

Machine Learning Methods for Neural Data Analysis

Lecture 14: Variational EM for SLDS

Scott Linderman

STATS 220/320 (*NBIO220, CS339N*). Winter 2021.

Announcements

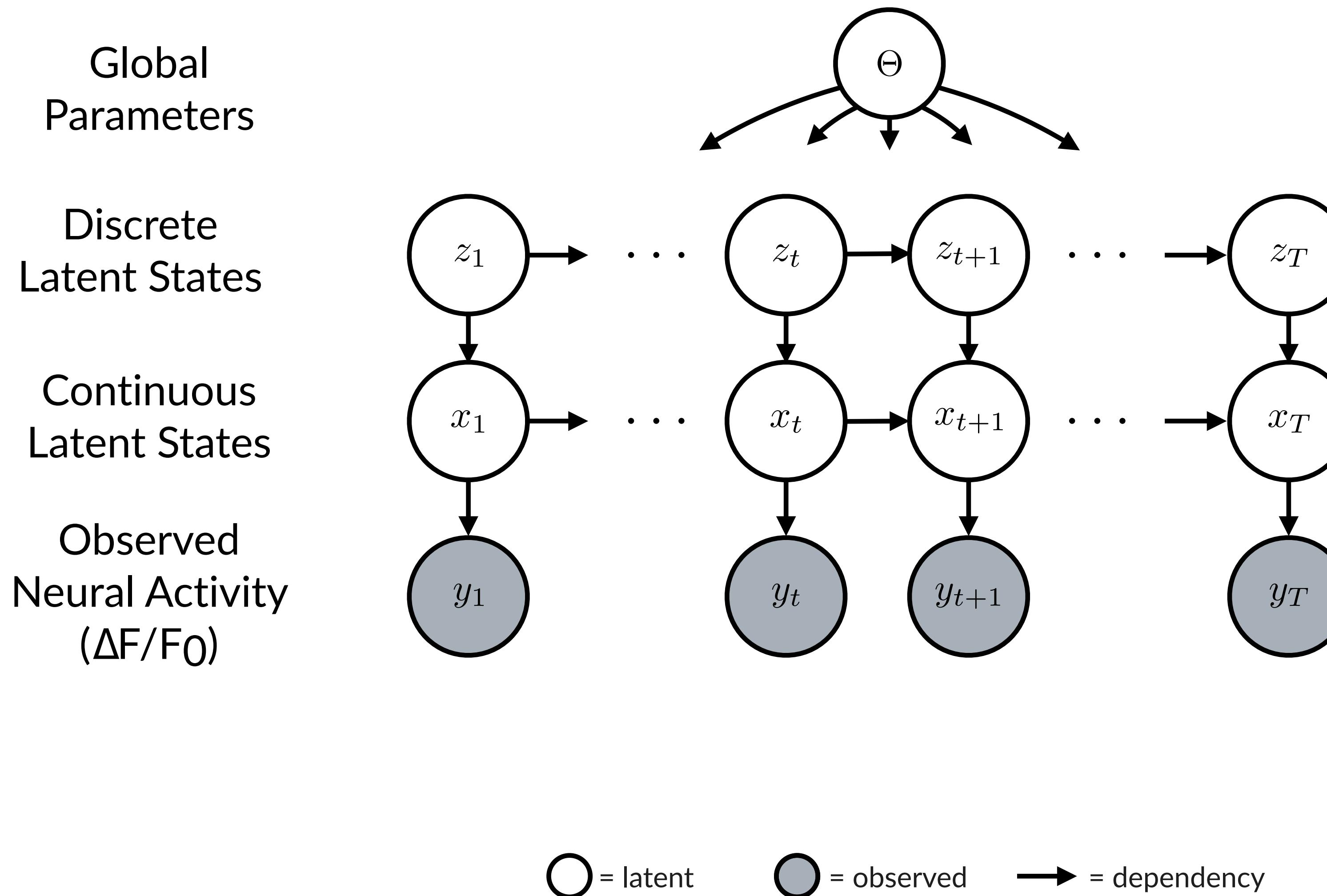
- **Updated proposals** due this Friday.
 - Looking for progress on downloading and extracting the data.

Agenda

- Variational EM
- Coordinate Ascent VI
- Variational EM for SLDS

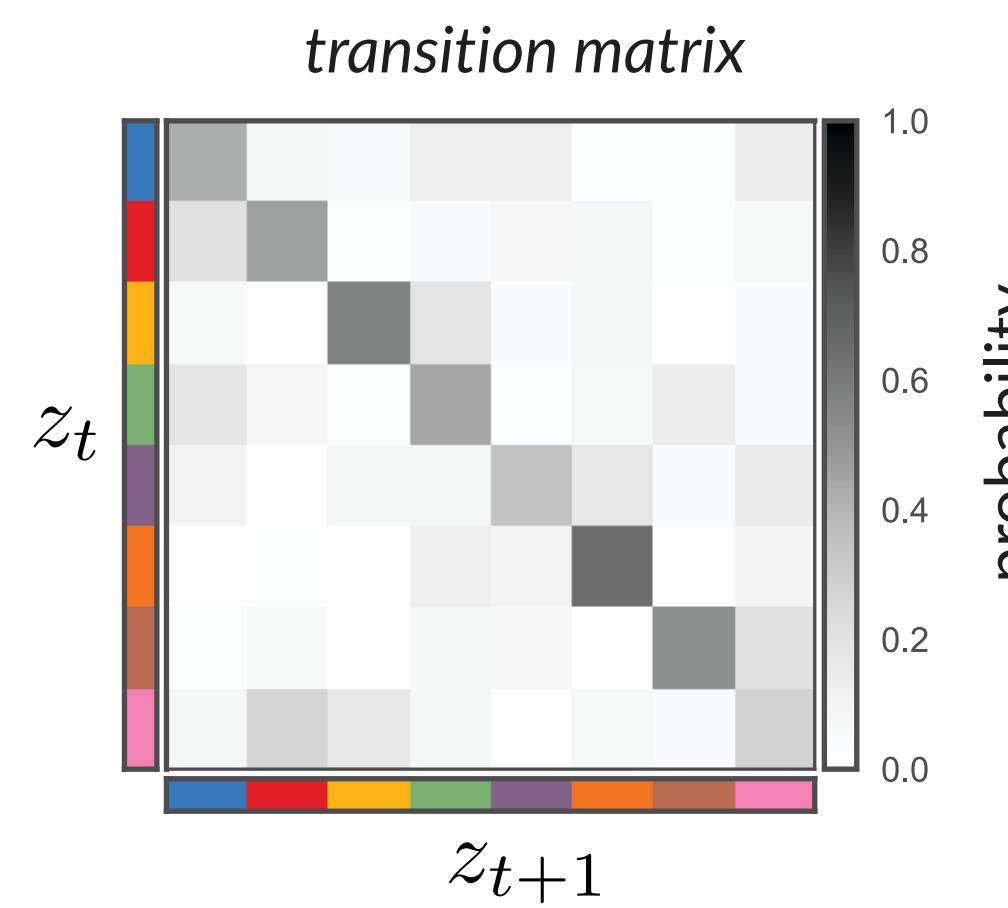
Recap

Switching Linear Dynamical Systems



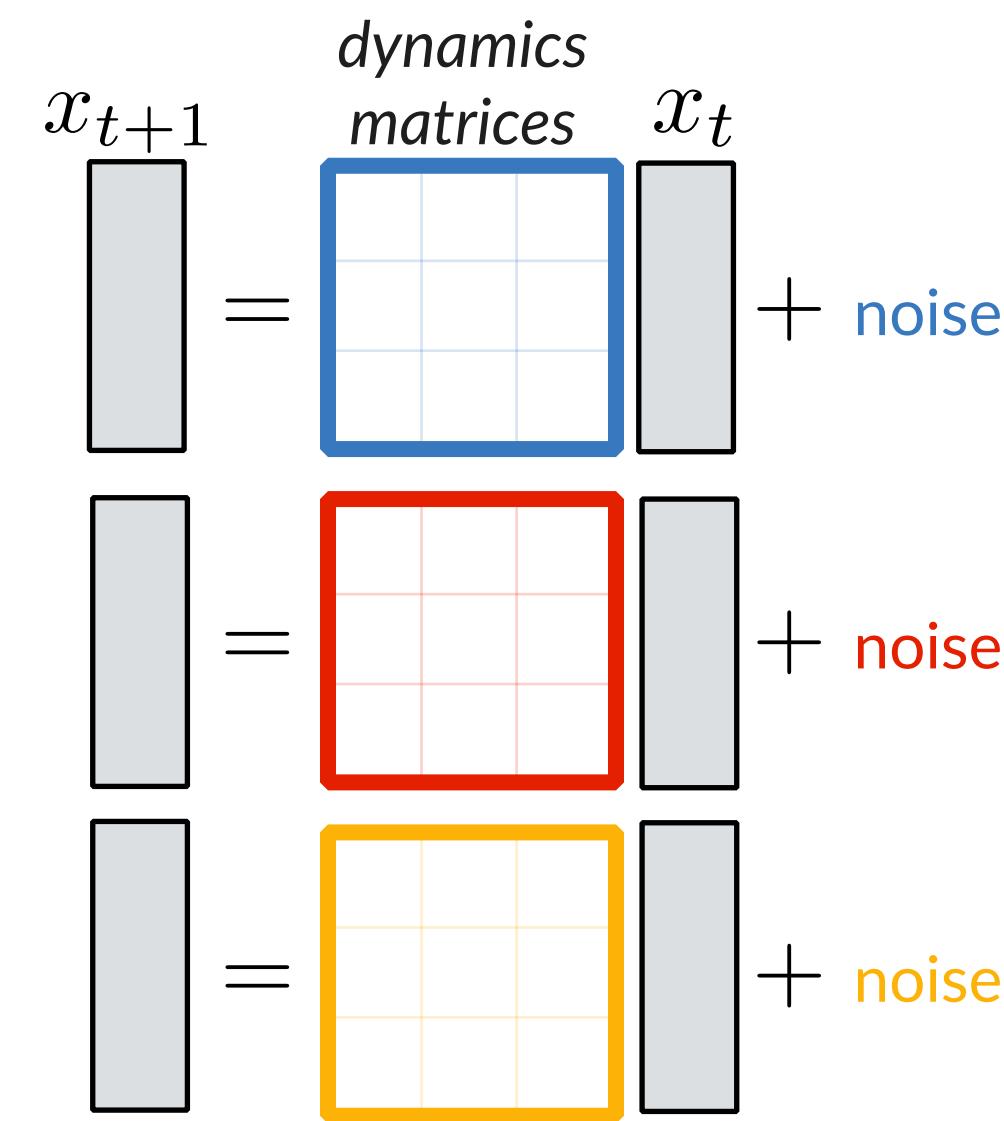
Specifying the form of the dependencies

State-dependent
switching probabilities



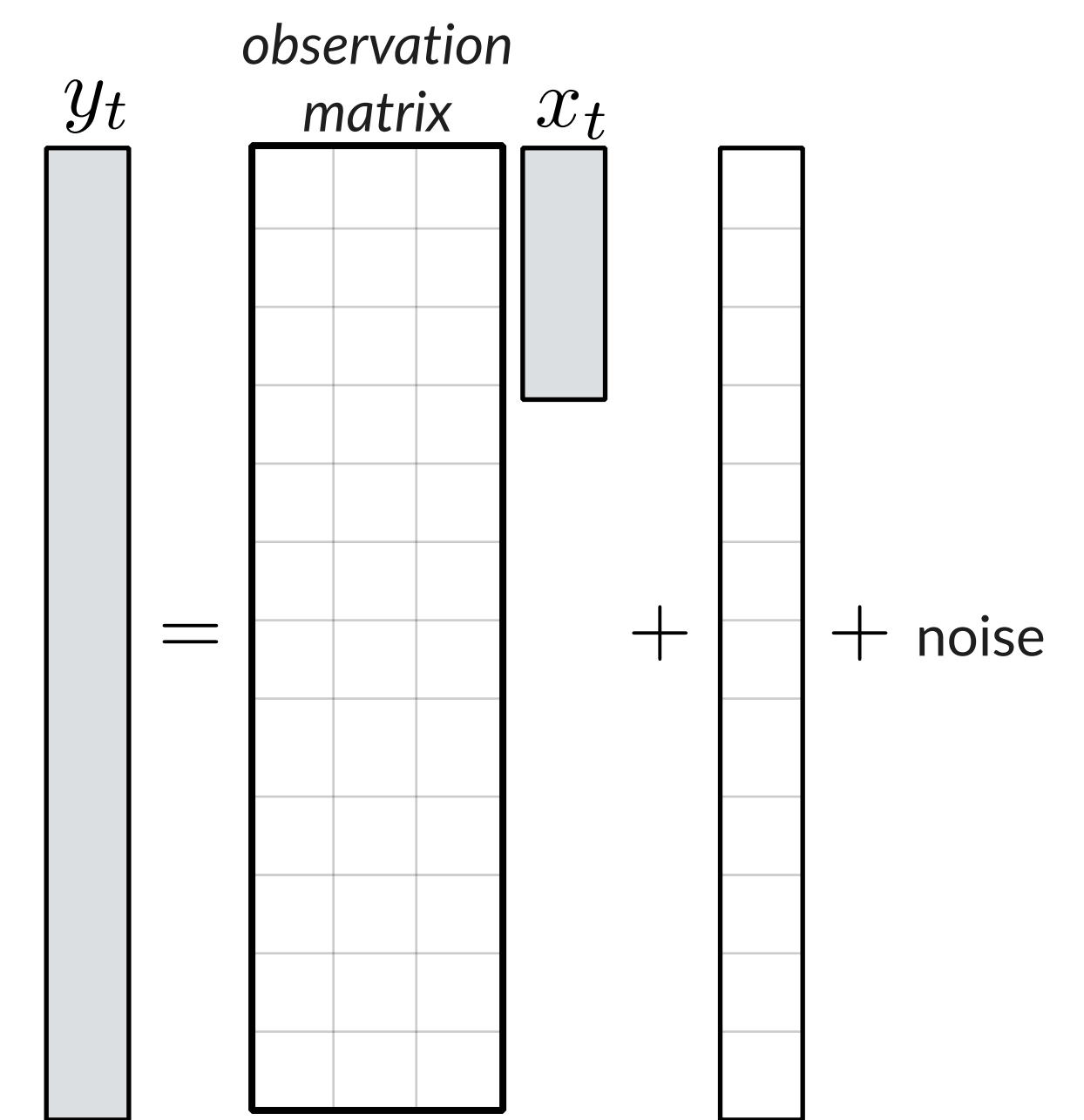
$$\Pr(z_{t+1} = j \mid z_t = i) = P_{ij}$$

