



Machine learning: lecture 19

Tommi S. Jaakkola

MIT CSAIL

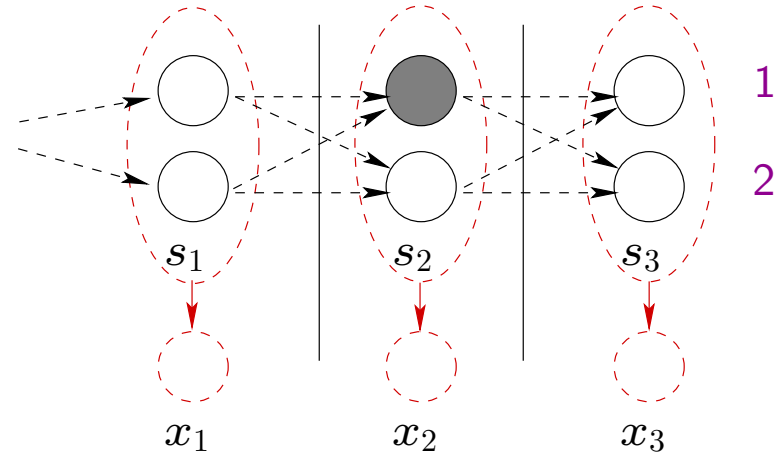
tommi@csail.mit.edu



Topics

- Hidden Markov models (HMMs)
 - the EM-algorithm
 - example
 - dynamic programming

Forward-backward probabilities: review



- Forward (predictive) probabilities $\alpha_t(i)$:

$$\alpha_t(i) = P(\mathbf{x}_1, \dots, \mathbf{x}_t, s_t = i)$$

- Backward (diagnostic) probabilities $\beta_t(i)$:

$$\beta_t(i) = P(\mathbf{x}_{t+1}, \dots, \mathbf{x}_n | s_t = i)$$

(evidence about the current state from future observations)

Uses of forward/backward probabilities

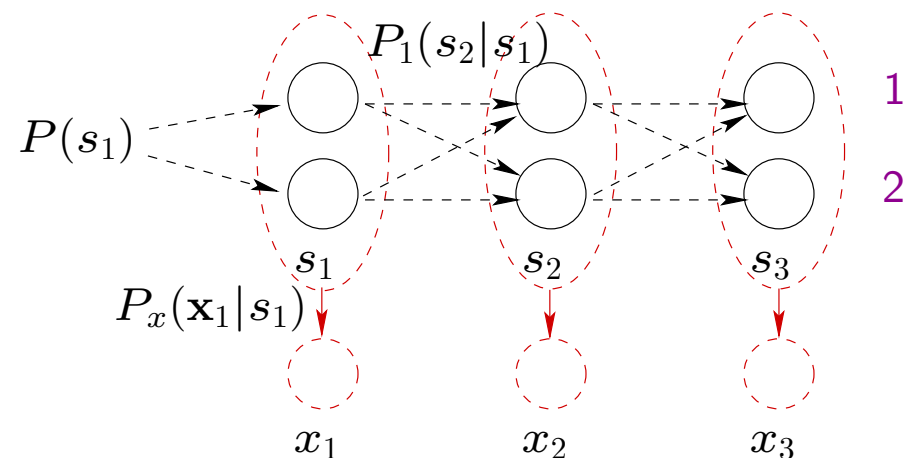
- The complementary forward/backward probabilities

$$\alpha_t(i) = P(\mathbf{x}_1, \dots, \mathbf{x}_t, s_t = i)$$

$$\beta_t(i) = P(\mathbf{x}_{t+1}, \dots, \mathbf{x}_n | s_t = i)$$

permit us to evaluate various probabilities:

- $P(\mathbf{x}_1, \dots, \mathbf{x}_n)$
- $\gamma_t(i) = P(s_t = i | \mathbf{x}_1, \dots, \mathbf{x}_n)$
- $\xi_t(i, j) = P(s_t = i, s_{t+1} = j | \mathbf{x}_1, \dots, \mathbf{x}_n)$



The EM algorithm for HMMs

Assume we have L observation sequences $\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_{n_l}^{(l)}$

E-step: compute the posterior probabilities

$$\gamma_t^{(l)}(i) \quad \text{for all } l, i, \text{ and } t \ (t = 1, \dots, n_l)$$

$$\xi_t^{(l)}(i, j) \quad \text{for all } l, i, \text{ and } t \ (t = 1, \dots, n_l - 1)$$

M-step: First, the initial state distribution can be updated according to the expected fraction of times the sequences started from a specific state i

$$\hat{P}(i) \leftarrow \frac{1}{L} \sum_{l=1}^L \gamma_1^{(l)}(i)$$

M-step cont'd

Second, the transition probabilities can be updated on the basis of the posterior counts:

$$\hat{P}_1(j|i) \leftarrow \frac{\hat{n}(i, j)}{\sum_{j'} \hat{n}(i, j')}$$

where

$$\hat{n}(i, j) = \sum_{l=1}^L \sum_{t=1}^{n-1} \xi_t^{(l)}(i, j)$$

defines the expected number of transitions from i to j

M-step cont'd

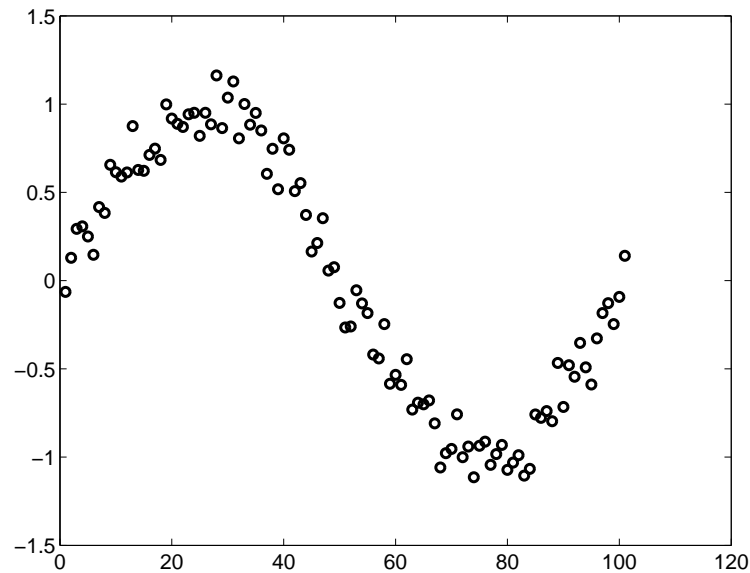
- Lastly, for the outputs we have to (in general) solve a weighted maximum likelihood estimation problem:

Separately for each state i we maximize:

$$J(\theta_i) = \sum_{l=1}^L \sum_{t=1}^{n_l} \gamma_t^{(l)}(i) \log P(\mathbf{x}_t^{(l)} | \theta_i)$$

with respect to the parameters θ_i (e.g, the mean and the covariance of a Gaussian).

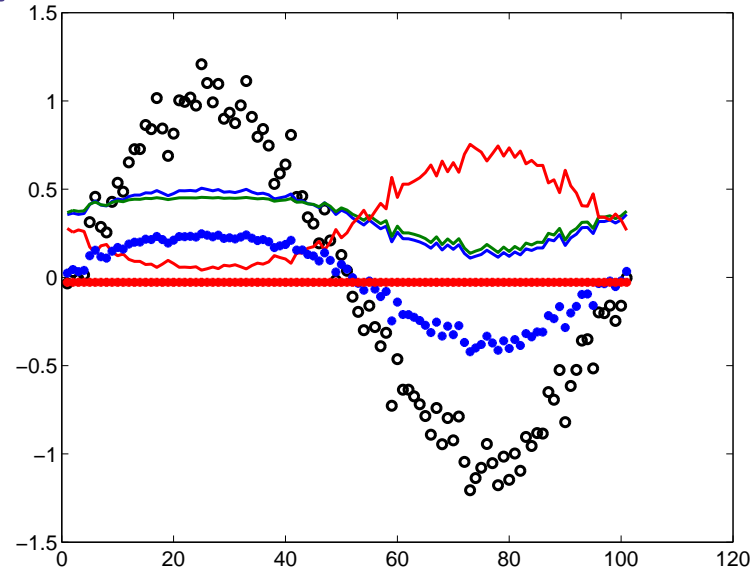
HMM example



Observed output as a function of time

- We will try to model this with a 3-state HMM with Gaussian outputs $p(x|s = i) = p(x|\mu_i, \sigma_i^2)$, $i = 1, 2, 3$.

HMM example cont'd



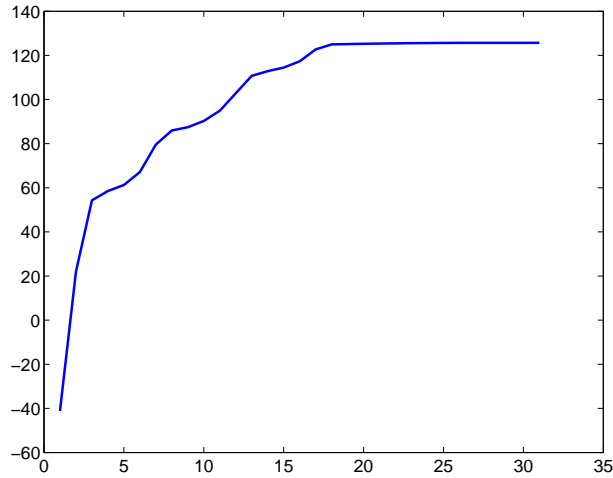
prior/posterior means and $\gamma_t(\cdot)$

$$\text{prior mean}(t) = \sum_i P_t(i) \hat{\mu}_i \quad ('*')$$

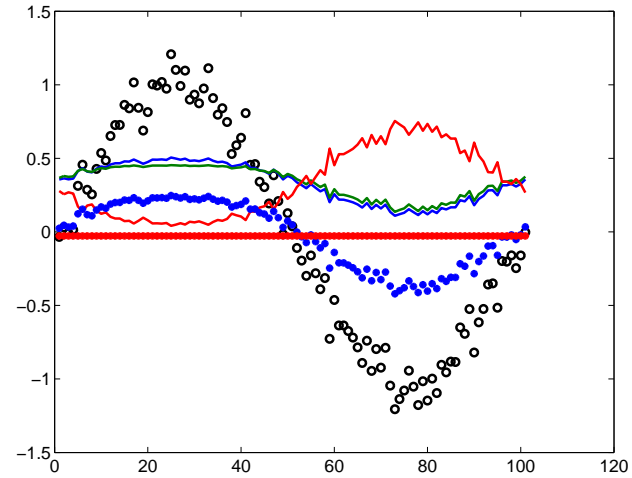
$$\text{posterior mean}(t) = \sum_i \gamma_t(i) \hat{\mu}_i \quad ('*')$$

where $P_t(i)$ is the probability of being in state i after t steps without observations; $\hat{\mu}_i$ is the mean output from the i^{th} state

