# Hidden Markov models

Wessel van Wieringen w.n.van.wieringen@vu.nl

Department of Epidemiology and Biostatistics, VUmc & Department of Mathematics, VU University Amsterdam, The Netherlands



VU medisch centrum

...

...

Consider (temporarily) a binary DNA sequence:

Question: anything striking?

...

...

Consider (temporarily) a binary DNA sequence:

Question: can the sequence be modeled by a Markov chain?

You would need two Markov chains to do this, e.g.:

one with a transition matrix like

$$\mathbf{P} = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$$

to generate something like

0101010100101010100100100010101001100

and one with a transition matrix like

$$\mathbf{P} = \begin{pmatrix} 0.0 & 1.0 \\ 0.0 & 1.0 \end{pmatrix}$$

to generate something like

Consider a charicature of the DNA that consists of introns and exons only.



Introns and exons are characterized by other nucleotide distributions.

In addition, assume that start and end base pair positions of introns and exons are unknown.

## Hidden Markov model

What does this look like?



This sequence:

- is unobserved.
- may be modelled by a 1<sup>st</sup> order Markov chain.

How to obtain the observed sequence of nucleotides?

Given the state (intron / exon):

• simply sample from (say) a multinomial.

## Hidden Markov model



## Hidden Markov model

Architecture of a hidden Markov model



#### Definition of a HMM

The processes  $\{X_t\}_{t=1,2,...}$  and  $\{Y_t\}_{t=1,2,...}$  form a HMM if: 1)  $\{X_t\}_{t=1,2,...}$  is a discrete, time-homogeneous Markov chain with state space  $S = \{E_1, ..., E_S\}$ , transition matrix  $\mathbf{P} = (p_{ij})$ and initial distribution  $\mathbf{\pi} = (\pi_1, ..., \pi_S)^T$ .

Axelson-Fisk (2010):

Section 2.1.2





#### Definition of a HMM (continued)

The processes  $\{X_t\}_{t=1,2,...}$  and  $\{Y_t\}_{t=1,2,...}$  form a HMM if: 2)  $\{Y_t\}_{t=1,2,...}$  is an observable stochastic process with state space  $V = \{V_1,...,V_m\}$ .

#### *In the intron-exon example:*

The DNA sequence is observable and  $V = \{A, C, G, T\}$ 





#### Definition of a HMM (continued)

The processes  $\{X_t\}_{t=1,2,...}$  and  $\{Y_t\}_{t=1,2,...}$  form a HMM if: 3)  $\{X_t\}_{t=1,2,...}$  and  $\{Y_t\}_{t=1,2,...}$  are related through the conditional probabilities  $P(Y_t=V_j | X_t = E_i) = b_i(V_j) = b_{ij}$ . The matrix  $\mathbf{B} = (b_{ij})$  is called the *emission matrix*.





#### Definition of a HMM (continued)

The processes  $\{X_t\}_{t=1,2,...}$  and  $\{Y_t\}_{t=1,2,...}$  form a HMM if: 4) Given the states  $\{X_t\}_{t=1,2,...}$  the observations  $\{Y_t\}_{t=1,2,...}$  are independent.



## Hidden Markov model



book Axelson-Fisk (2010): Section 2.1.2



## Hidden Markov model

A HMM may generate something like:

GTGGCACGGGTGCAGGTACGTCACCAACTCAGACTCAACG



#### **Practice**

Can the intron-exon sequence be recovered from that of the nucleotides? I.e.:

Consider the following parametrization of the intron-exon example.



$$\mathbf{P} = \mathbf{I} \begin{pmatrix} \mathbf{0}.80 & 0.20 \\ 0.10 & 0.90 \end{pmatrix} \text{ with } \boldsymbol{\pi} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \mathbf{I} \\ \mathbf{E} \end{pmatrix}$$
$$\mathbf{B} = \mathbf{I} \begin{pmatrix} \mathbf{0}.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

Question What is the probability of CGA? I.e. P(CGA)?

Consider the following parametrization of the intron-exon example.



$$\mathbf{P} = \mathbf{I} \begin{pmatrix} \mathbf{0}.80 & 0.20 \\ 0.10 & 0.90 \end{pmatrix} \text{ with } \boldsymbol{\pi} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \mathbf{I} \\ \mathbf{E} \end{pmatrix}$$
$$\mathbf{B} = \mathbf{I} \begin{pmatrix} \mathbf{0}.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

#### Question

What is the probability of **C** on the 1<sup>st</sup> position? I.e.  $P(Y_1 = C)$ ?

Consider the following parametrization of the intron-exon example.



$$\mathbf{P} = \mathbf{I} \begin{pmatrix} 0.80 & 0.20 \\ 0.10 & 0.90 \end{pmatrix} \text{ with } \boldsymbol{\pi} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \mathbf{I} \\ \mathbf{E} \\ \mathbf{B} = \mathbf{I} \begin{pmatrix} 0.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

How likely is the sequence: CGA?

1.1

We write down the likelihood of the observed sequence.