Different linear dynamics
in each discrete state

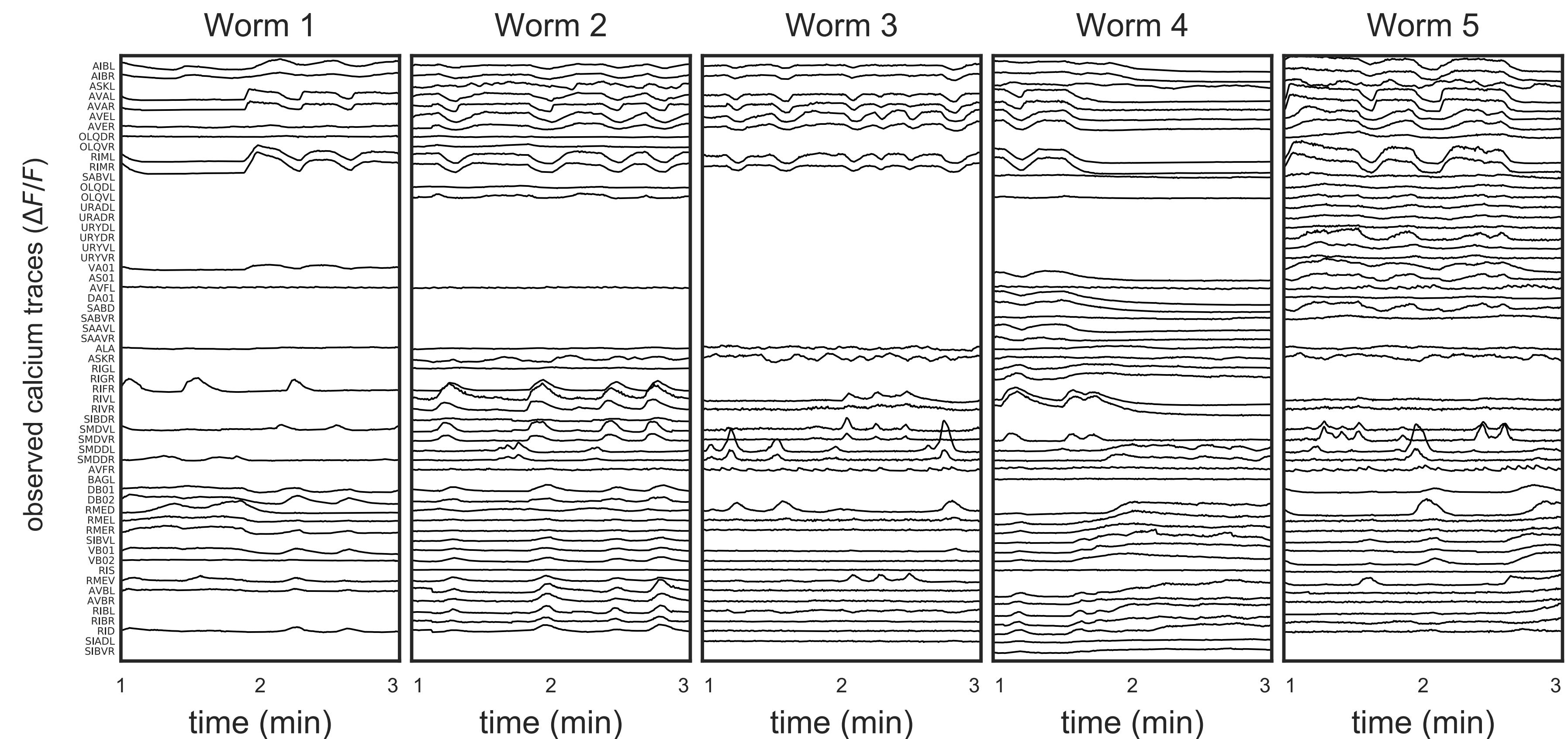
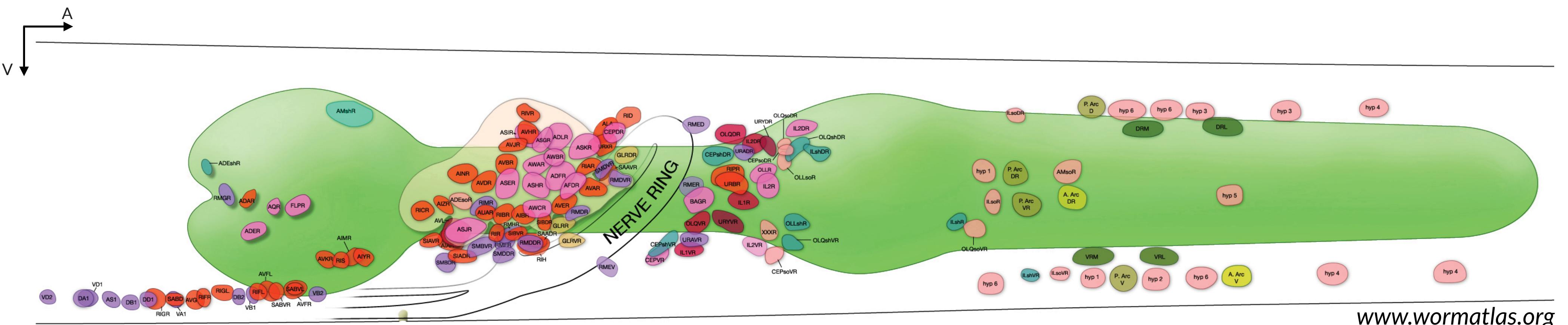


$$x_{t+1} = A_{z_{t+1}} x_t + b_{z_{t+1}} + \epsilon_{t+1}$$

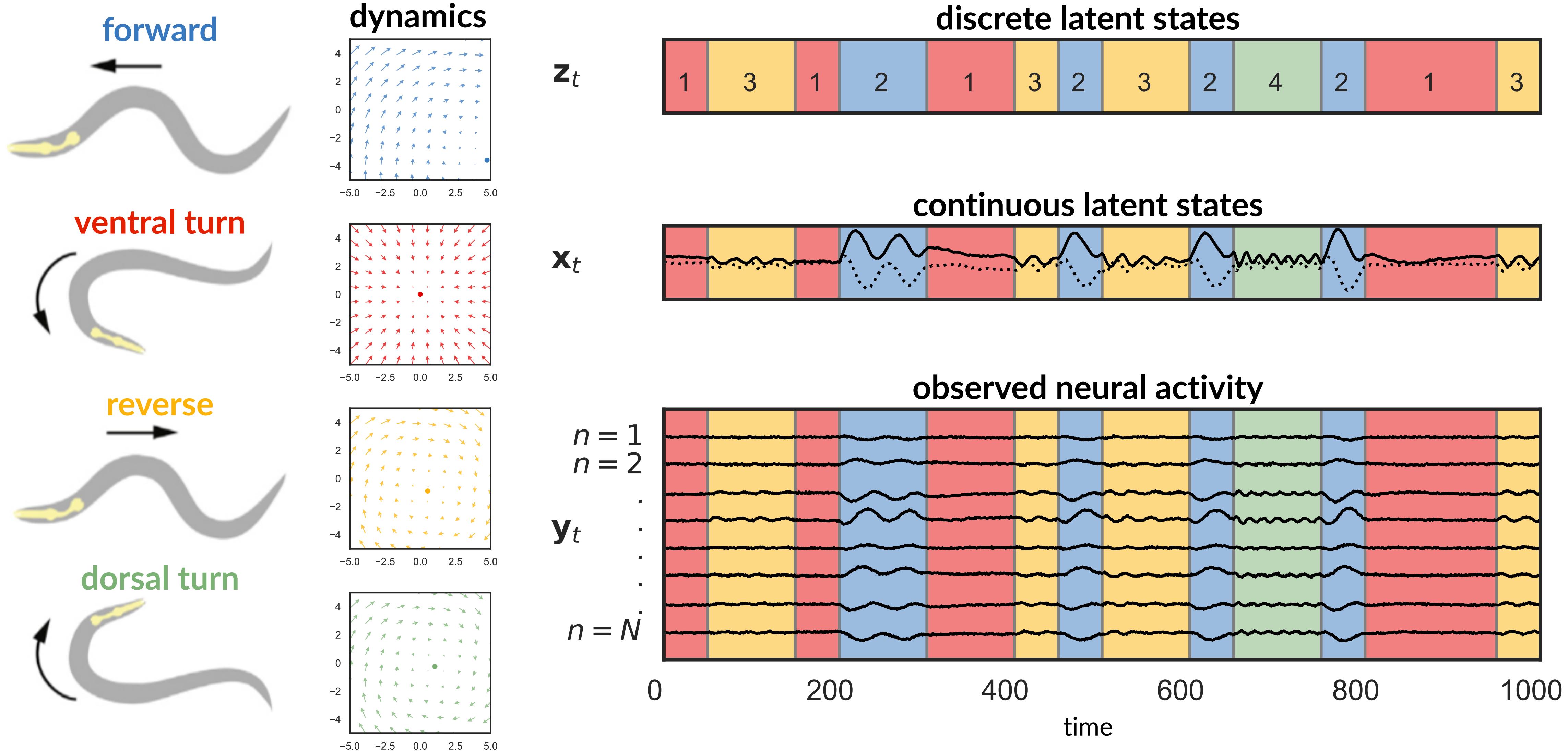
Linear mapping from continuous latent
states to observed neural activity



$$y_t = C x_t + d + \delta_t$$



Building a probabilistic model of neural data



Exact EM for the SLDS

- **E-step:** Update the posterior over latent variables,

$$q(z, x) \leftarrow p(z, x | y, \Theta) = \frac{p(z, x, y | \Theta)}{p(y | \Theta)}$$

- As before, we only need certain expectations under q ,

$$\mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t x_t^\top], \quad \mathbb{E}_{q(z,x)} [\mathbb{I}[z_t = k]x_t x_{t+1}^\top],$$

- **M-step:** Update the parameters,

$$\Theta \leftarrow \arg \max \mathbb{E}_{q(z,x)} [\log p(z, x, y | \Theta)]$$

- Unfortunately, computing the necessary expectations is a lot harder now!
- The posterior distribution is a mixture of Gaussians with K^T components.

Variational EM

Bayesian inference in latent variable models

Recall our derivation of the EM algorithm

- Goal: find parameters that maximize the **marginal likelihood** (aka the **model evidence**):

$$\log p(y \mid \Theta) = \log \int p(y, z \mid \Theta) dz$$

Bayesian inference in latent variable models

Recall our derivation of the EM algorithm

- Goal: find parameters that maximize the **marginal likelihood** (aka the **model evidence**):

$$\begin{aligned}\log p(y \mid \Theta) &= \log \int p(y, z \mid \Theta) dz \\ &= \log \int \frac{q(z)}{q(z)} p(y, z \mid \Theta) dz \quad \text{for any distribution } q(z)\end{aligned}$$

Bayesian inference in latent variable models

Recall our derivation of the EM algorithm

- Goal: find parameters that maximize the **marginal likelihood** (aka the **model evidence**):

$$\begin{aligned}\log p(y \mid \Theta) &= \log \int p(y, z \mid \Theta) dz \\ &= \log \int \frac{q(z)}{q(z)} p(y, z \mid \Theta) dz && \text{for any distribution } q(z) \\ &= \log \mathbb{E}_{q(z)} \left[\frac{p(y, z \mid \Theta)}{q(z)} \right]\end{aligned}$$

Bayesian inference in latent variable models

Recall our derivation of the EM algorithm

- Goal: find parameters that maximize the **marginal likelihood** (aka the **model evidence**):

$$\begin{aligned}\log p(y \mid \Theta) &= \log \int p(y, z \mid \Theta) dz \\ &= \log \int \frac{q(z)}{q(z)} p(y, z \mid \Theta) dz && \text{for any distribution } q(z) \\ &= \log \mathbb{E}_{q(z)} \left[\frac{p(y, z \mid \Theta)}{q(z)} \right] \\ &\geq \mathbb{E}_{q(z)} [\log p(y, z \mid \Theta) - \log q(z)] && \text{by Jensen's inequality}\end{aligned}$$

Bayesian inference in latent variable models

Recall our derivation of the EM algorithm

- Goal: find parameters that maximize the **marginal likelihood** (aka the **model evidence**):

$$\begin{aligned}\log p(y \mid \Theta) &= \log \int p(y, z \mid \Theta) dz \\ &= \log \int \frac{q(z)}{q(z)} p(y, z \mid \Theta) dz && \text{for any distribution } q(z) \\ &= \log \mathbb{E}_{q(z)} \left[\frac{p(y, z \mid \Theta)}{q(z)} \right] \\ &\geq \mathbb{E}_{q(z)} [\log p(y, z \mid \Theta) - \log q(z)] && \text{by Jensen's inequality} \\ &\triangleq \mathcal{L}[q, \Theta]\end{aligned}$$

- \mathcal{L} is called the **evidence lower bound** or the **ELBO** for short.

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$q \leftarrow \arg \max_q \mathcal{L}[q, \Theta]$$

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$\begin{aligned} q &\leftarrow \arg \max_q \mathcal{L}[q, \Theta] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(y, z | \Theta)}{q(z)} \right] \end{aligned}$$

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$\begin{aligned} q &\leftarrow \arg \max_q \mathcal{L}[q, \Theta] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(y, z | \Theta)}{q(z)} \right] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(z | y | \Theta)}{q(z)} \right] + \log p(y | \Theta) \end{aligned}$$

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$\begin{aligned} q &\leftarrow \arg \max_q \mathcal{L}[q, \Theta] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(y, z | \Theta)}{q(z)} \right] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(z | y | \Theta)}{q(z)} \right] + \log p(y | \Theta) \\ &= \arg \max_q -\text{KL} (q(z) \| p(z | y, \Theta)) + \log p(y | \Theta) \end{aligned}$$

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$\begin{aligned} q &\leftarrow \arg \max_q \mathcal{L}[q, \Theta] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(y, z | \Theta)}{q(z)} \right] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(z | y | \Theta)}{q(z)} \right] + \log p(y | \Theta) \\ &= \arg \max_q -\text{KL}(q(z) \| p(z | y, \Theta)) + \log p(y | \Theta) \\ &= \arg \min_q \text{KL}(q(z) \| p(z | y, \Theta)) \end{aligned}$$

Bayesian inference in latent variable models

Coordinate ascent on the ELBO

- **E Step:** Update the posterior distribution on latent variables,

$$\begin{aligned} q &\leftarrow \arg \max_q \mathcal{L}[q, \Theta] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(y, z | \Theta)}{q(z)} \right] \\ &= \arg \max_q \mathbb{E}_{q(z)} \left[\log \frac{p(z | y | \Theta)}{q(z)} \right] + \log p(y | \Theta) \\ &= \arg \max_q -\text{KL}(q(z) \| p(z | y, \Theta)) + \log p(y | \Theta) \\ &= \arg \min_q \text{KL}(q(z) \| p(z | y, \Theta)) \end{aligned}$$

- **Maximizing the ELBO** w.r.t. q is equivalent to **minimizing the Kullback-Leibler (KL) divergence**.
- The KL divergence is **non-negative**, and it equals zero iff $q(z) \equiv p(z | y, \Theta)$.

