

# Training Hidden Markov Models with Multiple Observations – A Combinatorial Method

Xiaolin Li, *Member, IEEE Computer Society*

CADlink Technology Corporation, 2440 Don Reid Drive, Suite 100,

Ottawa, Ontario, Canada K1H 1E1. xli@cadlink.com

Marc Parizeau, *Member, IEEE Computer Society*

Département de Génie Electrique et de Génie Informatique,

Université Laval, Ste-Foy, Québec, Canada G1K 7P4. parizeau@gel.ulaval.ca

and

Réjean Plamondon\*, *Fellow, IEEE*

École Polytechnique de Montréal,

Montréal, Québec, Canada H3C 3A7. rejean.plamondon@polymtl.ca

## Abstract

Hidden Markov models (HMMs) are stochastic models capable of statistical learning and classification. They have been applied in speech recognition and handwriting recognition because of their great adaptability and versatility in handling sequential signals. On the other hand, as these models have a complex structure, and also because the involved data sets usually contain uncertainty, it is difficult to analyze the multiple observation training problem without certain assumptions. For many years researchers have used Levinson's training equations in speech and handwriting applications simply assuming that all observations are independent of each other. This paper present a formal treatment of HMM multiple observation training without imposing the above assumption. In this treatment, the multiple observation probability is expressed as a combination of individual observation probabilities without losing generality. This combinatorial method gives one more freedom in making different dependence-independence assumptions. By generalizing Baum's auxiliary function into this framework and building up an associated objective function using Lagrange multiplier method, it is proved that the derived training equations guarantee the maximization of the objective function. Furthermore, we show that Levinson's training equations can be easily derived as a special case in this treatment.

Index Terms — Hidden Markov model, forward-backward procedure, Baum-Welch algorithm, multiple observation training

---

\*To whom correspondence should be addressed.

# 1 Introduction

Hidden Markov models (HMMs) are stochastic models which were introduced and studied in the late 1960s and early 1970s [1, 2, 3, 4, 5]. As the parameter space of these models is usually super-dimensional, the model training problem seems very difficult at the first glance. In 1970 Baum and his colleagues published their maximization method which gave a solution to the model training problem with a single observation [4]. In 1977 Dempster, Laird and Rubin introduced the Expectation-Maximization (EM) method for maximum likelihood estimates from incomplete data and later Wu proved some convergence properties of the EM algorithm [6], which made the EM algorithm a solid framework in statistical analysis. In 1983 Levinson, Rabiner and Sondhi presented a maximum likelihood estimation method for HMM multiple observation training, assuming that all observations are independent of each other [7]. Since then, HMMs have been widely used in speech recognition [7, 8, 9, 10, 11, 12]. More recently they have also been applied to handwriting recognition [18, 19, 20, 21, 22] as they are adaptive to random sequential signals and capable of statistical learning and classification.

Although the independence assumption of observations is helpful for problem simplification, it may not hold in some cases. For example, the observations of a syllable pronounced by a person are possibly highly correlated. Similar examples can also be found in handwriting: given a set of samples of a letter written by a person, it is difficult to assume or deny their independence properties when viewed from different perspectives. Based on these phenomena, it is better not to just rely on the independence assumption.

This paper presents a formal treatment for HMM multiple observation training without imposing the independence assumption. In this treatment, the multiple observation probability is expressed as a combination of individual observation probabilities rather than their product. The dependence-independence property of the observations is characterized by combinatorial weights. These weights give us more freedom in making different assumptions and hence in deriving corresponding training equations. By generalizing Baum's auxiliary function into this framework and building up an associated objective function using Lagrange multiplier method, it is proved that the derived training equations guarantee the maximization of the objective function and hence the convergence of the training process. Furthermore, as two special cases in this treatment, we show that Levinson's training equations can be easily derived with an independence assumption, and some other training equations can also be derived with a uniform dependence assumption.