Using the conditional independence assumption:

$$P(Y_1, \dots, Y_T | X_1, \dots, X_T; (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B})) = \prod_{t=1}^T P(Y_t | X_t; (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B})) = \prod_{t=1}^T b_{X_t}(Y_t)$$



We write down the likelihood of the observed sequence. Using the Markov property:

$$P(X_1, \ldots, X_T | (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B})) = \pi_{X_1} \prod_{t=1}^{T-1} (\mathbf{P})_{X_t; X_{t+1}}$$



Joint distribution  $\{X_t\}_{t=1,..,T}$  and  $\{Y_t\}_{t=1,..,T}$  is given by:  $P(Y_1, \dots, Y_T, X_1, \dots, X_T \mid (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B}))$   $= P(Y_1, \dots, Y_T \mid X_1, \dots, X_T; (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B}))$   $\times P(X_1, \dots, X_T \mid (\boldsymbol{\pi}, \mathbf{P}, \mathbf{B}))$ 



After summing over all possible choices of  $\{X_t\}_{t=1,..,T}$ , the likelihood of  $\{Y_t\}_{t=1,..,T}$  becomes:



Hence, calculate:  $P((Y_0, Y_1, Y_2) = CGA)$ 

Using the formula on the previous slide:

 $P(Y_0 = \mathsf{C}, Y_1 = \mathsf{G}, Y_2 = \mathsf{A}) = \sum_{X_0 \in \{\mathsf{I},\mathsf{E}\}} \sum_{X_1 \in \{\mathsf{I},\mathsf{E}\}} \sum_{X_2 \in \{\mathsf{I},\mathsf{E}\}} P(X_0, X_1, X_2) \times P(Y_0 = \mathsf{C}, Y_1 = \mathsf{G}, Y_2 = \mathsf{A} \mid X_0, X_1, X_2) = \sum_{X_0 \in \{\mathsf{I},\mathsf{E}\}} \sum_{X_1 \in \{\mathsf{I},\mathsf{E}\}} \sum_{X_2 \in \{\mathsf{I},\mathsf{E}\}} P(X_2 \mid X_1) P(X_1 \mid X_0) P(X_0)$ 

 $\times P(Y_0 = C | X_0) P(Y_1 = G | X_1) P(Y_2 = A | X_2)$ 

Hence, calculate:  $P((Y_0, Y_1, Y_2) = CGA)$ 

Generate all intron-exon sequences of length 3: III, IIE, IEI, IEE, EEE, EEI, EIE, EII

For each sequence calculate its likelihood, e.g.: P(III) = 0.9 \* 0.8 \* 0.8 = 0.576 ... = ...

For each sequence calculate the likelihood of CGA, e.g.: P(CGA | III) = 0.3 \* 0.4 \* 0.1 = 0.012 = ...

#### Combine these:

P(CGA	III)	*	P(III)	=	0.006912
P(CGA	IIE)	*	P(IIE)	=	0.008640
P(CGA	IEI)	*	P(IEI)	=	0.000096
P(CGA	IEE)	*	P(IEE)	=	0.004860
P(CGA	EEE)	*	P(EEE)	=	0.001620
P(CGA	EEI)	*	P(EEI)	=	0.00036
P(CGA	EIE)	*	P(EIE)	=	0.000080
P(CGA	EII)	*	P(EII)	=	0.000064

Sum the above probabilities: P(CGA) = 0.0223

All nice these HMMs, but ...

- ... does it not overly complicate matters?
- ... can we not make do with a regular Markov model?

A HMM and 1st order Markov chain model the same data:

1<sup>st</sup> order Markov chain

 $Y_t \longrightarrow Y_{t+1}$ 



Do these models yield identical results?

The Markov chain specifies the transition probabilities  $P(Y_{t+1} | Y_t)$  for the observed sequence.

Using the HMM we need to calculate (see next slide):

$$P(Y_{t+1}|Y_t) = \sum_{X_t, X_{t+1}} P(Y_{t+1} | X_{t+1}) P(X_{t+1} | X_t) P(X_t | Y_t)$$
  
conditional independence





$$= \sum_{X_t, X_{t+1}} P(Y_{t+1}, X_t, X_{t+1} | Y_t)$$

$$\sum_{X_t, X_{t+1}} P(Y_{t+1}, X_t, X_{t+1} | Y_t) = P(Y_t | Y_t) + Y_t$$

$$= \sum_{X_t, X_{t+1}} P(Y_{t+1} \mid X_t, X_{t+1}, Y_t) P(X_t, X_{t+1} \mid Y_t)$$

$$= \sum_{X_t, X_{t+1}} P(Y_{t+1} | X_{t+1}) P(X_{t+1} | X_t, Y_t) P(X_t | Y_t)$$

$$= \sum_{X_{t}, X_{t+1}} P(Y_{t+1} | X_{t+1}) P(X_{t+1} | X_{t}) P(X_{t} | Y_{t})$$

where the definition of conditional probability and the Markov properties of the HMM have been used.



