



# JOINT TRAINING OF MULTI-STREAM HMMs

Samy Bengio <sup>a</sup>

IDIAP-RR 05-22

MAY 2005

---

<sup>a</sup> IDIAP Research Institute, P.O. Box 592, CH-1920 Martigny, Switzerland



# JOINT TRAINING OF MULTI-STREAM HMMs

Samy Bengio

MAY 2005

**Abstract.** This report describes a novel technique to jointly train efficiently several streams of data describing the same sequence of events using a unified EM algorithm.

## 1 Introduction

Multi-stream HMMs is a very popular and efficient technique to decompose a complex problem such as automatic speech recognition into several streams of information, where each stream of data is trained independently of the others, while all streams are jointly used during decoding.

In this report, we propose a method to trained all streams simultaneously in order to maximize the joint probability of all data streams. This can be done by adding an additional, but reasonable assumption, which states that given the value of the state variable at time  $t$ , all streams up to time  $t$  are independent from each other (hence, the state variable contains all the correlation between all streams up to time  $t$ ).

## 2 Notation

Let

- $N$  be the number of streams.
- $x_{a:b}^n$  be the observations of stream  $n$  between time  $a$  and  $b$

## 3 Usual Derivation for Each Stream

Let us start by computing the normal forward and backward equations for each stream. First the forward equation:

$$\alpha_n(i, t) = p(x_{1:t}^n, q_t = i) \quad (1)$$

$$= p(x_t^n | q_t = i) \sum_j p(q_t = i | q_{t-1} = j) \alpha_n(j, t-1) \quad (2)$$

which can be used to compute the likelihood of one stream:

$$L_n = p(x_{1:T}^n) \quad (3)$$

$$= \sum_i p(x_{1:T}^n, q_T = i) = \sum_i \alpha_n(i, T) \quad (4)$$

and then the backward equation:

$$\beta_n(i, t) = p(x_{t+1:T}^n | q_t = i) \quad (5)$$

$$= \sum_j p(x_{t+1}^n | q_{t+1} = j) p(q_{t+1} = j | q_t = i) \beta_n(j, t+1) \quad (6)$$

which, together with the forward equation, can be used to compute the posterior probability of a state given the stream observation:

$$\gamma_n(i, t) = p(q_t = i | x_{1:T}^n) \quad (7)$$

$$\begin{aligned} &= \frac{p(x_{1:t}^n, x_{t+1:T}^n, q_t = i)}{p(x_{1:T}^n)} \\ &= \frac{\alpha_n(i, t) \beta_n(i, t)}{L_n} \end{aligned} \quad (8)$$

## 4 Derivation for Joint Streams

In this section, we derive the new EM training algorithm for multi-streams.

On top of the usual HMM assumptions accepted for the single stream case (emissions at time  $t$  only depend on the state, and the current state only depends on the previous state), we need to add one more assumption, which basically says that, given the value of the current state, the random variables of all streams up to time  $t$  are independent:

$$p(x_{1:t}^1, x_{1:t}^2, \dots, x_{1:t}^N | q_t = i) = \prod_{n=1}^N p(x_{1:t}^n | q_t = i) \quad (9)$$

### 4.1 Joint Likelihood and Auxiliary Function

As usual when deriving an EM algorithm, we start by defining the following indicator variable:

$$z_{i,t} = \begin{cases} 1 & \text{if } q_t = i \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

The joint complete log likelihood of the streams and the state  $l_c$  is:

$$l_c = \log p(X^1, X^2, \dots, X^N, Q) \quad (11)$$

$$= \log p(X^1, X^2, \dots, X^N | Q) + \log P(Q)$$

$$= \log p(X^1 | Q) + \log p(X^2 | Q) + \dots + \log p(X^N | Q) + \log P(Q)$$

$$= \sum_n \log p(X^n | Q) + \log P(Q)$$

$$= \sum_t \sum_n \sum_i z_{i,t} \log p(x_t^n | q_t = i) +$$

$$\sum_t \sum_i \sum_j z_{i,t} z_{j,t-1} \log p(q_t = i | q_{t-1} = j) \quad (12)$$

and the corresponding auxiliary variable becomes:

$$A = E_Q [\log p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N, Q) | X^{1:N}] \quad (13)$$

$$= \sum_t \sum_n \sum_i E_Q [z_{i,t} | X^{1:N}] \log p(x_t^n | q_t = i) +$$

$$\sum_t \sum_i \sum_j E_Q [z_{i,t} z_{j,t-1} | X^{1:N}] \log p(q_t = i | q_{t-1} = j) \quad (14)$$

