# Information Extraction Using Hidden Markov Models

2001    2

# Information Extraction Using Hidden Markov Models

2000    10

2000    12

_____    _____    \_\_

_____    _____    \_\_

\_\_  \_\_    _____    \_\_

(Hidden Markov Model, HMM)

(automata)

.

.

.

left-to-right HMM

. , HMM HMM

(Self-Organizing

Hidden Markov Model: S-HMM) .

S-HMM

Call-For-Papers , CMU

, LA (http://www.laweekly.com/)

.

S-HMM

.

12% 4%

.

: , , ,

: 99419-524

- iii -

# 1.

## 1.1

,

.

(information overload)

[Maes, 1994].

. ,

.

,

(filtering) ,

.

.

,

(grammar rule), (machine learning),

.

.              ,

HMM                    [David, 1999]                    .    ,

HMM

,        ,

(learner)

regression

(multistrategy)                                                      .

.

(Call  for  Papers)

,

.              ,

.

,

.                                        ,

.

**1.2**

(speech recognition)

.

.                    ,

,

HMM                        ,

Viterbi

[Rabiner, 1989].


HMM

Leek                              . Leek      HMM

[Leek, 1997]. Leek      HMM


(language modeling)

.          Leek                (syntax)

(network topology)                  ,                          (unigram)

.                              (token)                  (gap state)

. Leek                                                      HMM

.


Bikel

.  Bikel        HMM                    MUC-6(Message

Understanding Conference)                          (name entity)

Nymble                  [Bikel et al., 1997]. Nymble

Seymore[Seymore et al., 1999]

HMM　　　　　.　　,

(fully
connected)　　　　　　.　　,

(sparse training data)

(emission probabilities)　　　　(transition probabilities)　　　.

(shrinkage)

.　　,

, EM

(optimal mixture weight)　　　　　　　　.

HMM　　　　　　　. David[1999]

,

HMM　　　　　.　　　　　HMM

. David

TREC(Text REtrieval Conference)　　TREC-6, TREC-7 ad hoc retrieval

$tf \cdot idf$

[David et al., 1999].　　　　　HMM　　　　　　(blind
feedback)　　　, bigram

.

HMM

(profile)　　　　　　　. Lane[Lane 1990]　　HMM

(anormality)　　　　　. Lane

HMM

.

. Lane

，

．

HMM

．Seymore　[Semore et al., 1999]　HMM

．

Seymore

HMM　　．

，　　　　HMM　　　，

HMM　　．

，

，　　．

Freitag[Freitag, 1998]　HMM

．

(regression)

(multistrategy)

，　　．

．Riloff[Riloff,

1999]　(dictionary)　　(semantic　lexicon)

(extraction pattern)　　，

(Multi-level mutual bootstrapping)

(seed word)

，　．

．

**1.3**

(Self-Organizing Hidden Markov Model)

. EM (Expectation-Maximizati-

on) [Dempster et al., 1997] .

HMM

HMM

.

. ,

.

1 CFP

. 1 , , ,

URL ,

,

. 1

,

,

. ,

.

```
              SECOND CALL FOR PAPERS

       FIFTH ANNUAL INTERNATIONAL CONFERENCE ON
            COMPUTATIONAL MOLECULAR BIOLOGY

                   (RECOMB 2001)

                  April 21-24, 2001
                   Montrel, Canada

                    Organized by
           Centre de recherches math?atiques
                Universite de Montreal

                     Sponsored by
       Association for Computing Machinery (ACM-SIGACT)

                  with support from
                   Celera Genomics
                      Compugen
                   IBM Corporation
       International Society for Computational Biology (ISCB)
                   SLOAN Foundation
                  SmithKline Beecham
                 US Department of Energy
                 US National Science Foundation

                http://recomb2001.gmd.de


  The Fifth Annual Conference on Research in Computational Molecular
  Biology (RECOMB 2001), sponsored by the Association for Computing
  Machinery Special Interest Group on Algorithms and Computation Theory
  (ACM-SIGACT) with support from Celera Genomics, Compugen, IBM
  Corporation, SLOAN Foundation, International Society for Computational
```

1. RECOMB 2001 　　　　　CFP

　　　　　　　　　HMM

　　　　　　　　　　　　　　　　　　　．

　　　　　　　　　　．　　　，

HMM　　　　　　　　　HMM

　　　　　　　　　　　　　　　　　　．

**1.4**

. 2

(Forward Algorithm)           -                    (Forward-Backward Algorithm)

,                           EM              Viterbi

. 3

CFP(Call-For-Papers)

, 4

.                    5

.

## 2. (Hidden Markov Model)

### 2.1

HMM (hidden) (stochastic process)
(symbol)
(modeling) .
,
. HMM 2
3 5 .

- (hidden state set)

- (observable state set)

- $\pi$

$t = 1$

- 

-

, HMM

. HMM

.

$$P(q_t = j \mid q_{t-1} = i, q_{t-2} = k, \ldots) = P(q_t = j \mid q_{t-1} = i) \tag{1}$$

$$P(q_t = j \mid q_{t-1} = i) = P(q_{t+l} = j \mid q_{t+l-1} = i) \tag{2}$$

(1) ,

. , $t$ $q_t$ $j$ $t-1$

$q_{t-1}$ . **1** (first order

Markov assumption) HMM (2)

$t$ .

(observation symbol

sequence) (3) .

$$O = O_1, O_2, \ldots, O_{T-1}, O_T \tag{3}$$

$T$ . $N$

, $M$ .

$$Q = \{q_1, q_{2,} \ldots, q_n\} \tag{4}$$

(4)                                                                .

$$V = \{ v_1, v_2, \ldots, v_m \} \tag{5}$$

(5)                                                                .

( $\lambda$ )                    ( $\Pi, A, B$ )    3                    .

· $\Pi = (\pi_i)$, $\pi_i = P(q_1 = i)$,    $1 \leq i \leq N$

· $A = (a_{ij})$, $a_{ij} = P(q_t = j | q_{t-1} = i)$,    $1 \leq i, j \leq N$

· $B = (b_j(k))$, $b_j(k) = P(o_t = v_k | q_t = j)$,    $1 \leq k \leq M$, $1 \leq j \leq N$

$$P(o_i(k) | q_j)$$                        .

$A$                            .