Transition probability:

$$P(Y_{t+1}|Y_t) = \sum_{X_t, X_{t+1}} \frac{P(Y_{t+1} | X_{t+1}) P(X_{t+1} | X_t) P(X_t | Y_t)}{e^{\text{mission}}}$$

$$e^{\text{mission}}_{\text{matrix of HMM}} \text{transition}_{\text{matrix of HMM}}$$

$$obtained \text{ using Bayes' rule:}$$

$$P(X_t|Y_t) = \frac{P(Y_t | X_t) P(X_t)}{\sum_{X_t \in \{I, E\}} P(Y_t | X_t) P(X_t)}$$

Recall the parametrization of the intron-exon example.



The transition probabilities between the nucleotides using the HMM are (in the form of a transition matrix **P**):

	Α	C	G	Т
Α	0.435	0.216	0.233	0.116
C	0.340	0.240	0.280	0.140
G	0.320	0.245	0.290	0.145
т	0.320	0.245	0.290	0.145

Recall the parametrization of the intron-exon example.



The stationary distribution of this transition matrix  $P^n = 0.367 \ 0.233 \ 0.267 \ 0.133$ 

This is the same for the HMM, when we use:

 $\sum_{X^t \in \{\mathbf{I}, \mathbf{E}\}} P(X_t) \, b_{X_t}(Y_t) = \psi_{\mathbf{I}} \, b_{\mathbf{I}}(Y_t) + \psi_{\mathbf{E}} \, b_{\mathbf{E}}(Y_t)$ 

Recall the parametrization of the intron-exon example.



The initial distribution is sobtained from:

$$P(Y_1) = \sum_{X_1 \in \{I, E\}} P(Y_1 | X_1) P(X_1)$$

which yields:

$$\begin{array}{cccc} A & C & G & T \\ \pi &= 0.14 & 0.29 & 0.38 & 0.19 \end{array}$$

Recall the parametrization of the intron-exon example.



The couple  $(\pi, P)$  specifies a Markov chain. However, this Markov chain and the HMM from which it has been derived do not (necessarily) yield the same likelihood:

$$P_{HMM}(CGA) = 0.02232$$
  
vs.  
 $P_{Markov chain}(CGA) = 0.025984$
## Parsimony

#### 1<sup>st</sup> order Markov chain

The transition matrix **P** has 4x4 - 4 free parameters, while for its initial distribution has 4 - 1 parameters. In total: 15 parameters.

#### HMM

The HMM has 2x2-2 (matrix **P**), 2x4-2 (matrix **B**) and 2-1 (initial distribution) parameters. In total: 9.

The HMM is to be preferred from a parsimony perspective.

#### Question

Which model does the parsimony comparison favor when the HMM has a larger latent state space? E.g.:



Could you name other reasons why one model is then preferred?

Recall the parametrization of the intron-exon example.



Nonetheless, the HMM *can* be written as a Markov chain with the following state space:

 $S = \{IA, IC, IG, IT, EA, EC, EG, ET\}$   $intron \& \dots exon \&$ nucleotide A nucleotide T

Recall the parametrization of the intron-exon example.



This state space:

 $S = \{IA, IC, IG, IT, EA, EC, EG, ET\}$ 

has eight states. Consequently, the corresponding 1<sup>st</sup> Markov chain has in total: 63 parameters.

In comparison, the HMM has 9 parameters (in total).

*Question* Is this a fair comparison?

Recall the parametrization of the intron-exon example.



The transition matrix of this Markov chain contains, e.g.:

$$P(Y_{t+1} = IA | Y_t = EG)$$

To analyze a DNA sequence with Markov chain however requires that *both* the nucleotide and intron/exon sequence have been observed.

Canonical HMM problems

## Three basic problems for HMMs

- 1) For given parameters  $\lambda = (\pi, P, B)$ , how do we calculate  $P(\{Y_t\}_{t=1,2,...,T} | \lambda)$ , the probability of an observation, efficiently?
- 2) For given  $\{Y_t\}_{t=1,2,...,T}$  find the sequence  $\{X_t\}_{t=1,2,...,T}$  that maximizes P ( $\{X_t\}_{t=1,2,...,T} | \{Y_t\}_{t=1,2,...,T}$ ,  $\lambda$ ). Hence, which underlying state path is most probable?
- 3) What is the  $\lambda = (\pi, P, B)$  that maximizes the probability of an observed sequence  $\{Y_t\}_{t=1,2,...,T}$ ? How to maximize the likelihood?

Solutions for the basic problems for HMMs

- 1) The forward algorithm.
- 2) The Viterbi algorithm.

3) The Baum-Welch algorithm *(not discussed)*. Numbers corresponding to previous slide.

The first two problems could have been solved by direct calculation. This becomes rather cumbersome as T, N and M grow.

# The forward algorithm

Consider the following parametrization of the intron-exon example.



$$\mathbf{B} = \mathbf{I} \begin{pmatrix} \mathbf{A} & \mathbf{C} & \mathbf{G} & \mathbf{T} \\ 0.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

How likely is the sequence: CGA?

