

{sbpark,yhkim,btzhang}@cse.snu.ac.kr

# Automatic Text Classification by Learning from Unlabeled Data

Seong-Bae Park<sup>0</sup> Yu-Hwan Kim Byoung-Tak Zhang

School of Computer Science and Engineering, Seoul National University

(classifier) 가 , 가 , 가 가 , 2/3  
 NIPS 2000 9.2% 가 , WebKB 7.9% ,

## 1.

가 , 가 ,  $y \in \{-1, +1\}$  ,  $x$  가  
 [8]. , -1

Cramér-Rao (unbiased estimator)  $T(\mathbf{x})$  ,  $\theta$  (Fisher information)

$$\text{var}(T) \geq \frac{1}{I(\theta)}$$

$f(\cdot, \theta)$ 가

$$I(\theta) = E_{\theta} \left[ \frac{\partial}{\partial \theta} \ln f(\mathbf{x}; \theta) \right]^2$$

(classifier)가 가

가 , 가  
 Shahshahani Landgrebe

$I_{\text{labeled+unlabeled}}$

$$I_{\text{labeled+unlabeled}} = I_{\text{labeled}} + I_{\text{unlabeled}}$$

[6]. , 가

가 , Zhang Oles (semi-parametric [9])  
 model)  $I_{\text{unlabeled}} = 0$

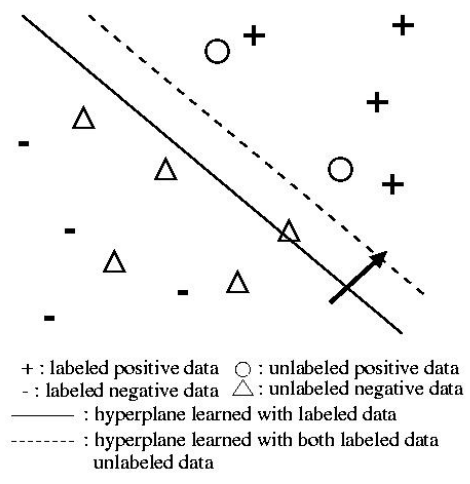
NIPS 2000 가 , WebKB 가 , 가

co-training[2]

1. 가
2. 가

## 2.

가 *bag-of-words* 가 , 1 (hyperplane)



1] 가 (negative example) 가

3. SEQUEL (SEQUENCE Learner)[1]  
 SEQUEL (ensemble)  
 $f_t$  가  $f_t(\mathbf{x})$  가  $\mathbf{x}$  가  $\tau_t$  가  
 $\mathbf{x} : f_t(\mathbf{x}) \geq \tau_t$

SEQUEL SEQUEL 가  $\tau_t$  가 가 가 가 가 가 0.5 가  $\mathbf{x}$  가 가 가 2 SEQUEL  $L$  Zhang Oles 가  $f_0$

$f_t(\mathbf{x})$   $\tau_t$

Given unlabeled example set  $U = \{\mathbf{x}_1, \dots, \mathbf{x}_u\}$   
 and labeled example set  $L = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$   
**Train** a classifier  $f_0$  with  $L_0 = L$ .  
**Set**  $t = 0$  and  $\tau_{-1} = 1$ .

**Train** each base classifier  $C_j (1 \leq j \leq M)$  from  $S_j^{(0)}$ .  
**Do**

1. **Calculate**  $\tau_t = \tau_{t-1} \times \tau$ , where  $\tau$  is the probability given to the negative example in  $L_t$  with the highest probability.
2. **Sort** data in  $L_t$  according to  $f_t(\mathbf{x} \in L_t)$ .
3. **Sort** data in  $U_t$  according to  $f_t(\mathbf{x} \in U_t)$ .
4. **Delete** data in  $L_t$  and  $U_t$  such that  $f_t(\mathbf{x}) > \tau_t$ .
5. **Set**  $s = |L_t|$ .
6. **Set**  $U_{add}$  such that  $U_{add} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-s}, y_{t-s}) \mid \mathbf{x} \in U_t, y = f_t(\mathbf{x})\}$ .
7. **Set**  $L_{t+1} = L_t + U_{add}$
8. **Train**  $f_{t+1}$  with  $L_{t+1}$ .
9. **Set**  $t = t + 1$ .