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1j} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2j} & \cdots & a_{2N} \\ \vdots & \vdots & \cdots & \vdots & \cdots & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{ij} & \cdots & a_{iN} \\ \vdots & \vdots & \cdots & \vdots & \cdots & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{Nj} & \cdots & a_{NN} \end{pmatrix} \tag{6}$$

(6)                        (7)                , $a_{ij}$        (8), (9)

.

$$a_{ij} = P(q_t = j \mid q_{t-1} = i) \qquad 1 \leq i,j \leq N \tag{7}$$

$$a_{ij} \geq 0, \qquad \forall i,j \tag{8}$$

$$\sum_{j=1}^{N} a_{ij} = 1, \qquad \forall i \tag{9}$$

$A$ ，$B$                                                         ，

．                                 ，

HMM

．

HMM

3                                                         ．

**1.        (probability estimation)**

$O = (o_1, o_2, \ldots, o_T)$                $\lambda = (\Pi, A, B)$

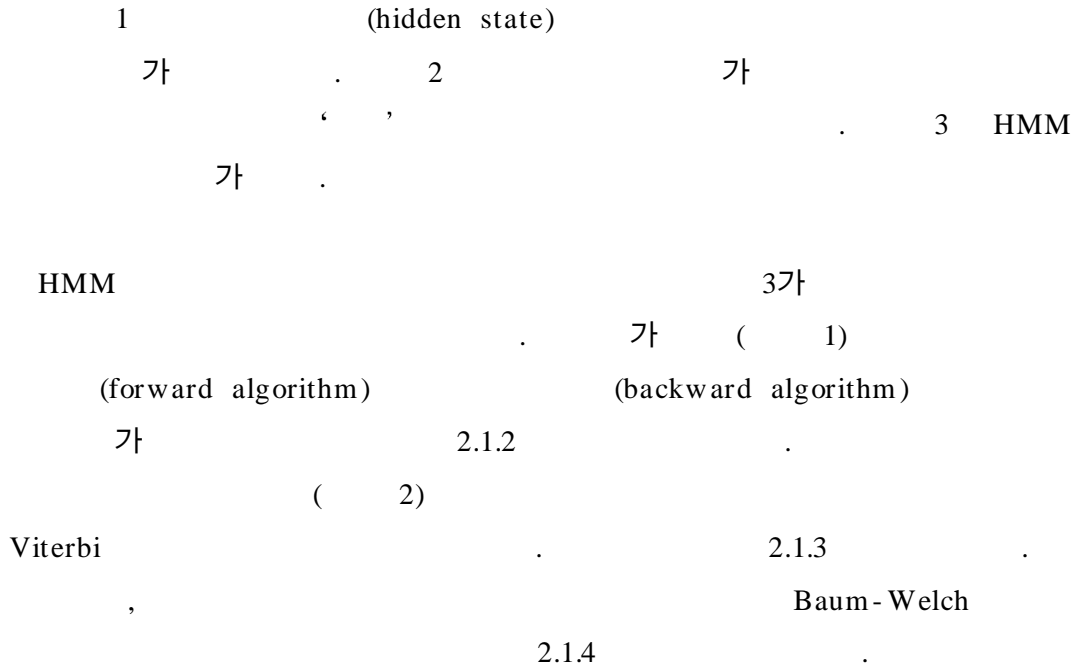HMM                                 $P(O \mid \lambda)$                                 ．

**2.                        (optimal sequence)**

$O = (o_1, o_2, \ldots, o_T)$                $\lambda$

$q = (q_1, q_2, \ldots, q_T)$

．

**3.                        (parameter estimation)**

$O = (o_1, o_2, \ldots, o_T)$                $P(O \mid \lambda)$

$\lambda = (\Pi, A, B)$                (parameter)

．

$$\max{}_{\lambda}\{P(O\,|\lambda)\} \tag{10}$$

1 (hidden state)

. 2

'   '                                                                    . 3    HMM

.

HMM                                                                3

.                                    (      1)

(forward algorithm)                        (backward algorithm)

2.1.2                            .

(      2)

Viterbi                                    .                    2.1.3                        .

,                                                        Baum‑Welch

2.1.4                                .

## 2.2

**(Forward Algorithm)**

$\lambda = (\Pi, A, B)$

$O = (o_1, o_2, \ldots, o_T)$ $\qquad P(o_1, o_2, \ldots, o_T \mid \lambda)$

.

.

---

Algorithm

- Let $q = (q_1, q_2, \ldots, q_T)$ be a state sequence.

- Assume the observations are independent:

$$P(O \mid q, \lambda) = \prod_{t=1}^{T} P(o_t \mid q_t, \lambda)$$

$$= b_{q1}(o_1) b_{q2}(o_2) \cdots b_{qT}(o_T)$$

- Probability of a particular state sequence is:

$$P(q \mid \lambda) = \pi_{q1} a_{q1q2} a_{q2q3} \cdots a_{qT-1qT}$$

- Also, $P(O, q \mid \lambda) = P(O \mid q, \lambda) P(q \mid \lambda)$

- Enumerate paths and sum probabilities:

$$P(O \mid \lambda) = \sum_{q} P(O \mid q, \lambda) P(q \mid \lambda)$$

---

1.

1 $\qquad N^T$

$O(T)$ $\qquad$ ,

$O(\,TN^{\,T})$ .

$T$ , $T$

.



2. (Forward Procedure)

.

.

(forward variable) $\alpha_t(\,i\,)$ .

$$\alpha_t(\,i\,) = P(\,o_1, o_2, \cdots , o_t, q_t = i\ |\lambda) \qquad (11)$$

(11) $\alpha_t(\,i\,)$ $q_t$ $i$

$(\,o_1, o_2, \cdots , o_t)$ . (2)

.

$O(N^2 T)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $T$

.

- Induction

  1. Initialization:

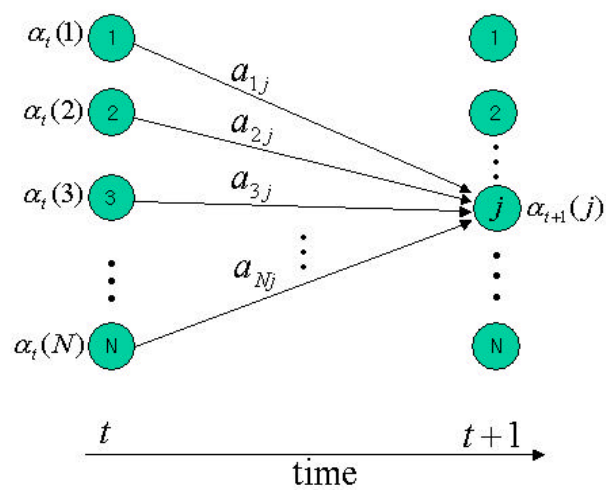     $$\alpha_1(i) = \pi_i b_i(o_1)$$

  2. Induction:

     $$\alpha_{t+1}(j) = \left[ \sum_{i=1}^{N} \alpha_t(i) a_{ij} \right] b_j(o_{t+1})$$

  3. Termination:

     $$P(O|\lambda) = \sum_{i=1}^{N} \alpha_T(i)$$

2.