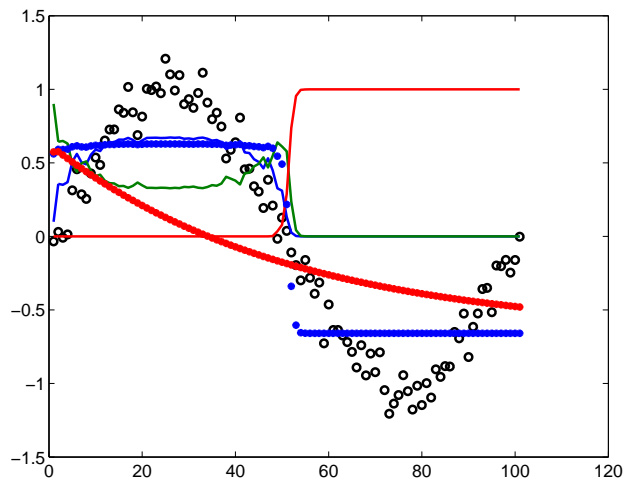
HMM example cont'd



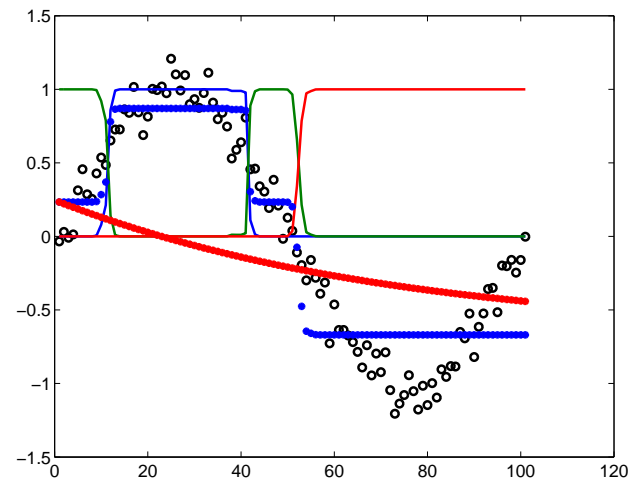
Log-prob. of data



after 0 iterations



after 7 iterations

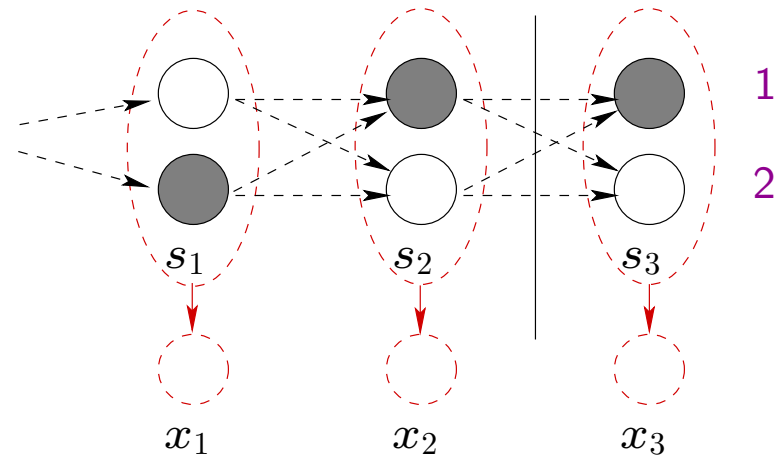


final

HMM problems

- There are several problems we have to solve
 1. How do we evaluate the probability of an observation sequence $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$?
 - forward-backward algorithm
 2. How do we adapt the parameters of the HMM to better account for the observations?
 - the EM-algorithm
 3. How do we uncover the most likely hidden state sequence corresponding to the observations?
 - dynamic programming (Viterbi algorithm)

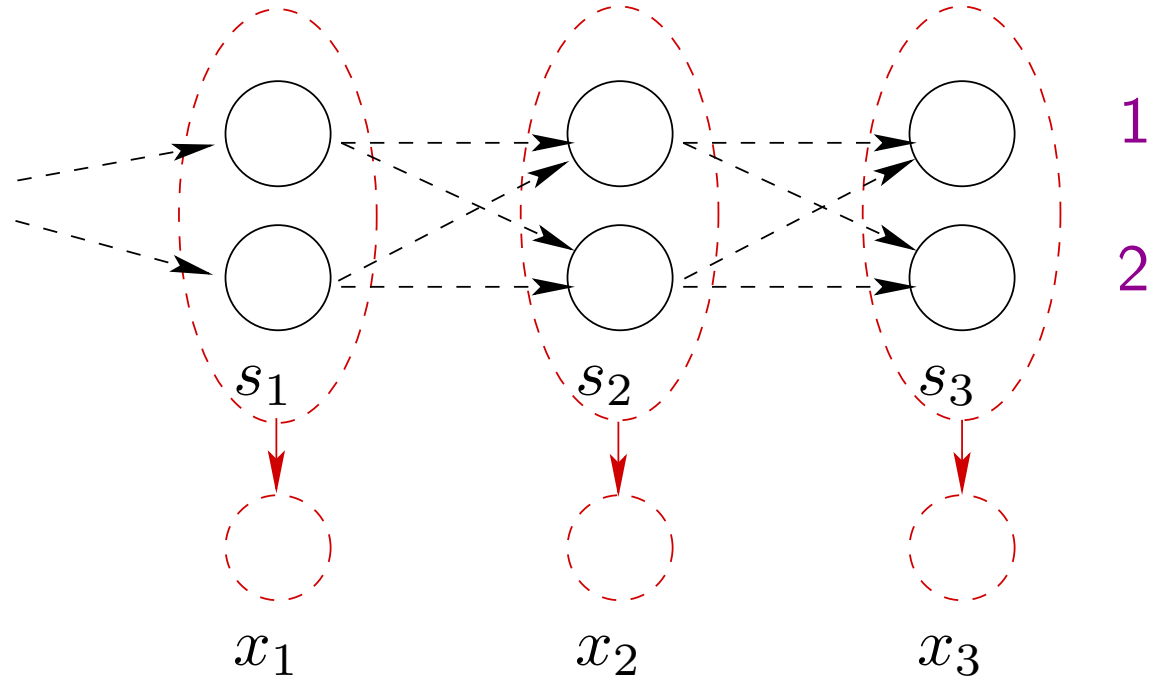
Max-probabilities



- We can recover the most likely hidden state sequence corresponding to a sequence of observations by evaluating the following max-probabilities:

$$\delta_t(i) = \max_{s_1, \dots, s_{t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_t, s_1, \dots, s_{t-1}, s_t = i)$$