Hence, calculate:  $P((Y_0, Y_1, Y_2) = CGA)$ 

- Generate all intron-exon sequences of length 3. III, IIE, IEI, IEE, EEE, EEI, EIE, EII
- For each sequence calculate its likelihood, e.g.:
   P(III) = 0.9 \* 0.8 \* 0.8 = 0.576
- For each sequence calculate the likelihood of CGA, e.g.:
   P(CGA | III) = 0.3 \* 0.4 \* 0.1 = 0.012
- Combine these:

P(CGA | III) \* P(III) = 0.006912

• Sum the above probabilities:

P(CGA) = 0.0223

Impractical for longer sequences: forward algorithm!





### *Idea behind forward algorithm* Likelihood of CGA:



equals (recall: total probability law)



These two quantities are "easy" to calculate.

Here we illustrate the calculation of  $P(Y_1, ..., Y_T | \lambda)$  using the forward algorithm.

First define the forward variable:

 $\alpha_t(i) = P(Y_1, \dots, Y_t, X_t = S_i | \boldsymbol{\lambda})$ 



hidden state path

Axelson-Fisk (2010):

Section 2.1.4

book

observed sequenc

This forward variable is of interest for:

$$\sum_{i=1}^{N} \alpha_{T}(i) = \sum_{i=1}^{N} P(Y_{1}, \dots, Y_{T}, X_{T} = S_{i} | \boldsymbol{\lambda})$$
$$= P(Y_{1}, \dots, Y_{T} | \boldsymbol{\lambda})$$

This uses total probability law, summing over all possible states of the unobserved hidden variable at a particular time.

In order to compute  $P(Y_1, ..., Y_T | \lambda)$  efficiently, we may derive (see SM) a recursive relationship for the forward variable:

$$\alpha_t(i) = b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji}$$

The recursive relationship for the forward variable (graphically)

AT.

$$\alpha_t(i) = b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji}$$

The sum over j washes out  $X_{t-1}$ .





# The forward algorithm

Apply the recursive relationship:  $\alpha_t(i) = b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji}$ 

to the example:

 $\frac{P((Y_1, Y_2, Y_3) = (C, G, A), X_3 = E)}{\alpha_3(E)} =$ 



The recursive relationship expresses the forward variable at time *t* in terms of those at time *t*-1.

book

Axelson-Fisk (2010):

Section 2.1.3

Walking forward in time, to calculate the probability of interest we only store the last probabilities. Dynamic programming!



The forward algorithm comprises of:

1) Initialization:

$$\alpha_1(i) = \pi_i \, b_i(Y_1)$$

2) Induction:  $\alpha_t(i) = b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji}$ for  $2 \le t \le T, 1 \le i \le N$ 

**A Y** 

Axelson-Fisk (2010):

Section 2.1.4

book

*3) Termination*:

$$P(Y_1, \dots, Y_T | \boldsymbol{\lambda}) = \sum_{i=1}^N \alpha_T(i)$$

The forward algorithm comprises of:

1) Initialization:

$$\alpha_1(i) = \pi_i b_i(Y_1)$$

Reconsider the intron-exon example. How likely is the sequence: CGA?

Initialization:  $\alpha_1(I) = 0.9 * 0.3 = 0.27$  $\alpha_1(E) = 0.1 * 0.2 = 0.02$  The forward algorithm comprises of:

2) Induction:  $\alpha_t(i) = b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji}$ for  $2 \le t \le T, 1 \le i \le N$ Induction:  $\alpha_2(\mathbf{I}) = 0.4 * (0.27 * 0.8 + 0.02 * 0.1)$  = 0.0872

 $\alpha_2(E) = 0.2 * (0.27 * 0.2 + 0.02 * 0.9)$ = 0.0144

$$\alpha_3(I) = 0.1 * (0.0872 * 0.8 + 0.0144 * 0.1)$$
  
= 0.00712

$$\alpha_3(E) = 0.5 * (0.0872 * 0.2 + 0.0144 * 0.9)$$
  
= 0.0152

# The forward algorithm

The forward algorithm comprises of:

3) Termination:

$$P(Y_1, \dots, Y_T | \boldsymbol{\lambda}) = \sum_{i=1}^{N} \alpha_T(i)$$

N

Reconsider the intron-exon example.



Termination:

 $\alpha_3(I) + \alpha_3(E) = 0.00712 + 0.0152$ = 0.0223

As before:

P(CGA) = 0.0223

т

Consider the following parametrization of the intron-exon sample.



$$\mathbf{P} = \mathbf{I} \begin{pmatrix} 0.80 & 0.20 \\ 0.10 & 0.90 \end{pmatrix} \text{ with } \boldsymbol{\pi} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \mathbf{I} \\ \mathbf{E} \\ \mathbf{B} = \mathbf{I} \begin{pmatrix} 0.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

What is the most likely intron-exon sequence to have generated the nucleotide sequence: CGA?

Г

What is the most likely intron-exon sequence to have generated the nucleotide sequence: CGA?

Question

How would you proceed to answer the above?

→ Would you base your decision on (e.g.):
P(IIE | CGA), P(CGA | IIE), P(IIE), Or P(CGA)?

 $\rightarrow$  How do you distinguish between (say) **IIE** and **EIE**?

