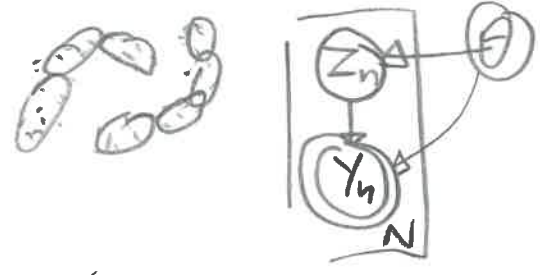


1. MOTIVATION: How to model complicated data?



APPROACH: GMM

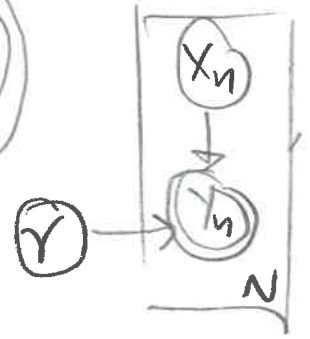
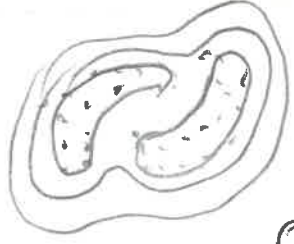


- Let  $\theta \triangleq \{ \pi, \mu_1, \Sigma_1, \dots, \mu_k, \Sigma_k \}$   
 DIR  $\sim$  NIW

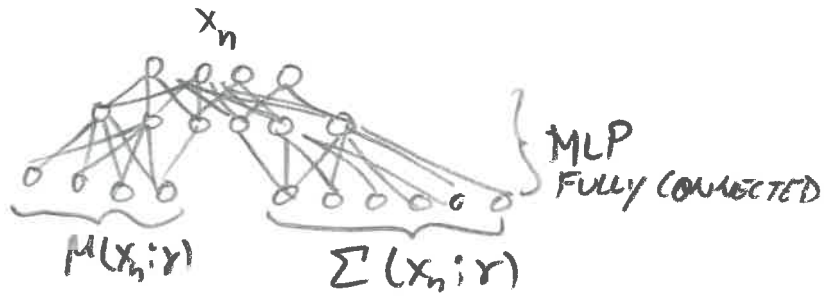
$Z_n | \pi \sim \text{CAT}(\pi)$   
 $Y_n | Z_n, \theta \sim N(\mu_{Z_n}, \Sigma_{Z_n})$

✓ INTERPRETABLE / STRUCTURED    ✗ RESTRICTIVE

APPROACH: DENSITY NET (VAE)



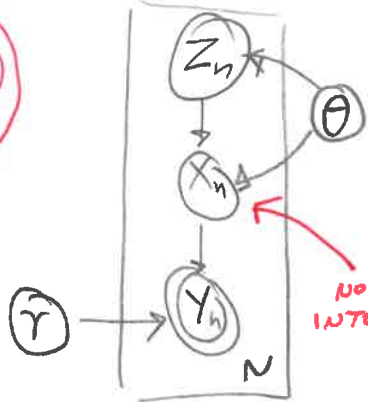
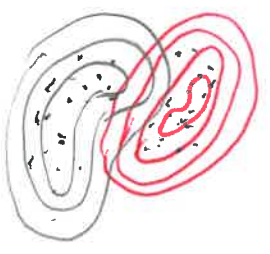
$\gamma \sim p(\gamma), X_n \sim N(0, I)$   
 $Y_n | X_n, \gamma \sim N(\mu(x_n; \gamma), \Sigma(x_n; \gamma))$



✓ FLEXIBLE

✗ NOT INTERPRETABLE

APPROACH: LATENT GMM (SVAE)



-  $\theta$ : SAME AS ABOVE  $\rightarrow$  GMM  
 $\gamma \sim p, X_n \sim N(\mu_{Z_n}, \Sigma_{Z_n}), Z_n \sim \text{CAT}(\pi)$   
 $Y_n | X_n, \gamma \sim N(\mu(x_n; \gamma), \Sigma(x_n; \gamma))$

NOT INTERPRETABLE

MLP

✓ FLEXIBLE    ✓ STRUCTURED / INTERPRETABLE

□ BUT HOW TO DO INFERENCE?

### 3: VAE RECAP

$Y$ : DATA,  $X$ : LATENT,  $\{y_n\}_{n=1}^N$ : Training Data

$P_r(Y|X) \Rightarrow N(\mu(x; \gamma), \Sigma(x; \gamma))$ : DECODER

$q_\phi(X|Y) \Rightarrow N(\mu(y; \phi), \Sigma(y; \phi))$ : ENCODER

#### MAX LIK. ESTIMATION

$$\max_{\gamma} \sum_{n=1}^N \log P_r(y_n) \geq \max_{\gamma, \phi} \sum_{n=1}^N \underbrace{-\text{KL}(q_\phi(X|y_i) \parallel P_r(X, y_i))}_{\equiv \mathcal{L}(\phi, \gamma; y_i)}$$

□ AVOIDS ISSUES w/ CONJUGACY

□ SGD optimization □ LACKS INTERPRETABLE X

$$\equiv \mathcal{L}(\phi, \gamma; y_i)$$

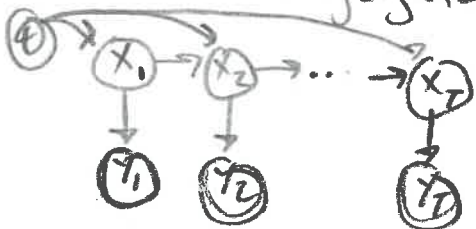
□ "Reparameterization Trick"  $\Rightarrow$  Lower variance gradient est.

