

# HMM: Introduction and segmentation

J. Lember

University of Tartu, Estonia

6.12.2019

## (Finite) Markov chain

Let  $S$  be a finite or countable set – **state space**; let for every  $t = 1, \dots, n$ ,  $Y_t$  be a random variable taking values in  $S$  (random state).

The random variables  $Y_1, \dots, Y_n$  form **Markov chain** when for every  $y_i \in S$  and  $t = 2, \dots, n - 1$

$$P(Y_{t+1} = y_{t+1} | Y_t = y_t, \dots, Y_1 = y_1) = P(Y_{t+1} = y_{t+1} | Y_m = y_t)$$

$$\text{equivalently } P(y_{t+1} | y_t, \dots, y_1) = P(y_{t+1} | y_t)$$

MC is sometimes denoted

$$Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_n.$$

Thus  $Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_n$  iff for every  $y_1, \dots, y_n$  such that  $y_t \in S$

$$P(y_1, \dots, y_n) = P(y_1)p(y_2|y_1)P(y_3|y_2) \cdots P(y_n|y_{n-1}).$$

## Elementary properties:

1) reversed MC is also a MC:

if  $Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_n$ , then  $Y_n \rightarrow Y_{n-1} \rightarrow \dots \rightarrow Y_1$

2) Every sub-chain of MC is a MC:

if  $Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_n$ , then  $Y_{n_1} \rightarrow Y_{n_2} \rightarrow \dots \rightarrow Y_{n_k}$ .

3) If  $Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_n$ , then for every  $t < n$  and  $y_i \in S_i$

$$P(y_n, \dots, y_{t+1} | y_t, \dots, y_1) = P(y_n, \dots, y_{t+1} | y_t).$$

4)  $Y_1 \rightarrow \dots \rightarrow Y_n$  iff for every  $t = 2, \dots, n - 1$  the random variables  $Y_1, \dots, Y_{t-1}$  and  $Y_{t+1}, \dots, Y_n$  are conditionally independent given  $Y_t$ :  
for every  $y_t \in S_t$ ,

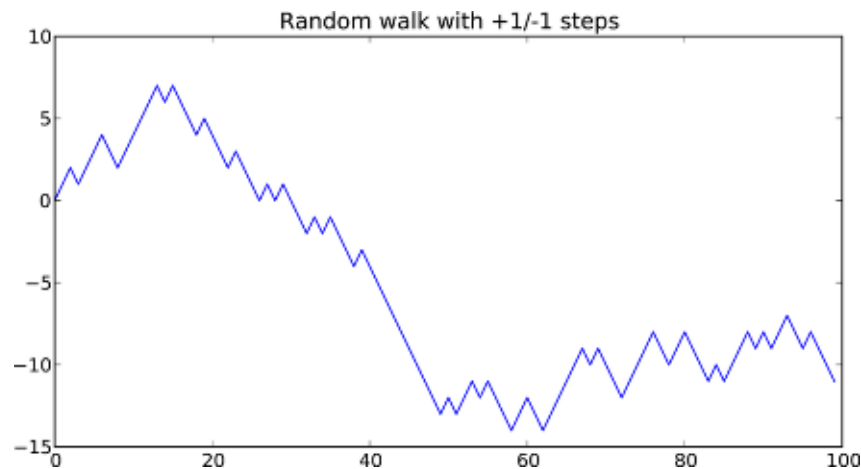
$$P(y_1, \dots, y_{t-1}, y_{t+1}, \dots, y_n | y_t) = P(y_1, \dots, y_{t-1} | y_t) P(y_{t+1}, \dots, y_n | y_t).$$

**Examples:** 1) Independent random variables  $Y_1, \dots, Y_n$ .

2) Simple random walk. Let  $\xi_1, \dots, \xi_n$  be independent tosses of fair coin, i.e.  $P(\xi_i = -1) = P(\xi_i = +1) = 0.5$ , let

$$Y_1 = \xi_1, \quad Y_2 = \xi_1 + \xi_2, \dots, Y_k = \sum_{i=1}^k \xi_i, \quad k = 1, \dots, n.$$

In general, sums of independent random variables is MC. Here  $S = \mathbb{Z}$ .



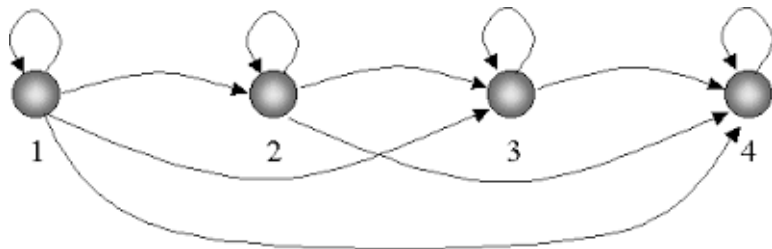
3) Polya urn. In a urn, there is  $m_b$  black balls and  $m_w$  white balls. Let  $Y_0 = m_b$ . Then take a random ball and return the ball along with another ball with the same color. Let  $Y_1$  be now the number of black balls; repeat the procedure etc. Here  $S = \mathbb{N}$ .

4) **Birth-death processes.**  $Y_t$  is the size of population at time  $t$ , every time  $t$  there is a birth and death probability  $p_t$  and  $q_t$ :

$$P(Y_{t+1} = k + 1 | Y_t = k) = p_t, \quad P(Y_{t+1} = k - 1 | Y_t = k) = q_t,$$
$$P(Y_{t+1} = k | Y_t = k) = 1 - q_t - p_t.$$

Polya urn is a birth-process

5) **Left-right MC.** Let  $a_0, a_1, \dots, a_m$  be states (letters) and if chain is at  $a_k$ , then it can: stay or jump to the next or second next (given it exists) state. Models duration from  $a_0$  to  $a_m$ . Used in speech recognition for modelling the length of a phonem.



6) **Bayesian model.** Random parameter  $\theta$  (prior probability), data  $X$  and a statistic  $T(X)$  (function of data). Thus  $\theta \rightarrow X \rightarrow T(X)$ .

## Homogenous MC

For any time  $t$ , the **transition matrix**  $\mathbb{P}_t$  is  $|S| \times |S|$ -matrix with entries being the **transition probabilities**

$$p_t(i, j) := P(Y_{t+1} = s_j | Y_t = s_i).$$

MC is **homogeneous** if transition probabilities do not depend on  $t$ , i.e.  $\mathbb{P}_t = \mathbb{P}$  for every  $t$ . In this case, we speak about **transition matrix of Markov chain**.

Examples: Polya urn. Not homogenous MC.

Simple random walk. Homogeneous MC,  $S = \mathbb{Z}$  and

$$p(i, j) = \begin{cases} \frac{1}{2}, & \text{if } j = i + 1; \\ \frac{1}{2}, & \text{if } j = i - 1; \\ 0, & \text{else;} \end{cases} \quad \mathbb{P} = \begin{pmatrix} \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & 0 & 1/2 & 0 & 1/2 & 0 & 0 & \dots \\ \dots & 0 & 0 & 1/2 & 0 & 1/2 & 0 & \dots \\ \dots & 0 & 0 & 0 & 1/2 & 0 & 1/2 & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

Two state MC.  $S = \{0, 1\}$ , transition matrix is

$$\begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix}$$

If  $\alpha < \beta < 0.5$  are small, then typical outcome is

0000000000**1111111111**000000000**1111**0000000000000000000**111111**00000...

Expected length of 0-block is  $\frac{1}{\alpha}$  and expected length of 1-block is  $\frac{1}{\beta}$ .

The proportion of ones is asymptotically?

The asymptotic proportion of 0 and 1 is according to **stationary distribution**  $(\pi(0), \pi(1))$ . MCMC method are based on that fact. In our case

$$\pi(1) = \frac{\alpha}{\alpha + \beta}.$$

**Three state "triangular" MC.**  $S = \{0, 1, 2\}$ , transition matrix is

$$\begin{pmatrix} 0.9 & 0.1 & 0 \\ 0 & 0.5 & 0.5 \\ 0.1 & 0 & 0.9 \end{pmatrix}$$

Being in a state, the chain either stays or jumps to the next state.  
Typical outcome

00000000**111**22222222220000000000000000**1**2222....



## Hidden Markov Model (HMM)

$Y_1, \dots, Y_n$  – homogeneous Markov Chain with finite alphabet  $S = \{1, \dots, K\}$ , and transition probabilities  $p(i, j)$  and initial probabilities  $\pi(i) := P(Y_1 = i)$ .

MC is sometimes called as the regime.

To each state  $y \in S$  corresponds an emission distribution  $P_y$  taking values on set  $\mathcal{X}$ . The set  $\mathcal{X}$  can be anything, typically either  $\mathbb{R}^d$  or a finite or countable set. In what follows, we assume that every  $P_y$  has a density  $f_y$  w. r. t. some reference measure  $dx$  (typically Lebesgue or counting measure).

