

1

Hidden Markov Model

Wei-Ta Chu

Markov Model

Chain rule

$$P(A, B, C) = P(A|B, C)P(B|C)P(C)$$

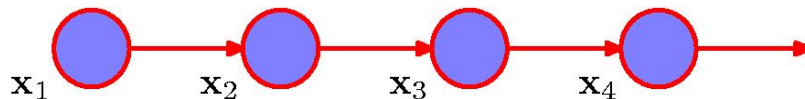
2

- Assume that each of the condition distributions is independent of all previous observations except the most recent, we obtain the *first-order Markov chain*.

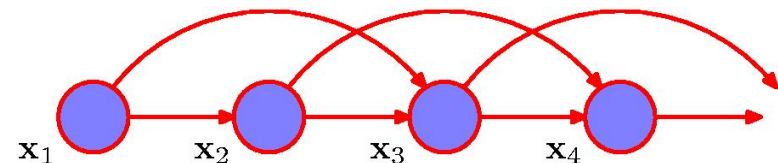
$$p(x_1, \dots, x_N) = \prod_{n=1}^N p(x_n | x_1, \dots, x_{n-1})$$

$$\Rightarrow p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n | x_{n-1})$$

$$p(x_n | x_1, \dots, x_{n-1}) = p(x_n | x_{n-1})$$



First-order Markov chain



Second-order Markov chain

Example

3

- What's the probability that the weather for eight consecutive days is “sun-sun-sun-rain-rain-sun-cloudy-sun”?

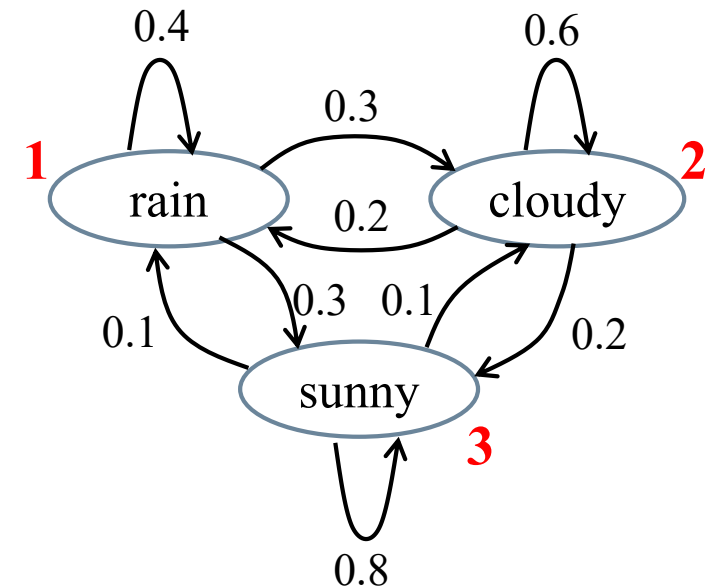
$$P(\mathbf{O}|\text{Model}) = P[3, 3, 3, 1, 1, 3, 2, 3|\text{Model}]$$

$$= P[3]P[3|3]^2P[1|3]P[1|1]P[3|1]P[2|3]P[3|2]$$

$$= \pi_3(a_{33})^2a_{31}a_{11}a_{13}a_{32}a_{23}$$

$$= (1.0)(0.8)^2(0.1)(0.4)(0.3)(0.1)(0.2)$$

$$= 1.536 \times 10^{-4}$$



$$A = a_{ij} = \begin{bmatrix} 0.4 & 0.3 & 0.3 \\ 0.2 & 0.6 & 0.2 \\ 0.1 & 0.1 & 0.8 \end{bmatrix}$$

L. Rabiner and B.-H. Juang, Fundamentals of Speech Recognition, Prentice-Hall, 1993

L. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” Proceedings of IEEE, vol. 77, no. 2, pp. 257-286, 1989.

