

自我來黃州已過三寒
食年、欲惜春、意不
寧惜今年又苦雨、月社
簫瑟、河海、棠花泥
污、遊支雪、閣中偷負
多夜半、真有力、何殊少
年、病起、頭白
春江欲入户、雨勢未
止、雨小屋如漁舟、濛
水雲裏、空庖煮寒菜
破灶燒濕草、那
知是寒食、但見烏
銜泥、君門深
九重、誰能去、在萬里、誰
哭、淪窮、那、更吹不
起

右黃州寒食二首

计算语言学

Computational Linguistics

教师：孙茂松

Tel: 62781286

Email: sms@tsinghua.edu.cn

TA：林衍凯

Email: linyankai423@qq.com

郑重声明

- 此课件仅供选修清华大学计算机系研究生课《计算语言学》(70240052)的学生个人学习使用，所以只允许学生将其下载、存贮在自己的电脑中。未经孙茂松本人同意，任何人不得以任何方式扩散之（包括不得放到任何服务器上）。否则，由此可能引起的一切涉及知识产权的法律责任，概由该人负责。
- 此课件仅限孙茂松本人讲课使用。除孙茂松本人外，凡授课过程中，PTT文件显示此《郑重声明》之情形，即为侵权使用。



第六章

词性自动标注

(隐Markov模型)

6.1. Part-of-Speech Tagging



Task: Select the most likely sequence of syntactic categories for the words in a sentence.

The simplest algorithm: to disambiguate words based solely on the largest probability that a word occurs with a particular tag.

90% success rate

Consider context! Use of some of the local context of the sentence in which the word appears

‘flies’ prefers V $\text{PROB}(0.0005)$ than N $\text{PROB}(0.0003)$
if without context

But given the prior word is ‘the’, ‘flies’ will prefer N

6.1. Part-of-Speech Tagging

Solution?

Let w_1, \dots, w_T be a sequence of words(sentence). The goal is to find the sequence of lexical categories C_1, \dots, C_T that maximizes

$$PROB(C_1, \dots, C_T \mid w_1, \dots, w_T) \quad \text{(formula 1)}$$

According to Bayes' rule:

$$\begin{aligned} & PROB(C_1, \dots, C_T \mid w_1, \dots, w_T) \\ &= \frac{PROB(C_1, \dots, C_T) \times PROB(w_1 \dots w_T \mid C_1, \dots, C_T)}{PROB(w_1, \dots, w_T)} \end{aligned} \quad \text{(formula 2)}$$

Thus the problem reduces to finding the sequence C_1, \dots, C_T that maximizes

$$PROB(C_1, \dots, C_T) \times PROB(w_1, \dots, w_T \mid C_1, \dots, C_T) \quad \text{(formula 3)}$$

Approximation: making some independence assumptions.

6.1. Part-of-Speech Tagging

▲ Approximation for the first expression in formula 3

Using the definition of conditional probability sequentially:

$$\begin{aligned} & \textit{PROB}(C_1, \dots, C_T) \\ &= \textit{PROB}(C_1, \dots, C_{T-1}) \times \textit{PROB}(C_T \mid C_1, \dots, C_{T-1}) \\ &= \textit{PROB}(C_1, \dots, C_{T-2}) \times \textit{PROB}(C_{T-1} \mid C_1, \dots, C_{T-2}) \times \textit{PROB}(C_T \mid C_1, \dots, C_{T-1}) \\ &= \dots \\ &= \textit{PROB}(C_1 C_2) \times \textit{PROB}(C_3 \mid C_1 C_2) \times \dots \times \textit{PROB}(C_{T-1} \mid C_1, \dots, C_{T-2}) \times \\ & \quad \textit{PROB}(C_T \mid C_1, \dots, C_{T-1}) \\ &= \textit{PROB}(C_1) \times \textit{PROB}(C_2 \mid C_1) \times \textit{PROB}(C_3 \mid C_1 C_2) \times \dots \times \\ & \quad \textit{PROB}(C_{T-1} \mid C_1, \dots, C_{T-2}) \times \textit{PROB}(C_T \mid C_1, \dots, C_{T-1}) \end{aligned}$$

n-gram models:

bigram $\textit{PROB}(C_i \mid C_{i-1})$

trigram $\textit{PROB}(C_i \mid C_{i-1} C_{i-2})$

Using bigram approximation:

$$\begin{aligned} \textit{PROB}(C_1, \dots, C_T) &\cong \textit{PROB}(C_1 \mid C_0) \times \textit{PROB}(C_2 \mid C_1) \times \dots \times \textit{PROB}(C_T \mid C_{T-1}) \\ &= \prod_{i=1, T} \textit{PROB}(C_i \mid C_{i-1}) \end{aligned}$$

6.1. Part-of-Speech Tagging

To account the beginning of a sentence: pseudo-category ϕ at position 0 as the value of C_0

Example:

$$\begin{aligned} & \textit{PROB}(\textit{ART } N \textit{ } V \textit{ } N) \\ &= \textit{PROB}(\textit{ART} \mid \phi) \times \textit{PROB}(N \mid \textit{ART}) \times \textit{PROB}(V \mid N) \times \textit{PROB}(N \mid V) \end{aligned}$$

▲ Approximation for the second expression in formula 3

$$\textit{PROB}(w_1, \dots, w_T \mid C_1, \dots, C_T) \cong \prod_{i=1, T} \textit{PROB}(w_i \mid C_i)$$

Example:

$$\begin{aligned} & \textit{PROB}(\textit{the fly likes flowers} \mid \textit{ART } N \textit{ } V \textit{ } N) \\ &= \textit{PROB}(\textit{the} \mid \textit{ART}) \times \textit{PROB}(\textit{fly} \mid N) \times \textit{PROB}(\textit{likes} \mid V) \times \textit{PROB}(\textit{flowers} \mid N) \end{aligned}$$

6.1. Part-of-Speech Tagging

▲ The combination of the above two approximations

With these two approximations, the problem changes into finding the sequence C_1, \dots, C_T that maximizes the value of

$$\begin{aligned} & \prod_{i=1, T} \text{PROB}(C_i | C_{i-1}) \times \prod_{i=1, T} \text{PROB}(w_i | C_i) \\ &= \prod_{i=1, T} \text{PROB}(C_i | C_{i-1}) \times \text{PROB}(w_i | C_i) \end{aligned} \quad (\text{formula 4})$$

▲ Train the statistical parameters by annotated corpus

bigram probabilities:

$$\text{PROB}(C_i = V | C_{i-1} = N) \cong \frac{\text{Count}(N \text{ at position } i-1 \text{ and } V \text{ at } i)}{\text{Count}(N \text{ at position } i-1)}$$

lexical-generation probabilities:

$$\text{PROB}(w_i | C_i) \cong \frac{\text{Count}(w_i \text{ with } C_i)}{\text{Count}(C_i)}$$

