

BIBLIOGRAPHY

- [1] Baum, Leonard E., Petrie, Ted, Soules, George, and Weiss, Norman, A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains, Annals of Mathematical Statistics, Vol. 41, No. 1, 1970, pp. 164-171.
- [2] Liporace, Louis A., Maximum likelihood analysis for multi-variate observations of Markov sources, (These Proceedings).
- [3] Liporace, Louis A., Statistical analysis of time series, (These Proceedings).
- [4] Ferguson, John D., Variable duration models for speech, (These Proceedings).
- [5] Markel, J.D., and Gray, A.H. Jr., Linear prediction of speech, Springer-Verlag, New York, (1976).

VARIABLE DURATION MODELS FOR SPEECH

John D. Ferguson

ABSTRACT

In this paper we discuss an extension of the basic theory of hidden Markov models to allow for the fact that speech sounds have varying durations. The extension introduces a distribution on the duration spent in a state, and allows a sequence of observations to be made while remaining in a single state.

VARIABLE DURATION MODELS FOR SPEECH

John D. Ferguson

0. Hidden Markov models

In the case of analyzing printed text using hidden Markov models ala Cave-Neuwirth [4], the observations are discrete and are always of exactly the same length: one character. In trying to analyze speech in a similar way, we have to allow for continuous measurements, and states persisting for various durations. The extension of the basic technique to continuous measurements was discussed in Laporace's paper [6]. Here we discuss the question of variable duration states.

We are interested in a stochastic process which produces observations in the following way. There is a Markov chain with S states, $i, i = 1, 2, \dots, S$, with initial distribution (a_i) , transition matrix (a_{ij}) , and output matrix $(b_j(k))$ which expresses the probability that the observation will be k , given state j . [We assume, for exposition, a finite output alphabet. We will remark later on the extension of the results obtained to the case of continuously distributed output.]

144

In the usual case [1], a sequence of observations $\sigma_1, \sigma_2, \dots, \sigma_T$ is produced by a hidden Markov model with one observation per state. The actual sequence of events is taken to be:

- (A) Choose an initial state i_1 according to the initial distribution a_i .
- (B) Choose an observation σ_1 according to the distribution $b_{i_1}(\sigma)$.
- (C) Choose the next state i_2 , according to the transition probability distribution $a_{i_1,j}$ ($j = 1, 2, \dots, S$).
- (D) Choose an observation σ_2 according to the output distribution $b_{i_2}(k)$, $k = 1, 2, \dots, K$, etc.

In the case of variable duration hidden Markov models, there is an additional ingredient: associated with each state, i , there is a duration distribution $P(d|i)$, ($d = 1, 2, \dots, D$) which determines the length of time spent in state i , i.e., the number of observations produced while in state i . The sequence of events is now:

- (A) Choose an initial state i_1 according to the initial distribution a_i .
- (B) Choose a duration d_1 according to the duration distribution $P(d|i_1)$.

- (C) choose observations $\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_{d_1}$ according to $b_{i_1}(\mathcal{O}_1, \dots, \mathcal{O}_{d_1})$. [This is often, for convenience, taken to be $\prod_{t=1}^{d_1} b_{i_1}(\mathcal{O}_t)$.]
- (D) Choose the next state i_2 according to the transition probabilities $a_{i_1, j}$ ($j = 1, 2, \dots, S$).
- (E) Choose a duration d_2 according to the duration distribution $P(d|i_2)$.

- (F) Choose observations $\mathcal{O}_{d_1+1}, \mathcal{O}_{d_1+2}, \dots, \mathcal{O}_{d_1+d_2}$ according to $b_{i_2}(\mathcal{O}_{d_1+1}, \dots, \mathcal{O}_{d_1+d_2})$, etc.

It is important to note that there is no external demarcation between the observations arising from state i_1 and those arising from i_2 .

Let us note that the variable duration model includes the usual model in either of the following two ways:

- (i) For each state j , set $P(d=1|j) = 1$.

Thus each state lasts exactly one time, and we return to the usual model.

- (ii) For each state, j , set $P(d|j) = \rho_j^{d-1}(1-\rho_j)$ for some ρ_j between 0 and 1. Thus the duration distribution is geometric. We also require $a_{i,j} = 0$ for each i . Such a model then mimics the usual case.

Conversely, with considerable labor, the variable duration case can be realized in the old framework. Take as states the triples (i, x, d) where i is one of the variable duration states, d a duration, and x a counter, $1 \leq x \leq d$, which indicates the number of observations produced so far while in state i .

Transitions are of the form $(i, x, d) \rightarrow (i, x+1, d)$, unless $x = d$, when $(i, d, d) \rightarrow (j, 1, d')$ for some j , some d' .

In the Poritz LPC model, described in an earlier paper in this collection, [5], it would be a definite advantage to allow for variable duration. It is clear that phones or phone categories have characteristic duration distributions which are not, in general, geometric. Thus the basic model, without variable duration, should be somewhat ill-fitting.

As usual in the study of hidden Markov models, there are three basic problems:

- (1) Given $\mathcal{O}_1, \dots, \mathcal{O}_T$ and a model, $(a_i, (a_{ij}), (b_j(\mathcal{O})), (P(d|i)))$ compute the likelihood of the sequence $\mathcal{O}_1, \dots, \mathcal{O}_T$.
- (2) Given $\mathcal{O}_1, \dots, \mathcal{O}_T$ and a state space of size S , find the maximum likelihood model.

(iii) Given $\mathcal{O}_1, \dots, \mathcal{O}_T$ and a model, find the most likely state sequence either in the sense of the single state sequence or maximum posterior probability, or the state sequence with maximum per state posterior probability.

