



# Tutorial on Hidden Markov Model

Loc Nguyen

Sunflower Soft Company, Ho Chi Minh city, Vietnam

**Email address:**

ng\_phloc@yahoo.com

**To cite this article:**

Loc Nguyen. Tutorial on Hidden Markov Model. *Applied and Computational Mathematics*. Special Issue: Some Novel Algorithms for Global Optimization and Relevant Subjects. Vol. 6, No. 4-1, 2017, pp. 16-38. doi: 10.11648/j.acm.s.2017060401.12

**Received:** September 11, 2015; **Accepted:** September 13, 2015; **Published:** June 17, 2016

---

**Abstract:** Hidden Markov model (HMM) is a powerful mathematical tool for prediction and recognition. Many computer software products implement HMM and hide its complexity, which assist scientists to use HMM for applied researches. However comprehending HMM in order to take advantages of its strong points requires a lot of efforts. This report is a tutorial on HMM with full of mathematical proofs and example, which help researchers to understand it by the fastest way from theory to practice. The report focuses on three common problems of HMM such as evaluation problem, uncovering problem, and learning problem, in which learning problem with support of optimization theory is the main subject.

**Keywords:** Hidden Markov Model, Optimization, Evaluation Problem, Uncovering Problem, Learning Problem

---

## 1. Introduction

There are many real-world phenomena (so-called states) that we would like to model in order to explain our observations. Often, given sequence of observations symbols, there is demand of discovering real states. For example, there are some states of weather: *sunny*, *cloudy*, *rainy* [1, p. 1]. Suppose you are in the room and do not know the weather outside but you are notified observations such as wind speed, atmospheric pressure, humidity, and temperature from someone else. Basing on these observations, it is possible for you to forecast the weather by using hidden Markov model (HMM). Before discussing about HMM, we should glance over the definition of Markov model (MM). First, MM is the statistical model which is used to model the stochastic process. MM is defined as below [2]:

- Given a finite set of state  $S = \{s_1, s_2, \dots, s_n\}$  whose cardinality is  $n$ . Let  $\Pi$  be the *initial state distribution* where  $\pi_i \in \Pi$  represents the probability that the stochastic process begins in state  $s_i$ . In other words  $\pi_i$  is the initial probability of state  $s_i$ , where

$$\sum_{s_i \in S} \pi_i = 1$$

- The stochastic process which is modeled gets only one state from  $S$  at all time points. This stochastic process is defined as a finite vector  $X = (x_1, x_2, \dots, x_T)$  whose element  $x_t$  is a state at time point  $t$ . The process  $X$  is called *state*

*stochastic process* and  $x_t \in S$  equals some state  $s_i \in S$ . Note that  $X$  is also called *state sequence*. Time point can be in terms of second, minute, hour, day, month, year, etc. It is easy to infer that the initial probability  $\pi_i = P(x_1 = s_i)$  where  $x_1$  is the first state of the stochastic process. The state stochastic process  $X$  must meet fully the *Markov property*, namely, given previous state  $x_{t-1}$  of process  $X$ , the conditional probability of current state  $x_t$  is only dependent on the previous state  $x_{t-1}$ , not relevant to any further past state  $(x_{t-2}, x_{t-3}, \dots, x_1)$ . In other words,  $P(x_t | x_{t-1}, x_{t-2}, x_{t-3}, \dots, x_1) = P(x_t | x_{t-1})$  with note that  $P(\cdot)$  also denotes probability in this report. Such process is called first-order Markov process.

- At time point, the process changes to the next state based on the *transition probability distribution*  $a_{ij}$ , which depends only on the previous state. So  $a_{ij}$  is the probability that the stochastic process changes current state  $s_i$  to next state  $s_j$ . It means that  $a_{ij} = P(x_t = s_j | x_{t-1} = s_i) = P(x_{t+1} = s_j | x_t = s_i)$ . The probability of transitioning from any given state to some next state is 1, we have

$$\forall s_i \in S, \sum_{s_j \in S} a_{ij} = 1$$

All transition probabilities  $a_{ij}(s)$  constitute the *transition probability matrix*  $A$ . Note that  $A$  is  $n$  by  $n$  matrix because there are  $n$  distinct states. It is easy to infer that matrix  $A$  represents state stochastic process  $X$ . It is possible to

understand that the initial probability matrix  $\Pi$  is degradation case of matrix  $A$ .

Briefly, MM is the triple  $\langle S, A, \Pi \rangle$ . In typical MM, states are observed directly by users and transition probabilities ( $A$  and  $\Pi$ ) are unique parameters. Otherwise, hidden Markov model (HMM) is similar to MM except that the underlying states become hidden from observer, they are hidden parameters. HMM adds more output parameters which are called observations. Each state (hidden parameter) has the conditional probability distribution upon such observations. HMM is responsible for discovering hidden parameters (states) from output parameters (observations), given the stochastic process. The HMM has further properties as below [2]:

- Suppose there is a finite set of possible observations  $\Phi = \{\varphi_1, \varphi_2, \dots, \varphi_m\}$  whose cardinality is  $m$ . There is the second stochastic process which produces *observations* correlating with hidden states. This process is called *observable stochastic process*, which is defined as a finite vector  $O = (o_1, o_2, \dots, o_T)$  whose element  $o_t$  is an observation at time point  $t$ . Note that  $o_t \in \Phi$  equals some  $\varphi_k$ . The process  $O$  is often known as *observation sequence*.
- There is a probability distribution of producing a given observation in each state. Let  $b_i(k)$  be the probability of observation  $\varphi_k$  when the state stochastic process is in state  $s_i$ . It means that  $b_i(k) = b_i(o_t = \varphi_k) = P(o_t = \varphi_k | x_t = s_i)$ . The sum of probabilities of all observations which observed in a certain state is 1, we have

$$\forall s_i \in S, \sum_{\theta_k \in \Phi} b_i(k) = 1$$

All probabilities of observations  $b_i(k)$  constitute the *observation probability matrix*  $B$ . It is convenient for us to use notation  $b_{ik}$  instead of notation  $b_i(k)$ . Note that  $B$  is  $n$  by  $m$  matrix because there are  $n$  distinct states and  $m$  distinct observations. While matrix  $A$  represents state stochastic process  $X$ , matrix  $B$  represents observable stochastic process  $O$ .

Thus, HMM is the 5-tuple  $\Delta = \langle S, \Phi, A, B, \Pi \rangle$ . Note that components  $S, \Phi, A, B$ , and  $\Pi$  are often called parameters of HMM in which  $A, B$ , and  $\Pi$  are essential parameters. Going back weather example, suppose you need to predict how weather tomorrow is: *sunny*, *cloudy* or *rainy* since you know only observations about the humidity: *dry*, *dryish*, *damp*, *soggy*. The HMM is totally determined based on its parameters  $S, \Phi, A, B$ , and  $\Pi$  according to weather example. We have  $S = \{s_1 = \text{sunny}, s_2 = \text{cloudy}, s_3 = \text{rainy}\}$ ,  $\Phi = \{\varphi_1 = \text{dry}, \varphi_2 = \text{dryish}, \varphi_3 = \text{damp}, \varphi_4 = \text{soggy}\}$ . Transition probability matrix  $A$  is shown in table 1.

**Table 1.** Transition probability matrix  $A$ .

		Weather current day (Time point $t$ )		
		sunny	cloudy	rainy
Weather previous day (Time point $t-1$ )	sunny	$a_{11}=0.50$	$a_{12}=0.25$	$a_{13}=0.25$
	cloudy	$a_{21}=0.30$	$a_{22}=0.40$	$a_{23}=0.30$
	rainy	$a_{31}=0.25$	$a_{32}=0.25$	$a_{33}=0.50$

From table 1, we have  $a_{11}+a_{12}+a_{13}=1$ ,  $a_{21}+a_{22}+a_{23}=1$ ,  $a_{31}+a_{32}+a_{33}=1$ .

Initial state distribution specified as uniform distribution is shown in table 2.

**Table 2.** Uniform initial state distribution  $\Pi$ .

sunny	cloudy	rainy
$\pi_1=0.33$	$\pi_2=0.33$	$\pi_3=0.33$

From table 2, we have  $\pi_1+\pi_2+\pi_3=1$ .

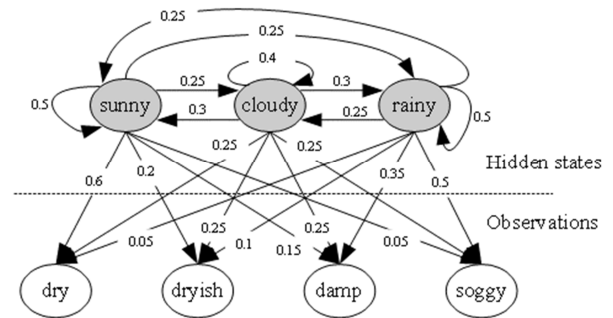
Observation probability matrix  $B$  is shown in table 3.

**Table 3.** Observation probability matrix  $B$ .

		Humidity			
		dry	dryish	damp	soggy
Weather	sunny	$b_{11}=0.60$	$b_{12}=0.20$	$b_{13}=0.15$	$b_{14}=0.05$
	cloudy	$b_{21}=0.25$	$b_{22}=0.25$	$b_{23}=0.25$	$b_{24}=0.25$
	rainy	$b_{31}=0.05$	$b_{32}=0.10$	$b_{33}=0.35$	$b_{34}=0.50$

From table 3, we have  $b_{11}+b_{12}+b_{13}+b_{14}=1$ ,  $b_{21}+b_{22}+b_{23}+b_{24}=1$ ,  $b_{31}+b_{32}+b_{33}+b_{34}=1$ .

The whole weather HMM is depicted in fig. 1.



**Figure 1.** HMM of weather forecast (hidden states are shaded).

There are three problems of HMM [2] [3, pp. 262-266]:

1. Given HMM  $\Delta$  and an observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  where  $o_t \in \Phi$ , how to calculate the probability  $P(O|\Delta)$  of this observation sequence. Such probability  $P(O|\Delta)$  indicates how much the HMM  $\Delta$  affects on sequence  $O$ . This is *evaluation problem* or *explanation problem*. Note that it is possible to denote  $O = \{o_1 \rightarrow o_2 \rightarrow \dots \rightarrow o_T\}$  and the sequence  $O$  is aforementioned observable stochastic process.
2. Given HMM  $\Delta$  and an observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  where  $o_t \in \Phi$ , how to find the sequence of states  $X = \{x_1, x_2, \dots, x_T\}$  where  $x_t \in S$  so that  $X$  is most likely to have produced the observation sequence  $O$ . This is *uncovering problem*. Note that the sequence  $X$  is aforementioned state stochastic process.
3. Given HMM  $\Delta$  and an observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  where  $o_t \in \Phi$ , how to adjust parameters of  $\Delta$  such as initial state distribution  $\Pi$ , transition probability matrix  $A$ , and observation probability matrix  $B$  so that the quality of HMM  $\Delta$  is enhanced. This is *learning problem*.

These problems will be mentioned in sections 2, 3, and 4, in turn.

## 2. HMM Evaluation Problem

The essence of evaluation problem is to find out the way to compute the probability  $P(O|\Delta)$  most effectively given the observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ . For example, given HMM  $\Delta$  whose parameters  $A, B$ , and  $\Pi$  specified in tables 1, 2, and 3, which is designed for weather forecast. Suppose we need to calculate the probability of event that humidity is *soggy* and *dry* in days 1 and 2, respectively. This is evaluation problem with sequence of observations  $O = \{o_1=\phi_4=soggy, o_2=\phi_1=dry, o_3=\phi_2=dryish\}$ . There is a complete set of  $3^3=27$  mutually exclusive cases of weather states for three days:  $\{x_1=s_1=sunny, x_2=s_1=sunny, x_3=s_1=sunny\}, \{x_1=s_1=sunny, x_2=s_1=sunny, x_3=s_2=cloudy\}, \{x_1=s_1=sunny, x_2=s_1=sunny, x_3=s_3=rainy\}, \{x_1=s_1=sunny, x_2=s_2=cloudy, x_3=s_1=sunny\}, \{x_1=s_1=sunny, x_2=s_2=cloudy, x_3=s_2=cloudy\}, \{x_1=s_1=sunny, x_2=s_2=cloudy, x_3=s_3=rainy\}, \{x_1=s_1=sunny, x_2=s_3=rainy, x_3=s_1=sunny\}, \{x_1=s_1=sunny, x_2=s_3=rainy, x_3=s_2=cloudy\}, \{x_1=s_1=sunny, x_2=s_3=rainy, x_3=s_3=rainy\}, \{x_1=s_2=cloudy, x_2=s_1=sunny, x_3=s_1=sunny\}, \{x_1=s_2=cloudy, x_2=s_1=sunny, x_3=s_2=cloudy\}, \{x_1=s_2=cloudy, x_2=s_1=sunny, x_3=s_3=rainy\}, \{x_1=s_2=cloudy, x_2=s_2=cloudy, x_3=s_1=sunny\}, \{x_1=s_2=cloudy, x_2=s_2=cloudy, x_3=s_2=cloudy\}, \{x_1=s_2=cloudy, x_2=s_2=cloudy, x_3=s_3=rainy\}, \{x_1=s_2=cloudy, x_2=s_3=rainy, x_3=s_1=sunny\}, \{x_1=s_2=cloudy, x_2=s_3=rainy, x_3=s_2=cloudy\}, \{x_1=s_2=cloudy, x_2=s_3=rainy, x_3=s_3=rainy\}, \{x_1=s_3=rainy, x_2=s_1=sunny, x_3=s_2=cloudy\}, \{x_1=s_3=rainy, x_2=s_1=sunny, x_3=s_3=rainy\}, \{x_1=s_3=rainy, x_2=s_2=cloudy, x_3=s_1=sunny\}, \{x_1=s_3=rainy, x_2=s_2=cloudy, x_3=s_2=cloudy\}, \{x_1=s_3=rainy, x_2=s_2=cloudy, x_3=s_3=rainy\}, \{x_1=s_3=rainy, x_2=s_3=rainy, x_3=s_1=sunny\}, \{x_1=s_3=rainy, x_2=s_3=rainy, x_3=s_2=cloudy\}, \{x_1=s_3=rainy, x_2=s_3=rainy, x_3=s_3=rainy\}$ .

According to total probability rule [4, p. 101], the probability  $P(O|\Delta)$  is:

$$\begin{aligned}
 P(O|\Delta) &= P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_2) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_3) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_2) \\
 &\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_3) \\
 &\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_2) \\
 &\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_3) \\
 &\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_2) \\
 &\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_2)
 \end{aligned}$$

$$\begin{aligned}
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_3) \\
 &\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_1) \\
 &\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_2) \\
 &\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_3) \\
 &\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_1) \\
 &\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_2) \\
 &\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
 &\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_2) \\
 &\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_3) \\
 &\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_1) \\
 &\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_2) \\
 &\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_3) \\
 &\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_3) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_1) \\
 &\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_1) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_2) \\
 &\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_2) \\
 &+ P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_3) \\
 &\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_3)
 \end{aligned}$$

We have:

$$\begin{aligned}
 &P(o_1 = \phi_4, o_2 = \phi_1, o_3 = \phi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &= P(o_1 = \phi_4 | x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(o_2 = \phi_1 | x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(o_3 = \phi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad \text{(Because observations } o_1, o_2, \text{ and } o_3 \text{ are mutually independent)} \\
 &= P(o_1 = \phi_4 | x_1 = s_1) * P(o_2 = \phi_1 | x_2 = s_1) \\
 &\quad * P(o_3 = \phi_2 | x_3 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_1) \\
 &\quad \text{(Because an observation is only dependent on the day when it is observed)} \\
 &= P(o_1 = \phi_4 | x_1 = s_1) * P(o_2 = \phi_1 | x_2 = s_1) \\
 &\quad * P(o_3 = \phi_2 | x_3 = s_1) \\
 &\quad * P(x_3 = s_1 | x_1 = s_1, x_2 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1) \\
 &\quad \text{(Due to multiplication rule [4, p. 100])} \\
 &= P(o_1 = \phi_4 | x_1 = s_1) * P(o_2 = \phi_1 | x_2 = s_1) \\
 &\quad * P(o_3 = \phi_2 | x_3 = s_1) \\
 &\quad * P(x_3 = s_1 | x_2 = s_1) \\
 &\quad * P(x_1 = s_1, x_2 = s_1) \\
 &\quad \text{(Due to Markov property, current state is only dependent on right previous state)}
 \end{aligned}$$

$$\begin{aligned}
&= P(o_1 = \varphi_4 | x_1 = s_1) * P(o_2 = \varphi_1 | x_2 = s_1) \\
&\quad * P(o_3 = \varphi_2 | x_3 = s_1) \\
&\quad * P(x_3 = s_1 | x_2 = s_1) \\
&\quad * P(x_2 = s_1 | x_1 = s_1) * P(x_1 = s_1) \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= b_{14}b_{11}b_{12}a_{11}a_{11}\pi_1 \\
&\text{(According to parameters A, B, and } \Pi \text{ specified in tables 1, 2, and 3)}
\end{aligned}$$