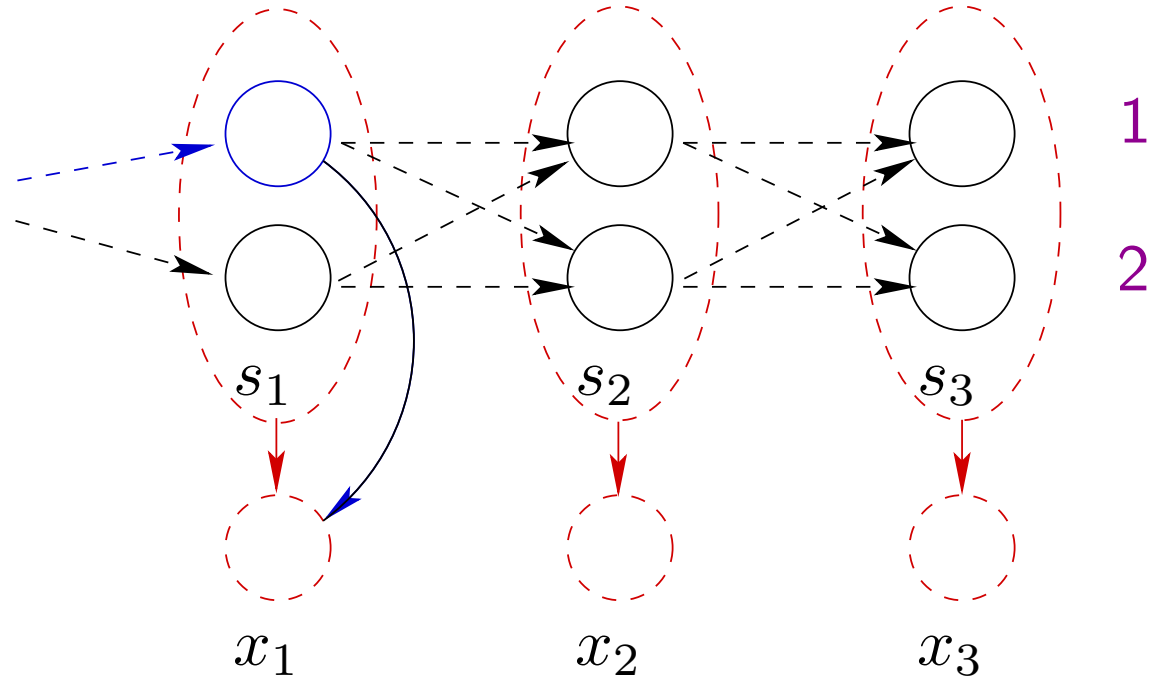
Viterbi algorithm



$$\delta_1(1) = P(x_1, s_1 = 1)$$

$$\delta_1(2) = P(x_1, s_1 = 2)$$

Viterbi algorithm

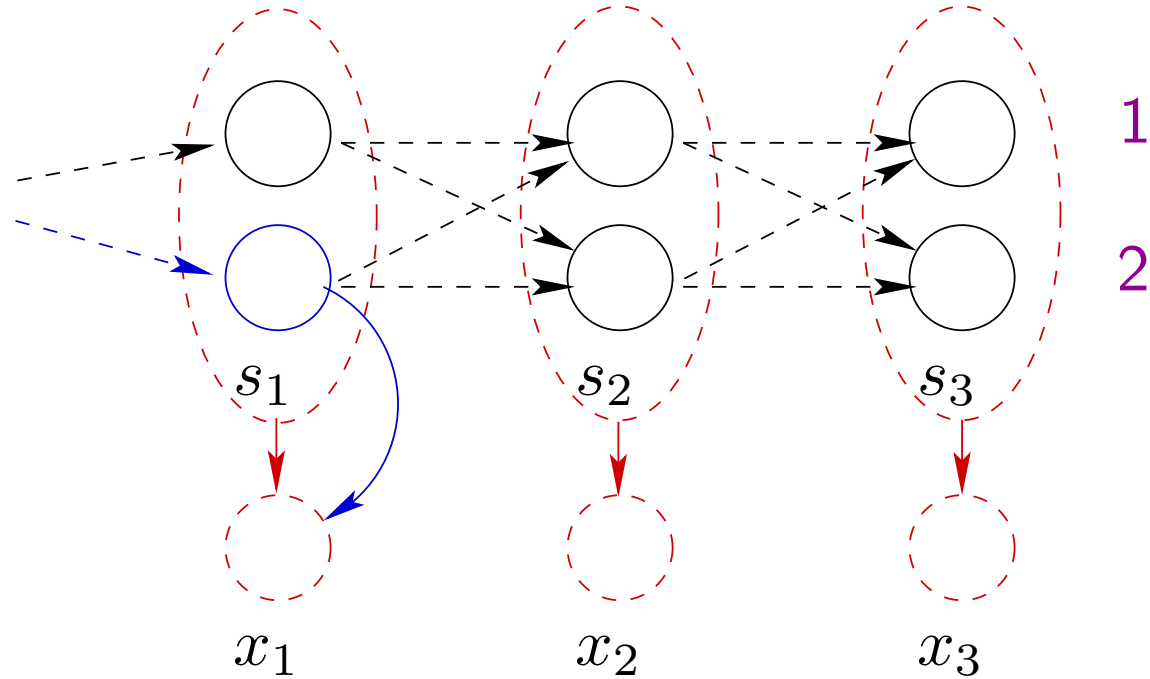


$$\delta_1(1) = P(x_1, s_1 = 1)$$

$$\delta_1(2) = P(x_1, s_1 = 2)$$

$$\delta_1(1) = P(1)P_x(\mathbf{x}_1|1)$$

Viterbi algorithm



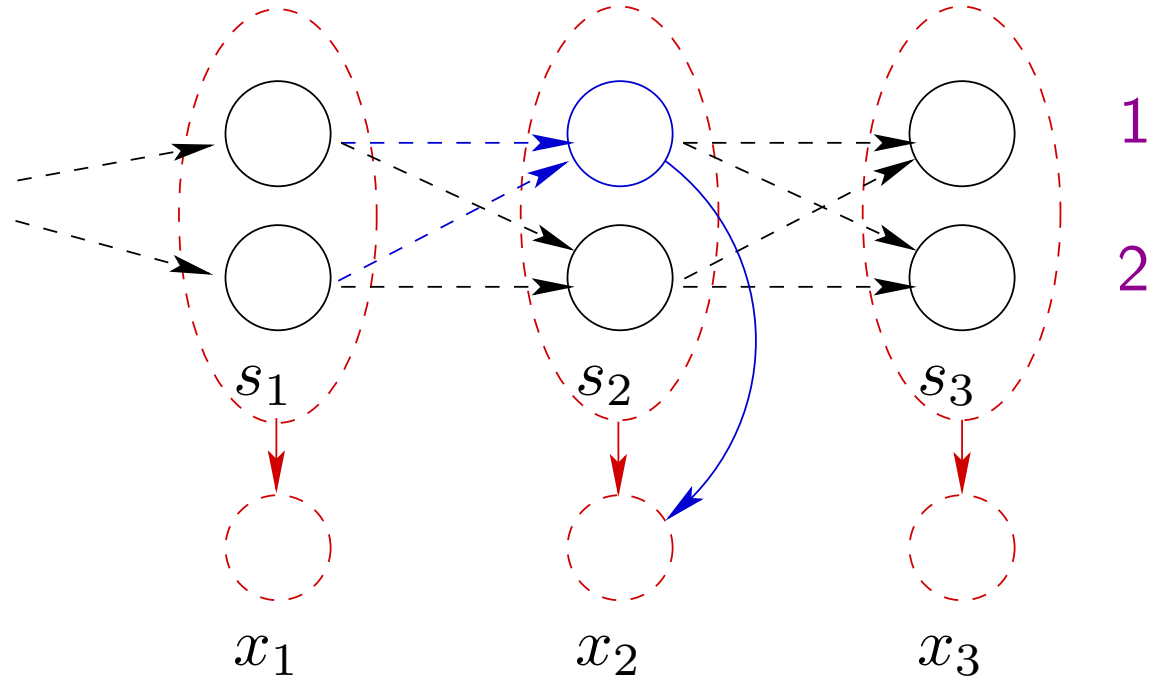
$$\delta_1(1) = P(x_1, s_1 = 1)$$

