

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

Text Classification By Boosting Na ve Bayes

Yu-Hwan Kim Byoung-Tak Zhang
School of Computer Science and Engineering, Seoul National University

AdaBoost
가
가
(confidence ratio)
(filtering track)
TREC-7 TREC-8

1. AdaBoost
C4.5 decision stump tf 가
[4] [6].

가 가 가 BayesBoost
가 TREC 가 가 가
[8]. TREC-7 TREC-8
가 가 [6].

2.

가 가 가 Rocchio
가 TREC (batch filtering) Rocchio
track) 가 dynamic feedback
pivot document
AdaBoost 가 optimization, query zoning, pivoted document
(weak learner)가 normalization (support vector
AdaBoost 가 k- machine, SVM), 가
가 (weighted voting)

Schapire
REUTER-21578 TREC-3

AdaBoost

Rocchio
, LF1

3.3 BayesBoost

3. BayesBoost

3.1 AdaBoost

AdaBoost 가

Schapire 가 $\{-1, 1\}$ 가 h_t 가 [1].

가 α_t 가 $\sum_{i=1}^m \alpha_t h_t$ 가 m 가

h_t 가 $|h_t|$ 가 h_t 가 C4.5
decision stump 가 [5].
[4], 가

tf 가

3.2

[3].
가
가 w_k 가 $(\theta_{w_k|c_k})$,
 c_k 가 (θ_{c_k})

$$\theta_{w_k|c_j} = \frac{1 + \sum_{i=1}^D N(w_k, d_i) P(c_j | d_i)}{|V| + \sum_{s=1}^{|V|} \sum_{i=1}^{|D|} N(w_k, d_i) P(c_j | d_i)}$$

$$\theta_{c_j} = \frac{\sum_{i=1}^{|D|} P(c_j | d_i)}{|D|}$$

$N(w_k, d_i)$ d_i w_k 가
 $|V|$, $|D|$

$$P(c_j | d_i; \hat{\theta}) = \frac{P(c_j | \hat{\theta}) P(d_i | c_j; \hat{\theta})}{P(d_i | \hat{\theta})} = \frac{\hat{\theta}_{w_i|c_j} \hat{\theta}_{c_j}}{P(d_i | \hat{\theta})}$$

$\arg \max_j P(c_j | d_i; \hat{\theta})$ 가

Decision Stump 가

BayesBoost

$$h_t(d_i; \hat{\theta}) = \tanh \left\{ \log \left(\frac{P(c_i = 1 | d_i; \hat{\theta})}{P(c_i = -1 | d_i; \hat{\theta})} \right) \right\}$$

가 α $[-1, 1]$
 α

$$\alpha = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right), r = \sum_i D(i) y_i h_t(x_i)$$

4.

4.1

가
precision/recall,
break-even point TREC linear
utility
[2]. linear
utility Linear
utility

$$\text{Linear Utility} = aR_+ + bN_+ + cR_- + dN_-$$

R_+ (relevant)