Similarly, we have:

$$\begin{aligned}
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_2) \\
&\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_2) \\
&= b_{14}b_{11}b_{22}a_{12}a_{11}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_1, x_3 = s_3) \\
&\quad * P(x_1 = s_1, x_2 = s_1, x_3 = s_3) \\
&= b_{14}b_{11}b_{32}a_{13}a_{11}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_1) \\
&\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_1) \\
&= b_{14}b_{21}b_{12}a_{21}a_{12}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_2) \\
&\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_2) \\
&= b_{14}b_{21}b_{22}a_{22}a_{12}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_2, x_3 = s_3) \\
&\quad * P(x_1 = s_1, x_2 = s_2, x_3 = s_3) \\
&= b_{14}b_{21}b_{32}a_{23}a_{12}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_1) \\
&\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_1) \\
&= b_{14}b_{31}b_{12}a_{31}a_{13}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_2) \\
&\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_2) \\
&= b_{14}b_{31}b_{22}a_{32}a_{13}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_1, x_2 = s_3, x_3 = s_3) \\
&\quad * P(x_1 = s_1, x_2 = s_3, x_3 = s_3) \\
&= b_{14}b_{31}b_{32}a_{33}a_{13}\pi_1 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_1) \\
&\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_1) \\
&= b_{24}b_{11}b_{12}a_{11}a_{21}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_2) \\
&\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_2) \\
&= b_{24}b_{11}b_{22}a_{12}a_{21}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_1, x_3 = s_3) \\
&\quad * P(x_1 = s_2, x_2 = s_1, x_3 = s_3) \\
&= b_{24}b_{11}b_{32}a_{13}a_{21}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_1) \\
&\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_1) \\
&= b_{24}b_{21}b_{12}a_{21}a_{22}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_2) \\
&\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_2) \\
&= b_{24}b_{21}b_{22}a_{22}a_{22}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_2, x_3 = s_3) \\
&\quad * P(x_1 = s_2, x_2 = s_2, x_3 = s_3) \\
&= b_{24}b_{21}b_{32}a_{23}a_{22}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_1) \\
&\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_1) \\
&= b_{24}b_{31}b_{12}a_{31}a_{23}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_2) \\
&\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_2) \\
&= b_{24}b_{31}b_{22}a_{32}a_{23}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
&\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
&= b_{24}b_{31}b_{32}a_{33}a_{23}\pi_2
\end{aligned}$$

$$\begin{aligned}
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
&\quad * P(x_1 = s_2, x_2 = s_3, x_3 = s_3) \\
&= b_{24}b_{31}b_{32}a_{33}a_{23}\pi_2 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_1) \\
&\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_1) \\
&= b_{34}b_{11}b_{12}a_{11}a_{31}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_2) \\
&\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_2) \\
&= b_{34}b_{11}b_{22}a_{12}a_{31}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_1, x_3 = s_3) \\
&\quad * P(x_1 = s_3, x_2 = s_1, x_3 = s_3) \\
&= b_{34}b_{11}b_{32}a_{13}a_{31}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_1) \\
&\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_1) \\
&= b_{34}b_{21}b_{12}a_{21}a_{32}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_2) \\
&\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_2) \\
&= b_{34}b_{21}b_{22}a_{22}a_{32}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_2, x_3 = s_3) \\
&\quad * P(x_1 = s_3, x_2 = s_2, x_3 = s_3) \\
&= b_{34}b_{21}b_{32}a_{23}a_{32}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_1) \\
&\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_1) \\
&= b_{34}b_{31}b_{12}a_{31}a_{33}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_2) \\
&\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_2) \\
&= b_{34}b_{31}b_{22}a_{32}a_{33}\pi_3 \\
P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2 | x_1 = s_3, x_2 = s_3, x_3 = s_3) \\
&\quad * P(x_1 = s_3, x_2 = s_3, x_3 = s_3) \\
&= b_{34}b_{31}b_{32}a_{33}a_{33}\pi_3
\end{aligned}$$

It implies

$$\begin{aligned}
P(O|\Delta) &= P(o_1 = \varphi_4, o_2 = \varphi_1, o_3 = \varphi_2) \\
&= b_{14}b_{11}b_{12}a_{11}a_{11}\pi_1 + b_{14}b_{11}b_{22}a_{12}a_{11}\pi_1 \\
&\quad + b_{14}b_{11}b_{32}a_{13}a_{11}\pi_1 \\
&\quad + b_{14}b_{21}b_{12}a_{21}a_{12}\pi_1 + b_{14}b_{21}b_{22}a_{22}a_{12}\pi_1 \\
&\quad + b_{14}b_{21}b_{32}a_{23}a_{12}\pi_1 \\
&\quad + b_{14}b_{31}b_{12}a_{31}a_{13}\pi_1 + b_{14}b_{31}b_{22}a_{32}a_{13}\pi_1 \\
&\quad + b_{14}b_{31}b_{32}a_{33}a_{13}\pi_1 \\
&\quad + b_{24}b_{11}b_{12}a_{11}a_{21}\pi_2 + b_{24}b_{11}b_{22}a_{12}a_{21}\pi_2 \\
&\quad + b_{24}b_{11}b_{32}a_{13}a_{21}\pi_2 \\
&\quad + b_{24}b_{21}b_{12}a_{21}a_{22}\pi_2 + b_{24}b_{21}b_{22}a_{22}a_{22}\pi_2 \\
&\quad + b_{24}b_{21}b_{32}a_{23}a_{22}\pi_2 \\
&\quad + b_{24}b_{31}b_{12}a_{31}a_{23}\pi_2 + b_{24}b_{31}b_{22}a_{32}a_{23}\pi_2 \\
&\quad + b_{24}b_{31}b_{32}a_{33}a_{23}\pi_2 \\
&\quad + b_{34}b_{11}b_{12}a_{11}a_{31}\pi_3 + b_{34}b_{11}b_{22}a_{12}a_{31}\pi_3 \\
&\quad + b_{34}b_{11}b_{32}a_{13}a_{31}\pi_3 \\
&\quad + b_{34}b_{21}b_{12}a_{21}a_{32}\pi_3 + b_{34}b_{21}b_{22}a_{22}a_{32}\pi_3 \\
&\quad + b_{34}b_{21}b_{32}a_{23}a_{32}\pi_3 \\
&\quad + b_{34}b_{31}b_{12}a_{31}a_{33}\pi_3 + b_{34}b_{31}b_{22}a_{32}a_{33}\pi_3 \\
&\quad + b_{34}b_{31}b_{32}a_{33}a_{33}\pi_3 \\
&= 0.012980859375
\end{aligned}$$

It is easy to explain that given weather HMM modeled by parameters  $A$ ,  $B$ , and  $\Pi$  specified in tables 1, 2, and 3, the event that it is *soggy*, *dry*, and *dryish* in three successive days is rare because the probability of such event  $P(O|\Delta)$  is low ( $\approx 1.3\%$ ). It is easy to recognize that it is impossible to browse all combinational cases of given observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  as we knew that it is necessary to survey  $3^3=27$

mutually exclusive cases of weather states with a tiny number of observations  $\{soggy, dry, dryish\}$ . Exactly, given  $n$  states and  $T$  observations, it takes extremely expensive cost to survey  $n^T$  cases. According to [3, pp. 262-263], there is a so-called *forward-backward procedure* to decrease computational cost for determining the probability  $P(O|\Delta)$ . Let  $\alpha_t(i)$  be the joint probability of partial observation sequence  $\{o_1, o_2, \dots, o_t\}$  and state  $x_t = s_i$  where  $1 \leq t \leq T$ , specified by (1).

$$\alpha_t(i) = P(o_1, o_2, \dots, o_t, x_t = s_i | \Delta) \quad (1)$$

The joint probability  $\alpha_t(i)$  is also called *forward variable* at time point  $t$  and state  $s_i$ . The product  $\alpha_t(i)a_{ij}$  where  $a_{ij}$  is the transition probability from state  $i$  to state  $j$  counts for probability of join event that partial observation sequence  $\{o_1, o_2, \dots, o_t\}$  exists and the state  $s_i$  at time point  $t$  is changed to  $s_j$  at time point  $t+1$ .

$$\begin{aligned} \alpha_t(i)a_{ij} &= P(o_1, o_2, \dots, o_t, x_t = s_i | \Delta) P(x_{t+1} = s_j | x_t = s_i) \\ &= P(o_1, o_2, \dots, o_t | x_t = s_i) P(x_t = s_i) P(x_{t+1} = s_j | x_t = s_i) \\ &\quad (\text{Due to multiplication rule [4, p. 100]}) \\ &= P(o_1, o_2, \dots, o_t | x_t = s_i) P(x_{t+1} = s_j | x_t = s_i) P(x_t = s_i) \\ &= P(o_1, o_2, \dots, o_t, x_{t+1} = s_j | x_t = s_i) P(x_t = s_i) \\ &\quad (\text{Because the partial observation sequence } \{o_1, o_2, \dots, o_t\} \text{ is} \\ &\quad \text{independent from next state } x_{t+1} \text{ given current state } x_t) \\ &= P(o_1, o_2, \dots, o_t, x_t = s_i, x_{t+1} = s_j) \end{aligned}$$

(Due to multiplication rule [4, p. 100])

Summing product  $\alpha_t(i)a_{ij}$  over all  $n$  possible states of  $x_t$  produces probability of join event that partial observation sequence  $\{o_1, o_2, \dots, o_t\}$  exists and the next state is  $x_{t+1} = s_j$  regardless of the state  $x_t$ .

$$\begin{aligned} \sum_{i=1}^n \alpha_t(i)a_{ij} &= \sum_{i=1}^n P(o_1, o_2, \dots, o_t, x_t = s_i, x_{t+1} = s_j) \\ &= P(o_1, o_2, \dots, o_t, x_{t+1} = s_j) \end{aligned}$$

The forward variable at time point  $t+1$  and state  $s_j$  is calculated on  $\alpha_t(i)$  as follows:

$$\begin{aligned} \alpha_{t+1}(j) &= P(o_1, o_2, \dots, o_t, o_{t+1}, x_{t+1} = s_j | \Delta) \\ &= P(o_{t+1} | o_1, o_2, \dots, o_t, x_{t+1} = s_j) P(o_1, o_2, \dots, o_t, x_{t+1} = s_j) \\ &\quad (\text{Due to multiplication rule}) \\ &= P(o_{t+1} | x_{t+1} = s_j) P(o_1, o_2, \dots, o_t, x_{t+1} = s_j) \\ &\quad (\text{Due to observations are mutually independent}) \\ &= b_j(o_{t+1}) \sum_{i=1}^n \alpha_t(i)a_{ij} \end{aligned}$$

Where  $b_j(o_{t+1})$  is the probability of observation  $o_{t+1}$  when the state stochastic process is in state  $s_j$ , please see an example of observation probability matrix shown in table 3. In brief, please pay attention to recurrence property of forward variable specified by (2).

$$\alpha_{t+1}(j) = \left( \sum_{i=1}^n \alpha_t(i)a_{ij} \right) b_j(o_{t+1}) \quad (2)$$

The aforementioned construction of forward recurrence equation (2) is essentially to build up Markov chain, illustrated by fig. 2 [3, p. 262].

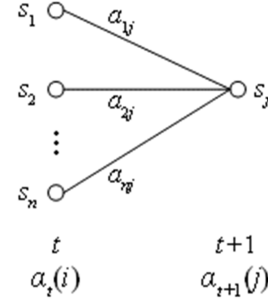


Figure 2. Construction of recurrence formula for forward variable.

According to the forward recurrence equation (2), given observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , we have:

$$\alpha_T(i) = P(o_1, o_2, \dots, o_T, x_T = s_i | \Delta)$$

The probability  $P(O|\Delta)$  is sum of  $\alpha_T(i)$  over all  $n$  possible states of  $x_T$ , specified by (3).

$$\begin{aligned} P(O|\Delta) &= P(o_1, o_2, \dots, o_T) = \\ &= \sum_{i=1}^n P(o_1, o_2, \dots, o_T, x_T = s_i | \Delta) = \sum_{i=1}^n \alpha_T(i) \end{aligned} \quad (3)$$

The forward-backward procedure to calculate the probability  $P(O|\Delta)$ , based on forward equations (2) and (3), includes three steps as shown in table 4 [3, p. 262].

Table 4. Forward-backward procedure based on forward variable to calculate the probability  $P(O|\Delta)$ .

1. Initialization step: Initializing  $\alpha_1(i) = b_i(o_1)\pi_i$  for all  $1 \leq i \leq n$
2. Recurrence step: Calculating all  $\alpha_{t+1}(j)$  for all  $1 \leq j \leq n$  and  $1 \leq t \leq T-1$  according to (2).

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

3. Evaluation step: Calculating the probability  $P(O|\Delta) = \sum_{i=1}^n \alpha_T(i)$

It is required to execute  $n+2n^2(T-1)+n-1 = 2n^2(T-1)+2n-1$  operations for forward-backward procedure based on forward variable due to:

- There are  $n$  multiplications at initialization step.
- There are  $n$  multiplications,  $n-1$  additions, and 1 multiplication over the expression  $\left( \sum_{i=1}^n \alpha_t(i)a_{ij} \right) b_j(o_{t+1})$ . There are  $n$  cases of values  $\alpha_{t+1}(j)$  for all  $1 \leq j \leq n$  at time point  $t+1$ . So, there are  $(n+n-1+1)n = 2n^2$  operations over values  $\alpha_{t+1}(j)$  for all  $1 \leq j \leq n$  at time point  $t$ . The recurrence step runs over  $T-1$  times and so, there are  $2n^2(T-1)$  operations at recurrence step.
- There are  $n-1$  additions at evaluation step.

Inside  $2n^2(T-1)+2n-1$  operations, there are  $n+(n+1)n(T-1) = n+(n^2+n)(T-1)$  multiplications and  $(n-1)n(T-1)+n-1 = (n^2+n)(T-1)+n-1$  additions.

Going back example with weather HMM whose parameters  $A$ ,  $B$ , and  $\Pi$  are specified in tables 1, 2, and 3. We need to re-calculate the probability of observation sequence  $O = \{o_1=\varphi_4=soggy, o_2=\varphi_1=dry, o_3=\varphi_2=dryish\}$  by forward-backward procedure shown in table 4 (based on forward variable). According to initialization step of forward-backward procedure based on forward variable, we have:

$$\alpha_1(1) = b_1(o_1 = \varphi_4)\pi_1 = b_{14}\pi_1 = 0.0165$$

$$\alpha_1(2) = b_2(o_1 = \varphi_4)\pi_2 = b_{24}\pi_2 = 0.0825$$

$$\alpha_1(3) = b_3(o_1 = \varphi_4)\pi_3 = b_{34}\pi_3 = 0.165$$

According to recurrence step of forward-backward procedure based on forward variable, we have:

$$\alpha_2(1) = \left( \sum_{i=1}^3 \alpha_1(i)a_{i1} \right) b_1(o_2 = \varphi_1) = \left( \sum_{i=1}^3 \alpha_1(i)a_{i1} \right) b_{11} = 0.04455$$

$$\alpha_2(2) = \left( \sum_{i=1}^3 \alpha_1(i)a_{i2} \right) b_2(o_2 = \varphi_1) = 0.01959375$$

$$\alpha_2(3) = \left( \sum_{i=1}^3 \alpha_1(i)a_{i3} \right) b_3(o_2 = \varphi_1) = 0.00556875$$

$$\alpha_3(1) = \left( \sum_{i=1}^3 \alpha_2(i)a_{i1} \right) b_1(o_3 = \varphi_2) = 0.0059090625$$

$$\alpha_3(2) = \left( \sum_{i=1}^3 \alpha_2(i)a_{i2} \right) b_2(o_3 = \varphi_2) = 0.005091796875$$

$$\alpha_3(3) = \left( \sum_{i=1}^3 \alpha_2(i)a_{i3} \right) b_3(o_3 = \varphi_2) = 0.00198$$

According to evaluation step of forward-backward procedure based on forward variable, the probability of observation sequence  $O = \{o_1=s_4=soggy, o_2=s_1=dry, o_3=s_2=dryish\}$  is:

$$P(O|\Delta) = \alpha_3(1) + \alpha_3(2) + \alpha_3(3) = 0.012980859375$$

The result from the forward-backward procedure based on forward variable is the same to the one from aforementioned brute-force method that browses all  $3^3=27$  mutually exclusive cases of weather states.

There is interesting thing that the forward-backward procedure can be implemented based on so-called *backward variable*. Let  $\beta_t(i)$  be the backward variable which is conditional probability of partial observation sequence  $\{o_t, o_{t+1}, \dots, o_T\}$  given state  $x_t=s_i$  where  $1 \leq t \leq T$ , specified by (4).

$$\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | x_t = s_i, \Delta) \quad (4)$$

We have

$$\begin{aligned} a_{ij}b_j(o_{t+1})\beta_{t+1}(j) &= P(x_{t+1} = s_j | x_t = s_i) \\ &\quad * P(o_{t+1} | x_{t+1} = s_j) \\ &\quad * P(o_{t+2}, o_{t+3}, \dots, o_T | x_{t+1} = s_j, \Delta) \\ &= P(x_{t+1} = s_j | x_t = s_i) \\ &\quad * P(o_{t+1}, o_{t+2}, o_{t+3}, \dots, o_T | x_{t+1} = s_j, \Delta) \\ &\quad \text{(Because observations } o_{t+1}, o_{t+2}, \dots, o_T \text{ are mutually independent)} \\ &= P(x_{t+1} = s_j | x_t = s_i) \\ &\quad * P(o_{t+1}, o_{t+2}, o_{t+3}, \dots, o_T | x_t = s_i, x_{t+1} = s_j, \Delta) \\ &\quad \text{(Because partial observation sequence } o_{t+1}, o_{t+2}, \dots, o_T \text{ is independent from state } x_t \text{ at time point } t) \\ &= P(o_{t+1}, o_{t+2}, o_{t+3}, \dots, o_T, x_{t+1} = s_j | x_t = s_i, \Delta) \\ &\quad \text{(Due to multiplication rule [4, p. 100])} \end{aligned}$$

Summing the product  $a_{ij}b_j(o_{t+1})\beta_{t+1}(j)$  over all  $n$  possible states of  $x_{t+1}=s_j$ , we have:

$$\begin{aligned} &\sum_{j=1}^n a_{ij}b_j(o_{t+1})\beta_{t+1}(j) \\ &= \sum_{j=1}^n P(o_{t+1}, o_{t+2}, o_{t+3}, \dots, o_T, x_{t+1} = s_j | x_t = s_i, \Delta) \\ &= P(o_{t+1}, o_{t+2}, o_{t+3}, \dots, o_T | x_t = s_i, \Delta) \\ &\quad \text{(Due to the total probability rule [4, p. 101])} \\ &= \beta_t(i) \end{aligned}$$

In brief, the recurrence property of backward variable specified by (5).

$$\beta_t(i) = \sum_{j=1}^n a_{ij}b_j(o_{t+1})\beta_{t+1}(j) \quad (5)$$

Where  $b_j(o_{t+1})$  is the probability of observation  $o_{t+1}$  when the state stochastic process is in state  $s_j$ , please see an example of observation probability matrix shown in table 3. The construction of backward recurrence equation (5) is essentially to build up Markov chain, illustrated by fig. 3 [3, p. 263].

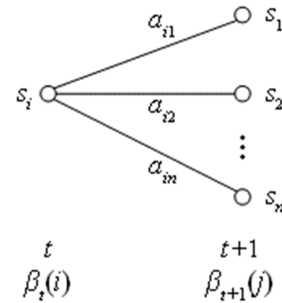


Figure 3. Construction of recurrence equation for backward variable.