### HMM:

To any realization  $y_1, y_2, \dots$  of  $Y$  corresponds a sequence of independent random variables  $X_1, X_2, \dots, X_n$ , where  $X_t \sim P_{y_t}$ .

Hence:

- 1) Given state sequence  $y_1, \dots, y_n$ , the random variables  $X_1, \dots, X_n$  are conditionally independent;
- 2) The distribution of  $X_t$  depends on  $y_t$  (independent of the rest of the states).

Hence for every state sequence  $y_1, \dots, y_n$  and (measurable) subsets  $A_1, \dots, A_n \subset \mathcal{X}$ , it holds:

$$\begin{aligned} & P(X_1 \in A_1, \dots, X_n \in A_n | Y_1 = y_1, \dots, Y_n = y_n) \stackrel{1)}{=} \\ & \prod_{t=1}^n P(X_t \in A_t | Y_1 = y_1, \dots, Y_n = y_n) \stackrel{2)}{=} \prod_{t=1}^n P(X_t \in A_t | Y_t = y_t) = \\ & \prod_{t=1}^n P_{y_t}(A_t) = \prod_{t=1}^n \int_{A_t} f_{y_t}(x) dx. \end{aligned}$$

**Exercise:** Write down full probability:

$$P(X_1 \in A_1, \dots, X_n \in A_n; Y_1 = y_1, \dots, Y_n = y_n).$$

The process  $X$  is sometimes called **observation process or hidden Markov process**. Can be taken as measurements of MC  $Y$  with error.

Some general properties (prove them as exercise):

1) In general,  $X_1, \dots, X_n$  is **not** a MC. It has longer memory. In particular, given a realization of  $X$ , it is difficult to prove/disprove that they are from a HMM;

2) 2D-process  $(X_1, Y_1), \dots, (X_n, Y_n)$  is a Markov process, when  $\mathcal{X}$  is countable then a Markov chain. Thus, given a realization of  $(X, Y)$ , it is very easy to prove/disprove that they are from a HMM;

3) Conditionally on a realization  $x_1, \dots, x_n$  of  $X$ , the  $Y$ -process (regime) is **inhomogeneous** MC:

$$P(y_{t+1}|y_t, \dots, y_1; x_1, \dots, x_n) = P(y_{t+1}|y_t; x_1, \dots, x_n).$$

Obviously the opposite is true as well: given  $y_1, \dots, y_n$ , the observations  $X_1, \dots, X_n$  are independent, thus Markov process.

## Example: GpC-islands.

Let the regime be a 2-state MC with transition matrix ( $S = \{0, 1\}$ )

$$\begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix},$$

where  $\alpha, \beta < 0.5$ . Let the random variables take values on  $\mathcal{X} = \{a, t, c, g\}$  with emission probabilities:

$$P_0(a) = P_0(t) = P_0(c) = P_0(g) = 0.25,$$
$$P_1(a) = P_1(t) = 0.1, \quad P_1(c) = P_1(g) = 0.4.$$

Thus, in regime 0 all letters are equiprobable, in regime 1  $g$  and  $c$  have much larger probability. Typical outcomes

*a c g t a c g t a a g t c c g t c c a g t t g c a a c g c t c g t g a a t c t t a g c a g g c c g a c g g t.*  
0000000000 **1111111111** 00000000 **1111** 000000000000000000 **111111** 00..

Example: pHMM – profile HMM (for pairwise alignments).

Three states: match (m), insertion (i), deletion (d).

Emissions: state match emits pairs (for example (a,g)); insertion emits residue against gap (for example (-,g)); deletion emits residue against gap (for example (a,-)).

An example of outcome is

a	t	t	g	g	a	t	-	g	t
a	a	-	g	c	-	t	g	c	t
m	m	d	m	m	d	m	i	m	m

Every outcome – alignment (table) – has a probability (score) . Hence pHMM is just a way to measure the goodness of alignments.

**Forward - backward recursions.** We are now study the probabilities

$$P(Y_t = j | X_1 = x_1, \dots, X_n = x_n) =: p(y_t = j | x_1, \dots, x_n).$$

The probabilities  $p(y_t = j | x_1, \dots, x_n)$  are called:

**smoothing** (posterior) probabilities, when  $t < n$ ;

**filtering** (posterior) probabilities, when  $t = n$ ;

**prediction** (posterior) probabilities, when  $t > n$ .

These probabilities are often the basis of the inference, but how to calculate the efficiently? For that, let us for every  $j \in S$  define

**forward** and **backward** variables

$$\alpha(j, x_1, \dots, x_t) := p(x_1, \dots, x_t, y_t = j),$$
$$\beta(x_{t+1}, \dots, x_n | j) := \begin{cases} 1, & \text{if } t = n \\ p(x_{t+1}, \dots, x_n | y_t = j), & \text{if } t < n \end{cases}$$

The  $\alpha$ - and  $\beta$ -variables can be calculated recursively:

$$\begin{aligned}\alpha(j, x_1, \dots, x_{t+1}) &= p(y_{t+1} = j, x_1, \dots, x_{t+1}) \\ &= \sum_{i \in S} p(y_{t+1} = j, y_t = i, x_1, \dots, x_{t+1}) \\ &= \sum_{i \in S} p(y_{t+1} = j, x_{t+1} | y_t = i, x_1, \dots, x_t) p(y_t = i, x_1, \dots, x_t) \\ &= \sum_{i \in S} p(y_{t+1} = j, x_{t+1} | y_t = i, x_t) \alpha(i, x_1, \dots, x_t) \\ &= \sum_{i \in S} p(y_{t+1} = j, x_{t+1} | y_t = i) \alpha(i, x_1, \dots, x_t) \\ &= \sum_{i \in S} p(i, j) f_i(x_t) \alpha(i, x_1, \dots, x_t).\end{aligned}$$

Here  $=$  follows from the fact that  $(X, Y)$  is a Markov process. The obtained recursion:

$$\alpha(j, x_1, \dots, x_{t+1}) = \sum_{i \in S} p(i, j) f_i(x_t) \alpha(i, x_1, \dots, x_t)$$

is called **forward recursion** (calculate  $\alpha$ -variables for  $t = 1$ , then update for  $t = 2$  and so on until  $t = n$ ).

Similarly, there is **backward recursion** for calculating  $\beta$ -variables:

$$\beta(x_t, \dots, x_n | j) = \sum_{i \in S} \beta(x_{t+1}, \dots, x_n | i) p(j, i) f_i(x_t).$$

With  $\alpha(j, x_1, \dots, x_t)$  and  $\beta(x_{t+1}, \dots, x_n | j)$ , obviously

$$\begin{aligned} p(y_t = j, x_1, \dots, x_n) &= p(x_1, \dots, x_t, y_t = j) p(x_{t+1}, \dots, x_n | y_t = j, x_1, \dots, x_t) = \\ &= p(x_1, \dots, x_t, y_t = j) p(x_{t+1}, \dots, x_n | y_t = j) = \alpha(j, x_1, \dots, x_t) \beta(x_{t+1}, \dots, x_n | j). \end{aligned}$$

Here  $=$  comes from Markov property of  $(X, Y)$  and definition of HMM.

Thus

$$\begin{aligned} p(y_t = j | x_1, \dots, x_n) &= \frac{p(y_t = j, x_1, \dots, x_n)}{p(x_1, \dots, x_n)} = \frac{p(y_t = j, x_1, \dots, x_n)}{\sum_{i \in S} p(y_t = i, x_1, \dots, x_n)} \\ &= \frac{\alpha(j, x_1, \dots, x_t) \beta(x_{t+1}, \dots, x_n | j)}{\sum_{i \in S} \alpha(i, x_1, \dots, x_t) \beta(x_{t+1}, \dots, x_n | j)}. \end{aligned}$$

Hence, to calculate smoothing probabilities, run forward and backward recursion (both linear in  $n$ ) and then apply the formula above.



In practice, there might happen numerical underflow – both  $\alpha$  and  $\beta$  values are too small for computer. Then several **rescaled versions** of the  $\alpha$  and  $\beta$ -variables can be used. Most important: the recursions still hold.

## Segmentation (decoding)

A realization  $x_1, \dots, x_n$  of  $X_1, \dots, X_n$  is observed. They are sample, data, observations.

The corresponding regime  $:= y_1, \dots, y_n$ , the realization of  $Y_1, \dots, Y_n$  is not observed ( $Y$  is hidden).

**The problem of segmentation:** To find out / prognose / estimate the hidden state sequence  $y_1, \dots, y_n$  – the hidden path.

**Warning:** There is no way to find out exact hidden path! All meaningful methods give a path which is the best in some sense.