Hence, we need to calculate:

$$\max_{X_0, X_1, X_2} P((X_0, X_1, X_2) | (Y_0, Y_1, Y_2))$$

$$= \max_{X_0, X_1, X_2} P((Y_0, Y_1, Y_2) | (X_0, X_1, X_2))$$

$$\times P((X_0, X_1, X_2)) / P((Y_0, Y_1, Y_2))$$

Bayes' rule

Intron/exon example, for any choice of  $(X_0, X_1, X_2)$  evaluate:  $P((X_0, X_1, X_2) \mid (Y_0, Y_1, Y_2) = CGA)$  Let us calculate, e.g:  $P((X_0, X_1, X_2) = III | (Y_0, Y_1, Y_2) = CGA)$ 

The likelihood of the latent intron-exon sequence is: P(III) = 0.9 \* 0.8 \* 0.8 = 0.576

Given the intron-exon sequence calculate the likelihood of the observed nucleotide sequence CGA:

P(CGA | III) = 0.3 \* 0.4 \* 0.1 = 0.012

Also, calculate the likelihood of the observed nucleotide sequence CGA using the forward algorithm:

P(CGA) = 0.0223

Calculate	e all conc	ditio	onal probabilities:			
P(III	CGA)	=	P(CGA   III) * P(	II	[) /	P(CGA)
		=	0.006912 / 0.0223	=	0.31	LO
P(IIE	CGA)	=	• • •	=	0.38	37
P(IEI	CGA)	=	• • •	=	0.00	)4
P(IEE	CGA)	=	• • •	=	0.21	L8
P(EEE	CGA)	=	• • •	=	0.07	73
P(EEI	CGA)	=	• • •	=	0.00	)2
P(EIE	CGA)	=	• • •	=	0.00	)4
P(EII	CGA)	=	• • •	=	0.00	)3

Select the one with the highest probability:

**Q**: Necessary to calculate the likelihood?

The preceeding calculation is impractical for long sequences.

The Viterbi algorithm efficiently calculates:  $P(X_1, ..., X_T | Y_1, ..., Y_T, \mathbf{\lambda}).$ 

First, define the variable:

$$\delta_t(i) = \max_{X_1, \dots, X_{t-1}} P(X_1, \dots, X_{t-1}, X_t = S_i, Y_1, \dots, Y_t | \boldsymbol{\lambda})$$

- The maximum is over all paths of length *t* ending in state  $S_i$ .
- The likelihood of this path together with the *t* associated observations is maximized.

Note:

$$P(X_1, ..., X_T | Y_1, ..., Y_T, λ) = \delta_t(i) / P(Y_1, ..., Y_T, λ)$$
likelihood



The recursive relationship:



As with the forward variable, the recursive relationship expresses the 'Viterbi' variable at time *t* in terms of those at time *t*-1. Hence, we do not need to keep track of probabilities at all time points.





```
\log[\delta_1(E)] = -3.912023
```

The log-scale is used for efficiency and accuracy.

The Viterbi algorithm comprises of:

2) Induction:

$$\delta_t(j) = \max_{1 \le i \le N} \delta_{t-1}(i) a_{ij} b_j(Y_t) \qquad 1 \le j \le N$$
  
$$\psi_t(j) = \arg \max_{1 \le i \le N} \delta_{t-1}(i) a_{ij} \qquad 2 \le t \le T$$

 $\psi_t(j)$  keeps track of which argument maximized  $\delta_t(j)$ 





- $= -0.9162907 + \max \{-1.532477, -6.214608\}$
- **= -2.448768**

Similarly:  $\log[\delta_2(E)] = -4.528209$   $\log[\delta_3(I)] = -4.974497$  $\log[\delta_3(E)] = -4.751353$  Parallel we get:  $\Psi_1(I) = 0, \Psi_1(E) = 0$  $\Psi_2(I) = \arg \max_{X_1} \{ \delta_1(X_1=I) * P(X_2=I | X_1=I), \}$  $\delta_1(\mathbf{X}_1 = \mathbf{E}) * \mathbf{P}(\mathbf{X}_2 = \mathbf{I} \mid \mathbf{X}_1 = \mathbf{E}) \}$  $= \arg \max_{X_1} \{ \exp(-1.309333) * 0.80,$  $\exp(-3.912023) * 0.10$ = I  $Ψ_2(E) = \arg \max_{X_1} \{ \delta_1(X_1 = I) * P(X_2 = E | X_1 = I), \}$  $\delta_1(\mathbf{X}_1 = \mathbf{E}) * \mathbf{P}(\mathbf{X}_2 = \mathbf{E} \mid \mathbf{X}_1 = \mathbf{E}) \}$ = arg max {exp(-1.309333) \* 0.20,  $X_1$  $\exp(-3.912023) * 0.90$ = I

= I

= I

Parallel we get:

$$\Psi_{3}(\mathbf{I}) = \arg \max_{X_{2}} \left\{ \delta_{2}(X_{1}=\mathbf{I}) * P(X_{3}=\mathbf{I} \mid X_{2}=\mathbf{I}), \\ \delta_{2}(X_{1}=\mathbf{E}) * P(X_{3}=\mathbf{I} \mid X_{2}=\mathbf{E}) \right\}$$
  