According to the backward recurrence equation (5), given observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , we have:

$$\beta_1(i) = P(o_2, o_3, \dots, o_T | x_1 = s_i, \Delta)$$

The product  $\pi_i b_i(o_1)\beta_1(i)$  is:

$$\begin{aligned} \pi_i b_i(o_1)\beta_1(i) &= P(x_1 = s_i)P(o_1 | x_1 = s_i)P(o_2, o_3, \dots, o_T | x_1 = s_i, \Delta) \\ &= P(x_1 = s_i)P(o_1, o_2, o_3, \dots, o_T | x_1 = s_i, \Delta) \\ &\quad \text{(Because observations } o_1, o_2, \dots, o_T \text{ are mutually independent)} \\ &= P(o_1, o_2, o_3, \dots, o_T, x_1 = s_i | \Delta) \end{aligned}$$

It implies that the probability  $P(O|\Delta)$  is:

$$\begin{aligned} P(O|\Delta) &= P(o_1, o_2, \dots, o_T) \\ &= \sum_{i=1}^n P(o_1, o_2, \dots, o_T, x_1 = s_i | \Delta) \\ &\quad \text{(Due to the total probability rule [4, p. 101])} \end{aligned}$$

$$= \sum_{i=1}^n \pi_i b_i(o_1)\beta_1(i)$$

Shortly, the probability  $P(O|\Delta)$  is sum of product  $\pi_i b_i(o_1)\beta_1(i)$  over all  $n$  possible states of  $x_1=s_i$ , specified by (6).

$$P(O|\Delta) = \sum_{i=1}^n \pi_i b_i(o_1)\beta_1(i) \quad (6)$$

The forward-backward procedure to calculate the probability  $P(O|\Delta)$ , based on backward equations (5) and (6), includes three steps as shown in table 5 [3, p. 263].

**Table 5.** Forward-backward procedure based on backward variable to calculate the probability  $P(O|\Delta)$ .

1. Initialization step: Initializing $\beta_t(i) = 1$ for all $1 \leq i \leq n$
2. Recurrence step: Calculating all $\beta_t(i)$ for all $1 \leq i \leq n$ and $t=T-1, t=T-2, \dots, t=1$ , according to (5).
3. Evaluation step: Calculating the probability $P(O \Delta)$ according to (6), $P(O \Delta) = \sum_{i=1}^n \pi_i b_i(o_1) \beta_1(i)$

It is required to execute  $(3n-1)n(T-1)+2n+n-1 = 3n^2(T-1)-n(T-4)-1$  operations for forward-backward procedure based on forward variable due to:

- There are  $2n$  multiplications and  $n-1$  additions over the sum  $\sum_{j=1}^n a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$ . So, there are  $(2n+n-1)n = (3n-1)n$  operations over values  $\beta_t(i)$  for all  $1 \leq i \leq n$  at time point  $t$ . The recurrence step runs over  $T-1$  times and so, there are  $(3n-1)n(T-1)$  operations at recurrence step.
- There are  $2n$  multiplications and  $n-1$  additions over the sum  $\sum_{i=1}^n \pi_i b_i(o_1) \beta_1(i)$  at evaluation step.

Inside  $3n^2(T-1)-n(T-4)-1$  operations, there are  $2n^2(T-1)+2n$  multiplications and  $(n-1)n(T-1)+n-1 = n^2(T-1)-n(T-2)-1$  additions.

Going back example with weather HMM whose parameters  $A$ ,  $B$ , and  $\Pi$  are specified in tables 1, 2, and 3. We need to re-calculate the probability of observation sequence  $O = \{o_1=\varphi_4=\text{soggy}, o_2=\varphi_1=\text{dry}, o_3=\varphi_2=\text{dryish}\}$  by forward-backward procedure shown in table 5 (based on backward variable). According to initialization step of forward-backward procedure based on backward variable, we have:

$$\beta_3(1) = \beta_3(2) = \beta_3(3) = 1$$

According to recurrence step of forward-backward procedure based on backward variable, we have:

$$\beta_2(1) = \sum_{j=1}^n a_{1j} b_j(o_3 = \varphi_2) \beta_3(j) = \sum_{j=1}^n a_{1j} b_{j2} \beta_3(j) = 0.1875$$

$$\beta_2(2) = \sum_{j=1}^n a_{2j} b_j(o_3 = \varphi_2) \beta_3(j) = 0.19$$

$$\beta_2(3) = \sum_{j=1}^n a_{3j} b_j(o_3 = \varphi_2) \beta_3(j) = 0.1625$$

$$\beta_1(1) = \sum_{j=1}^n a_{1j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.07015625$$

$$\beta_1(2) = \sum_{j=1}^n a_{2j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.0551875$$

$$\beta_1(3) = \sum_{j=1}^n a_{3j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.0440625$$

According to evaluation step of forward-backward procedure based on backward variable, the probability of observation sequence  $O = \{o_1=\varphi_4=\text{soggy}, o_2=\varphi_1=\text{dry}, o_3=\varphi_2=\text{dryish}\}$  is:

$$P(O|\Delta) = \sum_{i=1}^3 \pi_i b_i(o_1 = \varphi_4) \beta_1(i) = \sum_{i=1}^3 \pi_i b_{i4} \beta_1(i) = 0.012980859375$$

The result from the forward-backward procedure based on backward variable is the same to the one from aforementioned brute-force method that browses all  $3^3=27$  mutually exclusive cases of weather states and the one from forward-backward procedure based on forward variable.

The evaluation problem is now described thoroughly in this section. The uncovering problem is mentioned particularly in successive section.

### 3. HMM Uncovering Problem

Recall that given HMM  $\Delta$  and observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  where  $o_t \in \Phi$ , how to find out a state sequence  $X = \{x_1, x_2, \dots, x_T\}$  where  $x_t \in S$  so that  $X$  is most likely to have produced the observation sequence  $O$ . This is the uncovering problem: which sequence of state transitions is most likely to have led to given observation sequence. In other words, it is required to establish an *optimal criterion* so that the state sequence  $X$  leads to maximizing such criterion. The simple criterion is the conditional probability of sequence  $X$  with respect to sequence  $O$  and model  $\Delta$ , denoted  $P(X|O, \Delta)$ . We can apply brute-force strategy: “go through all possible such  $X$  and pick the one leading to maximizing the criterion  $P(X|O, \Delta)$ ”.

$$X = \underset{X}{\operatorname{argmax}} (P(X|O, \Delta))$$

This strategy is impossible if the number of states and observations is huge. Another popular way is to establish a so-called *individually optimal criterion* [3, p. 263] which is described right later.

Let  $\gamma_t(i)$  be joint probability that the stochastic process is in state  $s_i$  at time point  $t$  with observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , equation (7) specifies this probability based on forward variable  $\alpha_t$  and backward variable  $\beta_t$ .

$$\gamma_t(i) = P(o_1, o_2, \dots, o_T, x_t = s_i | \Delta) = \alpha_t(i) \beta_t(i) \quad (7)$$

The variable  $\gamma_t(i)$  is also called *individually optimal criterion* with note that forward variable  $\alpha_t$  and backward variable  $\beta_t$  are calculated according to (2) and (5), respectively.

Following is proof of (7).

$$\begin{aligned} \gamma_t(i) &= P(o_1, o_2, \dots, o_T, x_t = s_i | \Delta) \\ &= P(x_t = s_i, o_1, o_2, \dots, o_T | \Delta) \\ &\quad \text{(Due to Bayes' rule [4, p. 99])} \\ &= P(o_1, o_2, \dots, o_t, x_t = s_i, o_{t+1}, o_{t+2}, \dots, o_T | \Delta) \\ &= P(o_1, o_2, \dots, o_t, x_t = s_i | \Delta) \\ &\quad * P(o_{t+1}, o_{t+2}, \dots, o_T | o_1, o_2, \dots, o_t, x_t = s_i, \Delta) \\ &\quad \text{(Due to multiplication rule [4, p. 100])} \\ &= P(o_1, o_2, \dots, o_t, x_t = s_i | \Delta) P(o_{t+1}, o_{t+2}, \dots, o_T | x_t = s_i, \Delta) \\ &\quad \text{(Because observations } o_1, o_2, \dots, o_T \text{ are observed independently)} \\ &= \alpha_t(i) \beta_t(i) \\ &\quad \text{(According to (1) and (4) for determining forward variable)} \end{aligned}$$

and backward variable)

The state sequence  $X = \{x_1, x_2, \dots, x_T\}$  is determined by selecting each state  $x_t \in S$  so that it maximizes  $\gamma_t(i)$ .

$$\begin{aligned} x_t &= \operatorname{argmax}_i P(x_t = s_i | o_1, o_2, \dots, o_T, \Delta) \\ &= \operatorname{argmax}_i \frac{\alpha_t(i) \beta_t(i)}{P(o_1, o_2, \dots, o_T | \Delta)} \\ &= \operatorname{argmax}_i \frac{P(o_1, o_2, \dots, o_T, x_t = s_i | \Delta)}{P(o_1, o_2, \dots, o_T | \Delta)} \\ &\quad \text{(Due to Bayes' rule [4, p. 99])} \\ &= \operatorname{argmax}_i \frac{\gamma_t(i)}{P(o_1, o_2, \dots, o_T | \Delta)} \\ &\quad \text{(Due to (7))} \end{aligned}$$

Because the probability  $P(o_1, o_2, \dots, o_T | \Delta)$  is not relevant to state sequence  $X$ , it is possible to remove it from the optimization criterion. Thus, equation (8) specifies how to find out the optimal state  $x_t$  of  $X$  at time point  $t$ .

$$x_t = \operatorname{argmax}_i \gamma_t(i) = \operatorname{argmax}_i \alpha_t(i) \beta_t(i) \quad (8)$$

Note that index  $i$  is identified with state  $s_i \in S$  according to (8). The optimal state  $x_t$  of  $X$  at time point  $t$  is the one that maximizes product  $\alpha_t(i) \beta_t(i)$  over all values  $s_i$ . The procedure to find out state sequence  $X = \{x_1, x_2, \dots, x_T\}$  based on individually optimal criterion is called *individually optimal procedure* that includes three steps, shown in table 6.

**Table 6.** Individually optimal procedure to solve uncovering problem.

1. Initialization step:
- Initializing $\alpha_1(i) = b_1(o_1) \pi_i$ for all $1 \leq i \leq n$
- Initializing $\beta_T(i) = 1$ for all $1 \leq i \leq n$
2. Recurrence step:
- Calculating all $\alpha_{t+1}(i)$ for all $1 \leq i \leq n$ and $1 \leq t \leq T-1$ according to (2).
- Calculating all $\beta_t(i)$ for all $1 \leq i \leq n$ and $t=T-1, t=T-2, \dots, t=1$ , according to (5).
- Calculating all $\gamma_t(i) = \alpha_t(i) \beta_t(i)$ for all $1 \leq i \leq n$ and $1 \leq t \leq T$ according to (7).
- Determining optimal state $x_t$ of $X$ at time point $t$ is the one that maximizes $\gamma_t(i)$ over all values $s_i$ .
$x_t = \operatorname{argmax}_i \gamma_t(i)$
3. Final step: The state sequence $X = \{x_1, x_2, \dots, x_T\}$ is totally determined when its partial states $x_t(s)$ where $1 \leq t \leq T$ are found in recurrence step.

It is required to execute  $n + (5n^2 - n)(T-1) + 2nT$  operations for individually optimal procedure due to:

- There are  $n$  multiplications for calculating  $\alpha_1(i)$  (s).
- The recurrence step runs over  $T-1$  times. There are  $2n^2(T-1)$  operations for determining  $\alpha_{t+1}(i)$  (s) over all  $1 \leq i \leq n$  and  $1 \leq t \leq T-1$ . There are  $(3n-1)n(T-1)$  operations for determining  $\beta_t(i)$  (s) over all  $1 \leq i \leq n$  and  $t=T-1, t=T-2, \dots, t=1$ . There are  $nT$  multiplications for determining  $\gamma_t(i) = \alpha_t(i) \beta_t(i)$  over all  $1 \leq i \leq n$  and  $1 \leq t \leq T$ . There are  $nT$  comparisons for determining optimal state  $x_t = \operatorname{argmax}_i \gamma_t(i)$  over all  $1 \leq i \leq n$  and  $1 \leq t \leq T$ . In general, there are  $2n^2(T-1) + (3n-1)n(T-1) + nT + nT = (5n^2 - n)(T-1) + 2nT$  operations at the recurrence step.

Inside  $n + (5n^2 - n)(T-1) + 2nT$  operations, there are  $n + (n+1)n(T-1) + 2n^2(T-1) + nT = (3n^2 + n)(T-1) + nT + n$  multiplications and  $(n-1)n(T-1) + (n-1)n(T-1) = 2(n^2 - n)(T-1)$  ad-

ditions and  $nT$  comparisons.

For example, given HMM  $\Delta$  whose parameters  $A, B$ , and  $\Pi$  specified in tables 1, 2, and 3, which is designed for weather forecast. Suppose humidity is *soggy* and *dry* in days 1 and 2, respectively. We apply individual optimal procedure into solving the uncovering problem that finding out the optimal state sequence  $X = \{x_1, x_2\}$  with regard to observation sequence  $O = \{o_1 = \varphi_4 = \text{soggy}, o_2 = \varphi_1 = \text{dry}, o_3 = \varphi_2 = \text{dryish}\}$ . According to (2) and (5), forward variable and backward variable are calculated as follows:

$$\alpha_1(1) = b_1(o_1 = \varphi_4) \pi_1 = b_{14} \pi_1 = 0.0165$$

$$\alpha_1(2) = b_2(o_1 = \varphi_4) \pi_2 = b_{24} \pi_2 = 0.0825$$

$$\alpha_1(3) = b_3(o_1 = \varphi_4) \pi_3 = b_{34} \pi_3 = 0.165$$

$$\begin{aligned} \alpha_2(1) &= \left( \sum_{i=1}^3 \alpha_1(i) a_{i1} \right) b_1(o_2 = \varphi_1) = \left( \sum_{i=1}^3 \alpha_1(i) a_{i1} \right) b_{11} \\ &= 0.04455 \end{aligned}$$

$$\alpha_2(2) = \left( \sum_{i=1}^3 \alpha_1(i) a_{i2} \right) b_2(o_2 = \varphi_1) = 0.01959375$$

$$\alpha_2(3) = \left( \sum_{i=1}^3 \alpha_1(i) a_{i3} \right) b_3(o_2 = \varphi_1) = 0.00556875$$

$$\alpha_3(1) = \left( \sum_{i=1}^3 \alpha_2(i) a_{i1} \right) b_1(o_3 = \varphi_2) = 0.0059090625$$

$$\alpha_3(2) = \left( \sum_{i=1}^3 \alpha_2(i) a_{i2} \right) b_2(o_3 = \varphi_2) = 0.005091796875$$

$$\alpha_3(3) = \left( \sum_{i=1}^3 \alpha_2(i) a_{i3} \right) b_3(o_3 = \varphi_2) = 0.00198$$

$$\beta_3(1) = \beta_3(2) = \beta_3(3) = 1$$

$$\begin{aligned} \beta_2(1) &= \sum_{j=1}^n a_{1j} b_j(o_3 = \varphi_2) \beta_3(j) = \sum_{j=1}^n a_{1j} b_{j2} \beta_3(j) \\ &= 0.1875 \end{aligned}$$

$$\beta_2(2) = \sum_{j=1}^n a_{2j} b_j(o_3 = \varphi_2) \beta_3(j) = 0.19$$

$$\beta_2(3) = \sum_{j=1}^n a_{3j} b_j(o_3 = \varphi_2) \beta_3(j) = 0.1625$$

$$\beta_1(1) = \sum_{j=1}^n a_{1j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.07015625$$

$$\beta_1(2) = \sum_{j=1}^n a_{2j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.0551875$$

$$\beta_1(3) = \sum_{j=1}^n a_{3j} b_j(o_2 = \varphi_1) \beta_2(j) = 0.0440625$$

According to recurrence step of individually optimal procedure, individually optimal criterion  $\gamma_t(i)$  and optimal state  $x_t$  are calculated as follows:

$$\gamma_1(1) = \alpha_1(1) \beta_1(1) = 0.001157578125$$

$$\gamma_1(2) = \alpha_1(2) \beta_1(2) = 0.00455296875$$

$$\gamma_1(3) = \alpha_1(3) \beta_1(3) = 0.0072703125$$



$$\begin{aligned}
x_1 &= \underset{i}{\operatorname{argmax}}\{\gamma_1(i)\} = \underset{i}{\operatorname{argmax}}\{\gamma_1(1), \gamma_1(2), \gamma_1(3)\} = s_3 \\
&= \text{rainy} \\
\gamma_2(1) &= \alpha_2(1)\beta_2(1) = 0.008353125 \\
\gamma_2(2) &= \alpha_2(2)\beta_2(2) = 0.0037228125 \\
\gamma_2(3) &= \alpha_2(3)\beta_2(3) = 0.000904921875 \\
x_2 &= \underset{i}{\operatorname{argmax}}\{\gamma_2(i)\} = \underset{i}{\operatorname{argmax}}\{\gamma_2(1), \gamma_2(2), \gamma_2(3)\} = s_1 \\
&= \text{sunny} \\
\gamma_3(1) &= \alpha_3(1)\beta_3(1) = 0.0059090625 \\
\gamma_3(2) &= \alpha_3(2)\beta_3(2) = 0.005091796875 \\
\gamma_3(3) &= \alpha_3(3)\beta_3(3) = 0.00198 \\
x_3 &= \underset{i}{\operatorname{argmax}}\{\gamma_3(i)\} = \underset{i}{\operatorname{argmax}}\{\gamma_3(1), \gamma_3(2), \gamma_3(3)\} = s_1 \\
&= \text{sunny}
\end{aligned}$$

As a result, the optimal state sequence is  $X = \{x_1=\text{rainy}, x_2=\text{sunny}, x_3=\text{sunny}\}$ .