**Example:** In CpG-island example, thus only the DNA-sequence

*acgtacgtaagtcctccagttgcaacgctcgtgaatcttagcaggccgacggt*

is observed. The goal is to (find a good method to) estimate the islands.

**Example: DNA copy number alternations.** States: DNA copy numbers (typically 2).

Emissions: comparative genomic hybridization (CGH). Typically modeled via Gaussian mixtures. The observations: measurements CGH, the goal is to (find a good method to) estimate copy numbers.

**Example: Secondary structure of a protein.** States: types of foldings. Emissions: 20 amino-acids. The goal is to guess the folding (secondary) structure of a protein (sequence of amino acids).

Framework of statistical learning (pattern recognition). Formally, we are looking for a function – classifier –

$$g = (g_1, \dots, g_n) : \mathcal{X}^n \rightarrow S^n,$$

that maps the observed sequence into state sequence. What is the best  $g$  ?

**Notation:**  $x^n := x_1, \dots, x_n$ ,  $s^n = s_1, \dots, s_n$ ,  $y^n := y_1, \dots, y_n$ , so on.

Given  $x^n = x_1, \dots, x_n$  define the conditional risk

$$R(\cdot | x^n) : S^n \rightarrow [0, \infty]$$

so that  $R(s^n | x^n)$  measures the goodness of state sequence  $s^n$  given observations  $x^n$ . The best classifier is that with the minimal risk over all possible classifiers. It can be obtained by minimizing the conditional risk:

$$g^*(x^n) := \arg \min_{s^n \in S^n} R(s^n | x^n).$$

The conditional risk  $R(\cdot|x^n)$  depends on the task, often it is meaningful to define it via **loss function**

$$L : S^n \times S^n \rightarrow [0, \infty],$$

where  $L(a^n, b^n)$  is loss, when the actual state sequence is  $a^n$  and the prognose is  $b^n$ . The conditional expectation

$$R(s^n|x^n) := E[L(Y^n, s^n)|X^n = x^n]$$

is the conditional risk of  $y^n$ . Then minimal-risk (or best) classifier is

$$g^*(x^n) = \arg \min_{s^n \in S^n} \sum_{a^n \in S^n} L(y^n, s^n) P(Y^n = a^n | X^n = x^n).$$

---

Usually no (or few) training data - unsupervised (or semisupervised) learning.

## Standard alignments

**Viterbi alignment.** Let  $L$  be

$$L(a^n, b^n) = \begin{cases} 1, & \text{when } a^n \neq b^n; \\ 0, & \text{when } a^n = b^n. \end{cases} \quad (1)$$

Then the conditional risk (in this case denoted by  $R_\infty$ ) is

$$R_\infty(s^n|x^n) = 1 - P(Y^n = s^n|X^n = x^n)$$

and the minimal-risk classifier – for (1) denoted by  $v$  – maps every sequence of observations into sequence  $s^n$  with maximum (posterior) probability

$$v(x^n) := \arg \max_{s^n \in S^n} P(Y^n = s^n|X^n = x^n).$$

The maximum-prob state sequence  $v(x^n)$  is known as **Viterbi alignment (path)**. It inherits its name from the **Viterbi algorithm** – a dynamic programming algorithm for finding  $v(x^n)$ . Viterbi algorithm is very simple and Viterbi alignment is by far most popular classifier in segmentation.

**Viterbi algorithm.** How to find a path (or all of them) that maximizes the probability  $p(s^n, x^n)$  (or, equivalently,  $p(s^n|x^n)$ ) over all possible paths? For given  $s^n$ , the joint probability (or conditional probability) is easy to evaluate (linear in  $n$ ), but there are  $K^n$  possible paths. It is a big number even for a moderate  $n$  (there are about  $10^{80} \approx 2^{270}$  atoms in the universe).

Due to the Markov property (again) **(Bellman's) optimality principle** holds and a very elementary dynamic programming algorithm does the work.

Let, for every  $t = 1, \dots, n$  and every  $j \in S$ ,

$$\delta_t(j) := \max_{y_1, \dots, y_{t-1}} p(\overbrace{y_1, \dots, y_{t-1}}^{y^{t-1}}, y_t = j, \overbrace{x_1, \dots, x_t}^{x^t}) = \max_{y^{t-1} \in S^{t-1}} p(y^{t-1}, y_t = j, x^t).$$

$$\text{Thus } \delta_1(j) = p(y_1 = j, x_1) = \pi(j) f_j(x_1).$$

Now observe that

$$\begin{aligned} p(y^{t-1}, y_t = i, y_{t+1} = j, x^{t+1}) &= p(y^{t-1}, y_t = i, x^t) p(y_{t+1} = j, x_{t+1} | y_t = i) \\ &= p(y^{t-1}, y_t = i, x^t) p(i, j) f_j(x_{t+1}). \end{aligned}$$

Hence

$$\begin{aligned} \delta_{t+1}(j) &= \max_{y^t \in S^t} p(y^t, y_{t+1} = j, x^{t+1}) = \max_{i \in S} \left( \max_{y^{t-1} \in S^{t-1}} p(y^{t-1}, y_t = i, y_{t+1} = j, x^{t+1}) \right) \\ &= \max_{i \in S} \left( \overbrace{\max_{y^{t-1} \in S^{t-1}} p(y^{t-1}, y_t = i, x^t)}^{\delta_t(i)} p(i, j) f_j(x_{t+1}) \right) = \max_{i \in S} (\delta_t(i) p(i, j)) f_j(x_{t+1}) \end{aligned}$$

Let

$$i_t(j) := \arg \max_{i \in S} \delta_t(j) p(i, j).$$



## Viterbi algorithm:

1) **Initialize:** For every  $j \in S$ , define  $\delta_1(j) := \pi_j f_j(x_1)$ ;

2) **Do for**  $t = 1, \dots, n - 1$ :

Update  $\delta_{t+1}(j) = \max_i (\delta_t(i) p(i, j)) f_j(x_{t+1})$ ;

Record  $i_t(j) := \arg \max_i \delta_t(i) p(i, j)$

3) **Output:** Find Viterbi alignment  $v^n$  by backtracking:

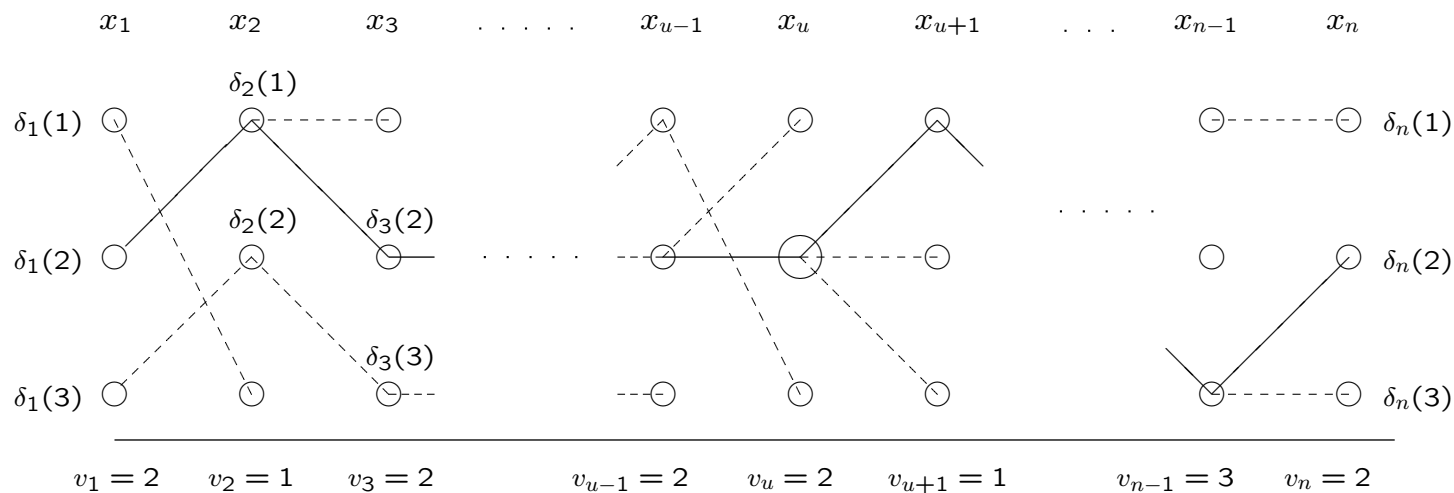
$$v_n := \arg \max_j \delta_n(j), \quad v_t = i_t(v_{t+1}), \quad t = n - 1, \dots, 1.$$

---

The numerical underflow can be the issue again. A possible solution is to use  $\ln \delta_t(j) =: \tilde{\delta}_t(j)$  so that the recursion is

$$\tilde{\delta}_{t+1}(j) = \max_i \left( \tilde{\delta}_t(i) + \ln p(i, j) \right) + \ln f_j(x_{t+1}).$$

The following picture illustrates Viterbi algorithm in action. The solid lines indicates the output, the dashed lines indicate  $i_t(j)$ .