The remainder of this paper is organized as follows. Section 2 summarizes the first order HMM. Section 3 describes the combinatorial method for HMM multiple observation training. Section 4 shows two special cases: an independence assumption versus a uniform dependence assumption. Finally, section 5 concludes this paper.

## 2 First Order Hidden Markov Model

### 2.1 Elements of HMM

A hidden Markov process is a doubly stochastic process: an underlying process which is hidden from observation, and an observable process which is determined by the underlying process. With respect to first order hidden Markov process, the model is characterized by the following elements [10]:

- *set of hidden states:*

$$S = \{S_1, S_2, \dots, S_N\} \quad (1)$$

where  $N$  is the number of states in the model.

- *state transition probability distribution<sup>1</sup>:*

$$A = \{a_{ij}\} \quad (2)$$

where for  $1 \leq i, j \leq N$ ,

$$a_{ij} = P[q_{t+1} = S_j | q_t = S_i] \quad (3)$$

$$\left\{ \begin{array}{l} 0 \leq a_{ij} \\ \sum_{j=1}^N a_{ij} = 1 \end{array} \right. \quad (4)$$

- *set of observation symbols:*

$$V = \{v_1, v_2, \dots, v_M\} \quad (5)$$

where  $M$  is the number of observation symbols per state.

- *observation symbol probability distribution<sup>2</sup>:*

$$B = \{b_j(k)\} \quad (6)$$

where for  $1 \leq j \leq N, 1 \leq k \leq M$ ,

$$b_j(k) = P[v_k \text{ at } t | q_t = S_j] \quad (7)$$

$$\left\{ \begin{array}{l} 0 \leq b_j(k) \\ \sum_{k=1}^M b_j(k) = 1 \end{array} \right. \quad (8)$$

---

<sup>1</sup> $A$  is also called transition matrix.

<sup>2</sup> $B$  is also called emission matrix.

- *initial state probability distribution:*

$$\pi = \{\pi_i\} \quad (9)$$

where for  $1 \leq i \leq N$ ,

$$\pi_i = P[q_1 = S_i] \quad (10)$$

$$\begin{cases} 0 \leq \pi_i \\ \sum_{i=1}^N \pi_i = 1 \end{cases} \quad (11)$$

For convenience, we denote an HMM as a triplet in all subsequent discussion:

$$\lambda = (A, B, \pi) \quad (12)$$

## 2.2 Ergodic model and left-right model

An HMM can be classified into one of the following types in the light of its state transition:

- *ergodic model:*

An ergodic model has full state transition.

- *left-right model*<sup>3</sup>:

A left-right model has only partial state transition such that  $a_{ij} = 0, \forall j < i$ .

## 2.3 Observation evaluation: forward-backward procedure

Let  $O = o_1 o_2 \cdots o_T$  be an observation sequence where  $o_t \in V$  is the observation symbol at time  $t$ , and let  $Q = q_1 q_2 \cdots q_T$  be a state sequence where  $q_t \in S$  is the state at time  $t$ . Given a model  $\lambda$  and an observation sequence  $O$ , the observation evaluation problem  $P(O|\lambda)$  can be solved using forward-backward procedure in terms of forward and backward variables (reference Figure 1):

- *forward variable*<sup>4</sup>:

$$\alpha_t(i) = P(o_1 o_2 \cdots o_t, q_t = S_i | \lambda) \quad (13)$$

$\alpha_t(i)$  can be solved inductively:

---

<sup>3</sup>This type of model is widely used in modeling sequential signals.

<sup>4</sup>i.e. the probability of the partial observation sequence  $o_1 o_2 \cdots o_t$  with state  $q_t = S_i$ , given model  $\lambda$ .