The individually optimal criterion  $\gamma_t(i)$  does not reflect the whole probability of state sequence  $X$  given observation sequence  $O$  because it focuses only on how to find out each partially optimal state  $x_t$  at each time point  $t$ . Thus, the indi-

$$\begin{aligned}
\delta_{t+1}(j) &= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, o_{t+1}, x_1, x_2, \dots, x_t, x_{t+1} = s_j | \Delta) \right) \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_{t+1} | o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t, x_{t+1} = s_j) * P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t, x_{t+1} = s_j) \right) \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_{t+1} | x_{t+1} = s_j) * P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t, x_{t+1} = s_j) \right) \\
&\quad \text{(Due to observations are mutually independent)} \\
&= \max_{x_1, x_2, \dots, x_t} \left( b_j(o_{t+1}) * P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t, x_{t+1} = s_j) \right) \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t, x_{t+1} = s_j) \right) * b_j(o_{t+1}) \\
&\quad \text{(The probability } b_j(o_{t+1}) \text{ is moved out of the maximum operation because it is independent from states } x_1, x_2, \dots, x_t) \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_{t+1} = s_j | x_t) * P(x_t) \right) * b_j(o_{t+1}) \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1} | x_{t+1} = s_j, x_t) * P(x_{t+1} = s_j | x_t) * P(x_t) \right) * b_j(o_{t+1}) \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1} | x_t) * P(x_{t+1} = s_j | x_t) * P(x_t) \right) * b_j(o_{t+1}) \\
&\quad \text{(Because observation } x_{t+1} \text{ is dependent from } o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}) \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1} | x_t) * P(x_t) * P(x_{t+1} = s_j | x_t) \right) * b_j(o_{t+1}) \\
&= \max_{x_1, x_2, \dots, x_t} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_t) * P(x_{t+1} = s_j | x_t) \right) * b_j(o_{t+1}) \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= \max_{x_t} \left( \max_{x_1, x_2, \dots, x_{t-1}} \left( P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_t) * P(x_{t+1} = s_j | x_t) \right) \right) * b_j(o_{t+1}) \\
&= \max_{x_t} \left( \left( \max_{x_1, x_2, \dots, x_{t-1}} P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_t) \right) * P(x_{t+1} = s_j | x_t) \right) * b_j(o_{t+1}) \\
&= \max_i \left( \left( \max_{x_1, x_2, \dots, x_{t-1}} P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_t = s_i) \right) * P(x_{t+1} = s_j | x_t = s_i) \right) * b_j(o_{t+1}) \\
&= \max_i \left( \left( \max_{x_1, x_2, \dots, x_{t-1}} P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_{t-1}, x_t = s_i) \right) * a_{ij} \right) * b_j(o_{t+1}) \\
&= \max_i (\delta_t(i) * a_{ij}) * b_j(o_{t+1}) \\
&= \left( \max_i (\delta_t(i) a_{ij}) \right) b_j(o_{t+1})
\end{aligned}$$

Given criterion  $\delta_{t+1}(j)$ , the state  $x_{t+1}=s_j$  that maximizes  $\delta_{t+1}(j)$  is stored in the backtracking state  $q_{t+1}(j)$  that is specified by (11).

vidually optimal procedure is heuristic method. Viterbi algorithm [3, p. 264] is alternative method that takes interest in the whole state sequence  $X$  by using joint probability  $P(X, O | \Delta)$  of state sequence and observation sequence as optimal criterion for determining state sequence  $X$ . Let  $\delta_t(i)$  be the maximum joint probability of observation sequence  $O$  and state  $x_t=s_i$  over  $t-1$  previous states. The quantity  $\delta_t(i)$  is called *joint optimal criterion* at time point  $t$ , which is specified by (9).

$$\begin{aligned}
\delta_t(i) &= \\
&\max_{x_1, x_2, \dots, x_{t-1}} (P(o_1, o_2, \dots, o_t, x_1, x_2, \dots, x_t = s_i | \Delta)) \quad (9)
\end{aligned}$$

The recurrence property of *joint optimal criterion* is specified by (10).

$$\delta_{t+1}(j) = \left( \max_i (\delta_t(i) a_{ij}) \right) b_j(o_{t+1}) \quad (10)$$

The semantic content of joint optimal criterion  $\delta_t$  is similar to the forward variable  $\alpha_t$ . Following is the proof of (10).

$$q_{t+1}(j) = \operatorname{argmax}_i (\delta_t(i) a_{ij}) \quad (11)$$

Note that index  $i$  is identified with state  $s_i \in S$  according to (11). The Viterbi algorithm based on joint optimal criterion  $\delta_t(i)$  includes three steps described in table 7.

**Table 7.** Viterbi algorithm to solve uncovering problem.

1. Initialization step:
- Initializing $\delta_1(i) = b_1(o_1)\pi_i$ for all $1 \leq i \leq n$
- Initializing $q_1(i) = 0$ for all $1 \leq i \leq n$
2. Recurrence step:
- Calculating all
$\delta_{t+1}(j) = \left( \max_i (\delta_t(i) a_{ij}) \right) b_j(o_{t+1})$
for all $1 \leq i, j \leq n$ and $1 \leq t \leq T-1$ according to (10).
- Keeping tracking optimal states
$q_{t+1}(j) = \arg \max_i (\delta_t(i) a_{ij})$
for all $1 \leq j \leq n$ and $1 \leq t \leq T-1$ according to (11).
3. State sequence backtracking step: The resulted state sequence $X = \{x_1, x_2, \dots, x_T\}$ is determined as follows:
- The last state $x_T = \arg \max_j (\delta_T(j))$
- Previous states are determined by backtracking: $x_t = q_{t+1}(x_{t+1})$ for $t=T-1, t=T-2, \dots, t=1$ .

The total number of operations inside the Viterbi algorithm is  $2n + (2n^2 + n)(T-1)$  as follows:

- There are  $n$  multiplications for initializing  $n$  values  $\delta_1(i)$  when each  $\delta_1(i)$  requires 1 multiplication.
- There are  $(2n^2 + n)(T-1)$  operations over the recurrence step because there are  $n(T-1)$  values  $\delta_{t+1}(j)$  and each  $\delta_{t+1}(j)$  requires  $n$  multiplications and  $n$  comparisons for maximizing  $\max_i (\delta_t(i) a_{ij})$  plus 1 multiplication.
- There are  $n$  comparisons for constructing the state sequence  $X$ ,  $x_T = \max_j (q_T(j))$ .

Inside  $2n + (2n^2 + n)(T-1)$  operations, there are  $n + (n^2 + n)(T-1)$  multiplications and  $n^2(T-1) + n$  comparisons. The number of operations with regard to Viterbi algorithm is smaller than the number of operations with regard to individually optimal procedure when individually optimal procedure requires  $(5n^2 - n)(T-1) + 2nT + n$  operations. Therefore, Viterbi algorithm is more effective than individually optimal procedure. Besides, individually optimal procedure does not reflect the whole probability of state sequence  $X$  given observation sequence  $O$ .

Going back the weather HMM  $\Delta$  whose parameters  $A$ ,  $B$ , and  $\Pi$  are specified in tables 1, 2, and 3. Suppose humidity is *soggy* and *dry* in days 1 and 2, respectively. We apply Viterbi algorithm into solving the uncovering problem that finding out the optimal state sequence  $X = \{x_1, x_2, x_3\}$  with regard to observation sequence  $O = \{o_1 = \varphi_4 = \text{soggy}, o_2 = \varphi_1 = \text{dry}, o_3 = \varphi_2 = \text{dryish}\}$ . According to initialization step of Viterbi algorithm, we have:

$$\begin{aligned}\delta_1(1) &= b_1(o_1 = \varphi_4)\pi_1 = b_{14}\pi_1 = 0.0165 \\ \delta_1(2) &= b_2(o_1 = \varphi_4)\pi_2 = b_{24}\pi_2 = 0.0825 \\ \delta_1(3) &= b_3(o_1 = \varphi_4)\pi_3 = b_{34}\pi_3 = 0.165 \\ q_1(1) &= q_1(2) = q_1(3) = 0\end{aligned}$$

According to recurrence step of Viterbi algorithm, we have:

$$\begin{aligned}\delta_1(1)a_{11} &= 0.00825 \\ \delta_1(2)a_{21} &= 0.02475 \\ \delta_1(3)a_{31} &= 0.04125\end{aligned}$$

$$\begin{aligned}\delta_2(1) &= \left( \max_i \{\delta_1(i) a_{i1}\} \right) b_1(o_2 = \varphi_1) \\ &= \left( \max_i \{\delta_1(i) a_{i1}\} \right) b_{11} = 0.04125 * 0.6 \\ &= 0.02475 \\ q_2(1) &= \arg \max_i \{\delta_1(i) a_{i1}\} \\ &= \arg \max_i \{\delta_1(1) a_{11}, \delta_1(2) a_{21}, \delta_1(3) a_{31}\} \\ &= s_3 = \text{rainy} \\ \delta_1(1)a_{12} &= 0.004125 \\ \delta_1(2)a_{22} &= 0.033 \\ \delta_1(3)a_{32} &= 0.04125 \\ \delta_2(2) &= \left( \max_i \{\delta_1(i) a_{i2}\} \right) b_2(o_2 = \varphi_1) \\ &= \left( \max_i \{\delta_1(i) a_{i2}\} \right) b_{21} = 0.04125 * 0.25 \\ &= 0.0103125 \\ q_2(2) &= \arg \max_i \{\delta_1(i) a_{i2}\} \\ &= \arg \max_i \{\delta_1(1) a_{12}, \delta_1(2) a_{22}, \delta_1(3) a_{32}\} \\ &= s_3 = \text{rainy} \\ \delta_1(1)a_{13} &= 0.004125 \\ \delta_1(2)a_{23} &= 0.02475 \\ \delta_1(3)a_{33} &= 0.0825 \\ \delta_2(3) &= \left( \max_i \{\delta_1(i) a_{i3}\} \right) b_3(o_2 = \varphi_1) \\ &= \left( \max_i \{\delta_1(i) a_{i3}\} \right) b_{31} = 0.0825 * 0.05 \\ &= 0.004125 \\ q_2(2) &= \arg \max_i \{\delta_1(i) a_{i3}\} \\ &= \arg \max_i \{\delta_1(1) a_{13}, \delta_1(2) a_{23}, \delta_1(3) a_{33}\} \\ &= s_3 = \text{rainy} \\ \delta_2(1)a_{11} &= 0.012375 \\ \delta_2(2)a_{21} &= 0.00309375 \\ \delta_2(3)a_{31} &= 0.00103125 \\ \delta_3(1) &= \left( \max_i \{\delta_2(i) a_{i1}\} \right) b_1(o_3 = \varphi_2) \\ &= \left( \max_i \{\delta_2(i) a_{i1}\} \right) b_{12} = 0.012375 * 0.2 \\ &= 0.002475 \\ q_3(1) &= \arg \max_i \{\delta_2(i) a_{i1}\} \\ &= \arg \max_i \{\delta_2(1) a_{11}, \delta_2(2) a_{21}, \delta_2(3) a_{31}\} \\ &= s_1 = \text{sunny} \\ \delta_2(1)a_{12} &= 0.0061875 \\ \delta_2(2)a_{22} &= 0.004125 \\ \delta_2(3)a_{32} &= 0.00103125 \\ \delta_3(2) &= \left( \max_i \{\delta_2(i) a_{i2}\} \right) b_2(o_3 = \varphi_2) \\ &= \left( \max_i \{\delta_2(i) a_{i2}\} \right) b_{22} \\ &= 0.0061875 * 0.25 = 0.001546875 \\ q_3(2) &= \arg \max_i \{\delta_2(i) a_{i2}\} \\ &= \arg \max_i \{\delta_2(1) a_{12}, \delta_2(2) a_{22}, \delta_2(3) a_{32}\} \\ &= s_1 = \text{sunny} \\ \delta_2(1)a_{13} &= 0.0061875 \\ \delta_2(2)a_{23} &= 0.00309375 \\ \delta_2(3)a_{33} &= 0.0020625\end{aligned}$$

$$\begin{aligned}\delta_3(3) &= \left(\max_i \{\delta_2(i)a_{i3}\}\right) b_3(o_3 = \varphi_2) \\ &= \left(\max_i \{\delta_2(i)a_{i3}\}\right) b_{32} \\ &= 0.0061875 * 0.1 = 0.00061875\end{aligned}$$

$$\begin{aligned}q_3(3) &= \operatorname{argmax}_i \{\delta_2(i)a_{i3}\} \\ &= \operatorname{argmax}_i \{\delta_2(1)a_{13}, \delta_2(2)a_{23}, \delta_2(3)a_{33}\} \\ &= s_1 = \text{sunny}\end{aligned}$$

According to state sequence backtracking of Viterbi algorithm, we have:

$$\begin{aligned}x_3 &= \operatorname{argmax}_j \{\delta_3(j)\} = \operatorname{argmax}_j \{\delta_3(1), \delta_3(2), \delta_3(3)\} = s_1 \\ &= \text{sunny}\end{aligned}$$

$$x_2 = q_3(x_3 = s_1) = q_3(1) = s_1 = \text{sunny}$$

$$x_1 = q_2(x_2 = s_1) = q_2(1) = s_3 = \text{rainy}$$

As a result, the optimal state sequence is  $X = \{x_1=\text{rainy}, x_2=\text{sunny}, x_3=\text{sunny}\}$ . The result from the Viterbi algorithm is the same to the one from aforementioned individually optimal procedure described in table 6.

The uncovering problem is now described thoroughly in this section. Successive section will mention the learning problem of HMM which is the main subject of this tutorial.

## 4. HMM Learning Problem

The learning problem is to adjust parameters such as initial state distribution  $\Pi$ , transition probability matrix  $A$ , and observation probability matrix  $B$  so that given HMM  $\Delta$  gets more appropriate to an observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  with note that  $\Delta$  is represented by these parameters. In other words, the learning problem is to adjust parameters by maximizing probability of observation sequence  $O$ , as follows:

$$(A, B, \Pi) = \operatorname{argmax}_{A, B, \Pi} P(O|\Delta)$$

The Expectation Maximization (EM) algorithm is applied successfully into solving HMM learning problem, which is equivalently well-known Baum-Welch algorithm [3]. Successive sub-section 4.1 describes EM algorithm in detailed before going into Baum-Welch algorithm.

### 4.1. EM Algorithm

Expectation Maximization (EM) is effective parameter estimator in case that incomplete data is composed of two parts: observed part and missing (or hidden) part. EM is iterative algorithm that improves parameters after iterations until reaching optimal parameters. Each iteration includes two steps: E(xpectation) step and M(aximization) step. In E-step the missing data is estimated based on observed data and current estimate of parameters; so the lower-bound of likelihood function is computed by the expectation of complete data. In M-step new estimates of parameters are determined by maximizing the lower-bound. Please see document [5] for short tutorial of EM. This sub-section focuses on practice general EM algorithm; the theory of EM algorithm is described comprehensively in article "Maximum Likelihood from Incomplete Data via the EM algorithm" by authors [6].

Suppose  $O$  and  $X$  are observed data and missing (hidden) data, respectively. Note  $O$  and  $X$  can be represented in any form such as discrete values, scalar, integer number, real number, vector, list, sequence, sample, and matrix. Let  $\Theta$  represent parameters of probability distribution. Concretely,  $\Theta$  includes initial state distribution  $\Pi$ , transition probability matrix  $A$ , and observation probability matrix  $B$  inside HMM. In other words,  $\Theta$  represents HMM  $\Delta$  itself. EM algorithm aims to estimate  $\Theta$  by finding out which  $\hat{\Theta}$  maximizes the likelihood function  $L(\Theta) = P(O|\Theta)$ .

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} L(\Theta) = \operatorname{argmax}_{\Theta} P(O|\Theta)$$

Where  $\hat{\Theta}$  is the optimal estimate of parameters which is called usually *parameter estimate*. Because the likelihood function is product of factors, it is replaced by the log-likelihood function  $\ln L(\Theta)$  that is natural logarithm of the likelihood function  $L(\Theta)$ , for convenience. We have:

$$\begin{aligned}\hat{\Theta} &= \operatorname{argmax}_{\Theta} \ln L(\Theta) = \operatorname{argmax}_{\Theta} \ln(L(\Theta)) \\ &= \operatorname{argmax}_{\Theta} \ln(P(O|\Theta))\end{aligned}$$

Where,

$$\ln L(\Theta) = \ln(L(\Theta)) = \ln(P(O|\Theta))$$

The method finding out the parameter estimate  $\hat{\Theta}$  by maximizing the log-likelihood function is called maximum likelihood estimation (MLE). Of course, EM algorithm is based on MLE.

Suppose the current parameter is  $\Theta_t$  after the  $t^{\text{th}}$  iteration. Next we must find out the new estimate  $\hat{\Theta}$  that maximizes the next log-likelihood function  $\ln L(\Theta)$ . In other words it maximizes the deviation between current log-likelihood  $\ln L(\Theta_t)$  and next log-likelihood  $\ln L(\Theta)$  with regard to  $\Theta$ .

$$\hat{\Theta} = \operatorname{argmax}_{\Theta} (\ln L(\Theta) - \ln L(\Theta_t)) = \operatorname{argmax}_{\Theta} (Q(\Theta, \Theta_t))$$

Where  $Q(\Theta, \Theta_t) = \ln L(\Theta) - \ln L(\Theta_t)$  is the deviation between current log-likelihood  $\ln L(\Theta_t)$  and next log-likelihood  $\ln L(\Theta)$  with note that  $Q(\Theta, \Theta_t)$  is function of  $\Theta$  when  $\Theta_t$  was determined.

Suppose the total probability of observed data can be determined by marginalizing over missing data:

$$P(O|\Theta) = \sum_X P(O|X, \Theta)P(X|\Theta)$$

The expansion of  $P(O|\Theta)$  is total probability rule [4, p. 101]. The deviation  $Q(\Theta, \Theta_t)$  is re-written:

$$\begin{aligned}Q(\Theta, \Theta_t) &= \ln L(\Theta) - \ln L(\Theta_t) \\ &= \ln(P(O|\Theta)) - \ln(P(O|\Theta_t)) \\ &= \ln \left( \sum_X P(O|X, \Theta)P(X|\Theta) \right) - \ln(P(O|\Theta_t)) \\ &= \ln \left( \sum_X P(O, X|\Theta) \right) - \ln(P(O|\Theta_t))\end{aligned}$$

(Due to multiplication rule [4, p. 100])

$$= \ln \left( \sum_X P(X|O, \theta_t) \frac{P(O, X|\theta)}{P(X|O, \theta_t)} \right) - \ln(P(O|\theta_t))$$

Because hidden  $X$  is the complete set of mutually exclusive variables, the sum of conditional probabilities of  $X$  is equal to 1 given  $O$  and  $\theta_t$ .

$$\sum_X P(X|O, \theta_t) = 1$$

Applying Jensen's inequality [5, pp. 3-4]

$$\ln \left( \sum_i \lambda_i x_i \right) \geq \sum_i \lambda_i \ln(x_i) \text{ where } \sum_i \lambda_i = 1$$

into deviation  $Q(\theta, \theta_t)$ , we have:

$$\begin{aligned} Q(\theta, \theta_t) &\geq \left( \sum_X P(X|O, \theta_t) \ln \left( \frac{P(O, X|\theta)}{P(X|O, \theta_t)} \right) \right) \\ &\quad - \ln(P(O|\theta_t)) \\ &= \left( \sum_X P(X|O, \theta_t) \left( \ln(P(O, X|\theta)) - \ln(P(X|O, \theta_t)) \right) \right) \\ &\quad - \ln(P(O|\theta_t)) \\ &= \left( \sum_X P(X|O, \theta_t) \ln(P(O, X|\theta)) \right) \\ &\quad - \left( \sum_X P(X|O, \theta_t) \ln(P(X|O, \theta_t)) \right) \\ &\quad - \ln(P(O|\theta_t)) \\ &= \sum_X P(X|O, \theta_t) \ln(P(O, X|\theta)) + C \end{aligned}$$

Where,

$$C = - \left( \sum_X P(X|O, \theta_t) \ln(P(X|O, \theta_t)) \right) - \ln(P(O|\theta_t))$$

Because  $C$  is constant with regard to  $\theta$ , it is possible to eliminate  $C$  in order to simplify the optimization criterion as follows:

$$\begin{aligned} \hat{\theta} &= \operatorname{argmax}_{\theta} (Q(\theta, \theta_t)) \\ &\approx \operatorname{argmax}_{\theta} \left( \sum_X P(X|O, \theta_t) \ln(P(O, X|\theta)) \right) \\ &\quad - C \\ &= \operatorname{argmax}_{\theta} \sum_X P(X|O, \theta_t) \ln(P(O, X|\theta)) \end{aligned}$$

The expression  $\sum_X P(X|O, \theta_t) \ln(P(O, X|\theta))$  is essentially expectation of  $\ln(P(O, X|\theta))$  given conditional probability distribution  $P(X|O, \theta_t)$  when  $P(X|O, \theta_t)$  is totally determined. Let  $E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \}$  denote this conditional expectation, equation (12) specifies EM optimization

criterion for determining the parameter estimate, which is the most important aspect of EM algorithm.

$$\hat{\theta} = \operatorname{argmax}_{\theta} E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \} \quad (12)$$

Where,

$$E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \} = \sum_X P(X|O, \theta_t) \ln(P(O, X|\theta))$$

If  $P(X|O, \theta_t)$  is continuous density function, the continuous version of this conditional expectation is:

$$E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \} = \int_X P(X|O, \theta_t) \ln(P(O, X|\theta))$$

Finally, the EM algorithm is described in table 8.

