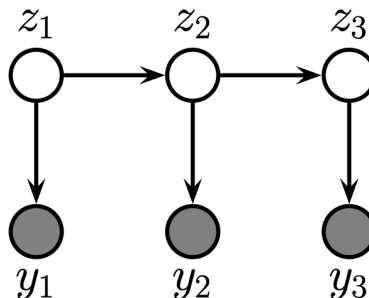


# 4 - Message Passing Algorithms

## Belief propagation on chains

Consider the **Hidden Markov Chain model** (all the latent variables are discrete):



The corresponding joint distribution has the form:

$$p(\mathbf{y}_{1:T}, \mathbf{z}_{1:T}) = \left[ p(\mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \right] \left[ \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{z}_t) \right]$$

Our goal is to *infer hidden states*, i.e. to compute the *posterior marginals*  $p(z_t | \mathbf{y}_{1:T})$ , using Bayes rule and Markov property we have:

$$p(z_t = j | \mathbf{y}_{t+1:T}, \mathbf{y}_{1:t}) \propto p(z_t = j | \mathbf{y}_{1:t}) p(\mathbf{y}_{t+1:T} | z_t = j, \cancel{\mathbf{y}_{1:t}})$$

The *given information* are:

1. Observations  $\mathbf{y}_{1:T}$
2. Initial distribution  $\pi_k = p(z_1 = k)$
3. Transition model  $A_{ij} = p(z_t = j | z_{t-1} = i)$
4. Local evidence  $p(y_t | z_t = k)$

Our algorithm breaks into two parts:

1. the forwards algorithm compute the filtering distribution  $p(z_t = j | \mathbf{y}_{1:t})$  by working forward in time,
2. the backward algorithm compute the  $p(\mathbf{y}_{t+1:T} | z_t = j)$  terms by working backward in time.

## The forwards algorithm (Bayes filter) (Online)

Define vectors:

- Belief states:  $\alpha_t(j) \triangleq p(z_t = j \mid \mathbf{y}_{1:t})$
- Local evidence:  $\lambda_t(j) \triangleq p(\mathbf{y}_t \mid z_t = j)$
- Transition matrix:  $A_{i,j} = p(z_t = j \mid z_{t-1} = i)$

The **predict step** is given by (**Chapman-Kolmogorov equation**):

$$\alpha_{t|t-1}(j) \triangleq p(z_t = j \mid \mathbf{y}_{1:t-1}) = \sum_i p(z_t = j \mid z_{t-1} = i) p(z_{t-1} = i \mid \mathbf{y}_{1:t-1}) = \sum_i A_{i,j} \alpha_{t-1}(i)$$

The **update step** is given by

$$\alpha_t(j) = \frac{1}{Z_t} p(\mathbf{y}_t \mid z_t = j) p(z_t = j \mid \mathbf{y}_{1:t-1}) = \frac{1}{Z_t} \lambda_t(j) \alpha_{t|t-1}(j) = \frac{1}{Z_t} \lambda_t(j) \left[ \sum_i \alpha_{t-1}(i) A_{i,j} \right]$$

where we used the law of total probability, Bayes' rule, and Markov property. In matrix notation we write  $\alpha_t = \text{normalize}(\lambda_t \odot (\mathbf{A}^\top \alpha_{t-1}))$ . The normalization constant for each time step is given by

$$Z_t \triangleq p(\mathbf{y}_t \mid \mathbf{y}_{1:t-1}) = \sum_{j=1}^K p(\mathbf{y}_t \mid z_t = j) p(z_t = j \mid \mathbf{y}_{1:t-1}) = \sum_{j=1}^K \lambda_t(j) \alpha_{t|t-1}(j)$$

## The forwards-backwards algorithm (Offline)

Define vectors:

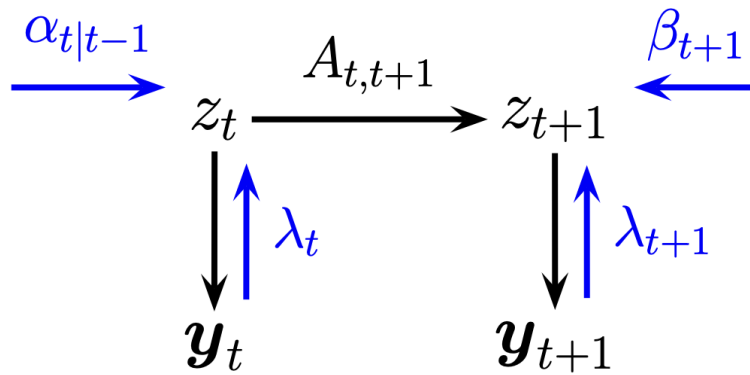
- the conditional likelihood  $\beta_t(j) \triangleq p(\mathbf{y}_{t+1:T} \mid z_t = j)$
- then  $\gamma_t(j) \triangleq p(z_t = j \mid \mathbf{y}_{t+1:T}, \mathbf{y}_{1:t}) = \alpha_t(j) \beta_t(j) = \text{normalize}(\alpha_t \odot \beta_t)$ .

We first have:

$$\beta_{t-1}(i) = p(\mathbf{y}_{t:T} \mid z_{t-1} = i) = \sum_j \beta_t(j) \lambda_t(j) A_{i,j}$$

in matrix-vector form:  $\beta_{t-1} = \mathbf{A} (\lambda_t \odot \beta_t)$ . The **base case** is

$$\beta_T(i) = p(\mathbf{y}_{T+1:T} \mid z_T = i) = p(\emptyset \mid z_T = i) = 1.$$



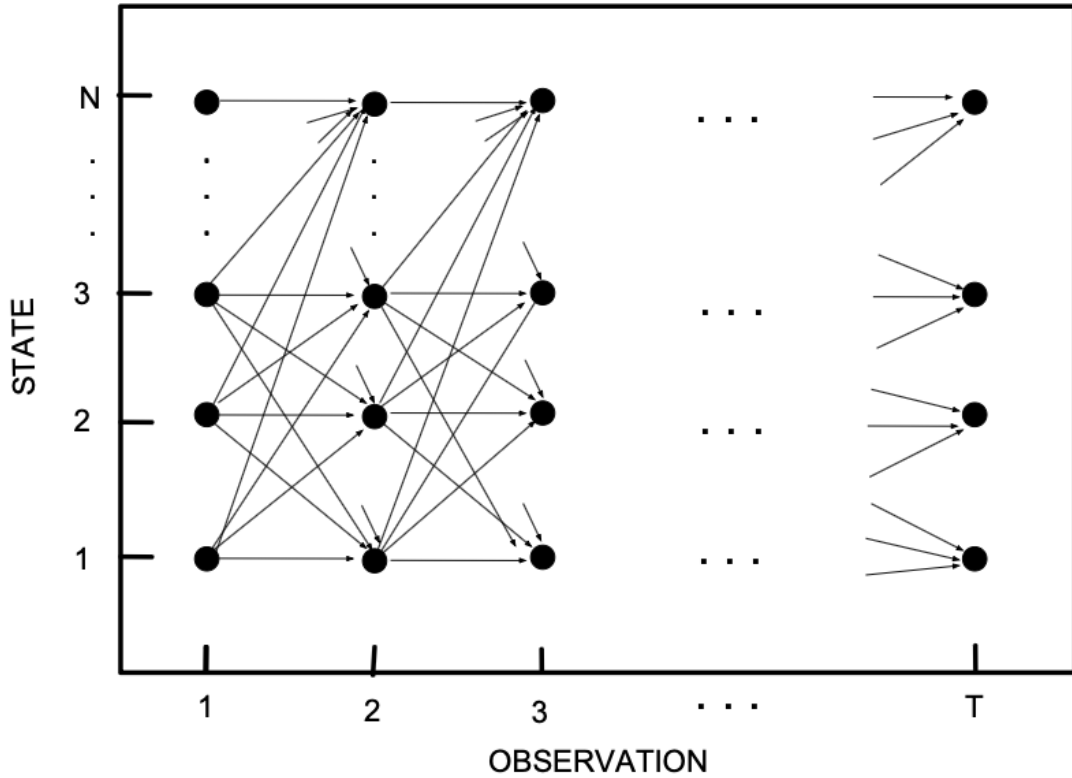
## The Viterbi algorithm

Now we consider the MAP estimate, finding the sequences with maximum posterior probability:

$$\begin{aligned} \mathbf{z}_{1:T}^* &= \operatorname{argmax}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} \mid \mathbf{y}_{1:T}) = \operatorname{argmax}_{\mathbf{z}_{1:T}} \log p(\mathbf{z}_{1:T} \mid \mathbf{y}_{1:T}) \\ &= \operatorname{argmax}_{\mathbf{z}_{1:T}} \log \pi_1(\mathbf{z}_1) + \log \lambda_1(\mathbf{z}_1) + \sum_{t=2}^T [\log A(\mathbf{z}_{t-1}, \mathbf{z}_t) + \log \lambda_t(\mathbf{z}_t)] \end{aligned}$$

This is equivalent to computing a *shortest path through the trellis diagram* below, where the nodes are possible states at each time step, and

- the nodes are  $\log \lambda_t(\mathbf{z}_t)$ , except the first layer nodes are  $\log \pi_1(\mathbf{z}_1) + \log \lambda_1(\mathbf{z}_1)$ ,
- the edges are weighted by  $\log A(\mathbf{z}_{t-1}, \mathbf{z}_t)$ .



The Viterbi algorithm is a *dynamical programming* algorithm, the algorithm is as follow:

**Forward Pass:** The maximum probability we can assign to the data up to  $z_t$  if we end up in state  $j$  is given by the path

$$\delta_t(j) \triangleq \max_{z_1, \dots, z_{t-1}} p(z_{1:t-1}, z_t = j, \mathbf{y}_{1:t})$$

The most probable path to state  $j$  at time  $t$  must consist of the most probable path to some other state  $i$  at time  $t - 1$ , so

$$\delta_t(j) = \lambda_t(j) \left[ \max_i \delta_{t-1}(i) A_{i,j} \right]$$

and we initialize by setting  $\delta_1(j) = \pi_j \lambda_1(j)$ . We also need to keep track of the most likely previous (ancestor) state, for each possible state that we end up in:

$$a_t(j) \triangleq \operatorname{argmax}_i \delta_{t-1}(i) A_{i,j}$$

At the end of forward pass, we will end up with, for each  $t$  and  $j$ , an index  $a_t(j)$ .

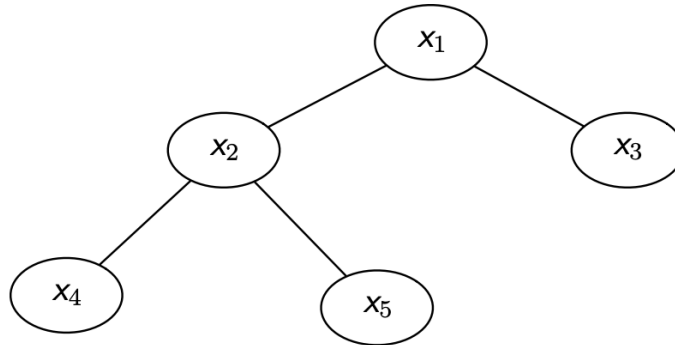
**Backward pass:** In the backwards pass, we *compute the most probable sequence* of states using a traceback procedure, as follows:  $z_t^* = a_{t+1}(z_{t+1}^*)$ , where we initialize using  $z_T^* = \operatorname{argmax}_i \delta_T(i)$ .

Overall, the time complexity of Viterbi is clearly  $O(K^2T)$  in general, and the space complexity is  $O(KT)$ .