3.

$t= 1$ $t= T$

$\alpha$ .

## (Backward Algorithm)

(backward

variable) . $\beta_t( i)$ .

$$\beta_t( i) = P( o_{t+ 1}, o_{t+ 2}, \cdots , o_T \mid q_t = i, \lambda) \tag{12}$$

(12) $\beta_t( i)$ $q_t$ $i$

$( o_{t+ 1}, o_{t+ 2}, \cdots , o_T)$ . (3)

.

- Induction
  1. Initialization:
     $$\beta_T( i) = 1$$
  2. Induction:
     $$\beta_t( i) = \sum_{=1}^{N} a_{ij} b_j( o_{t+ 1}) \beta_{t+ 1}(j)$$
     $1 \leq i \leq N$,
     $t= T - 1, \cdots , 1$

   3.

$\beta$ $t= T$

$t= 1$ .

图 4.

## 2.3        (state sequence)

，

.

(dynamic programming)

Viterbi        .

### (Viterbi algorithm)

Viterbi          $O$    $\lambda$        $O$

(state sequence) $(q_1, \cdots, q_t, \cdots, q_T)$

.     $O$     $\lambda$              (13)

.

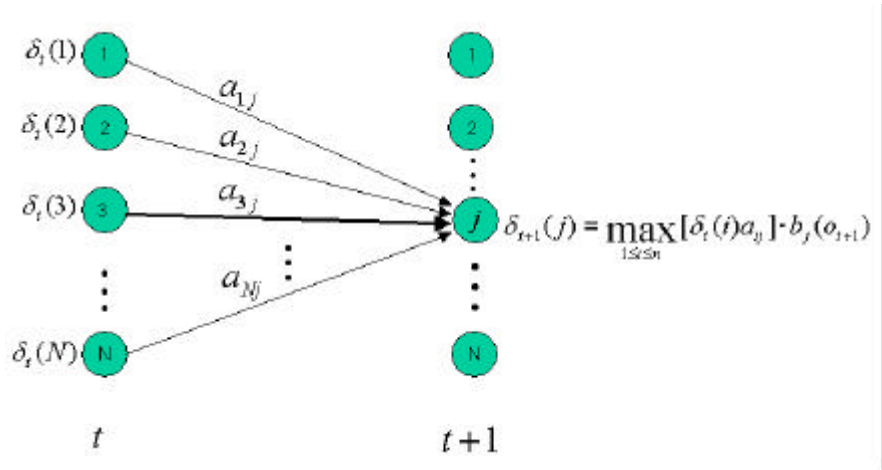$$P(q_1, q_2, \cdots, q_T \mid O, \lambda) \tag{13}$$

5           $t$                      $t+1$

.



5. Viterbi

5         $\delta_t(j)$           $j$

(14)                              .

$$\delta_t(i) = \max_{q_1, q_2, \cdots, q_{t-1}} P(q_1, q_2, \cdots, q_t = i, o_1, o_2, \cdots, o_t \mid \lambda) \tag{14}$$

$$\delta_{t+1}(j) = \max_i [\delta_t(i) a_{ij}] \cdot b_j(o_{t+1}) \tag{15}$$

(14)    (15)                              ,

$t$                                        $t+1$

.

---

- Initialization

$$\delta_1(i) = \pi_i b_i(o_1), \qquad 1 \leq i \leq N$$

$$\psi_1(i) = 0$$

- Recursion

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t)$$

$$\psi_t(j) = \arg\max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}]$$

$$2 \leq t \leq T, \ 1 \leq j \leq N$$

- Termination

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$$

$$q^*_T = \arg\max_{1 \leq i \leq N} [\delta_T(i)]$$

- Path (state sequence) backtracking

$$q^*_t = \psi_{t+1}(q^*_{t+1}), \qquad t = T-1, T-2, \cdots, 1$$

---

4. Viterbi

4        $\psi_t(i)$          $t$                    $i$

.  $\psi_t(i)$      $\psi_t(j) = \arg\max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}]$                        ($t-1$)

$\delta_{t-1}$                $t$                        $j$

.

## 2.4

,

$O = (o_1, o_2, \ldots, o_T)$                     $P(O \mid \lambda)$                     $\lambda = (\Pi, A, B)$

(parameter)                     (     3).

(analytic)                     .

EM

Baum - Welch                     .

## EM (Baum - Welch)

Baum - Welch                     $(\lambda_0)$                     ,                     $O$

$(\lambda)$                     .

.

5          .

---

Step 1. Let initial model be $\lambda_0$.

Step 2. Compute new $\lambda$ based on $\lambda_0$ and observation

$O$.

Step 3. If $\log P(O \mid \lambda) - \log P(O \mid \lambda_0) < DELTA$ then stop.

Step 4. Else set $\lambda_0$ $\lambda$ and goto step 2.

---

5. EM (Baum - Welch)

, Baum-Welch

. (16) $t$ $i$ , $t+1$

$j$ . , (17) $t$

$i$ .

$$\xi(i,j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(O \mid \lambda)} \qquad (16)$$

$$= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}$$

$$\gamma_t(i) = \sum_{j=1}^{N} \xi_t(i,j) \qquad (17)$$

(16) $\sum_{t=1}^{T-1} \xi(i,j)$ $O$ $i$ $j$

. , (17) $\sum_{t=1}^{T} \gamma_t(i)$

$O$ $i$ .

(16) $\alpha$ $\beta$ , - (forward-backward)

6

.

EM

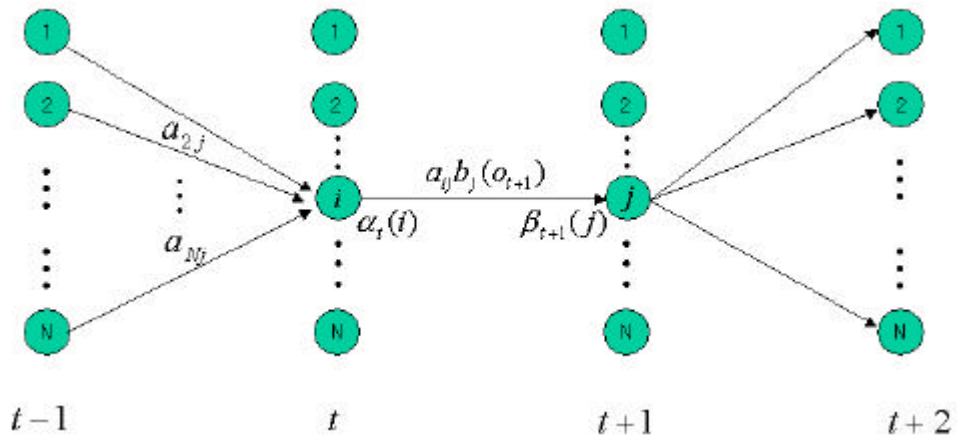(reestimate) , (18) (20) .

$$\widehat{\pi_i} = \gamma_1(i) \qquad (18)$$

$$\widehat{a_{ij}} = \frac{\sum \xi_t(i,j)}{\sum \gamma_t(i)} \qquad (19)$$

$$\widehat{b}_j(k) = \frac{\sum\limits_{t, o_t = k} \gamma_t(j)}{\sum\limits_t \gamma_t(j)} \tag{20}$$

,     (18)                 $t = 1$       ,         $i$

.     (19)                $i$                              ,

        $i$            $j$                            .     (20)

        $j$                       ,            $k$

$j$                             .



