

Statistical Bioinformatics, Makerere Hidden Markov Models Timo Koski

TK

02.08.2018



This lecture corresponds some of the sections of chapter 12 in Ewens and Grant

- 1) Definition and examples
- 2) Properties implied by conditional independence
- 3) Forward-Backward Algorithm and the Scoring/Evaluation Problem
- 4) Alignment (Viterbi algorithm), Learning (Baum - Welch algorithm)

A sequence of random variables $\{X_n\}_{n=0}^{\infty}$ is called a **Markov chain, (MC)**, if for all $n \geq 1$ and $j_0, j_1, \dots, j_n \in \mathcal{S}$,

$$P(X_n = j_n | X_0 = j_0, X_1 = j_1, \dots, X_{n-1} = j_{n-1}) =$$

$$P(X_n = j_n | X_{n-1} = j_{n-1}).$$

The Markov property

If $X_n = j_n$ is a future event, then the conditional probability of this event given the past history $X_0 = j_0, X_1 = j_1, \dots, X_{n-1} = j_{n-1}$ depends only upon the immediate past $X_{n-1} = j_{n-1}$ and not upon the remote past $X_0 = j_0, X_1 = j_1, \dots, X_{n-2} = j_{n-2}$.

Hidden Markov Models (HMM)

HMM is a model family for a sequence of symbols from an alphabet $\mathcal{O} = \{o_1, o_2, \dots, o_K\}$. The model uses the idea of a hidden sequence of state transitions.

HMM has a definition with parts I–III.



- (I) **Hidden Markov Chain** $\{X_n\}_{n=0}^{\infty}$ is a Markov chain assuming values in a finite state space $S = \{1, 2, \dots, J\}$ with J states. The time-homogeneous conditional probabilities are

$$a_{i|j} = P(X_n = j | X_{n-1} = i), n \geq 1, i, j \in S$$

and the transition probability matrix is

$$A = \left(a_{i|j} \right)_{i=1, j=1}^{J, J}$$

Hidden Markov Models (HMM) I

A matrix

$$A = (a_{i|j})_{i=1, j=1}^{J, J}$$

with the constraints

$$a_{i|j} \geq 0, \sum_{j=1}^J a_{i|j} = 1.$$

is called a stochastic matrix.



Hidden Markov Models (HMM) I

At time $n = 0$ the state X_0 is specified by the initial probability distribution $\pi_j(0) = P(X_0 = j)$ with

$$\pi(0) = (\pi_1(0), \dots, \pi_J(0)).$$

$$\pi_j(n) = P(X_n = j)$$

$$\pi(n) = (\pi_1(n), \dots, \pi_J(n)).$$



- (II) **Observable Random Process** A random process $\{Y_n\}_{n=0}^{\infty}$ with a finite state space $\mathcal{O} = \{o_1, o_2, \dots, o_K\}$, where K can be $\neq J$. The processes $\{Y_n\}_{n=0}^{\infty}$ and $\{X_n\}_{n=0}^{\infty}$ are for any fixed n related by the conditional probability distributions

$$b_j(k) = P(Y_n = o_k | X_n = j).$$

We set

$$B = \{b_j(k)\}_{j=1, k=1}^{J, K}$$

and call this the *emission probability matrix*. This is another stochastic matrix in the sense that

$$b_j(k) \geq 0, \sum_{k=1}^K b_j(k) = 1.$$

- (III) **Conditional independence** For any sequence of states $j_0 j_1 \dots j_n$ the probability of the sequence $o_0 o_1 \dots o_n$ is

$$P(Y_0 = o_0, \dots, Y_n = o_n \mid X_0 = j_0, \dots, X_n = j_n, B) = \prod_{l=0}^n b_{j_l}(l).$$

An HMM is designated by

$$\lambda = (A, B, \pi(0)).$$

UNDER THE HMM ASSUMPTIONS THE STRING $\mathbf{o} = o_0 \dots o_n$ HAS THE PROBABILITY

$$P(\mathbf{o}) = P(Y_0 = o_0, \dots, Y_n = o_n; \lambda) = \sum_{j_0=1}^J \dots \sum_{j_n=1}^J P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j_n; \lambda)$$

$$\sum_{j_0=1}^J \dots \sum_{j_n=1}^J P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j_n; \lambda)$$

where

$$P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j_n; \lambda) = \pi_{j_0}(0) \cdot \prod_{l=0}^n b_{j_l}(l) \prod_{l=1}^n a_{j_{l-1}|j_l}.$$

- [1] **The Evaluation or Scoring Problem** Compute $P(Y_0 = o_0, \dots, Y_n = o_n; \lambda)$. Since the marginalization involves J^{n+1} possible sequences, the total computational requirements are of the order $2(n+1) \cdot J^{n+1}$ operations. The solution is known as the *forward-backward procedure*

[2]] **The Decoding or Alignment Problem** Find the most probable state sequence that led to the observed sequence $(o_0 \dots o_n)$. This is an *alignment problem*. Find the sequence $j_0^* \dots j_n^*$ that maximizes

$$P(X_0 = j_0, \dots, X_n = j_n, Y_0 = o_0, \dots, Y_n = o_n; \lambda)$$

for a fixed observed sequence $o_0 \dots o_n$ (*Viterbi algorithm*).

- [3] **The Learning or Training Problem** Given an observed sequence $\mathbf{o} = o_0 \dots o_n$, find the 'right' model parameter values

$$\lambda = (A, B, \pi(0))$$

in a fixed topology that specify a model most likely to generate the given sequence

HIDDEN MARKOV MODELS: Conditional Independence

Markov property and the conditional independence property III imply useful expressions for smoothing, prediction, filtering and evaluation, and these yield the solutions to the three problems stated above.



The smoothing posterior probability is defined as

$$\hat{\pi}_j(n|m) = P(X_n = j | Y_0 = o_0, \dots, Y_m = o_m)$$

for a standard HMM.

