

**Ath. Kehagias.**

**"Optimal Control for Training: the Missing Link between Hidden Markov Models and Connectionist Networks".**

**This paper has appeared in the journal:  
Mathematical and Computer Modelling, Vol. 14, pp.284-289, 1990.**

# OPTIMAL CONTROL INTERPRETATION OF THE BACKWARD-FORWARD ALGORITHM

Athanasios Kehagias

Division of Electronics and Computer Engineering

Department of Electrical Engineering

Aristotle University of Thessaloniki

Thessaloniki, Greece

June 17, 1997

## Abstract

The Backward-Forward algorithm is a highly efficient method for training Hidden Markov Models. Its name derives (and its implementation depends) on two types of probabilities: the forward propagating  $\alpha$  probabilities and the backward propagating  $\beta$  probabilities. In this note we formulate the training problem as an optimal control problem and show that the  $\alpha$  and  $\beta$  probabilities are conjugate quantities, similar to position and momentum variables in the Hamiltonian formulation of optimal control problems.

This paper appeared in Mathematical and Computer Modelling, Vol.14, pp.284-289, 1990.

**Introduction:** Hidden Markov Models (HMM) are currently the most popular method of modelling speech signals and find wide application in speech recognition problems. Other applications of HMM include shape recognition, medical and biological problems. A significant reason for the popularity and

success of HMM's is the availability of a very efficient training algorithm, the Backward - Forward (BF) algorithm of Baum. A good exposition of this is found in [2].

The BF algorithm makes use of two types of probabilities: the forward propagating  $\alpha$  probabilities and the backward propagating  $\beta$  probabilities; these will be defined presently. The derivation of the algorithm and the use of the  $\alpha$  and  $\beta$  probabilities is ingenious, but also *ad hoc*. In this note we first introduce the forward probabilities and note that they satisfy a linear forward evolution equation; therefore we interpret them as state variables of a linear system. Then we note that the training problem can be formulated as an optimal control problem and solve it using the Pontryagin Maximum Principle [3]. In the course of the solution, the  $\beta$  probabilities emerge as Lagrange multipliers of the problem. This is a new interpretation of the  $\alpha$  and  $\beta$  probabilities; it has independent interest and also suggests that algorithms from the control literature may prove viable alternatives of the BF algorithm.

**Optimal Control Formulation of HMM Training:** The following review of HMM's follows [2]. A HMM is a pair of stochastic processes:  $X_t, t = 1, 2, \dots, T, O_t, t = 1, 2, \dots, T$ .  $X_t$  (the hidden process) is a Markov process: it takes values in the set  $\mathbf{x} = \{x_1, \dots, x_N\}$ ; for  $i = 1, \dots, N, t = 1, 2, \dots, T$  we have

$$Prob(X_1 = x_i) = c_i \tag{1}$$

$$Prob(X_{t+1} = x_j \mid X_t = x_i) = a_{ij} \tag{2}$$

Note that (1) and (2) are sufficient to compute  $Prob(X_t = x_j)$ .  $O_t$  is the observed process and takes values in  $\mathbf{o} = \{o_1, o_2, \dots, o_M\}$  locally depending on  $X_t$ :

$$Prob(O_t = o_k \mid X_t = x_i) = b_{ik} \tag{3}$$

We collect the  $a$ 's,  $b$ 's and  $c$ 's in appropriately dimensioned matrices  $A, B, C$  and define the triple  $\mathbf{M} \doteq (A, B, C)$ .  $\mathbf{M}$  specifies a Hidden Markov Model.

Now, suppose we observe  $O_1 = o_{i1}, O_2 = o_{i2}, \dots, O_T = o_{iT}$  (for some integers  $i1, i2, \dots, iT$ ). The following probabilities (dependent on the observations) are very useful for the recognition and training problems:

$$\alpha_t(i) \doteq Prob(O_1 = o_{i1}, O_2 = o_{i2}, \dots, O_T = o_{iT}, X_t = x_i) \quad (4)$$

$$\alpha_t \doteq [\alpha_t(i) \dots \alpha_t(N)]' \quad (5)$$

We will call these quantities  $\alpha$ -probabilities or forward probabilities. They are of interest, because

$$Prob(O_1 = o_{i1}, O_2 = o_{i2}, \dots, O_T = o_{iT}) = \sum_{i=1}^N \alpha_T(i). \quad (6)$$

