bigram

trigram

W_i

$(POS_{i-2}, POS_{i-1}, POS_i), (w_{i-1})$
 (C_{i-2}, C_{i-1}) 가
 가

[11] k -NN

(w_{i-2}, w_{i-1}, w_i)
 (E_i)

2

POS_{i-2}	W_{i-2}
POS_{i-1}	W_{i-1}
POS_i	W
POS_{i+1}	W_{i+1}
POS_{i+2}	W_{i+2}
W_i	W
C_{i-2}	W_{i-2}
C_{i-1}	W_{i-1}

1.

W_{i-2}	W_{i-2}
W_{i-1}	W_{i-1}
W_i	W
POS_{i-2}	W_{i-2}
POS_{i-1}	W_{i-1}
POS_i	W
C_{i-2}	W_{i-2}
C_{i-1}	W_{i-1}

2.

6.3

3 CoNLL-2000

96.72%

92.48

F-score

C4.5

F-score

2.34% 2.28

2,663

1,554

41.64%

가

가

가

4

[20]

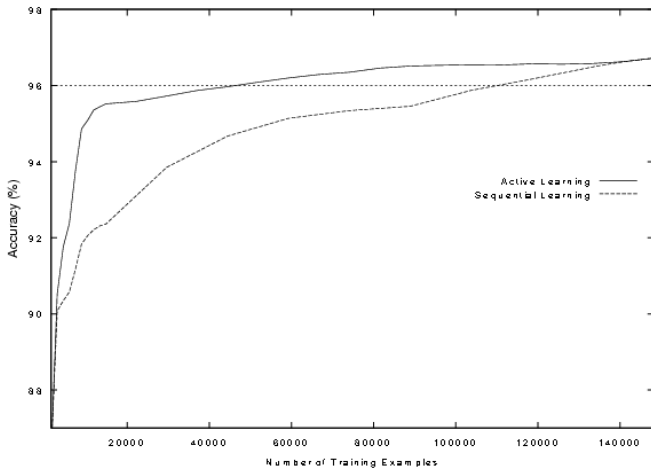
[21]

'Combined k -NN' [11]

k -Nearest Neighbor

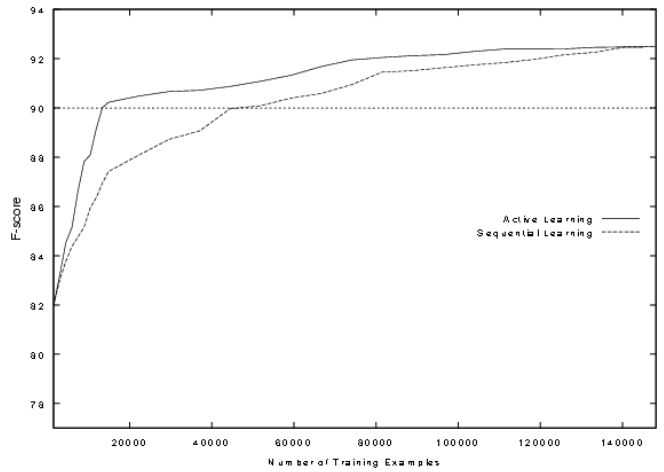
k -NN

k -NN



(a)

4.



(b) F-score

(a)

(b) F-score

가

F-score 90

1/4

		F-score	
C4.5	96.72%	92.48	1,554
	93.03%	90.20	2,663

3.

80,000

92.05

F-score

가

Combined <i>k</i> -NN	95.5 %
	97.8 %
	97.2 %

4.

6.4 AdaBoost

5 CoNLL-2000

AdaBoost

13

F-score

가

boosting

F-score

F-score 가

boosting

5 AdaBoost 가

, 96.88%

92.82 F-

score

3 F-score 0.34

가 0.16%

Koeling

[9]

가 CoNLL-2000

가 [19]

'Sequential Learning'

가

가

, 'Active Learning'