For $n = 0, \dots, N - 1$ it holds that

$$\hat{\pi}_j(n|N) = \hat{\pi}_i(n|n) \cdot \sum_{k=1}^J \frac{a_{j|k}}{\hat{\pi}_k(n+1|n)} \hat{\pi}_k(n+1|N).$$

We intend to explain this in some detail.

Here the typesetting is simplified e.g. by writing a conditional probability as

$$P(Y_m = o_m, \dots, Y_N = o_N \mid X_n = j_n, \dots, X_N = j_N)$$

simply as

$$P(Y_m, \dots, Y_N \mid X_n, \dots, X_N)$$

$$P(Y_0, \dots, Y_n \mid X_0, \dots, X_n) = \prod_{i=0}^n P(Y_i \mid X_i).$$

Proposition 1

Proposition

For all integers n and m such that $0 \leq n \leq m \leq N$

$$P(Y_m, \dots, Y_N \mid X_n, \dots, X_N) = P(Y_m, \dots, Y_N \mid X_m, \dots, X_N).$$



Proposition 1.

Proof: The left hand side of the asserted identity can be expressed as

$$\frac{1}{P(X_0, \dots, X_N)} \sum P(Y_0, \dots, Y_N | X_0, \dots, X_N) \cdot P(X_0, \dots, X_N),$$

where the summation is over j_0, \dots, j_{n-1} (i.e. the values of $X_{j_0}, \dots, X_{j_{n-1}}$). If $n = 0$, there is no summation. By conditional independence (and a marginalization argument)

$$P(Y_0, \dots, Y_N | X_0, \dots, X_N) = P(Y_0 | X_0) \cdot \dots \cdot P(Y_N | X_N).$$

This can be taken outside the summation sign \sum , since $m \geq n$.



Then we are dealing with

$$\prod_{l=m}^N P(Y_l|X_l) \frac{1}{P(X_n, \dots, X_N)} \sum P(X_0, \dots, X_N),$$

where the sum equals $P(X_n, \dots, X_N)$, since we are summing over j_0, \dots, j_{n-1} . Thus the whole last expression equals

$$= \prod_{l=m}^N P(Y_l|X_l),$$

which is independent of n . Since the right hand side of the above is a special case of the left hand side for $n = m$, this proves the assertion as claimed. □

Proposition 2.

Proposition

For all integers $n = 0, \dots, N - 1$

$$P(Y_{n+1}, \dots, Y_N \mid X_0, \dots, X_n) = \\ P(Y_{n+1}, \dots, Y_N \mid X_n).$$

Proof of proposition 2.

Set

$$X^{(t)} = (X_0, \dots, X_t), Y^{(t)} = (Y_0, \dots, Y_t),$$

Proof: The left hand side is

$$\frac{1}{P(X^{(n)})} \sum P(X^{(N)}) \cdot P(Y_{n+1}, \dots, Y_N | X^{(N)}),$$

where the summation is over j_{n+1}, \dots, j_N . By the first proposition, (with $m = n + 1$, $n = 0$), we have

$$P(Y_{n+1}, \dots, Y_N | X^{(N)}) = P(Y_{n+1}, \dots, Y_N | X_{n+1}, \dots, X_N)$$

and using the same proposition and equation once more (with $m = n + 1$) we have

$$\begin{aligned} P(Y_{n+1}, \dots, Y_N | X_{n+1}, \dots, X_N) = \\ P(Y_{n+1}, \dots, Y_N | X_n, \dots, X_N). \end{aligned}$$

Proof of proposition 2.

Thus

$$\begin{aligned} \sum P(X^{(N)}) \cdot P(Y_{n+1}, \dots, Y_N | X^N) &= \\ \sum P(X^N) \cdot P(Y_{n+1}, \dots, Y_N | X_n, \dots, X_N). \end{aligned}$$

By conditional probability $P(X^{(N)}) =$

$$P(X_{n+1}, \dots, X_N | X^{(n)}) \cdot P(X^{(n)})$$

and by a consequence of *Markov property* we have

$$P(X_{n+1}, \dots, X_N | X^{(n)}) = P(X_{n+1}, \dots, X_N | X_n).$$

Proof of proposition 2.

Thus the sum equals, since we are summing over j_{n+1}, \dots, j_N ,

$$\begin{aligned} & \sum P(X^{(N)}) \cdot P(Y_{n+1}, \dots, Y_N | X_n, \dots, X_N) = \\ & = \sum \frac{P(Y_{n+1}, \dots, Y_N, X_n, \dots, X_N) \cdot P(X^{(n)})}{P(X_n)}. \end{aligned}$$

Proof of proposition 2.

$$= P(X^{(n)}) \sum \frac{P(Y_{n+1}, \dots, Y_N, X_n, \dots, X_N)}{P(X_n)}.$$

And as we are summing over j_{n+1}, \dots, j_N , we have here that

$$\sum \frac{P(Y_{n+1}, \dots, Y_N, X_n, \dots, X_N)}{P(X_n)} = \frac{P(Y_{n+1}, \dots, Y_N, X_n)}{P(X_n)} =$$

Proof of proposition 2.

$$= P(Y_{n+1}, \dots, Y_N | X_n).$$

We have that

$$\begin{aligned} \frac{1}{P(X^{(n)})} \sum P(X^{(N)}) \cdot P(Y_{n+1}, \dots, Y_N | X^{(N)}) = \\ \frac{1}{P(X^{(n)})} P(Y_{n+1}, \dots, Y_N | X_n) \cdot P(X^{(n)}), \end{aligned}$$

which proves the assertion as claimed. □

Proposition 3.

Proposition

For all integers $n = 0, \dots, N$

$$P(Y_0, \dots, Y_n | X_0, \dots, X_N) =$$

$$P(Y_0, \dots, Y_n | X_0, \dots, X_n).$$



Proposition 4

Proposition

For all integers $n = 0, \dots, N$

$$P(Y_0, Y_1, \dots, Y_N | X_n) = \\ P(Y_0, Y_1, \dots, Y_n | X_n) \cdot P(Y_{n+1}, \dots, Y_N | X_n).$$

The conditional probability $P(Y_{n+1}, \dots, Y_N | X_n)$ is called the **backward variable**. The next proposition is used to find a recursion for this backward variable.

