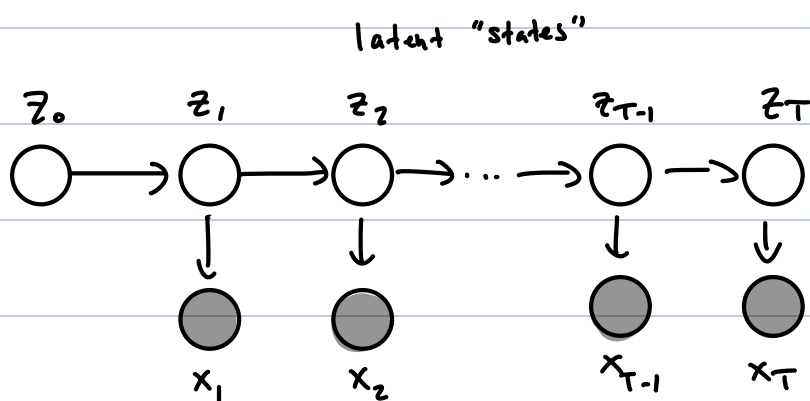


# State-space models (SSMs) (very general)



$$z_0 \sim P(z_0)$$

for  $t = 1 \dots T$ :

$$z_t \sim P(z_t | z_{t-1})$$

$$x_t \sim P(x_t | z_t)$$

- $P_{\phi}(x_t | z_t)$  "emission" distribution (likelihood)  
with fixed parameters  $\phi$
- $P_{\lambda}(z_t | z_{t-1})$  "transition" distribution  
with fixed parameters  $\lambda$
- $P(z_0)$  "initial" distribution  
 $\pi$  with fixed parameters  $\pi$

---

## Hidden Markov Models (HMMs) are SSMs

- $z_t \in [K] \equiv \{1 \dots K\}$  discrete states
- $P_{\lambda}(z_t = k | z_{t-1} = j) = \lambda(j, k)$   
"transition matrix"
- $\sum_{k=1}^K \lambda(j, k) = 1$  "row-stochastic" matrix
- $P_{\pi}(z_0 = k) = \pi_0(k)$  "initial dist"
- $x_t \in \mathbb{R}^d$  (no constraints on emission dist)

## Examples...

- speech recognition

$x_t \in \mathbb{R}^d$  acoustic signals

$z_t \in \text{WORDS}$  word taken

$p_{\lambda}(z_t | z_{t-1})$  language model

$p_{\theta}(x_t | z_t)$  acoustic model

- activity recognition

$x_t \in \mathbb{R}^d$  video frame of subject

$z_t \in \text{ACTIVITIES}$  activity of subject

e.g. is the zebrafish "hunting", "sleeping", ... at each  $t = 1 \dots T$

e.g. is the enemy aircraft "attacking", "escorting", etc. at  $t$

- gene finding :

$t = \text{locus}$  of nucleotide

$x_t \in \{A, C, G, T\}$

$z_t \in \text{REGIONS}$ , e.g. "gene-coding" region

# Types of inference in SSMs

filtering :  $P(z_t \mid X_{1:t})$  "online"

prediction :  $P(z_{t+1} \mid X_{1:t})$

smoothing :  $P(z_t \mid X_{1:T})$  "offline"

forecasting :  $P(x_{t+1} \mid X_{1:t})$

imputing :  $P(x_t \mid x_{1:T \setminus t})$   
(or denoising, or anomaly detection)

evidence :  $P(x_{1:T})$

## FORWARD pass

- Goal :  $P(x_{T+1} | x_{1:T})$  "forecasting"

- Requires :  $P(z_{T+1} | x_{1:T})$  "prediction"