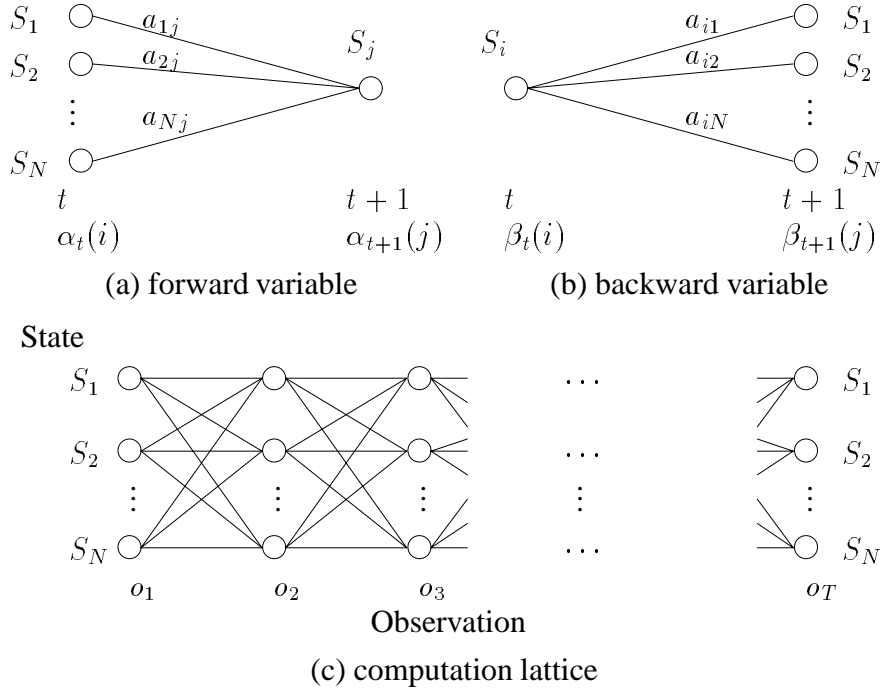


Figure 1: Illustration of forward-backward procedure

1. initialization:

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N \quad (14)$$

2. induction:

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad 1 \leq t \leq T-1, 1 \leq j \leq N \quad (15)$$

• *backward variable*<sup>5</sup>:

$$\beta_t(i) = P(o_{t+1} o_{t+2} \cdots o_T | q_t = S_i, \lambda) \quad (16)$$

$\beta_t(i)$  can be solved inductively:

1. initialization:

$$\beta_T(i) = 1, \quad 1 \leq i \leq N \quad (17)$$

2. induction:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), \quad 1 \leq t \leq T-1, 1 \leq i \leq N \quad (18)$$

<sup>5</sup>i.e. the probability of the partial observation sequence  $o_{t+1} o_{t+2} \cdots o_T$ , given state  $q_t = S_i$  and model  $\lambda$ .

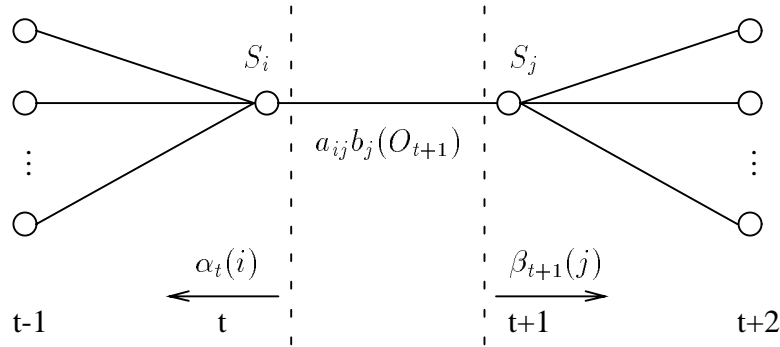


Figure 2: Illustration of the joint event

- *observation evaluation:*

$$P(O|\lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i), \quad \forall t \quad (19)$$

especially,

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i) \quad (20)$$

It is easy to see that the computational complexity of the forward-backward procedure is  $O(TN^2)$ .