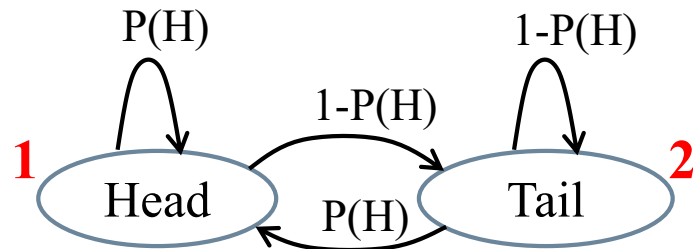
Coin-Toss Model

4

- You are in a room with a curtain through which you cannot see that is happening. On the other side of the curtain is another person who is performing a coin tossing experiment (using one or more coins). The person will not tell you which coin he selects at any time; he will only tell you the result of each coin flip.
- A typical observation sequence would be
$$\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \mathbf{o}_3 \cdots \mathbf{o}_T) = (HHTTTTHTTH \cdots H)$$
- The question is: how do we build an model to explain the observed sequence of head and tails?

Coin-Toss Model

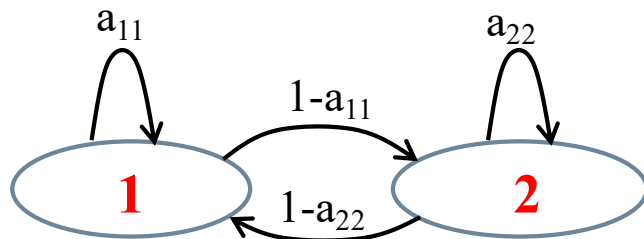
5



1-coin model
(Observable Markov Model)

$O = (HHTTTTHTTH)$

$S = (1122221221)$



$P(H) = P_1$

$P(T) = 1 - P_1$

$P(H) = P_2$

$P(T) = 1 - P_2$

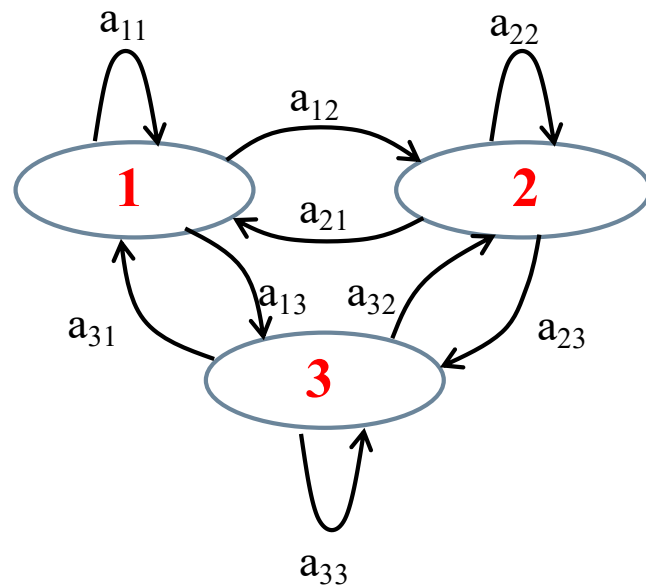
2-coins model
(Hidden Markov Model)

$O = (HHTTTTHTTH)$

$S = (2112221221)$

Coin-Toss Model

6



3-coins model
(Hidden Markov Model)

$O = (HHTTTTHTTH)$

$S = (3123311231)$

State 1

$P(H) = P_1$

$P(T) = 1 - P_1$

State 2

$P(H) = P_2$

$P(T) = 1 - P_2$

State 3

$P(H) = P_3$

$P(T) = 1 - P_3$

Elements of HMM

7

- N : the number of states in the model
- M : the number of distinct observation symbols per state
- The state-transition probability $A = \{a_{ij}\}$

$$a_{ij} = P[q_{t+1} = j | q_t = i] \quad 1 \leq i, j \leq N$$

- The observation symbol probability distribution $B = \{b_j(k)\}$

$$b_j(k) = P[\mathbf{o}_t = \mathbf{v}_k | q_t = j] \quad 1 \leq k \leq M$$

- The initial state distribution

$$\pi_i = P[q_1 = i] \quad 1 \leq i \leq N$$

- To describe an HMM, we usually use the compact notation

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

Three Basic Problems of HMM

8

□ Problem 1: Probability Evaluation

Given the observation sequence $\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \cdots \mathbf{o}_T)$, and a model $\lambda = (A, B, \pi)$, how do we efficiently compute $P(\mathbf{O}|\lambda)$, the probability of the observation sequence, given the model?