Bayesian inference in latent variable models

The Expectation-Maximization (EM) algorithm

- **M-step:** Maximize the expected log probability

$$\Theta \leftarrow = \arg \max_{\Theta} \mathbb{E}_{q(z)}[\log p(y, z, \Theta)]$$

- **E-step:** Update the posterior over latent variables

$$q \leftarrow \arg \max_q \mathcal{L}[q, \Theta]$$

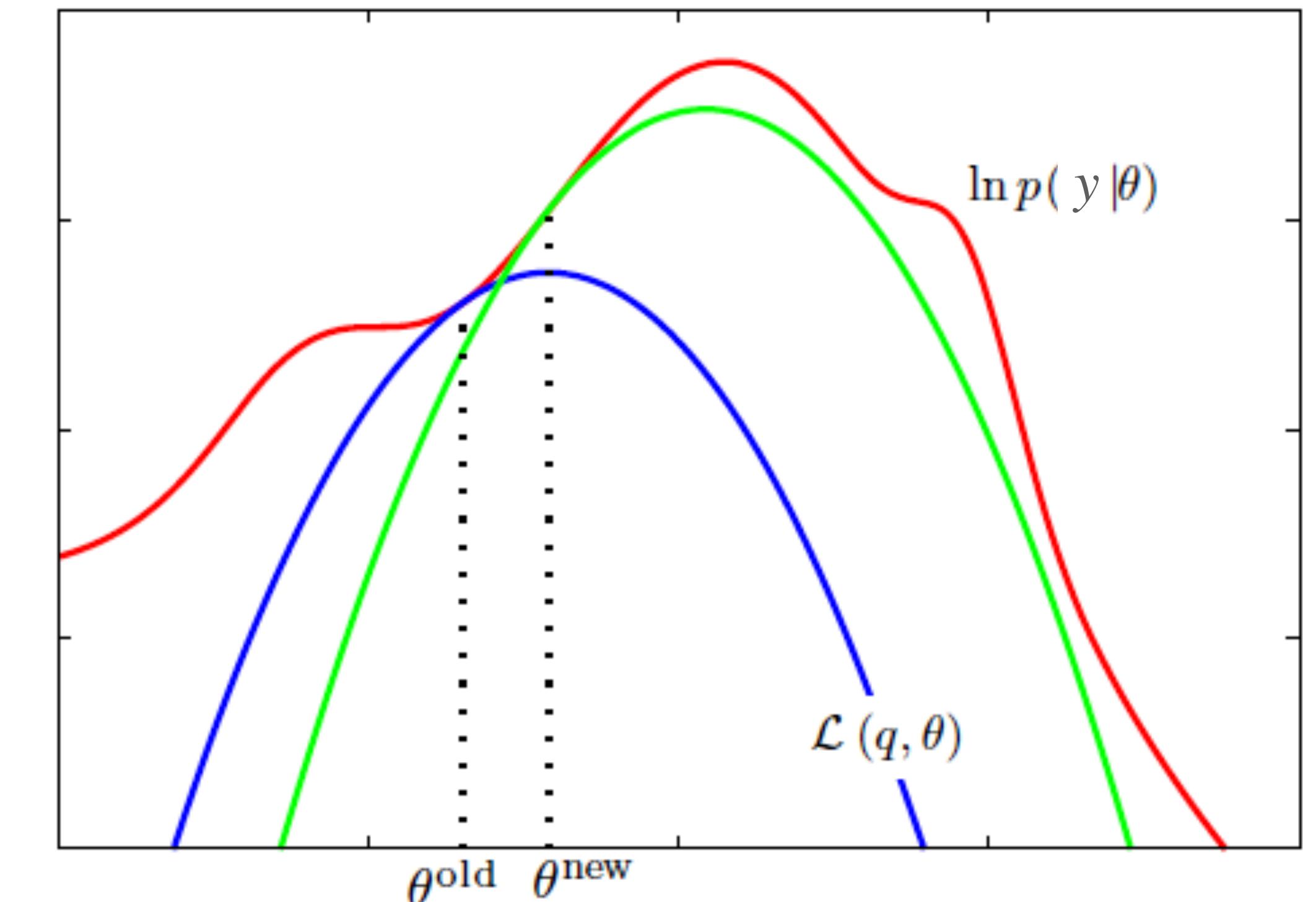
$$= \arg \min_q \text{KL}(q(z) \| p(z | y, \Theta))$$

$$= p(z | y, \Theta)$$

- After each E-step, the **ELBO is tight:**

$$\mathcal{L}[p(z | y, \Theta), \Theta] = \log p(y | \Theta)$$

- EM converges to **local optima** of the marginal distribution.



Bayesian inference in latent variable models

Variational Expectation-Maximization

- **M-step:** Maximize the expected log probability

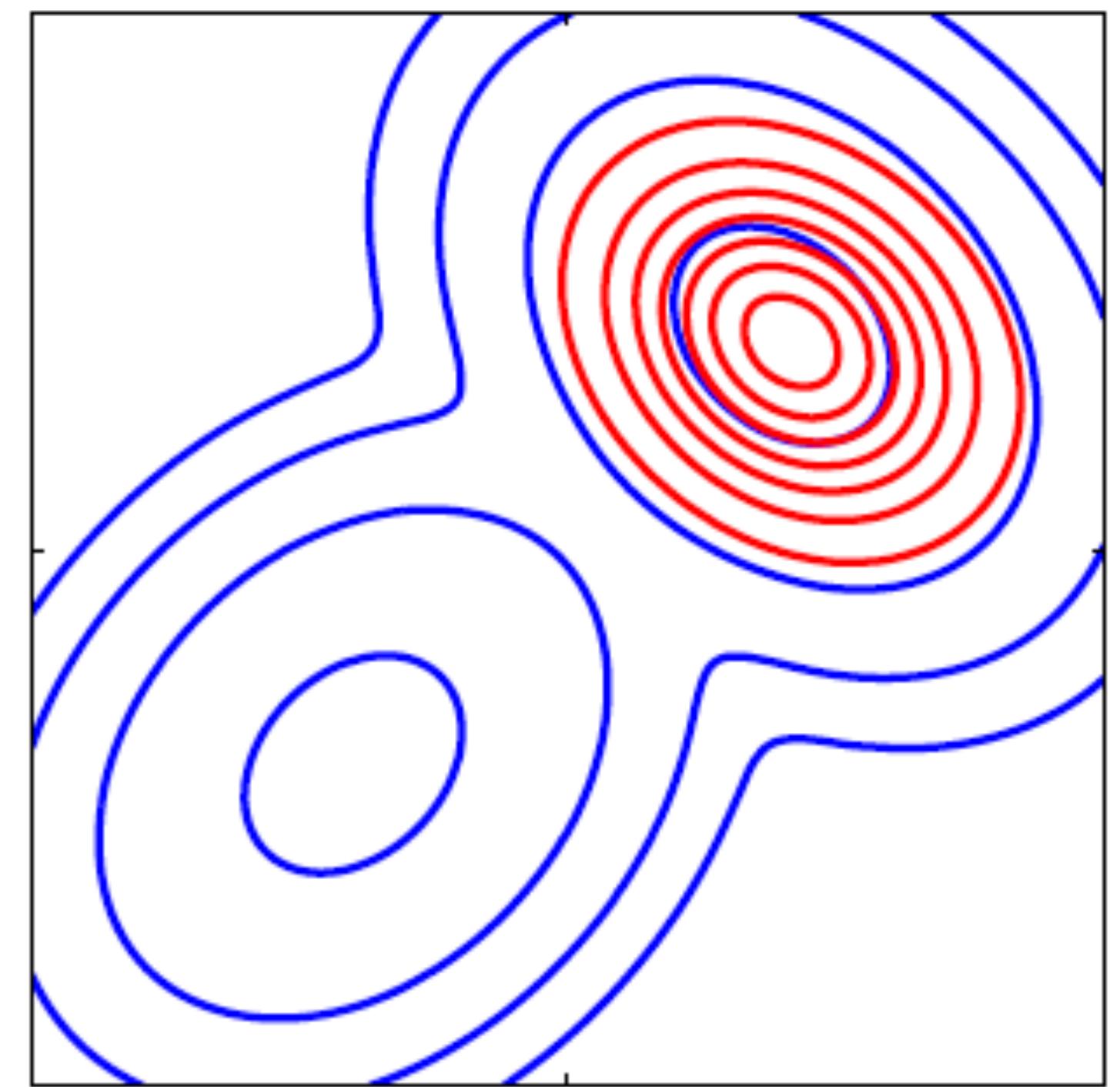
$$\Theta \leftarrow = \arg \max_{\Theta} \mathbb{E}_{q(z)}[\log p(y, z, \Theta)]$$

- **Variational E-step:** Update the posterior, subject to $q \in \mathcal{Q}$

$$\begin{aligned} q &\leftarrow \arg \max_{q \in \mathcal{Q}} \mathcal{L}[q, \Theta] \\ &= \arg \min_{q \in \mathcal{Q}} \text{KL}(q(z) \| p(z | y, \Theta)) \end{aligned}$$

where \mathcal{Q} is a set of **tractable approximate posteriors**.

- For example, \mathcal{Q} could **assume independence** or a **particular functional form**.
- If \mathcal{Q} does not contain the true posterior, the ELBO will be a strict lower bound on the marginal likelihood, $\mathcal{L}[q(z), \Theta] < \log p(y | \Theta)$.
- Optimizing over \mathcal{Q} to find the best approximation is called **variational inference**, hence the name Variational EM.



Coordinate Ascent Variational Inference (CAVI)

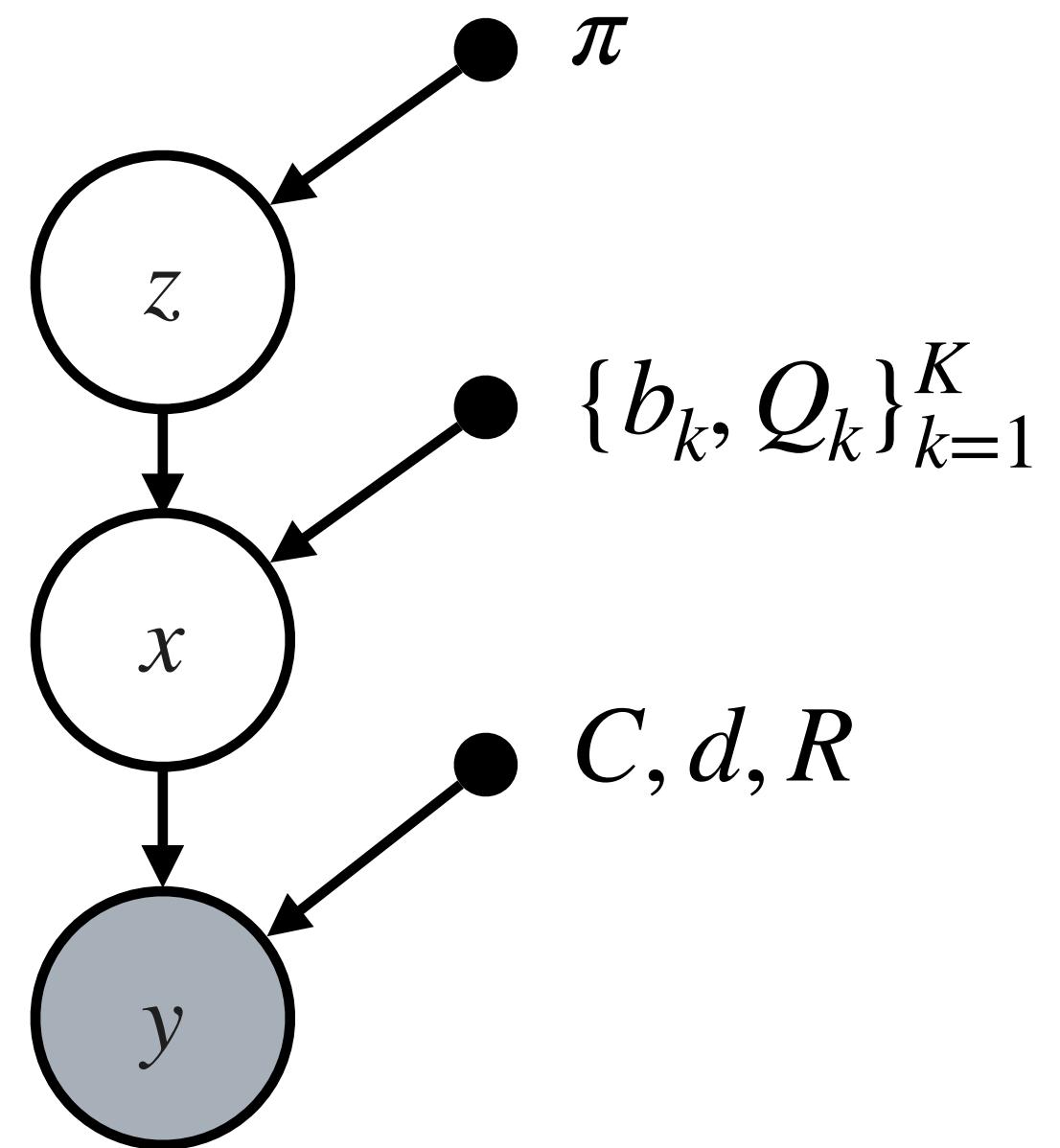
Coordinate ascent VI

Warm-up example

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Mean field variational family

- **Assume** \mathcal{Q} is the set of “factored” distributions

$$q(z, x) = q(z) q(x)$$

where z and x are independent.

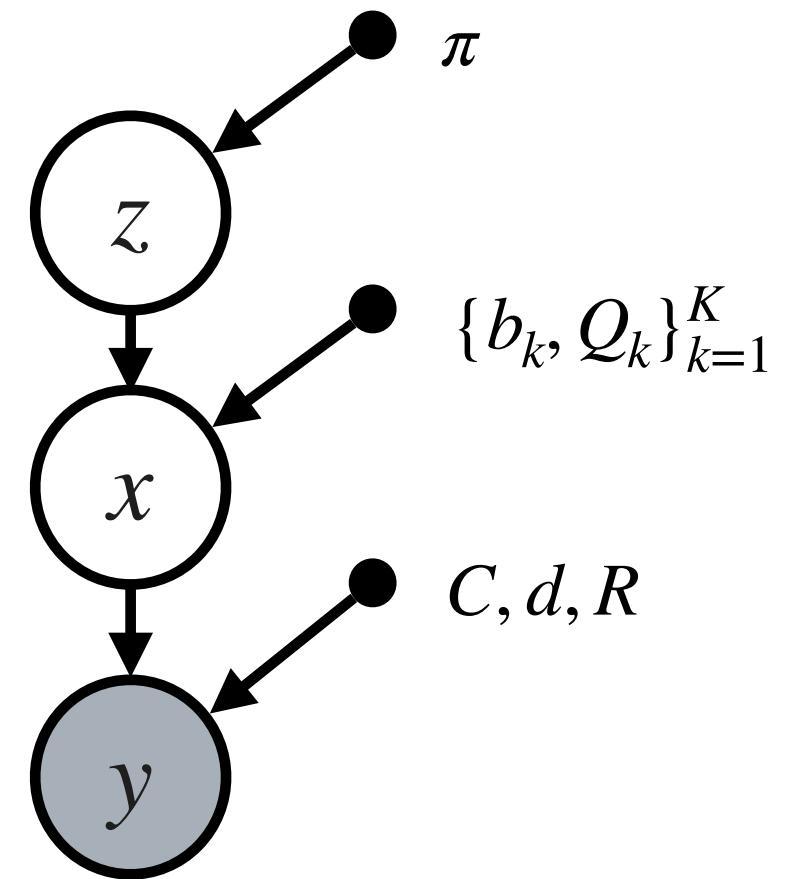
- This is called the **mean field family**.
- We want to find,

$$\arg \max_{q \in \mathcal{Q}} \mathcal{L}[q(z)q(x), \Theta] \equiv \arg \min_{q \in \mathcal{Q}} \text{KL} (q(z)q(x) \parallel p(z, x \mid y, \Theta))$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Coordinate updates