$$\delta_1(2) = P(x_1, s_1 = 2)$$

$$\delta_1(1) = P(1)P_x(\mathbf{x}_1|1)$$

$$\delta_1(2) = P(2)P_x(\mathbf{x}_1|2)$$

Viterbi algorithm

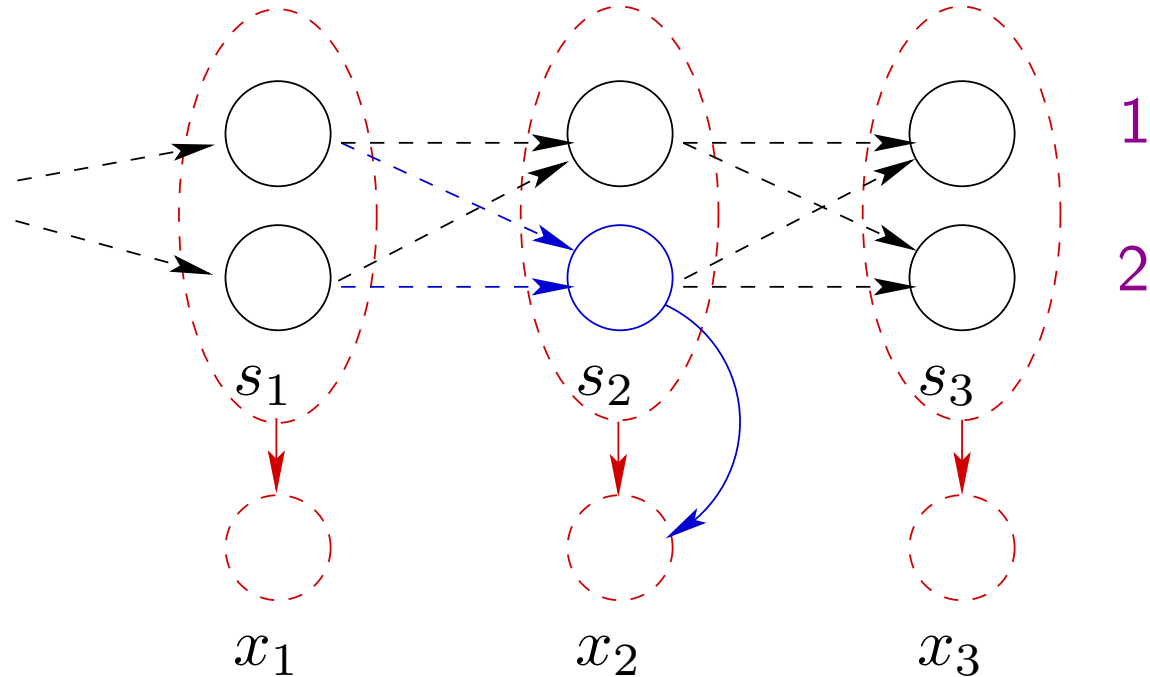


$$\delta_2(1) = \max_{s_1} P(x_1, x_2, s_1, s_2 = 1)$$

$$\delta_2(2) = \max_{s_1} P(x_1, x_2, s_1, s_2 = 2)$$

$$\delta_2(1) = \max \{ \delta_1(1)P_1(1|1), \delta_1(2)P_1(1|2) \} P_x(\mathbf{x}_2|1)$$