How do we compute the probability that the observed sequence was produced by the model?

Scoring how well a given model matches a given observation sequence.

Three Basic Problems of HMM

9

□ Problem 2: Optimal State Sequence

Given the observation sequence $\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \cdots \mathbf{o}_T)$, and the model λ , how do we choose a corresponding state sequence $\mathbf{q} = (q_1 q_2 \cdots q_T)$ that is optimal in some sense (i.e., best explains the observations)

Attempt to uncover the hidden part of the model – that is, to find the “correct” state sequence.

For practical situations, we usually use an optimality criterion to solve this problem as best as possible.

Three Basic Problems of HMM

10

□ **Problem 3: Parameter Estimation**

How do we adjust the model parameters $\lambda = (A, B, \pi)$ to maximize $P(\mathbf{O}|\lambda)$

Attempt to optimize the model parameters to best describe how a given observation sequence comes about.

The observation sequence used to adjust the model parameters is called a training sequence because it is used to “train” the HMM.

Solution to Problem 1

11

Given the observation sequence $\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \cdots \mathbf{o}_T)$, and a model $\lambda = (A, B, \pi)$, how do we efficiently compute $P(\mathbf{O}|\lambda)$, the probability of the observation sequence, given the model?

- There are N^T possible state sequences
- Consider one fixed-state sequence $\mathbf{q} = (q_1 q_2 \cdots q_T)$
- The prob. of the observation sequence given the sequence

$$P(\mathbf{O}|\mathbf{q}, \lambda) = \prod_{t=1}^T P(\mathbf{o}_t|q_t, \lambda)$$

- Where we have assumed statistical independence of observations, thus we get

$$P(\mathbf{O}|\mathbf{q}, \lambda) = b_{q_1}(\mathbf{o}_1) \cdot b_{q_2}(\mathbf{o}_2) \cdots b_{q_T}(\mathbf{o}_T)$$

Solution to Problem 1

12



Solution to Problem 1

13

- The prob. of such state sequence can be written as

$$P(\mathbf{q}|\lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \cdots a_{q_{T-1} q_T}$$

- The joint prob. of \mathbf{O} and \mathbf{q} , i.e., the prob. That \mathbf{O} and \mathbf{q} occur simultaneously, is simply the product of the above terms

$$P(\mathbf{O}, \mathbf{q}|\lambda) = P(\mathbf{O}|\mathbf{q}, \lambda)P(\mathbf{q}|\lambda)$$

- The prob. of \mathbf{O} is obtained by summing this joint prob. over all possible state sequences \mathbf{q} , giving

$$\begin{aligned} P(\mathbf{O}|\lambda) &= \sum_{all \mathbf{q}} P(\mathbf{O}|\mathbf{q}, \lambda)P(\mathbf{q}|\lambda) \\ &= \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} b_{q_1}(\mathbf{o}_1) a_{q_1 q_2} b_{q_2}(\mathbf{o}_2) \cdots a_{q_{T-1} q_T} b_{q_T}(\mathbf{o}_T) \end{aligned}$$

Solution to Problem 1

14

□ The Forward Procedure

- The prob. of the partial observation sequence $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_t$ (until time t) and state i at time t , given the model λ

$$\alpha_t(i) = P(\mathbf{o}_1 \mathbf{o}_2 \cdots \mathbf{o}_t, q_t = i | \lambda)$$

- We solve for it inductively

- 1. Initialization $\alpha_1(i) = \pi_i b_i(\mathbf{o}_1) \quad 1 \leq i \leq N$