## 2.4 Model training: Baum-Welch algorithm

Now let us consider the model training problem: given an observation sequence  $O$ , how to find the optimum model parameter vector  $\lambda \in \Lambda$  that maximizes  $P(O|\lambda)$ . To solve this problem, Baum and his colleagues defined an auxiliary function and proved the two propositions below [4]:

- *auxiliary function:*

$$Q(\lambda, \bar{\lambda}) = \sum_Q P(O, Q|\lambda) \log P(O, Q|\bar{\lambda}) \quad (21)$$

where  $\bar{\lambda}$  is the auxiliary variable that corresponds to  $\lambda$ .

- *proposition 1:*

If the value of  $Q(\lambda, \bar{\lambda})$  increases, then the value of  $P(O|\bar{\lambda})$  also increases, i.e.

$$Q(\lambda, \bar{\lambda}) \geq Q(\lambda, \lambda) \longrightarrow P(O|\bar{\lambda}) \geq P(O|\lambda) \quad (22)$$

- *proposition 2:*

$\lambda$  is a critical point of  $P(O|\lambda)$  if and only if it is a critical point of  $Q(\lambda, \bar{\lambda})$  as a function of  $\bar{\lambda}$ , i.e.

$$\frac{\partial P(O|\lambda)}{\partial \lambda_i} = \frac{\partial Q(\lambda, \bar{\lambda})}{\partial \bar{\lambda}_i} \Big|_{\bar{\lambda}=\lambda}, \quad 1 \leq i \leq D \quad (23)$$

where  $D$  is the dimension of  $\lambda$  and  $\lambda_i, 1 \leq i \leq D$ , are individual elements of  $\lambda$ .

In the light of the above propositions, the model training problem can be solved by the Baum-Welch algorithm in terms of joint events and state variables (reference Figure 2):

- *joint event*<sup>6</sup>:

$$\begin{aligned} \xi_t(i, j) &= P(q_t = S_i, q_{t+1} = S_j | O, \lambda) \\ &= \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(O|\lambda)} \end{aligned} \quad (24)$$

- *state variable*<sup>7</sup>:

$$\begin{aligned} \gamma_t(i) &= P(q_t = S_i | O, \lambda) \\ &= \sum_{j=1}^N \xi_t(i, j) \end{aligned} \quad (25)$$

- *parameter updating equations:*

1. state transition probability:

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}, \quad 1 \leq i \leq N, 1 \leq j \leq N \quad (26)$$

2. symbol emission probability:

$$\bar{b}_j(k) = \frac{\sum_{t=1, o_t=v_k}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}, \quad 1 \leq j \leq N, 1 \leq k \leq M \quad (27)$$

3. initial state probability:

$$\bar{\pi}_i = \gamma_1(i), \quad 1 \leq i \leq N \quad (28)$$

---

<sup>6</sup>i.e. the probability of being in state  $S_i$  at time  $t$ , and state  $S_j$  at time  $t + 1$ , given the observation sequence  $O$  and model  $\lambda$ .

<sup>7</sup>i.e. the probability of being in state  $S_i$  at time  $t$  given the observation sequence  $O$  and the model  $\lambda$ .

### 3 Multiple Observation Training

#### 3.1 Combinatorial method

Now let us consider a set of observation sequences from a pattern class:

$$\mathbf{O} = \{O^{(1)}, O^{(2)}, \dots, O^{(K)}\} \quad (29)$$

where

$$O^{(k)} = o_1^{(k)} o_2^{(k)} \dots o_{T_k}^{(k)}, \quad 1 \leq k \leq K \quad (30)$$

are individual observation sequences. Usually, one does not know if these observation sequences are independent of each other or not. And a contravercy can arise if one assumes the independence property while these observation sequences are statistically correlated. In either case, we have the following expressions without losing generality:

$$\begin{cases} P(\mathbf{O}|\lambda) = P(O^{(1)}|\lambda)P(O^{(2)}|O^{(1)}, \lambda) \dots P(O^{(K)}|O^{(K-1)} \dots O^{(1)}, \lambda) \\ P(\mathbf{O}|\lambda) = P(O^{(2)}|\lambda)P(O^{(3)}|O^{(2)}, \lambda) \dots P(O^{(1)}|O^{(K)} \dots O^{(2)}, \lambda) \\ \vdots \\ P(\mathbf{O}|\lambda) = P(O^{(K)}|\lambda)P(O^{(1)}|O^{(K)}, \lambda) \dots P(O^{(K-1)}|O^{(K)} O^{(K-2)} \dots O^{(1)}, \lambda) \end{cases} \quad (31)$$