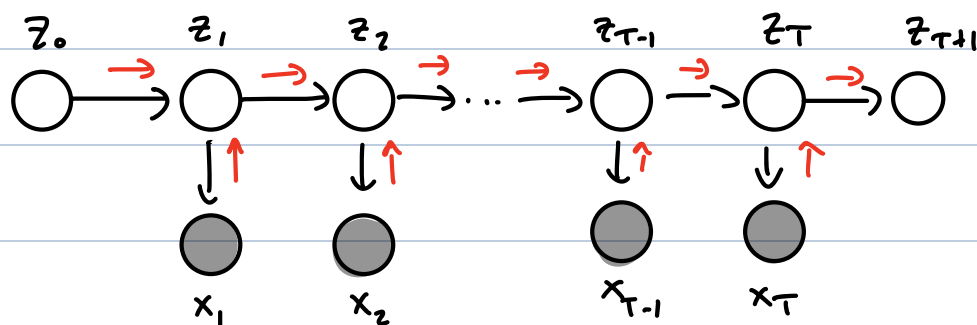
↑ target of inference

- HMMs are trees; variable elimination

- Set  $z_{T+1}$  as root of new tree

- point edges away

- ORDER by depth-first :  $\{x_1, \dots, x_T, z_1, \dots, z_{T+1}\}$



- Messages from observations  $x_1, \dots, x_T$

$$m_{x_t \rightarrow z_t}(z_t) = \sum_{x_t} \psi(x_t) \psi(x_t, z_t)$$

$$= \sum_{x_t} \delta_{\bar{x}_t}(x_t) P_{\bar{x}}(x_t | z_t)$$

$$= P_{\bar{x}}(\bar{x}_t | z_t)$$

- Vector notation

$$l_t = \begin{bmatrix} P_{\bar{x}}(\bar{x}_t | z_t = 1) \\ \vdots \\ P_{\bar{x}}(\bar{x}_t | z_t = K) \end{bmatrix} \equiv m_{x_t \rightarrow z_t}$$

## Messages between states

$$\bullet \mathcal{M}_{z_0 \rightarrow z_1}(z_1) = \sum_{z_0} \varphi(z_0) \psi(z_0, z_1)$$

$\swarrow$  "root"

$$= \sum_{z_0} p(z_0) p(z_1 | z_0) \equiv p(z_1)$$

$$= \sum_{z_0} \pi_0(z_0) \Lambda(z_0, z_1) \equiv \pi_0^T \Lambda(:, z_1)$$

$$\bullet \mathcal{M}_{z_1 \rightarrow z_2}(z_2) = \sum_{z_1} \varphi(z_1) \psi(z_1, z_2) \mathcal{M}_{z_0 \rightarrow z_1}(z_1) \mathcal{M}_{x_1 \rightarrow z_1}(z_1)$$

$$= \sum_{z_1} p(z_2 | z_1) p(z_1) p(\bar{x}_1 | z_1) \equiv p(z_2, \bar{x}_1)$$

$$= \sum_{z_1} \Lambda(z_1, z_2) \mathcal{M}_{z_0 \rightarrow z_1}(z_1) l_1(z_1)$$

All future messages follow the same recursion:

$$\bullet \text{ Define } \alpha_t(z_t) \stackrel{\circ}{=} \mathcal{M}_{z_{t-1} \rightarrow z_t}(z_t)$$

e.g.,  $\alpha_2(z_2)$

$$\bullet \alpha_t(z_t) = \sum_{z_{t-1}} \Lambda(z_{t-1}, z_t) \alpha_{t-1}(z_{t-1}) l_{t-1}(z_{t-1})$$

$$\alpha_t = \Lambda^T (\alpha_{t-1} \odot l_{t-1}), \quad t = 2 \dots T+1$$

$$\alpha_1 = \pi^T \Lambda \quad (\text{"base case"}) \quad t=1$$

- Each message computes :

$$\alpha_{t+1}(z_{t+1}) = P(z_{t+1}, \bar{x}_{1:t})$$

- e.g.  $\alpha_3(z_3) = \sum_{z_2} \Lambda(z_1, z_3) \alpha_2(z_2) \ell_2(z_2)$   
 $= \sum_{z_2} P(z_3 | z_2) P(z_2, \bar{x}_1) P(\bar{x}_2 | z_2)$   
 $= P(z_3, \bar{x}_{1:2})$

- last message  $\alpha_{T+1}(z_{T+1}) = P(z_{T+1}, \bar{x}_{1:T})$

- $P(z_{T+1} = k | \bar{x}_{1:T}) = \frac{\alpha_{T+1}(k)}{\sum_j \alpha_{T+1}(j)}$

- $P(\bar{x}_{1:T}) = \sum_{j=1}^k \alpha_{T+1}(j)$

- What about  $P(z_t | \bar{x}_{1:T})$ ?