= arg max {exp(-2.448768) \* 0.80,  
exp(-4.528209) \* 0.10 }

$$\Psi_{3}(\mathbf{E}) = \arg \max_{X_{2}} \{ \delta_{2}(X_{2}=\mathbf{I}) * P(X_{3}=\mathbf{E} \mid X_{2}=\mathbf{I}), \\ \delta_{2}(X_{2}=\mathbf{E}) * P(X_{3}=\mathbf{E} \mid X_{2}=\mathbf{E}) \}$$

Similarly:

 $\Psi_4(I) = I$  $\Psi_4(E) = E$ 

Results from induction step:



	$\psi_2$	$\psi_3$	$\psi_4$
Intron	I	I	I
Exon	I	I	Е
The Viterbi algorithm finalized with:

*3) Termination*:

$$p^* = \max_{1 \le i \le N} \delta_T(i)$$
$$q_T^* = \arg \max_{1 \le i \le N} \delta_T(i)$$

This yields:



Thus: 
$$q_3^* = E$$

The Viterbi algorithm finalized with: *4) Backtracking*:

$$q_{t-1}^* = \psi_t(q_t^*) \qquad t = T, T-1, \dots, 2$$

Backtracking:

an underlying sequence with an **exon** at *t*=3 is most likely.



Continue backtracking:

$$Y_1 = C$$
 $Y_2 = G$  $Y_3 = A$ Intron-1.309333-2.448768-4.974497Exon-3.912023-4.528209-4.751353

Thus:  $\psi_3(\mathbf{E}) = \mathbf{I}$  and  $\psi_2(\mathbf{I}) = \mathbf{I}$ 



The most likely underlying sequence is: **IIE** 

## The forward algorithm

#### Illustration

Sample from the following parametrization of the intron-exon example.



$$\mathbf{P} = \mathbf{I} \begin{pmatrix} \mathbf{I} & \mathbf{E} \\ 0.80 & 0.20 \\ 0.10 & 0.90 \end{pmatrix} \text{ with } \boldsymbol{\pi} = \begin{pmatrix} 0.9 \\ 0.1 \end{pmatrix} \mathbf{I} \\ \mathbf{E} \\ \mathbf{B} = \mathbf{I} \begin{pmatrix} 0.10 & 0.30 & 0.40 & 0.20 \\ 0.50 & 0.20 & 0.20 & 0.10 \end{pmatrix}$$

Data:

T C G C G C T G T T T G T C C T A A G T G T A T A C A

Illustration Most likely underlying state sequence?  $\rightarrow$  Viterbi!

True states (1<sup>st</sup> 100 positions)

Inferred states (1<sup>st</sup> 100 positions)

True vs. inferred (10000 positions):

	exon	intron	(true)
exon	5655	1633	
intron	776	1936	

Illustration Related is the posterior probability of states:  $P(X_{t} | Y_{1}, ..., Y_{T}, \boldsymbol{\lambda}) = P(X_{t}, Y_{1}, ..., Y_{T}, \boldsymbol{\lambda}) / P(Y_{1}, ..., Y_{T}, \boldsymbol{\lambda}),$ likelihood 400 forward variable 300 Frequency 200 Distribution of  $P(E|Y_1, ..., Y_T, \lambda)$ for intron bases 100 0 0.2 0.0 0.4 0.6 0.8

probability



```
# load library
library(HMM)
```

```
# specify HMM
```

```
hmm <- initHMM(c("I","E"),c("A","C","G", "T"),
    transProbs=matrix(c(0.8, 0.2, 0.1, 0.9), 2, 2, byrow=TRUE),
    emissionProbs=matrix(c(0.1, 0.3, 0.4, 0.2, 0.5, 0.2, 0.2, 0.1),
    2, 4, byrow=TRUE))</pre>
```

```
# simulate from the HMM
seqLength <- 10000
DNAseq <- simHMM(hmm, seqLength)</pre>
```

# apply viterbi and compare to true underlying sequence table(viterbi(hmm, DNAseq\$observation), DNAseq\$states)

# obtain posterior probabilities of latent states (given the sequence)
posterior(hmm, DNAseq\$observation)

# histogram of post. probabilities for exons. hist(posterior(hmm, obsSequence\$observation)[1,obsSequence\$states=="E"]) Note

Apart from the backtracking step, the Viterbi algorithm is similar to the forward algorithm. The key difference is the maximization in the former instead of summation in the latter.

## The Baum-Welch algorithm (maximum likelihood) (sketch only)

## The Baum-Welch algorithm

Intuitively, one may combine:

- $\rightarrow$  efficient likelihood evaluation,
- $\rightarrow$  (probabilistic) recovery of underlying state, and
- $\rightarrow$  ML estimation of Markov chain (lecture 1),
- to conceive an algorithm along the following lines:
- 0) Initiate parameters **P**, **B** and  $\pi$ .
- 1) Apply Viterbi: most likely underlying sequence.
- 2) Given sequence, update estimates of **P**, **B** and  $\pi$ .
- 3) Iterate between 1) and 2) until convergence.