Based on the above equations, the multiple observation probability given the model can be expressed as a summation:

$$P(\mathbf{O}|\lambda) = \sum_{k=1}^K w_k P(O^{(k)}|\lambda) \quad (32)$$

where

$$\begin{cases} w_1 = \frac{1}{K} P(O^{(2)}|O^{(1)}, \lambda) \dots P(O^{(K)}|O^{(K-1)} \dots O^{(1)}, \lambda) \\ w_2 = \frac{1}{K} P(O^{(3)}|O^{(2)}, \lambda) \dots P(O^{(1)}|O^{(K)} \dots O^{(2)}, \lambda) \\ \vdots \\ w_K = \frac{1}{K} P(O^{(1)}|O^{(K)}, \lambda) \dots P(O^{(K-1)}|O^{(K)} O^{(K-2)} \dots O^{(1)}, \lambda) \end{cases} \quad (33)$$

are weights. These weights are conditional probabilities and hence they can characterize the dependence-independence property.

Based on the above expression, we can construct an auxiliary function below for model training:

$$Q(\lambda, \bar{\lambda}) = \sum_{k=1}^K w_k Q_k(\lambda, \bar{\lambda}) \quad (34)$$

where  $\bar{\lambda}$  is the auxiliary variable corresponding to  $\lambda$ , and

$$Q_k(\lambda, \bar{\lambda}) = \sum_Q P(O^{(k)}, Q|\lambda) \log P(O^{(k)}, Q|\bar{\lambda}), \quad 1 \leq k \leq K \quad (35)$$

are Baum's auxiliary functions related to individual observations. Since  $w_k, 1 \leq k \leq K$ , are not functions of  $\bar{\lambda}$ , we have the following theorem related to the maximization of  $P(\mathbf{O}|\lambda)$ [23]:



- *theorem 1:*

If the value of  $Q(\lambda, \bar{\lambda})$  increases, then the value of  $P(\mathbf{O}|\bar{\lambda})$  also increases, i.e.

$$Q(\lambda, \bar{\lambda}) \geq Q(\lambda, \lambda) \longrightarrow P(\mathbf{O}|\bar{\lambda}) \geq P(\mathbf{O}|\lambda) \quad (36)$$

Furthermore, as  $w_k, 1 \leq k \leq K$  are weights that characterize the dependence-independence property of the observations, if one assumes that these weights are constants, one has the following theorem[23]:

- *theorem 2:*

For fixed  $w_k, 1 \leq k \leq K$ ,  $\lambda$  is a critical point of  $P(\mathbf{O}|\lambda)$  if and only if it is a critical point of  $Q(\lambda, \bar{\lambda})$  as a function of  $\bar{\lambda}$ , i.e.

$$\left. \frac{\partial P(\mathbf{O}|\lambda)}{\partial \lambda_i} = \frac{\partial Q(\lambda, \bar{\lambda})}{\partial \bar{\lambda}_i} \right|_{\bar{\lambda}=\lambda} \quad (37)$$

In such a case, the maximization of  $Q(\lambda, \bar{\lambda})$  is equivalent to the maximization of  $P(\mathbf{O}|\lambda)$ .