**Table 8.** General EM algorithm.

Starting with initial parameter $\theta_0$ , each iteration in EM algorithm has two steps:
1. <i>E-step</i> : computing the conditional expectation $E_{X O, \theta_t} \{ \ln(P(O, X \theta)) \}$ based on the current parameter $\theta_t$ according to (12).
2. <i>M-step</i> : finding out the estimate $\hat{\theta}$ that maximizes such conditional expectation. The next parameter $\theta_{t+1}$ is assigned by the estimate $\hat{\theta}$ , we have:
$\theta_{t+1} = \hat{\theta}$
Of course $\theta_{t+1}$ becomes current parameter for next iteration. How to maximize the conditional expectation is optimization problem which is dependent on applications. For example, the popular method to solve optimization problem is Lagrangian duality [7, p. 8].
EM algorithm stops when it meets the terminating condition, for example, the difference of current parameter $\theta_t$ and next parameter $\theta_{t+1}$ is smaller than some pre-defined threshold $\varepsilon$ .
$ \theta_{t+1} - \theta_t  < \varepsilon$
In addition, it is possible to define a custom terminating condition.

General EM algorithm is simple but please pay attention to the concept of lower-bound and what the essence of EM is. Recall that the next log-likelihood function  $LnL(\theta)$  is current likelihood function  $LnL(\theta_t)$  plus the deviation  $Q(\theta, \theta_t)$ . We have:

$$\begin{aligned} LnL(\theta) &= LnL(\theta_t) + Q(\theta, \theta_t) \\ &\geq LnL(\theta_t) + E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \} + C \end{aligned}$$

Where,

$$C = - \left( \sum_X P(X|O, \theta_t) \ln(P(X|O, \theta_t)) \right) - \ln(P(O|\theta_t))$$

Let  $lb(\theta, \theta_t)$  denote the lower-bound of the log-likelihood function  $LnL(\theta)$  given current parameter  $\theta_t$  [5, pp. 7-8]. The lower-bound  $lb(\theta, \theta_t)$  is the function of  $\theta$  as specified by (13):

$$lb(\theta, \theta_t) = LnL(\theta_t) + E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \} + C \quad (13)$$

Determining  $lb(\theta, \theta_t)$  is to calculate the EM conditional expectation  $E_{X|O, \theta_t} \{ \ln(P(O, X|\theta)) \}$  because terms  $LnL(\theta_t)$  and  $C$  were totally determined. The lower-bound  $lb(\theta, \theta_t)$  has a feature where its evaluation at  $\theta = \theta_t$  equals the log-likelihood function  $LnL(\theta)$ .

$$lb(\theta_t, \theta_t) = LnL(\theta_t)$$

In fact,

$$\begin{aligned}
lb(\Theta, \Theta_t) &= LnL(\Theta_t) + E_{X|O, \Theta_t} \{ \ln(P(O, X|\Theta)) \} + C \\
&= LnL(\Theta_t) + \left( \sum_X P(X|O, \Theta_t) \ln(P(O, X|\Theta)) \right) \\
&\quad - \left( \sum_X P(X|O, \Theta_t) \ln(P(X|O, \Theta_t)) \right) \\
&\quad - \ln(P(O|\Theta_t)) \\
&= LnL(\Theta_t) + \left( \sum_X P(X|O, \Theta_t) \ln \left( \frac{P(O, X|\Theta)}{P(X|O, \Theta_t)} \right) \right) \\
&\quad - \ln(P(O|\Theta_t)) \\
&= LnL(\Theta_t) + \left( \sum_X P(X|O, \Theta_t) \ln \left( \frac{P(X|O, \Theta) P(O|\Theta)}{P(X|O, \Theta_t)} \right) \right) \\
&\quad - \ln(P(O|\Theta_t)) \\
&\quad \text{(Due to multiplication rule [4, p. 100])}
\end{aligned}$$

It implies

$$\begin{aligned}
lb(\Theta_t, \Theta_t) &= LnL(\Theta_t) + \left( \sum_X P(X|O, \Theta_t) \ln \left( \frac{P(X|O, \Theta_t) P(O|\Theta_t)}{P(X|O, \Theta_t)} \right) \right) \\
&\quad - \ln(P(O|\Theta_t)) \\
&= LnL(\Theta_t) + \left( \sum_X P(X|O, \Theta_t) \ln(P(O|\Theta_t)) \right) \\
&\quad - \ln(P(O|\Theta_t)) \sum_X P(X|O, \Theta_t) \\
&= LnL(\Theta_t) + \ln(P(O|\Theta_t)) \sum_X P(X|O, \Theta_t) - \ln(P(O|\Theta_t)) \\
&= LnL(\Theta_t) + \ln(P(O|\Theta_t)) - \ln(P(O|\Theta_t)) \\
&\quad \left( \text{due to } \sum_X P(X|O, \Theta_t) = 1 \right) \\
&= LnL(\Theta_t)
\end{aligned}$$

Fig. 4. [8, p. 7] shows relationship between the log-likelihood function  $LnL(\Theta)$  and its lower-bound  $lb(\Theta, \Theta_t)$ .

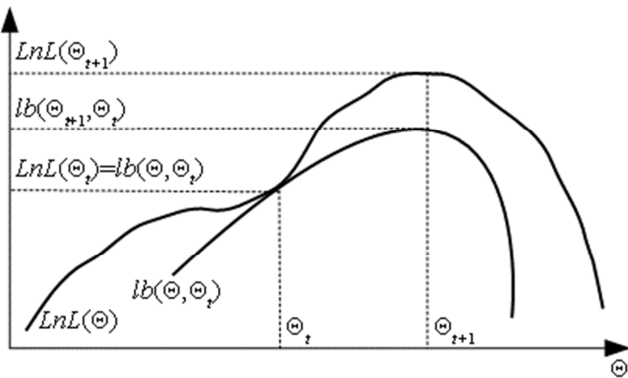


Figure 4. Relationship between the log-likelihood function and its lower-bound.

The essence of maximizing the deviation  $Q(\Theta, \Theta_t)$  is to maximize the lower-bound  $lb(\Theta, \Theta_t)$  with respect to  $\Theta$ . For each iteration the new lower-bound and its maximum are computed based on previous lower-bound. A single iteration in EM algorithm can be understood as below:

1. E-step: the new lower-bound  $lb(\Theta, \Theta_t)$  is determined based on current parameter  $\Theta_t$  according to (13). Of course, determining  $lb(\Theta, \Theta_t)$  is to calculate the EM conditional expectation  $E_{X|O, \Theta_t} \{ \ln(P(O, X|\Theta)) \}$ .
2. M-step: finding out the estimate  $\hat{\Theta}$  so that  $lb(\Theta, \Theta_t)$  reaches maximum at  $\hat{\Theta}$ . The next parameter  $\Theta_{t+1}$  is assigned by the estimate  $\hat{\Theta}$ , we have:

$$\Theta_{t+1} = \hat{\Theta}$$

Of course  $\Theta_{t+1}$  becomes current parameter for next iteration. Note, maximizing  $lb(\Theta, \Theta_t)$  is to maximize the EM conditional expectation  $E_{X|O, \Theta_t} \{ \ln(P(O, X|\Theta)) \}$ .

In general, it is easy to calculate the EM expectation  $E_{X|O, \Theta_t} \{ \ln(P(O, X|\Theta)) \}$  but finding out the estimate  $\hat{\Theta}$  based on maximizing such expectation is complicated optimization problem. It is possible to state that the essence of EM algorithm is to determine the estimate  $\hat{\Theta}$ . Now the EM algorithm is introduced with full of details. How to apply it into solving HMM learning problem is described in successive sub-section.

#### 4.2. Applying EM Algorithm into Solving Learning Problem

Now going back the HMM learning problem, the EM algorithm is applied into solving this problem, which is equivalently well-known Baum-Welch algorithm [3]. The parameter  $\Theta$  becomes the HMM model  $\Delta = (A, B, \Pi)$ . Recall that the learning problem is to adjust parameters by maximizing probability of observation sequence  $O$ , as follows:

$$\hat{\Delta} = (\hat{A}, \hat{B}, \hat{\Pi}) = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j) = \underset{\Delta}{\operatorname{argmax}} P(O|\Delta)$$

Where  $\hat{a}_{ij}$ ,  $\hat{b}_j(k)$ ,  $\hat{\pi}_j$  are parameter estimates and so, the purpose of HMM learning problem is to determine them.

The observation sequence  $O = \{o_1, o_2, \dots, o_T\}$  and state sequence  $X = \{x_1, x_2, \dots, x_T\}$  are observed data and missing (hidden) data within context of EM algorithm, respectively. Note  $O$  and  $X$  is now represented in sequence. According to EM algorithm, the parameter estimate  $\hat{\Delta}$  is determined as follows:

$$\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j) = \underset{\Delta}{\operatorname{argmax}} E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \}$$

Where  $\Delta_r = (A_r, B_r, \Pi_r)$  is the known parameter at the current iteration. Note that we use notation  $\Delta_r$  instead of popular notation  $\Delta_t$  in order to distinguish iteration indices of EM algorithm from time points inside observation sequence  $O$  and state sequence  $X$ . The EM conditional expectation in accordance with HMM is:

$$\begin{aligned}
E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \} &= \sum_X P(X|O, \Delta_r) \ln(P(O, X|\Delta)) \\
&= \sum_X P(X|O, \Delta_r) \ln(P(O|X, \Delta) P(X|\Delta)) \\
&= \sum_X P(X|O, \Delta_r) \ln(P(o_1, o_2, \dots, o_T | x_1, x_2, \dots, x_T, \Delta) \\
&\quad * P(x_1, x_2, \dots, x_T | \Delta))
\end{aligned}$$

$$\begin{aligned}
&= \sum_X P(X|O, \Delta_r) \ln(P(o_1|x_1, x_2, \dots, x_T, \Delta) \\
&\quad * P(o_2|x_1, x_2, \dots, x_T, \Delta) * \dots \\
&\quad * P(o_T|x_1, x_2, \dots, x_T, \Delta) \\
&\quad * P(x_1, x_2, \dots, x_T|\Delta)) \\
&\text{(Because observations } o_1, o_2, \dots, o_T \text{ are mutually independent)} \\
&= \sum_X P(X|O, \Delta_r) \ln(P(o_1|x_1, \Delta) * P(o_2|x_2, \Delta) * \dots \\
&\quad * P(o_T|x_T, \Delta) * P(x_1, x_2, \dots, x_T|\Delta)) \\
&\text{(Because each observations } o_i \text{ is only dependent on state } x_i) \\
&= \sum_X P(X|O, \Delta_r) \ln\left(\left(\prod_{t=1}^T P(o_t|x_t, \Delta)\right) \right. \\
&\quad \left. * P(x_1, x_2, \dots, x_T|\Delta)\right) \\
&= \sum_X P(X|O, \Delta_r) \ln\left(\left(\prod_{t=1}^T P(o_t|x_t, \Delta)\right) \right. \\
&\quad \left. * P(x_T|x_1, x_2, \dots, x_{T-1}, \Delta) \right. \\
&\quad \left. * P(x_1, x_2, \dots, x_{T-1}|\Delta)\right) \\
&= \sum_X P(X|O, \Delta_r) \ln\left(\left(\prod_{t=1}^T P(o_t|x_t, \Delta)\right) * P(x_T|x_{T-1}, \Delta) \right. \\
&\quad \left. * P(x_1, x_2, \dots, x_{T-1}|\Delta)\right) \\
&\text{(Because each state } x_i \text{ is only dependent on previous state } x_{i-1}) \\
&= \sum_X P(X|O, \Delta_r) \ln\left(\left(\prod_{t=1}^T P(o_t|x_t, \Delta)\right) * P(x_T|x_{T-1}, \Delta) \right. \\
&\quad \left. * P(x_{T-1}|x_{T-2}, \Delta) * \dots * P(x_2|x_1, \Delta) \right. \\
&\quad \left. * P(x_1|\Delta)\right) \\
&\text{(Due to recurrence on probability } P(x_1, x_2, \dots, x_i)) \\
&= \sum_X P(X|O, \Delta_r) \ln\left(\left(\prod_{t=1}^T P(o_t|x_t, \Delta)\right) \right. \\
&\quad \left. * \left(\prod_{t=2}^T P(x_t|x_{t-1}, \Delta)\right) * P(x_1|\Delta)\right)
\end{aligned}$$

It is conventional that  $P(x_1|x_0, \Delta) = P(x_1|\Delta)$  where  $x_0$  is pseudo-state, equation (14) specifies general EM conditional expectation for HMM:

$$\begin{aligned}
&E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \} = \\
&\sum_X P(X|O, \Delta_r) \ln(\prod_{t=1}^T P(x_t|x_{t-1}, \Delta) P(o_t|x_t, \Delta)) = \\
&\sum_X P(X|O, \Delta_r) \sum_{t=1}^T (\ln(P(x_t|x_{t-1}, \Delta)) + \ln(P(o_t|x_t, \Delta))) \\
&\quad (14)
\end{aligned}$$

Let  $I(x_{t-1} = s_i, x_t = s_j)$  and  $I(x_t = s_j, o_t = \varphi_k)$  are

two index functions so that

$$I(s_i = x_{t-1}, s_j = x_t) = \begin{cases} 1 & \text{if } s_i = x_{t-1} \text{ and } s_j = x_t \\ 0 & \text{otherwise} \end{cases}$$

$$I(x_t = s_j, o_t = \varphi_k) = \begin{cases} 1 & \text{if } x_t = s_j \text{ and } o_t = \varphi_k \\ 0 & \text{otherwise} \end{cases}$$

We have:

$$\begin{aligned}
&E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \} \\
&= \sum_X P(X|O, \Delta_r) \left( \sum_{t=1}^T \ln(P(x_t|x_{t-1}, \Delta)) \right. \\
&\quad \left. + \sum_{t=1}^T \ln(P(o_t|x_t, \Delta)) \right) \\
&= \sum_X P(X|O, \Delta_r) \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \ln(P(s_j|s_i, \Delta)) \right. \\
&\quad \left. + \sum_{j=1}^n \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \ln(P(\varphi_k|s_j, \Delta)) \right) \\
&= \sum_X P(X|O, \Delta_r) \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \ln(a_{ij}) \right. \\
&\quad \left. + \sum_{j=1}^n \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \ln(b_j(k)) \right)
\end{aligned}$$

Because of the convention  $P(x_1|x_0, \Delta) = P(x_1|\Delta)$ , matrix  $\Pi$  is degradation case of matrix  $A$  at time point  $t=1$ . In other words, the initial probability  $\pi_j$  is equal to the transition probability  $a_{ij}$  from pseudo-state  $x_0$  to state  $x_1=s_j$ .

$$P(x_1 = s_j|x_0, \Delta) = P(x_1 = s_j|\Delta) = \pi_j$$

Note that  $n=|S|$  is the number of possible states and  $m=|\Phi|$  is the number of possible observations.

Shortly, the EM conditional expectation for HMM is specified by (15).

$$\begin{aligned}
&E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \} = \\
&\sum_X P(X|O, \Delta_r) \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \ln(a_{ij}) \right. \\
&\quad \left. + \sum_{j=1}^n \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \ln(b_j(k)) \right) \\
&\quad (15)
\end{aligned}$$

Where,

$$I(x_{t-1} = s_i, x_t = s_j) = \begin{cases} 1 & \text{if } x_{t-1} = s_i \text{ and } x_t = s_j \\ 0 & \text{otherwise} \end{cases}$$

$$I(x_t = s_j, o_t = \varphi_k) = \begin{cases} 1 & \text{if } x_t = s_j \text{ and } o_t = \varphi_k \\ 0 & \text{otherwise} \end{cases}$$

$$P(x_1 = s_j | x_0, \Delta) = P(x_1 = s_j | \Delta) = \pi_j$$

Note that the conditional expectation  $E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}$  is function of  $\Delta$ . There are two constraints for HMM as follows:

$$\sum_{j=1}^n a_{ij} = 1, \forall i = \overline{1, n}$$

$$\sum_{k=1}^m b_j(k) = 1, \forall k = \overline{1, m}$$

Maximizing  $E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}$  with subject to these constraints is optimization problem that is solved by Lagrangian duality theorem [7, p. 8]. Original optimization problem mentions minimizing target function but it is easy to infer that maximizing target function shares the same methodology. Let  $l(\Delta, \lambda, \mu)$  be Lagrangian function constructed from  $E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}$  together with these constraints [9, p. 9], we have (16) for specifying HMM Lagrangian function as follows:

$$l(\Delta, \lambda, \mu) = l(a_{ij}, b_j(k), \lambda_i, \mu_j) = E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \} + \sum_{i=1}^n \lambda_i (1 - \sum_{j=1}^n a_{ij}) + \sum_{j=1}^m \mu_j (1 - \sum_{k=1}^m b_j(k)) \quad (16)$$

Where  $\lambda$  is  $n$ -component vector  $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)$  and  $\mu$  is  $m$ -component vector  $\mu = (\mu_1, \mu_2, \dots, \mu_m)$ . Factors  $\lambda_i \geq 0$  and  $\mu_j \geq 0$  are called Lagrange multipliers or Karush-Kuhn-Tucker multipliers [10] or dual variables. The expectation  $E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}$  is specified by (15).

