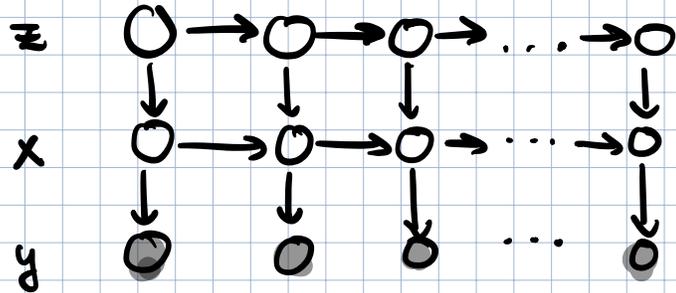
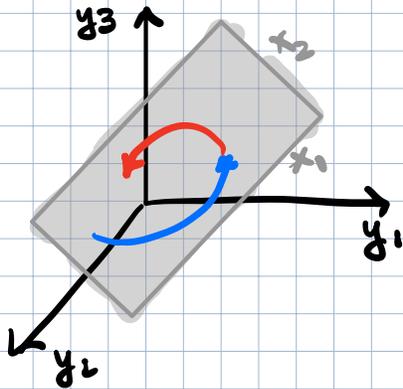


## Recall: SLDS



$$z_t \sim \text{Cat}(\pi_{z_{t+1}})$$

$$x_t \sim N(A_{z_t} x_{t-1} + b_{z_t}, Q_{z_t})$$

$$y_t \sim N(Cx_t + d, R)$$

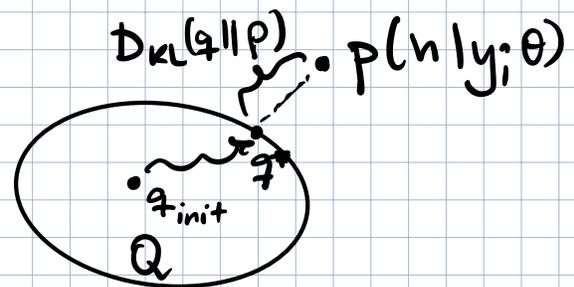
## EM for SLDS

E Step:  $q(z, x) \leftarrow p(z, x | y; \theta)$  **HARD!**

M Step:  $\theta \leftarrow \underset{\theta}{\text{argmax}} \mathbb{E}_{q} [\log p(z, x, y; \theta)]$

IDEA: Approximate  $q$  w/ VI  
Let  $h = (z, x)$

$$D_{KL}(q || p) = \mathbb{E}_{q} [\log q(h) - \log p(h | y; \theta)]$$



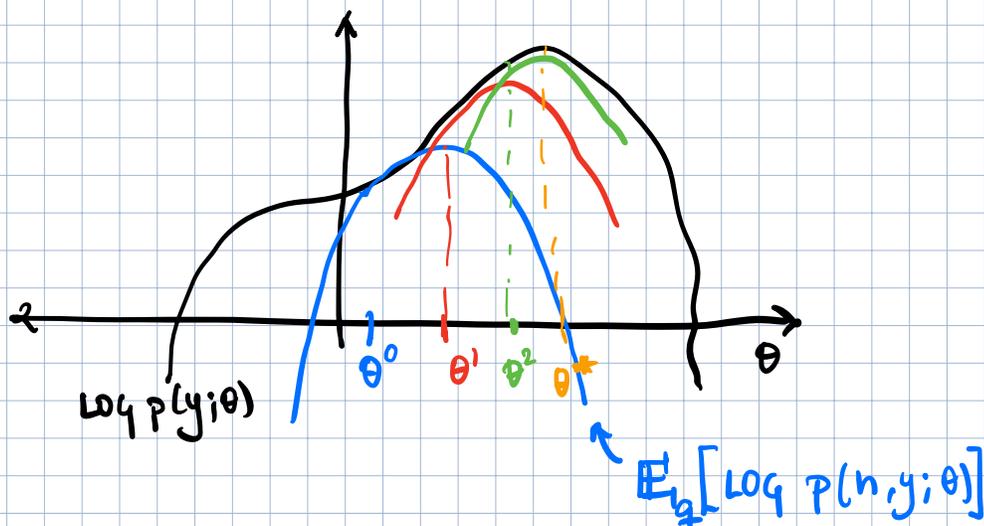
\* for certain  $Q$ , we can minimize KL via coordinate descent

## Variational EM:

E:  $q(h) \leftarrow \underset{Q}{\text{argmin}} D_{KL}(q(h) || p(h | y; \theta))$