6. EM                 ,

# 3.

## 3.1

, HMM

(mapping) . ,

. HMM

,

.

### 3.1.1 HMM

, HMM HMM

.

, HMM $O$ ,

. ,

.

6 9 .

3 .

, ,

$+ 1$

. ,

$$Number\ of\ Initial\ Model\ States\ =\ 2\times (Number\ of\ Target\ Fields)\ +\ 1 \qquad (21)$$

(21)

, (21)

.

6 9 3

.

.

| | | |
|---|---|---|
| | {n-NAME_SYM + year}, <br> {n-NAME_SYM + "'" + year}, <br> {n-NAME_SYM + "-" + year}, <br> {n-NAME_SYM + " " + year}, etc. <br> NAME_SYM := Capital character <br>         \| Begin with Capital Character <br> year := 4-digit \| 2-digit <br> $2 \leq n \leq 15$ | ICONIP 2000 <br> SIGIR '2001 <br> PRICAI'01 <br> KDD01 <br> EC-Web 2000 <br> CL2000 <br> GECCO 2000 <br> CEC 2001 |
| | {day + month + "," + year}, <br> {day + "," + month + "," + year}, <br> {day + th + D_SYM + day + "th + month + year}, <br> {day + D_SYM + day + "," + year}, <br> {month + day + D_SYM + day + "," + year}, etc. <br> D_SYM := "-" \| "to" | 23 September 2000 <br> 24th to 28th July, 2000 <br> 8-12 July 1999 <br> August 20-23, 2000 |

6. CFP (Call-For-Papers)

| | | |
|---|---|---|
| | {PLACE + \<NL\> + CITY + "," + COUNTRY}, <br> {PLACE + "," + CITY + "," + COUNTRY}, etc. <br> NL := new line character | Riviera Hotel <br> Las Vegas, Nevada USA |
| URL | {URL_MARK + URL}, etc. <br> URL_MARK := "http://" \| "www." | http://www.genetic-algorithm.org/GECCO2000/gecco2000mainpage.htm |
| | {DUE_SYM + \<NL\> + date}, <br> {DUE_SYM + ":" + date}, <br> {DUE_SYM + date}, etc. <br> DUE_SYM := "Important Dates" \| <br>　　　"Deadline" \| "Due Dates" \| <br>　　　"Submission Deadline" \| <br>　　　"Paper Submission" <br>　　　"Dates" \| "Schedule" \| <br>　　　"Conference Schedule" \| <br>　　　"CALENDAR" | Important Dates <br><br> December 20, 1999: Deadline for the submissions of the proposals. <br><br> SUBMISSION DEADLINE: January 26, 2000 |
| | {CON_SYM1 + ":" + number}, <br> {CON_SYM2 + ":" + e-mail}, etc. <br> CON_SYM1 := "Phone" \| "Tel" <br>　　　　　\| "Fax" <br> CON_SYM2 := "E-MAIL" | Phone:　　650-328-3123 <br> FAX:　　　650-321-4457 <br> E-MAIL: gecco@aaai.org |

7. CFP (Call-For-Papers)　　　　　　　　　　　　( 　 )

6

.

| | | |
|---|---|---|
| | {TOPIC_SYM + ":" + topic}, etc.<br>TOPIC_SYM := "Topic" + "Title" | Topic:<br>　　SKVORETZ Seminar |
| | {DATE_SYM + ":" + date}, etc.<br>DATE_SYM := "Dates"<br>date := day + "-" + month + "-" +<br>　　+year<br>day := 1-digit \| 2-digit<br>month := month<br>year := 4-digit \| 2-digit | Dates:　　4-May-95<br>Dates:　　15-April-1997 |
| | {LOC_SYM + location}, etc.<br>LOC_SYM := "in" \| "at" | 5:30 in PH 223D.<br>at Porter Hall A18C. |
| | {PER_SYM + name}, etc.<br>PER_SYM := "Mr." \| "Dr." \|<br>　　　　"President" | DR. WILLIAM FISH<br>MR. Jill Fain Lehman |
| | {TIME_SYM + ":" + time}, etc.<br>TIME_SYM := "Time"<br>time := number + ":" + number | 3:30-5:00 |
| | {"-" + time}, etc.<br>time := number + ":" + number | 3:30-5:00 |

8. CMU

CMU　　　　　　　　　　　　8
.
,
.
.

| | | |
|---|---|---|
| | {NAME_SYM + name + RST_SYM}, etc.<br>NAME_SYM := "Named after"<br>RST_SYM := "RESTAURANT" | EL FLORIDITA RESTAURANT<br>1253 N. Vine St., L.A.<br>(213) 871-8612 ... |
| | {STREET NO + STREET + "," + CITY}, etc. | 1253 N. Vine St., L.A. |
| | {"(" + number + ")" + number + "-" + number}, etc. | (213) 871-8612 |
| | {SECTION_SYM + Review + SECTION_SYM}, etc.<br>SECTION_SYM := new line \| new line + tab | Named after Hemingway-'s favorite hangout ... weekend reservations suggested. |
| | {n-CD + card + SECTION_SYM}, etc.<br>CD = {card + ","}<br>card := AE \| CB \| V \| MC \| DC \| DIC \| BC<br>$0 \leq n \leq 7$ | AE, CB, DC, DIS, MC, V. |

9. LA

6 9

HMM $O$

. ,

7 , CFP

6 7 CFP

$O$ ( 6 7 1 6

. $S:$ , $E:$ ).