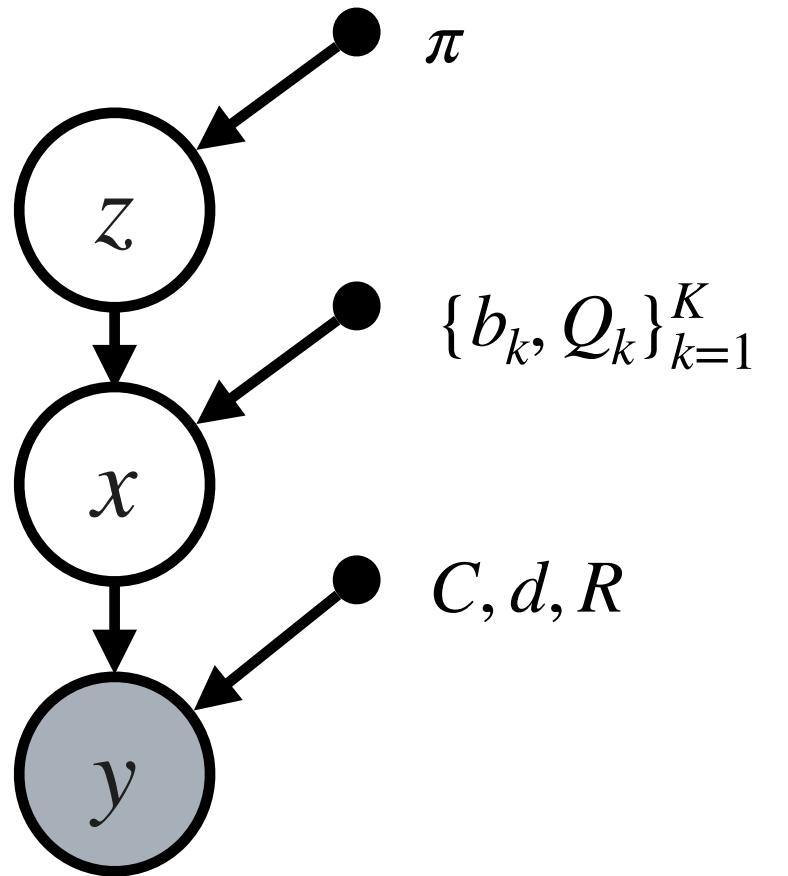
Hold $q(x)$ fixed and optimize w.r.t. $q(z)$:

$$\mathcal{L}[q(z)q(x), \Theta] = \mathbb{E}_{q(z)q(x)} [\log p(z, x, y | \Theta) - \log q(z) - \log q(x)]$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Coordinate updates

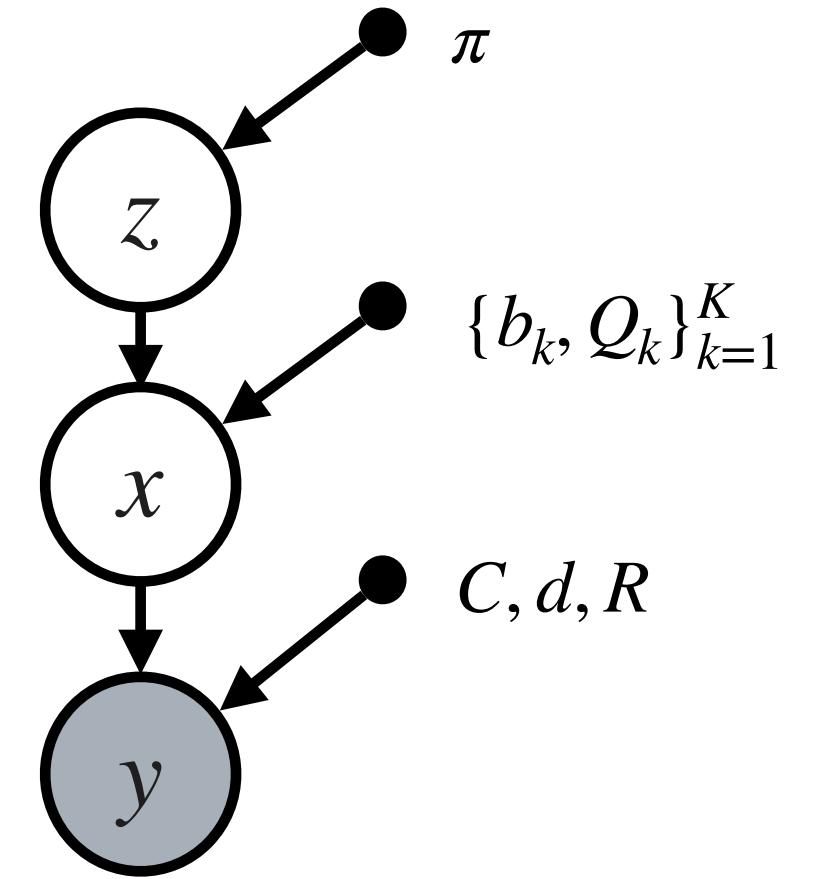
Hold $q(x)$ fixed and optimize w.r.t. $q(z)$:

$$\begin{aligned}\mathcal{L}[q(z)q(x), \Theta] &= \mathbb{E}_{q(z)q(x)} [\log p(z, x, y | \Theta) - \log q(z) - \log q(x)] \\ &= \mathbb{E}_{q(z)} \left[\mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] - \log q(z) \right] + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Coordinate updates

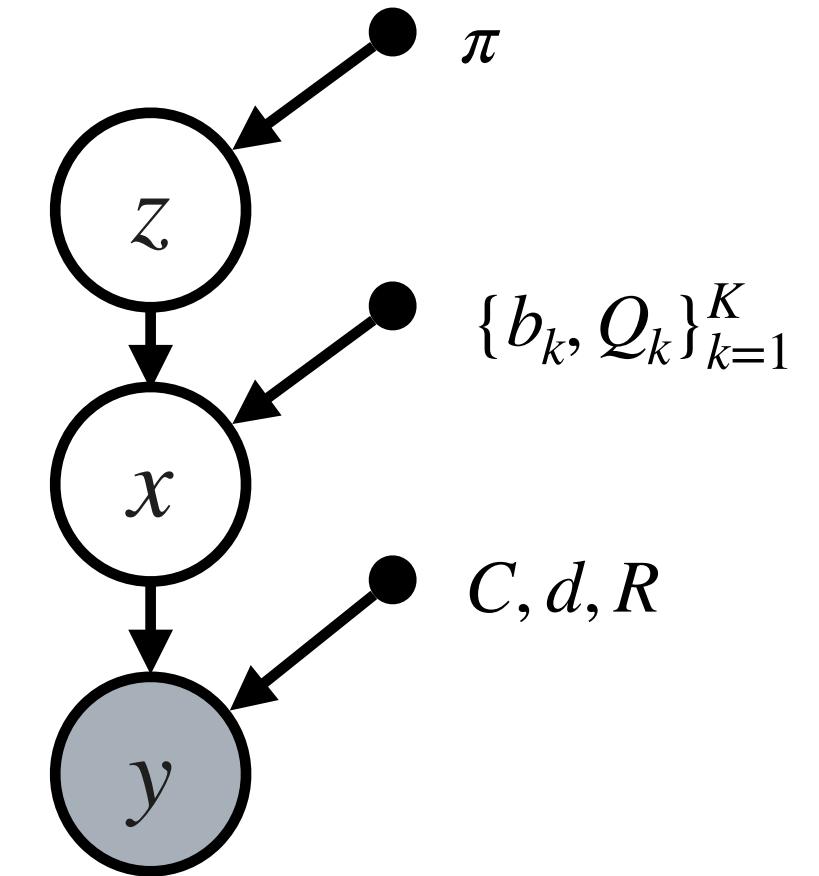
Hold $q(x)$ fixed and optimize w.r.t. $q(z)$:

$$\begin{aligned}\mathcal{L}[q(z)q(x), \Theta] &= \mathbb{E}_{q(z)q(x)} [\log p(z, x, y | \Theta) - \log q(z) - \log q(x)] \\ &= \mathbb{E}_{q(z)} \left[\mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] - \log q(z) \right] + c \\ &= \mathbb{E}_{q(z)} \left[\log \exp \left\{ \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] \right\} - \log q(z) \right] + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Coordinate updates

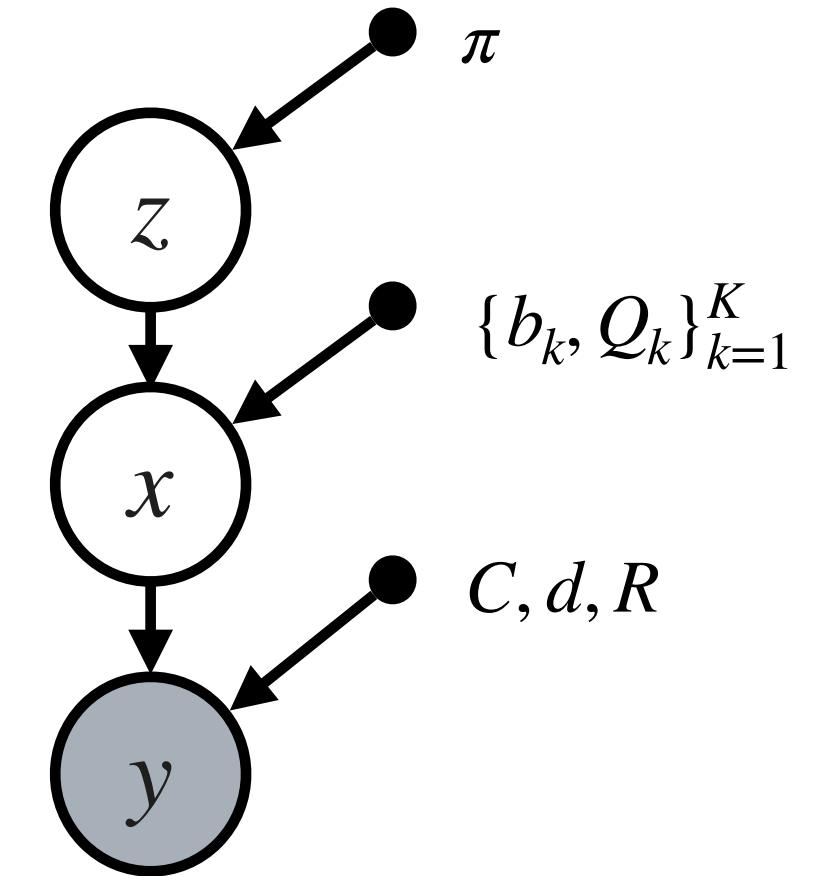
Hold $q(x)$ fixed and optimize w.r.t. $q(z)$:

$$\begin{aligned}\mathcal{L}[q(z)q(x), \Theta] &= \mathbb{E}_{q(z)q(x)} [\log p(z, x, y | \Theta) - \log q(z) - \log q(x)] \\ &= \mathbb{E}_{q(z)} \left[\mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] - \log q(z) \right] + c \\ &= \mathbb{E}_{q(z)} \left[\log \exp \left\{ \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] \right\} - \log q(z) \right] + c \\ &= \mathbb{E}_{q(z)} [\log \tilde{p}(z) - \log q(z)] + c'\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Coordinate updates

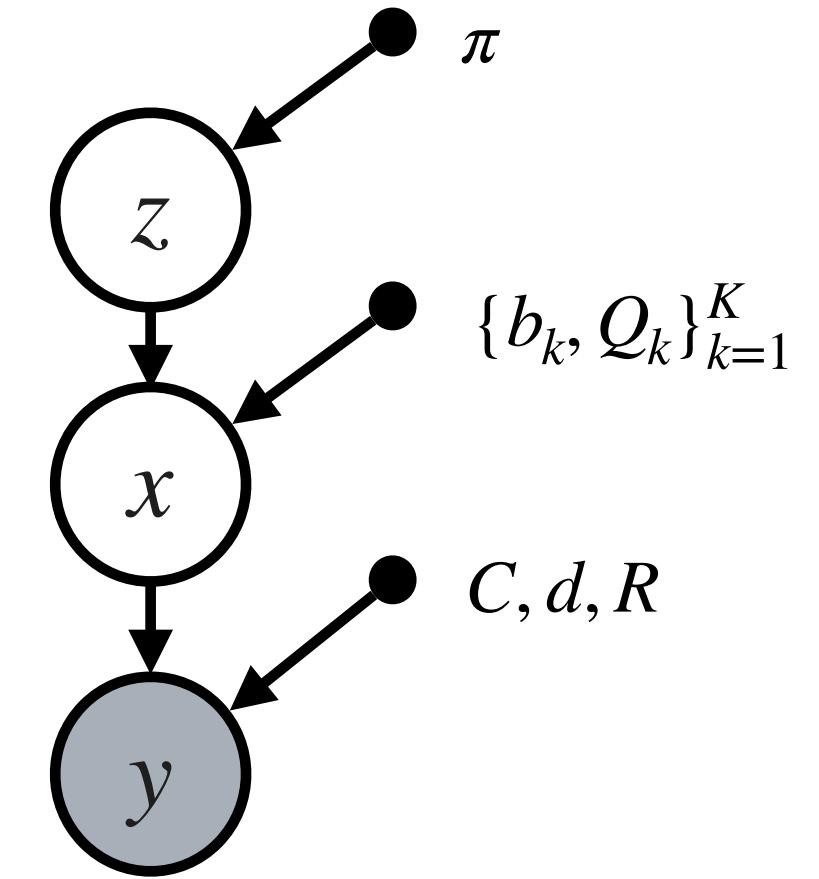
Hold $q(x)$ fixed and optimize w.r.t. $q(z)$:

$$\begin{aligned}\mathcal{L}[q(z)q(x), \Theta] &= \mathbb{E}_{q(z)q(x)} [\log p(z, x, y | \Theta) - \log q(z) - \log q(x)] \\ &= \mathbb{E}_{q(z)} \left[\mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] - \log q(z) \right] + c \\ &= \mathbb{E}_{q(z)} \left[\log \exp \left\{ \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] \right\} - \log q(z) \right] + c \\ &= \mathbb{E}_{q(z)} [\log \tilde{p}(z) - \log q(z)] + c' \\ &= - \text{KL} (q(z) \| \tilde{p}(z)) + c'\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

General form

- Thus, as a function of $q(z)$:

$$\mathcal{L}[q(z)q(x), \Theta] = -\text{KL} \left(q(z) \parallel \tilde{p}(z) \right) + c'$$

- This is minimized when

$$q(z) = \tilde{p}(z) \propto \exp \left\{ \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] \right\}$$

- By symmetry, the optimal update for $q(x)$ is

$$q(x) = \tilde{p}(x) \propto \exp \left\{ \mathbb{E}_{q(z)} [\log p(z, x, y | \Theta)] \right\}$$