- 2. Induction

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(\mathbf{o}_{t+1}) \quad \begin{array}{l} 1 \leq t \leq T-1 \\ 1 \leq j \leq N \end{array}$$

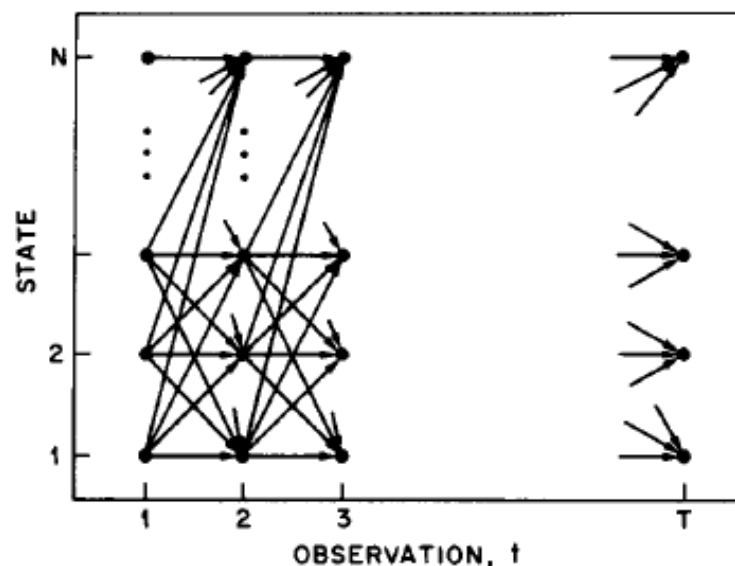
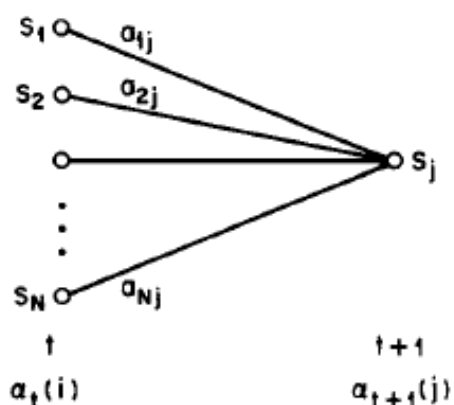
- 3. Termination

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

Solution to Problem 1

15

□ The Forward Procedure



- Require on the order of N^2T calculations, rather than $2TN^T$ as required by the direction calculation.

□ The Backward Procedure

- The prob. of partial observation sequence from $t+1$ to the end, given state i at time t and the model λ

$$\beta_t(i) = P(\mathbf{o}_{t+1}\mathbf{o}_{t+2} \cdots \mathbf{o}_T | q_t = i, \lambda)$$

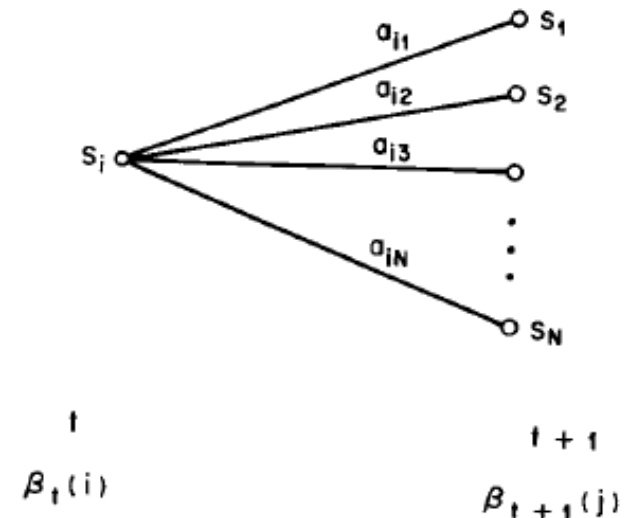
- 1. Initialization $\beta_T(i) = 1 \quad 1 \leq i \leq N$