### 3.2 Maximization: Lagrange multiplier method

Based on theorem 1, one can always maximize  $Q(\lambda, \bar{\lambda})$  to increase the value of  $P(\mathbf{O}|\bar{\lambda})$ , regardless of 1) if the individual observations are independent of one another or not, and 2) whether the combinatorial weights are constants or not. Let us consider the auxiliary function with boundary conditions:

$$\begin{aligned} Q(\lambda, \bar{\lambda}) &= \sum_{k=1}^K w_k Q_k(\lambda, \bar{\lambda}) \\ 1 - \sum_{j=1}^N \bar{a}_{ij} &= 0, & 1 \leq i \leq N \\ 1 - \sum_{k=1}^M \bar{b}_j(k) &= 0, & 1 \leq j \leq N \\ 1 - \sum_{i=1}^N \bar{\pi}_i &= 0 \end{aligned} \quad (38)$$

we can construct an objective function using Lagrange multiplier method:

$$F(\bar{\lambda}) = Q(\lambda, \bar{\lambda}) + \sum_{i=1}^N c_{ai} [1 - \sum_{j=1}^N \bar{a}_{ij}] + \sum_{j=1}^M c_{bj} [1 - \sum_{k=1}^M \bar{b}_j(k)] + c_{\pi} [1 - \sum_{i=1}^N \bar{\pi}_i] \quad (39)$$

where  $c_{ai}$ ,  $c_{bj}$  and  $c_{\pi}$  are Lagrange multipliers. Differentiating the objective function with respect to individual parameters and finding solutions to corresponding Lagrange multipliers, we obtain the following training equations that guarantee the maximization of the objective function (see appendix for detailed derivation):

1. state transition probability:

$$\bar{a}_{mn} = \frac{\sum_{k=1}^K w_k P(O^{(k)}|\lambda) \sum_{t=1}^{T_k-1} \xi_t^{(k)}(m, n)}{\sum_{k=1}^K w_k P(O^{(k)}|\lambda) \sum_{t=1}^{T_k-1} \gamma_t^{(k)}(m)}, \quad 1 \leq m \leq N, 1 \leq n \leq N \quad (40)$$

2. symbol emission probability:

$$\bar{b}_n(m) = \frac{\sum_{k=1}^K w_k P(O^{(k)}|\lambda) \sum_{t=1, o_t^{(k)}=v_m}^{T_k} \gamma_t^{(k)}(n)}{\sum_{k=1}^K w_k P(O^{(k)}|\lambda) \sum_{t=1}^{T_k} \gamma_t^{(k)}(n)}, \quad 1 \leq n \leq N, 1 \leq m \leq M \quad (41)$$

3. initial state probability:

$$\bar{\pi}_n = \frac{\sum_{k=1}^K w_k P(O^{(k)}|\lambda) \gamma_1^{(k)}(n)}{\sum_{k=1}^K w_k P(O^{(k)}|\lambda)}, \quad 1 \leq n \leq N \quad (42)$$

### 3.3 Convergence property

The training equations derived by Lagrange multiplier method guarantee the convergence of the training process. Firstly, these training equations give the zero points of the first order Jacobi differential matrix  $\frac{\partial F(\bar{\lambda})}{\partial \lambda}$ . Secondly, the second order Jacobi differential matrix  $\frac{\partial^2 F(\bar{\lambda})}{\partial \lambda^2}$  is diagonal and all its diagonal elements are negative. Thus, the algorithm guarantees local maxima and hence the convergence of the training process (See [23] for detailed proofs).

The above training equations are adaptive to both the ergodic model and the left-right model since we do not put any constraints on the model type during the derivation.

## 4 Two Special cases - Independence versus Uniform Dependence

### 4.1 Independence assumption

Now let us assume that the individual observations are independent of each other, i.e.

$$P(O|\lambda) = \prod_{k=1}^K P(O^{(k)}|\lambda) \quad (43)$$

In this case, the combinatorial weights become:

$$w_k = \frac{1}{K} P(\mathbf{O}|\lambda) / P(O^{(k)}|\lambda), \quad 1 \leq k \leq K \quad (44)$$

Substituting the above weights into equations (40) to (42), we obtain Levinson's training equations:

1. state transition probability:

$$\bar{a}_{mn} = \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \xi_t^{(k)}(m, n)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^{(k)}(m)}, \quad 1 \leq m \leq N, 1 \leq n \leq N \quad (45)$$