A Proposition for the Backward variable

Proposition

For all integers $n = 0, \dots, N$

$$P(Y_n, Y_{n+1}, \dots, Y_N | X_n) = \\ P(Y_n | X_n) \cdot P(Y_{n+1}, \dots, Y_N | X_n).$$

Proposition

For all integers $n = 0, \dots, N - 1$

$$P(Y_0, Y_1, \dots, Y_N \mid X_n, X_{n+1}) = \\ P(Y_0, Y_1, \dots, Y_n \mid X_n) \cdot P(Y_{n+1}, \dots, Y_N \mid X_{n+1}).$$

Smoothing probability

$$\hat{\pi}_j(n|N) = P(X_n = j | Y^{(N)}).$$

By another marginalization we get

$$\begin{aligned}\hat{\pi}_j(n|N) &= \sum_{k=1}^J P(X_n = j, X_{n+1} = k | Y^{(N)}) = \\ &= \sum_{k=1}^J P(X_n = j | X_{n+1} = k, Y^{(N)}) P(X_{n+1} = k | Y^{(N)}) =\end{aligned}$$

$$\begin{aligned} &= \sum_{k=1}^J P\left(X_n = j | X_{n+1} = k, Y^{(N)}\right) \hat{\pi}_k(n+1|N) \\ &= \sum_{k=1}^J \frac{P\left(X_n = j, X_{n+1} = k, Y^{(N)}\right)}{P\left(X_{n+1} = k, Y^{(N)}\right)} \hat{\pi}_k(n+1|N). \end{aligned}$$

In the numerator inside the summation above we have

$$\begin{aligned} P\left(X_n = j, X_{n+1} = k, Y^{(N)}\right) &= \\ P\left(Y^{(N)} | X_n = j, X_{n+1} = k\right) P\left(X_n = j, X_{n+1} = k\right) &= \\ = P\left(Y^{(n)} | X_n = j\right) P\left(Y_{n+1}, \dots, Y_N | X_{n+1} = k\right) \cdot P\left(X_n = j\right) \cdot a_{j|k} \end{aligned}$$

using the factorization in proposition 4 and the definition of $a_{j|k}$.

Then, since

$$P\left(Y^{(n)}|X_n = j\right) P\left(X_n = j\right) = \\ P\left(X_n = j|Y^{(n)}\right) P\left(Y^{(n)}\right),$$

we have obtained

$$P\left(X_n = j, X_{n+1} = k, Y^{(N)}\right) = \\ P\left(X_n = j|Y^{(n)}\right) P\left(Y^{(n)}\right) P\left(Y_{n+1}, \dots, Y_N|X_{n+1} = k\right) \cdot a_{j|k}.$$

Smoothing probability

For the denominator inside the summation for smoothing probability it holds that

$$\begin{aligned}P\left(X_{n+1} = k, Y^{(N)}\right) &= P\left(Y^{(N)}|X_{n+1} = k\right) P\left(X_{n+1} = k\right) = \\&= P\left(Y^{(n)}|X_{n+1} = k\right) P\left(Y_{n+1}, \dots, Y_N|X_{n+1} = k\right) P\left(X_{n+1} = k\right)\end{aligned}$$

Next by $P\left(Y^{(n)}|X_{n+1} = k\right) P\left(X_{n+1} = k\right) = P\left(X_{n+1} = k|Y^{(n)}\right) P\left(Y^{(n)}\right)$ we get

$$\begin{aligned}P\left(X_{n+1} = k, Y^{(N)}\right) &= \\&= P\left(X_{n+1} = k|Y^{(n)}\right) P\left(Y^{(n)}\right) P\left(Y_{n+1}, \dots, Y_N|X_{n+1} = k\right).\end{aligned}$$

Therefore

$$\begin{aligned}\hat{\pi}_j(n|N) &= \sum_{k=1}^J \frac{P(X_n = j|Y^{(n)}) \cdot a_{j|k}}{P(X_{n+1} = k|Y^{(n)})} \hat{\pi}_k(n+1|N) = \\ &= \hat{\pi}_i(n|n) \cdot \sum_{k=1}^J \frac{a_{j|k}}{\hat{\pi}_k(n+1|n)} \hat{\pi}_k(n+1|N),\end{aligned}$$

which is the result as claimed. □

(a) Filtering posterior

$$\hat{\pi}_j(n|n) = \frac{\left[\sum_{i=1}^J \hat{\pi}_i(n-1|n-1) \cdot a_{ij} \right] \cdot b_j(o_n)}{\sum_{j=1}^J \sum_{i=1}^J \hat{\pi}_i(n-1|n-1) \cdot a_{ij} \cdot b_j(o_n)},$$

for $j = 1, \dots, J$.

(b) Prediction posterior probability for an HMM:

$$\hat{\pi}_j(n|n-1) = \sum_{i=1}^J \hat{\pi}_i(n-1|n-1) \cdot a_{ij}, j = 1, \dots, J,$$

Thus

$$\hat{\pi}_j(n|n) = \frac{\hat{\pi}_j(n|n-1) \cdot b_j(o_n)}{\sum_{j=1}^J \hat{\pi}_j(n|n-1) \cdot b_j(o_n)}, j = 1, \dots, J.$$

The log likelihood function for the sequence $\mathbf{o} = o_{j_0} \dots o_{j_n}$ with respect to the HMM model family is

$$\log P \left(Y_0 = o_{j_0}, \dots, Y_n = o_{j_n} \right) =$$

$$\sum_{i=0}^n \log f \left(Y_i = o_{j_i} \mid o_{j_0}, \dots, o_{j_{i-1}} \right),$$

where

$$f \left(Y_i = o_{j_i} \mid o_{j_0}, \dots, o_{j_{i-1}} \right) = \sum_{l=1}^J \hat{\pi}_l (i \mid i-1) b_l \left(o_{j_i} \right).$$

The problem is to compute the simultaneous probability for the a sequence of emitted symbols, $\mathbf{o} = o_0 \dots o_N$, conditioned on some model $\lambda = (A, B, \pi(0))$,

$$L_N = P(Y_0 = o_0 \dots, Y_n = o_N; \lambda) = \sum_{j_0=1}^J \dots \sum_{j_N=1}^J \pi_{j_0}(0) b_{j_0}(0) \prod_{l=1}^N a_{j_{l-1}|j_l} b_{j_l}(l),$$

so that the exponential growth of operations in N involved in the marginalization is avoided. In order to simplify the notation, the reference to the model λ is omitted.