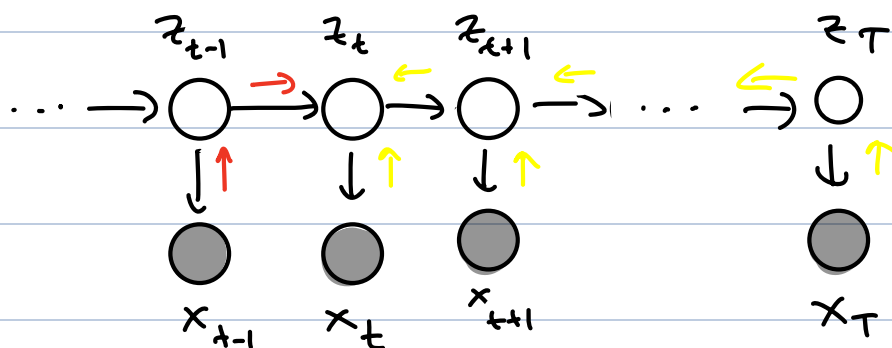
- Can we compute it for  $t < T+1$  using only the forward messages?

- No. why not?  $z_t$  needs "backward" info from  $\bar{x}_{t+1:T}$ .

## BACKWARD PASS :

$$\underline{Goal} : P(z_t | \bar{x}_{1:T})$$

- Set  $z_t$  as root of new tree
- point edges away
- ORDER by depth-first:  $\{x, \dots, x_T, z, \dots, z_{t+1}\}$



$$\bullet \mathcal{M}(z_{t-1}) = \sum_{\substack{z \rightarrow z_{t-1} \\ z_T}} \varphi(z_T) \varphi(z_{t-1}, z_T) \mathcal{M}(z_T)$$

$$= \sum_{z_T} P(z_T | z_{t-1}) P(\bar{x}_T | z_T) \equiv P(\bar{x}_T | z_{t-1})$$

$$\bullet \beta_t(z_t) \stackrel{\circ}{=} \mathcal{M}(z_t)_{z \rightarrow z_t, t-1}, \text{ e.g. } \beta_{T-1}(z_{T-1})$$

$$\bullet \beta_{T-2}(z_{T-2}) = \sum_{z_{T-1}} \varphi(z_{T-1}) \varphi(z_{T-2}, z_{T-1}) \mathcal{M}(z_{T-1})_{x \rightarrow z_{T-1}, T-1} \beta_{T-1}(z_{T-1})$$

$$= \sum_{z_{T-1}} P(z_{T-1} | z_{T-2}) P(\bar{x}_{T-1} | z_{T-1}) \beta_{T-1}(z_{T-1})$$

$$= P(\bar{x}_{T-1:T} | z_{T-2})$$

$$\begin{aligned} \bullet \quad \beta_t(z_t) &= \sum_{z_{t+1}} \Lambda(z_t, z_{t+1}) l_{t+1}(z_{t+1}) \beta_{t+1}(z_{t+1}) \\ &= P(\bar{x}_{t+1:T} | z_t) \end{aligned}$$

$$\beta_T = \frac{1}{K} \quad \text{"base case"} \quad t = T$$

$$\beta_t = \Lambda(l_{t+1} \odot \beta_{t+1}) \quad \text{for } t = T-1 \dots 0$$

Putting it all together...