- 2. Induction

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(\mathbf{o}_{t+1}) \beta_{t+1}(j)$$

$$t = T - 1, T - 2, \dots, 1$$

$$1 \leq i \leq N$$



Solution to Problem 2

17

Given the observation sequence $\mathbf{O} = (\mathbf{o}_1 \mathbf{o}_2 \cdots \mathbf{o}_T)$, and the model λ , how do we choose a corresponding state sequence $\mathbf{q} = (q_1 q_2 \cdots q_T)$ that is optimal in some sense (i.e., best explains the observations)

- We define the quantity

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P[q_1 q_2 \dots q_{t-1}, q_t = i, \mathbf{o}_{t+1} \mathbf{o}_{t+2} \cdots \mathbf{o}_t | \lambda]$$

- Which is the best score (highest probability) along a single path, at time t , which accounts for the first t observations and ends in state i . By induction we have

$$\delta_{t+1}(i) = [\max_j \delta_t(j) a_{ji}] \cdot b_i(\mathbf{o}_{t+1})$$

Solution to Problem 2

Solution to Problem 2

19

□ The Viterbi Algorithm

- 1. Initialization $\delta_1(i) = \pi_i b_i(\mathbf{o}_1) \quad \psi_1(i) = 0$
- 2. Recursion $\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(\mathbf{o}_t) \quad 1 \leq j \leq N$
 $\psi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] \quad 2 \leq t \leq T$
- 3. Termination $P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$
 $q_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)]$
- 4. Path (state sequence) backtracking
 $q_t^* = \psi_{t+1}(q_{t+1}^*) \quad t = T - 1, T - 2, \dots, 1$
- The major difference between Viterbi and the forward procedure is the maximization over previous states.

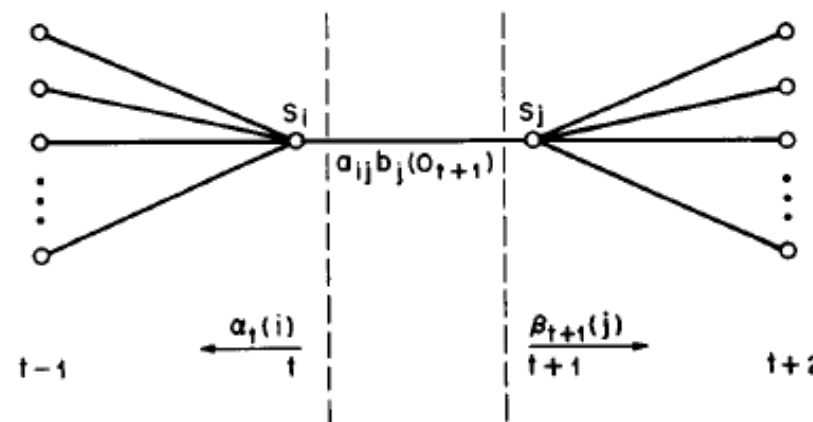
Solution to Problem 3

20

How do we adjust the model parameters $\lambda = (A, B, \pi)$ to maximize $P(\mathbf{O}|\lambda)$

- Choose $\lambda = (A, B, \pi)$ such that its likelihood, $P(\mathbf{O}|\lambda)$, is locally maximized using an iterative procedure such as the Baum-Welch algorithm (also known as EM algorithm or forward-backward algorithm)
- Define the prob. of being in state i at time t , and state j at time $t+1$, given the model and the observation sequence.