**While** ( $|U_{add}| > 0$  and  $\tau_t > 0.5$ )  
**Output** the final classifier:

$$f^*(\mathbf{x}) = \left( \prod_{i=1}^{k_x-1} \tau_i \right) f_{k_x}(\mathbf{x})$$

2]  $k_x$   $\mathbf{x}$  가  $f_t$  가  $\tau$   $f_t$   $\tau_t = \tau_{t-1} + 1$   $L_t$   $U_t$  (margin)  $(\mathbf{x}, y)$   $y \cdot f_t(\mathbf{x})$  가  $\mathbf{x} \in U$   $y \cdot f_t(\mathbf{x}) > 0$ .  $\mathbf{x}$   $y$   $f_t$   $L_t$   $U_t$   $\tau_t$

(t+1) 가 t  $f_{t+1}$  가  $\tau_t$  가 0.5  $\mathbf{x}$  가  $f^*(\mathbf{x}) = \left( \prod_{i=1}^{k_x-1} \tau_i \right) f_{k_x}(\mathbf{x})$   $k_x$   $\mathbf{x}$  가

4. NIPS 2000 "Using Unlabeled Data for Supervised Learning"  
 가 (P2)  
<http://www.microsoft.com> <http://www.linux.org>  
 (P6)  
<http://www.mit.edu> <http://www.uoguelph.ca>  
 'MIT', 'Institute',

'Guelph'

1 words bag-of-

*tf•idf*

[ 1] NIPS 2000

Data Set	P2	P6
No. of Labeled Data	500	50
No. of Unlabeled Data	5,481	3,952
No. of Test Data	1,000	100
No. of Terms	200	1,000

co-training CMU text learning group "The 4 Universities Data Set" Cornell, University of Washington, University of Wisconsin, University of Texas 1,051 1,051 course non-course 2 non-course

[ 2] WebKB

Data Set	Course	Non-Course	Baseline
Cornell	40	203	83.5%
Texas	38	216	85.0%
Washington	74	220	71.1%
Wisconsin	78	220	73.8%
Total	230	821	78.1%

5.

MLP, Pentium III 550MHz 256MB 가 PC Linux 3 NIPS 2000

[4] Transductive SVM 가

TSVM

0.7%, P6 15.0%

[ 3] NIPS 2000

Method	Using Only Labeled Data		Using Both Labeled and Unlabeled Data	
	P2	P6	P2	P6
Our Method	98.8%	60.0%	99.5%	75.0%
TSVM	N/A	N/A	99.7%	80.0%

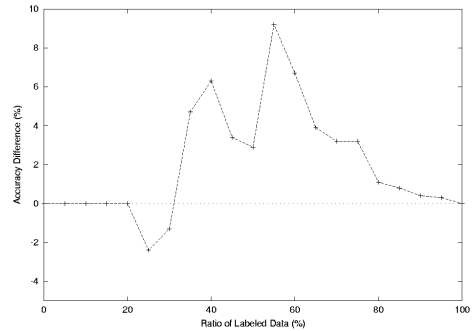
WebKB

1.5% ( 4). 16.2% , co-training 0.4% 3 가

9.2%

[ 4] WebKB

Data Set	Using Partially Labeled Data	Using All Labeled Data	Co-Training
Cornell	94.2%	93.4%	N/A
Texas	97.1%	96.5%	N/A
Washington	91.4%	89.9%	N/A
Wisconsin	94.2%	91.3%	N/A
Average	94.2%	92.8%	93.8%



[ 3]

6.

가 가

BK 21

( : 00-023)

[1] L. Asker and R. Maclin, "Ensembles as a Sequence of Classifiers," In *Proceedings of IJCAI-99*, pp. 860-865, 1999.

[2] A. Blum and T. Mitchell, "Combining Labeled and Unlabeled Data with Co-Training," In *Proceedings of COLT-98*, pp. 209-214, 1998.

[3] T. Joachims, "Transductive Inference for Text Classification Using Support Vector Machines," In *Proceedings of ICML-99*, pp. 200-209, 1999.

[4] B. Shahshahani and D. Landgrebe, "The Effect of Unlabeled Samples in Reducing the Small Sample Size Problem and Mitigating the Hughes Phenomenon," *IEEE Transactions on Geoscience and Remote Sensing*, Vol. 32, pp. 1087-1095, 2000.

[5] Y. Yang and J. Pederson, "Feature Selection in Statistical Learning of Text Categorization," In *Proceedings of ICML-97*, pp. 412-420, 1997.

[6] T. Zhang and F. Oles, "A Probability Analysis on the Value of Unlabeled Data for Classification," In *Proceedings of ICML-2000*, pp. 1191-1198, 2000.