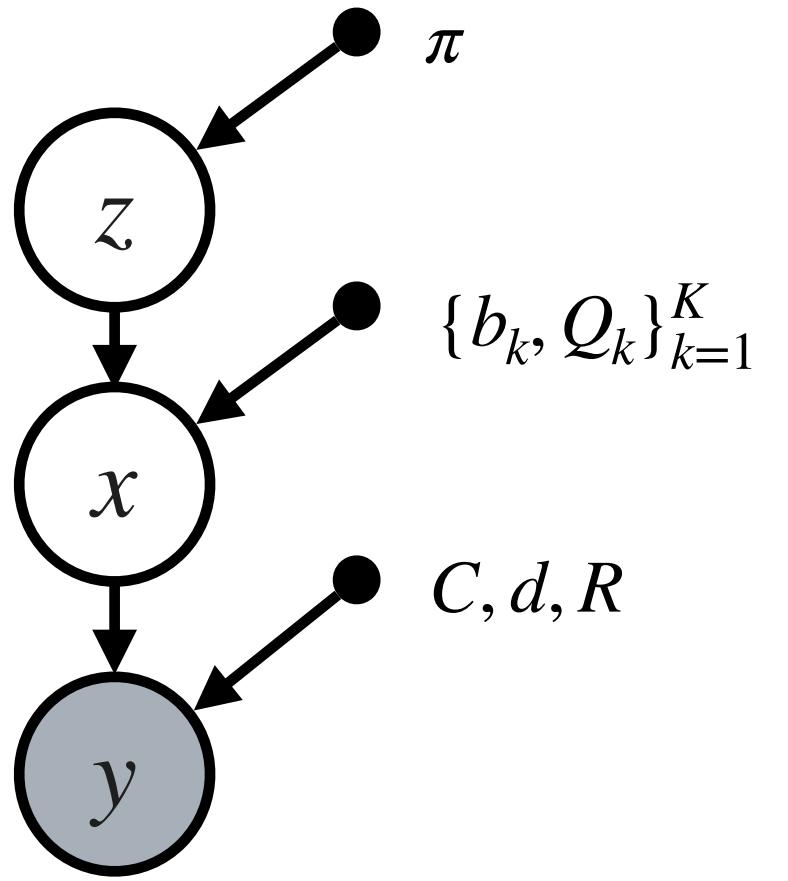
- In general, coordinate ascent for mean field variational families takes the form,

$$q(x_i) \propto \exp \left\{ \mathbb{E}_{q(x_{-i})} [\log p(x_1, \dots, x_i, \dots, x_D, y | \Theta)] \right\}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(z)$

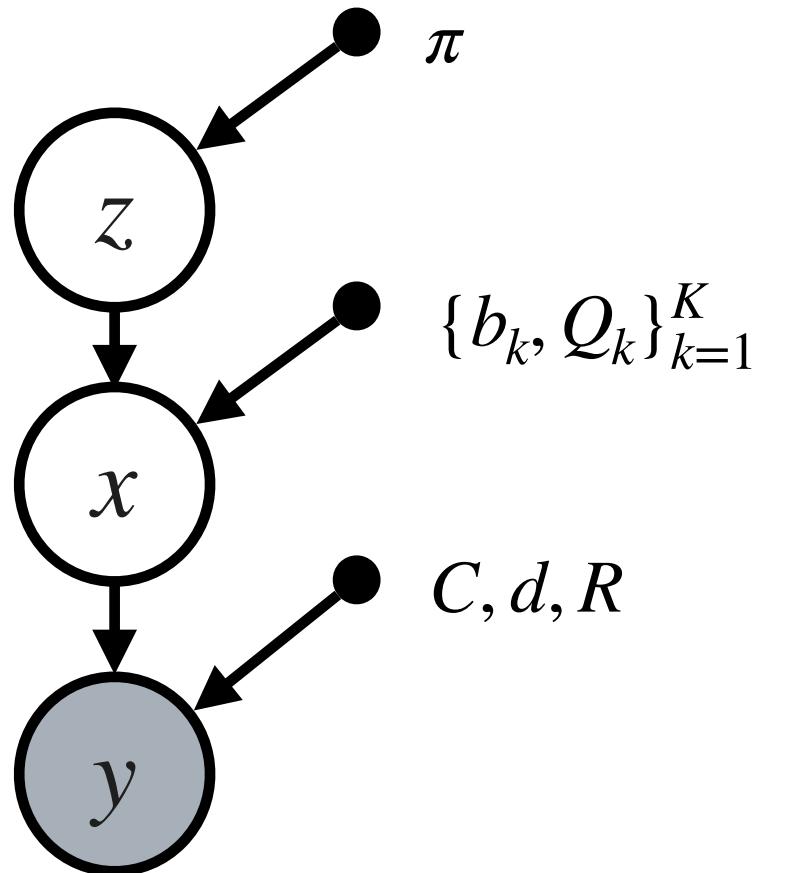
The optimal updates are often available in closed form,

$$\log q(z) = \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] + c$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(z)$

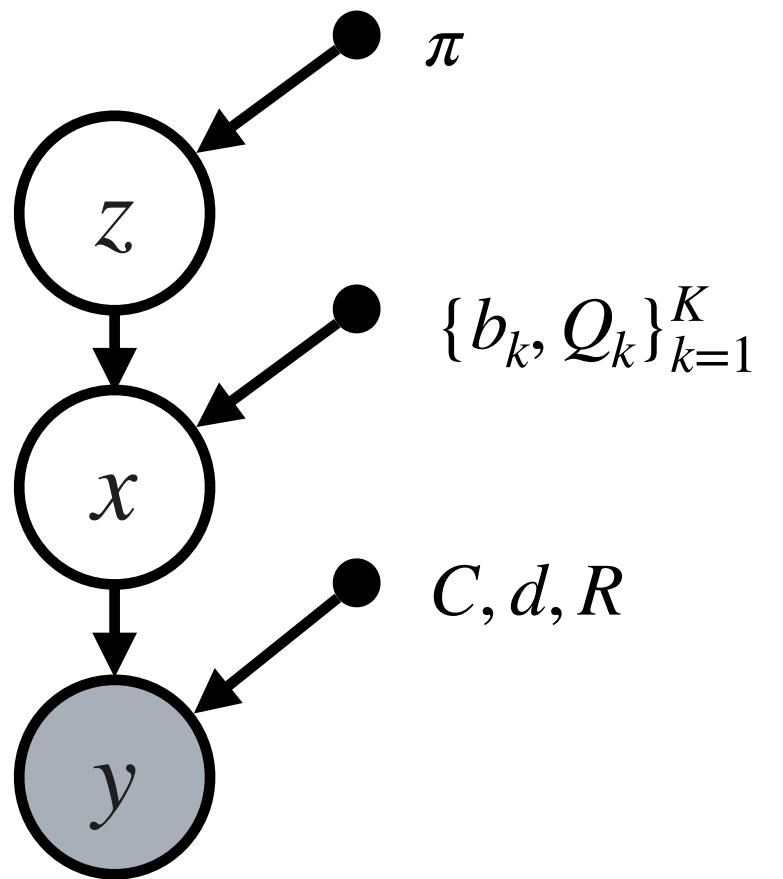
The optimal updates are often available in closed form,

$$\begin{aligned}\log q(z) &= \mathbb{E}_{q(x)} [\log p(z, x, y | \Theta)] + c \\ &= \mathbb{E}_{q(x)} [\log p(z | \Theta) + \log p(x | z, \Theta) + \log p(y | x, \Theta)] + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x | z \sim \mathcal{N}(b_z, Q_z)$$

$$y | x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(z)$

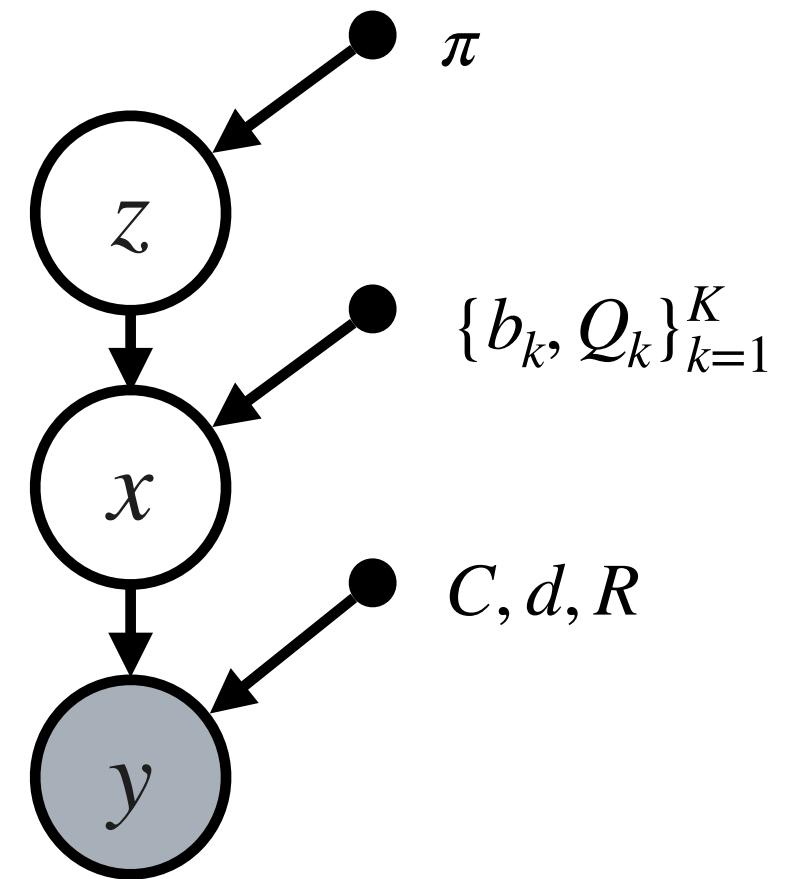
The optimal updates are often available in closed form,

$$\begin{aligned}\log q(z) &= \mathbb{E}_{q(x)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(x)} [\log p(z \mid \Theta) + \log p(x \mid z, \Theta) + \log p(y \mid x, \Theta)] + c \\ &= \log \text{Cat}(z \mid \pi) + \mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_z, Q_z)] + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(z)$

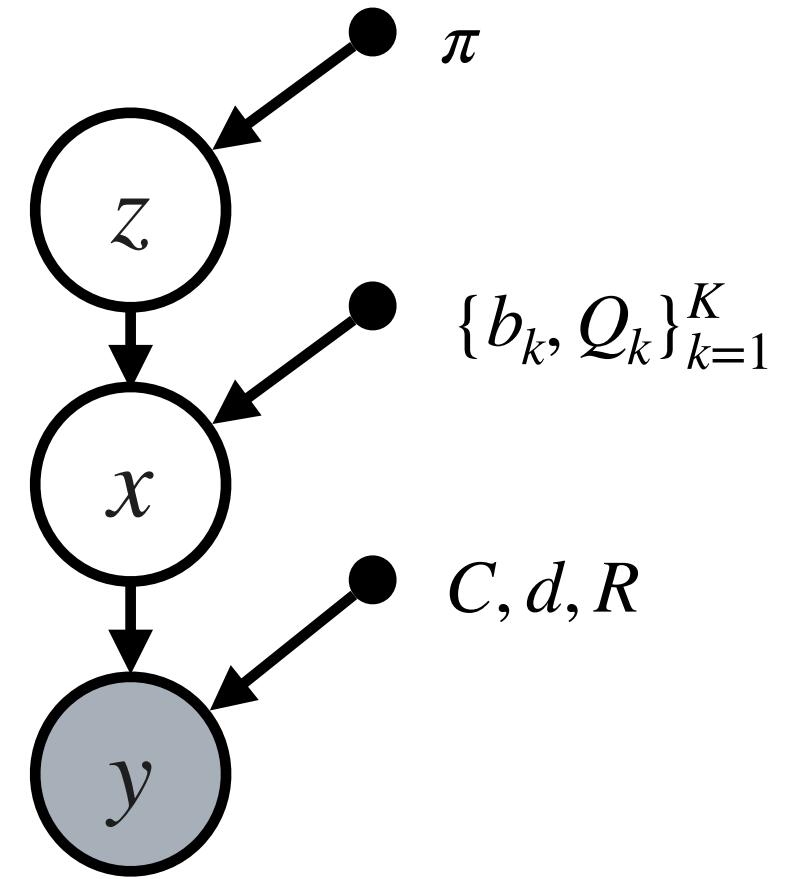
The optimal updates are often available in closed form,

$$\begin{aligned}\log q(z) &= \mathbb{E}_{q(x)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(x)} [\log p(z \mid \Theta) + \log p(x \mid z, \Theta) + \log p(y \mid x, \Theta)] + c \\ &= \log \text{Cat}(z \mid \pi) + \mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_z, Q_z)] + c \\ &= \sum_{k=1}^K \mathbb{I}[z = k] \left(\log \pi_k + \mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_k, Q_k)] \right) + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(z)$

The optimal updates are often available in closed form,

$$\begin{aligned}
 \log q(z) &= \mathbb{E}_{q(x)} [\log p(z, x, y \mid \Theta)] + c \\
 &= \mathbb{E}_{q(x)} [\log p(z \mid \Theta) + \log p(x \mid z, \Theta) + \log p(y \mid x, \Theta)] + c \\
 &= \log \text{Cat}(z \mid \pi) + \mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_z, Q_z)] + c \\
 &= \sum_{k=1}^K \mathbb{I}[z = k] \left(\log \pi_k + \mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_k, Q_k)] \right) + c \\
 &= \log \text{Cat}(z \mid \tilde{\pi})
 \end{aligned}$$

where

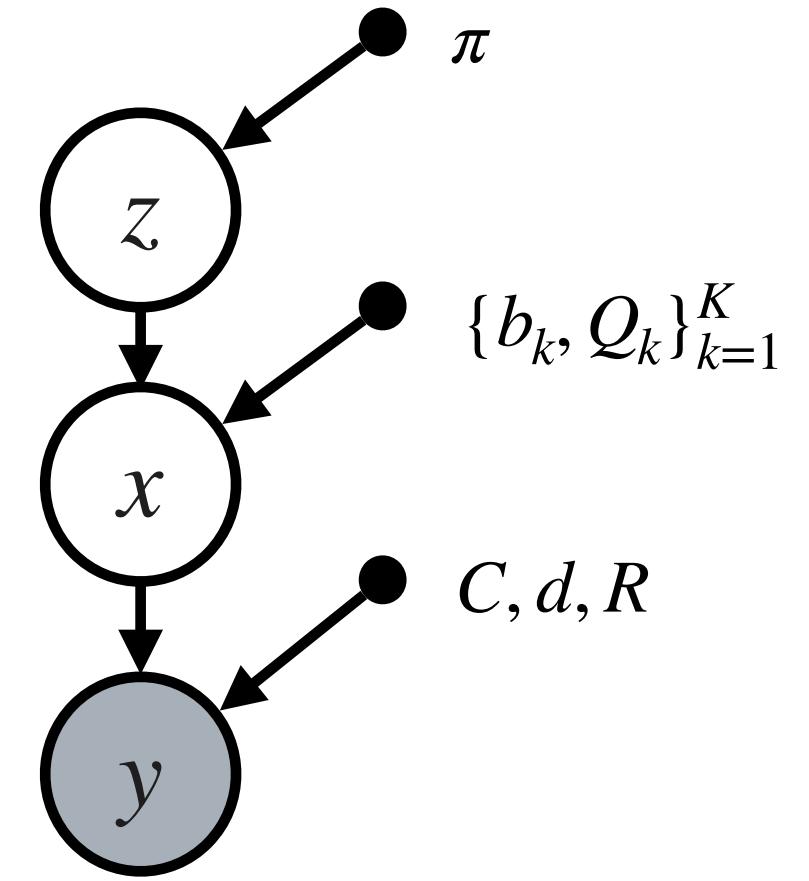
$$\log \tilde{\pi}_k = \log \pi_k + \underbrace{\mathbb{E}_{q(x)} [\log \mathcal{N}(x \mid b_k, Q_k)]}_{\text{"expected log likelihood"}} + c$$

The **expected log likelihood** is also called the **cross entropy** between $q(x)$ and $p(x \mid z)$.