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | \mathbf{O}, \lambda)$$

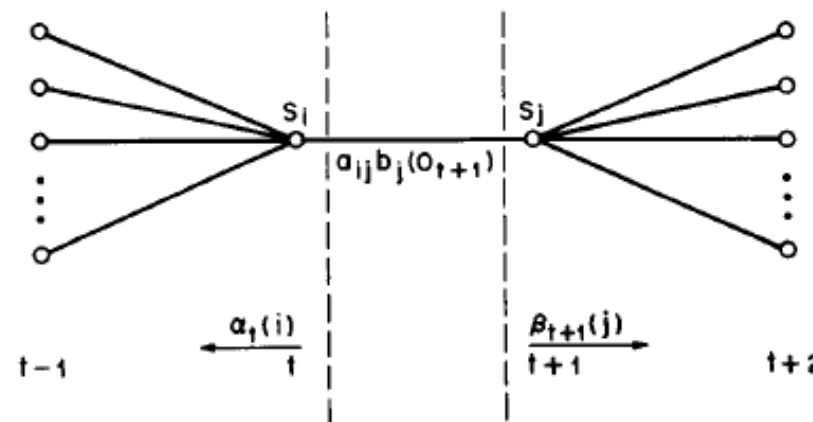


Solution to Problem 3

21

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | \mathbf{O}, \lambda)$$

$$\begin{aligned} \xi_t(i, j) &= \frac{P(q_t = i, q_{t+1} = j, \mathbf{O} | \lambda)}{P(\mathbf{O} | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(\mathbf{o}_{t+1}) \beta_{t+1}(j)}{P(\mathbf{O} | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(\mathbf{o}_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(\mathbf{o}_{t+1}) \beta_{t+1}(j)} \end{aligned}$$



$$\gamma_t(i) = P(q_t = i | \mathbf{O}, \lambda)$$

The prob. of being in state i at time t , given the observation sequence \mathbf{O} and the model

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

$$\begin{aligned} \gamma_t(i) &= P(q_t = i | \mathbf{O}, \lambda) = \frac{P(\mathbf{O}, q_t = i | \lambda)}{P(\mathbf{O} | \lambda)} \\ &= \frac{P(\mathbf{O}, q_t = i | \lambda)}{\sum_{i=1}^N P(\mathbf{O}, q_t = i | \lambda)} = \frac{\alpha_t(i) \beta_t(i)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)} \end{aligned}$$

Solution to Problem 3

22

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

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

$$\bar{\pi}_i = \text{expected frequency (number of times) in state } i \text{ at time } (t = 1) = \gamma_1(i)$$

$$\bar{a}_{ij} = \frac{\text{expected number of transitions from state } i \text{ to state } j}{\text{expected number of transitions from state } i} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

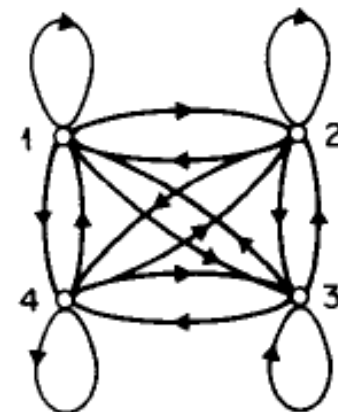
$$\bar{b}_j(k) = \frac{\text{expected number of times in state } j \text{ and observing symbol } \mathbf{v}_k}{\text{expected number of times in state } j}$$

$$= \frac{\sum_{\substack{t=1 \\ s.t. \mathbf{o}_t = \mathbf{v}_k}}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad \rightarrow \text{Usually modeled by GMM}$$

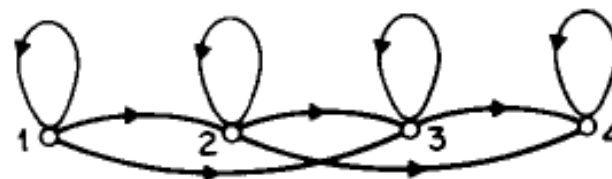
Types of HMMs

23

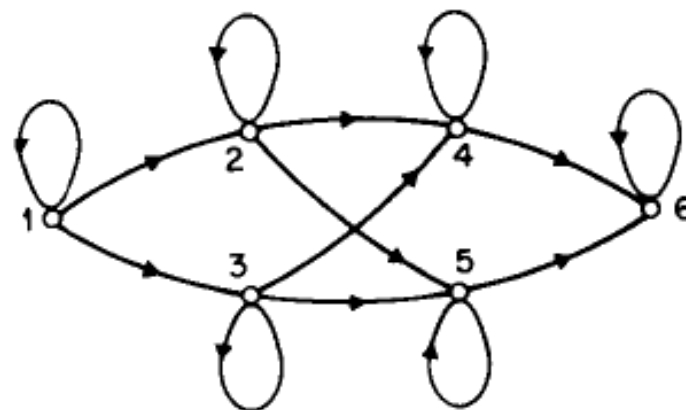
- Ergodic
- Left-right
- Parallel path left-right