*Baum-Welch* does something similar, but formalized to ensure convergence to the maximum likelihood estimator.



#### Consider two DNA sequences of two mammals: human: CTATACG mouse: CGATCG

In order to reconstruct the sequence in their common ancestor, the two sequences are aligned.

An alignment can be thought of as the result of sequential executions of the following operations:

- match,
- substitution,
- insertion, and
- deletion.

Hidden Markov models aid in alignment.

Formulate the alignment process in terms of a HMM. The operations *match*, *deletion*, and *insertion* are hidden states.

Suppose the two sequences match perfectly: human: CTATACG mouse: CTATACG (really: CGATCG)

Only one hidden state is needed:



As the two sequences do not match perfectly, the HMM is extended with the insertion state:

human: CTATACG
mouse: CGATCG -> CGTATACG



Now add the deletion state:

human : CTATACG

mouse





d

m

Remarks:

 Deletion state does not emit: affects the dimension of the emission matrix.



2) Consequently, the underlying sequence may be longer than the observed.



3) First hidden state is defined as observed start state.



DNA copy number of a genomic segment is simply the number of copies of that segment present in the cell under study.

Healthy normal cell: chr 1 : 2 ...

```
chr 22 : 2
chr X : 1 or 2
chr Y : 0 or 1
```



#### Chromosomes of a tumor cell



Technique: SKY

DNA copy number may be measured genome-wide in highthroughput fashion using microarrays.



- Test and reference (assumed to have DNA copy number 2) samples are labeled differently and hybridized together.
- Under ideal circumstances, the intensity of an array element is linearly proportional to the abundance of the corresponding DNA sequence.
- Log<sub>2</sub> ratios of test and reference intensities reflect the relative copy number in the test sample compared to that in the reference sample.

Principle of array CGH.



Copy numbers of each genomic segment are assessed simultaneously. Down-syndrome: extra copy of chr. 21

22

11

12

13 14

15

16

17

18

19

20

#### DNA copy number profile of a cancer sample



*Problem*: determine the aberrated genomic segments from the observed  $log_2$  ratios.

Solution: model the DNA copy number profile by a HMM.

*Ingredients for the HMM*: three hidden states:

- L : *loss* : < 2 copies
- N : normal : 2 copies
- G : *gain* : > 2 copies

Given the hidden states, a *continuous* value is emitted.

HMM architecture for array CGH



Sampling from this HMM:

	t =	1 : normal			
	t =	2/: gain			
	t=	3 : normal			
	t =	4 : normal			
	t =	5 : normal			
	t =	6 : normal			
	t =	7: gain			
	t =	8 : gain			
	t =	9 : normal			
	t =	10 : normal			
		1			
unobserved					
sequence					
	of states				



#### DNA copy number profiles of two cervix cancer samples.





#### Analysis flow

#### • Fit HMM to data by means of ML (Baum-Welch).

```
> summary(hmmFit)
$delta
[1] 1.000000e+00 1.801274e-45 2.257732e-58
```

\$Pi

	[,1]	[,2]	[,3]
[1,]	0.979794054	0.01695689	0.003249055
[2,]	0.011131489	0.98481411	0.004054398
[3,]	0.005139211	0.01587148	0.978989311

\$pm
\$pm\$mean
[1] -0.32802689 0.04630948 0.53721098
\$pm\$sd
\$pm\$sd

[1] 0.1850411 0.1335572 0.3447424

• Given ML estimates, calculate hidden states (Viterbi)

#### HMM result with three states.

profile 2 (HMM, 3 states)

profile 3 (HMM, 3 states)



#### HMM result with five states.

profile 2 (HMM, 5 states)

profile 3 (HMM, 5 states)



# R

#### In R:

- > # activate libraries
- > library(CGHcall)
- > library(HiddenMarkov)
- > # load data
- > data(Wilting)

```
> # specify profile to analyze
> profNo <- 3</pre>
```

> # put in particular format
> cgh <- make\_cghRaw(Wilting)</pre>

```
> # first preprocess the data
> rawCGH <- preprocess(cgh)</pre>
```

> # normalize the data (global median normalization)
> normCGH <- normalize(rawCGH)
> plot(copynumber(normCGH)[,profNo], pch=".", cex=2)



#### R code (continued)

- > # fit HMM
- > hmmFit <- BaumWelch(x)</pre>

> # estimate hidden states
> states <- Viterbi(hmmFit)</pre>

- > # overlay data with estimated hidden states
- > stateValues <- hmmFit\$pm\$mean</pre>
- > plot(copynumber(normCGH)[, profNo], pch=".", cex=2)
- > lines(stateValues[states], col="red", lwd=2)

#### Note

The HMM is fitted to data from the whole genome. Effectively, this assumes the chromosomes are glued together (ordered by their numbers). This is of course nonsense.

#### Exercise

- Modify the code to fit the HMM to each chromosome separate, and overlay the data with the resulting fit.
- Moreover, also investigate the number of states to be used.

## Hidden semi-Markov model

#### Hidden semi-Markov model

The HMM analysis suggest that within only a few steps multiple hidden states have been visited. This may be unrealistic biologically.