Let

$$P(Y_0 = o_0, \dots, Y_N = o_N, X_n = j) = \\ P(X_n = j) \cdot P(Y_0 = o_0, \dots, Y_N = o_N | X_n = j).$$

But the right hand side is factorized as

$$P(Y_0 = o_0, \dots, Y_N = o_N | X_n = j) = \\ P(Y_0 = o_0, \dots, Y_n = o_n | X_n = j) \cdot P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j)$$

Forward Algorithm

This gives

$$P(Y_0 = o_0, \dots, Y_n = o_n, X_n = j) \cdot P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j).$$

Since

$$\begin{aligned} &P(Y_0 = o_0, \dots, Y_N = o_N) \\ &= \sum_{j=1}^J P(Y_0 = o_0, \dots, Y_N = o_N, X_n = j), \end{aligned}$$

we get

$$P(Y_0 = o_0, \dots, Y_N = o_N) = \sum_{j=1}^J \alpha_n(j) \cdot \beta_n(j),$$



$$P(Y_0 = o_0, \dots, Y_N = o_N) = \sum_{j=1}^J \alpha_n(j) \cdot \beta_n(j),$$

where

$$\alpha_n(j) = P(Y_0 = o_0, \dots, Y_n = o_n, X_n = j).$$

$$\beta_n(j) = P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j).$$

We take $\beta_N(j) = 1$ for every j arbitrarily.

First

$$\begin{aligned} \alpha_{n+1}(j) &= P(Y_0 = o_0, \dots, Y_{n+1} = o_{n+1}, X_{n+1} = j) = \\ &= \sum_{i=1}^J P(Y_0 = o_0, \dots, Y_{n+1} = o_{n+1}, X_n = i, X_{n+1} = j) \\ &= \sum_{i=1}^J P(X_n = i, X_{n+1} = j) \cdot P(Y_0 = o_0, \dots, Y_{n+1} = o_{n+1} | X_n = i, X_{n+1} = j). \end{aligned}$$

Forward Algorithm

Here, by the properties derived ('One More Proposition')

$$\begin{aligned} &= \sum_{i=1}^J P(X_n = i, X_{n+1} = j) \cdot P(Y_0 = o_0, \dots, Y_n = o_n | X_n = i) \cdot P(Y_{n+1} = o_{n+1} | X_{n+1} = j). \\ &P(X_n = i, X_{n+1} = j) \cdot P(Y_{n+1} = o_{n+1} | X_{n+1} = j) = a_{ij} \cdot b_j(o_{n+1}) \cdot P(X_n = i). \end{aligned}$$

Hence we have

$$\begin{aligned} &\sum_{i=1}^J P(X_n = i, X_{n+1} = j) \cdot P(Y_0 = o_0, \dots, Y_n = o_n | X_n = i) \cdot P(Y_{n+1} = o_{n+1} | X_{n+1} = j) = \\ &= \sum_{i=1}^J P(Y_0 = o_0, \dots, Y_n = o_n, X_n = i) \cdot a_{ij} \cdot b_j(o_{n+1}). \end{aligned}$$



This last expression is by our definition of the forward variable equal to

$$\begin{aligned} &= \sum_{i=1}^J \alpha_n(i) \cdot a_{i|j} \cdot b_j(o_{n+1}) = \\ &\quad \left[\sum_{i=1}^J \alpha_n(i) \cdot a_{i|j} \right] \cdot b_j(o_{n+1}). \end{aligned}$$

This completes the derivation of the forward algorithm. We summarize the result in a formal way.

The Forward Recursion

Consider the forward variable $\alpha_n(j)$ defined as

$$\alpha_n(j) = P(Y_0 = o_0, \dots, Y_n = o_n, X_n = j),$$

which is the probability of the emitted subsequence $\mathbf{o} = o_0 \dots o_n$ and of the hidden chain being in the state j at time n (given the model λ).

Start:

$$\alpha_0(j) = b_j(o_0) \pi_j(0), \quad j = 1, \dots, J.$$

Recursion:

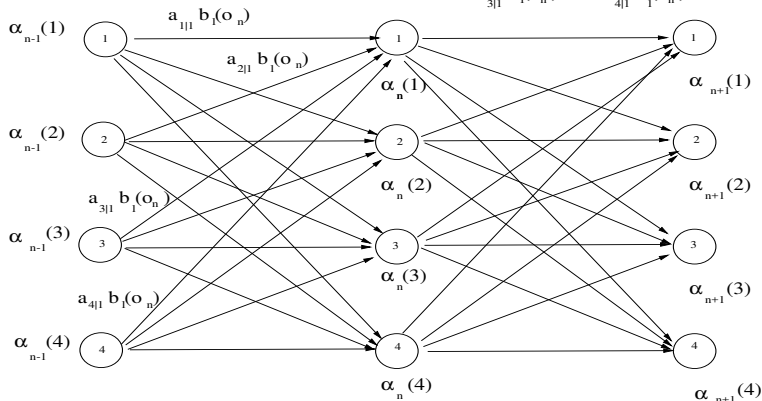
$$\alpha_{n+1}(j) = \left[\sum_{i=1}^J \alpha_n(i) \cdot a_{ij} \right] \cdot b_j(o_{n+1}).$$

$$j = 1, \dots, J, \quad 1 \leq n \leq N-1.$$



The Forward Recursion Trellis

$$\alpha_n(1) = a_{11} b_1(o_n) + a_{21} b_1(o_n) + a_{31} b_1(o_n) + a_{41} b_1(o_n)$$



