

# Deep Learning

- what is a deep neural network?
- neural conditional models (supervised)
- deep generative models (unsupervised)
- deep exponential families

## Neural network

parameterized function from features  $x \in \mathbb{R}^P$  to output  $w \in \mathbb{R}^d$

notation:  $\mathcal{L}(x; \theta)$

"hidden unit"

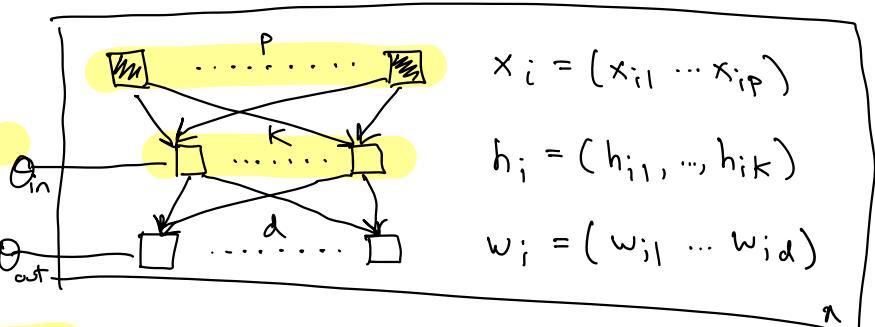
$$x_i \in \mathbb{R}^P$$

$$h_i = g(\theta \cdot x_i) \quad \text{e.g., } g \text{ is logistic}$$

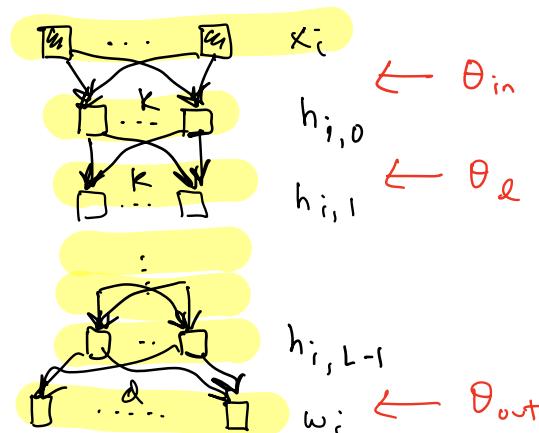
one-layer NN

$$h_{i,k} = g(\theta_{in,k} \cdot x_i) \quad k=1 \dots K$$

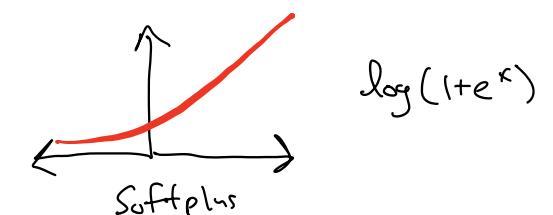
$$w_{ij} = \theta_{out,j} \cdot h_i \quad j=1 \dots d$$



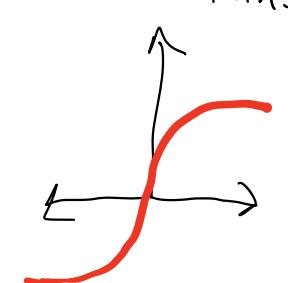
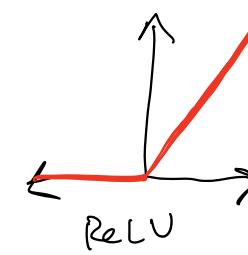
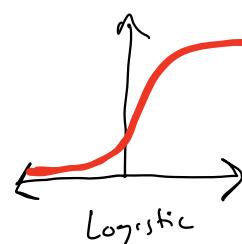
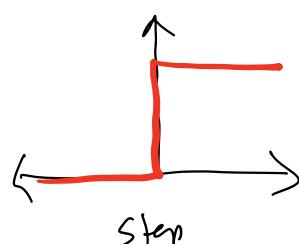
## deep NN, aka feed-forward NN



$$\begin{aligned} &\nabla_\theta \mathcal{L}(x; \theta) \\ &\mathcal{L}(x; \theta) \end{aligned}$$



### ① activation function $g(\cdot)$



## Neural regression

$$\theta \sim N_B(0, \lambda^2)$$

$$y_i | x_i \sim N(\mu(x_i; \theta), \sigma^2)$$

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

MAP estimation of  $\theta$ :

$$\mathcal{L}(\theta) = -\frac{1}{2\lambda^2} \sum_{j=1}^w \theta_j^2 + \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - \mu(x_i; \theta))^2$$

$$\nabla \mathcal{L}(\theta) = -\underbrace{\lambda \theta}_{\text{w-vector}} + \underbrace{\sum_{i=1}^n \frac{1}{\sigma^2} (y_i - \mu(x_i; \theta)) \nabla_\theta \mu(x_i; \theta)}_{\text{w-vector}}$$

