

## Black Box Variational Inference

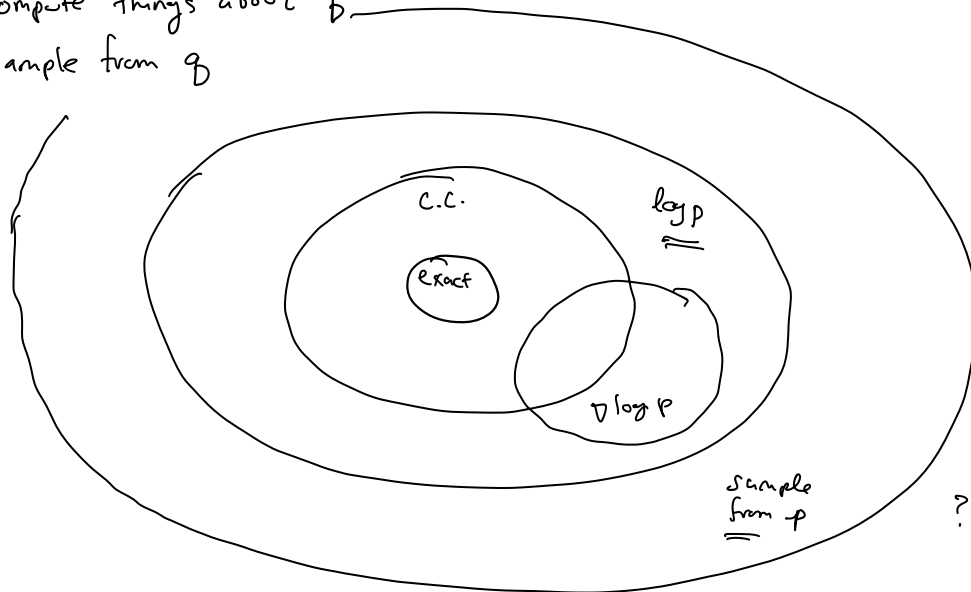
complete conditionals  $\rightarrow$  Gibbs, CAVI

$$p(\theta, z_{1:n}, x_{1:n}) \quad \text{global } \theta \quad \text{local } z_{1:n} \\ = p(\theta) \prod_{i=1}^n p(z_i | \theta) p(x_i | z_i, \theta)$$

$p(\theta, z_{1:n} | x_{1:n})$  - posterior. Goal: approximate it with  $q(\theta, z_{1:n}; \nu)$

Black box criteria:

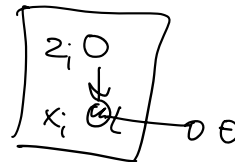
- compute the log joint (compute  $\nabla_{\theta, z} \log p(z, \theta, x)$ )
- compute things about  $q$
- sample from  $q$



Recall: Deep generative model.