(a)



(b)



(c)

Fig. 7. Illustration of 3 distinct types of HMMs. (a) A 4-state ergodic model. (b) A 4-state left-right model. (c) A 6-state parallel path left-right model.

Case Study

24

- Features
 - ▣ Field descriptor
 - ▣ Edge descriptor
 - ▣ Grass and sand
 - ▣ Player height
- For each highlight model, compute the optimal state sequence by the Viterbi algorithm.
- Check the prob. derived from the state sequence.

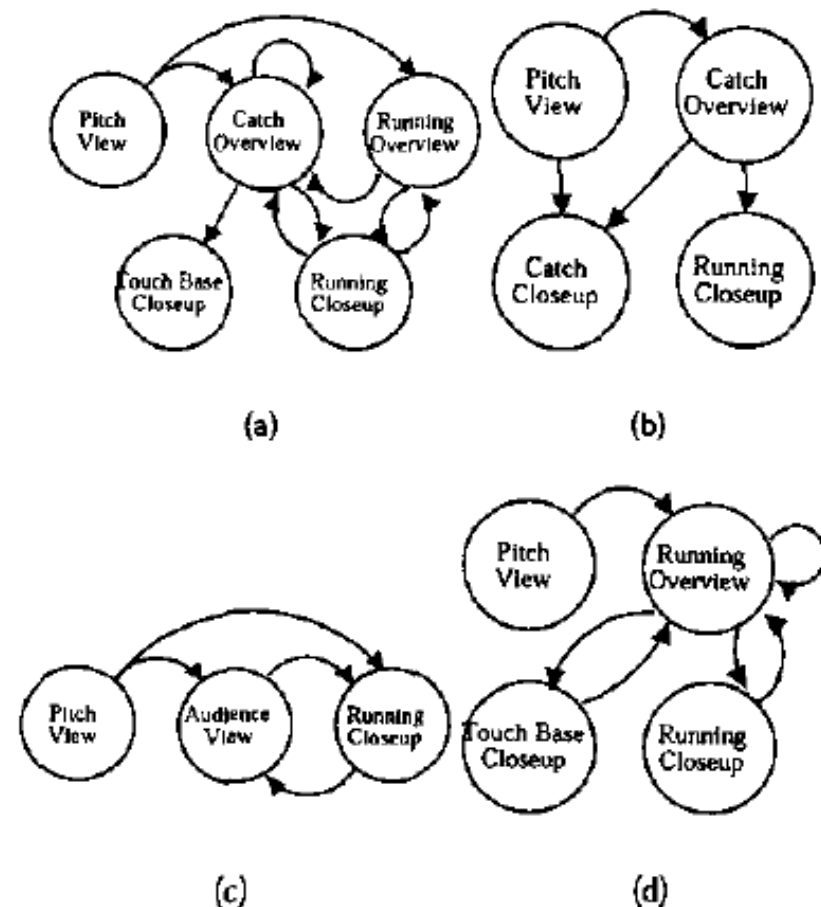


Fig. 2. (a) HMM model for nice hits (b) HMM model for nice catches (c) HMM model for home runs (d) HMM model for plays within the diamond

Peng, et al. "Extract highlights from baseball game video with hidden Markov models," In Proc. of ICIP, vol. 1, pp. 609-612, 2002.

Related Resources

25

- Hidden Markov Model (HMM) Toolbox for Matlab
 - ▣ <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
- The General Hidden Markov Model library (GHMM)
 - ▣ <http://ghmm.sourceforge.net/>
- HTK Speech Recognition Toolkit
 - ▣ <http://htk.eng.cam.ac.uk>