

Co-Trained Support Vector Machines

o

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

Text Categorization Using Co-Trained Support Vector Machines

Seong-Bae Park^o Byoung-Tak Zhang

School of Computer Science and Engineering, Seoul National University

co-training

가

Reuters-21578

TREC-7 filtering

1.

(text categorization)

SVM

confidence margin

2. Co-Training

2.1

가

(machine learning)
bag of words

[1,2].

co-training

가

가

가

가

가

(i)

(ii)

가

(formal definition)가

Stamatatos
가

가

가

(full parsing)

[5].

co-training

가

(text chunk)[3]

, *text chunk* 가

가

Co-training

[4]

●

tf-idf

●

co-trained Support Vector

가

Machine

2.2

Co-training

가

가

(chunk)

가

CoNLL-2000

SVM
CoNLL-2000 12 가 ,
가 : NP, VP, ADVP, PP O.
O NP, VP, ADVP, PP 가 O
(B-XP I-XP)
가 11 가 .
SVM
[3]. w_i
 w_j, POS_j ($j = i - 2, i - 1, i, i + 1, i + 2$)
 c_j ($j = i - 2, i - 1$)
, POS c
SVM 11
가 , (pairwise classification)[6]

1
5
5

SF1	Detected NPs / total detected chunks
SF2	Detected VPs / total detected chunks
SF3	Detected PPs / total detected chunks
SF4	Detected ADVPs / total detected chunks
SF5	Detected Os / total detected chunks
SF6	Words included in NPs / detected NPs
SF7	Words included in VPs / detected VPs
SF8	Words included in PPs / detected PPs
SF9	Words included in ADVPs / detected ADVPs
SF10	Words included in Os / detected Os
SF11	Sentences / words

1. (feature).

3.

3.1

Reuters-21578

Reuters-21578 135 가 ,
10
“ModLewis”, “ModApte”,
“ModHayes”가 가
“ModApte” , 9,603
3,299

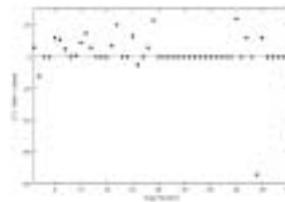
TREC-7

TREC-7 filtering track 1988
1990 AP . 1988
, 1989 1990
79,898 가
12% 9572
9572
88%

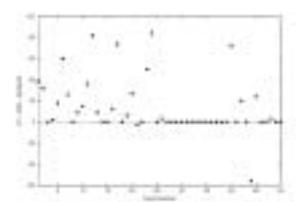
Class			
Earn	2876	2504	2895
Acq	1587	1028	1642
Money-fx	188	147	193
Grain	26	14	26
Crude	292	150	312
Trade	153	114	155
Interest	130	112	141
Ship	113	88	116
Wheat	98	70	108
Corn	86	68	87

2.

. LF1



(a) LF1 : Both - Lexical



(b) LF1 : Both - Syntactic

0.

3.2 가

가 ,
utility measure . R_+ relevant
, N_+ irrelevant , R_- relevant
, N_- irrelevant
, linear utility

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

, a, b, c, d . LF1 LF2

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

$$\text{LF2} = 3R_+ - N_+$$

Linear utility

가

scaled linear utility

$$\text{Scaled Linear Utility} = \frac{\max\{u(S, T), U(s)\} - U(s)}{\text{Max}U(T) - U(s)}$$

, $u(S, T)$ S T linear utility
, $\text{Max}U(T)$ T utility
 $U(s)$ s irrelevant utility

3.3

Reuters-21578

TREC-7

1

, Y

. X