$$O = (o_1, o_2, \cdots, o_T) = \quad S\ 1\ a\ 2\ b\ 3\ c\ 4\ d\ 5\ e\ 6\ E \qquad (22)$$

```
SECOND CALL FOR PAPERS

FIFTH ANNUAL INTERNATIONAL CONFERENCE ON
COMPUTATIONAL MOLECULAR BIOLOGY

(RECOMB 2001)

April 21-24, 2001
Montrel, Canada


US National Science Foundation

http://recomb2001.gmd.de

The Fifth Annual Conference on Research in Computational Molecular

...
CALENDAR:

Deadline for submission of papers:      Sep 30, 2000
Notification of acceptance/rejection:   Dec  5, 2000
Deadline for reception of final papers: Jan  5, 2001

STEERING COMMITTEE:
...

INFORMATION:
...
...
CANADA H3C 3J7

Tel: (514) 343-7501
Fax: (514) 343-2254
email: recomb01@CRM.UMontreal.CA
```

7. CFP

**3.1.2** **HMM** **HMM**

,                        HMM                         .
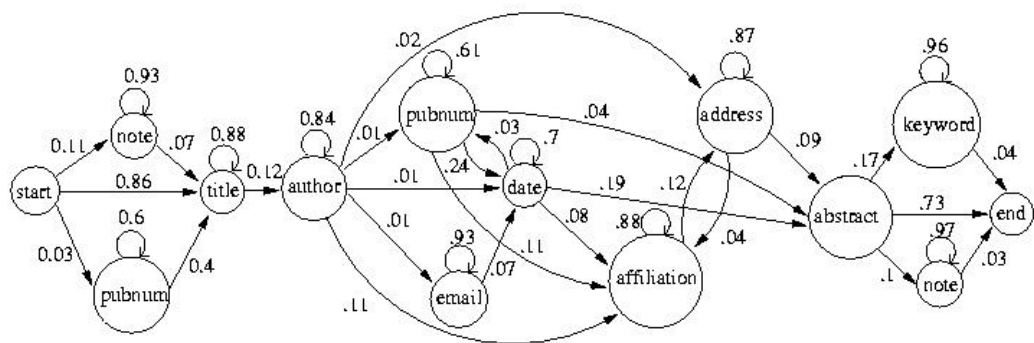
(21)                                    .        (21)

( $S$ )                    ( $E$ )                    .

.

.

HMM . CFP

, HMM 8

[Seymore et al., 1999].  ,

HMM

.  , 1

,

.  -

(forward-backward algorithm) Viterbi , EM

.



그림 8. HMM [Seymore et al., 1999]

,

,

HMM .

9 7 .

9 $\lambda_0$ . $\lambda_0$

(MID),

.



그림 9. 표 7의 CFP 문서에 대한 은닉 마코프 모델 ($\lambda_0$)

그림 9에서 Date와
Location MID

.

.

. 그림 9

Data Location, Conf. Name

.

HMM

관측열 $O$ .

그림 10

(pseudo code) .

- 31 -

Step 0. Do data preprocessing and generate a field generation

rule $R$ for each field.

$\tau$ = initial model state by equation (21).

Step 1. Construct initial model $\lambda_0$ with internal state $\tau$ and

rule $R$.

Step 2. If ( $nState_{\min} <= nState_{\lambda} < nState_{\max}$ ) Goto Step 5.

Step 3. Compute field distance of model $\lambda_0$.

Step 4. if (exist(state pair within state distance $\theta$)) then

Merge nearest two states(one pair).

$nState_{\lambda}$    $nState_{\lambda}$ - 1, Goto Step 2.

else Goto Step 5.

Step 5. Compute $P(O \mid \lambda)$ with observation $O = (o_1, o_2, o_{2,} \cdots, o_T)$
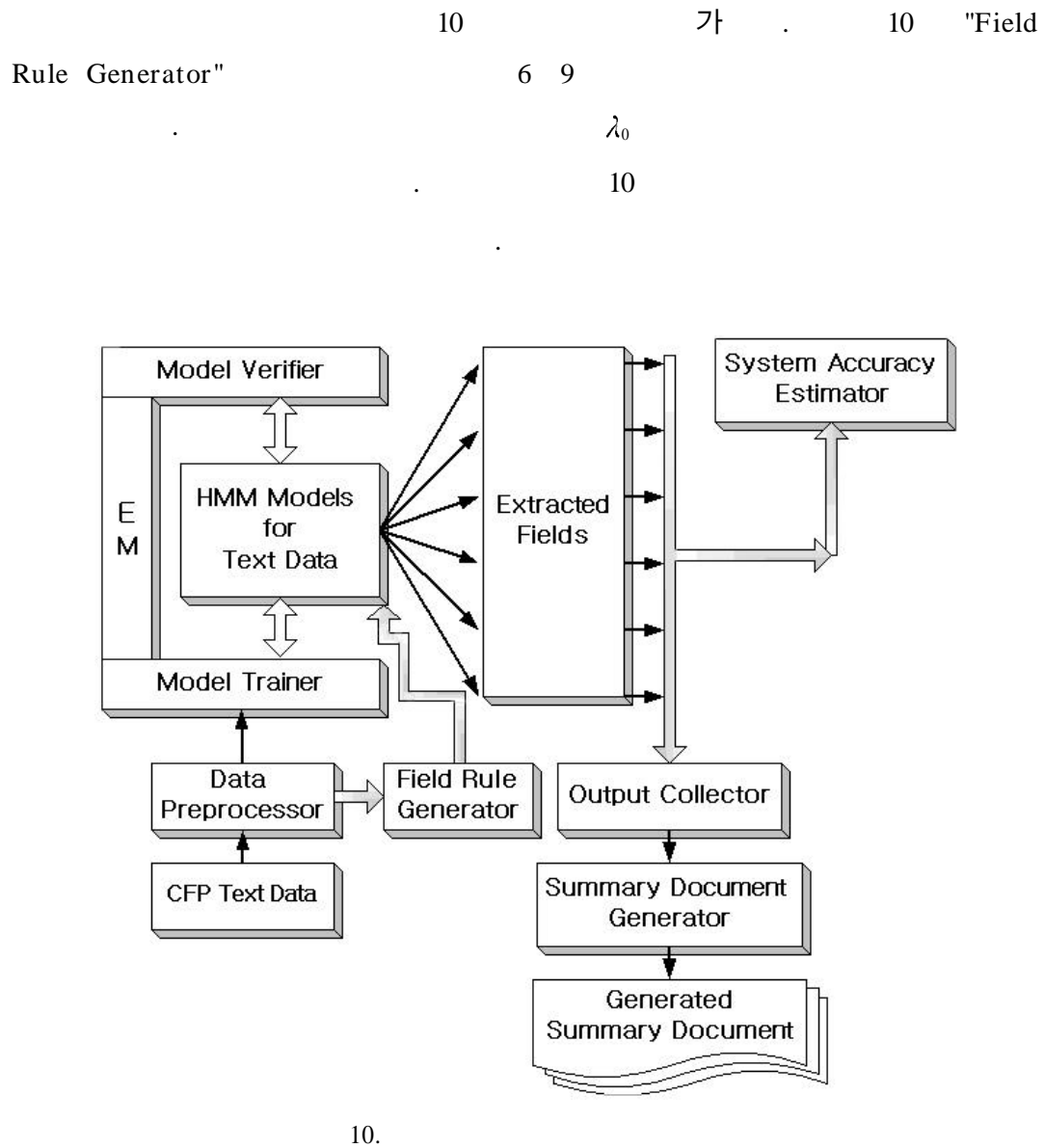
and model $\lambda = (\pi, A, B)$

Step 6. Estimate optimal model parameters with the EM

algorithm in Table 5.