6.1. Part-of-Speech Tagging

$$\begin{aligned} & \text{PROB}(N V \text{ART} N) \\ &= \text{PROB}(N | \phi) \times \text{PROB}(V | N) \times \text{PROB}(\text{ART} | V) \times \text{PROB}(N | \text{ART}) \\ &= 0.29 \times 0.43 \times 0.65 \times 1 = 0.081 \end{aligned}$$

$$\begin{aligned} & \text{PROB}(\text{flies likes a flower} | N V \text{ART} N) \\ &= \text{PROB}(\text{flies} | N) \times \text{PROB}(\text{likes} | V) \times \text{PROB}(a | \text{ART}) \times \text{PROB}(\text{flower} | N) \\ &= 0.025 \times 0.1 \times 0.36 \times 0.063 = 0.000054 = 5.4 \times 10^{-5} \end{aligned}$$

=>HMM (Hidden Markov Model)

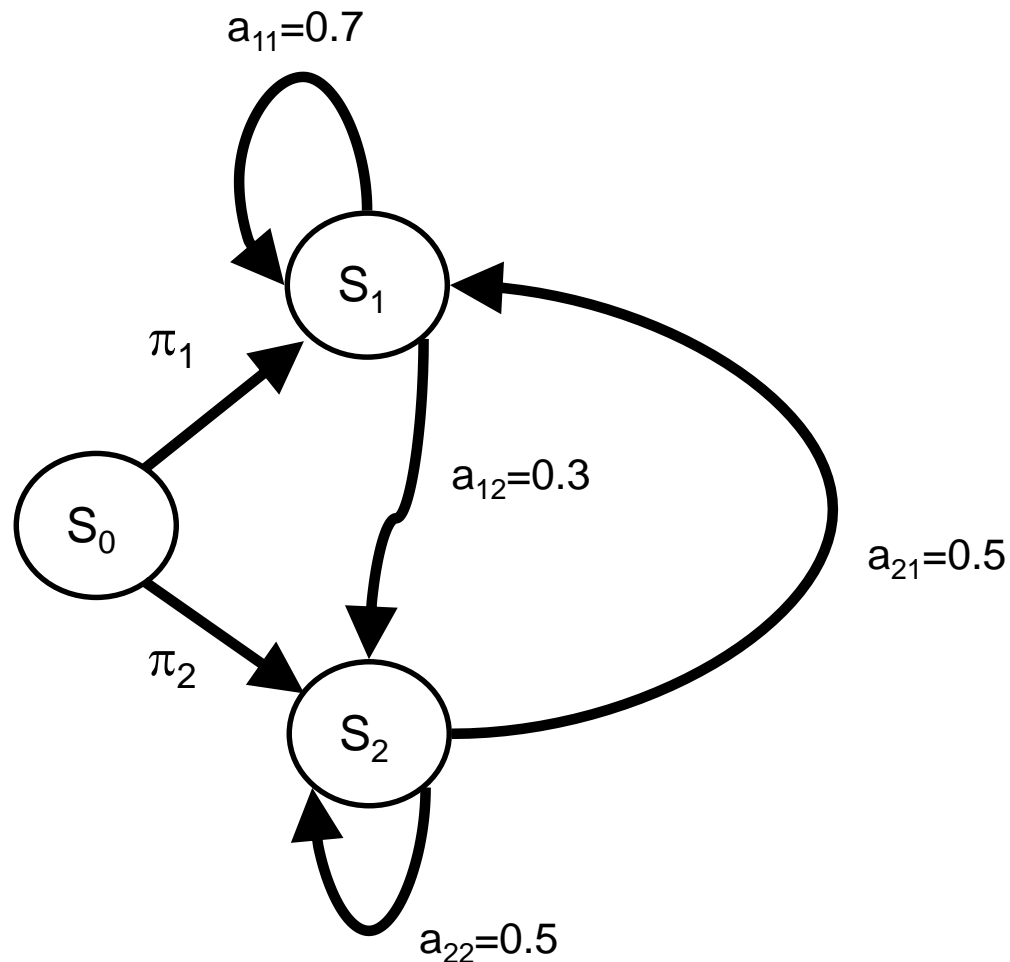
Solution? To find the sequence C_1, \dots, C_T that maximizes the value of

$$\begin{aligned} & \prod_{i=1, T} \text{PROB}(C_i | C_{i-1}) \times \prod_{i=1, T} \text{PROB}(w_i | C_i) \\ &= \prod_{i=1, T} \text{PROB}(C_i | C_{i-1}) \times \text{PROB}(w_i | C_i) \end{aligned}$$

Brute force method: combinatorial explosion problem

6.1. Part-of-Speech Tagging

VMM (Visible Markov Model)