The 3 messages sent to  $z_t$ :

$$\textcircled{1} \quad \alpha_t(z_t) = P(z_t, \bar{x}_{1:t-1})$$

$$\textcircled{2} \quad \beta_t(z_t) = P(\bar{x}_{t+1:T} | z_t)$$

$$\textcircled{3} \quad l_t(z_t) = P(\bar{x}_t | z_t)$$

$$\alpha_t(z_t) \beta_t(z_t) l_t(z_t) = P(z_t, \bar{x}_{1:T})$$

$$P(z_t = k | \bar{x}_{1:T}) = \frac{\alpha_t(k) \beta_t(k) l_t(k)}{\sum_j \alpha_t(j) \beta_t(j) l_t(j)}$$



- The "Forwards - Backwards" algo (1986):

$$\alpha_t = \dots \text{ for } t = 1 \dots T+1$$

$$\beta_t = \dots \text{ for } t = T \dots 0$$

- Invented in 1980s; special case of belief prop.
- Estimates all singleton posterior marginals exactly

$$P(z_t | \bar{x}_{1:T}) \quad t = 1 \dots T$$

↑ "beliefs"

$$\bullet O(k^2)$$


---

What about  $P(z_1 \dots z_T | \bar{x}_{1:T})$ ?

- Many applications motivate inference of sequences
- Instantiating entire joint distribution is  $O(k^T)$
- Working with individual sequences is tractable:

① Evaluating probability of a sequence

$$P(z_1 \dots z_T | \bar{x}_{1:T}) = \frac{\prod_{t=1}^T P(z_t | z_{t-1}) P(\bar{x}_t | z_t)}{P(\bar{x}_{1:T})}$$

② Sampling a sequence:

$$z_1 \dots z_T \sim P(z_1 \dots z_T | \bar{x}_{1:T})$$

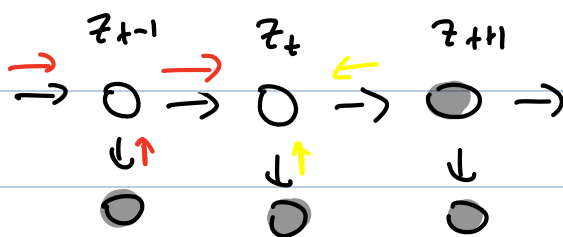
③ Maximizing

$$z_{1:T}^* = \underset{z_{1:T}}{\operatorname{argmax}} P(z_{1:T} | \bar{x}_{1:T})$$

## Forward-filtering backward sampling (FFBS)

- ① forward pass  $d_t = \dots$
- ②  $z_T \sim p(z_T | \bar{x}_{1:T})$
- ③  $z_t \sim p(z_t | \bar{z}_{t+1}, \bar{x}_{1:T}) \quad t = T-1 \dots 0$

$$p(z_t | \bar{z}_{t+1}, \bar{x}_{1:T}) \propto \alpha_t(z_t) q_t(z_t) p(\bar{z}_{t+1} | z_t)$$



MAP

$$\max_z p(z | \bar{x}) = \max_{z_1} p(z_1 | \bar{x}) \cdots \max_{z_T} p(z_T | z_{T-1}, \bar{x})$$

$$\hat{\beta}_T(z_T) = 1$$

$$\tilde{\beta}_t(z_t) = \max_{z_{t+1}} p(z_{t+1} | z_t) p(\bar{z}_{t+1} | z_{t+1}) \tilde{\beta}_{t+1}(z_{t+1}) \quad t = T-1 \dots 0$$