The parameter estimate  $\hat{\Delta}$  is extreme point of the Lagrangian function. According to Lagrangian duality theorem [11, p. 216] [7, p. 8], we have:

$$\hat{\Delta} = (\hat{A}, \hat{B}) = (\hat{a}_{ij}, \hat{b}_j(k)) = \operatorname{argmax}_{A, B} l(\Delta, \lambda, \mu)$$

$$(\hat{\lambda}, \hat{\mu}) = \operatorname{argmin}_{\lambda, \mu} l(\Delta, \lambda, \mu)$$

The parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k))$  is determined by setting partial derivatives of  $l(\Delta, \lambda, \mu)$  with respect to  $a_{ij}$  and  $b_j(k)$  to be zero. The partial derivative of  $l(\Delta, \lambda, \mu)$  with respect to  $a_{ij}$  is:

$$\begin{aligned} \frac{\partial l(\Delta, \lambda, \mu)}{\partial a_{ij}} &= \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}}{\partial a_{ij}} \\ &\quad + \frac{\partial}{\partial a_{ij}} \left( \sum_{i=1}^n \lambda_i \left( 1 - \sum_{j=1}^n a_{ij} \right) \right) \\ &\quad + \frac{\partial}{\partial a_{ij}} \left( \sum_{j=1}^m \mu_j \left( 1 - \sum_{k=1}^m b_j(k) \right) \right) \\ &= \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}}{\partial a_{ij}} - \lambda_i \end{aligned}$$

$$\begin{aligned} &= \frac{\partial}{\partial a_{ij}} \left( \sum_X P(X|O, \Delta_r) \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \ln(a_{ij}) \right. \right. \\ &\quad \left. \left. + \sum_{j=1}^m \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \ln(b_j(k)) \right) \right) - \lambda_i \\ &= \left( \sum_X P(X|O, \Delta_r) \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \frac{\partial \ln(a_{ij})}{\partial a_{ij}} \right) - \lambda_i \\ &= \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \frac{1}{a_{ij}} \right) - \lambda_i \\ &= \frac{1}{a_{ij}} \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \right) - \lambda_i \end{aligned}$$

Setting the partial derivative  $\frac{\partial l(\Delta, \lambda, \mu)}{\partial a_{ij}}$  to be zero:

$$\frac{\partial l(\Delta, \lambda, \mu)}{\partial a_{ij}} = 0 \Leftrightarrow \frac{1}{a_{ij}} \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \right) - \lambda_i = 0$$

The parameter estimate  $\hat{a}_{ij}$  is solution of equation  $\frac{\partial l(\Delta, \lambda, \mu)}{\partial a_{ij}} = 0$ , we have:

$$\hat{a}_{ij} = \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\lambda_i}$$

It is required to estimate the Lagrange multiplier  $\lambda_i$ . The multiplier estimate  $\hat{\lambda}_i$  is determined by setting the partial derivative of  $l(\Delta, \lambda, \mu)$  with respect to  $\lambda_i$  to be zero as follows:

$$\begin{aligned} \frac{\partial l(\Delta, \lambda, \mu)}{\partial \lambda_i} &= 0 \\ \Rightarrow \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}}{\partial \lambda_i} &+ \frac{\partial}{\partial \lambda_i} \left( \sum_{i=1}^n \lambda_i \left( 1 - \sum_{j=1}^n a_{ij} \right) \right) \\ &\quad + \frac{\partial}{\partial \lambda_i} \left( \sum_{j=1}^m \mu_j \left( 1 - \sum_{k=1}^m b_j(k) \right) \right) = 0 \end{aligned}$$

$$\Rightarrow 1 - \sum_{j=1}^n a_{ij} = 0$$

Substituting  $\hat{a}_{ij}$  for  $a_{ij}$ , we have:

$$\begin{aligned} 1 - \sum_{j=1}^n \hat{a}_{ij} &= 1 - \sum_{j=1}^n \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\lambda_i} \end{aligned}$$

$$= 1 - \frac{1}{\lambda_i} \sum_{j=1}^n \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) = 0$$

It implies:

$$\begin{aligned} \hat{\lambda}_i &= \sum_{j=1}^n \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \\ &= \sum_X P(X|O, \Delta_r) \sum_{t=1}^T \sum_{j=1}^n I(x_{t-1} = s_i, x_t = s_j) \\ &= \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i) \end{aligned}$$

Where,  $I(s_i = x_{t-1})$  is index function.

$$I(x_{t-1} = s_i) = \begin{cases} 1 & \text{if } x_{t-1} = s_i \\ 0 & \text{otherwise} \end{cases}$$

Substituting  $\hat{\lambda}_i$  for  $\lambda_i$  inside

$$\hat{a}_{ij} = \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\hat{\lambda}_i}$$

We have:

$$\begin{aligned} \hat{a}_{ij} &= \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\hat{\lambda}_i} \\ &= \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i)} \end{aligned}$$

Evaluating the numerator, we have:

$$\begin{aligned} &\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \\ &= \sum_{t=1}^T \sum_X I(x_{t-1} = s_i, x_t = s_j) P(X|O, \Delta_r) \\ &= \sum_{t=1}^T \sum_X I(x_{t-1} = s_i, x_t = s_j) P(x_1, \dots, x_{t-1}, x_t, \dots, x_T | O, \Delta_r) \\ &= \sum_{t=1}^T P(x_{t-1} = s_i, x_t = s_j | O, \Delta_r) \\ &\quad \text{(Due to total probability rule [4, p. 101])} \\ &= \sum_{t=1}^T \frac{P(O, x_{t-1} = s_i, x_t = s_j | \Delta_r)}{P(O | \Delta_r)} \\ &\quad \text{(Due to multiplication rule [4, p. 100])} \\ &= \frac{1}{P(O | \Delta_r)} \sum_{t=1}^T P(O, x_{t-1} = s_i, x_t = s_j | \Delta_r) \\ &\quad \text{Evaluating the denominator, we have:} \\ &\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i) \\ &= \sum_{t=1}^T \sum_X I(x_{t-1} = s_i) P(X|O, \Delta_r) \\ &= \sum_{t=1}^T \sum_X I(x_{t-1} = s_i) P(x_1, \dots, x_{t-1}, x_t, \dots, x_T | O, \Delta_r) \\ &= \sum_{t=1}^T P(x_{t-1} = s_i | O, \Delta_r) \end{aligned}$$

$$\begin{aligned} &\text{(Due to total probability rule [4, p. 101])} \\ &= \sum_{t=1}^T \frac{P(O, x_{t-1} = s_i | \Delta_r)}{P(O | \Delta_r)} \end{aligned}$$

(Due to multiplication rule [4, p. 100])

$$\begin{aligned} &= \frac{1}{P(O | \Delta_r)} \sum_{t=1}^T P(O, x_{t-1} = s_i | \Delta_r) \\ &\quad \text{It implies} \end{aligned}$$

$$\begin{aligned} \hat{a}_{ij} &= \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j)}{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_{t-1} = s_i)} \\ &= \frac{\sum_{t=1}^T P(O, x_{t-1} = s_i, x_t = s_j | \Delta_r)}{\sum_{t=1}^T P(O, x_{t-1} = s_i | \Delta_r)} \end{aligned}$$

Because of the convention  $P(x_1 | x_0, \Delta) = P(x_1 | \Delta)$ , the estimate  $\hat{a}_{ij}$  is fixed as follows:

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T P(O, x_{t-1} = s_i, x_t = s_j | \Delta_r)}{\sum_{t=2}^T P(O, x_{t-1} = s_i | \Delta_r)}$$

The estimate of initial probability  $\hat{\pi}_j$  is known as specific estimate  $\hat{a}_{ij}$  from pseudo-state  $x_0$  to state  $x_1 = s_j$ . It means that

$$\hat{\pi}_j = \frac{P(O, x_1 = s_j | \Delta_r)}{\sum_{i=1}^n P(O, x_1 = s_i | \Delta_r)}$$

Recall that the parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k))$  is determined by setting partial derivatives of  $l(\Delta, \lambda, \mu)$  with respect to  $a_{ij}$  and  $b_j(k)$  to be zero. The parameter estimate  $\hat{a}_{ij}$  was determined. Now it is required to calculate the parameter estimate  $\hat{b}_j(k)$ . The partial derivative of Lagrangian function  $l(\Delta, \lambda, \mu)$  with respect to  $b_j(k)$  is:

$$\begin{aligned} \frac{\partial l(\Delta, \lambda, \mu)}{\partial b_j(k)} &= \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}}{\partial b_j(k)} \\ &\quad + \frac{\partial}{\partial b_j(k)} \left( \sum_{i=1}^n \lambda_i \left( 1 - \sum_{j=1}^n a_{ij} \right) \right) \\ &\quad + \frac{\partial}{\partial b_j(k)} \left( \sum_{j=1}^n \mu_j \left( 1 - \sum_{k=1}^m b_j(k) \right) \right) \\ &= \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X | \Delta)) \}}{\partial b_j(k)} - \mu_j \\ &= \frac{\partial}{\partial b_j(k)} \left( \sum_X P(X|O, \Delta_r) \left( \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T I(x_{t-1} = s_i, x_t = s_j) \ln(a_{ij}) \right. \right. \\ &\quad \left. \left. + \sum_{j=1}^n \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \ln(b_j(k)) \right) \right) - \mu_j \end{aligned}$$



$$\begin{aligned}
&= \left( \sum_X P(X|O, \Delta_r) \sum_{j=1}^n \sum_{k=1}^m \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \right. \\
&\quad \left. = \varphi_k) \frac{\partial \ln(b_j(k))}{\partial b_j(k)} \right) - \mu_j \\
&= \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \frac{1}{b_j(k)} \right) - \mu_j \\
&= \frac{1}{b_j(k)} \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \right) - \mu_j \\
&\quad \text{Setting the partial derivative } \frac{\partial l(\Delta, \lambda, \mu)}{\partial b_j(k)} \text{ to be zero:} \\
&\frac{\partial l(\Delta, \lambda, \mu)}{\partial b_j(k)} = 0 \Leftrightarrow \frac{1}{b_j(k)} \left( \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \right. \\
&\quad \left. = \varphi_k) \right) - \mu_j = 0
\end{aligned}$$

The parameter estimate  $\hat{a}_{ij}$  is solution of equation  $\frac{\partial l(\Delta, \lambda, \mu)}{\partial b_j(k)} = 0$ , we have:

$$\hat{b}_j(k) = \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\mu_j}$$

It is required to estimate the Lagrange multiplier  $\mu_j$ . The multiplier estimate  $\hat{\mu}_j$  is determined by setting the partial derivative of  $l(\Delta, \lambda, \mu)$  with respect to  $\mu_j$  to be zero as follows:

$$\begin{aligned}
&\frac{\partial l(\Delta, \lambda, \mu)}{\partial \mu_j} = 0 \\
&\Rightarrow \frac{\partial E_{X|O, \Delta_r} \{ \ln(P(O, X|\Delta)) \}}{\partial \mu_j} + \frac{\partial}{\partial \mu_j} \left( \sum_{i=1}^n \lambda_i \left( 1 - \sum_{j=1}^n a_{ij} \right) \right) \\
&\quad + \frac{\partial}{\partial \mu_j} \left( \sum_{j=1}^n \mu_j \left( 1 - \sum_{k=1}^m b_j(k) \right) \right) = 0 \\
&\Rightarrow 1 - \sum_{k=1}^m b_j(k) = 0
\end{aligned}$$

Substituting  $\hat{b}_j(k)$  for  $b_j(k)$  we have:

$$\begin{aligned}
1 - \sum_{k=1}^m \hat{b}_j(k) &= 1 - \sum_{k=1}^m \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\mu_j} \\
&= 1 - \frac{1}{\mu_j} \sum_{k=1}^m \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \\
&= 0
\end{aligned}$$

It implies:

$$\begin{aligned}
\hat{\mu}_j &= \sum_{k=1}^m \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \\
&= \sum_X P(X|O, \Delta_r) \sum_{t=1}^T \sum_{k=1}^m I(x_t = s_j, o_t = \varphi_k) \\
&= \sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j)
\end{aligned}$$

Where,  $I(s_j = x_t)$  is index function.

$$I(x_t = s_j) = \begin{cases} 1 & \text{if } x_t = s_j \\ 0 & \text{otherwise} \end{cases}$$

Substituting  $\hat{\mu}_j$  for  $\mu_j$  inside

$$\hat{b}_j(k) = \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\hat{\mu}_j}$$

We have:

$$\begin{aligned}
\hat{b}_j(k) &= \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\hat{\mu}_j} \\
&= \frac{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j)}
\end{aligned}$$

Evaluating the numerator, we have:

$$\begin{aligned}
&\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k) \\
&= \sum_{t=1}^T \sum_X I(x_t = s_j, o_t = \varphi_k) P(X|O, \Delta_r) \\
&= \sum_{t=1}^T \sum_X I(x_t = s_j, o_t = \varphi_k) P(x_1, \dots, x_t, \dots, x_T | O, \Delta_r) \\
&= \sum_{t=1}^T P(x_t = s_j | O, \Delta_r) \\
&\quad \text{(Due to total probability rule [4, p. 101])} \\
&= \sum_{t=1}^T \frac{P(O, x_t = s_j | \Delta_r)}{P(O | \Delta_r)} \\
&\quad \text{(Due to multiplication rule [4, p. 100])} \\
&= \frac{1}{P(O | \Delta_r)} \sum_{t=1}^T P(O, x_t = s_j | \Delta_r)
\end{aligned}$$

Note, the expression  $\sum_{o_t = \varphi_k}^T P(O, x_t = s_j | \Delta_r)$  expresses the sum of probabilities  $P(O, x_t = s_j | \Delta_r)$  over  $T$  time points with condition  $o_t = \varphi_k$ .

Evaluating the denominator, we have:

$$\begin{aligned}
&\sum_X P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j) \\
&= \sum_{t=1}^T \sum_X I(x_t = s_j) P(X|O, \Delta_r) \\
&= \sum_{t=1}^T P(x_t = s_j | O, \Delta_r) \\
&\quad \text{(Due to total probability rule [4, p. 101])} \\
&= \sum_{t=1}^T \frac{P(O, x_t = s_j | \Delta_r)}{P(O | \Delta_r)}
\end{aligned}$$

(Due to multiplication rule [4, p. 100])

$$= \frac{1}{P(O|\Delta_r)} \sum_{t=1}^T P(O, x_t = s_j | \Delta_r)$$

It implies

$$\hat{b}_j(k) = \frac{\sum_x P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j, o_t = \varphi_k)}{\sum_x P(X|O, \Delta_r) \sum_{t=1}^T I(x_t = s_j)} = \frac{\sum_{o_t=\varphi_k}^T P(O, x_t = s_j | \Delta_r)}{\sum_{t=1}^T P(O, x_t = s_j | \Delta_r)}$$

In general, the parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  is totally determined as follows:

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T P(O, x_{t-1} = s_i, x_t = s_j | \Delta_r)}{\sum_{t=2}^T P(O, x_{t-1} = s_i | \Delta_r)}$$

$$\hat{b}_j(k) = \frac{\sum_{o_t=\varphi_k}^T P(O, x_t = s_j | \Delta_r)}{\sum_{t=1}^T P(O, x_t = s_j | \Delta_r)}$$

$$\hat{\pi}_j = \frac{P(O, x_1 = s_j | \Delta_r)}{\sum_{i=1}^n P(O, x_1 = s_i | \Delta_r)}$$

As a convention, we use notation  $\Delta$  instead of  $\Delta_r$  for denoting known HMM at current iteration of EM algorithm. We have (17) for specifying HMM parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  given current parameter  $\Delta = (a_{ij}, b_j(k), \pi_j)$  as follows:

$$\hat{a}_{ij} = \frac{\sum_{t=2}^T P(O, x_{t-1} = s_i, x_t = s_j | \Delta)}{\sum_{t=2}^T P(O, x_{t-1} = s_i | \Delta)}$$

$$\hat{b}_j(k) = \frac{\sum_{o_t=\varphi_k}^T P(O, x_t = s_j | \Delta)}{\sum_{t=1}^T P(O, x_t = s_j | \Delta)} \quad (17)$$

$$\hat{\pi}_j = \frac{P(O, x_1 = s_j | \Delta)}{\sum_{i=1}^n P(O, x_1 = s_i | \Delta)}$$

The parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  is the ultimate solution of the learning problem. As seen in (17), it is necessary to calculate probabilities  $P(O, x_{t-1}=s_i, x_t=s_j)$  and  $P(O, x_{t-1}=s_i)$  when other probabilities  $P(O, x_t=s_j)$ ,  $P(O, x_1=s_i)$ , and  $P(O, x_1=s_j)$  are represented by the joint probability  $\gamma_t$  specified by (7).

$$P(O, x_t = s_j | \Delta) = \gamma_t(j) = \alpha_t(j) \beta_t(j)$$

$$P(O, x_1 = s_i | \Delta) = \gamma_1(i) = \alpha_1(i) \beta_1(i)$$

$$P(O, x_1 = s_j | \Delta) = \gamma_1(j) = \alpha_1(j) \beta_1(j)$$

Let  $\xi_t(i, j)$  is the joint probability that the stochastic process receives state  $s_i$  at time point  $t-1$  and state  $s_j$  at time point  $t$  given observation sequence  $O$  [3, p. 264].

$$\xi_t(i, j) = P(O, x_{t-1} = s_i, x_t = s_j | \Delta)$$

Given forward variable  $\alpha_t$  and backward variable  $\beta_t$ , if

$t \geq 2$ , we have:

$$\begin{aligned} & \alpha_{t-1}(i) a_{ij} b_j(o_t) \beta_t(j) \\ &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i | \Delta) * P(x_t = s_j | x_{t-1} = s_i) \\ & \quad * b_j(o_t) * \beta_t(j) \\ &= P(o_1, o_2, \dots, o_t | x_{t-1} = s_i, \Delta) * P(x_{t-1} = s_i | \Delta) \\ & \quad * P(x_t = s_j | x_{t-1} = s_i) * b_j(o_t) * \beta_t(j) \\ & \quad \text{(Due to multiplication rule [4, p. 100])} \\ &= P(o_1, o_2, \dots, o_t | x_{t-1} = s_i, \Delta) * P(x_t = s_j | x_{t-1} = s_i) \\ & \quad * P(x_{t-1} = s_i | \Delta) * b_j(o_t) * \beta_t(j) \\ &= P(o_1, o_2, \dots, o_t, x_t = s_j | x_{t-1} = s_i, \Delta) * P(x_{t-1} = s_i | \Delta) \\ & \quad * b_j(o_t) * \beta_t(j) \end{aligned}$$

(Because the partial observation sequence  $\{o_1, o_2, \dots, o_t\}$  is independent from current state  $x_t$  given previous state  $x_{t-1}$ )

$$\begin{aligned} &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i, x_t = s_j | \Delta) * b_j(o_t) * \beta_t(j) \\ &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i | x_t = s_j, \Delta) * P(x_t = s_j | \Delta) \\ & \quad * b_j(o_t) * \beta_t(j) \\ & \quad \text{(Due to multiplication rule [4, p. 100])} \\ &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i | x_t = s_j, \Delta) * P(x_t = s_j | \Delta) \\ & \quad * P(o_t | x_t = s_j) \\ & \quad * P(o_{t+1}, o_{t+2}, \dots, o_T | x_t = s_j, \Delta) \\ &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i | x_t = s_j, \Delta) * P(x_t = s_j | \Delta) \\ & \quad * P(o_t, o_{t+1}, o_{t+2}, \dots, o_T | x_t = s_j, \Delta) \\ & \quad \text{(Because observations } o_t, o_{t+1}, o_{t+2}, \dots, o_T \text{ are mutually independent)} \\ &= P(o_1, o_2, \dots, o_t, x_{t-1} = s_i | x_t = s_j, \Delta) \\ & \quad * P(o_t, o_{t+1}, o_{t+2}, \dots, o_T | x_t = s_j, \Delta) \\ & \quad * P(x_t = s_j | \Delta) \end{aligned}$$

$$\begin{aligned} &= P(o_1, o_2, \dots, o_t, o_{t+1}, o_{t+2}, \dots, o_T, x_{t-1} = s_i | x_t = s_j, \Delta) \\ & \quad * P(x_t = s_j | \Delta) \\ & \quad \text{(Due to multiplication rule [4, p. 100])} \\ &= P(o_1, o_2, \dots, o_T, x_t = s_i, x_{t+1} = s_j | \Delta) \\ & \quad \text{(Due to multiplication rule [4, p. 100])} \\ &= P(O, x_{t-1} = s_i, x_t = s_j | \Delta) = \xi_t(i, j) \end{aligned}$$

In general, equation (18) determines the joint probability  $\xi_t(i, j)$  based on forward variable  $\alpha_t$  and backward variable  $\beta_t$ .

$$\xi_t(i, j) = \alpha_{t-1}(i) a_{ij} b_j(o_t) \beta_t(j) \text{ where } t \geq 2 \quad (18)$$

Where forward variable  $\alpha_t$  and backward variable  $\beta_t$  are calculated by previous recurrence equations (2) and (5).