Step 7. Find the optimal state sequence $q = (q_1, q_2, q_{2,} \cdots, q_T)$

with given observation $O = (o_1, o_2, o_{2,} \cdots, o_T)$ and model

$\lambda$.

10.          HMM                    (pseudo-code)

10    Step  3            $i$    $j$                  ,

.

$\theta$

.

$\theta$

.

## 3.2 HMM

10                                    .        10      "Field

Rule  Generator"                          6   9

.                                        $\lambda_0$

.                  10

.



10.

HMM                          CFP

10                              ,        10      CFP  Text

Data                                                    .

                                        .

                                .

                                                                        .

                                                                ,

                                                                .

                                .

                                .

HTML                                        .

                                        .

# 4.

## 4.1 (Call-For-Paper Data)

CFP

.

, , , , (topic), ,

, .

CFP 11 SGML (tagging)

. SGML

. CFP (Computer Science)

(Conference) CFP 200 CFP 100

300 CFP ,

2051K .

1 RECOMB 2001 CFP

. 1 CFP 11 SGML

11 . 11 "<NL>"

(new line character) .

300 CFP 200 CFP

, 100

.

```
<NL>
<paragraph><sentence> SECOND CALL FOR PAPERS </sentence></paragraph>
<NL>
<paragraph><sentence> FIFTH ANNUAL INTERNATIONAL CONFERENCE ON
COMPUTATIONAL MOLECULAR BIOLOGY </sentence></paragraph>
<NL>
<sentence><paragraph>(<c_name>RECOMB 2001</c_name>)</sentence></paragraph>
<NL>
<paragraph><date>April 21-24, 2001</date></paragraph>
<paragraph><location>Montrel, Canada</location></paragraph>
<NL>
<paragraph><sentence>Organized by <NL> Centre de recherches math?atiques
<NL> Universite de Montreal <NL></sentence></paragraph>
<NL>
<paragraph><sentence>Sponsored by <NL>
Association for Computing Machinery (ACM-SIGACT)</sentence></paragraph>
<NL>
<paragraph><sentence>with support from <NL> Celera Genomics <NL>
Compugen <NL> IBM Corporation <NL>
International Society for Computational Biology (ISCB) <NL>
SLOAN Foundation <NL>
SmithKline Beecham <NL>
US Department of Energy <NL>
US National Science Foundation </sentence></paragraph><NL>
<NL>
<c_url>http://recomb2001.gmd.de</c_url>
<NL><NL>
<paragraph><sentence>The Fifth Annual Conference on Research in
Computational Molecular <NL> Biology (<c_name>RECOMB 2001</c_name>),
sponsored by the Association for Computing <NL>
Machinery Special Interest Group on Algorithms and Computation Theory <NL>
(ACM-SIGACT) with support from Celera Genomics, Compugen, IBM <NL>
```

11. CFP (SGML )

| | | (Tag) |
|---|---|---|
| | | <c_name> ... </c_name> |
| | | <date> ... </date> |
| | | <location> ... </location> |
| URL | | <c_url> ... </c_url> |
| | | <due_date> ... </due_date> |
| | , , | <c_contact> ... </c_contact> |
| | CFP . | <paragraph> ... </pagragraph> |
| | CFP . | <sentence> ... </sentence> |

11. CFP

## 4.2 CMU

CMU         CMU(Carnegie Mellon University)

[CMU Data]   ,

            .

1149K   ,   488

  .               X

          SGML              .

             SGML

      . *

      . **           SGML

     ,

               .

          .    12

   ,    13    12        12    SGML

       (tagging)             .

```
<0.2.5.95.11.00.22.cd01+@andrew.cmu.edu.0>
Type:    cmu.andrew.academic.sds.seminars
Topic:   SKVORETZ Seminar
Dates:   4-May-95
Time:    4:00 - 5:30
PostedBy: Carole Deaunovich on 2-May-95 at 11:00 from andrew.cmu.edu
Abstract:

Professor John Skvoretz, U. of South Carolina, Columbia, will present
a seminar entitled "Embedded Commitment," on Thursday, May 4th from
4-5:30 in PH 223D.
```

12. CMU        (    )

```
<0.2.5.95.11.00.22.cd01+@andrew.cmu.edu.0>
Type:    cmu.andrew.academic.sds.seminars
Topic:   SKVORETZ Seminar
Dates:   4-May-95
Time:    <stime>4:00</stime> - <etime>5:30</etime>
PostedBy: Carole Deaunovich on 2-May-95 at 11:00 from andrew.cmu.edu
Abstract:

<paragraph><sentence><speaker>Professor John Skvoretz</speaker>, U. of South Carolina, Columbia, will present
a seminar entitled "Embedded Commitment," on Thursday, May 4th from
<stime>4</stime>-<etime>5:30</etime> in <location>PH 223D</sentence></location>.</paragraph>
```

13. CMU                                (SGML                    )

| | | (Tag) |
|---|---|---|
| (Topic)* | | <topic> ... </topic> |
| * | "25-Oct-93"  . | <date> ... </date> |
| | Room      Hall  . | <location> ... </location> |
| | . | <speaker> ... </speaker> |
| | . | <stime> ... </stime> |
| | . | <etime> ... </etime> |
| ** | . | <paragraph> ...  </pagragraph> |
| ** | . | <sentence> ...  </sentence> |

12. CMU

CMU                                                        488

250                                    ,                238
.

,                                12

,                              SGML

.


## 4.3  LA


LA                                      LA  Weekly (http://www.laweekly.com/)


.                  HTML                                        435K

200                                    .                                14

,            ,      ,            ,                                    ,


.

**EL FLORIDITA RESTAURANT**
1253 N. Vine St., L.A.
(213) 871-8612

Named after Hemingway's favorite hangout in
Havana, El Floridita offers food, fun and festivities.
Try the yellow-rice paella complete with lobster,
mussels and shrimp ($19.50) or the tender roasted
pork marinated in a savory lemon garlic sauce
($10.95). Wash it all down with a Mojito, the
popular Cuban drink made from mint leaves, lime
juice and sugar. Live bands, including salsa, entertain
Thursday through Saturday, and Monday nights.
Prices vary on weekends. Lunch and dinner seven
days. Full bar; takeout; delivery; valet parking;
weekend reservations suggested. AE, CB, DC, DIS,
MC, V.