6.1. Part-of-Speech Tagging

HMM

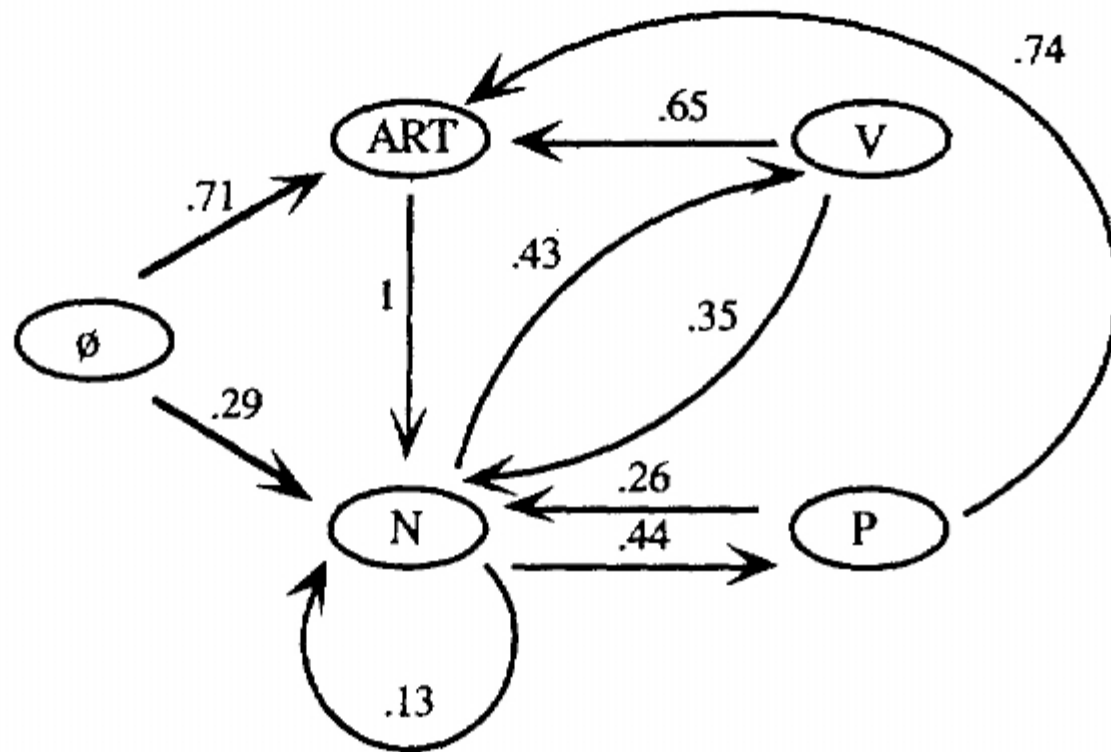
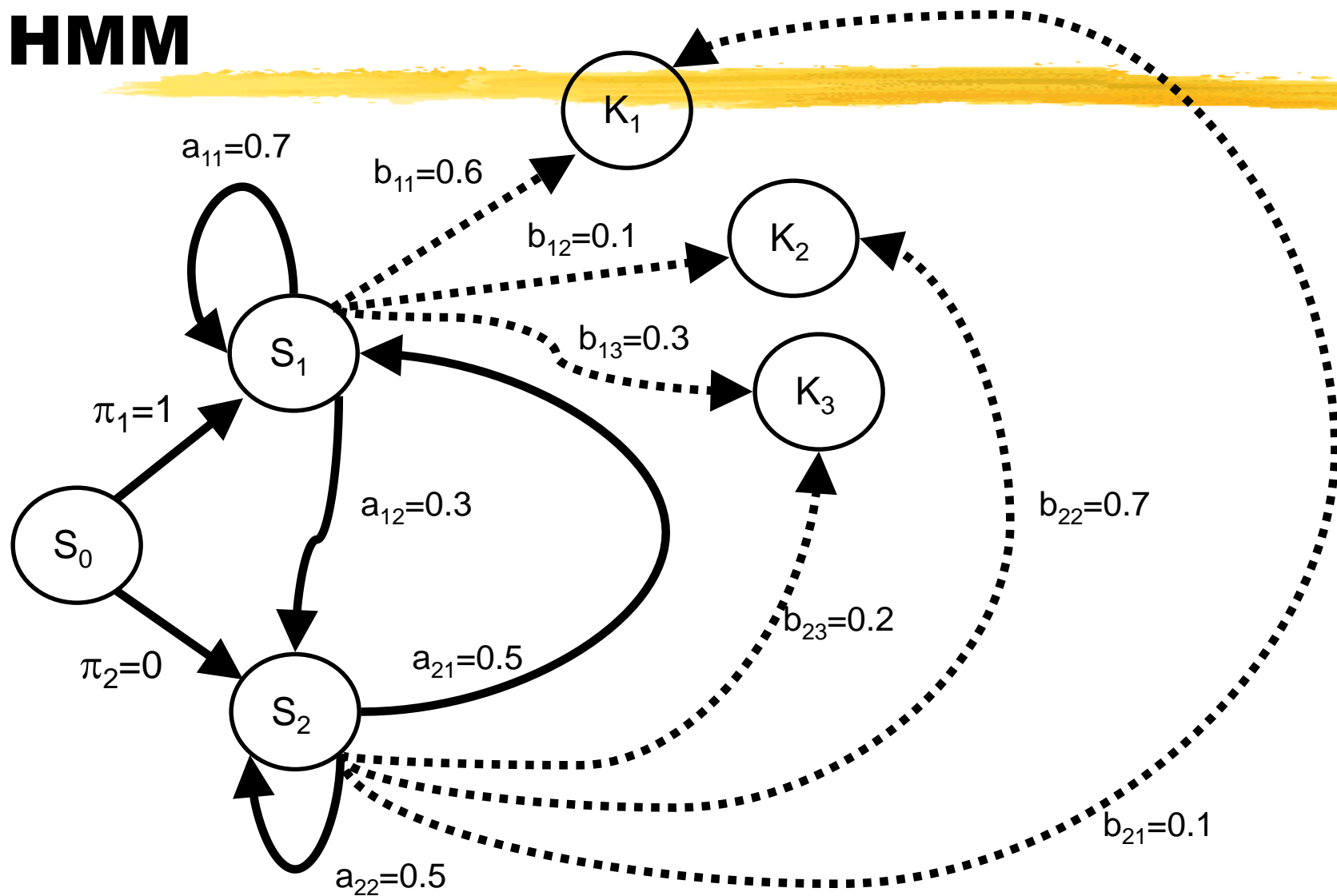


Figure 7.7 A Markov chain capturing the bigram probabilities

6.1. Part-of-Speech Tagging

HMM



6.1. Part-of-Speech Tagging

The Viterbi Algorithm

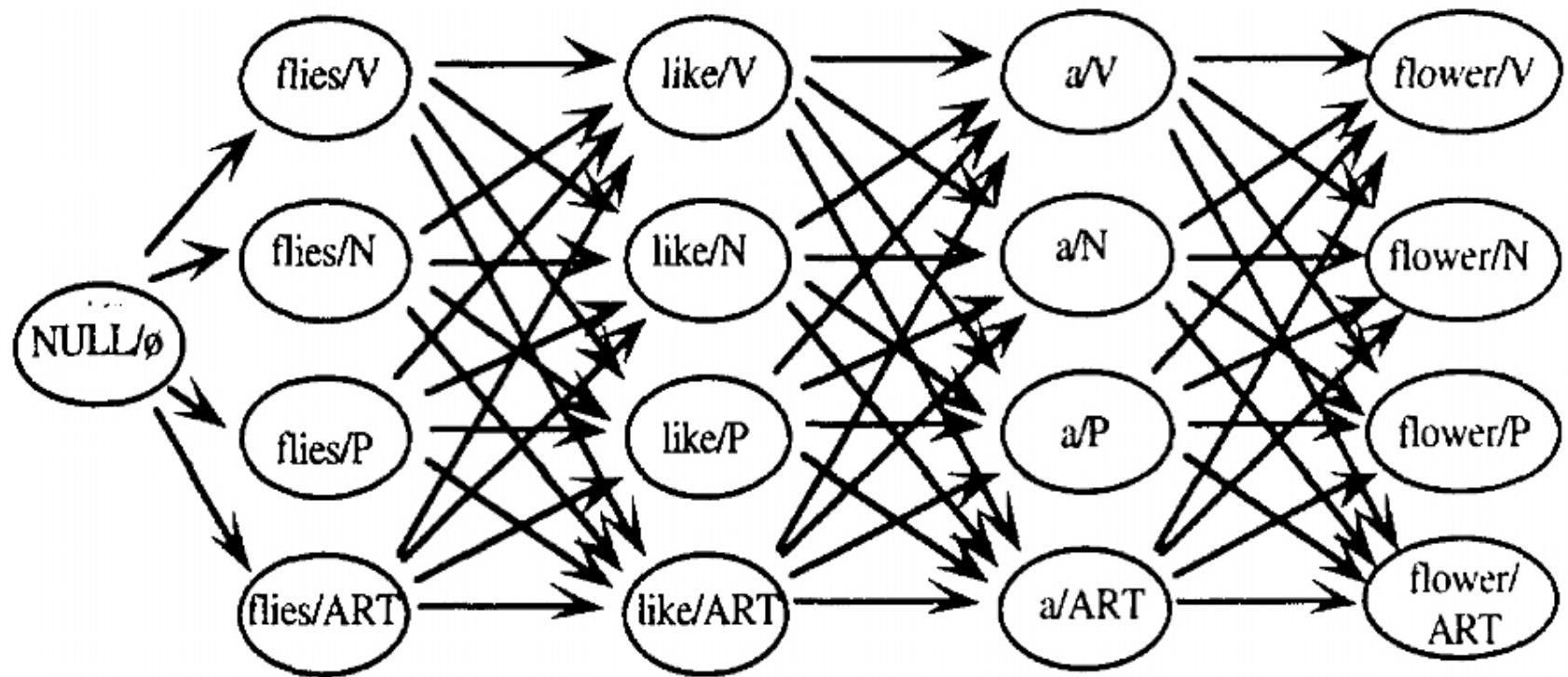


Figure 7.8 Encoding the 256 possible sequences exploiting the Markov assumption

6.1. Part-of-Speech Tagging

The Viterbi Algorithm

Data structure used in the algorithm:

- ▲ A $N \times T$ array SEQSCORE(n,t): record the probability for the best sequence up to the position t that ends with a word in category L_n , where N is the number of lexical categories (L_1, \dots, L_N) and T is the number of words in the sentence (w_1, \dots, w_T).

SEQSCORE	1	2	T
1				
.....				
N				

- ▲ A $N \times T$ array BACKPTR: indicate for each category at each position t what the preceding category is in the best sequence at position t-1.

BACKPTR	1	2	T
1				
.....				
N				

6.1. Part-of-Speech Tagging

The Viterbi algorithm:

Given word sequence w_1, \dots, w_T , lexical categories L_1, \dots, L_N , lexical probabilities $PROB(w_t | L_i)$, and bigram probabilities $PROB(L_i | L_j)$, find the most likely sequence of lexical categories C_1, \dots, C_T for the word sequence.

Initialization Step

For $i = 1$ to N do

$$SEQSCORE(i, 1) = PROB(w_1 | L_i) * PROB(L_i | \phi)$$

$$BACKPTR(i, 1) = 0$$

Iteration Step

For $t = 2$ to T do

For $i = 1$ to N do

$$SEQSCORE(i, t) =$$

$$MAX_{j=1, \dots, N} (SEQSCORE(j, t-1) * PROB(L_i | L_j)) * PROB(w_t | L_i)$$

$$BACKPTR(i, t) = \text{index of } j \text{ that gave the max above}$$

Sequence Identification Step

$$C(T) = i \text{ that maximizes } SEQSCORE(i, T)$$

$$\text{For } i = T-1 \text{ to } 1 \text{ do } C(i) = BACKPTR(C(i+1), i+1)$$

6.1. Part-of-Speech Tagging

Example: *Flies like a flower.*

Bigram probabilities:

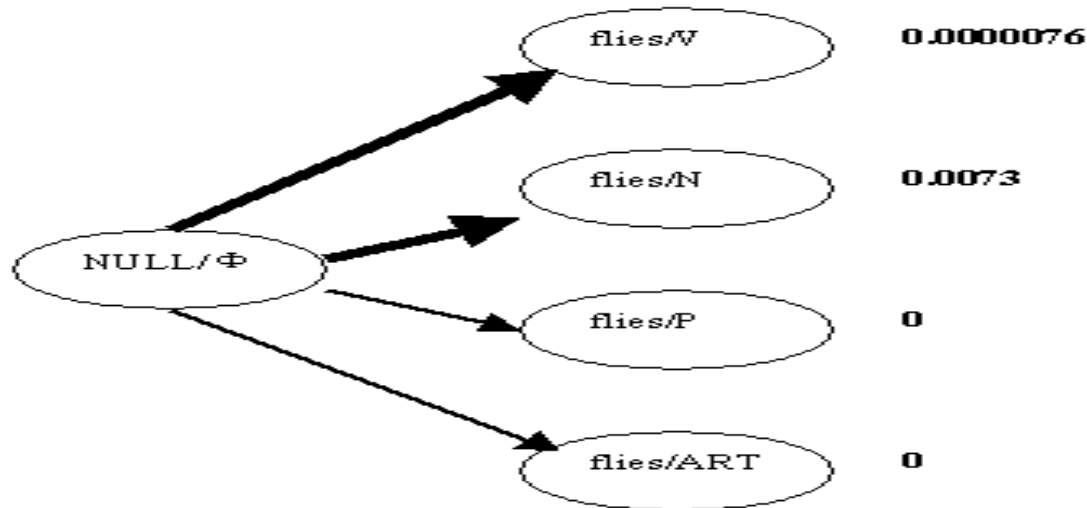
$PROB(ART \phi) = 0.71$	$PROB(N V) = 0.35$
$PROB(N \phi) = 0.29$	$PROB(ART V) = 0.65$
$PROB(N ART) = 1$	$PROB(ART P) = 0.74$
$PROB(V N) = 0.43$	$PROB(N P) = 0.26$
$PROB(N N) = 0.13$	$PROB(\text{other bigram}) = 0.0001$
$PROB(P N) = 0.44$	

The lexical-generation probabilities:

$PROB(\text{flies} N) = 0.025$	$PROB(a ART) = 0.36$
$PROB(\text{flies} V) = 0.076$	$PROB(a N) = 0.001$
$PROB(\text{like} V) = 0.10$	$PROB(\text{flower} N) = 0.063$
$PROB(\text{like} P) = 0.068$	$PROB(\text{flower} V) = 0.05$
$PROB(\text{like} N) = 0.012$	

6.1. Part-of-Speech Tagging

(1) Initialization Step



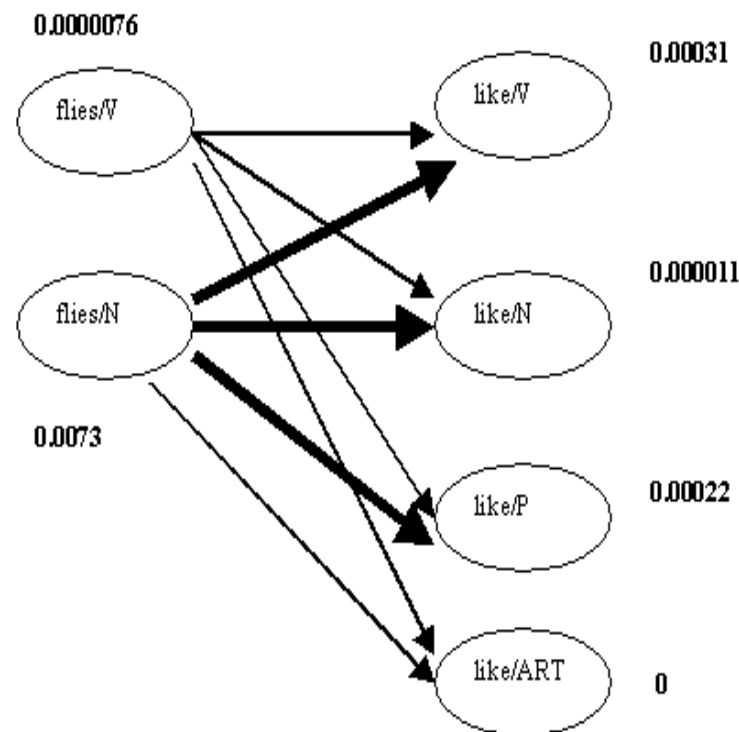
SEQSCORE	1(flies)	2(like)	3(a)	4(flower)
1(V)	0.0000076			
2(N)	0.0073			
3(P)	0			
4(ART)	0			