Backward Algorithm

By definition

$$\begin{aligned}\beta_n(j) &= \sum_{i=1}^J \frac{P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N, X_n = j, X_{n+1} = i)}{P(X_n = j)} \\ &= \sum_{i=1}^J \frac{P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j, X_{n+1} = i) P(X_n = j, X_{n+1} = i)}{P(X_n = j)}.\end{aligned}$$

Here $\frac{P(X_n=j, X_{n+1}=i)}{P(X_n=j)} = a_{j|i}$ and

$$\begin{aligned}P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j, X_{n+1} = i) &= \\ &= P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_{n+1} = i).\end{aligned}$$

Backward Recursion

We apply furthermore one of the previous properties

$$\begin{aligned} P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_{n+1} = i) &= \\ &= P(Y_{n+1} = o_{n+1} | X_{n+1} = i) \cdot P(Y_{n+2} = o_{n+2}, \dots, Y_N = o_N | X_{n+1} = i). \end{aligned}$$



Backward Recursion

Hence it follows that

$$\beta_n(j) = \sum_{i=1}^J P(Y_{n+1} = o_{n+1} | X_{n+1} = i) \cdot P(Y_{n+2} = o_{n+2}, \dots, Y_N = o_N | X_{n+1} = i) \cdot a_{ji}$$

Recalling the definition of the backward variable and the emission probability $b_i(o_{n+1}) = P(Y_{n+1} = o_{n+1} | X_{n+1} = i)$ we have

$$\beta_n(j) = \sum_{i=1}^J b_i(o_{n+1}) \cdot \beta_{n+1}(i) \cdot a_{ji}$$

The Backward Procedure

Consider the backward variable $\beta_n(j)$ defined as

$$\beta_n(j) = P(Y_{n+1} = o_{n+1}, \dots, Y_N = o_N | X_n = j),$$

which is the probability of the emitted subsequence $o_{n+1} \dots o_N$ (from $n+1$ till the end) conditioned on the hidden chain being in the state j at time n (conditional on the model λ).

Start:

$$\beta_N(j) = 1 \quad j = 1, \dots, J.$$

Recursion :

$$\beta_n(j) = \sum_{i=1}^J b_i(o_{n+1}) \cdot \beta_{n+1}(i) \cdot a_{ji}. \quad j = 1, \dots, J, \quad n = N-1, N-2, \dots, 0.$$



The Scoring (Evaluation) Problem

$$L_N = P(Y_0 = o_0, \dots, Y_N = o_N) = \sum_{j=1}^J \alpha_n(j) \cdot \beta_n(j).$$

Hence we have for any $n = 0, \dots, N$ a respective way of computing L_N .
For example with $n = N$ we have

$$L_N = \sum_{j=1}^J \alpha_N(j)$$

by the convention $\beta_N(j) = 1$.



Filtering, Smoothing and Prediction

All probabilistic information about X_n given a sequence of observations $o_0 \dots o_N$ is contained in the conditional probabilities

$$\hat{\pi}_j(n|N) = P(X_n = j | Y_0 = o_0, \dots, Y_N = o_N).$$

This is conditioned on the model λ .

For $n > N$ the probability $\hat{\pi}_j(n|N)$ deals with **prediction**, $n = N$ the probability $\hat{\pi}_j(n|N)$ is a **filtering** probability. This is the standard phrase for reconstruction of a hidden variable from observations. For $n < N$ we talk about a **smoothing** probability.



For $n < N$

$$\hat{\pi}_{j|k}(n|N) = P(X_n = j, X_{n+1} = k | Y_0 = o_0, \dots, Y_N = o_N),$$

is the conditional posterior probability that a transition has taken place between any two states j and k at time $n + 1$.

Let us first find $\hat{\pi}_j(n|n)$. We use the definition of conditional probability to write

$$\begin{aligned}\hat{\pi}_j(n|n) &= P(X_n = j | Y_0 = o_0, \dots, Y_n = o_n) = \\ &= \frac{P(Y_0 = o_0, \dots, Y_m = o_m, X_n = j)}{P(Y_0 = o_0, \dots, Y_n = o_n)}.\end{aligned}$$

Using the definition and the evaluation formula we have

$$= \frac{\alpha_N(j)}{\sum_{j=1}^J \alpha_N(j)}.$$

Let

$$\hat{\pi}_j(n|N) = P(X_n = j | Y_0 = o_0, \dots, Y_n = o_n) =$$

and simplify the notation

$$\begin{aligned} &= \frac{P(Y_0, \dots, Y_N, X_n)}{P(Y_0, \dots, Y_N)} = \\ &= \frac{P(Y_0, \dots, Y_n, Y_{n+1}, \dots, Y_N, X_n)}{P(Y_0, \dots, Y_n) P(Y_{n+1}, \dots, Y_N | Y_0, \dots, Y_n)} = \end{aligned}$$

$$= \frac{P(Y_0, Y_1, \dots, Y_n, X_n) \cdot P(Y_{n+1}, \dots, Y_N | X_n)}{P(Y_0, \dots, Y_n) P(Y_{n+1}, \dots, Y_N | Y_0, \dots, Y_n)},$$

by one of the factorizations derived earlier and the definition of conditional probability.

Scalings in Filtering, Smoothing and Prediction

Here Devijver (1985) (no reference included) introduces

$$\hat{\pi}_j(n|N) = \bar{\alpha}_n(j) \cdot \bar{\beta}_n(j)$$

with

$$\bar{\alpha}_n(j) = \frac{P(Y_0, \dots, Y_n, X_n = j)}{P(Y_0, \dots, Y_n)}$$

and

$$\bar{\beta}_n(j) = \frac{P(Y_{n+1}, \dots, Y_N | X_n = j)}{P(Y_{n+1}, \dots, Y_N | Y_0, \dots, Y_n)},$$

These are evidently scalings of the forward and backward variables. The important property of these particular scalings is that there exist recursions for the scaled forward and backward variables $\bar{\alpha}_n(j)$ and $\bar{\beta}_n(j)$.