Viterbi algorithm



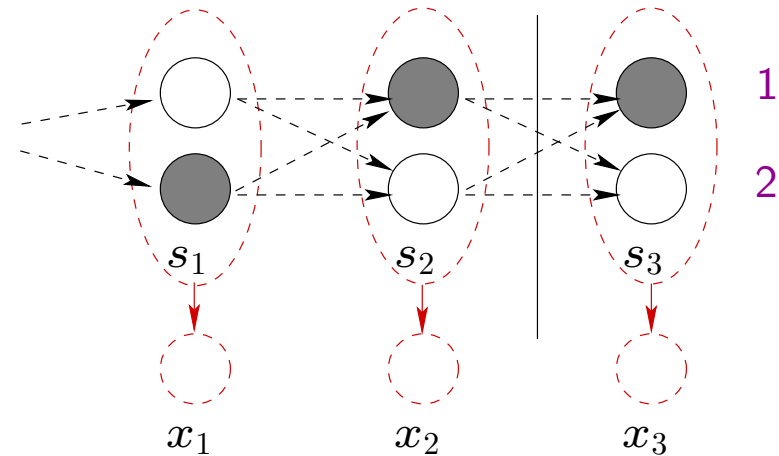
$$\delta_2(1) = \max_{s_1} P(x_1, x_2, s_1, s_2 = 1)$$

$$\delta_2(2) = \max_{s_1} P(x_1, x_2, s_1, s_2 = 2)$$

$$\delta_2(1) = \max \{ \delta_1(1)P_1(1|1), \delta_1(2)P_1(1|2) \} P_x(\mathbf{x}_2|1)$$

$$\delta_2(2) = \max \{ \delta_1(1)P_1(2|1), \delta_1(2)P_1(2|2) \} P_x(\mathbf{x}_2|2)$$

Viterbi algorithm



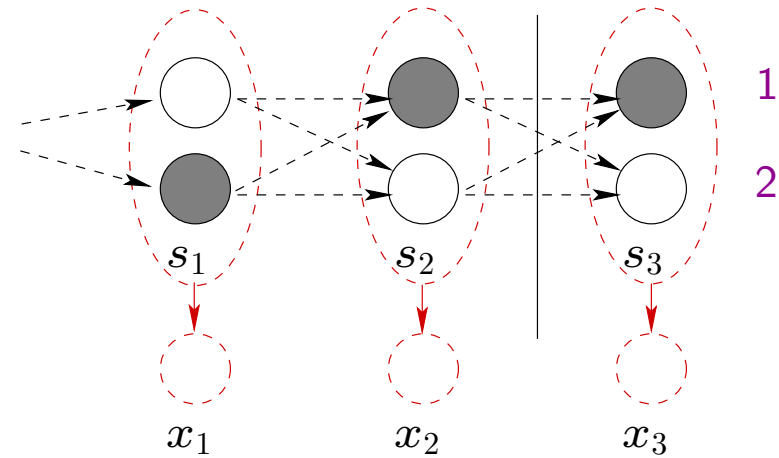
$$\delta_t(i) = \max_{s_1, \dots, s_{t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_t, s_1, \dots, s_{t-1}, s_t = i)$$

- We get the following recursive equation for calculating the max probabilities:

$$\delta_1(i) = P(i)P_x(\mathbf{x}_1|i)$$

$$\delta_t(i) = \max_j \{ \delta_{t-1}(j)P_1(i|j) \} P_x(\mathbf{x}_t|i)$$

Viterbi algorithm: back-tracking

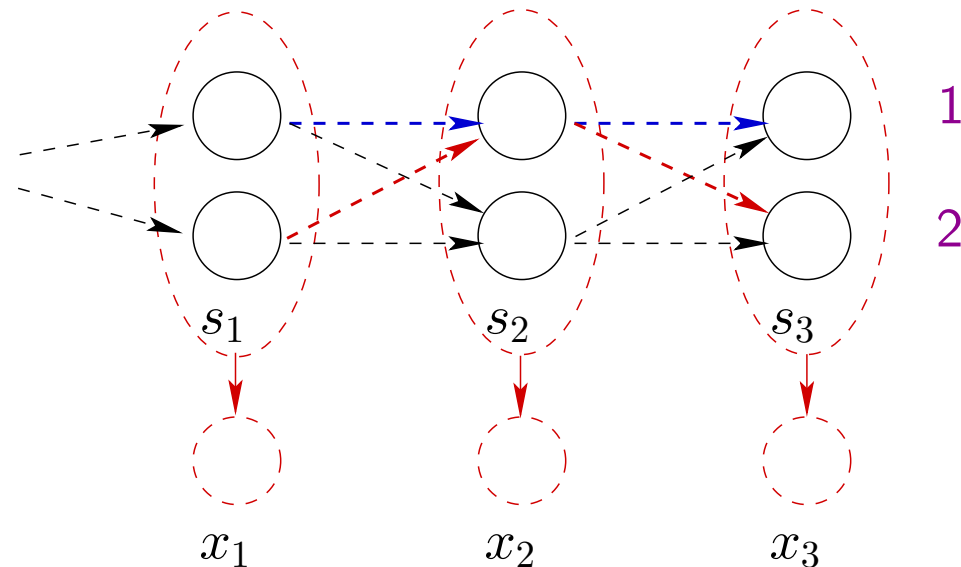


- We can recover the most likely state sequence by working backwards:

$$s_n^* = \operatorname{argmax}_i \delta_n(i)$$

$$s_t^* = \operatorname{argmax}_j \{ \delta_t(j) P_1(s_{t+1}^* | j) \}$$

Viterbi algorithm: properties

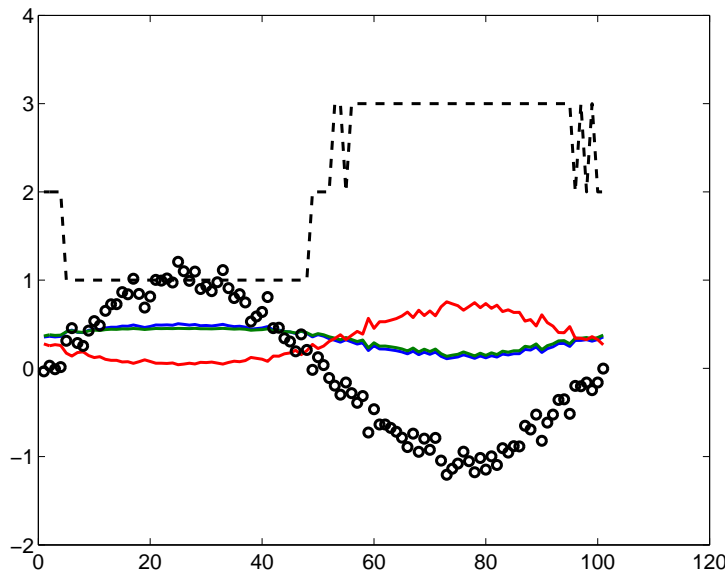


- The most likely path has the property that any partial path is also optimal:

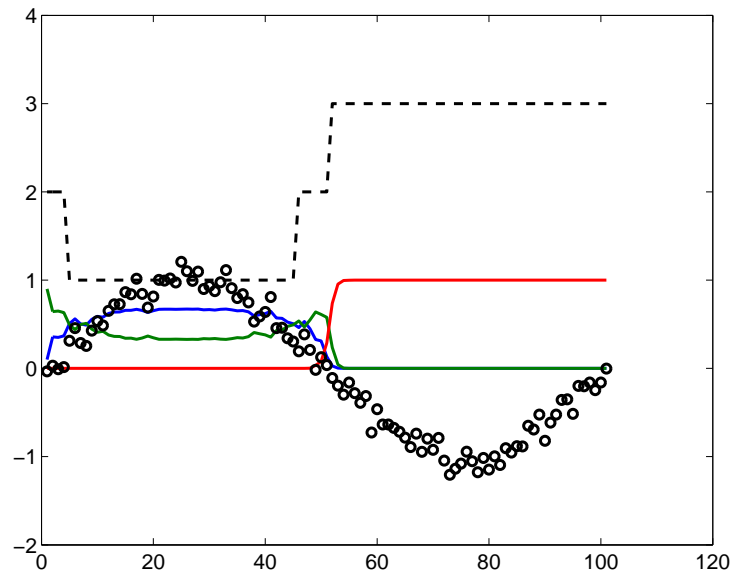
If $s_t^* = i$ then $\{s_1^*, \dots, s_t^*\}$ is also the most likely state sequence forced to end up in $s_t = i$ at time t given only $\mathbf{X}_1, \dots, \mathbf{X}_t$.

Viterbi algorithm: example

- Same example as in the EM case (3 states, Gaussian outputs)



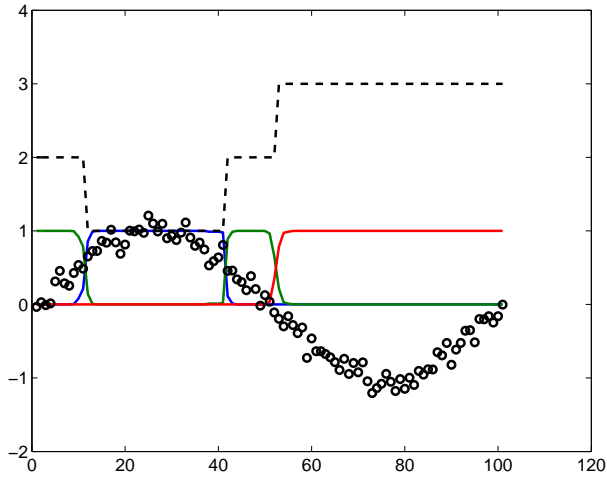
after 0 iterations



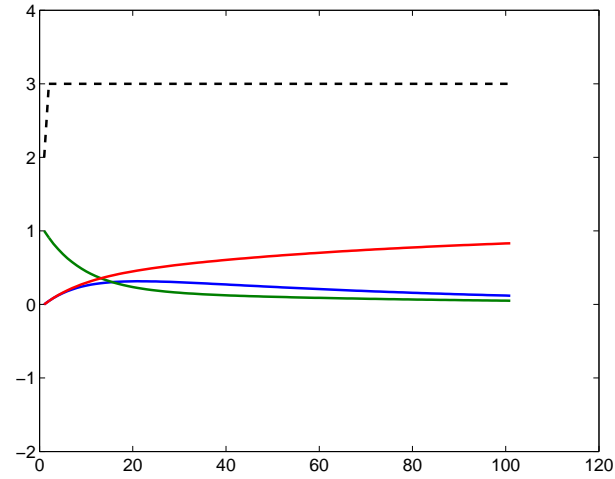
after 7 iterations

- The most likely hidden state sequence $\{s_0^*, \dots, s_n^*\}$ need not agree with the most likely states derived from the posterior marginals $\gamma_t(i)$