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

$$\beta_t(i) = \sum_{j=1}^n a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$

Shortly, the joint probability  $\xi_t(i, j)$  is constructed from forward variable and backward variable, as seen in fig. 5 [3, p. 264].

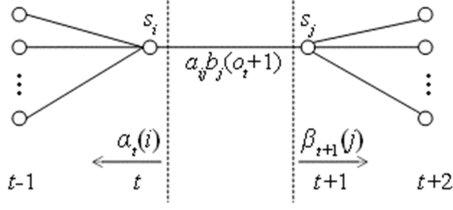


Figure 5. Construction of the joint probability  $\xi_t(i, j)$ .

Recall that  $\gamma_t(j)$  is the joint probability that the stochastic process is in state  $s_j$  at time point  $t$  with observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , specified by (7).

$$\gamma_t(j) = P(O, x_t = s_j | \Delta) = \alpha_t(j) \beta_t(j)$$

According to total probability rule [4, p. 101], it is easy to infer that  $\gamma_t$  is sum of  $\xi_t$  over all states with  $t \geq 2$ , as seen in (19).

$$\forall t \geq 2, \gamma_t(j) = \sum_{i=1}^n \xi_t(i, j) \text{ and } \gamma_{t-1}(i) = \sum_{j=1}^n \xi_t(i, j) \quad (19)$$

Deriving from (18) and (19), we have:

$$\begin{aligned} P(O, x_{t-1} = s_i, x_t = s_j | \Delta) &= \xi_t(i, j) \\ P(O, x_{t-1} = s_i | \Delta) &= \sum_{j=1}^n \xi_t(i, j), \forall t \geq 2 \\ P(O, x_t = s_j | \Delta) &= \gamma_t(j) \\ P(O, x_1 = s_j | \Delta) &= \gamma_1(j) \end{aligned}$$

By extending (17), we receive (20) for specifying HMM parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$  given current parameter  $\Delta = (a_{ij}, b_i(k), \pi_i)$  in detailed.

$$\begin{aligned} \hat{a}_{ij} &= \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_{t=2}^T \sum_{l=1}^n \xi_t(i, l)} \\ \hat{b}_j(k) &= \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \\ \hat{\pi}_j &= \frac{\gamma_1(j)}{\sum_{i=1}^n \gamma_1(i)} \end{aligned} \quad (20)$$

Followings are interpretations relevant to the joint probabilities  $\xi_t(i, j)$  and  $\gamma_t(j)$  with observation sequence  $O$ .

- The sum  $\sum_{t=2}^T \xi_t(i, j)$  expresses expected number of transitions from state  $s_i$  to state  $s_j$  [3, p. 265].
- The double sum  $\sum_{t=2}^T \sum_{l=1}^n \xi_t(i, l)$  expresses expected number of transitions from state  $s_i$  [3, p. 265].
- The sum  $\sum_{t=1}^T \gamma_t(j)$  expresses expected number of times in state  $s_j$  and in observation  $\phi_k$  [3, p. 265].
- The sum  $\sum_{t=1}^T \gamma_t(j)$  expresses expected number of times in state  $s_j$  [3, p. 265].

Followings are interpretations of the parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$ :

- The transition estimate  $\hat{a}_{ij}$  is expected frequency of

transitions from state  $s_i$  to state  $s_j$ .

- The observation estimate  $\hat{b}_j(k)$  is expected frequency of times in state  $s_j$  and in observation  $\phi_k$ .
- The initial estimate  $\hat{\pi}_j$  is (normalized) expected frequency of state  $s_j$  at the first time point ( $t=1$ ).

It is easy to infer that the parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$  is based on joint probabilities  $\xi_t(i, j)$  and  $\gamma_t(j)$  which, in turn, are based on current parameter  $\Delta = (a_{ij}, b_i(k), \pi_i)$ . The EM conditional expectation  $E_{X|O, \Delta_r} \{\ln(P(O, X | \Delta))\}$  is determined by joint probabilities  $\xi_t(i, j)$  and  $\gamma_t(j)$ ; so, the main task of E-step in EM algorithm is essentially to calculate the joint probabilities  $\xi_t(i, j)$  and  $\gamma_t(j)$  according to (18) and (7). The EM conditional expectation  $E_{X|O, \Delta_r} \{\ln(P(O, X | \Delta))\}$  gets maximal at estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$  and so, the main task of M-step in EM algorithm is essentially to calculate  $\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i$  according to (20). The EM algorithm is interpreted in HMM learning problem, as shown in table 9.

Table 9. EM algorithm for HMM learning problem.

Starting with initial value for $\Delta$ , each iteration in EM algorithm has two steps:
1. <i>E-step</i> : Calculating the joint probabilities $\xi_t(i, j)$ and $\gamma_t(j)$ according to (18) and (7) given current parameter $\Delta = (a_{ij}, b_i(k), \pi_i)$ . $\xi_t(i, j) = \alpha_{t-1}(i) a_{ij} b_j(o_t) \beta_t(j)$ where $t \geq 2$ $\gamma_t(j) = P(O, x_t = s_j   \Delta) = \alpha_t(j) \beta_t(j)$ Where forward variable $\alpha_t$ and backward variable $\beta_t$ are calculated by (2) and (5).
2. <i>M-step</i> : Calculating the estimate $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$ based on the joint probabilities $\xi_t(i, j)$ and $\gamma_t(j)$ determined at E-step, according to (20).
$\hat{a}_{ij} = \frac{\sum_{t=2}^T \xi_t(i, j)}{\sum_{t=2}^T \sum_{l=1}^n \xi_t(i, l)}$ $\hat{b}_j(k) = \frac{\sum_{t=1}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}$ $\hat{\pi}_j = \frac{\gamma_1(j)}{\sum_{i=1}^n \gamma_1(i)}$
The estimate $\hat{\Delta}$ becomes the current parameter for next iteration.
EM algorithm stops when it meets the terminating condition, for example, the difference of current parameter $\Delta$ and next parameter $\hat{\Delta}$ is insignificant. It is possible to define a custom terminating condition.

The algorithm to solve HMM learning problem shown in table 9 is known as Baum-Welch algorithm [3]. Please see document “Hidden Markov Models Fundamentals” by [9, pp. 8-13] for more details about HMM learning problem. As aforementioned in sub-section 4.1, the essence of EM algorithm applied into HMM learning problem is to determine the estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_i(k), \hat{\pi}_i)$ .

As seen in table 9, it is not difficult to run E-step and M-step of EM algorithm but how to determine the terminating condition is considerable problem. It is better to establish a computational terminating criterion instead of applying the general statement “EM algorithm stops when it meets the terminating condition, for example, the difference of current parameter  $\Delta$  and next parameter  $\hat{\Delta}$  is insignificant”. Going back the learning problem that EM algorithm solves, the EM algorithm aims to maximize probability  $P(O | \Delta)$  of given observation sequence  $O = (o_1, o_2, \dots, o_T)$  so as to find out the estimate  $\hat{\Delta}$ . Maximizing the probability  $P(O | \Delta)$  is equivalent to max-

imizing the conditional expectation. So it is easy to infer that EM algorithm stops when probability  $P(O|\Delta)$  approaches to maximal value and EM algorithm cannot maximize  $P(O|\Delta)$  any more. In other words, the probability  $P(O|\Delta)$  is terminating criterion. Calculating criterion  $P(O|\Delta)$  is evaluation problem described in section 2. Criterion  $P(O|\Delta)$  is determined according to forward-backward procedure; please see tables 4 and 5 for more details about forward-backward procedure.

At the end of M-step, the next criterion  $P(O|\hat{\Delta})$  that is calculated based on the next parameter (also estimate)  $\hat{\Delta}$  is compared with the current criterion  $P(O|\Delta)$  that is calculated in the previous time. If these two criterions are the same or there is no significantly difference between them then, EM algorithm stops. This implies EM algorithm cannot maximize  $P(O|\Delta)$  any more. However, calculating the next criterion  $P(O|\hat{\Delta})$  according to forward-backward procedure causes EM algorithm to run slowly. This drawback is overcome by following comment and improvement. The essence of forward-backward procedure is to determine forward variables  $\alpha_t$  while EM algorithm must calculate all forward variables and backward variables in its learning process (E-step). Thus, the evaluation of terminating condition is accelerated by executing forward-backward procedure inside the E-step of EM algorithm. In other words, when EM algorithm results out forward variables in E-step, the forward-backward procedure takes advantages of such forward variables so as to determine criterion  $P(O|\Delta)$  the at the same time. As a result, the speed of EM algorithm does not decrease. However, there is always a redundant iteration; suppose that the terminating criterion approaches to maximal value at the end of the  $r^{th}$  iteration but the EM algorithm only stops at the E-step of the  $(r+1)^{th}$  iteration when it really evaluates the terminating criterion. In general, the terminating criterion  $P(O|\Delta)$  is calculated based on the current parameter  $\Delta$  at E-step instead of the estimate  $\hat{\Delta}$  at M-step. Table 10 shows the proposed implementation of EM algorithm with terminating criterion  $P(O|\Delta)$ . Pseudo-code like programming language C is used to describe the implementation of EM algorithm. Variables are marked as *italic words* and comments are followed by the signs `//` and `/*`. Please pay attention to programming language keywords: *while*, *for*, *if*, `[]`, `==`, `!=`, `&&`, `//`, `/*`, `*/`, etc. For example, notation `[]` denotes array index operation; concretely,  $\alpha[t][i]$  denotes forward variable  $\alpha_t(i)$  at time point  $t$  with regard to state  $s_i$ .

**Table 10.** Proposed implementation of EM algorithm for learning HMM with terminating criterion  $P(O|\Delta)$ .

```
/*
Input:
  HMM with current parameter  $\Delta = \{a_{ij}, \pi_j, b_{jk}\}$ 
  Observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ 
Output:
  HMM with optimized parameter  $\Delta = \{a_{ij}, \pi_j, b_{jk}\}$ 
*/

Allocating memory for two matrices  $\alpha$  and  $\beta$  representing forward variables
and backward variables.
previous_criterion = -1
current_criterion = -1
iteration = 0
/*Pre-defined number MAX_ITERATION is used to prevent from infinite
```

```
loop.*/
MAX_ITERATION = 10000
While (iteration < MAX_ITERATION)

  //Calculating forward variables and backward variables
  For t = 1 to T
    For i = 1 to n
      Calculating forward variables  $\alpha[t][i]$  and backward variables
       $\beta[T-t+1][i]$  based on observation sequence  $O$  according to (2) and
      (5).
    End for i
  End for t

  //Calculating terminating criterion  $current\_criterion = P(O|\Delta)$ 
  current_criterion = 0
  For i = 1 to n
    current_criterion = current_criterion +  $\alpha[T][i]$ 
  End for i

  //Terminating condition
  If previous_criterion >= 0 && previous_criterion == current_criterion
    then
      break //breaking out the loop, the algorithm stops
  Else
    previous_criterion = current_criterion
  End if

  //Updating transition probability matrix
  For i = 1 to n
    denominator = 0
    Allocating numerators as a 1-dimension array including  $n$  zero elements.
    For t = 2 to T
      For k = 1 to n
         $\zeta = \alpha[t-1][i] * a_{ik} * b_k(o_t) * \beta[t][k]$ 
        numerators[k] = numerators[k] +  $\zeta$ 
        denominator = denominator +  $\zeta$ 
      End for k
    End for t

    If denominator != 0 then
      For j = 1 to n
         $a_{ij} = \text{numerators}[j] / \text{denominator}$ 
      End for j
    End if
  End for i

  //Updating initial probability matrix
  Allocating g as a 1-dimension array including  $n$  elements.
  sum = 0
  For j = 1 to n
     $g[j] = \alpha[1][j] * \beta[1][j]$ 
    sum = sum +  $g[j]$ 
  End for j

  If sum != 0 then
    For j = 1 to n
       $\pi_j = g[j] / \text{sum}$ 
    End for j
  End if

  //Updating observation probability distribution
  For j = 1 to n
    Allocating  $\gamma$  as a 1-dimension array including  $T$  elements.
    denominator = 0
    For t = 1 to T
       $\gamma[t] = \alpha[t][j] * \beta[t][j]$ 
      denominator = denominator +  $\gamma[t]$ 
    End for t

    Let  $m$  be the columns of observation distribution matrix  $B$ .
    For k = 1 to  $m$ 
      numerator = 0
      For t = 1 to T
```

```

    If  $o_t = k$  then
         $numerator = numerator + \gamma[t]$ 
    End if
End for  $t$ 

 $b_{jk} = numerator / denominator$ 
End for  $k$ 
End for  $j$ 

 $iteration = iteration + 1$ 
End while

```

According to table 10, the number of iterations is limited by a pre-defined maximum number, which aims to solve a so-called infinite loop optimization. Although it is proved that EM algorithm always converges, maybe there are two different estimates  $\hat{\Delta}_1$  and  $\hat{\Delta}_2$  at the final convergence. This situation causes EM algorithm to alternate between  $\hat{\Delta}_1$  and  $\hat{\Delta}_2$  in infinite loop. Therefore, the final estimate  $\hat{\Delta}_1$  or  $\hat{\Delta}_2$  is totally determined but the EM algorithm does not stop. This is the reason that the number of iterations is limited by a pre-defined maximum number.

Going back given weather HMM  $\Delta$  whose parameters  $A$ ,  $B$ , and  $\Pi$  are specified in tables 1, 2, and 3, suppose observation sequence is  $O = \{o_1=o_4=soggy, o_2=o_1=dry, o_3=o_2=dryish\}$ , the EM algorithm and its implementation described in tables 9 and 10 are applied into calculating the parameter estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  which is the ultimate solution of the learning problem, as below.