2. symbol emission probability:

$$\bar{b}_n(m) = \frac{\sum_{k=1}^K \sum_{t=1, o_t^{(k)}=v_m}^{T_k} \gamma_t^{(k)}(n)}{\sum_{k=1}^K \sum_{t=1}^{T_k} \gamma_t^{(k)}(n)}, \quad 1 \leq n \leq N, 1 \leq m \leq M \quad (46)$$

3. initial state probability:

$$\bar{\pi}_n = \frac{1}{K} \sum_{k=1}^K \gamma_1^{(k)}(n), \quad 1 \leq n \leq N \quad (47)$$

## 4.2 Uniform dependence assumption

If we assume that the individual observations are uniformly dependent on one another, i.e.

$$w_k = \text{const}, \quad 1 \leq k \leq K \quad (48)$$

Substituting the above weights into equations (40) to (42), it readily follows that

1. state transition probability:

$$\bar{a}_{mn} = \frac{\sum_{k=1}^K P(O^{(k)}|\lambda) \sum_{t=1}^{T_k-1} \xi_t^{(k)}(m, n)}{\sum_{k=1}^K P(O^{(k)}|\lambda) \sum_{t=1}^{T_k-1} \gamma_t^{(k)}(m)}, \quad 1 \leq m \leq N, 1 \leq n \leq N \quad (49)$$

2. symbol emission probability:

$$\bar{b}_n(m) = \frac{\sum_{k=1}^K P(O^{(k)}|\lambda) \sum_{t=1, o_t^{(k)}=v_m}^{T_k} \gamma_t^{(k)}(n)}{\sum_{k=1}^K P(O^{(k)}|\lambda) \sum_{t=1}^{T_k} \gamma_t^{(k)}(n)}, \quad 1 \leq n \leq N, 1 \leq m \leq M \quad (50)$$

3. initial state probability:

$$\bar{\pi}_n = \frac{\sum_{k=1}^K P(O^{(k)}|\lambda) \gamma_1^{(k)}(n)}{\sum_{k=1}^K P(O^{(k)}|\lambda)}, \quad 1 \leq n \leq N \quad (51)$$

## 5 Conclusions

A formal treatment for HMM multiple observation training has been presented in this paper. In this treatment, the multiple observation probability is expressed as a combination of individual observation probabilities without losing generality. The independence-dependence property of the observations are characterized by the combinatorial weights, and hence it gives us more freedom in making different assumptions and also in deriving corresponding training equations.

The well known Baum's auxiliary function has been generalized into the case of multiple observation training, and two theorems related to the maximization have been presented in this paper. Based on the auxiliary function and its boundary conditions, an objective function has been constructed using Lagrange multiplier method, and a set of training equations have been derived by maximizing the objective function. Similar to the EM algorithm, this algorithm guarantees the local maxima and hence the convergence of the training process.

We have also shown, through two special cases, that the above training equations are general enough to include different situations. Once the independence assumption is made, one can readily obtain Levinson's training equations. On the other hand, if the uniform dependence is assumed, one can also have the corresponding training equations.

## Acknowledgements

This research work was supported by NSERC Grant OGP0155389 to M. Parizeau, and NSERC Grant OGP00915 and FCAR Grant ER-1220 to R. Plamondon.