BACKPTR	1(flies)	2(like)	3(a)	4(flower)
1(V)	0			
2(N)	0			
3(P)	0			
4(ART)	0			

6.1. Part-of-Speech Tagging

(2) Iteration Step

(2.1) t=2

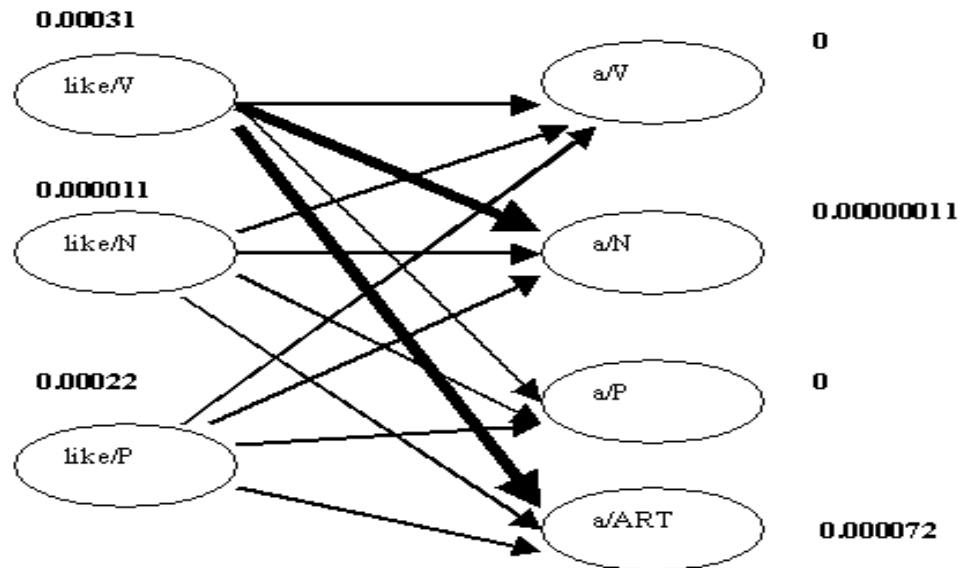


SEQSCORE	1(flies)	2(like)	3(a)	4(flower)
1(V)	0.0000076	0.00031		
2(N)	0.0073	0.000011		
3(P)	0	0.00022		
4(ART)	0	0		

BACKPTR	1(flies)	2(like)	3(a)	4(flower)
1(V)	0	2		
2(N)	0	2		
3(P)	0	2		
4(ART)	0			

6.1. Part-of-Speech Tagging

(2.2) t=3

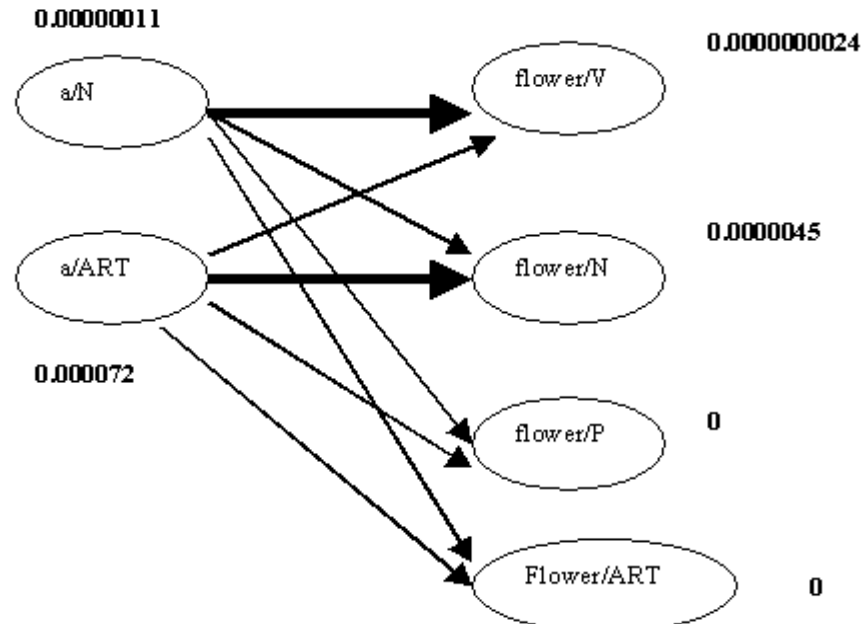


SEQSCORE	1(flies)	2(like)	3(a)	4(flower)
1(V)	0.0000076	0.00031	0	
2(N)	0.0073	0.000011	0.00000011	
3(P)	0	0.00022	0	
4(ART)	0	0	0.000072	

BACKPTR	1(flies)	2(like)	3(a)	4(flower)
1(V)	0	2		
2(N)	0	2	1	
3(P)	0	2		
4(ART)	0		1	

6.1. Part-of-Speech Tagging

(2.3) t=4



SEQSCORE	1(flies)	2(like)	3(a)	4(flower)
1(V)	0.0000076	0.00031	0	0.0000000024
2(N)	0.00725	0.000011	0.00000011	0.00000045
3(P)	0	0.00022	0	0
4(ART)	0	0	0.000072	0

BACKPTR	1(flies)	2(like)	3(a)	4(flower)
1(V)	0	2		2
2(N)	0	2	1	4
3(P)	0	2		
4(ART)	0		1	

6.1. Part-of-Speech Tagging



(3) Sequence Identification Step

$C(4)=2$

$i=3: C(3)=\text{BACKPTR}(C(4), 4)=\text{BACKPTR}(2, 4)=4$

$i=2: C(2)=\text{BACKPTR}(C(3), 3)=\text{BACKPTR}(4, 3)=1$

$i=1: C(1)=\text{BACKPTR}(C(2), 2)=\text{BACKPTR}(1, 2)=2$

Solution: flies(N) like(V) a(ART) flower(N)