The Viterbi algorithm is extremely easy to program and understand, therefore it is very popular. The output of it is the path with (conditional) probability. How does it look in practice?

Consider a MC with  $S = \{0, 1\}$   $\pi(0) = 3/13$ ,  $\pi(1) = 10/13$  (stationary distribution) and transition matrix

$$\begin{pmatrix} 0.6 & 0.4 \\ 0.3 & 0.7 \end{pmatrix}$$

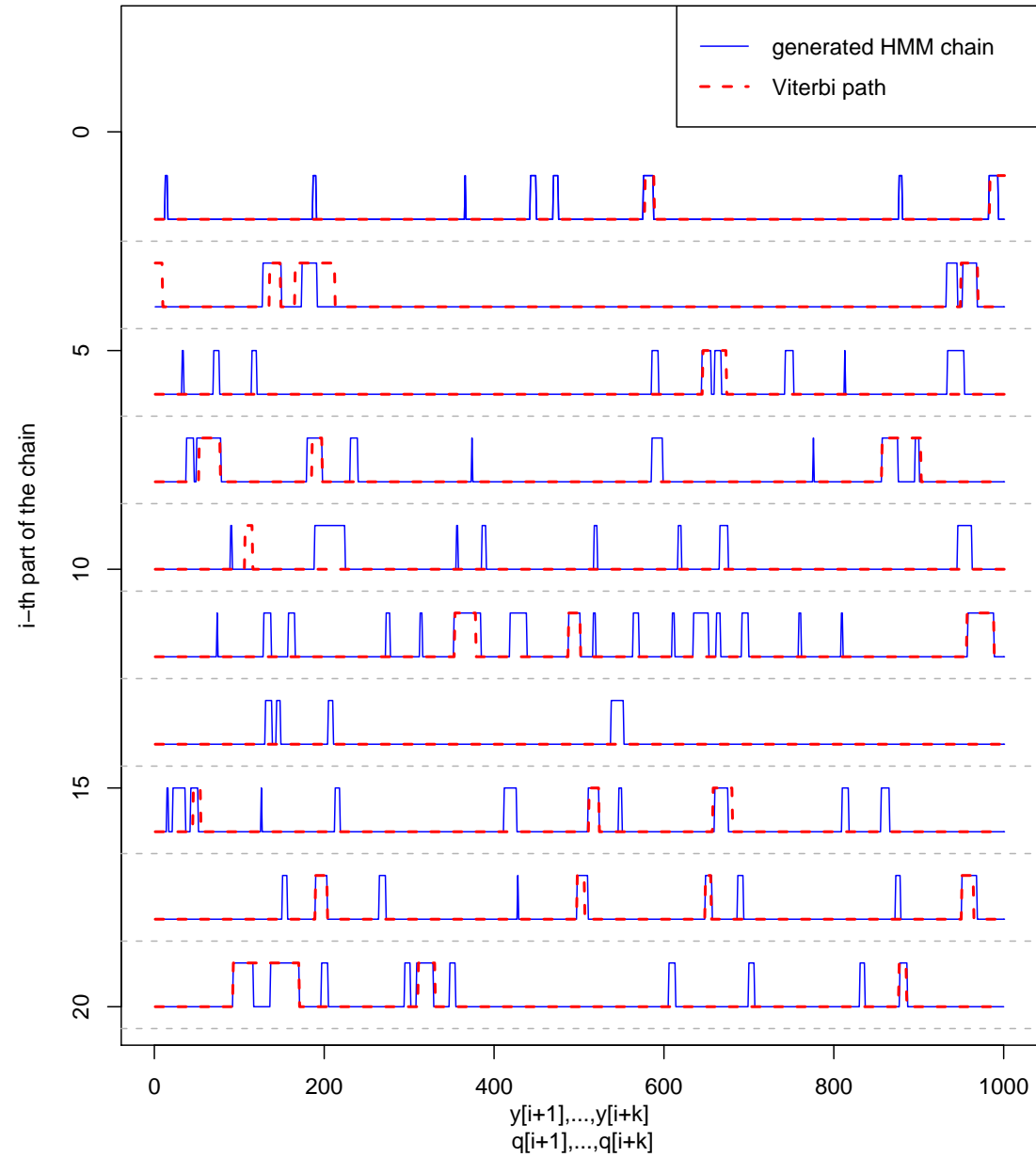
The typical realization of this chain consists of independent blocks of 0's and 1's having lengths distributed as  $G(0.4)$  and  $G(0.3)$ , respectively

$$\overbrace{00}^{G(0.4)} \quad \overbrace{111}^{G(0.3)} \quad 00000 \quad 11 \quad 000 \quad 1111\dots$$

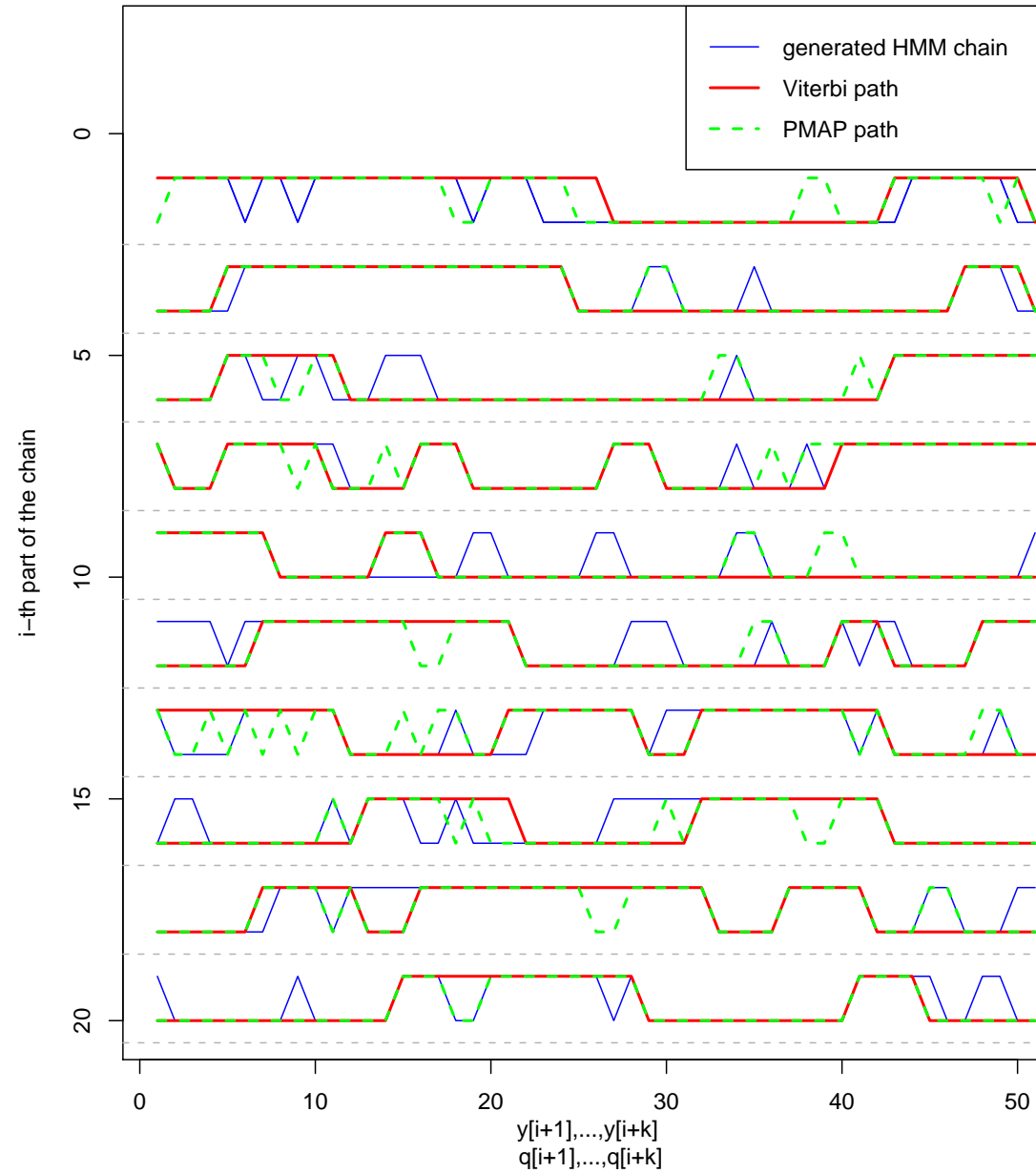
How does the path with maximal probability look like?