가

X-

, Y-

F-score

가

96%

50,000

,

34%

가

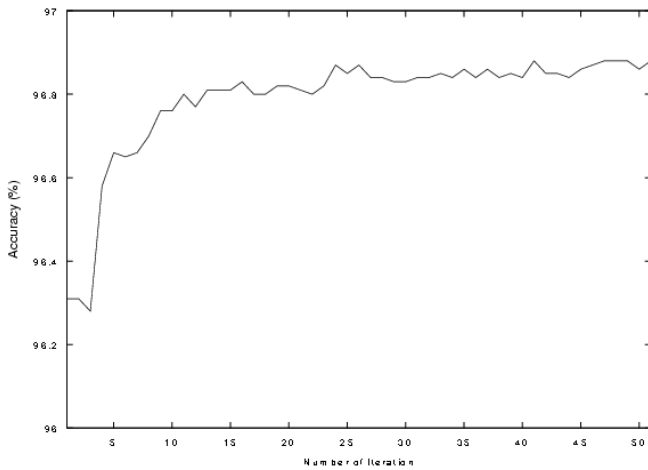
54% 80,000

가

가

96.46%

F-score

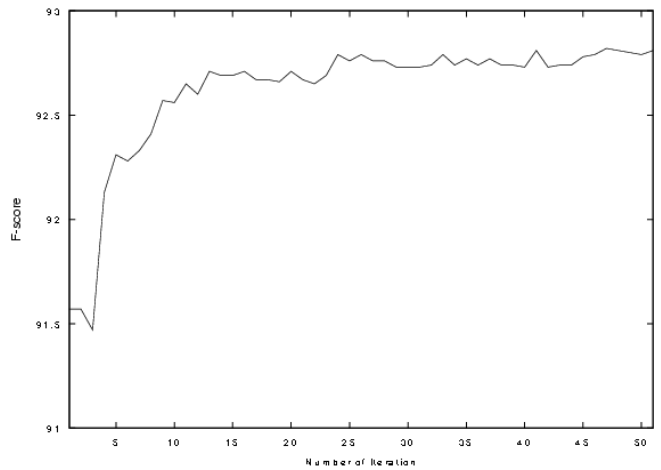


(a)

5.

AdaBoost

(a)



(b) F-score

, (b) F-score

Type	Precision	Recall	F-score
ADJP	62.23%	65.07%	63.62
ADVP	74.48%	78.87%	76.61
CONJP	40.00%	66.67%	50.00
INTJ	100.00%	50.00%	66.67
LST	0.00%	0.00%	0.00
NP	92.49%	94.75%	93.61
PP	96.63%	97.71%	97.17
PRT	73.74%	68.87%	71.22
SBAR	90.89%	87.66%	89.25
VP	92.67%	93.58%	93.12
All	91.96%	93.69%	92.82

5. CoNLL-2000

96.88% . [9] F-score 91.97

7.

97.2%

, CoNLL-2000

, 92.48 F-score

가

가 , 가
 가
 54%
 AdaBoost 가
 AdaBoost
 F-score
 가
 92.82
 가
 (KOSEF)
 (AITrc) BK 21

[1] S. Abney, "Parsing by Chunks," In *Principle-Based Parsing*, Kluwer Academic Publishers, 1991.
 [2] S. Argamon, I. Dagan, and Y. Krymolowski, "A Memory-based Approach to Learning Shallow Natural Language Patterns," In *Proceedings of COLING/ACL 1998*, pp. 67-73, 1998.
 [3] A. Berger, S. Pietra, and V. Pietra, "A Maximum Entropy Approach to Natural Language Processing," *Computational Linguistics*, Vol. 22, No. 1, 1996.
 [4] S. Chen and R. Rosenfeld, "Efficient Sampling and Feature Selection in Whole Sentence Maximum Entropy Language Models," In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, pp. 549-552, 1999.

- [5] CoNLL, *Shared Task for Computational Natural Language Learning (CoNLL)*, <http://lcg-www.uia.ac.be/conll2000/chunking>, 2000.
- [6] J. Darroch and D. Ratcliff, "Generalized Iterative Scaling for Log-linear Models," *The Annals of Mathematical Statistics*, Vol. 43, No. 5, pp. 1470-1480, 1972.
- [7] R. Florian, J. Henderson, and G. Ngai, "Coaxing Confidences from an Old Friend: Probabilistic Classification from Transformation Rule Lists," In *Proceedings of EMNLP/VLC-2000*, pp. 26-34, 2000.
- [8] Y. Freund, S. Seung, E. Shamir, and N. Tishby, "Information, Prediction, and Query by Committee," In *Proceedings of NIPS-92*, pp. 483-490, 1992.
- [9] R. Koeling, "Chunking with Maximum Entropy Models," In *Proceedings of CoNLL-2000 and LLL-2000*, pp. 139-141, 2000.
- [10] G. Ngai and D. Yarowsky, "Rule Writing or Annotation: Cost-efficient Resource Usage for Base Noun Phrase Chunking," In *Proceedings of ACL-2000*, pp. 547-554, 2000.
- [11] S.-B. Park and B.-T. Zhang, "Combining a Rule-based Method and a k -NN for Chunking Korean Text," In *Proceedings of ICCPOL 2001*, pp. 225-230, 2001.
- [12] S. Pietra, V. Pietra, and J. Lafferty, "Inducing Features of Random Fields," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 19, No. 4, pp. 380-393, 1997.
- [13] R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kaufmann Publishers, 1993.
- [14] L. Ramshaw and M. Marcus, "Text Chunking Using Transformation-based Learning," In *Proceedings of VLC-95*, pp. 82-94, 1995.
- [15] S. Seung, M. Opper, and H. Sompolinsky, "Query by Committee," In *Proceedings of COLT-92*, pp. 287-294, 1992.
- [16] B.-T. Zhang, "Accelerated Learning by Active Example Selection," *International Journal of Neural Systems*, Vol. 5, No. 1, pp. 67-75, 1994.
- [17] G. Zhou and J. Su, "Error-driven HMM-based Chunk Tagger with Context-dependent Lexicon," In *Proceedings of EMNLP/VLC-2000*, pp.71-79, 2000.
- [18] J.-M. O and B.-T. Zhang, "Boosting Linear Perceptrons for Unbalanced Data," In *Proceedings of International Conference on Neural Information Processing*, pp. 642-645, 2000.
- [19] T. Zhang, F. Damerau, and D. Johnson, "Text Chunking Using Regularized Winnow," In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, 2001.
- [20] , " , " 11 , pp. 242-247, 1999.
- [21] , " , " 27 , " , pp. 327-329, 2000.