At the first iteration ( $r=1$ ) we have:

$$\begin{aligned}
 \alpha_1(1) &= b_1(o_1 = \varphi_4)\pi_1 = b_{14}\pi_1 = 0.0165 \\
 \alpha_1(2) &= b_2(o_1 = \varphi_4)\pi_2 = b_{24}\pi_2 = 0.0825 \\
 \alpha_1(3) &= b_3(o_1 = \varphi_4)\pi_3 = b_{34}\pi_3 = 0.165 \\
 \alpha_2(1) &= \left(\sum_{i=1}^3 \alpha_1(i)a_{i1}\right)b_1(o_2 = \varphi_1) = \left(\sum_{i=1}^3 \alpha_1(i)a_{i1}\right)b_{11} = 0.04455 \\
 \alpha_2(2) &= \left(\sum_{i=1}^3 \alpha_1(i)a_{i2}\right)b_2(o_2 = \varphi_1) = 0.01959375 \\
 \alpha_2(3) &= \left(\sum_{i=1}^3 \alpha_1(i)a_{i3}\right)b_3(o_2 = \varphi_1) = 0.00556875 \\
 \alpha_3(1) &= \left(\sum_{i=1}^3 \alpha_2(i)a_{i1}\right)b_1(o_3 = \varphi_2) = 0.0059090625 \\
 \alpha_3(2) &= \left(\sum_{i=1}^3 \alpha_2(i)a_{i2}\right)b_2(o_3 = \varphi_2) = 0.005091796875 \\
 \alpha_3(3) &= \left(\sum_{i=1}^3 \alpha_2(i)a_{i3}\right)b_3(o_3 = \varphi_2) = 0.00198 \\
 \beta_3(1) &= \beta_3(2) = \beta_3(3) = 1 \\
 \beta_2(1) &= \sum_{j=1}^n a_{1j}b_j(o_3 = \varphi_2)\beta_3(j) = \sum_{j=1}^n a_{1j}b_{j2}\beta_3(j) = 0.1875 \\
 \beta_2(2) &= \sum_{j=1}^n a_{2j}b_j(o_3 = \varphi_2)\beta_3(j) = 0.19 \\
 \beta_2(3) &= \sum_{j=1}^n a_{3j}b_j(o_3 = \varphi_2)\beta_3(j) = 0.1625 \\
 \beta_1(1) &= \sum_{j=1}^n a_{1j}b_j(o_2 = \varphi_1)\beta_2(j) = 0.07015625 \\
 \beta_1(2) &= \sum_{j=1}^n a_{2j}b_j(o_2 = \varphi_1)\beta_2(j) = 0.0551875 \\
 \beta_1(3) &= \sum_{j=1}^n a_{3j}b_j(o_2 = \varphi_1)\beta_2(j) = 0.0440625
 \end{aligned}$$

Within the E-step of the first iteration ( $r=1$ ), the terminating criterion  $P(O|\Delta)$  is calculated according to forward-backward procedure (see table 4) as follows:

$$P(O|\Delta) = \alpha_3(1) + \alpha_3(2) + \alpha_3(3) \approx 0.013$$

Within the E-step of the first iteration ( $r=1$ ), the joint probabilities  $\xi_t(i,j)$  and  $\gamma_t(j)$  are calculated based on (18) and (7) as follows:

$$\begin{aligned}
 \xi_2(1,1) &= \alpha_1(1)a_{11}b_1(o_2 = \varphi_1)\beta_2(1) = \alpha_1(1)a_{11}b_{11}\beta_2(1) \\
 &= 0.000928125 \\
 \xi_2(1,2) &= \alpha_1(1)a_{12}b_2(o_2 = \varphi_1)\beta_2(2) = 0.0001959375 \\
 \xi_2(1,3) &= \alpha_1(1)a_{13}b_3(o_2 = \varphi_1)\beta_2(3) = 0.000033515625 \\
 \xi_2(2,1) &= \alpha_1(2)a_{21}b_1(o_2 = \varphi_1)\beta_2(1) = 0.002784375 \\
 \xi_2(2,2) &= \alpha_1(2)a_{22}b_2(o_2 = \varphi_1)\beta_2(2) = 0.0015675 \\
 \xi_2(2,3) &= \alpha_1(2)a_{23}b_3(o_2 = \varphi_1)\beta_2(3) = 0.00020109375 \\
 \xi_2(3,1) &= \alpha_1(3)a_{31}b_1(o_2 = \varphi_1)\beta_2(1) = 0.004640625 \\
 \xi_2(3,2) &= \alpha_1(3)a_{32}b_2(o_2 = \varphi_1)\beta_2(2) = 0.001959375 \\
 \xi_2(3,3) &= \alpha_1(3)a_{33}b_3(o_2 = \varphi_1)\beta_2(3) = 0.0006703125 \\
 \xi_3(1,1) &= \alpha_2(1)a_{11}b_1(o_3 = \varphi_2)\beta_3(1) = 0.004455 \\
 \xi_3(1,2) &= \alpha_2(1)a_{12}b_2(o_3 = \varphi_2)\beta_3(2) = 0.002784375 \\
 \xi_3(1,3) &= \alpha_2(1)a_{13}b_3(o_3 = \varphi_2)\beta_3(3) = 0.00111375 \\
 \xi_3(2,1) &= \alpha_2(2)a_{21}b_1(o_3 = \varphi_2)\beta_3(1) = 0.001175625 \\
 \xi_3(2,2) &= \alpha_2(2)a_{22}b_2(o_3 = \varphi_2)\beta_3(2) = 0.001959375 \\
 \xi_3(2,3) &= \alpha_2(2)a_{23}b_3(o_3 = \varphi_2)\beta_3(3) = 0.0005878125 \\
 \xi_3(3,1) &= \alpha_2(3)a_{31}b_1(o_3 = \varphi_2)\beta_3(1) = 0.0002784375 \\
 \xi_3(3,2) &= \alpha_2(3)a_{32}b_2(o_3 = \varphi_2)\beta_3(2) = 0.000348046875 \\
 \xi_3(3,3) &= \alpha_2(3)a_{33}b_3(o_3 = \varphi_2)\beta_3(3) = 0.0002784375 \\
 \gamma_1(1) &= \alpha_1(1)\beta_1(1) = 0.001157578125 \\
 \gamma_1(2) &= \alpha_1(2)\beta_1(2) = 0.00455296875 \\
 \gamma_1(3) &= \alpha_1(3)\beta_1(3) = 0.0072703125 \\
 \gamma_2(1) &= \alpha_2(1)\beta_2(1) = 0.008353125 \\
 \gamma_2(2) &= \alpha_2(2)\beta_2(2) = 0.0037228125 \\
 \gamma_2(3) &= \alpha_2(3)\beta_2(3) = 0.000904921875 \\
 \gamma_3(1) &= \alpha_3(1)\beta_3(1) = 0.0059090625 \\
 \gamma_3(2) &= \alpha_3(2)\beta_3(2) = 0.005091796875 \\
 \gamma_3(3) &= \alpha_3(3)\beta_3(3) = 0.00198
 \end{aligned}$$

Within the M-step of the first iteration ( $r=1$ ), the estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  is calculated based on joint probabilities  $\xi_t(i,j)$  and  $\gamma_t(j)$  determined at E-step.

$$\begin{aligned}
 \hat{a}_{11} &= \frac{\sum_{t=2}^3 \xi_t(1,1)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(1,l)} \approx 0.5660 \\
 \hat{a}_{12} &= \frac{\sum_{t=2}^3 \xi_t(1,2)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(1,l)} \approx 0.3134 \\
 \hat{a}_{13} &= \frac{\sum_{t=2}^3 \xi_t(1,3)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(1,l)} \approx 0.1206 \\
 \hat{a}_{21} &= \frac{\sum_{t=2}^3 \xi_t(2,1)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(2,l)} \approx 0.4785 \\
 \hat{a}_{22} &= \frac{\sum_{t=2}^3 \xi_t(2,2)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(2,l)} \approx 0.4262 \\
 \hat{a}_{23} &= \frac{\sum_{t=2}^3 \xi_t(2,3)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(2,l)} \approx 0.0953 \\
 \hat{a}_{31} &= \frac{\sum_{t=2}^3 \xi_t(3,1)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(3,l)} \approx 0.6017 \\
 \hat{a}_{32} &= \frac{\sum_{t=2}^3 \xi_t(3,2)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(3,l)} \approx 0.2822
 \end{aligned}$$

$$\begin{aligned}\hat{a}_{33} &= \frac{\sum_{t=2}^3 \xi_t(3,3)}{\sum_{t=2}^3 \sum_{l=1}^3 \xi_t(3,l)} \approx 0.1161 \\ \hat{b}_2(1) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_1}}^3 \gamma_t(2)}{\sum_{t=1}^3 \gamma_t(2)} = \frac{\gamma_2(2)}{\gamma_1(2) + \gamma_2(2) + \gamma_3(2)} \approx 0.2785 \\ \hat{b}_2(2) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_2}}^3 \gamma_t(2)}{\sum_{t=1}^3 \gamma_t(2)} = \frac{\gamma_3(2)}{\gamma_1(2) + \gamma_2(2) + \gamma_3(2)} \approx 0.3809 \\ \hat{b}_2(3) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_3}}^3 \gamma_t(2)}{\sum_{t=1}^3 \gamma_t(2)} = \frac{0}{\gamma_1(2) + \gamma_2(2) + \gamma_3(2)} = 0 \\ \hat{b}_2(4) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_4}}^3 \gamma_t(2)}{\sum_{t=1}^3 \gamma_t(2)} = \frac{\gamma_1(2)}{\gamma_1(2) + \gamma_2(2) + \gamma_3(2)} \approx 0.3406 \\ \hat{b}_3(1) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_1}}^3 \gamma_t(3)}{\sum_{t=1}^3 \gamma_t(3)} = \frac{\gamma_2(3)}{\gamma_1(3) + \gamma_2(3) + \gamma_3(3)} \approx 0.0891 \\ \hat{b}_3(2) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_2}}^3 \gamma_t(3)}{\sum_{t=1}^3 \gamma_t(3)} = \frac{\gamma_3(3)}{\gamma_1(3) + \gamma_2(3) + \gamma_3(3)} \approx 0.1950 \\ \hat{b}_3(3) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_3}}^3 \gamma_t(3)}{\sum_{t=1}^3 \gamma_t(3)} = \frac{0}{\gamma_1(3) + \gamma_2(3) + \gamma_3(3)} = 0 \\ \hat{b}_3(4) &= \frac{\sum_{\substack{t=1 \\ o_t=\varphi_4}}^3 \gamma_t(3)}{\sum_{t=1}^3 \gamma_t(3)} = \frac{\gamma_1(3)}{\gamma_1(3) + \gamma_2(3) + \gamma_3(3)} \approx 0.7159 \\ \hat{\pi}_1 &= \frac{\gamma_1(1)}{\gamma_1(1) + \gamma_1(2) + \gamma_1(3)} \approx 0.0892 \\ \hat{\pi}_2 &= \frac{\gamma_1(2)}{\gamma_1(1) + \gamma_1(2) + \gamma_1(3)} \approx 0.3507 \\ \hat{\pi}_3 &= \frac{\gamma_1(3)}{\gamma_1(1) + \gamma_1(2) + \gamma_1(3)} \approx 0.5601\end{aligned}$$

At the second iteration ( $r=2$ ), the current parameter  $\Delta = (a_{ij}, b_j(k), \pi_j)$  is received values from the estimate  $\hat{\Delta} = (\hat{a}_{ij}, \hat{b}_j(k), \hat{\pi}_j)$  above. By repeating the similar calculation, it is easy to determine HMM parameters at the second iteration. Table 11 summarizes HMM parameters resulted from the first iteration and the second iteration of EM algorithm.

**Table 11.** HMM parameters resulted from the first iteration and the second iteration of EM algorithm.

Iteration	HMM parameters			
1 <sup>st</sup>	$\hat{a}_{11} = 0.5660$	$\hat{a}_{12} = 0.3134$	$\hat{a}_{13} = 0.1206$	
	$\hat{a}_{21} = 0.4785$	$\hat{a}_{22} = 0.4262$	$\hat{a}_{23} = 0.0953$	
	$\hat{a}_{31} = 0.6017$	$\hat{a}_{32} = 0.2822$	$\hat{a}_{33} = 0.1161$	
	$\hat{b}_1(1) = 0.5417$	$\hat{b}_1(2) = 0.3832$	$\hat{b}_1(3) = 0$	$\hat{b}_1(4) = 0.0751$
	$\hat{b}_2(1) = 0.2785$	$\hat{b}_2(2) = 0.3809$	$\hat{b}_2(3) = 0$	$\hat{b}_2(4) = 0.3406$
	$\hat{b}_3(1) = 0.0891$	$\hat{b}_3(2) = 0.1950$	$\hat{b}_3(3) = 0$	$\hat{b}_3(4) = 0.7159$
	$\hat{\pi}_1 = 0.0892$	$\hat{\pi}_2 = 0.3507$	$\hat{\pi}_3 = 0.5601$	
	Terminating criterion $P(O \Delta) = 0.013$			
	$\hat{a}_{11} = 0.6053$	$\hat{a}_{12} = 0.3299$	$\hat{a}_{13} = 0.0648$	
	$\hat{a}_{21} = 0.5853$	$\hat{a}_{22} = 0.3781$	$\hat{a}_{23} = 0.0366$	
2 <sup>nd</sup>	$\hat{a}_{31} = 0.7793$	$\hat{a}_{32} = 0.1946$	$\hat{a}_{33} = 0.0261$	
	$\hat{b}_1(1) = 0.5605$	$\hat{b}_1(2) = 0.4302$	$\hat{b}_1(3) = 0$	$\hat{b}_1(4) = 0.0093$

$\hat{b}_2(1) = 0.2757$	$\hat{b}_2(2) = 0.4517$	$\hat{b}_2(3) = 0$	$\hat{b}_2(4) = 0.2726$
$\hat{b}_3(1) = 0.0283$	$\hat{b}_3(2) = 0.0724$	$\hat{b}_3(3) = 0$	$\hat{b}_3(4) = 0.8993$
$\hat{\pi}_1 = 0.0126$	$\hat{\pi}_2 = 0.2147$	$\hat{\pi}_3 = 0.7727$	
Terminating criterion $P(O \Delta) = 0.0776$			

As seen in table 11, the EM algorithm does not converge yet when it produces two different terminating criterions (0.013 and 0.0776) at the first iteration and the second iteration. It is necessary to run more iterations so as to gain the most optimal estimate. Within this example, the EM algorithm converges absolutely after 10 iterations when the criterion  $P(O|\Delta)$  approaches to the same value 1 at the 9<sup>th</sup> and 10<sup>th</sup> iterations. Table 12 shows HMM parameter estimates along with terminating criterion  $P(O|\Delta)$  at the 9<sup>th</sup> and 10<sup>th</sup> iterations of EM algorithm.

**Table 12.** HMM parameters along with terminating criterions after 10 iterations of EM algorithm.

Iteration	HMM parameters			
9 <sup>th</sup>	$\hat{a}_{11} = 0$	$\hat{a}_{12} = 1$	$\hat{a}_{13} = 0$	
	$\hat{a}_{21} = 0$	$\hat{a}_{22} = 1$	$\hat{a}_{23} = 0$	
	$\hat{a}_{31} = 1$	$\hat{a}_{32} = 0$	$\hat{a}_{33} = 0$	
	$\hat{b}_1(1) = 1$	$\hat{b}_1(2) = 0$	$\hat{b}_1(3) = 0$	$\hat{b}_1(4) = 0$
	$\hat{b}_2(1) = 0$	$\hat{b}_2(2) = 1$	$\hat{b}_2(3) = 0$	$\hat{b}_2(4) = 0$
	$\hat{b}_3(1) = 0$	$\hat{b}_3(2) = 0$	$\hat{b}_3(3) = 0$	$\hat{b}_3(4) = 1$
	$\hat{\pi}_1 = 0$	$\hat{\pi}_2 = 0$	$\hat{\pi}_3 = 1$	
	Terminating criterion $P(O \Delta) = 1$			
	$\hat{a}_{11} = 0$	$\hat{a}_{12} = 1$	$\hat{a}_{13} = 0$	
	$\hat{a}_{21} = 0$	$\hat{a}_{22} = 1$	$\hat{a}_{23} = 0$	
10 <sup>th</sup>	$\hat{a}_{31} = 1$	$\hat{a}_{32} = 0$	$\hat{a}_{33} = 0$	
	$\hat{b}_1(1) = 1$	$\hat{b}_1(2) = 0$	$\hat{b}_1(3) = 0$	$\hat{b}_1(4) = 0$
	$\hat{b}_2(1) = 0$	$\hat{b}_2(2) = 1$	$\hat{b}_2(3) = 0$	$\hat{b}_2(4) = 0$
	$\hat{b}_3(1) = 0$	$\hat{b}_3(2) = 0$	$\hat{b}_3(3) = 0$	$\hat{b}_3(4) = 1$
	$\hat{\pi}_1 = 0$	$\hat{\pi}_2 = 0$	$\hat{\pi}_3 = 1$	
	Terminating criterion $P(O \Delta) = 1$			

As a result, the learned parameters  $A$ ,  $B$ , and  $\Pi$  are shown in table 13:

**Table 13.** HMM parameters of weather example learned from EM algorithm.

		Weather current day (Time point $t$ )			
		<i>sunny</i>	<i>cloudy</i>	<i>rainy</i>	
Weather previous day (Time point $t - 1$ )	<i>sunny</i>	$a_{11}=0$	$a_{12}=1$	$a_{13}=0$	
	<i>cloudy</i>	$a_{21}=0$	$a_{22}=1$	$a_{23}=0$	
	<i>rainy</i>	$a_{31}=1$	$a_{32}=0$	$a_{33}=0$	
sunny $\pi_1=0$		cloudy $\pi_2=0$	rainy $\pi_3=1$		
		Humidity			
		<i>dry</i>	<i>dryish</i>	<i>damp</i>	<i>soggy</i>
Weather	<i>sunny</i>	$b_{11}=1$	$b_{12}=0$	$b_{13}=0$	$b_{14}=0$
	<i>cloudy</i>	$b_{21}=0$	$b_{22}=1$	$b_{23}=0$	$b_{24}=0$
	<i>rainy</i>	$b_{31}=0$	$b_{32}=0$	$b_{33}=0$	$b_{34}=1$

Such learned parameters are more appropriate to the training observation sequence  $O = \{o_1=\phi_4=soggy, o_2=\phi_1=dry, o_3=\phi_2=dryish\}$  than the original ones shown in tables 1, 2, and 3 when the terminating criterion  $P(O|\Delta)$  corresponding to its optimal state sequence is 1.

Now three main problems of HMM are described; please see an excellent document “A tutorial on hidden Markov models and selected applications in speech recognition” written by the author Rabiner [3] for advanced details about HMM.

## 5. Conclusion

In general, there are three main problems of HMM such as evaluation problem, uncovering problem, and learning problem. For evaluation problem and uncovering problem, researchers should pay attention to forward variable and backward variable. Most computational operations are relevant to them. They reflect unique aspect of HMM. The Viterbi algorithm is very effective to solve the uncovering problem. The Baum-Welch algorithm is often used to solve the learning problem. It is easier to explain Baum-Welch algorithm by combination of EM algorithm and optimization theory, in which the Lagrangian function is maximized so as to find out optimal parameters of EM algorithms when such parameters are also learned parameters of HMM.

Observations of normal HMM described in this report are quantified by discrete probability distribution which is observation probability matrix  $B$ . In the most general case, observation is represented by continuous variable and matrix  $B$  is replaced by probability density function. At that time the normal HMM becomes continuously observational HMM. Readers are recommended to research continuously observational HMM, an enhanced variant of normal HMM.

---

## References

- [1] E. Fosler-Lussier, "Markov Models and Hidden Markov Models: A Brief Tutorial," 1998.
- [2] J. G. Schmolze, "An Introduction to Hidden Markov Models," 2001.
- [3] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257-286, 1989.
- [4] L. Nguyen, "Mathematical Approaches to User Modeling," *Journals Consortium*, 2015.
- [5] B. Sean, "The Expectation Maximization Algorithm - A short tutorial," Sean Borman's Homepage, 2009.
- [6] A. P. Dempster, N. M. Laird and D. B. Rubin, "Maximum Likelihood from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society, Series B (Methodological)*, vol. 39, no. 1, pp. 1-38, 1977.
- [7] Y.-B. Jia, "Lagrange Multipliers," 2013.
- [8] S. Borman, "The Expectation Maximization Algorithm - A short tutorial," Sean Borman's Home Page, South Bend, Indiana, 2004.
- [9] D. Ramage, "Hidden Markov Models Fundamentals," 2007.
- [10] Wikipedia, "Karush–Kuhn–Tucker conditions," Wikimedia Foundation, 4 August 2014. [Online]. Available: [http://en.wikipedia.org/wiki/Karush–Kuhn–Tucker\\_conditions](http://en.wikipedia.org/wiki/Karush–Kuhn–Tucker_conditions). [Accessed 16 November 2014].
- [11] S. Boyd and L. Vandenberghe, *Convex Optimization*, New York, NY: Cambridge University Press, 2009, p. 716. G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," *Phil. Trans. Roy. Soc. London*, vol. A247, pp. 529–551, April 1955. (references).