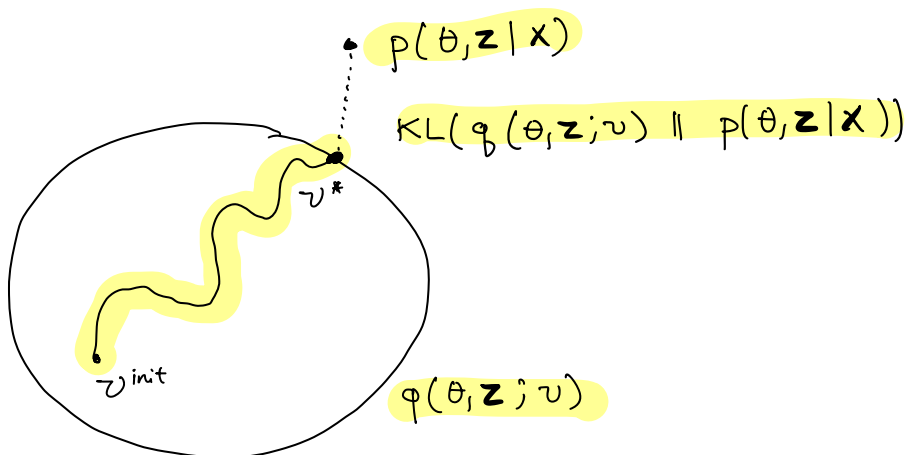
$$z_i \sim N_p(0, 1)$$

$$x_i | z_i, \theta \sim N(\Sigma_\mu(z_i), \Sigma_\sigma^z(z_i))$$



$$p(z_i | x_i, \theta) = \frac{p(z_i, x_i | \theta)}{p(x_i | \theta)}$$

$$p(z_i, x_i | \theta) = p(z_i) p(x_i | z_i, \theta)$$



Objective function

ELBO

EM: maximum likelihood estimation alg.  
when there are latent variables in the model.

$$\mathcal{L}(v) = \mathbb{E}_v [\log p(\theta, Z, x) - \log q(\theta, Z; v)]$$

Strategy:

- ① Write  $\nabla_v \mathcal{L}$  as an expectation  $\mathbb{E}_q[\sim]$
- ② Take a Monte Carlo approximation of  $\nabla_v \mathcal{L} \approx \frac{1}{B} \sum_b \sim_b$
- ③ Use stochastic optimization

Score gradient

$$\nabla_v \mathcal{L} = \mathbb{E} \left[ \underbrace{\nabla_v \log q(\theta, Z; v)}_{\text{score function}} \underbrace{(\log p(\theta, Z, x) - \log q(\theta, Z; v))}_{\text{instantaneous ELBO}} \right]$$

$$\left. \begin{array}{l} \nabla_v \log q(\theta, Z; v) \\ \nabla_{\theta, Z} \log q(\theta, Z; v) \\ \nabla \log q(\theta, Z; v) \end{array} \right\} \text{"score function"}$$

## Score VI

For iteration  $t$ :

$$\theta_b, z_b \sim q(\theta, z; v_t) \quad b = 1 \dots B$$

$$g_t = \frac{1}{B} \sum_{b=1}^B \nabla_v \log q(\theta_b, z_b; v_t) (\log p(\theta_b, z_b, x) - \log q(\theta_b, z_b; v_t))$$

$$v_{t+1} = v_t + \rho_t g_t$$

Example: DGM

$$\log p(\theta, z, x) = \log p(\theta) + \sum_{i=1}^n \log p(z_i) + \log p(x; z, \theta)$$

$$q(\theta, z; v) = q(\theta; v_\theta) \prod_{i=1}^n q(z_i; v_i)$$

$\uparrow \quad \quad \quad \uparrow$   
 $N \quad \quad \quad N$

In practice:

Variance is high.  $\int$  Andy Miller, Justin Domke, ...  
Control variates.

Let  $g(\theta, z; v) \triangleq$  <sup>stochastic</sup> Score gradient,  $\mathbb{E}[g(\theta, z; v)] = \nabla_v \log p(\theta, z, x)$

Choose  $h(\theta, z; v)$  s.t.  $\mathbb{E}[h(\theta, z; v)] = \nabla_v \log p(\theta, z, x)$

Define  $\tilde{g}(\theta, z; v) = g(\theta, z; v) - \xi (h(\theta, z; v) - \mathbb{E}[h])$   
 $\uparrow$   
scalar

$$\text{Var}(\tilde{g}) = \text{Var}(g) + \xi^2 \text{Var}(h) - 2\xi \text{Cov}(g, h)$$

$$\xi^* = \frac{\text{Cov}(g, h)}{\text{Var}(h)}$$

$$\begin{aligned}
\mathbb{E} [\nabla_v \log q(\mathbf{z}; v)] &= \int q(\mathbf{z}; v) \nabla_v \log q(\mathbf{z}; v) d\mathbf{z} \\
&= \int \nabla_v q(\mathbf{z}; v) d\mathbf{z} \\
&= \nabla_v \int q(\mathbf{z}; v) d\mathbf{z} \\
&= \nabla_v 1 = \emptyset
\end{aligned}$$

FACT:

$$\begin{aligned}
\nabla_v q(\mathbf{z}; v) \\
= q(\mathbf{z}; v) \nabla_v \log q(\mathbf{z}; v)
\end{aligned}$$

Reparameterization gradient

transform the variables in  $q$ .