First, $\tilde{\alpha}_n(j)$ equals the filtering posterior probability $\tilde{\alpha}_n(j) = \hat{\pi}_j(n|n)$. By the preceding

$$\tilde{\alpha}_n(j) = \frac{\left[\sum_{i=1}^J \tilde{\alpha}_{n-1}(i) \cdot a_{i|j} \right] \cdot b_j(o_n)}{\sum_{j=1}^J \sum_{i=1}^J \tilde{\alpha}_{n-1}(i) \cdot a_{i|j} \cdot b_j(o_n)},$$

which is the desired forward recursion.

Next we get

$$\tilde{\beta}_n(j) = \sum_{k=1}^J \frac{P(Y_{n+1}, Y_{n+2}, \dots, Y_N, X_{n+1} = k \mid X_n = j)}{P(Y_{n+1}, Y_{n+2}, \dots, Y_N \mid Y_0, \dots, Y_n)}.$$

By some rearranging we have

$$\tilde{\beta}_n(j) = \sum_{k=1}^J \frac{P\left(Y_{(n+1)}^{(N)} \mid X_{n+1} = k, X_n = j\right) P(X_{n+1} = k \mid X_n = j)}{P(Y_{n+1}, Y_{n+2}, \dots, Y_N \mid Y_0, \dots, Y_n)}.$$

From the preceding

$$P(Y_{n+1}, Y_{n+2}, \dots, Y_N \mid X_{n+1} = k, X_n = j) =$$

$$P(Y_{n+1}, Y_{n+2}, \dots, Y_N \mid X_{n+1} = k)$$

and

Scaling Recursions

$$P(Y_{n+1}, Y_{n+2}, \dots, Y_N | X_{n+1} = k) = \\ P(Y_{n+1} | X_{n+1} = k) \cdot P(Y_{n+2}, \dots, Y_N | X_{n+1}).$$

Since $P(Y_{n+1} | X_{n+1} = k) = b_k(o_{jk})$ and $P(X_{n+1} = k | X_n = j) = a_{j|k}$ we have

$$P(Y_{n+1}, Y_{n+2}, \dots, Y_N | X_{n+1} = k, X_n = j) P(X_{n+1} = k | X_n = j) = \\ b_k(o_{jk}) P(Y_{n+2}, \dots, Y_N | X_{n+1}) a_{j|k}.$$



Scaling Recursions

In the denominators we have

$$\begin{aligned} &P(Y_{n+1}, \dots, Y_N \mid Y_0, \dots, Y_n) = \\ &= P(Y_{n+1} \mid Y_0, \dots, Y_n) \cdot P(Y_{n+2}, \dots, Y_N \mid Y_0, \dots, Y_{n+1}). \end{aligned}$$

Thus we have

$$\tilde{\beta}_n(j) = \sum_{k=1}^J \frac{b_k(\sigma_{jk}) P(Y_{n+2}, \dots, Y_N \mid X_{n+1}) a_{j|k}}{P(Y_{n+1} \mid Y_0, \dots, Y_n) \cdot P(Y_{n+2}, \dots, Y_N \mid Y_0, \dots, Y_{n+1})},$$

By definition of $\tilde{\beta}_n(j)$ the last equality gives

$$\tilde{\beta}_n(j) = \sum_{k=1}^J \frac{b_k(o_{j_k}) a_{j|k}}{P(Y_{n+1} | Y_0, \dots, Y_n)} \cdot \frac{P(Y_{n+2}, \dots, Y_N | X_{n+1} = k)}{P(Y_{n+2}, \dots, Y_N | Y_0, \dots, Y_{n+1})},$$

which equals

$$\tilde{\beta}_n(j) = \frac{1}{P(Y_{n+1} | Y_0, \dots, Y_n)} \sum_{k=1}^J b_k(o_{j_k}) a_{j|k} \cdot \tilde{\beta}_{n+1}(k).$$

Scaling Recursions

Finally,

$$\begin{aligned}P(Y_{n+1} | Y_0, \dots, Y_n) &= \sum_{j=1}^J \sum_{k=1}^J P(Y_{n+1}, X_n = j, X_{n+1} = k | Y_0, \dots, Y_n) = \\ &= \sum_{j=1}^J \sum_{k=1}^J a_{j|k} b_j(o_{j_j}) \bar{\alpha}_{n-1}(j)\end{aligned}$$

using the preceding results. We set

$$N_n = \frac{1}{\sum_{j=1}^J \sum_{k=1}^J a_{j|k} b_j(o_{j_j}) \bar{\alpha}_{n-1}(j)}$$

and obtain

$$\bar{\beta}_n(j) = N_n \sum_{k=1}^J b_k(o_{j_k}) a_{j|k} \bar{\beta}_{n+1}(k).$$



Posterior Smoothing

The forward variable $\tilde{\alpha}_n(j)$ and the backward variable $\tilde{\beta}_n(j)$ are defined as

$$\tilde{\alpha}_n(j) = \frac{P(Y_0, Y_1, \dots, Y_n, X_n = j)}{P(Y_0, \dots, Y_n)}$$

and

$$\tilde{\beta}_n(j) = \frac{P(Y_{n+1}, \dots, Y_N | X_n = j)}{P(Y_{n+1}, \dots, Y_N | Y_0, \dots, Y_n)}.$$

Start:

$$\tilde{\alpha}_0(j) = N_0 b_j(o_0) \pi_j(0), \quad j = 1, \dots, J.$$

$$\tilde{\beta}_N(j) = 1$$

Recursions:

$$\tilde{\alpha}_{n+1}(j) = N_n \left[\sum_{i=1}^J \tilde{\alpha}_n(i) \cdot a_{i|j} \right] \cdot b_j(o_{n+1}).$$

and

$$\tilde{\beta}_n(j) = N_n \sum_{k=1}^J b_k(o_{j_k}) a_{j|k} \cdot \tilde{\beta}_{n+1}(k).$$

where in both cases

$$N_n = \frac{1}{\sum_{j=1}^J \sum_{k=1}^J a_{j|k} b_j(o_{l_j}) \tilde{\alpha}_{n-1}(j)}.$$



The scaled recursions above are immune to *underflow problems*.