Here are no data. With data the path changes but the long block structure prevails. Hence no wonder that Viterbi path (for models like this) overestimates the lengths of the blocks. It is highly untypical.

Comparison of Viterbi alignment and the actual HMM chain, splitted into parts with length  $l$ .  
 Parameters:  $p_{11} = 0.90$ ,  $p_{22} = 0.99$ , emission distributions are discrete and non-synchronous.  
 Note: read the graph like normal text.



on of Viterbi alignment, PMAP alignment and the actual HMM chain, splitted into parts w  
 odel parameters:  $p_{11} = 0.80$ ,  $p_{22} = 0.80$ , emssion distributions are discrete and symr  
 Note: read the graph like normal text.



**PMAP alignment.** The loss (1) penalizes all differences alike: no matter whether  $a^n$  and  $b^n$  differ from one entry or all entries, the penalty is one. Often it is natural to penalize every different entry. This can be done by **pointwise loss-function**

$$l : S \times S \rightarrow [0, \infty) \quad \text{where } l(s, s) = 0, \quad \forall s \in S. \quad (2)$$

Using  $l$ , we can define loss-function  $L$  as follows

$$L(a^n, b^n) := \frac{1}{n} \sum_{t=1}^n l(a_t, b_t).$$

Then the minimal-risk classifier is obtained pointwise:  $g^* = (g_1^*, \dots, g_n^*)$ , where

$$g_t^*(x^n) = \arg \min_{s \in S} E[l(Y_t, s) | X^n = x^n]. \quad (3)$$

**Exercise:** Show (3).

Counting errors. If

$$l(a, a') = \begin{cases} 0, & \text{if } a = a'; \\ 1, & \text{if } a \neq a'. \end{cases},$$

then the loss-function  $L$  counts the differences between  $a^n$  and  $b^n$  when compared pairwise. Thus the conditional risk measures the **expected number of misclassification errors** of  $s^n$  given the observations are  $x^n$  and can be calculated as follows:

$$R_1(s^n|x^n) := 1 - \frac{1}{n} \sum_{t=1}^n P(Y_t = y_t|X^n = x^n).$$

So the minimal-risk classifier – let us denote it by  $u$  – minimizes the expected number of misclassification errors, from above we see that it can be calculated pointwise, indeed:

$$u_t(x^n) = \arg \max_{s \in S} P(Y_t = s|X^n = x^n) = \arg \max_{s \in S} p(y_t = s|x^n), \quad t = 1, \dots, n.$$

We shall call  $u$  as **PMAP (pointwise maximum a posteriori)**-alignment (path). Other names encountered: *marginal posterior mode*, *maximum posterior marginals*, *optimal symbol-by-symbol detection*, *symbol-by-symbol MAP estimation*, *MAP-state estimation*.

Typically in practice  $K$  is small so that maximizing the smoothing probabilities  $p(y_t = s|x^n)$  over  $S$  is not a problem. We also know that calculating the smoothing probabilities  $p(y_t = s|x^n)$  is not a problem, because it can be done by forward-backward recursions. Thus, finding PMAP path is computationally feasible and from computational point of view it has the same complexity as Viterbi path.

In general, PMAP and Viterbi paths are not equal. Thus **Viterbi path is not the best for misclassification errors point of view!**

Since PMAP path is purely local, then in the presence of 0-s in transition matrix, it might have probability 0. Thus, it might happen that  $p(u^n|x^n) = 0$ . Such paths are **inadmissible** or **forbidden** and (again) highly untypical (never occur).

Consider our MC again:  $S = \{0, 1\}$ ,  $\pi(0) = 3/13$ ,  $\pi(1) = 10/13$  (stationary distribution) and transition matrix  $\begin{pmatrix} 0.6 & 0.4 \\ 0.3 & 0.7 \end{pmatrix}$ . What is PMAP path in this case (without data)?



Logarithmic versions of standard risks:

$$\begin{aligned}\bar{R}_\infty(s^n|x^n) &:= -\frac{1}{n} \ln P(Y^n = s^n|x^n) = -\frac{1}{n} \ln p(s^n|x^n) \\ \bar{R}_1(s^n|x^n) &:= -\frac{1}{n} \sum_{t=1}^n \ln P(Y_t = s_t|x^n).\end{aligned}$$

Clearly

Viterbi alignment  $v(x^n)$  minimizes  $\bar{R}_\infty(\cdot|x^n)$

PMAP alignment  $u(x^n)$  minimizes  $\bar{R}_1(\cdot|x^n)$ .

**Between PMAP and Viterbi.** If the aim is to minimize the expected number of errors, then one should use PMAP-alignment  $u^n$ . Unfortunately, in the presence of forbidden transitions and/or emissions, it might be that  $p(u^n|x^n) = 0$ , i.e.  $u$  is **inadmissible**.

The simplest solution: **restricted optimization**:

$$\min_{s^n: p(s^n|x^n) > 0} R_1(s^n|x^n) \Leftrightarrow \max_{s^n: p(s^n|x^n) > 0} \sum_{i=1}^n P(Y_i = s_i|x^n). \quad (4)$$

Because of the restrictions, (4) is not equivalent to problem:

$$\min_{s^n: p(s^n|x^n) > 0} \bar{R}_1(s^n|x^n) \Leftrightarrow \max_{s^n: p(s^n|x^n) > 0} \sum_{t=1}^n \ln P(Y_t = s_t|x^n). \quad (5)$$

The solution of (5): **posterior Viterbi decoding (PVD)**.

**Note:** the solution of (5) does not change if  $p(s^n|x^n) > 0$  is replaced by  $p(s^n) > 0$ , but (4) is in general not equivalent to

$$\min_{s^n: p(s^n) > 0} R_1(s^n|x^n) \Leftrightarrow \max_{s^n: p(s^n) > 0} \sum_{t=1}^n P(Y_t = s_t|x^n). \quad (6)$$

In fact, **the solution of (6) can still be inadmissible!**

Algorithm for restricted optimization (4):

**1) Initialize:** For every  $1 \leq t \leq n$  and  $j \in S$  compute  $p(y_t = j|x^n)$  with forward-backward recursions. Set  $\delta_1(j) := p(y_1 = j|x^n)$ .

**2) Do for**  $t = 1, \dots, n - 1$ :

For every  $j \in S$  update

$$\delta_{t+1}(j) = \left( \max_i \delta_t(i) r_{ij} + p(y_{t+1} = j|x^n) \right) r_j^{t+1},$$

$$\text{where } r_{ij} := \mathbb{I}_{\{p(i,j) > 0\}}, \quad r_j^t := \mathbb{I}_{\{p(y_t=j|x^n) > 0\}}.$$

Record  $i_t(j) := \arg \max_i \delta_t(i) r_{ij}$ ,  $t = 1, \dots, n - 1$ .

**3) Output:** Find the optimal alignment  $s^n$  by backtracking:

$$s_n := \arg \max_j \delta_n(j), \quad s_t = i_t(s_{t+1}), \quad t = n - 1, \dots, 1.$$

## Algorithm for PVD (5):

**1) Initialize:** For every  $1 \leq t \leq n$  and  $j \in S$  compute  $p(y_t = j|x^n)$  with forward-backward recursions. Set  $\delta_1(j) := \ln(p(y_1 = j|x^n))$ .

**2) Do for**  $t = 1, \dots, n - 1$ :

For every  $j \in S$  update

$$\delta_{t+1}(j) = \left( \max_i \delta_t(i) + \ln(r_{ij}) \right) + \ln(p(y_{t+1} = j|x^n));$$

Record  $i_t(j) := \arg \max_i \delta_t(i) + \ln(r_{ij})$ ,  $t = 1, \dots, n - 1$ .

**3) Output:** Find the optimal alignment  $s^n$  by backtracking:

$$s_n := \arg \max_j \delta_n(j), \quad s_t = i_t(s_{t+1}), \quad t = n - 1, \dots, 1.$$

The solutions of constrained problems can still have small probability. This suggests to consider instead of (4) the following more general **penalized optimization problem**:

$$\min_{s^n} [R_1(s^n|x^n) + Ch(s^n, x^n)], \quad (7)$$

where  $C > 0$  and  $h(s^n, x^n)$  or  $h(s^n)$  is some penalty term. Some possibilities for  $h$ :

\*  $h(s^n) = I_{\{p(s^n)=0\}}$ . For  $C$  large, we get (6) .