$$\begin{aligned}
\varepsilon &\sim s(\varepsilon) \rightarrow \mathbf{z} \sim q(\mathbf{z}; v) \\
\mathbf{z} &= t(\varepsilon, v)
\end{aligned}$$

e.g. for a normal:

$$\begin{aligned}
\varepsilon &\sim \mathcal{N}(0, 1) \rightarrow \mathbf{z} \sim \mathcal{N}(v_\mu, v_\sigma^2) \\
\mathbf{z} &= \varepsilon v_\sigma + v_\mu
\end{aligned}$$

Reparameterization gradient

$$\nabla \mathcal{L}(v) =$$

$$\mathcal{L}(v) = \mathbb{E}_{s(\varepsilon)} [\log p(\theta, \mathbf{z}, x) - \log q(\theta, \mathbf{z}; v)]$$

$$\text{where: } \theta = t(\varepsilon_\theta, v_\theta), \mathbf{z} = t(\varepsilon_z, v_z)$$

$$\begin{aligned}
\nabla_v \mathcal{L}(v) &= \mathbb{E}_{s(\varepsilon)} \left[ \nabla_{\theta, \mathbf{z}} [\log p(\theta, \mathbf{z}, x) - \log q(\theta, \mathbf{z}; v)] \nabla_v t(\varepsilon, v) \right. \\
&\quad \left. - \nabla_v \log q(\theta, \mathbf{z}; v) \right]
\end{aligned}$$

$$= \mathbb{E}_{s(\varepsilon)} \left[ \nabla_{\theta, \mathbf{z}} [\log p(\theta, \mathbf{z}, x) - \log q(\theta, \mathbf{z}; v)] \nabla_v t(\varepsilon, v) \right]$$

Score gradient  
discrete + cont LV's  
larger variance  
cheaper

vs. Reparameterization gradient  
continuous LV's  
smaller variance  
more expensive

ADVI  
Stan HMC

ELBO surgery

$$\mathcal{L} = \mathbb{E}[\log p(x|z, \theta)] - \underbrace{\text{KL}(q(z, \theta) \parallel p(z, \theta))}_{\text{Sticking the Landing}}$$

Sticking the Landing

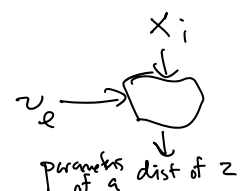
Amortized inference and the variational autoencoder

DGM:  $\theta \sim p(\theta)$

for each data point  $i$ :

$$z_i \sim \mathcal{N}_p(0, 1)$$

$$x_i \sim \mathcal{N}(\mu(z_i; \theta), \Sigma_\sigma(z_i; \theta)^2)$$



observe:  $\{x_i\}_{i=1}^n$  goal:  $p(\theta, z_{1:n} | x_{1:n})$

inference network:  $\Phi(x_i; v) \rightarrow q(z_i)$

$$q(z_i; v_e, x_i) = \Phi(x_i; v_e)$$

$$q(\theta, z_{1:n}; v) = \frac{q(\theta; v_\theta) \prod_{i=1}^n q(z_i; v_e, x_i)}{q(\theta; v_\theta) \prod_{i=1}^n q(z_i; v_i)}$$

Amortized mean field  
Classical mean field

$$q(z_i; v_e, x_i) = \mathcal{N}(\Phi_\mu(x_i; v_e), \Phi_\sigma(x_i; v_e)^2)$$

$$\mathcal{L}(v) = \mathbb{E} \left[ \log p(\theta) - \log q(\theta; v_\theta) + \sum_{i=1}^n (\log p(z_i, x_i | \theta) - \log q(z_i; v_e, x_i)) \right]$$

Transformation of local latent variables.

$$\varepsilon_i \sim \mathcal{N}_p(0, 1)$$

$$t(\varepsilon_i, x_i, v_e) = \varepsilon_i \cdot \Phi_\sigma(x_i; v_e) + \Phi_\mu(x_i; v_e) \rightarrow z_i \sim q(z_i; x_i, v_e)$$