For a given observation sequence  $o_{i1}, \dots, o_{iT}$

$$J_0(\mathbf{M}) \doteq \sum_{i=1}^N \alpha_T(i) \quad (7)$$

is the probability of this particular sequence occurring, having assumed a particular model  $\mathbf{M}$  to be true. In recognition, we assume a fixed model and try to find the sequence of highest probability. In training, we fix a sequence  $o_{i1}, \dots, o_{iT}$  and look for the model  $\mathbf{M}$  that maximizes the probability  $J_0(\mathbf{M})$ .

Now we will formulate the training problem as an optimal control problem. It is easy to see that the  $\alpha$ -probabilities evolve recursively. For  $i = 1, 2, \dots, N, t = 1, 2, \dots$

$$\alpha_1(i) = g_i(O_1)c_i \quad (8)$$

$$\alpha_{t+1}(i) = \sum_{j=1}^N a_{ji} g_i(O_{t+1}) \alpha_t(j). \quad (9)$$

Here  $g_i : \mathbf{o} \rightarrow [0, 1]$  is a **nonlinear** (decision) function that is defined as follows:  $g_i(v) = b_{ik}$  iff  $v = o_k$ .

That is,  $g_i$  takes as input an observation and gives as an output the probability that this observation

occurred while  $X_t = x_i$ . Using matrix notation, we construct appropriate matrices  $\Phi(t, \mathbf{M})$  and (for a fixed observation sequence) we write (8, 9), for  $t = 1, 2, \dots, T$

$$\alpha_1 = \Phi(1, \mathbf{M}) \quad (10)$$

$$\alpha_{t+1} = \Phi(t, \mathbf{M}) \cdot \alpha_t. \quad (11)$$

Equations (10,11) describe a dynamical system, controlled by  $\mathbf{M} = (A, B, C)$ . The object of the control is to maximize

$$J_0(\mathbf{M}) = \sum_{i=1}^N \alpha_T(i). \quad (12)$$

We proceed to solve the optimal control problem by Pontryagin's Maximum Principle [3]. We must take into account two types of constraints:

1. The controls  $\mathbf{M} = (A, B, C)$  are not free quantities. Since their elements are probabilities, we must have the following conditions:

$$\sum_j a_{ij} = 1, \quad \sum_j b_{ij} = 1, \quad \sum_j c_j = 1 \quad (13)$$

$$a_{ij} \geq 0, \quad b_{ij} \geq 0, \quad c_j \geq 0. \quad (14)$$

However, by the Maximum Principle can free and constrained controls are treated in the same manner, so we will ignore this constraint for the time being.

2. The second kind of constraints is expressed by (8,11)

$$\alpha_1 = \Phi(1, \mathbf{M})$$

$$\alpha_{t+1} = \Phi(t, \mathbf{M}) \cdot \alpha_t.$$

We will treat these constraints explicitly, by use of the Maximum Principle.

We augment the cost function  $J_0$  with (10,11), using Lagrange multipliers  $\beta_t$ ,  $t = 1, 2, \dots, T$ :

$$J = J_0 - \sum_{t=1}^{T-1} \beta_{t+1} \cdot (a_{t+1} - \Phi(t)\alpha_t) - \beta_1 \cdot (\alpha_1 - \Phi(1, \mathbf{M})). \quad (15)$$

By the Maximum Principle, the introduction of Lagrange multipliers allows us to proceed as if the problem is unconstrained. At a maximizing point the following conditions must be necessarily satisfied:

$$\frac{\partial J}{\partial \alpha_t} = 0 \quad (16)$$

$$\frac{\partial J}{\partial A} = 0 \quad (17)$$

$$\frac{\partial J}{\partial B} = 0 \quad (18)$$

$$\frac{\partial J}{\partial C} = 0 \quad (19)$$

Let us first look at (16)). For  $t = 1$ , it implies  $\beta_1 = \beta_2 \cdot \Phi(2\mathbf{M})$ . For  $t \geq 2$ , define the *Hamiltonian* function