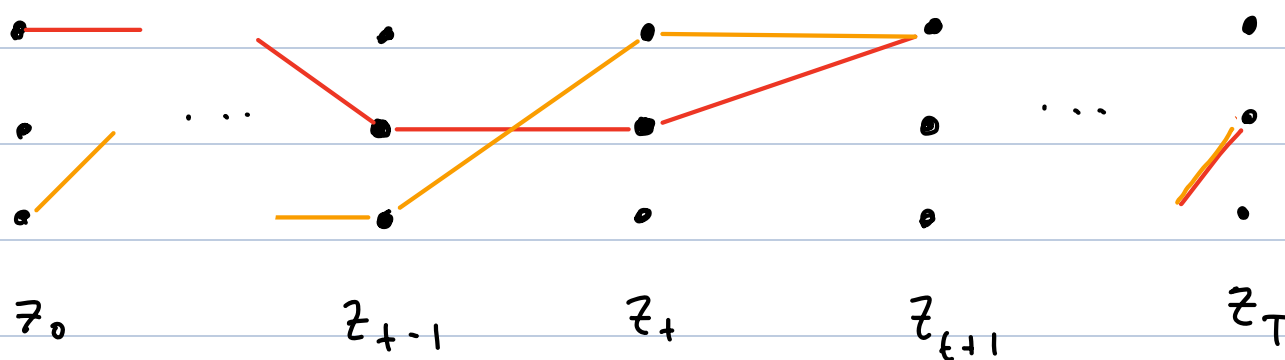
$$\tilde{\beta}_t^*(z_t) = \underset{z_{t+1}}{\operatorname{argmax}} \quad " \quad " \quad t = T-1 \dots 0$$

$$z_0^* = \arg \max_{z_0} \pi_*(z_0) \hat{\beta}_0(z_0)$$

$$z_t^* = \tilde{\beta}_{t-1}^* (z_{t-1}^*) \quad \text{für } t = 1 \dots T$$

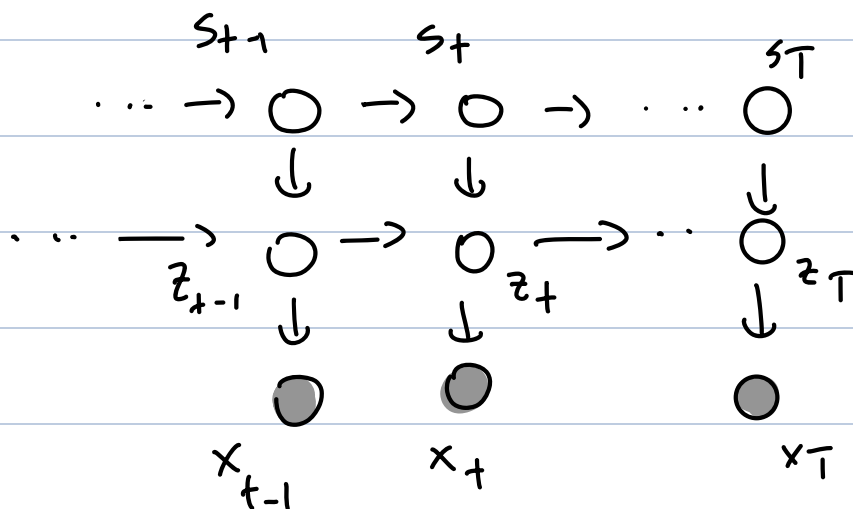
This is the Viterbi algorithm <sup>(1967)</sup> (special case of max-prod).

Often visualized with a lattice:



## "Switching" HMMs

- $P(z_t = k \mid z_{t-1} = j, s_t = s) = \Lambda^s(j, k)$
- $P(s_t = s \mid s_{t-1} = r) = \Delta(r, s)$