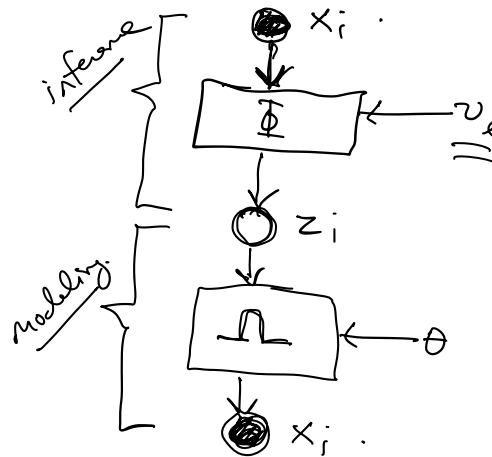
$$\nabla_v \mathcal{L} = \mathbb{E} \left[ \nabla_\theta (\log p(\theta) + \sum_{i=1}^n \log p(z_i, x_i | \theta) - \log q(\theta; v_\theta)) \nabla_v \theta + (\varepsilon_\theta, v_\theta) \right]$$

$$\nabla_{\nu} \mathcal{L} = \sum_{i=1}^n \mathbb{E} \left[ \nabla_{z_i} (\log p(z_i, x_i | \theta) - \log q(z_i, x_i, \nu_e)) \nabla_{\nu_e} \epsilon(\epsilon_i, x_i, \nu_e) \right]$$

$\Omega(z_i; \theta) = \text{"decoder"}$

$\Phi(x_i; \nu_e) = \text{"encoder"}$

$$\frac{1}{B} \sum$$



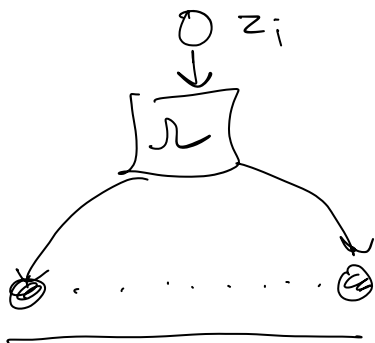
Zemel.  
Helmholtz  
Machine  
Autoencoder

SCVI

$$z_i \sim \mathcal{N}_p(0, 1)$$

$$x_i \sim \text{NB}(\Omega(z_i; \theta))$$

$$q(z_i; x_i, \nu_e) = \mathcal{N}(\Phi_\mu(z_i; \nu_e), \Phi_\sigma(z_i; \nu_e)^2)$$



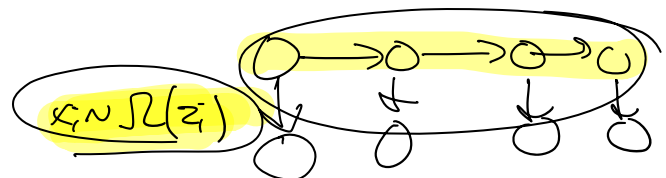
counts, sparse.

$$\nabla \mathcal{L}(\nu) = \mathbb{E}_q \left[ \nabla_{\nu} \log q(\cdot) (\log p(\cdot) - \log q(\cdot)) \right]$$

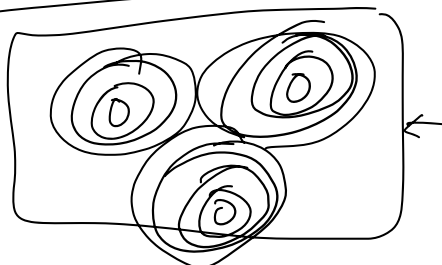
Galaxy properties (from sim)

$z_i$ 's.

low-dim.



Structured VAEs — "old" (2016)



$$z \rightarrow x \sim \mathcal{P}(z) \quad \Omega_z(\cdot)$$

EDM.

Conditional models

Linear + logistic regression

Mixtures - cluster the data.

Mixed-membership - grouped data

Matrix factorization - interaction data.

Exponential family + conjugacy

Generalized linear models

Deep learning - regression  
generalized neural models  
deep generative models

= Deep exp. families.

exact inference.

MAP inference.

Gibbs sampler.

CAVI

generic inference w/ BBVI.