14.  LA  Weekly

200                                100

100                                          .

.        ,      9                                              SGML

HTML                                      .

## 4.4

15   16                   HMM              HMM              CFP

.        17

.

15   16                                  CFP

(        15)



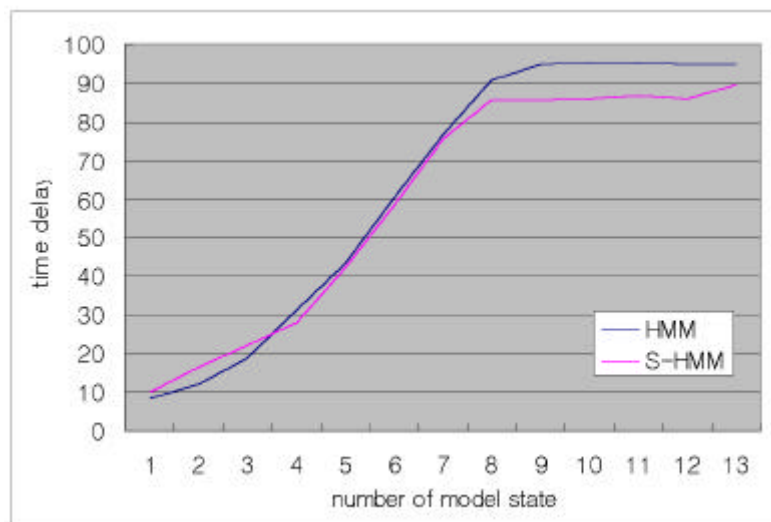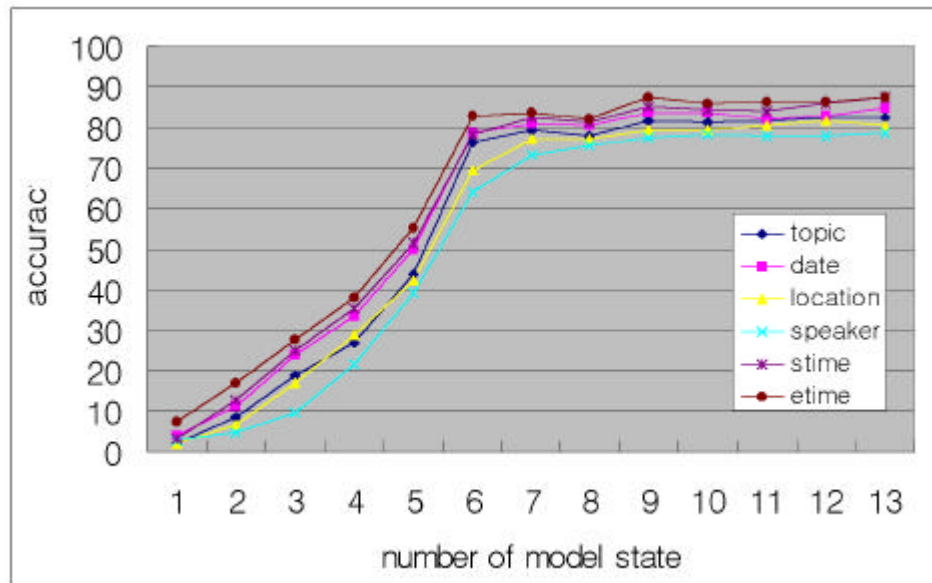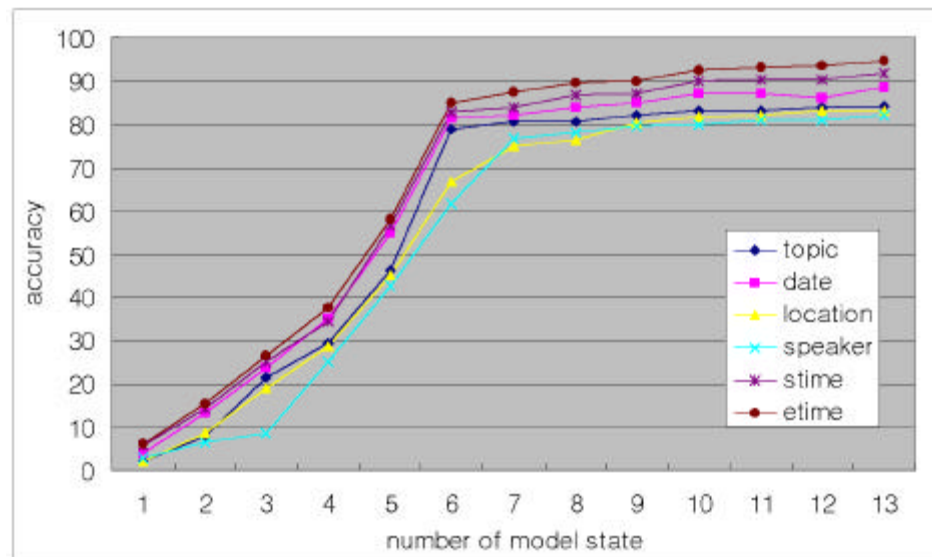15.              HMM              CFP

( 16)                                                                                  .

HMM

.