In the case of variable duration hidden Markov models, this problem is complicated by decisions about the durations of the individual states.

We will discuss all three of these problems for our case.

1. Computation of likelihood for variable duration models

Suppose we have observations $\mathcal{O} = \mathcal{O}_1, \dots, \mathcal{O}_T$ and a variable duration model. We wish to compute the likelihood of the sequence \mathcal{O} given the model. We will make the simplifying assumption that the first state commences at time 1, and the last state ends at time T [no observations are "lost" on either end]. It is messy, but possible, to handle the situation where the first state begins at or before observation \mathcal{O}_1 , and the last state ends at or after observation \mathcal{O}_T . We will not discuss this complication here.

To compute the likelihood, and for later use, we define an analogue of the usual alpha calculation [1].

Let

$$\alpha_t(i) = \Pr(\mathcal{O}_1 \dots \mathcal{O}_t \text{ and state } i \text{ terminates at } t).$$

This probability is a sum over paths

$$p = (i_1, d_1), \dots, (i_r, d_r) \text{ such that } i_r = i \text{ and } d_1 + \dots + d_r = t,$$

of the joint probability of the path and the observations. Thus

$$\alpha_t(i) = \sum_p \Pr(p \& \mathcal{O}_1 \mathcal{O}_2 \dots \mathcal{O}_t)$$

$$\begin{aligned} &= \sum_p a_{i_1} \cdot \Pr(d_1 | i_1) \cdot \Pr(\mathcal{O}_1 \dots \mathcal{O}_{d_1} | i_1) \cdot \\ &\quad a_{i_1 i_2} \cdot \Pr(d_2 | i_2) \cdot \Pr(\mathcal{O}_{d_1+1} \dots \mathcal{O}_{d_1+d_2} | i_2) \dots \\ &\quad a_{i_{r-1} i_r} \cdot \Pr(d_r | i_r) \cdot \\ &\quad \Pr(\mathcal{O}_{d_1+\dots+d_{r-1}+1} \dots \mathcal{O}_t | i_r) \end{aligned}$$

We note that, as in the usual case, there is an effectively computable inductive process for calculating the alphas. The likelihood we want is simply $\sum_i \alpha_T(i)$.

We have, inductively:

$$\alpha_t(j) = \sum_{d=1}^S \sum_{i=1}^D \alpha_{t-d}(i) \cdot a_{ij} \cdot P(d|j)$$

$$\cdot P(\sigma_{t-d+1} \dots \sigma_t | j)$$

ALPHA INDUCTIVE FORMULA

To complete the calculation of likelihood, given a model, we have to show how to initialize the alphas.

Since we have assumed that some state must commence at time 1, we can compute:

$$\alpha_1(i) = a_i P(d=1|i) \cdot P(\sigma_1|i), \quad i = 1, 2, \dots, S,$$

$$\alpha_2(1) = a_1 P(d=2|1) \cdot P(\sigma_1 \sigma_2|1)$$

$$+ \sum_j \alpha_1(j) \cdot a_{j1} \cdot P(d=1|i) P(\sigma_2|1),$$

$$\alpha_3(1) = a_1 P(d=3|1) \cdot P(\sigma_1 \sigma_2 \sigma_3|1)$$

$$+ \sum_{d < 3} \sum_j \alpha_{3-d}(j) a_{ji} P(d|1) P(\sigma_{4-d} \dots \sigma_3|1)$$

etc. until $\alpha_D(1)$ is completed, after which we use the basic inductive formula. [There is an alternative approach involving a dummy array $\alpha_0(1), \alpha_{-1}(1), \dots, \alpha_{-D+1}(1)$ and dummy observations $\sigma_0, \sigma_{-1}, \dots, \sigma_{-D+1}$.]

For speech, the results of this section can be used as follows: given several competing models, say for different talkers, or for different phrases, and given a segment of speech, we can compute the probability of the segment given the model. Together with the prior probability of the model, this yields the posterior probability of the model given the segment.

2. Iterative method for obtaining the maximum likelihood model

In this section we show how the Baum algorithm [1] can be adapted to the case of variable duration hidden Markov models to provide an iterative method for finding the maximum likelihood model.

In the Poritz scheme, [5], extended, we would thus be able to estimate the LPC parameters for each state as well as the distribution on the duration for which these LPC parameters are to be held fixed.

Given any model M , with likelihood $L(M)$, the Baum algorithm produces another model $M' = \mathcal{T}(M)$ such that $L(M') > L(M)$, unless $M = M'$ in which case M must be a critical point of the likelihood function. The algorithm is used iteratively. We start with a

crude model M_0 , and compute, successively, $M_1 = \mathcal{T}(M_0)$, $M_2 = \mathcal{T}(M_1)$, etc. Since the likelihood increases at each step, (unless a critical point is reached), then either the likelihood converges, or, in the continuous case, possibly approaches $+\infty$.

3. Computations

We discuss how we can use this algorithm in our case.

In the following, the dependence of all probabilities on the current model is suppressed.

We will need three other inductively computable

arrays: $\alpha_t^*(i)$, $\beta_t^*(i)$, $\beta_t(i)$ for $t = 0, 1, \dots, T$, $i = 1, \dots, S$.