6.2. Tag Set

1.	CC	Coordinating conjunction	19.	PP\$	Possessive pronoun
2.	CD	Cardinal number	20.	RB	Adverb
3.	DT	Determiner	21.	RBR	Comparative adverb
4.	EX	Existential <i>there</i>	22.	RBS	Superlative Adverb
5.	FW	Foreign word	23.	RP	Particle
6.	IN	Preposition / subord. conj	24.	SYM	Symbol (math or scientific)
7.	JJ	Adjective	25.	TO	to
8.	JJR	Comparative adjective	26.	UH	Interjection
9.	JJS	Superlative adjective	27.	VB	Verb, base form
10.	LS	List item marker	28.	VBD	Verb, past tense
11.	MD	Modal	29.	VBG	Verb, gerund/pres. participle
12.	NN	Noun, singular or mass	30.	VBN	Verb, past participle
13.	NNS	Noun, plural	31.	VBZ	Verb, non-3s, present
14.	NNP	Proper noun, singular	32.	VBZ	Verb, 3s, present
15.	NNPS	Proper noun, plural	33.	WDT	Wh-determiner
16.	PDT	Predeterminer	34.	WP	Wh-pronoun
17.	POS	Possessive ending	35.	WPZ	Possessive wh-pronoun
18.	PRP	Personal pronoun	36.	WRB	Wh-adverb

The Penn Treebank tag set

6.3. Related Factors



(1) Handling Unknown Words:

$\mathbf{w}_1 \mathbf{w}_2 \mathbf{w}_3$

$\mathbf{C}_1 \mathbf{C}_2$

pick the category \mathbf{C} for the unknown word that maximizes

$\text{PROB}(\mathbf{C} \mid \mathbf{C}_1 \mathbf{C}_2)$

(2) n-grams? $n=3?$

(3) 中文：分词与词性标注一体化？

6.4. The Three Fundamental Questions for HMM

General form of an HMM:

Set of states $S = \{s_1, \dots, s_N\}$

Output alphabet $K = \{k_1, \dots, k_M\} = \{1, \dots, M\}$

Initial state probabilities $\Pi = \{\pi_i\}, i \in S$

State transition probabilities $A = \{a_{ij}\}, i, j \in S$

Symbol emission probabilities $B = \{b_{ijk}\}, i, j \in S, k \in K$

State sequence $x = (X_1, \dots, X_{T+1}) \quad X_t : s \mapsto \{1, \dots, N\}$

Output sequence $o = (o_1, \dots, o_T) \quad o_t \in K$


```
1  $t := 1$ ;  
2 Start in state  $s_i$  with probability  $\pi_i$  (i.e.,  $X_1 = i$ )  
3 forever do  
4     Move from state  $s_i$  to state  $s_j$  with probability  $a_{ij}$  (i.e.,  $X_{t+1} = j$ )  
5     Emit observation symbol  $o_t = k$  with probability  $b_{ijk}$   
6      $t := t + 1$   
7 end
```

6.4. The Three Fundamental Questions for HMM

The Three Fundamental Questions for HMMs

There are three fundamental questions that we want to know about an HMM:

1. Given a model $\mu = (A, B, \Pi)$, how do we efficiently compute how likely a certain observation is, that is $P(O|\mu)$?
2. Given the observation sequence O and a model μ , how do we choose a state sequence (X_1, \dots, X_{T+1}) that best explains the observations?
3. Given an observation sequence O , and a space of possible models found by varying the model parameters $\mu = (A, B, \pi)$, how do we find the model that best explains the observed data?

Fundamental question 1:

Finding the probability of an observation

Given the observation sequence $O = (o_1, \dots, o_T)$ and a model $\mu = (A, B, \Pi)$, we wish to know how to efficiently compute $P(O|\mu)$ – the probability of the observation given the model. This process is often referred to as *decoding*.

For any state sequence $X = (X_1, \dots, X_{T+1})$,

$$\begin{aligned} P(O|X, \mu) &= \prod_{t=1}^T P(o_t | X_t, X_{t+1}, \mu) \\ &= b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \dots b_{X_T X_{T+1} o_T} \end{aligned}$$

and,

$$P(X|\mu) = \pi_{X_1} a_{X_1 X_2} a_{X_2 X_3} \dots a_{X_T X_{T+1}}$$

Now,

$$P(O, X|\mu) = P(O|X, \mu) P(X|\mu) \quad (2T+1) \cdot N^{T+1} \text{ multiplications}$$

Therefore,

$$\begin{aligned} P(O|\mu) &= \sum_{\mathbf{X}} P(O|X, \mu) P(X|\mu) \\ &= \sum_{X_1 \dots X_{T+1}} \pi_{X_1} \prod_{t=1}^T a_{X_t X_{t+1}} b_{X_t X_{t+1} o_t} \end{aligned}$$

The forward procedure

The form of caching that is indicated in these diagrams is called the forward *procedure*. We describe it in terms of forward variables:

$$\alpha_i(t) = P(o_1 o_2 \dots o_{t-1}, X_t = i | \mu)$$

The forward variable $\alpha_i(t)$ is stored at (s_i, t) in the trellis and expresses the total probability of ending up in state s_i at time t (given that the observations $o_1 \dots o_{t-1}$ were seen). It is calculated by summing probabilities for all incoming arcs at a trellis node. We calculate the forward variables in the trellis left to right using the following procedure:

1. Initialization

$$\alpha_i(1) = \pi_i, \quad 1 \leq i \leq N$$

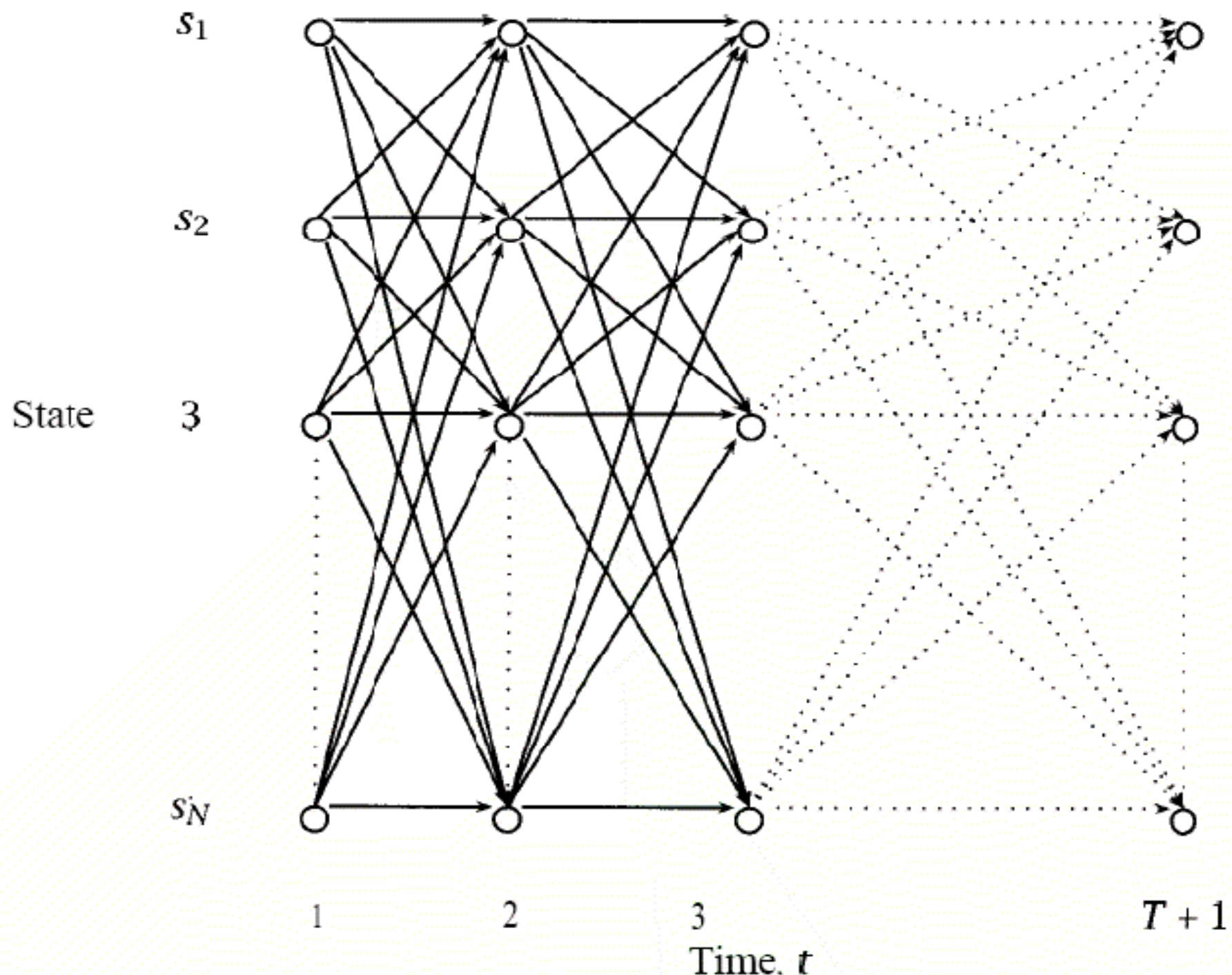
2. Induction

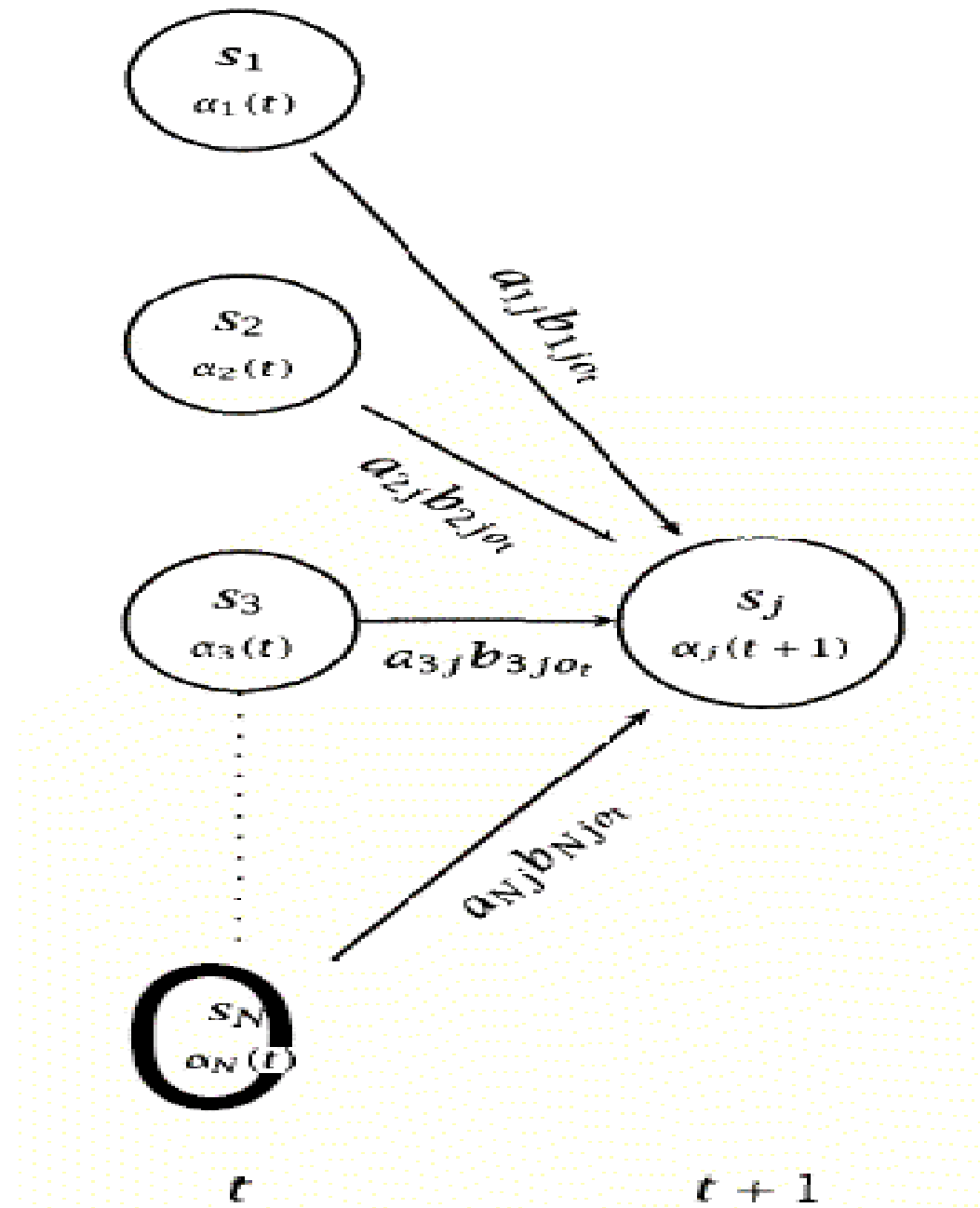
$$\alpha_j(t+1) = \sum_{i=1}^N \alpha_i(t) a_{ij} b_{ij} o_t, \quad 1 \leq t \leq T, 1 \leq j \leq N$$

3. Total

$$P(O | \mu) = \sum_{i=1}^N \alpha_i(T+1)$$

This is a much cheaper algorithm that requires only $2N^2T$ multiplications.





6.4. The Three Fundamental Questions for HMM

The backward procedure

$$\beta_i(t) = P(o_t \cdot \dots \cdot o_T | X_t = i, \mu)$$

1. Initialization

$$\beta_i(T + 1) = 1, \quad 1 \leq i \leq N$$

2. Induction

$$\beta_i(t) = \sum_{j=1}^N a_{ij} b_{ij} o_t \beta_j(t + 1), \quad 1 \leq t \leq T, 1 \leq i \leq N$$

3. Total

$$P(O|\mu) = \sum_{i=1}^N \pi_i \beta_i(1)$$

Combining them

More generally, in fact, we can use any combination of forward and backward caching to work out the probability of an observation sequence. Observe that:

$$\begin{aligned}P(O, X_t = i | \mu) &= P(o_1 \dots o_T, X_t = i | \mu) \\&= P(o_1 \dots o_{t-1}, X_t = i, o_t \dots o_T | \mu) \\&= P(o_1 \dots o_{t-1}, X_t = i | \mu) \\&\quad \times P(o_t \dots o_T | o_1 \dots o_{t-1}, X_t = i, \mu) \\&= P(o_1 \dots o_{t-1}, X_t = i | \mu) P(o_t \dots o_T | X_t = i, \mu) \\&= \alpha_i(t) \beta_i(t)\end{aligned}$$

Therefore:

$$P(O | \mu) = \sum_{i=1}^N \alpha_i(t) \beta_i(t), \quad 1 \leq t \leq T + 1$$

6.4. The Three Fundamental Questions for HMM

Fundamental question 2:

Finding the best state sequence that best explains the observations.