16.            HMM                CFP



17. HMM      S -HMM                         -  CFP

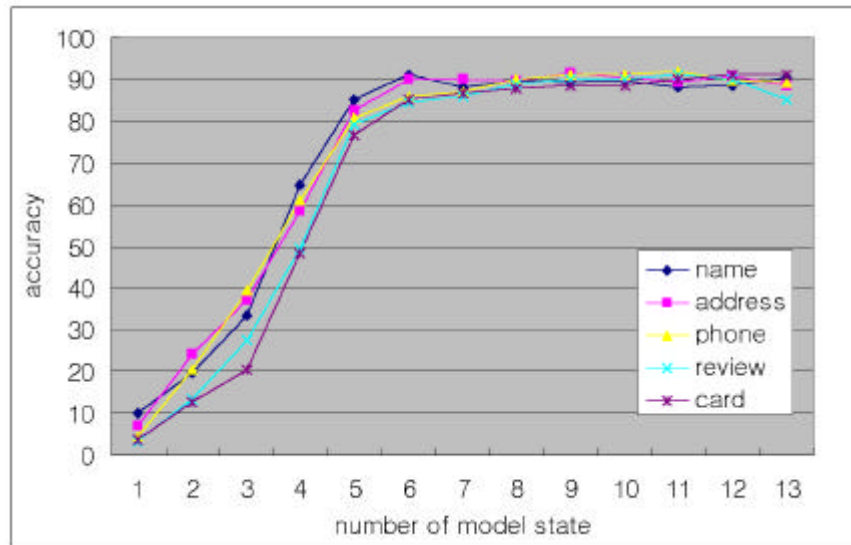17            HMM                          HMM
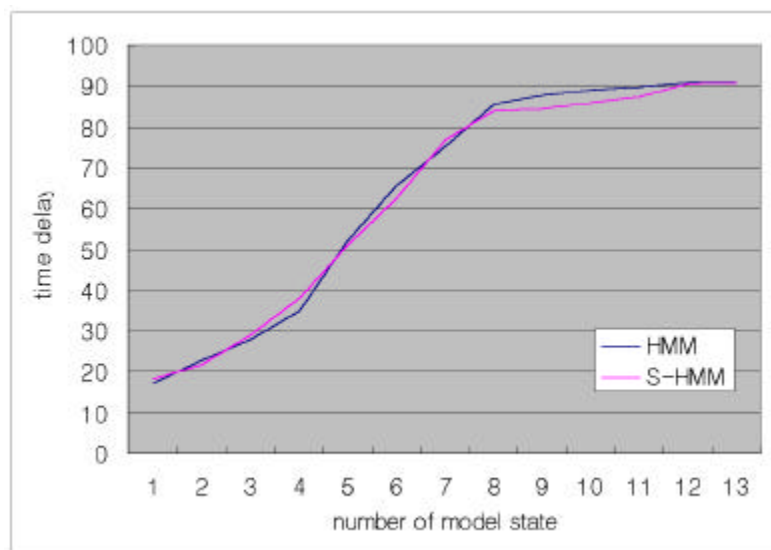10%                                    .



18.          HMM          CMU



19.          HMM          CMU

20. HMM    S-HMM           -  CMU



21.             HMM              LA Weekly

그림 22.                HMM                   LA  Weekly



그림 23.  HMM      S-HMM                    - LA  Weekly


18   19    CMU

.                                            HMM                      HMM

.          ,                    HMM                    20

　　　HMM              .


21   22                                        LA  Weekly

　　　　　　.

　　　　　　　　　　　　　　　　　　　　　　　.        23

　　　　HMM                                        .

**5.**

HMM

,

.

,

.

,

.

HMM

HMM

3

.

,

.

. XML

XML

.

[Baldi et al, 1998] Baldi, P. and Brunak, S., *BIOINFORMATICS - The Machine Learning Approach*, MIT Press, 1998.

[Bikel et al., 1997] Bikel, D. M., Miller, S., Schwartz, R., and Weischedel, R., "Nymble: a high-performance learning", *Proceesings of ANLP-97*, pp. 194-201, 1997.

[Cherkassky et al., 1998] Cherkassky, V. and Mulier, F. *LEARNING FROM DATA - Concepts, Theory, and Methods"*, Wiley-Interscience, 1998.

[CMU Data] Repository of Test Domains for Information Extraction, http://www.isi.edu/~muslea/RISE/repository.html.

[Cohen, 1998] Cohen, W, "A Web-based Information System that Reasons with Structured Collections of Text," *Proceedings of Second International Conference on Autonomous Agents*, pp. 400-407, 1998.

[David et al., 1999] David, R. H. Miller, Tim Leek, and Richard M. Schwartz, "A hidden Markov model information retrieval system", *Proceedings on the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 214-221, 1999.

[Dempster et al., 1997] Dempster, A. P., Laird, N. M., and Rubin, D. B.,

"Maximum likelihood from incomplete data via the EM algorithm (with discussion)", *Journal of Royal Statistic Society B*, Vol. 39, pp. 1-38, 1977.

[Freitag, 1997] Freitag, D., "Using grammatical inference to improve precision in information extraction", *Working Notes of the ICML-97 Workshop on Automata Induction, Grammatical Inference, and Language Acquisition*, 1997.

[Freitag, 1998] Freitag, D., "Multistrategy learning for information extraction," *15th International Conference on Machine Learning*, 1998.

[Freitag & McCallum, 2000] Freitag, D., McCallum, A., "Information extraction with HMM structures learned by stochastic Optimization", *AAAI-2000*, 2000.

[Kushmerick et al., 1997] Kushmerick, N., Weld, D., and Doorenbos, R., "Wrapper Induction for Information Extraction," *International Joint Conference on Artficial Intelligent*, pp. 729-735, 1997.

[Lane, 1999] Lane, T., "Hidden Markov models for human/computer interface modeling", *In Proceedings of the IJCAI-99 Workshop on Learning about Users*, pp. 35-44, 1999.

[Leek, 1997] Leek, T. R., "Information extraction using hidden Markov models", Master's thesis, UC San Diego, 1996.

[Maes, 1994] Maes, P., "Agents that reduce work and information overload",

*Communications of the ACM*, Vol. 37, No. 7, pp. 31-40, 1994.

[Mitchell, 1997] Mitchell, T. M., *Machine Learning*, McGraw-Hill, 1997.

[Muslea et. al, 1998] Muslea, I., Minton, S., and Knoblock, C., "STALKER: Learning extraction rules for semistructured, Web-based information sources", *AAAI 98*, 1998.

[Rabiner, 1989] Rabiner, L. R., "A tutorial on hidden Markov models and selected applications in speech recognition", *Proceedings of IEEE*, Vol. 77, No. 2, February 1989.

[Seymore et al., 1999] Seymore, K., McCallum, A., and Rosenfeld, R., "Learning hidden Markov model structure for information extraction", *AAAI'99 Workshop on Machine Learning for Information Extraction,* pp. 37-42, 1999.

# Abstract

HMM is a kind of automata and it's internal state transitions are decided with some probabilistic values. HMM is used widely for an application with temporal data of sequential characteristics, for example, speech data. This thesis presents a new effective method for building HMM structure for information extraction tasks.

For information extraction tasks, we used HMM based on left-to-right structure. Traditional HMM are used with pre-constructed static model structure and trained its model parameter after model construction. We present here a new HMM called S-HMM (Self-Organizing Hidden Markov Model) that constructs it's structure with the rules that are obtained from training dataset.

In this paper, we used S-HMM for information extraction tasks from Call-For-Papers data, CMU online seminar announcement data, and LA restaurants review and recommendation data (http://www.laweekly.com/). We construct model structure using S-HMM from initial abstract model structure to more detailed structure with the set of rules learned from training data. We could find more appropriate structure with this set of rules. The experimental results show improved average extraction accuracy of 12% increase in average extraction speed in comparison with fixed-state HMM.

Keywords: Information Extraction, Hidden Markov Model, Model Structure Learning, Self-Organizing Hidden Markov Model