## References

- [1] L.E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains". *The Annals of Mathematical Statistics*, Vol.37, 1554-1563 (1966)
- [2] L.E. Baum and J.A. Egon, "An inequality with applications to statistical estimation for probabilistic functions of a Markov process and to a model for ecology". *Bull. Amer. Meteorol Soc.*, Vol.73, 360-363 (1967)
- [3] L.E. Baum and G.R. Sell, "Growth functions for transformations on manifolds". *Pac. J. Math.*, Vol.27, No.2, 211-227 (1968)
- [4] L.E. Baum, T. Petrie, G. Soules, and N. Weiss, "A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains". *The Annals of Mathematical Statistics*, Vol.41, No.1, 164-171 (1970)
- [5] L.E. Baum, "An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes". *Inequalities*, Vol.3, 1-8 (1970)
- [6] C. F. J. Wu, "On the convergence properties of the EM algorithm". *The Annals of Statistics*, Vol. 11, No. 1, 95-103 (1983)
- [7] S.E. Levinson, L.R. Rabiner and M.M. Sondhi, "An introduction to the application of the theory of probabilistic functions of Markov process to automatic speech recognition". *Bell System Technical Journal*, Vol.62, No.4, 1035-1074 (1983)
- [8] L. R. Bahl, F. Jelinek and R. L. Mercer, "A maximum likelihood approach to continuous speech recognition". *IEEE Trans. Pattern Anal. Machine Intell.* Vol. PAMI-5, 179-190 (1983)
- [9] L. R. Rabiner and S. E. Levinson, "A speaker-independent, syntax-directed, connected word recognition system based on hidden Markov models and level building". *IEEE trans. Acoust. Speech Signal Processing*, Vol. ASSP 33, No. 3, 561-573 (1985)
- [10] Lawrence R. Rabiner, "A Tutorial on hidden Markov models and selected applications in speech recognition". *Proceedings of the IEEE*, 77(2). 257-286 (1989)
- [11] Kai-Fu Lee, "Automatic Speech Recognition - the Development of SPHINX System". *Kluwer Academic Publishers* (1989)
- [12] L. Rabiner and B.H. Juang, "Fundamentals of Speech Recognition". Prentice Hall, Englewood Cliffs, N.J., 1993
- [13] Makoto Iwayama, Nitin Indurkha and Hiroshi Motoda, "A new algorithm for automatic configuration of hidden Markov models". *Proc. 4th International Workshop on Algorithmic Learning Theory (ALT'93)*, 237-250 (1993)
- [14] A. Kaltenmeier, T. Caesar, J.M. Gloger and E. Mandler, "Sophisticated topology of hidden Markov models for cursive script recognition". *Proceedings of the 2nd International Conference on Document Analysis and Recognition*, 139-142 (1993)
- [15] P. Baldi and Y. Chauvin, "Smooth on-line learning algorithms for hidden Markov models". *Neural Computation*, Vol.6, 307-318 (1994)

- [16] S.B. Cho and J.H. Kim, "An HMM/MLP architecture for sequence recognition". *Neural Computation*, Vol.7, 358–369 (1995)
- [17] J. Dai, "Robust Estimation of HMM parameters using fuzzy vector quantization and Parzen's window". *Pattern Recognition*, Vol.28, No.1, 53–57 (1995)
- [18] A. Kundu, Y. He and P. Bahl, "Recognition of handwritten word: first and second order hidden Markov model based approach". *Pattern Recognition*, Vol.22, No.3, 283–297 (1989)
- [19] S.R. Veltman and R. Prasad, "Hidden Markov models applied to on-line handwritten isolated character recognition". *IEEE Trans. Image Processing*, Vol.3, No.3, 314–318 (1994)
- [20] E.J. Bellegarda, J.R. Bellegarda, D. Nahamoo and K.S. Nathan, "A fast statistical mixture algorithm for on-line handwriting recognition". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol.16, No.12, 1227–1233 (1994)
- [21] G. Rigoll, A. Kosmala, J. Rottland, and Ch. Neukirchen, "A Comparison between continuous and discrete density hidden Markov models for cursive handwriting recognition". *Proceedings of ICPR'96*, 205–209 (1996)
- [22] J. Hu, M.K. Brown and W.Turin, "HMM based on-line handwriting recognition". *IEEE transactions on Pattern Analysis and Machine Intelligence*, Vol.18, No.10, 1039-1045 (1996)
- [23] X. Li, M. Parizeau and R. Plamondon, "Hidden Markov model multiple observation training". *Technical Report EPM/RT-99/16*, November 1999, 16 pages.