#### Hidden semi-Markov model

Consider the intron-exon charicature of the DNA again.



How long do we spend in a particular state:

P(d time steps time spent in state k)

= 
$$P(X_t \neq k, X_{t+1} = k, ..., X_{t+d} = k, X_{t+d+1} \neq k)$$
  
=  $P(X_t \neq k) (\mathbf{P})_{\neg k, k} (\mathbf{P})_{kk}^{d-1} (\mathbf{P})_{k, \neg k}$ 

Short durations more likely than long durations.



## Hidden semi-Markov model

Modify the HMM to allow more realistic duration times.



where  $P_{I}(d)$  and  $P_{E}(d)$  are, e.g. Poisson distributions describing the duration time.
## Hidden semi-Markov model

**Poisson distribution** 



Index

# Hidden semi-Markov model

Sampling from this modified HMM:



# Hidden semi-Markov model

The modified HMM is called a *hidden semi-Markov model*.

Within this model the Markov property holds for state transitions, but no longer for sequential observations.



*Note:* with two states, the unobserved sequence of states is deterministic, alternating between the two states.

Recall: array CGH measures DNA copy number in high - throughput fashion. Below two cervix cancer profiles.



#### HSMM result with three states.

profile 2 (HSMM, 3 states)

profile 3 (HSMM, 3 states)



genomic order

genomic order

#### HSMM result with five states.

profile 2 (HSMM, 5 states)

profile 3 (HSMM, 5 states)



genomic order

genomic order

The HSMM analysis yields longer duration times than the HMM for the DNA copy number profile of the cervix cancer samples.



Supplementary Material: Proofs of recursive relations of the forward and Viterbi variable



# The forward algorithm

$$\begin{aligned} \alpha_t(i) &= P(Y_1, \dots, Y_t, X_t = S_i | \boldsymbol{\lambda}) \\ &= \sum_{j=1}^N P(Y_1, \dots, Y_t, X_t = S_i, X_{t-1} = S_j | \boldsymbol{\lambda}) \\ &= \sum_{j=1}^N P(Y_t, X_t = S_i | Y_1, \dots, Y_{t-1}, X_{t-1} = S_j; \boldsymbol{\lambda}) \\ &\times P(Y_1, \dots, Y_{t-1}, X_{t-1} = S_j | \boldsymbol{\lambda}) \\ &= \sum_{j=1}^N P(Y_t | X_t = S_i; \boldsymbol{\lambda}) P(X_t = S_i | X_{t-1} = S_j; \boldsymbol{\lambda}) \\ &\times P(Y_1, \dots, Y_{t-1}, X_{t-1} = S_j | \boldsymbol{\lambda}) \\ &= b_i(Y_t) \sum_{j=1}^N \alpha_{t-1}(j) (\mathbf{P})_{ji} \end{aligned}$$

For this variable too there exists a recursive relationship:

 $\delta_{t+1}(j) = \max_{X_1, \dots, X_t} P(X_1, \dots, X_t, X_{t+1} = S_j, Y_1, \dots, Y_{t+1} | \boldsymbol{\lambda})$  $= |\max_{X_1,\ldots,X_t} P(X_1,\ldots,X_t,Y_1,\ldots,Y_t | \boldsymbol{\lambda})|$  $P(X_{t+1} = S_i | X_t, \boldsymbol{\lambda})]$  $\times P(Y_{t+1} \mid X_{t+1} = S_i, \boldsymbol{\lambda})$  $= \max_{i} \max_{X_1,\ldots,X_{t-1}} P(X_1,\ldots,X_t = S_i,Y_1,\ldots,Y_t \mid \boldsymbol{\lambda})$  $P(X_{t+1} = S_i \mid X_t = S_i, \boldsymbol{\lambda})]$  $\times P(Y_{t+1} \mid X_{t+1} = S_i, \boldsymbol{\lambda})$  $= \left[\max_{i} \delta_t(i) a_{ij}\right] b_j(Y_{t+1})$ 

References & further reading

# References and further reading

- Durbin, R., Eddy, S., Krogh, A., Mitchison, G. (1998), *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*, Cambridge University Press.
- Ewens, W.J, Grant, G. (2006), *Statistical Methods for Bioinformatics*, Springer, New York.
- Fridlyand, J., Snijders, A.M., Pinkel, D., Albertson, D.G., Jain, A.N. (2004), "Hidden Markov models approach to the analysis of CGH data", *Journal of Multivariate Analysis*, 90(1), 132-153.
- Rabiner, L.R. (1989), "A tutorial on hidden Markov models and selected applications in speech recognition", *Proceedings of the IEEE*, 77(2), 257-286.
- Van de Wiel, M.A., Picard, F., Van Wieringen, W.N., Ylstra, B. (2011), "Preprocessing and downstream analysis of microarray DNA copy number profiles", *Briefings in Bioinformatics*, 12(1), 10-21.
- Yu, S.-Z. (2010), "Hidden semi-Markov models", Artificial Intelligence, 174, 215-243.



This material is provided under the Creative Commons Attribution/Share-Alike/Non-Commercial License.

See http://www.creativecommons.org for details.