\*  $h(s^n) = -\frac{1}{n} \ln p(s^n)$  (log-prior). Then (7) is

$$\max_{s^n} \left[ \sum_{i=1}^n P(Y_i = s_i|x^n) + C \ln p(s^n) \right]$$

and for  $C > 0$  small we get (6)

\*  $h(s^n, x^n) = \bar{R}_\infty(s^n|x^n)$  (log-posterior). Then (7) is

$$\min_{s^n} [R_1(s^n|x^n) + C\bar{R}_\infty(s^n|x^n)] \Leftrightarrow \max_{s^n} \left[ \sum_{i=1}^n P(Y_i = s_i|x^n) + C \ln p(s^n|x^n) \right] \quad (8)$$

and for  $C > 0$  small we get (4);

---

With  $\bar{R}_1$ -risk instead of  $R_1$ -risk, the  $\bar{R}_1$  counterpart of (8) is

$$\min_{s^n} [\bar{R}_1(s^n|x^n) + C\bar{R}_\infty(s^n|x^n)] \Leftrightarrow \max_{s^n} \left[ \sum_{t=1}^n \ln P(Y_t = s_t|x^n) + C \ln p(s^n|x^n) \right] \quad (9)$$

Small  $C$  gives the solution of (5): PVD.

The solutions of (8) and (9) are in general not the same, but both problems naturally interpolate between the two standard alignments: the case  $C = 0$  corresponds to PMAP and  $C$  very large – Viterbi.

If the alignment  $s^n$  has likelihood 0, then  $\bar{R}_\infty(s^n|x^n) = \infty$ . Thus for every  $C > 0$ , the solution has positive (posterior, hence prior) likelihood.

## Block PMAP alignment

**The idea (Rabiner):** Instead of maximizing (over all  $s^n \in S^n$ )

$$\sum_{t=1}^n P(Y_t = s_t | x^n)$$

(PMAP alignment) maximize (here  $a_r^l := (a_r, \dots, a_l)$ )

$$P(Y_1^k = s_1^k | x^n) + P(Y_2^{k+1} = s_2^{k+1} | x^n) + \dots + P(Y_{n-k+1}^n = s_{n-k+1}^n | x^n), \quad (10)$$

where  $k$  is the size of block. The case  $k = 1$  corresponds to the PMAP-alignment, the bigger  $k$ , the "closer" to the Viterbi alignment.

Minimizes the expected number of correctly guessed  $k$ -blocks. For example, when  $k = 2$ , then the minimizer gives a path that with highest expected number of correct transitions.

**Warning:** Can be of 0 likelihood even for  $k > 1$ . (To prove that is an easy exercise).

**A modification:** Replace  $+$  by  $\times$ . Denote, for  $k = 1, 2, \dots$

$$\bar{R}_k(s^n | x^n) := -\frac{1}{n} \sum_{t=1-k}^n \ln P\left(Y_{(t+1)\vee 1}^{(t+k)\wedge n} = s_{(t+1)\vee 1}^{(t+k)\wedge n} | x^n\right).$$

So, for  $k > 1$ ,  $-n\bar{R}_k(s^n | x^n)$  is basically the sum

$$\ln P(Y_1^k = s_1^k | x^n) + \ln P(Y_2^{k+1} = s_2^{k+1} | x^n) + \dots + \ln P(Y_{n-k+1}^n = s_{n-k+1}^n | x^n).$$

- \* PMAP alignment  $u(x^n)$  minimizes  $\bar{R}_1(\cdot | x^n)$ ;
- \* Let  $u_k(x^n)$  minimize  $\bar{R}_k(\cdot | x^n)$ , hence  $u_1(x^n) = u(x^n)$ .



**Lemma** Let  $n \geq k > 1$ . Then, for every state sequence  $s^n$ :

$$\bar{R}_k(s^n|x^n) = (k - 1)\bar{R}_\infty(s^n|x^n) + \bar{R}_1(s^n|x^n).$$

Thus  $u_k$  is the solution of (9):

$$\min_{s^n} [\bar{R}_1(s^n|x^n) + C\bar{R}_\infty(s^n|x^n)]$$

for  $C = k - 1$ . Gives a nice interpretation to (9) and relates it to Rabiner's  $k$ -block idea.

When  $k$  increases, then

$\bar{R}_\infty(u_k)$  decreases (likelihood increases)

$\bar{R}_1(u_k)$  risk increases.

for every  $k > 1$ , the alignment  $u_k(x^n)$  has positive likelihood.

## Generalizations

The risk

$$\bar{R}_1(s^n|x^n) + C\bar{R}_\infty(s^n|x^n)$$

can be obviously generalized to

$$C_1\bar{R}_\infty(s^n|x^n) + C_2\bar{R}_1(s^n|x^n)$$

( $C_1 \geq 0, C_2 \geq 0$ ) and, to incorporate prior information more,

$$C_1\bar{R}_\infty(s^n|x^n) + C_2\bar{R}_1(s^n|x^n) + C_3\bar{R}_\infty(s^n) + C_4\bar{R}_1(s^n),$$

where

$$\bar{R}_\infty(s^n) := -\frac{1}{n} \ln P(Y^n = s^n) \quad - \text{prior likelihood}$$

$$\bar{R}_1(s^n) := -\frac{1}{n} \sum_{t=1}^n \ln P(Y_t = s_t) \quad - \text{prior pseudolikelihood.}$$

Consider

$$\arg \min_{s^n} [C_1 \bar{R}_\infty(s^n|x^n) + C_2 \bar{R}_1(s^n|x^n) + C_3 \bar{R}_\infty(s^n) + C_4 \bar{R}_1(s^n)].$$

$C_1 > 0, C_2 = C_3 = C_4 = 0$  – Viterbi alignment;

$C_1 = 0, C_2 > 0, C_3 = C_4 = 0$  – PMAP alignment;

$C_1 = C_2 = 0, C_3 > 0, C_4 = 0$  – the alignment with maximal (prior) probability;

$C_1 = C_2 = C_3 = 0, C_4 > 0$  – the alignment with (average) minimal number of errors when no data;

$C_1 > 0, C_2 > 0, C_3 = C_4 = 0$  – generalized  $k$ -block alignment;

$C_1 > 0, C_2 = 0, C_3 > 0, C_4 = 0$  – generalized Viterbi alignment (with prior probabilities);

$C_1 = 0, C_2 > 0, C_3 > 0, C_4 = 0$  – generalized PMAP alignment (with prior probabilities – has always admissible solution!)

...

**How to find the solution?** There is a recursion that can be nicely implemented by a (Viterbi like) dynamic programming algorithm.

## Algorithm for optimization:

**1) Initialize:** For every  $1 \leq t \leq n$  and  $j \in S$  compute  $p(y_t = j|x^n)$  with forward-backward recursions and set

$$g_t(j) := C_1 \log p(y_t = j|x^n) + C_2 \log f_j(x_t) + C_3 \log p(y_t = j).$$

Set  $\delta_1(j) = C_1 \log p(y_1 = j|x^n) + (C_2 + C_3 + C_4) \log \pi_j + C_2 \log f_j(x_1)$ .

**2) Do for**  $t = 1, \dots, n - 1$ :

For every  $j \in S$  update

$$\delta_{t+1}(j) = \max_i \left( \delta_t(i) + (C_2 + C_4) \log(p(i, j)) \right) + g_{t+1}(j).$$

Record  $i_t(j) := \arg \max_i \delta_t(i) r_{ij}$ ,  $t = 1, \dots, n - 1$ .

**3) Output:** Find  $s^n$  by backtracking ( $t = n - 1, \dots, 1$ ):

$$s_n := \arg \max_j \max_i \left( \delta_t(i) + (C_2 + C_4) \log(p(i, j)) \right), \quad s_t = i_t(s_{t+1}).$$

Example.

$$\mathbb{P} = \begin{pmatrix} 0.99 & 0.01 & 0 \\ 0.3 & 0.3 & 0.4 \\ 0 & 0.02 & 0.98 \end{pmatrix}, \quad \pi = \begin{pmatrix} 0.5882 \\ 0.0196 \\ 0.3922 \end{pmatrix},$$

Emission distributions given by

$$p_1 = (0.3, 0.2, 0.2, 0.3), \quad p_2 = (0.1, 0.3, 0.3, 0.3), \quad p_3 = (1/6, 1/6, 1/6, 1/2).$$

Note: transitions  $1 \leftrightarrow 3$  not allowed.

Chain mostly stays in state "one", producing "islands" of other states, mostly of state "three".

An outcome is generated and different alignment compared.



## Experiments with real data from Protein Data Bank

Data:  $N = 25713$  proteins. Each of them is modeled as a (non-stationary) HMM. Hidden states are 6 type of structure motifs ( $\alpha$ -helix,  $\beta$ -helix, coil and their modifications). Many impossible transitions. Emissions: 20 amino-acids.