We define

$$\begin{aligned} \alpha_t^*(i) &= P(\sigma_1 \dots \sigma_t \text{ & } i \text{ commences at } t+1), \\ \beta_t(i) &= P(\sigma_{t+1} \dots \sigma_T | i \text{ terminates at } t), \\ \beta_t^*(i) &= P(\sigma_{t+1} \dots \sigma_T | i \text{ commences at } t+1). \end{aligned}$$

Then we have:

$$(i) \quad \alpha_t^*(j) = \sum_i \alpha_t(i) a_{ij},$$

$$(ii) \quad \alpha_t(i) = \sum_d \alpha_{t-d}^*(i) \cdot P(d|i) \cdot P(\sigma_{t-d+1} \dots \sigma_t | i),$$

= Expected number of times state i started at time t or before, given σ ,

$$(iii) \quad E_t(i) = \sum_{\tau < t} \alpha_\tau(i) \beta_\tau(i) / P$$

= Expected number of times state i terminated before time t , given σ ,

$$(iv) \quad \beta_t^*(i) = \sum_d \beta_{t+d}(i) P(d|i) \cdot P(\sigma_{t+1} \dots \sigma_{t+d} | i).$$

$$(E) \quad r_t(i) = S_t(i) - E_t(i)$$

= Prob(state i was used to produce observation σ_t , given σ),
 [Hence $\forall t$, $\sum_i r_t(i) = 1$.]

$$(F) \quad \sum_t r_t(i) = \text{Expected total duration spent in state } i, \text{ given } \sigma,$$

(G) $S_T(i) = \text{Expected total number of times state } i \text{ commenced, given } \sigma$
 $= E_{T+1}(i) = \text{expected number of times state } i \text{ terminated, given } \sigma$

$$(H) \quad \sum_t r_t(i)/S_T(i) = \text{estimated average duration of state } i, \text{ given } \sigma,$$

$$(I) \quad c(i) = \frac{a_i \beta_0^*(i)}{P} = \text{Probability that state } i$$

was the first state, given σ ,

$$(J) \quad C(i,j) = \sum_t \frac{\alpha_t(i)a_{ij}\beta_t^*(j)}{P} = \text{Expected number of transitions from } i \text{ to } j, \text{ given } \sigma,$$

since it is not necessarily true that some state must end at time t .

For (E), we argue as follows: the probability that state i occurred with duration d , given σ ,

$$(K) \quad C(i,d) = \sum_t \alpha_t^*(i)P(d|i) \cdot P(\sigma_{t+1} \dots \sigma_{t+d}|i)$$

= Expected number of times

state i occurred with duration d , given σ ,

$$(L) \quad C(i,k) = \sum_{t \ni \sigma_t=k} r_t(i)$$

= Expected number of times that state i occurred with observation $\sigma_t = k$, given σ .

Most of these relations are quite simple. Let us discuss only (A) and (E).

For (A), $\alpha_t(i)\beta_t(i)$ is the sum of the probabilities for all paths (i.e., state sequences with durations) which have the property that state i terminates at time t . This holds since $\alpha_t(i)$ accounts for all paths to the left of t , ending in i , and $\beta_t(i)$ accounts for all paths to the right, given i ends at t , and any pairs can be combined.

This is in direct analogy to [1], but note that

$$\sum_i \frac{\alpha_t(i)\beta_t(i)}{P} \leq 1 ,$$

$$\sum_i \frac{\alpha_t(i)\beta_t(i)}{P} \leq 1 ,$$

for (E), we argue as follows: the probability that state i is used at time t is

$$\sum_{\tau < t} \sum_{d \geq t \leq \tau+d} \frac{\alpha_{\tau}^*(1) P(d|1) \cdot P(O_{\tau+1} \dots O_{\tau+d} | 1) \beta_{\tau+d}(1)}{P},$$

by examining paths, as above.

This can be rewritten:

$$\sum_{\tau < t} \sum_{\text{all } d} \frac{\alpha_{\tau}^*(1) P(d|1) \cdot P(O_{\tau+1} \dots O_{\tau+d} | 1) \beta_{\tau+d}(1)}{P}$$

$$- \sum_{\tau < t} \sum_{\substack{\text{all } d \\ t > \tau+d}} \frac{\alpha_{\tau}^*(1) P(d|1) \cdot P(O_{\tau+1} \dots O_{\tau+d} | 1) \beta_{\tau+d}(1)}{P}$$

The first quantity is

$$\sum_{\tau < t} \frac{\alpha_{\tau}^*(1) \cdot \beta_{\tau}^*(1)}{P},$$

using the inductive formula (iv). The second quantity, replacing $\tau + d$ by θ , τ by $\theta - d$, is

$$\sum_{\theta < t} \sum_{d} \frac{\alpha_{\theta-d}^*(1) \cdot P(d|1) \cdot P(O_{\theta-d+1} \dots O_{\theta} | 1) \beta_{\theta}(1)}{P}$$

= (using inductive formula (ii))

$$= \sum_{\theta < t} \frac{\alpha_{\theta}(1) \beta_{\theta}(1)}{P},$$

as asserted.

Now we show how to use these various computed quantities to reestimate the parameters of the model.

The new initial probabilities are

$$\hat{a}_1 = c(1) \quad (\text{see (I)}).$$

The new transition probabilities are

$$\hat{a}_{1,j} = c(1,j) / \sum_j c(1,j) \quad (\text{see (J)}).$$

The new duration probabilities are

$$\hat{P}(d|1) = c(1,d) / \sum_d c(1,d) \quad (\text{see (K)}).$$

For the reestimates of observation probabilities, we make the assumption of independence, although other hypotheses can be dealt with.

The new observation probabilities are

$$\hat{P}(O_k | 1) = c(1,k) / \sum_k c(1,k) \quad (\text{see (L)}).$$

4. Proof of increasing likelihood

We show how the proof of [1] that the likelihood must increase can be adapted to fit our case.