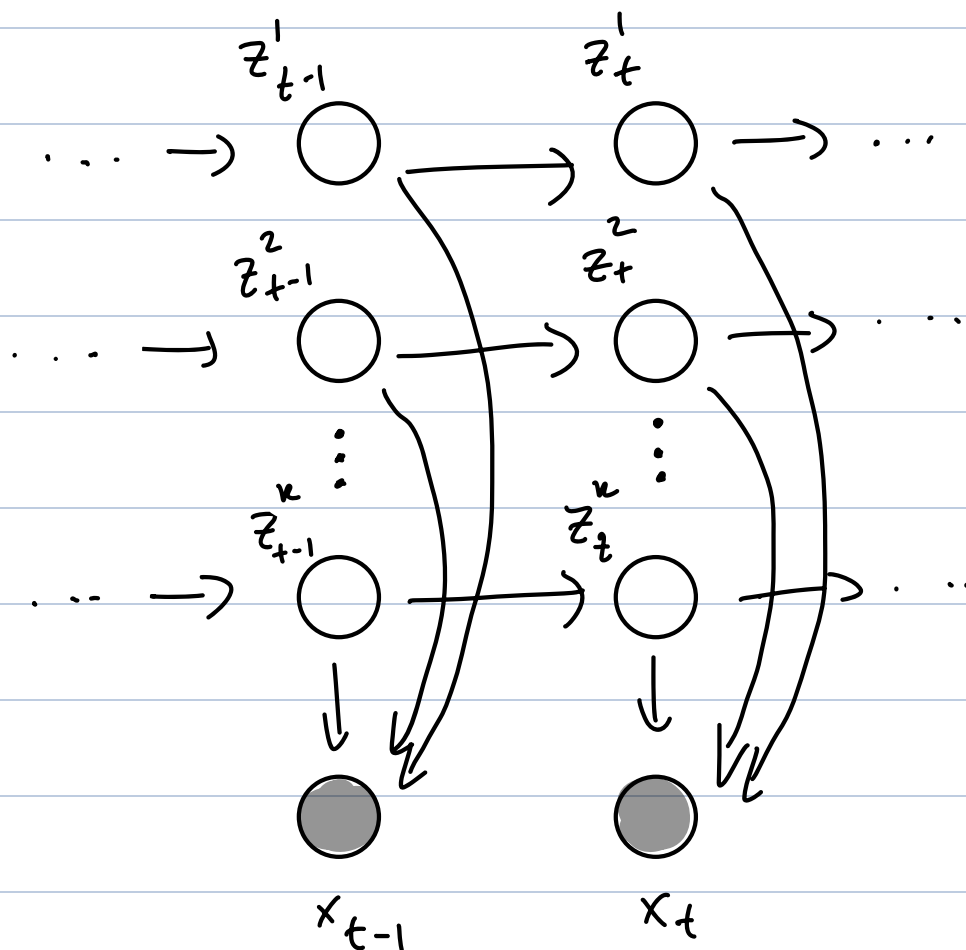
No longer a tree. How do we do inference?

$$\dots \sum_{s_T} P(s_T \mid s_{T-1}) \sum_{z_T} P(z_T \mid z_{T-1}, s_T) P(\bar{x}_T \mid z_T)$$

$$= \dots \sum_{s_T, z_T} P(s_T, z_T \mid s_{T-1}, z_{T-1}) P(\bar{x}_T \mid s_T, z_T)$$

Redefine  $(s_t, z_t)$  as a "super node"

## "factorial" HMMs



$$\vec{z}_t = [z_t^1 \dots z_t^k]^T, \quad z_t^k \in [p]$$

"distributed representation"

- e.g. to represent 30 bits of info from the history  $\bar{x}_{1:t}$ , a vanilla HMM would require  $z_{t+1} \in [2^{30}]$ , whereas, a factorial HMM could do this with  $p=2, k=30$ .
- Connection to general state-space models

# Learning in HMMs

- We have only thus far considered inference in graphical models (e.g.  $P(x_F | \bar{x}_E)$ ) assuming that the parameters are known

- Parameters in an HMM:

$$\Theta = \{ \Lambda, \pi_0, \Phi \}$$

                    ↑                    ↑                    ↑  
                    "transition"      "initial"      "emission"

- Type-II MLE would be: ← "evidence"

$$\hat{\Theta} \leftarrow \underset{\Theta}{\operatorname{argmax}} P(\bar{x}_{1:T}; \Theta)$$

- For a given  $\Theta$ , we can compute the evidence using just the forward pass
- One option is gradient ascent ...
- Another: EM (next class)