We have  $N$  realizations of  $(x^{n_j}, y^{n_j})$ ,  $j = 1, \dots, N$ . For every  $i$ , the parameters are estimated from remaining  $N - 1$  sequences and using these parameters, the sequence  $x^{n_j}$  is segmented using different methods (alignments). For every method, we calculated: the overall averaged error rate

$$\frac{\sum_{j=1}^N \sum_{i=1}^{n_j} I_{\{g_i(x^{n_j}) \neq y_i^{n_j}\}}}{\sum_{j=1}^N n_j}$$

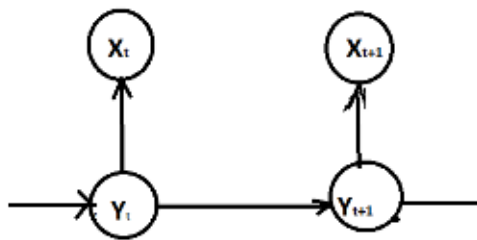
and averaged alignment probability

$$\frac{1}{N} \sum_{i=1}^n \left( P(g(x^{n_i}) | x^{n_i}) \right)^{\frac{1}{n_i}}.$$

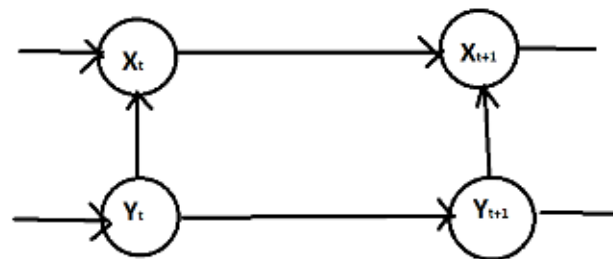
Pairwise Markov models (PMM) is a two-dimensional Markov chain

$$(X_1, Y_1), (X_2, Y_2), \dots,$$

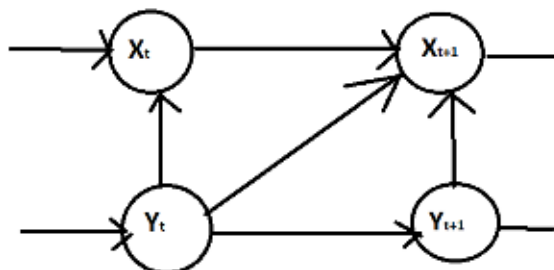
where, as previously,  $Y_t$  takes values in a finite set  $\mathcal{Y}$  and  $X_t$  takes values in  $\mathcal{X}$  (in general  $\mathbb{R}^d$ ). Introduced by W. Piezcyński. PMM is a very rich class of models, allowing dependence between observations (given the state sequence), HMM is just a narrow subclass.



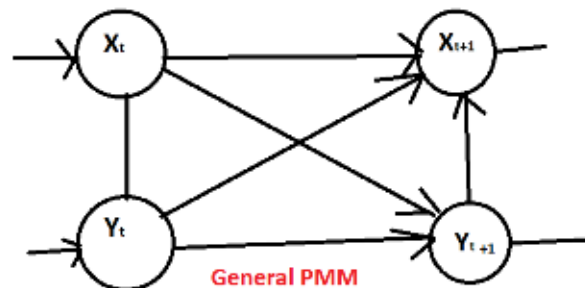
HMM



Markov switching model (Y is MC)



HMM-DN (Y is MC)



General PMM



In general neither  $Y$  nor  $X$  is a Markov chain, but in special cases one of the might be. Given  $X$ , the  $Y$ -sequence is a MC and vice versa.

**Example:** A Markov chain can be modeled as follows: start with a state  $y_1$  chosen according to to initial probabilities, then generate **duration** from Geometrical distribution  $G_{y_1}$ , where for every  $y \in S$

$$G_y(k) = (p_{yy})^k(1 - p_{yy}), \quad k = 0, 1, \dots$$

Then change the state ( $y_1 \rightarrow y_2$ ) with transition probabilities

$$q_{yy'} = \frac{p_{yy'}}{1 - p_{yy}}, \quad y \neq y', \quad q_{yy} = 0,$$

then generate a duration from  $G_{y_2}$  and so on.

When instead of Geometrical distribution, the duration is modeled via another distribution  $P_y$ , one obtains a **semi-Markov chain**.

It is not hard to see that semi-Markov process is a PMM.

Allows, among others, to model dependence between random sequences (DNA-sequences, for example). The following example shows how to jointly model two binary (1/2 valued) MC.

**Example:** Let  $\mathcal{X} = \mathcal{Y} = \{1, 2\}$ , and we want that  $X$  and  $Y$  were both Markov chains with the same transition matrix  $\begin{pmatrix} p & 1-p \\ q & 1-q \end{pmatrix}$

It can be achieved by considering PMM  $Z$  with transition matrix

$$\mathbb{Q} = \begin{array}{cc} & \begin{array}{cccc} (1,1) & (1,2) & (2,1) & (2,2) \end{array} \\ \begin{array}{c} (1,1) \\ (1,2) \\ (2,1) \\ (2,2) \end{array} & \begin{array}{cccc} p\lambda_1 & p(1-\lambda_1) & p(1-\lambda_1) & 1+p\lambda_1-2p \\ p\lambda_2 & p(1-\lambda_2) & q-p\lambda_2 & 1+p\lambda_2-q-p \\ q\mu_1 & q(1-\mu_1) & p-q\mu_1 & 1+q\mu_1-p-q \\ q\mu_2 & q(1-\mu_2) & q(1-\mu_2) & 1+q\mu_2-2q \end{array} \end{array}$$

where  $\lambda_1, \lambda_2, \mu_1, \mu_2$  are the parameters to control the dependence between  $X$  and  $Y$ . They are independent, if  $\lambda_1 = \mu_1 = p$  and  $\lambda_2 = \mu_2 = q$ ; the bigger  $\lambda_i$  and  $\mu_i$ , the more they are positively corr; the smaller parameters, the bigger neg. corr.

Examples of generated realizations of  $(X, Y)$  with different dependence (different  $\lambda_i$  and  $\mu_i$ ,  $p = q = 0.7$ )

$p=q=0.7$

Max

```
10011111110101101100001011111110111100111100111110110110110010010110
101111111101011011001010111111110111100111110111110111110111110110000010100
```

Min

```
111111111111100100111001111111111011000111100111011101110011101111011
111001100001111111000111000110011100111100011101101101111101111000100
```

Ind

```
001110110101111010110011111101011111111110111111001011010101101101110
001111011111111111101101100010110111111011101111111101111010011110111
```

Triplet Markov model (TMM) is a three-dimensional Markov chain

$$(X_1, Y_1, U_1), (X_2, Y_2, U_2), \dots,$$

where typically  $U_t$  takes values in a finite set. Introduced by W. Piez-cynski, allows **deep modeling**. Now, the pair  $(X, Y)$  – observations and classes – might not be a MC any more.

The third component allows to change the transition matrix and so the triplet model is very flexible to model inhomogeneous dynamics.

**Example:** A HMM with changing transition matrix can be modeled as TMM. Suppose we have two transition matrices

$$\mathbb{P}_A = \begin{pmatrix} p_A & 1 - p_A \\ q_A & 1 - q_A \end{pmatrix}, \quad \mathbb{P}_B = \begin{pmatrix} p_B & 1 - p_B \\ q_B & 1 - q_B \end{pmatrix},$$

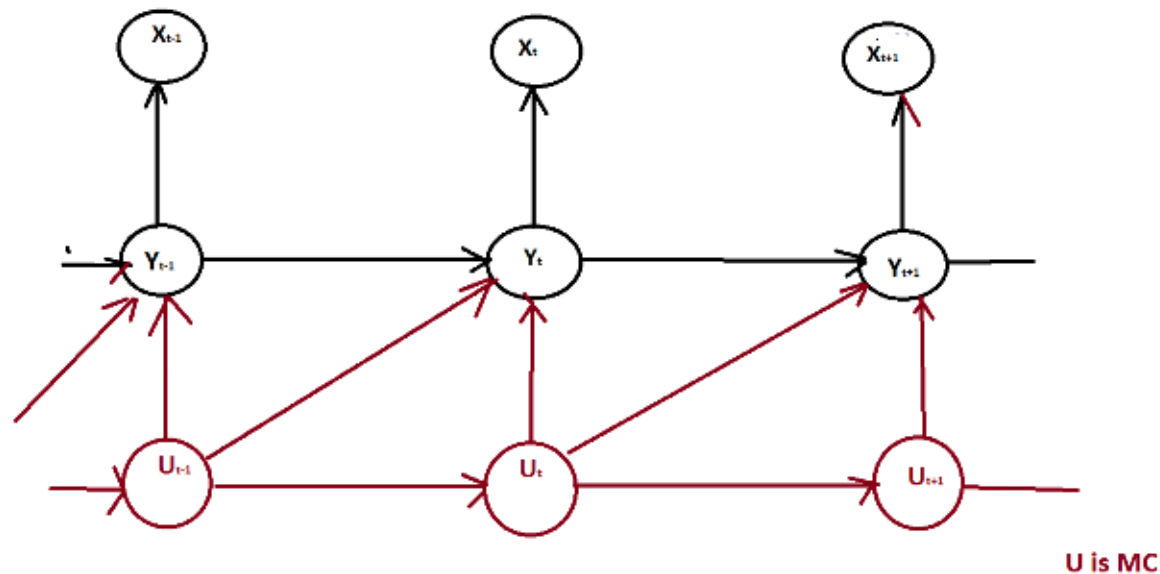
and with probability  $\epsilon > 0$  they change. This can be modelled as a PMM  $(Y, U)$  as follows  $U$  takes values  $A$  and  $B$ ,  $Y$  takes values 0 and 1 and the transition matrix of  $(Y, U)$  is