N_+
 R_-

N_-

a, b, c, d

TREC

LF1

$$\text{LF1} = 3R_+ - 2N_-$$

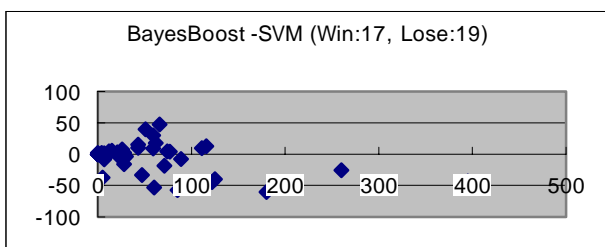
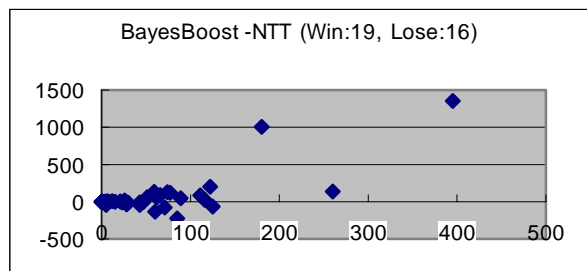
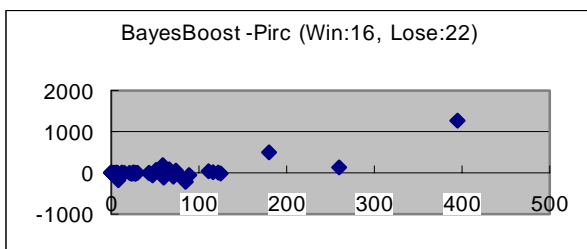
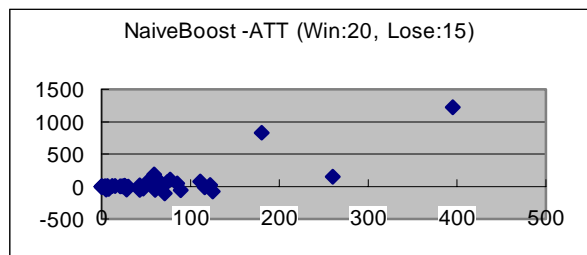
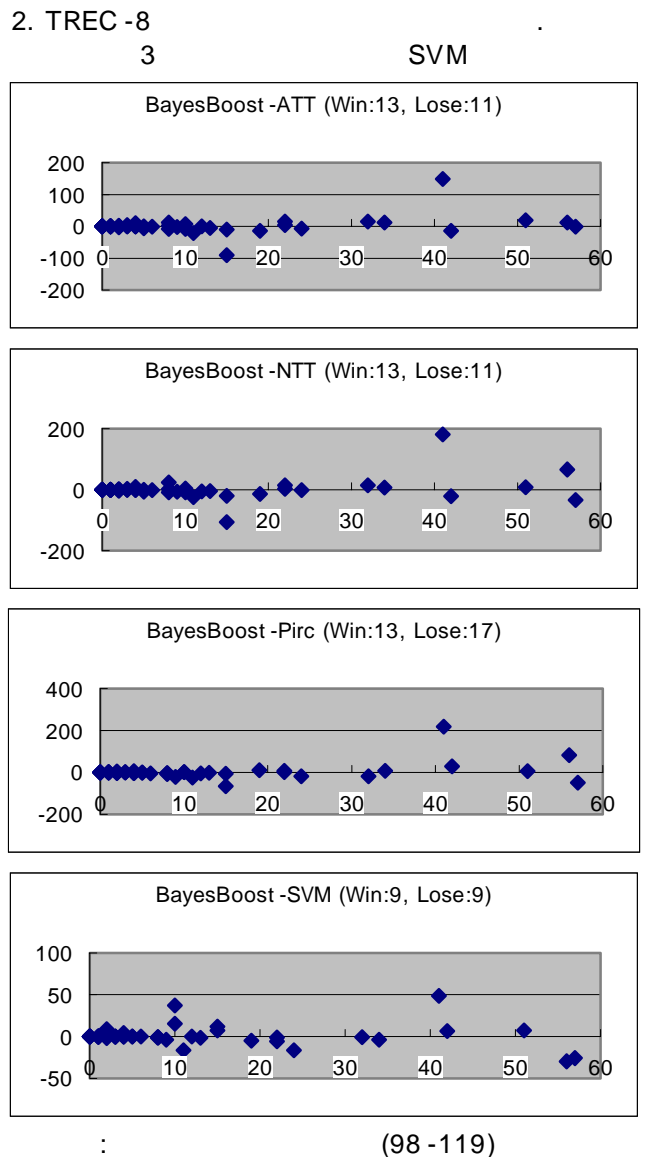
4.2 TREC-7 TREC-8

TREC-7 1988-1990 AP
1-50
1988
TREC-8 1992-
1994 351-400
, 1992 , 1993-1994
'description'
TREC-7 5250 TREC-8 4079
. 100

가

TREC-7 TREC-8 가
 SVM
 TREC-7 1, TREC-8
 2 x
 positive example y BayesBoost

5. BayesBoost
 가 positive example
 가
 pirc
 TREC positive
 example 가
 AP 가
 2 가
 SVM 가
 SVM 가
 가
 가
 1. TREC -7
 3 SVM



[1] Y. Freund and R. E. Schapire, "Experiments with a new boosting algorithm", In *Proc. 13th Int. Conf. On Machine Learning*, pp. 148 -156, 1996

[2] D. Hull. "The TREC-8 filtering track: Description and analysis", In *Proc. 7th Text Retrieval Conf. (TREC -7)*, pp. 33 -56, 1998.

[3] T. Joachims, "A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization", In *Proc. Int. Conf. on Machine Learning (ICML -97)*, pp. 143 -151, 1997.

[4] J. R. Quinlan, "bagging, boosting and C4.5", In *Proc. AAAI-96*, pp. 725 -730, 1996.

[5] R. E. Schapire and Y. Singer, "Improved boosting algorithms using confidence-rated predictions", *Machine Learning*, 37(3):297 -336, 1999.

[6] R. E. Schapire, Y. Singer, and A. Singhal, "Boosting and Rocchio applied to text filtering", In *Proc. SIGIR-98*, pp. 251 -223, 1998.

[7] D. K. Harman, "Overview of the 8th Text Retrieval Conference (TREC-8)", In *Proc. 8th Text Retrieval Conf. (TREC -8)*, pp 1 -10, 1999.