## HMM with Poisson observations

$$P(x_t \mid z_t = k) = \text{Pois}(\mu_k)$$

Goal:  $\hat{M} \leftarrow \underset{M}{\text{argmax}} P(X_{1:T}; M)$

This is easier w.r.t the complete data likelihood:

$$P(X_{1:T}, z_{1:T}; M) \propto_M \prod_t \text{Pois}(x_t; \mu_{z_t})$$

Define  $f_{tk} \stackrel{\Delta}{=} \mathbb{1}(z_t = k)$

$\propto_M \prod_t \prod_k \text{Pois}(x_t; \mu_k)^{f_{tk}}$

$$\frac{\partial}{\partial \mu_k} \log \prod_t \prod_k \text{Pois}(x_t; \mu_k)^{f_{tk}}$$

$$\propto_{\mu_k} \frac{\partial}{\partial \mu_k} \sum_t \sum_k f_{tk} [x_t \log \mu_k - \mu_k]$$

$$\propto_{\mu_k} \sum_t f_{tk} x_t \frac{1}{\mu_k} - \sum_t f_{tk}$$

$$= \frac{\sum_t f_{tk} x_t}{\mu_k} - N_{tk} = 0 \rightarrow \hat{\mu}_k = \frac{1}{N_{tk}} \sum_t f_{tk} x_t$$

Not surprising:  $\hat{\mu}_k = \hat{\mathbb{E}}[x_t \mid z_t = k]$

$P(x, z; \mu)$  is easily optimizable.

We don't observe  $z_{1:T}$ . Instead:

$$\log P_\theta(x) = \log \sum_z P_\theta(x, z)$$

$$= \log \sum_z \underbrace{\frac{Q(z)}{Q(z)}}_{\text{any surrogate dist.}} P_\theta(x, z)$$

$$= \log \mathbb{E}_Q \left[ \frac{P_\theta(x, z)}{Q(z)} \right]$$

$$\geq \mathbb{E}_Q \left[ \log \frac{P_\theta(x, z)}{Q(z)} \right] \quad (\text{Jensen})$$

$$\stackrel{\Delta}{=} B(Q, \theta) \quad (\text{evidence lower bound})$$

"ELBO"

$$B(Q, \theta) \propto_{\mu_k} \mathbb{E}_Q \left[ \log P(x_{1:T} | z_{1:T}; \mu) \right]$$

$$\propto_{\mu_k} \mathbb{E}_Q \left[ \sum_t \sum_k f_{tk} [x_t \log \mu_k - \mu_k] \right]$$

$$\propto_{\mu_k} \sum_t \sum_k \mathbb{E}_Q [f_{tk}] [x_t \log \mu_k - \mu_k]$$

$$\frac{\partial}{\partial \mu_k} \dots = 0 \rightarrow \hat{\mu}_k = \frac{\sum_t \mathbb{E}_Q [f_{tk}] x_t}{\sum_t \mathbb{E}_Q [f_{tk}]}$$

So: if we have  $Q(z)$ , we compute the beliefs

$$\mathbb{E}_Q[y_{t_k}] = \mathbb{E}_Q[\mathbb{1}(z_t = k)] = Q(z_t = k),$$

and then maximize ELBO WRT  $\mu_v$ .

Where do we get  $Q(z)$ ?

$$\hat{Q}(z) \leftarrow \underset{Q(z)}{\operatorname{argmax}} B(Q, \Theta)$$

$$B(Q, \theta) = \mathbb{E}_Q \left[ \log \frac{p(x, z)}{Q(z)} \right]$$

$$\quad \quad \quad \text{"} \quad \quad \quad \text{"} \quad \quad \quad -\log p(x)$$

$$\mathcal{L}_Q$$

$$= -\text{KL} \left( Q(z) \parallel p(z|x) \right)$$

Therefore :

$$\underset{Q(z)}{\operatorname{argmax}} \mathcal{B}(Q, \Theta) = p_{\Theta}^*(z | x)$$

(the exact posterior)

Can we set  $Q(z_{1:T}) = p(z_{1:T} | x_{1:T})$ ?  
(tractably)

No. But, do we need the whole distribution?

No! Only the beliefs, which we can compute with forwards-backwards.