We wish to find the state sequence that maximizes the probability

$$P(Y_0 = o_0, \dots, Y_N = o_N, X_0 = j_0, \dots, X_N = j_N)$$

by selection of $j_0 \dots j_N$, when the sequence $o_0 \dots o_N$ is fixed and the model λ is known and omitted in the notation.

Alignment Problem

Let us set

$$\delta_n(j) = \max_{j_0 \dots j_{n-1}} P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j),$$

which is the highest probability along a single subsequence of states that at time n is in state j and accounts for the first $n + 1 \leq N$ emitted symbols.



Proposition

$$\delta_n(j) = \left[\max_{i=1, \dots, J} \delta_{n-1}(i) \cdot a_{i|j} \right] \cdot b_j(o_n).$$

Alignment Problem

Proof: Using the notational conventions we set

$$P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j) = P(Y^{(n)} | X^{(n)}) \cdot P(X^{(n)}).$$

By the conditional independence of the emitted symbols and the definition of conditional probability

$$= \prod_{j=1}^n P(Y_j | X_j) \cdot P(X_n | X^{(n-1)}) P(X^{(n-1)}).$$



Alignment Problem

But the Markov property of the hidden chain and some reorganization give

$$\begin{aligned} &= b_j(o_n) \prod_{j=1}^{n-1} P(Y_j | X_j) \cdot P(X_n | X_{n-1}) P(X^{(n-1)}) = \\ &= a_{j_{n-1}|j} \cdot b_j(o_n) \cdot \prod_{j=1}^{n-1} P(Y_j | X_j) \cdot P(X^{(n-1)}). \end{aligned}$$

Reverting back to the unabridged notation this equals

$$b_j(o_n) \left[a_{j_{n-1}|j} P(Y_0 = o_0, \dots, Y_{n-1} = o_{n-1}, X_0 = j_0, \dots, X_{n-1} = j_{n-1}) \right].$$

Alignment Problem

For each $j \in S$ at time n we have to find the transition to that state from every state $i \in S$ at time $n-1$ giving the best score (in the sense above). There are many paths leading to i at time $n-1$. But we see that the score for the current transition is factorized as the product

$$a_{j_{n-1}|j} P(o_0, \dots, o_{n-1}, X_0, \dots, X_{n-1}).$$

But this shows again that if we do not choose at time $n-1$ for every i that special subsequence leading to i with maximal probability, we cannot obtain

$$\max_{j_0 \dots j_{n-1}} a_{j_{n-1}|j} P(o_0, \dots, o_{n-1}, X_0, \dots, X_{n-1} = i).$$



Alignment Problem

Hence

$$\max_{j_0 \dots j_{n-1}} a_{j_{n-1}|j} P(o_0, \dots, o_{n-1}, X_0 = j_0, \dots, X_{n-1} = j_{n-1}) =$$
$$\left[\max_{i=1, \dots, J} \delta_{n-1}(i) \cdot a_{i|j} \right]$$

as was to be proved. □

Alignment Problem

The subsequence yielding $\delta_n(j)$ is called a *survivor* and denoted by

$$\psi_n(j) = \operatorname{argmax}_{i=1,\dots,J} \delta_{n-1}(i) \cdot a_{ij}$$

and consists of the prefix yielding $\delta_{n-1}(i)$ concatenated by the best scoring transition between times $n - 1$ and n . Hence we need at any j and any n only remember the survivor and no other path leading to this state through the trellis.

Alignment Problem

The complete procedure yielding the best decoded state sequence (path) is now formalizable in the following manner.

Viterbi Algorithm (VA)

Storage: n time index, for each $j \in S$ the survivor $\psi_n(j)$ and the corresponding scores $\delta_n(j)$, $j \in S$.

Start: $n = 0$. Compute for each $j \in S$

$$\delta_0(j) = \pi_j(o_0) \cdot b_j(o_0),$$

$$\psi_0(j) = \emptyset.$$

Recursion: Compute

$$\delta_{n+1}(j) = \left[\max_{i=1, \dots, J} \delta_n(i) \cdot a_{ij} \right] \cdot b_j(o_{n+1}).$$

Store the survivors

$$\psi_n(j) = \operatorname{argmax}_{i=1, \dots, J} \delta_{n-1}(i) \cdot a_{ij}.$$

for $j = 1, \dots, N$. Set $n+1$ to n unless $n = N$ and repeat.



Termination:

$$P^* = \max_{i=1,\dots,J} \delta_N(i)$$

$$j_*(N) = \operatorname{argmax}_{i=1,\dots,J} \delta_N(i).$$

Backtracking: The best path is found by

$$j_*(n) = \psi_{n+1}(j_*(n+1)), n = N-1, N-2, \dots, 0.$$



Let now t denote the number of state sequences $\mathbf{x} = j_0 j_1 \dots j_n$ of length $n + 1$ that have positive probability with regard to the model λ with the given sequence of emission symbols \mathbf{o} .

We enumerate the state sequences (j_0, \dots, j_n) by the index s , $s = 1, \dots, t$. Then we set

$$u_s = P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j_n; \lambda)$$

if $(j_0 \dots j_n) \mapsto s$.

Quasilikelihood for HMM

For any other model λ^* we set

$$v_s = P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0, \dots, X_n = j_n; \lambda^*)$$

if $(j_0 \dots j_n) \mapsto s$.

Note that some v_s may be in fact be equal to zero, since we are checking state paths with positive probability with regard to λ .

We have to exclude the converse situation and thus make the following assumption.