Let (cf. [1]) λ be a generic symbol for the parameters of our model. Thus $\lambda = ((a_i), (a_{ij}), (P(d|i)), (P(k|i)))$. Define

$$Q(\lambda, \lambda') = \sum_{\text{Paths } p} P_\lambda(p \& \mathcal{O}) \log P_{\lambda'}(p \& \mathcal{O}).$$

The sum is over all paths

$$p : (i_1, d_1), (i_2, d_2), \dots, (i_N, d_N) \ni \sum_{T=1}^N d_T = T.$$

The first probability, $P_\lambda(p \& \mathcal{O})$, denotes the probability of the path p and the observation sequence \mathcal{O} , given the parameter set λ . The second probability $P_{\lambda'}(p \& \mathcal{O})$ denotes the probability of the same path and \mathcal{O} , given a second parameter set λ' . This auxiliary function greatly facilitates the proof of increasing likelihood.

Lemma 1. $Q(\lambda, \lambda') - Q(\lambda, \lambda) \leq P_{\lambda'}(\mathcal{O}) - P_\lambda(\mathcal{O})$. The inequality is strict unless, for each path, p ,

$$P_\lambda(p \& \mathcal{O}) = P_{\lambda'}(p \& \mathcal{O}).$$

Proof. We use the elementary inequality $\log x \leq x - 1$, which is strict unless $x = 1$.

Thus,

$$\begin{aligned} Q(\lambda, \lambda') - Q(\lambda, \lambda) &= \sum_p P_\lambda(p \& \mathcal{O}) \ln \frac{P_{\lambda'}(p \& \mathcal{O})}{P_\lambda(p \& \mathcal{O})} \\ &\leq \sum_p P_\lambda(p \& \mathcal{O}) \left(\frac{P_{\lambda'}(p \& \mathcal{O})}{P_\lambda(p \& \mathcal{O})} - 1 \right) = P_{\lambda'}(\mathcal{O}) - P_\lambda(\mathcal{O}). \end{aligned}$$

As shown in [1], a stronger inequality holds:

$$Q(\lambda, \lambda') - Q(\lambda, \lambda) \leq P_{\lambda'}(\mathcal{O}) \ln \frac{P_{\lambda'}(\mathcal{O})}{P_\lambda(\mathcal{O})},$$

but Lemma 1 will suffice for our purposes.

This lemma relates the Q function to the likelihood. It indicates that, for fixed λ , if we can find any $\lambda' \ni Q(\lambda, \lambda') > Q(\lambda, \lambda)$ then $P_{\lambda'}(\mathcal{O}) > P_\lambda(\mathcal{O})$. One natural method is therefore to maximize $Q(\lambda, \lambda')$ as a function of λ' .

Lemma 2.

$$\begin{aligned} \frac{Q(\lambda, \lambda') - Q(\lambda, \lambda)}{P_\lambda(\mathcal{O})} &= \sum_i C_\lambda(i) \ln \frac{a'_i}{a_i} + \sum_{i,j} C_\lambda(i, j) \ln \frac{a'_{ij}}{a_{ij}} \\ &+ \sum_{i,d} C_\lambda(i, d) \ln \frac{P'(d|i)}{P(d|i)} + \sum_{i,k} C_\lambda(i, k) \ln \frac{P'(k|i)}{P(k|i)}. \end{aligned}$$

[Here the C 's refer to the counts made for reestimation, based on parameters λ , described in Section 3, (I), (J), (K), (L).]

Proof. A path can be defined by a state sequence,

with associated durations:

$$(i_1, d_1), (i_2, d_2), \dots, (i_N, d_N) \Rightarrow \sum_{\tau} d_{\tau} = T .$$

Given a path, p , the observation string is decomposed into segments $\sigma_p^{(1)} = \sigma_1 \dots \sigma_{d_1}$, $\sigma_p^{(2)} = \sigma_{d_1+1} \dots \sigma_{d_1+d_2}, \dots, \sigma_p^{(N)} = \sigma_{d_1+\dots+d_{N-1}+1} \dots \sigma_T$. Using this notation, we have:

$$Q(\lambda, \lambda') = \sum_p P_{\lambda}(p \& \sigma) \ln P_{\lambda'}(p \& \sigma)$$

$$= \sum_p P_{\lambda}(p \& \sigma) \ln \left\{ a'_{i_1} \cdot \prod_{\tau=2}^N a'_{i_{\tau-1} i_{\tau}} \cdot \right.$$

$$\left. \prod_{\tau=1}^N P'(d_{\tau} | i_{\tau}) \cdot \prod_{\tau=1}^N P'(\sigma_p^{(\tau)} | i_{\tau}) \right\}$$

$$= \sum_p P_{\lambda}(p \& \sigma) \ln a'_{i_1} + \sum_{\tau=2}^N \sum_p P_{\lambda}(p \& \sigma) \ln a'_{i_{\tau-1} i_{\tau}}$$

$$+ \frac{1}{P_{\lambda}(\sigma)} \sum_p \sum_{\tau=1}^N P_{\lambda}(p \& \sigma) \ln (P'(d_{\tau} | i_{\tau}) / P(d_{\tau} | i_{\tau}))$$

Each of the four summands here will be related to the corresponding summand in the conclusion of the lemma.

First of all,

$$\begin{aligned} & \frac{1}{P_{\lambda}(\sigma)} \sum_p P_{\lambda}(p \& \sigma) \ln (a'_{i_1} / a_{i_1}) \\ &= \frac{1}{P_{\lambda}(\sigma)} \sum_{i=1}^S \ln (a'_{i_1} / a_{i_1}) \sum_{p: i_1=i} P_{\lambda}(p \& \sigma) \\ &= \sum_{i=1}^S \ln (a'_{i_1} / a_{i_1}) \cdot \frac{a_{i_1} \cdot \beta_0^*(i)}{P_{\lambda}(\sigma)} = \sum_{i=1}^S \ln (a'_{i_1} / a_{i_1}) \cdot c_{\lambda}(i), \end{aligned}$$

according to (I).