Example cont'd



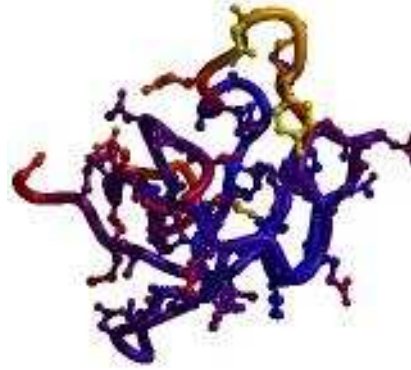
final



ML model, no observations

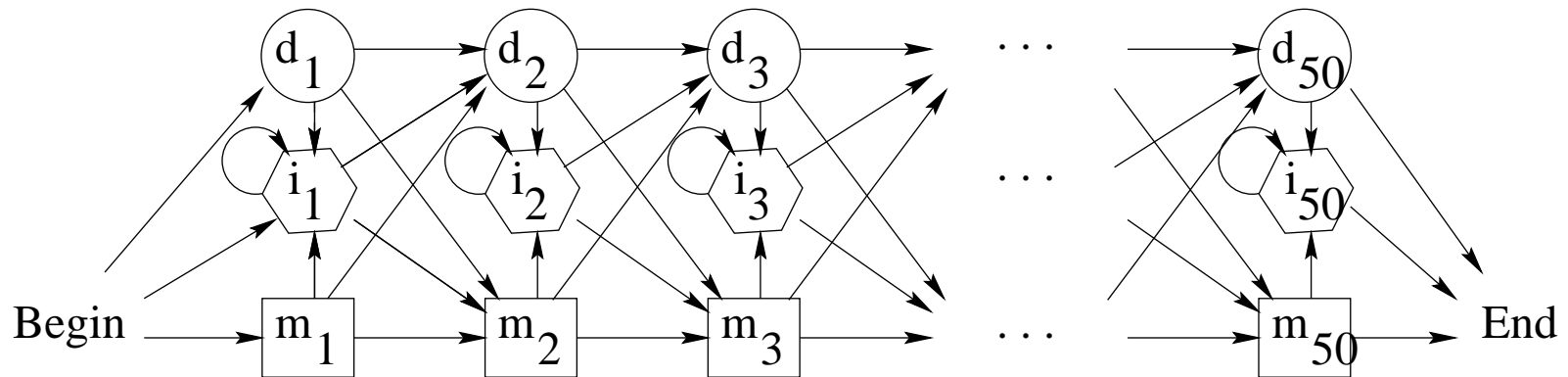
Linear HMMs, alignment

- Proteins



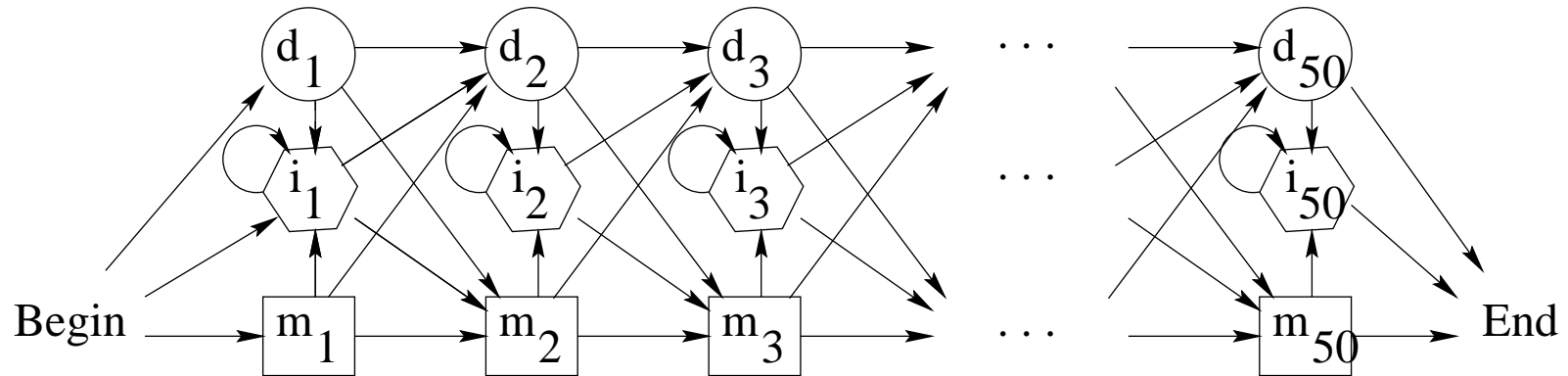
GVAALGAKVLAQIGVAVSHLGDEGKMVAQMKA VGRHKG YGNKH IKAQYFEPLGASLLS AMEHRIG

- A Linear HMM model for protein sequences



GVAALGAKVLAQIGVAVSHLGD~~egk~~MVAQMKA VGRHKG~~gyg~~NK-HIKAQYFEPLGASLLS AMEHRIG

Linear HMMs, multiple alignment



-VKGHGKKVADALTNVAHVDDMPNALSALSDLHA . . .HKLRVDPV.NFKLLSHCLLVTLAAHLP
 KVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHC . . .DKLHVDPE.NFRLLGNVLCVLAHHFG
 DLKKHGVTVLTALGAILKKKKGHHEAELKPLAQSHA . . .TK-HKIPikYLEFISEAIIHVLHSRHP
 PFETHANRIVGFFSKIIGELPNIEADVNTFVASHK . . .PR-GVTHD.QLNNFRAGFVSYMKAH--
 DVRWHAERIINAVNDAVASMDDtek . .MSMKLRDLSGKHA . . .KSFQVDPQ.YFKVLA AVIADTVAA---
 ELQAHAGKVFKLVEAAIQLQVtgvvvTDATLKNLGSVHV . . .SK-GVADA.HFPVVKEAILKTIKEVVG
 GVAALGAKVLAQIGVAVSHLGDegk . .MVAQMKAVGVRHKgygNK-HIKAQ.YFEPLGASLLSAMEHRIG