Quasilikelihood for HMM

We assume that the model λ^* does not assign a positive probability, conditioned on the given \mathbf{o} , to a state path in S^{n+1} that has probability zero with regard to the model λ or, if we have $\mathbf{x}^\dagger = j_0^\dagger \dots j_n^\dagger$ such that

$$P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0^\dagger, \dots, X_n = j_n^\dagger; \lambda) = 0$$

then

$$P(Y_0 = o_0, \dots, Y_n = o_n, X_0 = j_0^\dagger, \dots, X_n = j_n^\dagger; \lambda^*) = 0.$$



$$\ln \frac{P(Y_0 = o_0 \dots, Y_n = o_n; \lambda^*)}{P(Y_0 = o_0 \dots, Y_n = o_n; \lambda)},$$

which is comparing the plausibility of the two models for the fixed sequence of emitted symbols.

Under the assumptions above for $s = 1, \dots, t$ we have

$$u_s > 0$$

and

$$\ln \frac{P(Y_0 = o_0 \dots, Y_n = o_n; \lambda^*)}{P(Y_0 = o_0 \dots, Y_n = o_n; \lambda)} \geq \frac{Q(\lambda, \lambda^*) - Q(\lambda, \lambda)}{P(Y_0 = o_0 \dots, Y_n = o_n; \lambda)},$$

$$Q(\lambda, \lambda^*) = Q(\lambda, \lambda^* | \mathbf{o}) = \sum_{s=1}^t u_s \ln v_s$$

and

$$Q(\lambda, \lambda) = Q(\lambda, \lambda | \mathbf{o}) = \sum_{s=1}^t u_s \ln u_s.$$

Quasiloglikelihood for HMM

$$\begin{aligned} Q(\lambda, \lambda^*) &= \sum_{s=1}^t u_s \ln v_s = \\ &= \sum_{s=1}^t u_s \left[\sum_{j=1}^J r_j(s) \ln \pi_j^*(0) + \right. \\ &\quad \left. \sum_{j=1}^J \sum_{k=1}^K m_{j|k}(s) \ln b_j^*(o_k) + \sum_{j=1}^J \sum_{i=1}^J n_{ij}(s) \ln a_{ij}^* \right] = \end{aligned}$$

(interchanging the order of the finite summations)

We maximize

$$\begin{aligned} & \sum_{j=1}^J \left[\sum_{s=1}^t u_s r_j(s) \right] \ln \pi_j^*(0) + \\ & \sum_{j=1}^J \sum_{k=1}^K \left[\sum_{s=1}^t u_s m_{j|k}(s) \right] \ln b_j^*(o_k) + \\ & \sum_{j=1}^J \sum_{i=1}^J \left[\sum_{s=1}^t u_s n_{i|j}(s) \right] \ln a_{ij}^*. \end{aligned}$$

as function of the unknown parameters. This gives:

1. For $j = 1, \dots, J$,

$$\pi_j^*(0) = \frac{e_j}{P(Y_0 = \alpha_0, \dots, Y_n = \alpha_n; \lambda)}.$$

2. For $j = 1, \dots, J$ and for $k = 1, \dots, K$

$$b_j^*(\alpha_k) = \frac{d_{j|k}}{n_j}.$$

3. For $j = 1, \dots, J$ and for $i = 1, \dots, J$

$$a_{ij}^* = \frac{c_{ij}}{\sum_{j=1}^J c_{ij}}.$$

1. For $j = 1, \dots, J$,

$$\pi_j^*(0) = \frac{\alpha_0(j) \cdot \beta_0(j)}{P(Y_0 = o_0, \dots, Y_n = o_n)}. \quad (1)$$

π_j^* is the expected frequency of j at starting time given $o_0 \dots o_n$ and conditioned on the current model λ .

2. For $j = 1, \dots, J$ and for $k = 1, \dots, K$

$$b_j^*(o_k) = \frac{\sum_{l=0}^n I_{\{Y_l=o_k\}} \alpha_l(j) \cdot \beta_l(j)}{\sum_{l=0}^n \alpha_l(j) \cdot \beta_l(j)}. \quad (2)$$

$b_j^*(o_k)$ is the expected number of visits in state j and emitting the symbol o_k divided by the expected number of transitions from state j , given $o_0 \dots o_n$ and conditioned on the current model λ .

3. For $j = 1, \dots, J$ and for $i = 1, \dots, J$

$$a_{i|j}^* = \frac{a_{i|j} \cdot \sum_{l=0}^{n-1} \alpha_l(i) \cdot b_j(o_{l+1}) \cdot \beta_{l+1}(j)}{\sum_{l=0}^{n-1} \alpha_l(i) \cdot \beta_l(i)}. \quad (3)$$

$a_{i|j}^*$ is the ratio of the expected number of transitions from state i to state j divided by the expected number of transitions from state i given $o_0 \dots o_n$ and conditioned on the current model λ .

Consider a Markov chain $(X_k)_{k=0}^{\infty}$ with the state space $\{0, 1\}$ and with the transition probability matrix

$$A = \begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix}$$

and with the initial distribution

$$\pi(0) = (1-a, a).$$

The emitted sequence $(Y_k)_{k=0}^{\infty}$ is given by

$$Y_k = \begin{cases} 1 & \text{if } X_k + V_k \geq 1 \\ 0 & \text{if } X_k + V_k \leq 0, \end{cases}$$

where $(V_k)_{k=0}^{\infty}$ is a sequence of independent, identically distributed discrete random variables, which are independent of $(X_k)_{k=0}^{\infty}$, too. The variables V_k assume values in the alphabet $\{-1, 0, 1\}$ with the probabilities

$$1 - p_0 - p_1, p_0, p_1$$

respectively.



A Problem

- (a) Show that this is a hidden Markov model in the sense of our definition. Give the emission probability matrix B .
- (b) Let for $j = 0, 1$

$$\hat{\pi}_j(n|m) = P(X_n = j | Y_0, \dots, Y_m)$$

be the prediction ($n > m$) or filtering ($n = m$) probability. Show that

$$\hat{\pi}_1(n+1|n) = p - (p - (1 - q)) \cdot \hat{\pi}_1(n|n).$$