Secondly,

$$\begin{aligned} & \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\tau=2}^N P_\lambda(p \& \mathcal{O}) \ln \left(\frac{a'_{i_{\tau-1} i_\tau}}{a_{1_{\tau-1} i_\tau}} \right) \\ & = \frac{1}{P_\lambda(\mathcal{O})} \sum_{i,j} \sum_{N} \sum_{\tau=2}^N \sum_{\substack{\text{p of length } N, \\ \exists i_{\tau-1} = i \\ i_\tau = j}} P_\lambda(p \& \mathcal{O}) \ln \left(\frac{a'_{i_{\tau-1} i_\tau}}{a_{i_{\tau-1} i_\tau}} \right) \end{aligned}$$

$$\begin{aligned} & = \sum_{i,j} \ln \left(\frac{a'_{i,j}}{a_{i,j}} \right) \cdot \frac{1}{P_\lambda(\mathcal{O})} \sum_{t=1}^{T-1} \sum_{\substack{\text{p} \ni t \text{ ends at } t \\ \text{j starts at } t+1}} P_\lambda(p \& \mathcal{O}) \end{aligned}$$

$$\begin{aligned} & = \sum_{i,j} \ln \left(\frac{a'_{i,j}}{a_{i,j}} \right) \cdot c_\lambda(i,j), \quad (\text{from (J)}) . \end{aligned}$$

Thirdly,

$$\begin{aligned} & \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\tau=1}^N P_\lambda(p \& \mathcal{O}) \ln \left(\frac{P'(\mathcal{O}_\tau | i_\tau)}{P(\mathcal{O}_\tau | i_\tau)} \right) \\ & = \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\tau=1}^N P_\lambda(p \& \mathcal{O}) \sum_{t \in U_\tau} \ln \frac{P'(\mathcal{O}_t | i_\tau)}{P(\mathcal{O}_t | i_\tau)} \end{aligned}$$

$$= \frac{1}{P_\lambda(\mathcal{O})} \sum_{i,d} \ln \left(\frac{P'(d|i)}{P(d|i)} \right) \sum_{\substack{p, \tau \\ i_\tau = i \\ d_\tau = d}} P_\lambda(p \& \mathcal{O})$$

$$= \sum_{i,k} \ln \frac{P'(k|i)}{P(k|i)} \cdot \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\substack{\tau \\ i_\tau = i \\ \exists \mathcal{O}_t = k}} \sum_{t \in U_\tau} P_\lambda(p \& \mathcal{O})$$

$$= \sum_{i,d} \ln \left(\frac{P'(d|i)}{P(d|i)} \right) \cdot c_\lambda(i,d), \quad (\text{from (K)}) .$$

Finally we consider the last sum. To complete the proof we recall the assumption of independence of the observations within a state. We also need some new notation. If p is a path, $p = (i_1, d_1), (i_2, d_2), \dots, (i_N, d_N)$, then the τ th segment refers to the observations $\mathcal{O}_\tau = \mathcal{O}_{d_1} + \dots + \mathcal{O}_{d_{\tau-1}} + \dots + \mathcal{O}_{d_\tau}$. Let U_τ denote the set of integers from $d_1 + \dots + d_{\tau-1} + 1$ to $d_1 + \dots + d_\tau$, inclusive. Then

$$P(\mathcal{O}_\tau | i_\tau) = \prod_{t \in U_\tau} P(\mathcal{O}_t | i_\tau) .$$

Thus

$$\frac{1}{P_\lambda(\mathcal{O})} \cdot \sum_p \sum_{\tau=1}^N P_\lambda(p \& \mathcal{O}) \ln \frac{P'(\mathcal{O}_\tau | i_\tau)}{P(\mathcal{O}_\tau | i_\tau)}$$

$$= \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\tau=1}^N P_\lambda(p \& \mathcal{O}) \sum_{t \in U_\tau} \ln \frac{P'(\mathcal{O}_t | i_\tau)}{P(\mathcal{O}_t | i_\tau)}$$

$$= \sum_{i,k} \ln \frac{P'(k|i)}{P(k|i)} \cdot \frac{1}{P_\lambda(\mathcal{O})} \sum_p \sum_{\substack{\tau \\ i_\tau = i \\ \exists \mathcal{O}_t = k}} \sum_{t \in U_\tau} P_\lambda(p \& \mathcal{O})$$

$$= \sum_{i,k} \ln \frac{P'(k|i)}{P(k|i)} \cdot c_\lambda(i,k), \quad \text{according to (L)} .$$

Lemma 3. The reestimates (\hat{a}_1) , (\hat{a}_{ij}) , $(\hat{P}(d|1))$, and $(\hat{P}(k|1))$ given in Section 3 respectively maximize the four summands comprising

$$\frac{Q(\lambda, \lambda') - Q(\lambda, \lambda)}{P}$$

as a function of λ' .

Proof. All four parts are similar. We prove part 4.

We are to maximize

$$\sum_{i,k} c(i,k) \ln \frac{P'(k|i)}{P(k|i)}$$

subject to

$$\sum_k P'(k|i) = 1, \quad \text{for } i = 1, 2, \dots, s.$$

Equivalently, maximize

$$\sum_{i,k} c(i,k) \ln P'(k|i) \quad \text{with} \quad \sum_k P'(k|i) = 1.$$

The maximum clearly occurs when the P' are proportional to the c , (e.g., by Lagrange multipliers). Thus

$$P'(k|i) = c(i,k) \left/ \sum_k c(i,k) \right. ,$$

as asserted.