$$H_t \doteq \beta_{t+1} \cdot \Phi(t, \mathbf{M})\alpha_t. \quad (20)$$

Now (16) becomes

$$\beta_t = \frac{\partial H_t}{\partial \alpha_t} \Rightarrow \quad (21)$$

$$\beta_t = \beta_{t+1} \cdot \Phi(t+1, \mathbf{M}). \quad (22)$$

We also have

$$\frac{\partial J}{\partial a_T} = 0 \Rightarrow \quad (23)$$

$$\beta_T = [1 \ 1 \dots 1]. \quad (24)$$

Equations (22,24) allow us to solve for  $\beta_t$  backwards. Equation (22) can be written for  $t = 1, 2, \dots, T-1$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} g_j(O_t) \beta_{t+1}(j) \quad (25)$$

Using (25) and the final condition (24) we can that  $\beta_t(i) = \text{Prob}(O_{t+1} = o_{i,t+1}, \dots, O_T = o_{iT})$ . In [2] things are done in the reverse way. Starting from the above definition of  $\beta_t$ , it is proven that eq.(25) has to be satisfied. The reader familiar with the Backward-Forward algorithm will recognize  $\beta_t$  as the backward probabilities, or  $\beta$ -probabilities.

We also see that  $\alpha$  and  $\beta$  have a Hamiltonian structure. Equations (10,11) can be written for  $t = 2, \dots, T$

$$\alpha_t = \frac{\partial H_{t-1}}{\partial \beta_t}. \quad (26)$$

Equations (21,26) show clearly the Hamiltonian structure:  $\alpha$  and  $\beta$  are conjugate quantities (like position and momentum in classical mechanics).

**Choice of Control:** So far we have not said anything about the actual selection of the control variables (whether these are given in terms of  $A, B, C$  or  $Q, R, S$ ). Equations (17-19) are necessary conditions that the optimal control has to satisfy, but obtaining  $A, B, C$  from them is not trivial. As a matter of fact, there is no canonical way to solve the general Optimal Control problem. Many optimization algorithms are reported in the optimal control literature (see [1]). On the other hand, in the HMM community, the method of choice is the Backward-Forward algorithm which consists of the following iteration:

Given  $\mathbf{M}^n = (A^n, B^n, C^n)$ ,

$$a_{ij}^{n+1} = \frac{\sum_{t=1}^{T-1} \alpha_t^n(i) a_{ij}^n g_j^n(O_{t+1}) \beta_{t+1}^n(j)}{\sum_{t=1}^{T-1} \alpha_t^n(i) \beta_t^n(i)} \quad (27)$$

$$b_{ij}^{n+1} = \frac{\sum_{t:O_t=o_j} \alpha_t^n(i) \beta_{t+1}^n(i)}{\sum_{t=1}^{T-1} \alpha_t^n(i) \beta_t^n(i)} \quad (28)$$

$$c_i^{n+1} = \frac{\alpha_i^n(i) \beta_i^n(i)}{\sum_{i=1}^N \alpha_T^n(i)} \quad (29)$$

This is an ascent procedure (not necessarily steepest ascent). As proven in [2], for all  $n \geq 1$

$$J(\mathbf{M}^{n+1}) \geq J(\mathbf{M}^n) \quad (30)$$

and in fact we have guaranteed convergence to a local maximum.

**Conclusion:** The optimal control formulation of the HMM training problem is independent of the optimization method. Hence the new formulation achieves two things. First, it offers an alternative interpretation of the forward and backward probabilities. Second, it suggests a comparison of the efficiency of optimal control optimization algorithms (such as the Kelly - Bryson method [1]) to the BF algorithm. This comparison may yield practically significant results. However this will be investigated elsewhere.

## References

- [1] A.E. Bryson and Y.-C. Ho, *Applied Optimal Control*, Blaisdell, Waltham, 1969.
- [2] S.E. Levinson, L.R. Rabiner and M.M. Sondhi, An Introduction to the Application of the Theory of Probabilistic Functions of a Markov Process to Automatic Speech Recognition, *Bell Sys. Tech. J.*, Vol. 62, No. 4, April 1983.
- [3] A.P. Sage, *Optimum Systems Control*, Prentice-Hall, Englewood Cliffs, 1968.