## Belief propagation on trees

### Sum-product algorithm

We assume that our model is an undirected tree, we pick an arbitrary node as a root, and orient all the edges downwards away from this root, so that each node has a unique parent:



In the graph we define the message passing:

$$\begin{aligned}
 m_{4 \rightarrow 2}(z_2) &= \sum_{z_4} \psi_4(z_4) \psi_{24}(z_2, z_4) \\
 m_{5 \rightarrow 2}(z_2) &= \sum_{z_5} \psi_5(z_5) \psi_{25}(z_2, z_5) \\
 m_{3 \rightarrow 1}(z_1) &= \sum_{z_3} \psi_3(z_3) \psi_{13}(z_1, z_3) \\
 m_{2 \rightarrow 1}(z_1) &= \sum_{z_2} \psi_2(z_2) \psi_{12}(z_1, z_2) m_{4 \rightarrow 2}(z_2) m_{5 \rightarrow 2}(z_2)
 \end{aligned}$$

Finally, we can obtain the marginal distribution over  $z_1$ :

$$p(z_1) \propto \psi_1(z_1) m_{2 \rightarrow 1}(z_1) m_{3 \rightarrow 1}(z_1).$$

**Upward Passing (Collect to root):** *In general*, we define the **belief state** for each variable:

$$\text{bel}_s(z_s) \propto \psi_s(z_s) \prod_{t \in \text{ch}_s} m_{t \rightarrow s}(z_s)$$

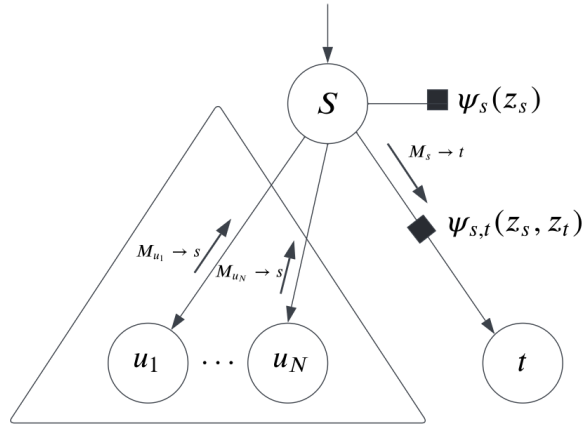
To compute the outgoing message that  $s$  should send to its parent  $t$ , we pass the local belief through the pairwise potential linking  $s$  and  $t$ , and then marginalize out  $s$  to get

$$m_{s \rightarrow t}(z_t) = \sum_{z_s} \psi_{st}(z_s, z_t) \text{bel}_s(z_s)$$

Then, at the root of the tree,  $\text{bel}_t(z_t) = p(z_t)$  will have seen all the evidence.

**Downward Passing (Distribute from root):** Send messages back down to the leaves. The message that  $s$  sends to its child  $t$  is given by

$$m_{s \rightarrow t}(z_t) = \sum_{z_s} \left( \psi_s(z_s) \psi_{st}(z_s, z_t) \prod_{u \in \text{ch}_s \setminus t} m_{u \rightarrow s}(z_s) \right) = \sum_{z_s} \psi_{st}(z_s, z_t) \frac{\text{bel}_s(z_s)}{m_{t \rightarrow s}(z_s)}$$



We use this to update the belief state at node  $t$  to get:

$$\text{bel}_t(z_t) \propto \text{bel}_t(z_t) m_{s \rightarrow t}(z_t)$$

Now, after both passing, we have that for each  $z_t$ :

$$\text{bel}_t(z_t) = p(z_t)$$

## Max-product algorithm

We can replace the sum operation with the max operation to get max-product belief propagation. The result of this computation are a set of max marginals for each node, so we have that for each  $z_i$ :

$$\text{bel}_i^*(z_i) \propto \max_{z_{-i}} p(z_i = k, z_{-i} | \mathbf{y})$$

Now we can derive two different kinds of “MAP” estimates from these local quantities:

- **Maximizer of the posterior marginal (MPM):**  $\hat{z}_i = \text{argmax}_{z_i} \text{bel}_i(z_i)$ ,
- **Maximizer of the max marginal or (MMM):**  $\tilde{z}_i = \text{argmax}_{z_i} \text{bel}_i^*(z_i)$ .

Let:

- $\hat{\mathbf{z}} = [\hat{z}_1, \dots, \hat{z}_{N_z}]$  be the MPM estimate,
- $\tilde{\mathbf{z}} = [\tilde{z}_1, \dots, \tilde{z}_{N_z}]$  be the MMM estimate,
- $\mathbf{z}^* = \operatorname{argmax}_{\mathbf{z}} p(\mathbf{z} | \mathbf{y})$  be the true MAP estimate.