	(0, A)	(1, A)	(0, B)	(1, B)
(0, A)	$(1 - \epsilon)p_A$	$(1 - \epsilon)(1 - p_A)$	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$
(1, A)	$(1 - \epsilon)q_A$	$(1 - \epsilon)(1 - q_A)$	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$
(0, B)	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$	$(1 - \epsilon)p_B$	$(1 - \epsilon)(1 - p_B)$
(1, B)	$\frac{\epsilon}{2}$	$\frac{\epsilon}{2}$	$(1 - \epsilon)q_B$	$(1 - \epsilon)(1 - q_B)$

Conditionally on  $U_1 = u_1, \dots, U_n = u_n$  ( $u_t \in \{A, B\}$ ), the process  $Y_1, \dots, Y_n$  is a inhomogeneous MC with transition matrix

$$P(Y_{t+1} = j | Y_t = i, U_1 = u_1, \dots, U_n = u_n) = \begin{cases} p_A(i, j), & \text{if } u_t = u_{t+1} = A; \\ p_B(i, j), & \text{if } u_t = u_{t+1} = B; \\ 1/2, & \text{if } u_t \neq u_{t+1}; \end{cases}$$

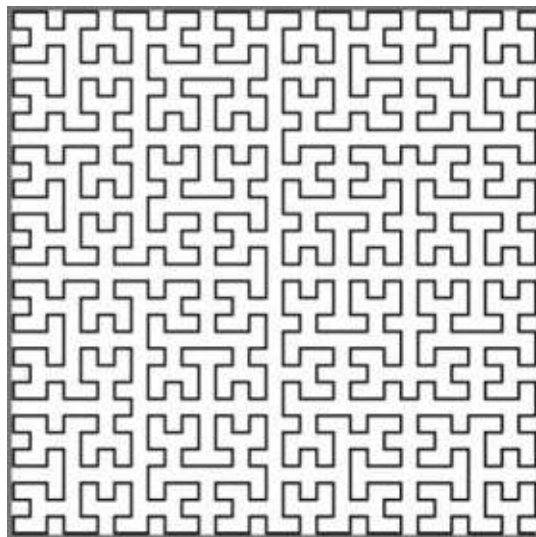
Like for HMM, let the observations be conditionally (given  $U$  and  $Y$  values) independent and depending only on  $Y_t$ . We have the following fairly simple TMM:



Many more examples possible, a way larger class of models than PMM. It can be that none of the processes:  $X$ ,  $Y$ ,  $U$ ,  $(X, Y)$ ,  $(Y, U)$ ,  $(U, X)$  MC, only  $(X, Y, U)$  is. Conditionally on  $U$ ,  $(X, Y)$  is a MC and so on.

**2D case.** 2D generalization of MC is **Markov random field (MRF)**, with independent noise **hidden MRF (HMRF)**. PMM and TMM are just two or three-component MRFs. **Problem with fields: no dynamics! In particular – no Viterbi algorithm, no forward-backward algorithms, no EM-algorithms.**

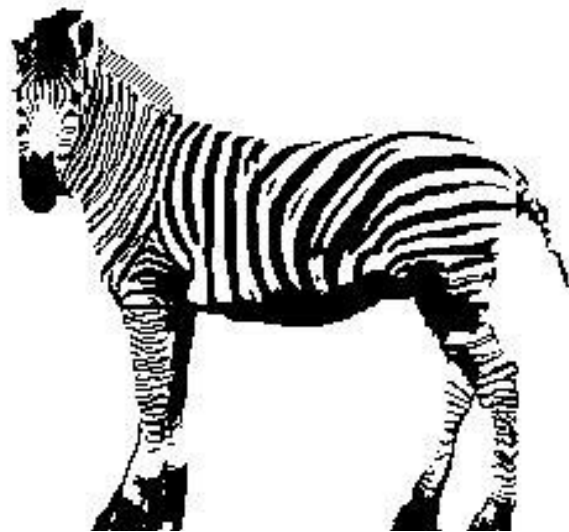
A solution: map 2D (or higher dim) field into 1D via **Hilbert-Peano curve**. Then model (HMM, PMM, TMM), estimate parameters, segment and reconstruct picture.



## Denoising a zebra with HMM and PMM (by W. Pieczynski).

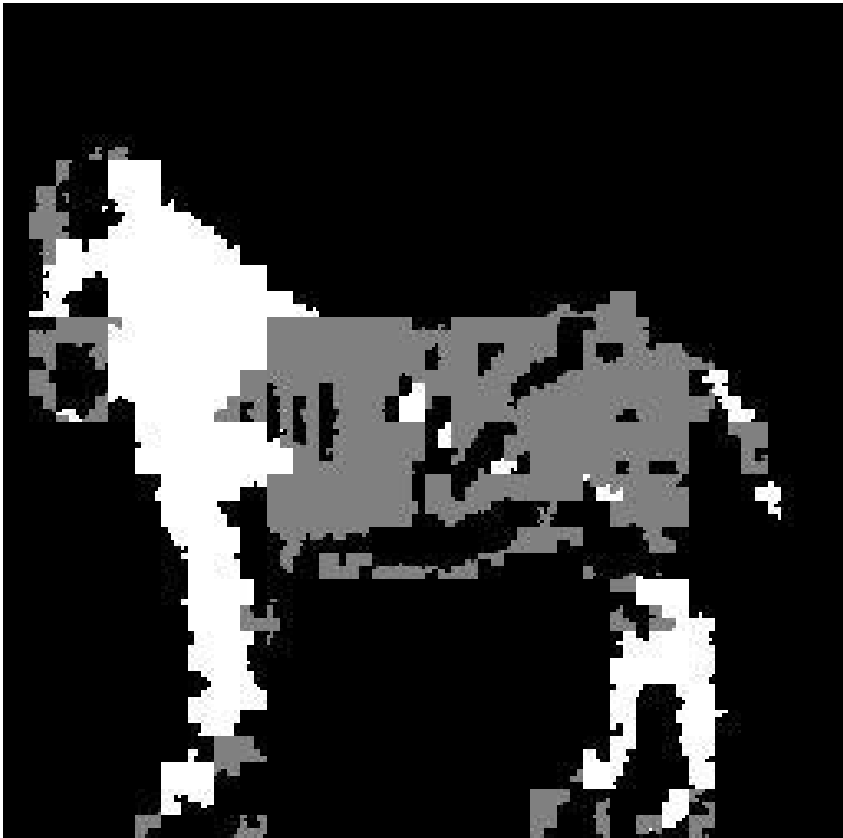
The noise is independent Gaussian, thus HMM or TMM with changing transitions (as above) are feasible. The picture is scanned with HP curve, then model parameters are estimated and PMAP segmentation is performed. After that the picture is reconstructed.





HMM reconstruction

TMM reconstruction



Estimates of U

## Some (selected) references

About (basic of) HMM:

L. Rabiner, "A Tutorial in Hidden Markov Models and Selected Applications in Speech Recognition", Proceedings of the IEEE, 1989;

T. Koski, "Hidden Markov models for bioinformatics", Kluwer, 2002;

O. Cappe, E. Moulines, T. Ryden, "Inferences with Hidden Markov models", Springer, 2005;

...

About segmentation:

J. Lember, K. Kuljus, A. Koloydenko, "Theory of segmentation", In: Hidden Markov Models: Theory and Applications (P. Dymarksi (Ed)), InTech, 2011;

J. Lember, A. Koloydenko, "Bridging Viterbi and Posterior Decoding: A Generalized Risk Approach to Hidden Path Inference Based on Hidden Markov Models", Journal of Machine Learning Research, 2014;

## About PMM and TMM:

W. Pieczynski, "Pairwise Markov chains", IEEE Trans. on Pattern Analysis and Machine Intelligence, 2003.

S. Derrode and W. Pieczynski, "Unsupervised data classification using pairwise Markov chains with automatic copulas selection", Comp. Stat. and Data Anal., 2013.

D. Benboudjema and Benboudjema, "Unsupervised Statistical Segmentation of Nonstationary Images Using Triplet Markov Fields", IEEE Transactions on Pattern analysis and Machine Inteliignece, 2007.

...