

* * * **
*
,
{chjeong, cokeman, joo}@kisti.re.kr,
myaeng@icu.ac.kr

A Probabilistic Method for Recognizing Unlabeled Text on Web Pages

Chang-Hoo Jeong*, Min-Ho Lee*, Won-Kyun Joo*, Sung-Hyon Myaeng**
Korea Institute of Science and Technology Information*,
Information and Communications University**

HTML

가

1.

HTML

가

가

가

[1].

가

가

가

(가),
_____ (Token Set)

2.

Sequence)

_____ (Token Set

[1]

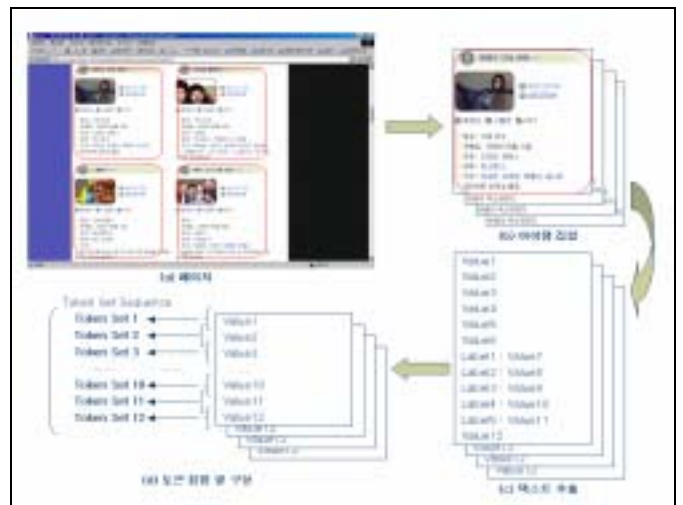
가

가

가

()

가



[1]

[1] (a) HTML

, 가

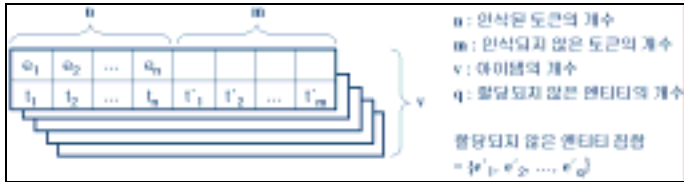
(b) 가

(c)

(b)

(d)

[1]
[2]



[2]

[2]

m

가가

1. n

{t₁, t₂, ..., t_n}

2. n 가 {e₁, e₂, ..., e_n}

3. m

{t₁, t₂, ..., t_m}

4. q 가 {e₁, e₂, ..., e_q}

e_k

E

5. v

6. n

{T₁, T₂, ..., T_n}, T_i = {t_{i1}, t_{i2}, ..., t_{iv}}

7. m

{T₁, T₂, ..., T_m}, T_j = {t_{j1}, t_{j2}, ..., t_{jv}}

8. (n + q) 가

2.

ERM(Entity Recognition Model)

HMM(Hidden Markov Model)

HMM

(Word)

(Category)

[2].

ERM

(Token)

(Entity)

HMM

HMM

HMM

, ERM

Viterbi Algorithm

가

가

ERM

HMM

, ERM

HMM

[1]

	HMM	ERM
Class	Category	Entity
Object	Word	Token
Probability	Word가 Category	Token Entity
Corpus Statistics	Lexical Generation Probability	Model 1 Probability
Context Information	Bigram Probability	Model 2 Probability
Difference	Category 가	Entity 가

[1] HMM ERM

HMM

$$= \text{PROB}(C_1, \dots, C_T | W_1, \dots, W_T)$$

$$\cong \prod_{i=1, T} \text{PROB}(W_i | C_i) * \text{PROB}(C_i | C_{i-1})$$

ERM

$$= \text{PROB}(e'_1, \dots, e'_q | T'_1, \dots, T'_m)$$

$$\cong \alpha * \{P(e'_i) * \frac{1}{v} \sum_{k=1}^v P(t'_{jk} | e'_i)\} +$$

$$(1 - \alpha) * \{ \frac{1}{v} \sum_{k=1}^v \sum_{h=1}^n P(e'_i = t'_{jk} | e_h = t_{hk}) * P(e_h = t_{hk}) \}$$

(, 1 <= i <= q and 1 <= j <= m)

HMM

Word가 Category가

Lexical Generation Probability: $\text{PROB}(W_i | C_i)$

ERM

Model 1 Probability:

$$P(e'_i) * \frac{1}{v} \sum_{k=1}^v P(t'_{jk} | e'_i)$$

$$1. P(t_{jk} | e_i)$$

e_i t_{jk} /

$$. P(t_{jk} | e_i) =$$

$$2. P(e_i)$$

e_i가

/

가

HMM

Category가

Bigram Probability: $PROB(C_i | C_{i-1})$

ERM

Model 2 Probability:

$$\frac{1}{V} \sum_{k=1}^V \sum_{h=1}^n P(e'_i = t'_{jk} | e_h = t_{hk}) * P(e_h = t_{hk})$$

$$3. P(e_i = t_{jk} | e_h = t_{hk}) = \frac{P(e_i = t_{jk}, e_h = t_{hk})}{P(e_h = t_{hk})}$$

$$4. P(e_h = t_{hk}) = \frac{P(e_h = t_{hk})}{P(e_h = t_{hk})}$$

5. $PROB(e_1, \dots, e_q | T_1, \dots, T_m)$ 가 가

가

6. T_1, T_2, \dots, T_{m-1} , T_j
1, 2, 3, 4, 5

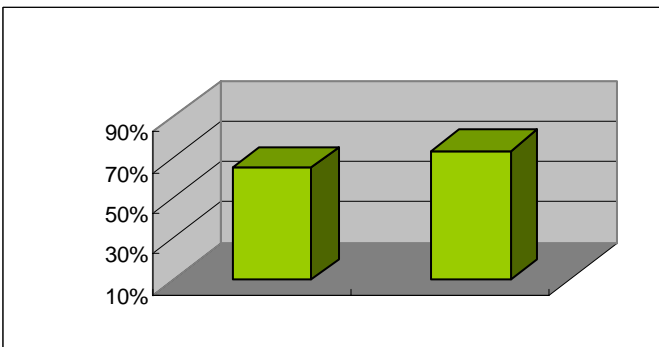
가 , e_i , $P(e_i)$
가 , e_i
ERM HMM
가 α
가 α

4.

7 (Site A, Site B, ..., Site G)

가

가



[3]

사이트 구분	레이블 개수	Model 1	확률 비교	Model 2
Site D	2	0.0007380	>	0.0004231
Site E	5	0.0004548	<	0.0010652
Site G	2	0.0014716	>	0.0006424

[4] Model 1 Model 2

[3]

가 가

가

Model 1 2 [4] 가

7
Model 1

Model 2

E

Site

가

([4]

가 5

가 , Site E Site G

가 2

Model 1

Model 2

가

Model 2

Model

, Model 1 Model 2

가

α

Model

5.

가 가

6.

[1] H. Seo, J. Yang, and J. Choi, "Knowledge-based Wrapper Generation by Using XML", IJCAI-2001 Workshop on Adaptive Text Extraction and Mining (ATEM 2001), pp. 1-8, Seattle, USA, 2001.

[2] James Allen, "Natural Language Understanding (2nd Edition)", Addison-Wesley Publishing Co, pp. 189-204, 1995