□ Ammortized inference  $q_\phi(X|Y)$  learns mapping for any input  $y$

#### START HERE

### 2: STRUCTURED MEAN FIELD VARIATIONAL

□ Consider Conjugate linear-Gaussian dynamical system.



$$X_t | X_{t-1} \sim N(AX_{t-1}, \Sigma), \quad Y_t | X_t \sim N(BX_t, R)$$

$$\theta \equiv \{A, B, \Sigma, R\} \sim \text{CONJUGATE PRIOR } p(\theta)$$

- VARIATIONAL LOWER BOUNDS

$$\log P(Y) \geq \max_{q(\theta)q(x)} -\text{KL}(q(\theta)q(x) \parallel P(X, Y, \theta))$$

$$q(x) = q(x_1) \prod_{t=2}^T q(x_t | x_{t-1})$$

$\swarrow$  GAUSSIAN  $\searrow$  GAUSSIAN

- Complete conditionals are Conjugate

- Closed-form coordinate ascent updates

$$q(\theta) \propto \exp \left\{ \mathbb{E}_{q(x)} [\log P(\theta | x, y)] \right\}$$

$\Rightarrow$  Expfam w/ expected natural params  $\rightarrow$  Same family as  $p(\theta)$

~~When  $p(Y|X)$  more general closed-form updates not possible~~

# 4. STRUCTURED VAE

- MAIN IDEA: Combine VAE encoder for portions of structured model so closed-form SVI updates possible via conditional conjugacy.
- Assume conjugate pair:  $p(\theta) p(x|\theta)$
- Likelihood non-conjugate: e.g.  $p(y|x, r) = N(\mu(x; r), \Sigma(x; r))$

→ Variational Bounds

~~$$\log p(y) \geq \max_{q(\theta), q(r), q(x)} -\text{KL} [q(\theta) q(r) q(x) \| p(\theta) p(r) p(x|\theta) p(y|x, r)]$$~~

MLP

- Parametric variational bounds:

$$\log p(y) \geq \max_{\eta_\theta, \eta_r, \eta_x} -\text{KL} [q_{\eta_\theta}(\theta) q_{\eta_r}(r) q_{\eta_x}(x) \| p(\theta) p(r) p(x|\theta) p(y|x, r)]$$

$$\equiv \mathcal{L}(\eta_\theta, \eta_r, \eta_x)$$

- One issue:  $p(x|\theta)$  not conjugate to  $p(y|x, r)$  so update not closed-form

⇒ Sol'n: Surrogate objectives:

$$\hat{\mathcal{L}}(\eta_\theta, \eta_x, \phi) \triangleq -\text{KL} [q(\theta) q(x) \| p(\theta) p(x|\theta) \tilde{p}_\phi(y|x)]$$

where  $\tilde{p}_\phi(y|x) \propto \exp\{\langle r(y; \phi), t_x(x) \rangle\}$

MLP ↑

↑  
suff. stats.  
of  $p(x|\theta)$

⇒  $p(x|\theta)$  conjugate to  $\tilde{p}_\phi(y|x)$

⇒ Closed-form update:

$$\eta_x^*(\eta_\theta, \phi) = \underset{\eta_x}{\text{argmin}} \hat{\mathcal{L}}(\eta_\theta, \eta_x, \phi), \quad q^*(x) \propto \exp\{\langle \eta_x^*(\eta_\theta, \phi), t_x(x) \rangle\}$$

⇒ Closed-form updates of  $\eta_\theta$  &  $\eta_r$  minimizing  $\mathcal{L}(\eta_\theta, \eta_r, \eta_x^*(\eta_\theta, \phi))$

$$\mathcal{L}_r \triangleq \mathcal{L}^{\text{VAE}}(\eta_\theta, \eta_r, \phi) \quad \text{③}$$

## ADDE: CONDITIONALLY CONJUGATE MF UPDATES

- Conjugate pair:

$$p(x|\theta) = \exp\{\eta(\theta)^T T(x) - A(\theta)\}$$

$$p(\theta) \propto \exp\left\{\begin{pmatrix} \tau \\ \nu \end{pmatrix}^T \begin{pmatrix} \eta(\theta) \\ -A(\theta) \end{pmatrix}\right\}$$

- Posterior conjugacy:

$$\begin{aligned} p(\theta|x) &\propto p(x, \theta) \propto \exp\left\{\begin{pmatrix} \tau \\ \nu \end{pmatrix}^T \begin{pmatrix} \eta(\theta) \\ -A(\theta) \end{pmatrix} + \eta(\theta)^T T(x) - A(\theta)\right\} \\ &= \exp\left\{\begin{pmatrix} \tau + T(x) \\ \nu + 1 \end{pmatrix}^T \begin{pmatrix} \eta(\theta) \\ -A(\theta) \end{pmatrix}\right\} \Rightarrow \text{SAME EXPFORM AS } p(\theta) \end{aligned}$$

- MF FIXED-POINT:

$$q(\theta) \propto \exp\left\{\mathbb{E}_{q_x}[\log p(\theta, x)]\right\}$$

$$= \exp\left\{\begin{pmatrix} \tau + \mathbb{E}_{q_x}[T(x)] \\ \nu + 1 \end{pmatrix}^T \begin{pmatrix} \eta(\theta) \\ -A(\theta) \end{pmatrix}\right\} \Rightarrow \text{SAME EXPFORM AS } p(\theta)$$