Corollary. The likelihood increases if we replace our current parameters (a_1) , (a_{ij}) , $(P(d|i))$, $(P(k|i))$ by (\hat{a}_1) , (\hat{a}_{ij}) , $(\hat{P}(d|i))$, $(\hat{P}(k|i))$.

More is in fact true: any convex combination of old and new parameters also increases likelihood.

Lemma 4. Let

$$a'_1 = \theta a_1 + (1-\theta)\hat{a}_1, \quad 0 \leq \theta \leq 1,$$

$$a'_{ij} = \eta_i a_{ij} + (1-\eta_i)\hat{a}_{ij}, \quad 0 \leq \eta_i \leq 1,$$

$$P'(d|i) = \mu_i P(d|i) + (1-\mu_i)\hat{P}(d|i), \quad 0 \leq \mu_i \leq 1,$$

and

$$P'(k|i) = \nu_i P(k|i) + (1-\nu_i)\hat{P}(k|i), \quad 0 \leq \nu_i \leq 1.$$

Then $P_{\lambda'}(\sigma) \geq P_\lambda(\sigma)$ for any choice of θ , (η_i) , (μ_i) , (ν_i) .

Proof.

$$\frac{P_{\lambda'}(\sigma) - P_\lambda(\sigma)}{P_\lambda(\sigma)} \geq \frac{Q(\lambda, \lambda') - Q(\lambda, \lambda)}{P_\lambda(\sigma)}$$

$$= \sum_i c(i) \ln \frac{a'_1}{a_1} + \sum_{i,j} c(i,j) \ln \frac{a'_{ij}}{a_{ij}}$$

$$+ \sum_{i,d} c(i,d) \ln \frac{P'(d|i)}{P(d|i)} + \sum_{i,k} c(i,k) \ln \frac{P'(k|i)}{P(k|i)} .$$

The argument is the same for each sum.

For example,

$$\begin{aligned} & \sum_{i,j} c(i,j) \ln \frac{a'_{ij}}{a_{ij}} \\ &= \sum_i \sum_j c(i,j) \ln \left(\frac{\eta_i a_{ij} + (1-\eta_i) \hat{a}_{ij}}{a_{ij}} \right) \\ &\geq \sum_i \sum_j c(i,j) \left(\eta_i \ln \frac{a_{ij}}{a_{ij}} + (1-\eta_i) \ln \frac{\hat{a}_{ij}}{a_{ij}} \right) \end{aligned}$$

(by concavity of \ln),

$$= \sum_i (\eta_i) \sum_j c(i,j) \ln \left(\frac{\hat{a}_{ij}}{a_{ij}} \right) \geq 0$$

by maximality of the \hat{a}_{ij} .

5. Posterior estimation procedures

One standard method of finding the "best" state sequence in the usual case of hidden Markov models, is dynamic programming ("the Viterbi algorithm" [2]). We consider how that technique can be adapted to our case.

We wish to find the state sequence

$$p_{\max} = (i_1, d_1), \dots, (i_N, d_N)$$

which has maximum posterior probability. Equivalently,

we maximize $P(p|\sigma)$ over p . We assume a fixed model, and the path found will be the most likely, for this model, given the observations.

Let $\delta_t(i) = \max \text{Prob}(\sigma_1 \dots \sigma_t)$ over all partial paths p for which state i terminates at time t .
 $\delta_t^*(i) = \max \text{Prob}(\sigma_1 \dots \sigma_t)$ over all partial paths p for which state i starts at time $t+1$. Then, analogously to the a, a^* formulas,

$$\delta_t^*(j) = \max_i \delta_t(i)a_{ij},$$

$$\delta_t(i) \approx \max_d \delta_{t-d}(i)P(d|i)\text{P}(\sigma_{t-d+1} \dots \sigma_t | i).$$

As usual in dynamic programming, we keep a record of best predecessors, so that we can retrace our steps when we have finished the inductive calculation of $\delta_T(i)$, all i .

An alternative procedure gives (in the case where the model is exactly correct) the path with the maximum expected number of correct state-time assignments. This is also described in [2]. The analogous method here is simply to choose, at time t , that state i_t which maximizes $\gamma_t(i)$. The resulting state sequence (which

may have sudden discontinuities and impossible durations) gives the maximum expected number of correct state-time assignments.

Other procedures can be used to produce decisions with a variety of optimality properties. For example, if we define

$$\gamma_t(i,d) = \frac{\alpha_t^*(i)P(d|i)P(\mathcal{G}_{t+1} \dots \mathcal{G}_{t+d} | i) \beta_{t+d}(i)}{P}$$

then γ represents the posterior probability that state i occupied the time segment from $t+1$ to $t+d$. If we choose triplets t, i, d with the largest values of $\gamma_t(i, d)$ we expect the maximum number of correctly matched states with segments. Of course, there may be overlaps, and near-duplicates. If a complete path is desired, with maximum expected number of correctly matched state-segment pairs, this can be arranged using a dynamic program (variable duration) to find the path

$$\begin{aligned} p = (i_1, d_1), \dots, (i_N, d_N) &\ni \gamma_0(i_1, d_1) \\ &+ \gamma_{d_1}(i_2, d_2) + \gamma_{d_1+d_2}(i_3, d_3) + \dots \end{aligned}$$

is maximal.

This path, while complete, may have impossible state transitions $i_\tau \rightarrow i_{\tau+1}$ in it. If it is desired to eliminate these, then a dynamic program can again be employed to find the path with all transitions having nonzero probability, and maximal sum of gammas for segments. The ideas of this section will allow us to estimate, in the Poritz scheme [5] extended, which phone class is (probably) operative at each moment, and when transitions occur (probably). This amounts to probabilistic labeling and segmentation.