Viterbi algorithm

Commonly we want to find the most likely complete path, that is:

$$\arg \max_{\mathbf{x}} P(\mathbf{X}|\mathbf{O}, \mu)$$

To do this, it is sufficient to maximize for a fixed \mathbf{O} :

$$\arg \max_{\mathbf{x}} P(\mathbf{X}, \mathbf{O}|\mu)$$

An efficient trellis algorithm for computing this path is the *Viterbi algorithm*. Define:

$$\delta_j(t) = \max_{x_1 \dots x_{t-1}} P(X_1 \dots X_{t-1}, o_1 \dots o_{t-1}, X_t = j|\mu)$$

This variable stores for each point in the trellis the probability of the most probable path that leads to that node. The corresponding variable $\psi_j(t)$ then records the node of the incoming arc that led to this most probable path. Using dynamic programming, we calculate the most probable path through the whole trellis as follows:

1. Initialization

$$\delta_j(1) = \pi_j, \quad 1 \leq j \leq N$$

2. Induction

$$\delta_j(t+1) = \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij} o_t, \quad 1 \leq j \leq N$$

Store backtrace

$$\psi_j(t+1) = \arg \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij} o_t, \quad 1 \leq j \leq N$$

3. Termination and path readout (by backtracking). The most likely state sequence is worked out from the right backwards:

$$\hat{X}_{T+1} = \arg \max_{1 \leq i \leq N} \delta_i(T+1)$$

$$\hat{X}_t = \psi_{\hat{X}_{t+1}}(t+1)$$

$$P(\hat{X}) = \max_{1 \leq i \leq N} \delta_i(T+1)$$

6.4. The Three Fundamental Questions for HMM

Fundamental question 3: Parameter estimation

Given a certain observation sequence, we want to find the values of the model parameters $\mu = (A, B, \pi)$ which best explain what we observed. Using Maximum Likelihood Estimation, that means we want to find the values that maximize $P(O|\mu)$:

$$\arg \max_{\mu} P(O_{\text{training}}|\mu)$$

There is no known analytic method to choose μ to maximize $P(O|\mu)$. But we can locally maximize it by an iterative hill-climbing algorithm. This algorithm is the *Baum-Welch* or *Forward-Backward algorithm*, which is a special case of the Expectation Maximization method which we will cover

Define $p_t(i, j)$, $1 \leq t \leq T$, $1 \leq i, j \leq N$ as shown below. This is the probability of traversing a certain arc at time t given observation sequence O ; see figure 9.7.

$$\begin{aligned} p_t(i, j) &= \frac{P(X_t = i, X_{t+1} = j | O, \mu)}{P(O | \mu)} \\ &= \frac{P(X_t = i, X_{t+1} = j, O | \mu)}{P(O | \mu)} \end{aligned}$$

$$= \frac{\alpha_i(t) a_{ij} b_{ij|o_t} \beta_j(t+1)}{\sum_{m=1}^N \alpha_m(t) \beta_m(t)}$$

$$= \frac{\alpha_i(t) a_{ij} b_{ij|o_t} \beta_j(t+1)}{\sum_{m=1}^N \sum_{n=1}^N \alpha_m(t) a_{mn} b_{mn|o_t} \beta_n(t+1)}$$

Note that $\gamma_i(t) = \sum_{j=1}^N p_t(i, j)$.

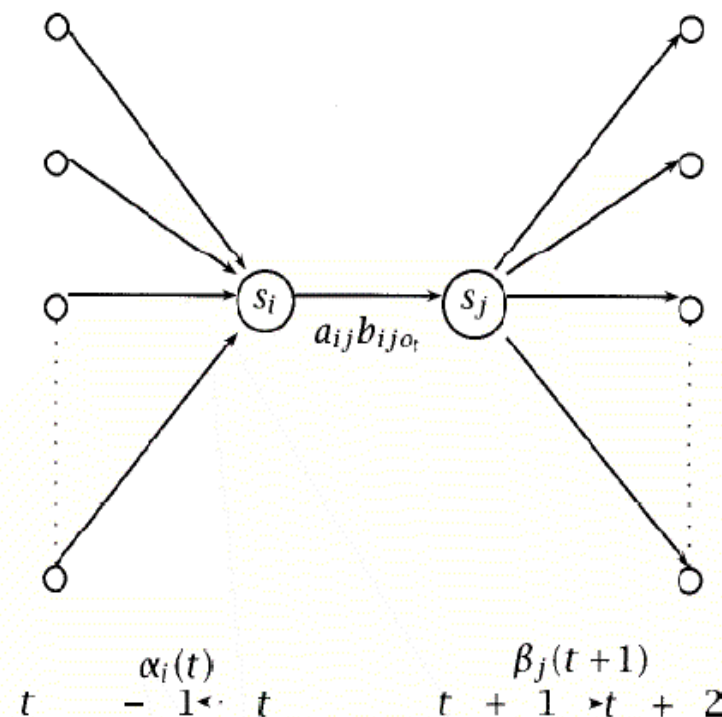


Figure 9.7 The probability of traversing an arc. Given an observation sequence O and a model, we can work out the probability that the Markov process went from state s_i to s_j at time t .

Now, if we sum over the time index, this gives us expectations (counts):

$$\sum_{t=1}^T y_i(t) = \text{expected number of transitions from state } i \text{ in } O$$

$$\sum_{t=1}^T p_t(i, j) = \text{expected number of transitions from state } i \text{ to } j \text{ in } O$$

$$\begin{aligned} \hat{\pi}_i &= \text{expected frequency in state } i \text{ at time } t = 1 \\ &= y_i(1) \end{aligned}$$

$$\begin{aligned} \hat{a}_{ij} &= \frac{\text{expected number of transitions from state } i \text{ to } j}{\text{expected number of transitions from state } i} \\ &= \frac{\sum_{t=1}^T p_t(i, j)}{\sum_{t=1}^T y_i(t)} \end{aligned}$$

$$\begin{aligned} \hat{b}_{ijk} &= \frac{\text{expected number of transitions from } i \text{ to } j \text{ with } k \text{ observed}}{\text{expected number of transitions from } i \text{ to } j} \\ &= \frac{\sum_{\{t: o_t=k, 1 \leq t \leq T\}} p_t(i, j)}{\sum_{t=1}^T p_t(i, j)} \end{aligned}$$

$$P(O|\hat{\mu}) \geq P(O|\mu)$$

This is a general property of the EM algorithm