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

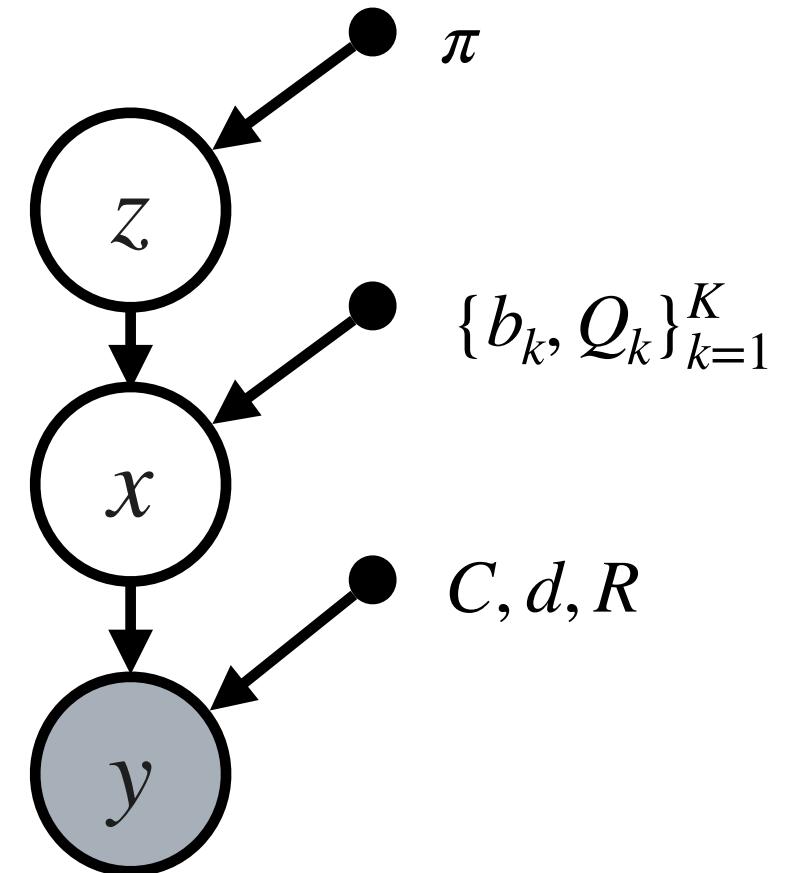
Now do the same for $q(x)$:

$$\log q(x) = \mathbb{E}_{q(z)} [\log p(z, x, y \mid \Theta)] + c$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

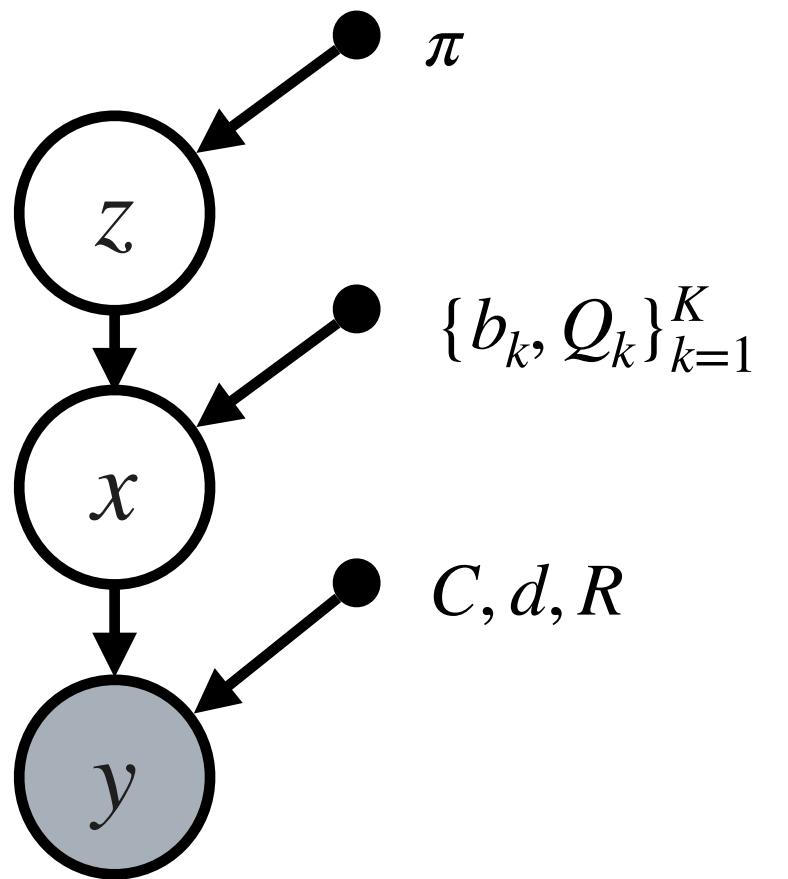
Now do the same for $q(x)$:

$$\begin{aligned}\log q(x) &= \mathbb{E}_{q(z)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(z)} [\log p(x \mid z, \Theta)] + \log p(y \mid x, \Theta) + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

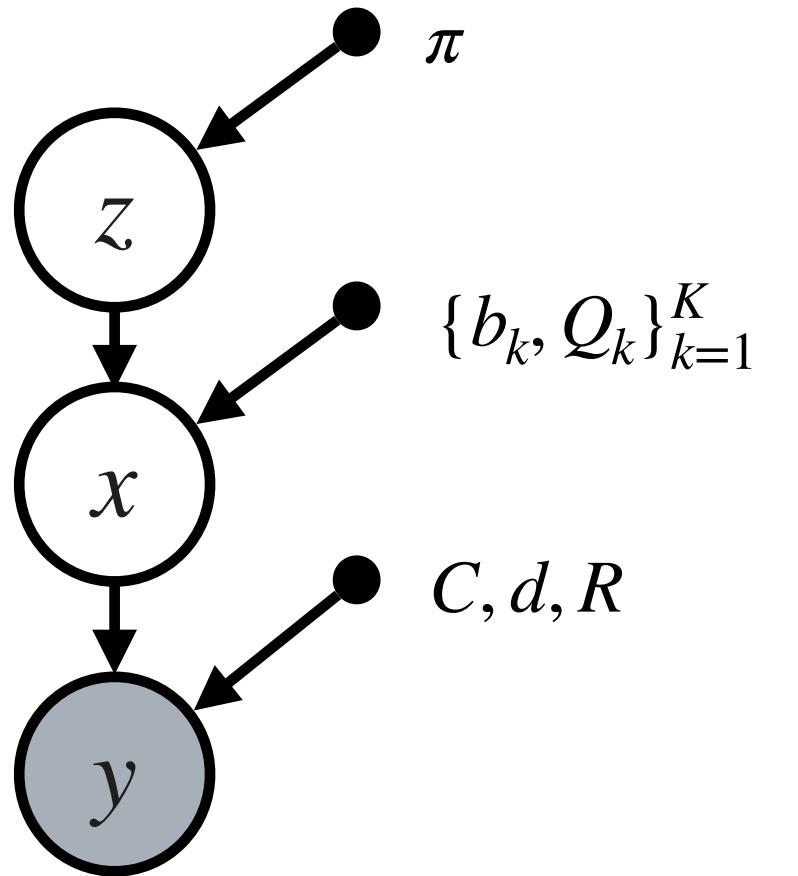
Now do the same for $q(x)$:

$$\begin{aligned}\log q(x) &= \mathbb{E}_{q(z)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(z)} [\log p(x \mid z, \Theta)] + \log p(y \mid x, \Theta) + c \\ &= \mathbb{E}_{q(z)} [\log \mathcal{N}(x \mid b_z, Q_z)] + \log p(y \mid x, \Theta) + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

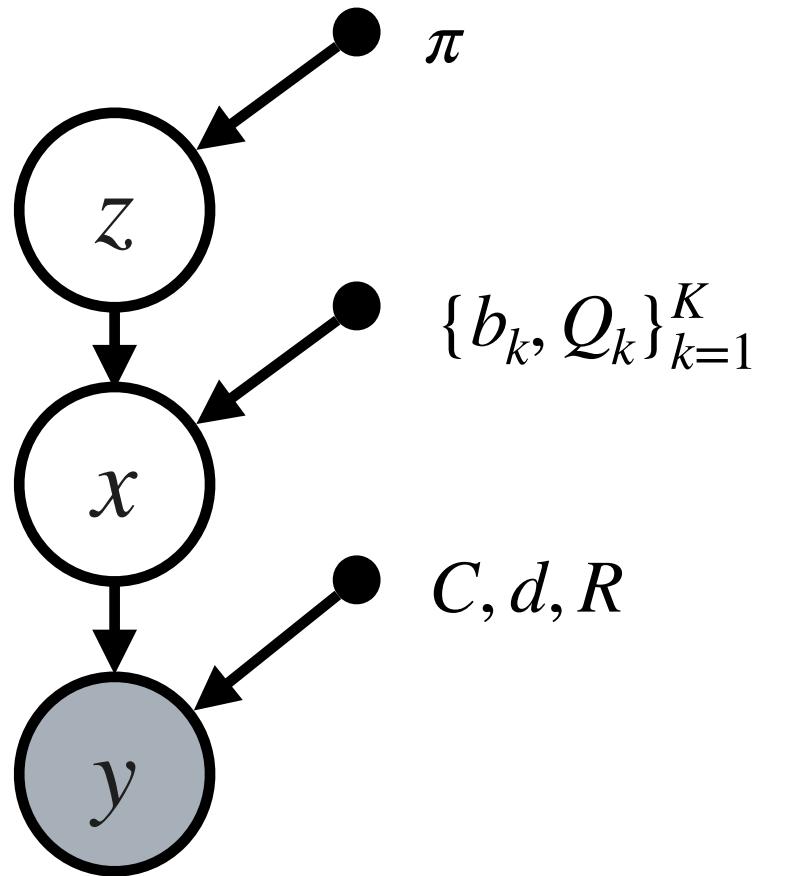
Now do the same for $q(x)$:

$$\begin{aligned}\log q(x) &= \mathbb{E}_{q(z)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(z)} [\log p(x \mid z, \Theta)] + \log p(y \mid x, \Theta) + c \\ &= \mathbb{E}_{q(z)} [\log \mathcal{N}(x \mid b_z, Q_z)] + \log p(y \mid x, \Theta) + c \\ &= -\frac{1}{2}x^\top \mathbb{E}_{q(z)}[Q_z^{-1}]x + x^\top \mathbb{E}_{q(z)}[Q_z^{-1}b_z] - \frac{1}{2}x^\top C^\top R^{-1}Cx + x^\top C^\top R^{-1}(y - d) + c\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

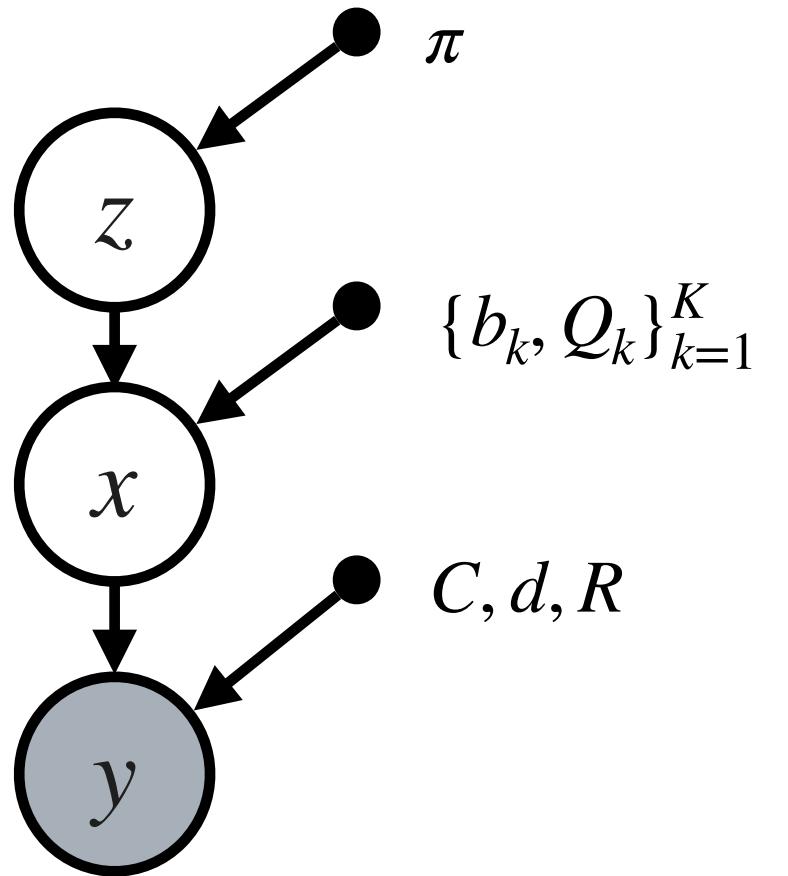
Now do the same for $q(x)$:

$$\begin{aligned}\log q(x) &= \mathbb{E}_{q(z)} [\log p(z, x, y \mid \Theta)] + c \\ &= \mathbb{E}_{q(z)} [\log p(x \mid z, \Theta)] + \log p(y \mid x, \Theta) + c \\ &= \mathbb{E}_{q(z)} [\log \mathcal{N}(x \mid b_z, Q_z)] + \log p(y \mid x, \Theta) + c \\ &= -\frac{1}{2}x^\top \mathbb{E}_{q(z)}[Q_z^{-1}]x + x^\top \mathbb{E}_{q(z)}[Q_z^{-1}b_z] - \frac{1}{2}x^\top C^\top R^{-1}Cx + x^\top C^\top R^{-1}(y - d) + c \\ &= \log \mathcal{N}(x \mid \tilde{\mu}, \tilde{\Sigma})\end{aligned}$$

$$z \sim \text{Cat}(\pi)$$

$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Coordinate ascent VI

Closed form updates for $q(x)$

The final result is,

$$\log q(x) = \log \mathcal{N}(x \mid \tilde{\mu}, \tilde{\Sigma})$$

where

$$\tilde{\mu} = \tilde{J}^{-1}\tilde{h}$$

$$\tilde{h} = \mathbb{E}_{q(z)}[Q_z^{-1}b_z] + C^\top R^{-1}(y - d)$$

$$= \sum_{k=1}^K [q(z=k)Q_k^{-1}b_k] + C^\top R^{-1}(y - d)$$

$$\tilde{\Sigma} = \tilde{J}^{-1}$$

$$\tilde{J} = \mathbb{E}_{q(z)}[Q_z^{-1}] + C^\top R^{-1}C$$

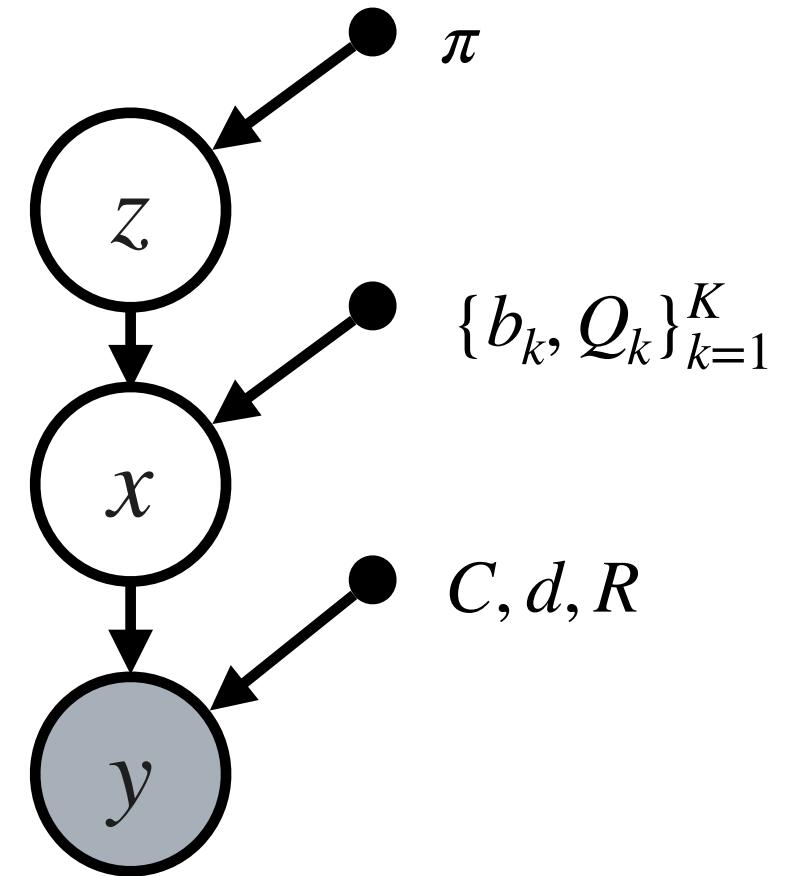
$$= \sum_{k=1}^K [q(z=k)Q_k^{-1}] + C^\top R^{-1}C$$