6. Special duration distributions

There are a few special duration distributions for which reestimation is particularly simple. For example, suppose that, for state i , the duration distribution is Poisson, with mean $\mu_i + 1$, starting at 1. Thus

$$P(\text{Duration} = 1 + d|i) = e^{-\mu_i} \cdot \frac{\mu_i^d}{d!}, \quad d = 0, 1, \dots$$

This distribution is potentially infinite, but for the problem at hand, we know $d \leq T$. Also, if we impose an arbitrary upper bound $d \leq D_{\max}$, there may be very little loss. In any case, we can show that the reestimate for the duration distribution (by maximizing the

Q -function) depends only on the expected average duration spent in state i :

$$\hat{\mu}_1 = \frac{\Sigma_t \gamma_t(1)}{S_T(1)} - 1 \quad (\text{see (F) \& (G)}) .$$

Lemma 5. If

$$P'(1+d|1) = e^{-\mu'_1} \frac{\mu'^d}{d!} ,$$

then $Q(\lambda, \lambda')$ is maximized, as a function of μ'_1 , at

$$\mu'_1 = \frac{\Sigma_t \gamma_t(1)}{S_T(1)} - 1 .$$

Proof.

$$\frac{Q(\lambda, \lambda') - Q(\lambda, \lambda)}{P(\mathcal{O})} = \sum_1 c(1) \ln \frac{a'_1}{a_1} + \dots$$

(according to Lemma 2). We are only interested in maximizing

$$\sum_d c(1, d) \ln \frac{P'(d|1)}{P(d|1)}$$

since the other terms either do not involve state 1 , or do not refer to duration. Thus we wish to maximize

Here $\Sigma_d d \cdot c(1, d)$ is the expected total duration spent in state $1 = \Sigma_t \gamma_t(1)$, and $\Sigma_d c(1, d)$ is the expected number of segments in which 1 was the state = $S_T(1) = E_{T+1}(1)$.

Of course, similar arguments apply whenever the maximum of the Q -function just depends on the average duration, for example:

$P(d|i)$ can be geometric:

$P(d|i) = \rho_1^{d-1} (1-\rho_1)$ for some ρ_1 ,
or binomial on the range 1 to $M_1 + 1$, (where M_1 is
fixed and known):

$$P(1+d|i) = \binom{M_1}{d} \rho_1^d (1-\rho_1)^{M_1-d}, \quad d = 0, \dots, M_1.$$

We observe, following [3], that whenever the maximum likelihood parameters for $P(d|i)$ can be determined, given the number of occurrences $n(i,d)$ for each d , then the maximum of the Q -function can be similarly found, simply by replacing the numbers $n(i,d)$ by the expected numbers $c(i,d)$. This maximization of Q need not, of course, be merely algebraic, as above, but might entail an iterative approximative scheme.

7. Extension to continuously distributed output

As described in [1], and extended in [3], it is not hard to replace $P(k|i)$ by a probability density function on continuous, possibly multivariate observations.

For application to the Poritz scheme, [5], or to the ideas proposed in the next paper, by Liporace [7], this extension of the variable duration idea to continuously distributed output is necessary.

We illustrate by allowing the output to be Gaussian distributed (in one-dimension for simplicity) with mean μ_1 , variance σ_1^2 , in state i . We again assume independence of the observations occurring while in state

- Our conclusion is that the reestimate which maximizes the Q -function is

$$\begin{aligned} \hat{\mu}_1 &= \frac{\sum_t \sigma_t \gamma_t^i(1)}{\sum_t \gamma_t^i(1)}, \\ \hat{\sigma}_1^2 &= \frac{\sum_t \sigma_t^2 \gamma_t^i(1)}{\sum_t \gamma_t^i(1)} - \hat{\mu}_1^2. \end{aligned}$$

For the proof, we need an analogue of Lemma 2 for the case of continuous output.

Lemma 2'.

$$\begin{aligned} \frac{Q(\lambda, \lambda') - Q(\lambda, \lambda)}{P_\lambda(\mathcal{O})} &= \sum_i c(i) \ln \frac{a'_i}{a_i} \\ &+ \sum_{i,j} c(i,j) \ln \frac{a'_{ij}}{a_{ij}} \\ &+ \sum_{i,d} c(i,d) \ln \frac{P'(d|i)}{P(d|i)} \\ &+ \sum_{i,t} \gamma_t^i(1) \ln \frac{P'(\sigma_t^i|i)}{P(\sigma_t^i|i)}. \end{aligned}$$

(This formula holds in both the continuous and discrete cases interpreting $P'(\sigma_t | i)$ and $P(\sigma_t | i)$ as likelihoods, but the Lemma 2 form is simpler in the discrete case.)

Proof. We can take up the proof of Lemma 2 in the consideration of the last sum:

$$H = \frac{1}{P_\lambda(\phi)} \cdot \sum_p \sum_{\tau=1}^N P_\lambda(p \& \phi) \sum_{t \in U_\tau} \ln \frac{P'(\sigma_t | i_\tau)}{P(\sigma_t | i_\tau)} .$$

We have:

$$H = \frac{1}{P_\lambda(\phi)} \sum_i \sum_t \sum_{\substack{p \ni \text{state} \\ i \text{ is used} \\ \text{at time } t}} P_\lambda(p \& \phi) \ln \frac{P'(\sigma_t | i)}{P(\sigma_t | i)}$$

$$= \frac{1}{P_\lambda(\phi)} \sum_{i,t} \gamma_t(i) \ln \frac{P'(\sigma_t | i)}{P(\sigma_t | i)}$$

Hence

$$H = \sum_{i,t} \gamma_t(i) \left[-\ln \sigma'_i - \frac{1}{2} \frac{(\sigma_t - \mu'_i)^2}{\sigma'^2_i} + \ln \sigma_i \right]$$

$$\frac{\partial H}{\partial \mu'_i} = \sum_t \gamma_t(i) \frac{(\sigma_t - \mu'_i)}{\sigma'^2_i} ,$$

and

$$\frac{\partial H}{\partial \sigma'_i} = - \sum_t \frac{\gamma_t(i)}{\sigma'_i} + \sum_t \frac{(\sigma_t - \mu'_i)^2}{\sigma'^3_i} \gamma_t(i) .$$

The maximum occurs at

$$\hat{\mu}_i = \sum_t \gamma_t(i) \sigma_t / \sum_t \gamma_t(i) ,$$

as asserted.