The forward recursion becomes:

$$\alpha(i, t) = p(x_{1:t}^1, x_{1:t}^2, \dots, x_{1:t}^N, q_t = i) \quad (15)$$

$$= p(x_{1:t}^1, x_{1:t}^2, \dots, x_{1:t}^N | q_t = i) p(q_t = i)$$

$$= p(x_{1:t}^1 | q_t = i) p(x_{1:t}^2 | q_t = i) \dots p(x_{1:t}^N | q_t = i) p(q_t = i)$$

$$= \frac{p(x_{1:t}^1 | q_t = i) p(q_t = i)}{p(q_t = i)} \frac{p(x_{1:t}^2 | q_t = i) p(q_t = i)}{p(q_t = i)} \dots \frac{p(x_{1:t}^N | q_t = i) p(q_t = i)}{p(q_t = i)} p(q_t = i)$$

$$= \frac{p(x_{1:t}^1 | q_t = i)}{p(q_t = i)} \frac{p(x_{1:t}^2 | q_t = i)}{p(q_t = i)} \dots \frac{p(x_{1:t}^N | q_t = i)}{p(q_t = i)} p(q_t = i)$$

$$= \frac{\alpha_1(i, t)}{p(q_t = i)} \frac{\alpha_2(i, t)}{p(q_t = i)} \dots \frac{\alpha_N(i, t)}{p(q_t = i)} p(q_t = i)$$

$$= \frac{\prod_{n=1}^N \alpha_n(i, t)}{p(q_t = i)^{N-1}} \quad (16)$$

We can recover the joint likelihood of the streams as:

$$L = p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N) \quad (17)$$

$$\begin{aligned} &= \sum_i p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N, q_T = i) \\ &= \sum_i \alpha(i, T) \end{aligned} \quad (18)$$

Let us also compute the joint posterior probability of being in a state given all data of all streams, the so-called gamma:

$$E_Q [z_{i,t} | X^{1:N}] = \gamma(i, t) \quad (19)$$

$$\begin{aligned} \gamma(i, t) &= p(q_t = i | x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N) \quad (20) \\ &= \frac{p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N, q_t = i)}{p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N)} \\ &= \frac{p(x_{t+1:T}^1 | q_t = i) \cdots p(x_{t+1:T}^N | q_t = i) p(x_{1:t}^1, \dots, x_{1:t}^N, q_t = i)}{L} \\ &= \frac{\beta_1(i, t) \cdots \beta_N(i, t) \alpha(i, t)}{L} \\ &= \frac{\alpha(i, t) \prod_n \beta_n(i, t)}{L} \end{aligned} \quad (21)$$

Note that we need to compute  $p(q_t = i)$ , which can be done recursively as follows:

$$p(q_t = i) = \sum_j p(q_t = i, q_{t-1} = j) \quad (22)$$

$$= \sum_j p(q_t = i | q_{t-1} = j) p(q_{t-1} = j) \quad (23)$$

We can also similarly compute the posterior of a transition as follows:

$$\begin{aligned} E_Q [z_{i,t} z_{j,t-1} | X^{1:N}] &= p(q_t = i, q_{t-1} = j | x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N) \quad (24) \\ &= \frac{p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N, q_t = i, q_{t-1} = j)}{p(x_{1:T}^1, x_{1:T}^2, \dots, x_{1:T}^N)} \\ &= \frac{\alpha(j, t-1) p(q_t = i | q_{t-1} = j) \prod_n p(x_t^n | q_t = i) \beta_n(i, t)}{L} \end{aligned}$$

## 5 Conclusion

While most multi-stream HMM techniques assume that each stream is first trained independently, we have shown in this report a principled way to jointly trained all streams in order to optimize the joint probability of all streams.