In other words, the natural parameters are the **expected natural parameters** under $q(z)$.

$$z \sim \text{Cat}(\pi)$$

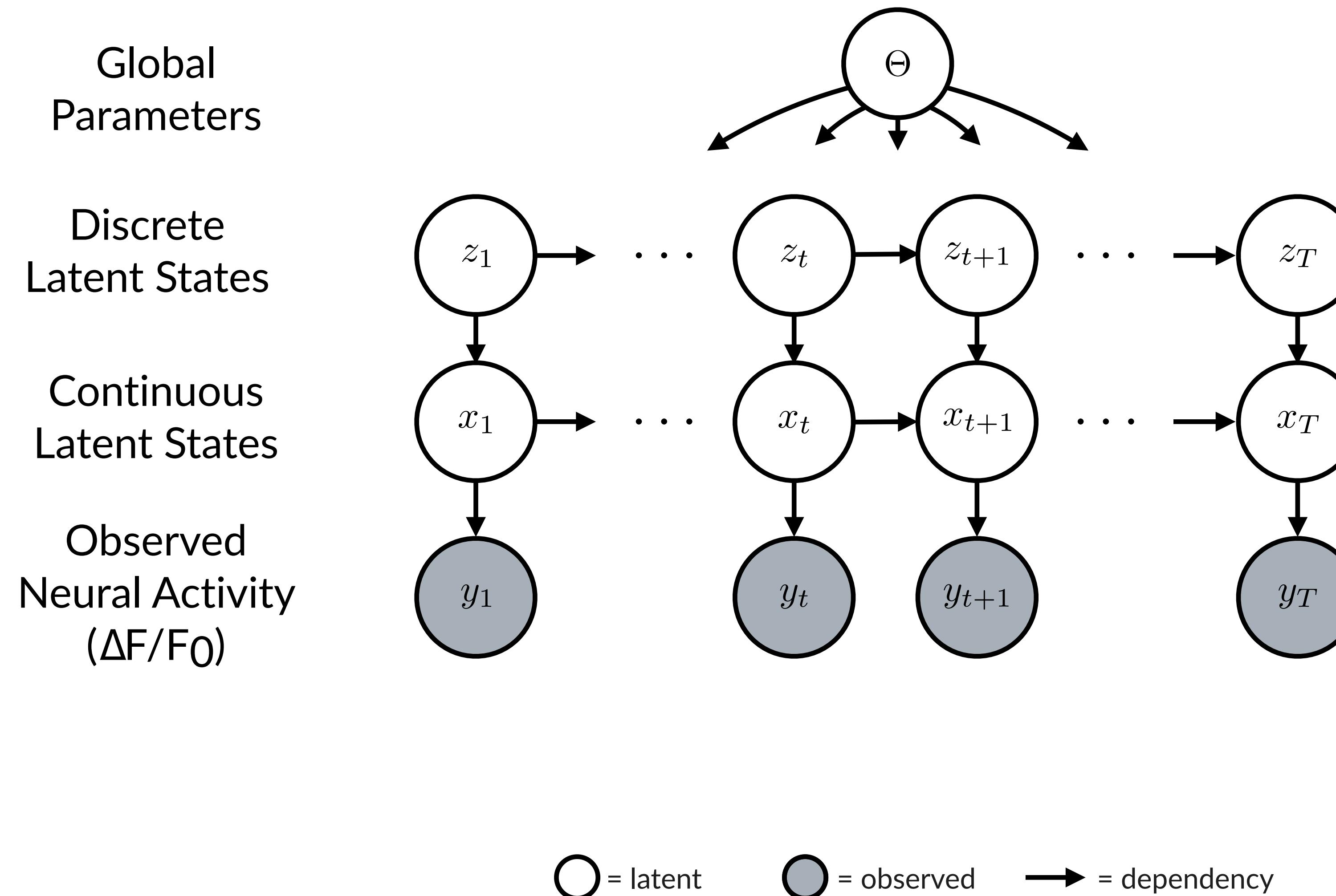
$$x \mid z \sim \mathcal{N}(b_z, Q_z)$$

$$y \mid x \sim \mathcal{N}(Cx + d, R)$$



Variational EM for SLDS

Variational EM for SLDS



Variational EM for SLDS

Structured mean field variational family

- Assume the variational posterior factors over discrete and continuous states.

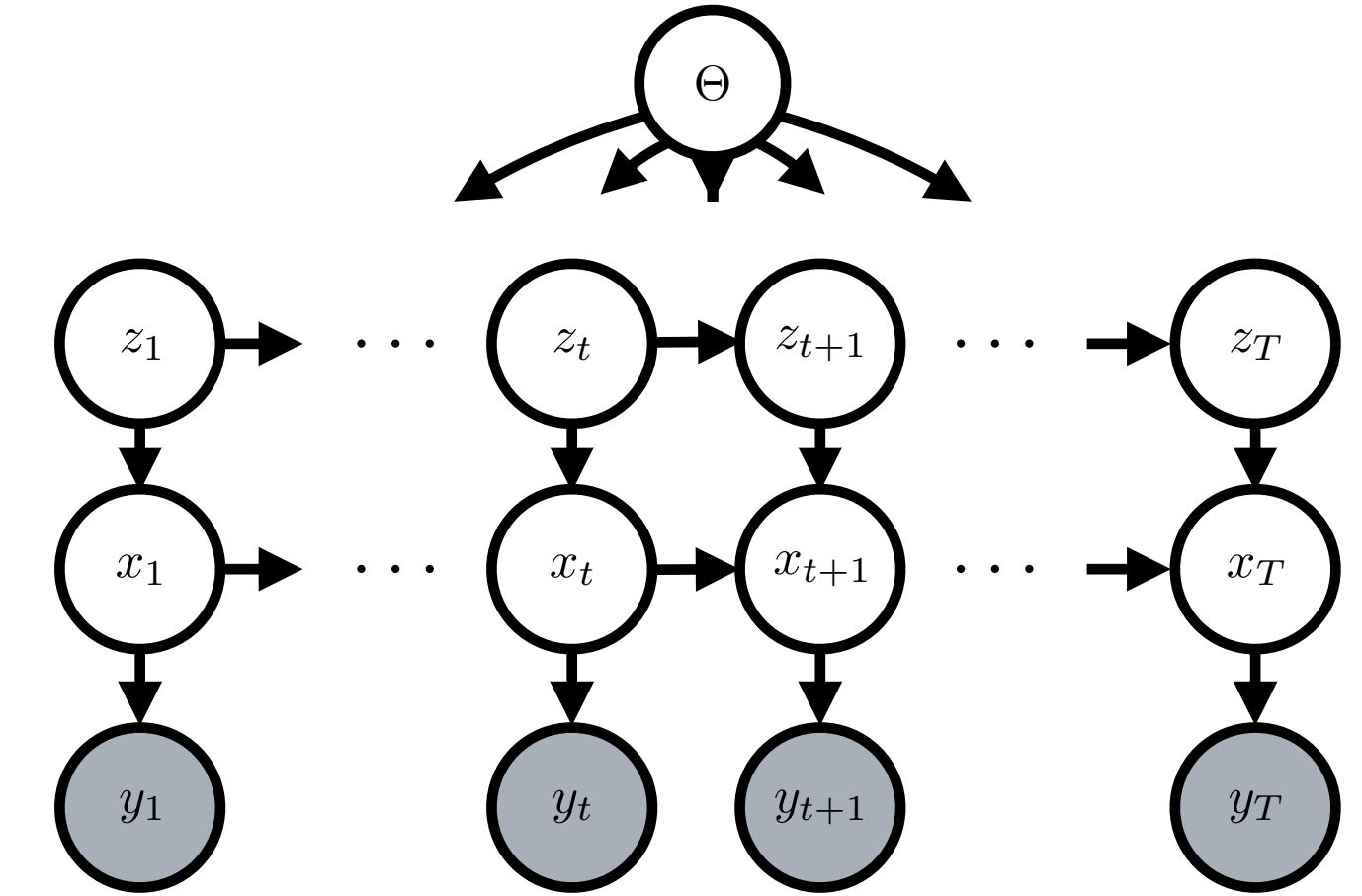
$$q(z_{1:T}, x_{1:T}) = q(z_{1:T}) q(x_{1:T})$$

where $z_{1:T}$ and $x_{1:T}$ are independent.

- There can still be dependencies within each factor though!
- This is called a **structured mean field family**.
- We want to find,

$$\arg \max_{q \in \mathcal{Q}} \mathcal{L}[q(z_{1:T})q(x_{1:T}), \Theta]$$

$$\equiv \arg \min_{q \in \mathcal{Q}} \text{KL} (q(z_{1:T})q(x_{1:T}) \parallel p(z_{1:T}, x_{1:T} \mid y_{1:T}, \Theta))$$



Variational EM for SLDS

Updating the discrete state posterior $q(z_{1:T})$

The optimal update for the discrete states takes the form

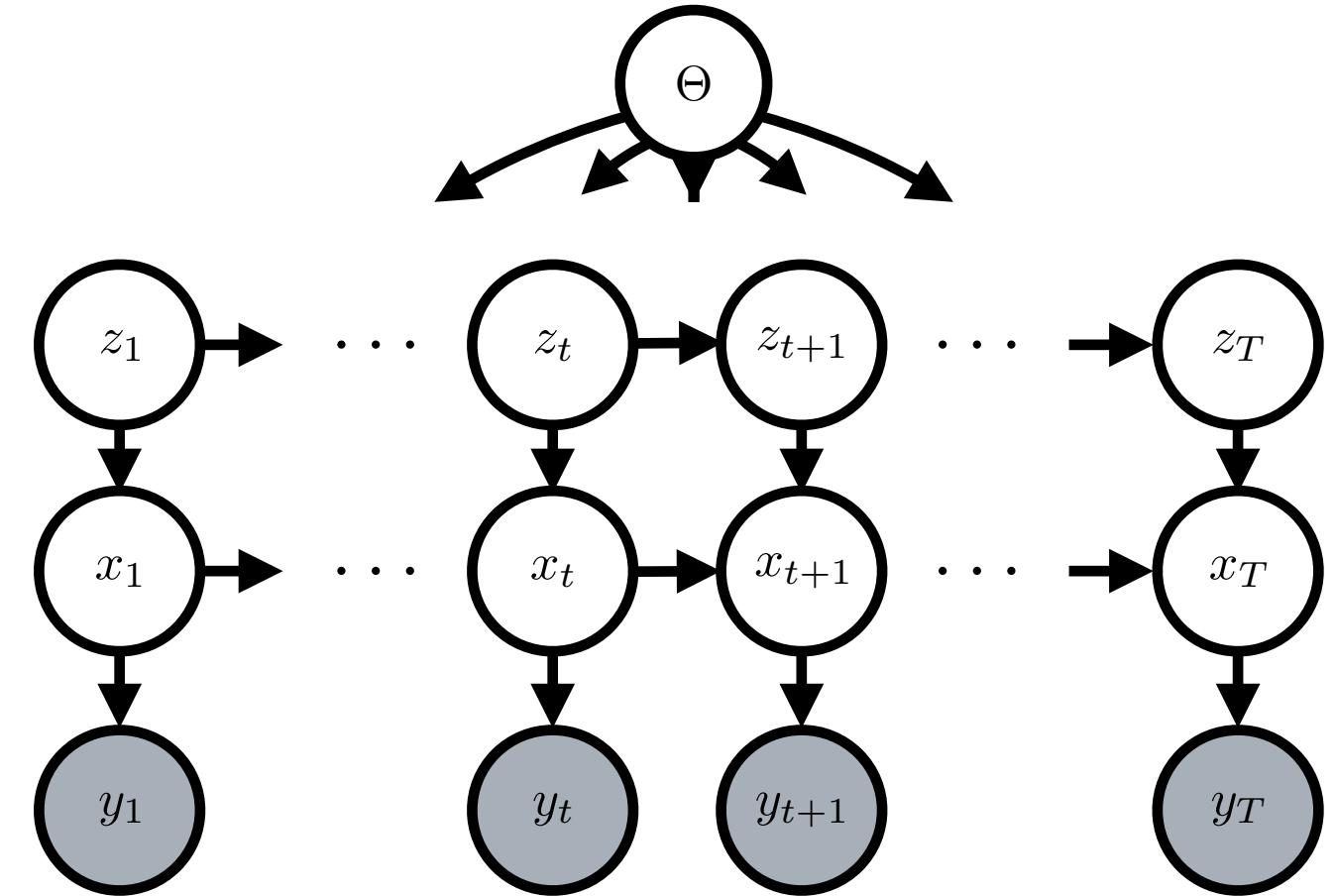
$$\begin{aligned} \log q(z_{1:T}) &= \log \text{Cat}(z_1 | \pi) + \sum_{t=2}^T \log \text{Cat}(z_t | P_{z_{t-1}}) \\ &\quad + \sum_{t=1}^T \sum_{k=1}^K \mathbb{I}[z_t = k] \log \tilde{\ell}_{tk} + c \end{aligned}$$

where

$$\log \tilde{\ell}_{tk} = \mathbb{E}_{q(x)} [\log \mathcal{N}(x_t | A_k x_{t-1} + b_k, Q_k)]$$

This is the **same form as the posterior in a hidden Markov model!**

But here, the log likelihoods are replaced with **expected log likelihoods** under $q(x)$.