We now employ Lemma 2' for the Gaussian example at hand. Here

$$P'(\sigma_t | i) = \frac{1}{\sqrt{2\pi} \sigma'_i} e^{-\frac{1}{2} \frac{(\sigma_t - \mu'_i)^2}{\sigma'^2_i}} .$$

$$\hat{\sigma}_1^2 = \sum_t (\sigma_t - \hat{\mu}_1)^2 \cdot r_t(1) / \sum_t r_t(1)$$

$$= \sum_t \sigma_t^2 \cdot r_t(1) / \sum_t r_t(1) - \hat{\mu}_1^2,$$

as asserted.

This analytic approach to maximizing Q is applicable to a variety of simple continuous distributions, including the multi-variate Gaussian, but, as shown in [3], all that is really necessary is the ability to find maximum likelihood estimators, given actual observations, analytically, or otherwise.

8. Duration 0

In the preceding analyses, durations were required to be at least one. In some special cases it is useful to allow a state to have duration 0. In these cases, the inductive formulas for $\alpha, \alpha^*, \beta, \beta^*$ go awry, and special methods must be used. For example, if we have a hidden Markov model for a specific word or sentence, then a definite state order is assumed. But it is quite possible for a speaker to omit, entirely, one or more of the phones in a sequence, especially when excited. For this reason we may wish to allow a state or states to

have 0 duration. In this particular example, it is possible to accommodate 0 duration in a simple way.

Since the states are ordered, there is no way to return to an earlier state in the assumed sequence, thus obtaining an infinite loop of 0 duration states. The argument below describes how to handle 0 duration even in the general case, where infinite loops can occur.

Suppose we have a variable duration hidden Markov model with parameters $(\alpha_1), (\alpha_{1j}), (P(d|1), d = 0, 1, \dots, D), (b_j(k))$. Suppose we also have a sequence of observations $\mathcal{O} = \mathcal{O}_1 \mathcal{O}_2 \dots \mathcal{O}_T$. We show how to compute the probability of \mathcal{O} , given the model. Let $\bar{P}(d|1) = P(d|1)/(1-P(0|1))$ for each i , and $d > 0$. Thus $\bar{P}(d|i) = \text{Prob}(\text{duration is } d, \text{ given that it is not 0})$. Similarly, let $\bar{a}_{ij} = \text{Prob}(j \text{ is the next state after } i \text{ for which the duration is not 0})$. Then

$$\bar{a}_{1j} = a_{1j} \cdot (1 - P(0|j))$$

$$+ \sum_k a_{1k} P(0|k) \cdot a_{kj} (1 - P(0|j))$$

$$+ \sum_{k,l} a_{1k} P(0|k) \cdot a_{kj} P(0|l) \cdot a_{lj} (1 - P(0|j)) + \dots$$

This infinite sum can either be approximated, or

computed exactly by summing an infinite series of powers of the matrix $M = (a_{k,l}P(o|l))$, obtaining $(I-M)^{-1}$, etc.

Lastly, we define $\bar{a}_1 = \text{Prob(first state with nonzero duration is } i)$

$$\begin{aligned} &= a_1(1-P(o|1)) + \sum_j a_j P(o|j) a_{ji}(1-P(o|1)) \\ &+ \sum_{k,j} a_k P(o|k) \cdot a_{kj} P(o|j) \cdot a_{ji}(1-P(o|1)) \\ &+ \dots, \end{aligned} \quad [3]$$

which is handled similarly.

Then the hidden Markov model with parameters $(\bar{a}_1), (\bar{a}_{1j}), (\bar{P}(d|1)), d = 1, 2, \dots, D$, and $(b_j(k))$ can be treated in the standard manner, and allow the computation of the probability of \mathcal{O} .

Alternative approaches might involve structuring the Markov chain so that no infinite sequences of 0-duration states can occur, or approximating the model by truncating all such chains to a small finite length.

REFERENCES

- [1] Baum, Leonard, E., Petrie, Ted, Soules, George, and Weiss, Norman, A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains, *Annals of Mathematical Statistics*, Vol. 41, No. 1, pp. 164-171, 1970.
- [2] Forney, G. David, Jr., The Viterbi Algorithm, *Proceedings of the IEEE*, Vol. 61, No. 3, pp. 268-278, March 1973.
- [3] Dempster, A. P., Laird, N. M., Rubin, D. B., Maximum likelihood from incomplete data via the EM algorithm, *Journal of the Royal Statistical Society, Series B (Methodological)*, Vol. 39, pp. 1-38, 1977.
- [4] Cave, Robert L., Neuwirth, Lee P., Hidden Markov models for English, *these Proceedings*.
- [5] Poritz, Alan B., Linear predictive hidden Markov models, *these Proceedings*.
- [6] Liporace, Louis A., Maximum likelihood analysis for multivariate observations of Markov sources, *these Proceedings*.
- [7] Liporace, Louis A., Statistical analysis of time series, *these Proceedings*.