M: same as above

# What is (Variational) EM doing?



Consider minimizing  $D_{KL}(q \parallel p)$

$$D_{KL}(q(h) \parallel p(h | y; \theta))$$

Bayes' Rule

$$p(h | y; \theta) = \frac{p(h, y; \theta)}{p(y; \theta)}$$

$$= -\mathbb{E}_q[\log p(h | y; \theta) - \log q(h)]$$

$$= -\mathbb{E}_q[\log p(h, y; \theta) - \log p(y; \theta) - \log q(h)]$$

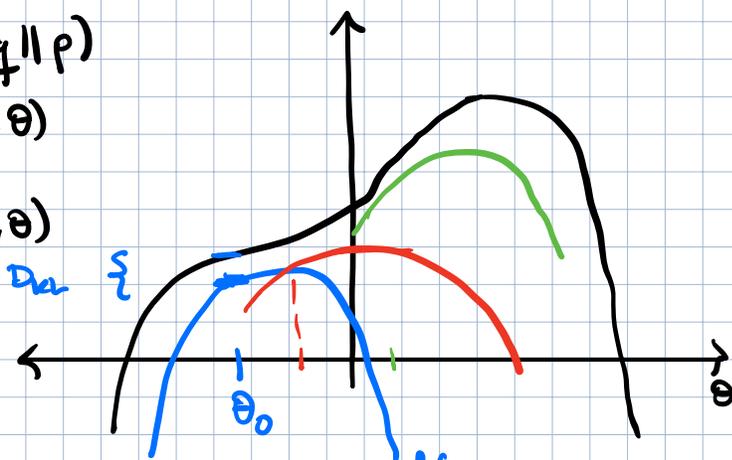
$$= \log p(y; \theta) - \underbrace{\mathbb{E}_q[\log p(h, y; \theta) - \log q(h)]}_{\text{EVIDENCE LOWER BOUND (ELBO) } \mathcal{L}(q, \theta)}$$

EVIDENCE LOWER BOUND (ELBO)  $\mathcal{L}(q, \theta)$

$$\mathcal{L}(q, \theta) = \log p(y, \theta) - D_{KL}(q(h) \parallel p(h | y; \theta))$$

Var. E step:  $q \leftarrow \underset{Q}{\operatorname{argmin}} D_{KL}(q \parallel p)$   
 $\equiv \underset{Q}{\operatorname{argmax}} \mathcal{L}(q, \theta)$

M step:  $\theta \leftarrow \underset{\Theta}{\operatorname{argmax}} \mathcal{L}(q, \theta)$



## Fixed form VI

- last time, we simply assumed  $q(z, x)$  factorized

$$Q_{MF} = \{q: q(z, x) = q(z)q(x)\}$$

This is called the mean field assumption

- A stricter assumption is to fix the functional form

$$\text{eg: } Q = \{q: q(h; \phi) = \text{Cat}(z; \phi_1) N(x | \phi_2, \phi_3)\}$$

$\phi = (\phi_1, \phi_2, \phi_3)$  are the variational params.

Then optimizing over  $Q \equiv$   
optimizing over params  $\phi$

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{q(z, x; \phi)} [\log p(h, y; \theta) - \log q(h; \phi)]$$

IDEA: maximize  $\mathcal{L}(\phi, \theta)$  via stochastic gradient ascent

$$\nabla_{\theta} \mathcal{L}(\phi, \theta) = \mathbb{E}_{q(h; \phi)} [\nabla_{\theta} \log p(h, y; \theta)]$$

$$\approx \frac{1}{M} \sum_m \nabla_{\theta} \log p(h^{(m)}, y; \theta) \quad \text{where } h^{(m)} \stackrel{\text{iid}}{\sim} q(h; \phi)$$

$$\nabla_{\phi} \mathcal{L}(\phi, \theta) \neq \mathbb{E}_{q(h; \phi)} [\nabla_{\phi} (\log p(h, y; \theta) - \log q(h; \phi))]$$

## Reparameterization Trick

Suppose  $\varepsilon \sim N(0, 1)$   $\leftarrow$  noise indep of  $\phi$   
 $h = r(\varepsilon, \phi) \Rightarrow h \sim q(h; \phi)$   
 $\leftarrow$  reparam. function

eg  $q(h; \phi) = N(h | \phi, 1)$  then  $h = \phi + \varepsilon = r(\varepsilon, \phi)$

Then

$$\nabla_{\phi} \alpha(\theta, \phi) = \mathbb{E}_{r(\varepsilon)} \left[ \nabla_{\phi} (\log q(r(\varepsilon, \phi), y; \theta) - \log q(r(\varepsilon, \phi); \phi)) \right]$$

"Law of the unconscious statistician"