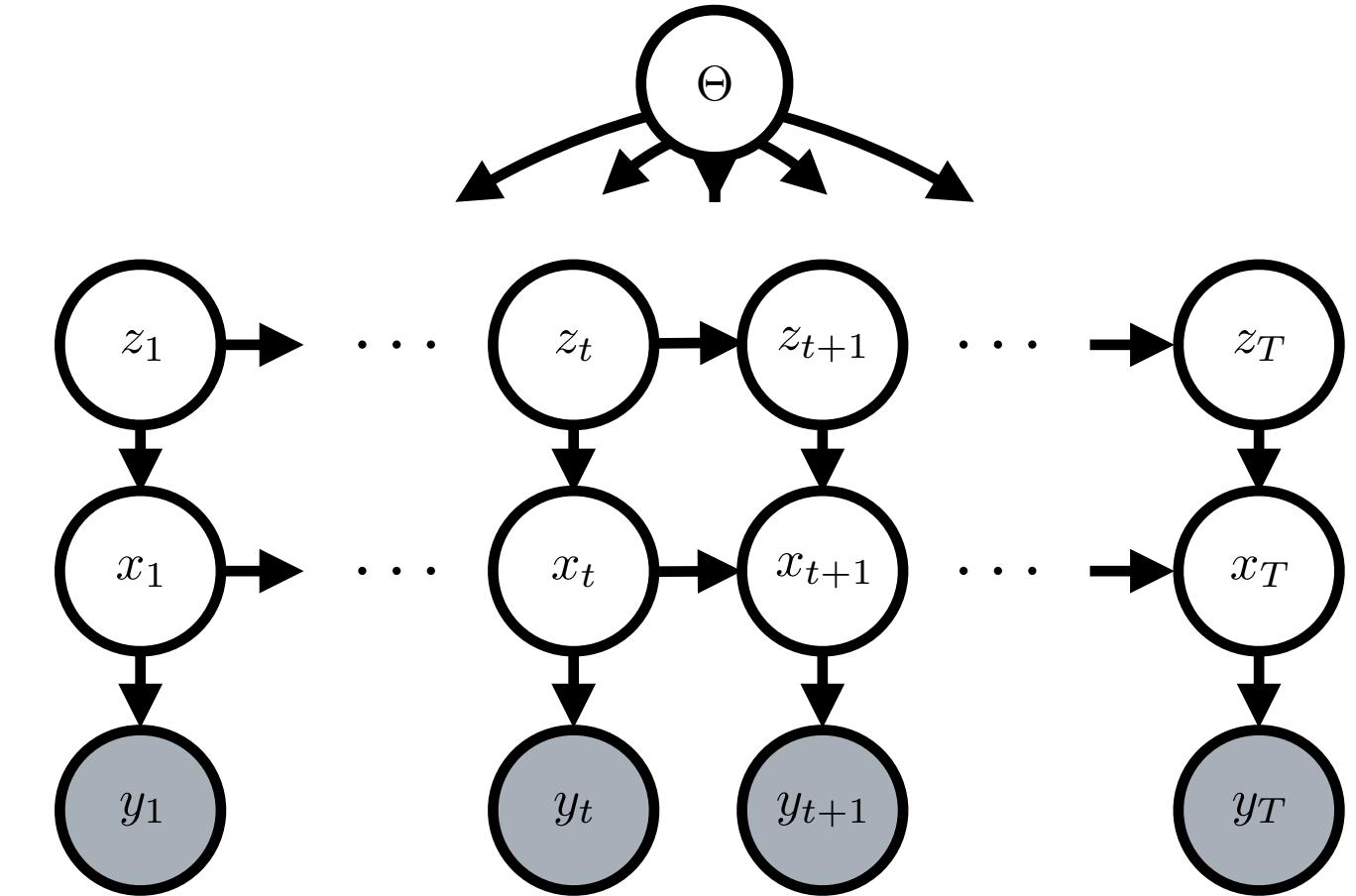


Variational EM for SLDS

Updating the continuous state posterior $q(x_{1:T})$

The optimal update for the continuous states takes the form

$$\begin{aligned} \log q(x_{1:T}) &= \mathbb{E}_{q(z)} \left[\log \mathcal{N}(x_1 \mid b_{z_1}, Q_{z_1}) \right] \\ &\quad + \sum_{t=2}^T \mathbb{E}_{q(z)} \left[\log \mathcal{N}(x_t \mid A_{z_t}x_{t-1} + b_{z_t}, Q_{z_t}) \right] \\ &\quad + \sum_{t=1}^T \log \mathcal{N}(y_t \mid Cx_t + d, R) + c \end{aligned}$$



Question: can you see what form this will take?

Variational EM for SLDS

Updating the continuous state posterior $q(x_{1:T})$

The optimal update for the continuous states is the same form as the posterior in **linear dynamical system**.

$$\log q(x_{1:T}) = \mathcal{N}(\text{vec}(x_{1:T}) \mid \tilde{J}^{-1}\tilde{h}, \tilde{J}^{-1})$$

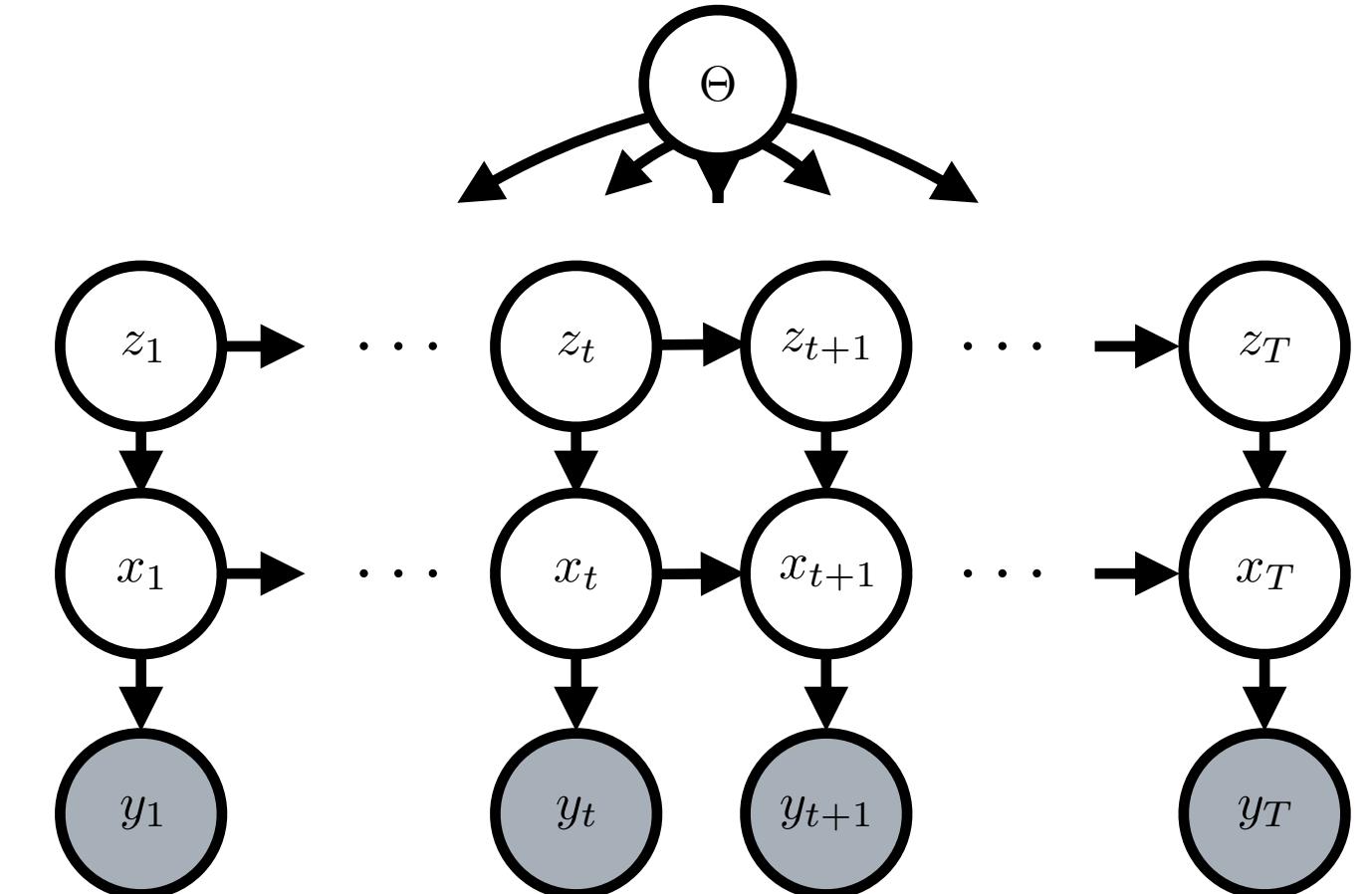
where

$$\tilde{J}_{tt} = \mathbb{E}_{q(z)}[Q_{z_t}^{-1}] + \mathbb{E}_{q(z)}[A_{z_{t+1}} Q_{z_{t+1}}^{-1} A_{z_{t+1}}] + C^\top R^{-1} C$$

$$\tilde{J}_{t,t-1} = \mathbb{E}_{q(z)}[Q_{z_t}^{-1} A_{z_t}]$$

$$\tilde{h}_t = \mathbb{E}_{q(z)}[Q_{z_t}^{-1} b_{z_t}] - \mathbb{E}_{q(z)}[A_{z_{t+1}} Q_{z_{t+1}}^{-1} b_{z_{t+1}}] + C^\top R^{-1} (y_t - d)$$

But here, the **natural parameters are expectations** under $q(z)$.



Variational EM for SLDS

Putting it all together

- **M-step:** Maximize the expected log probability

$$\Theta \leftarrow = \arg \max_{\Theta} \mathbb{E}_{q(z)q(x)}[\log p(y, x, z, \Theta)]$$

using expected sufficient statistics.

- **Variational E-step:**

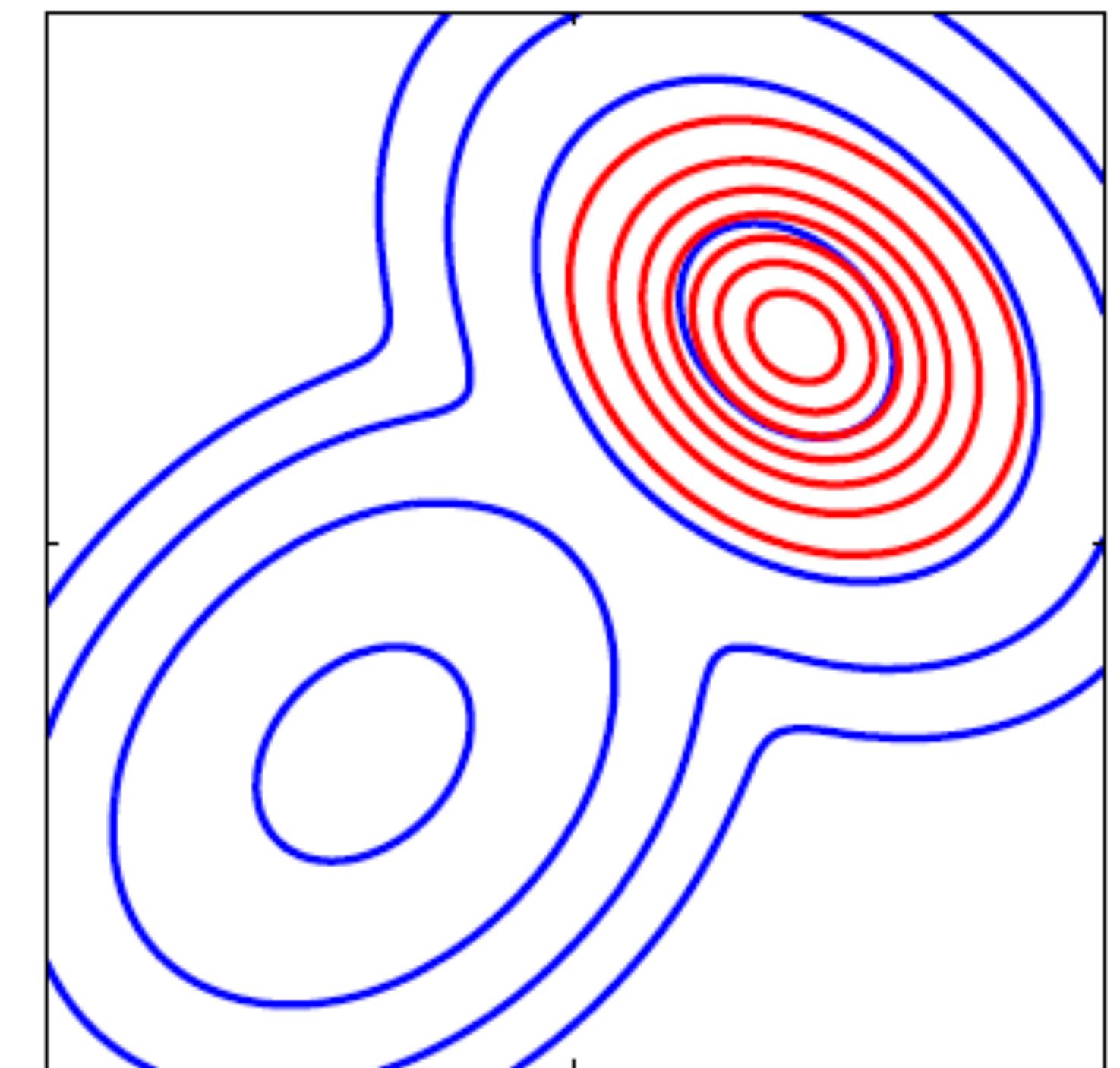
- Repeat until convergence:

1. $q(z) \leftarrow \text{HMM}(\pi, P, \tilde{\ell})$
where $\tilde{\ell}$ are the expected log likelihoods under $q(x)$

2. $q(x) \leftarrow \text{LDS}(\tilde{J}, \tilde{h})$
where \tilde{J} and \tilde{h} are the expected natural parameters under $q(z)$

- Compute the ELBO

$$\mathcal{L}[q(z) q(x), \Theta] \leq \log p(y | \Theta)$$



Conclusion

- In many latent variable models, including SLDS, the posterior distribution is intractable.
- Variational EM replaces the E step with a tractable variational approximation that minimizes the divergence to the true but intractable posterior.
- Mean field approximations are commonly used, as they often admit simple coordinate updates.
- In the SLDS, we can use a structured mean field approximation that retains dependencies across time while assuming independence of the discrete and continuous states.

Further Reading

- Barber, David. 2012. Bayesian Reasoning and Machine Learning. Cambridge University Press. **Chapter 25.**
- Blei, David M., Alp Kucukelbir, and Jon D. McAuliffe. 2017. “Variational Inference: A Review for Statisticians.” *Journal of the American Statistical Association* 112 (518): 859–77.
- Ghahramani, Z., and G. E. Hinton. 2000. “Variational Learning for Switching State-Space Models.” *Neural Computation* 12 (4): 831–64.
- Zoltowski, David, Jonathan Pillow, and Scott Linderman. 2020. “A General Recurrent State Space Framework for Modeling Neural Dynamics during Decision-Making.” In *Proceedings of the 37th International Conference on Machine Learning (ICML)*.