Generalized neural conditional model (GNCM)

$$\theta \sim N_w(0, \lambda^2)$$

$$y_i | x_i \sim \text{expfam}(\eta_i) \quad \eta_i \triangleq \eta_\mu(f(\mu(x_i; \theta)))$$

$$x_i \xrightarrow{\mu} w_i \xrightarrow{f} \mu_i \xrightarrow{\eta_\mu} \eta_i \rightarrow y_i$$

$$\text{canonical link: } f = \eta_\mu^{-1} = \nabla a(\eta)$$

in a GLM,  $w_i = \theta \cdot x$ ;  
here,  $w_i = \mu(x_i; \theta)$

$$y_i | x_i \sim \text{expfam}(\mu(x_i; \theta))$$

cross entropy loss : Bernoulli response w/ canonical link

soft max : Categorical response w/ " "

$$y_i | x_i \sim N(\mu(x_i; \theta), \sigma^2)$$

$$\mathcal{L}(\theta) = -\frac{1}{2\lambda^2} \sum_{j=1}^w \theta_j^2 + \sum_{i=1}^n \underbrace{\eta_\mu(f(\mu(x_i; \theta))) \cdot y_i - a(\eta_\mu(f(\mu(x_i; \theta))))}_{\eta_i}$$

$y_i$  is scalar

with canonical link:

$$\mathcal{L}(\theta) = -\frac{1}{2\lambda^2} \sum_{j=1}^w \theta_j^2 + \sum_{i=1}^n (\mu(x_i; \theta) \cdot y_i - a(\mu(x_i; \theta)))$$

$$\nabla \mathcal{L}(\theta) = \frac{1}{N} \theta + \sum_{i=1}^n (y_i - \mathbb{E}[y|x=x_i, \theta]) \nabla_\theta \mathcal{L}(x_i; \theta)$$

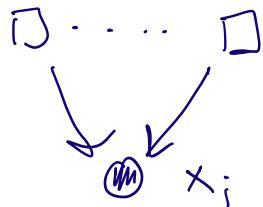
Deep generative models.  $\theta \sim p(\theta)$



$$z_i \sim N_k(0, \lambda^2)$$

$$x_i | z_i \sim \text{expfam}(\mathcal{L}(z_i; \theta))$$

Deep generative model.



$x_i$  : high dimensional

$z_i$  : lower dimension.

dimension reduction.

$$z_i \sim N_k(0, \lambda^2)$$

$$x_i \sim N_p(\mathcal{L}_\mu(z_i; \theta), \mathcal{L}_{\sigma^2}(z_i; \theta))$$

"decoder"  $\xrightarrow{\quad}$

$p(z_i | x_i, \theta)$  —  $K$ -dim embedding of  $x_i$

$$p(x|\theta) = \int p(x|z, \theta) p(z) dz$$

$$z_i \sim N_k(0, \lambda^2)$$

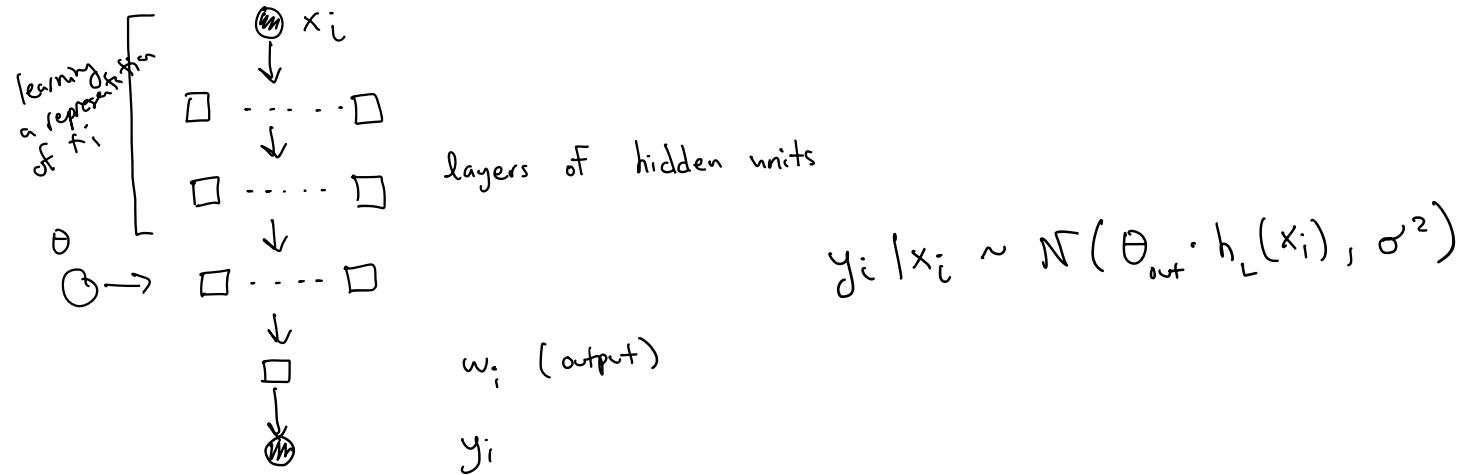
$$x_i | z_i \sim N(\beta \cdot z_i, I_{\sigma^2}) \quad \beta = K \cdot p$$

prob. PCA.

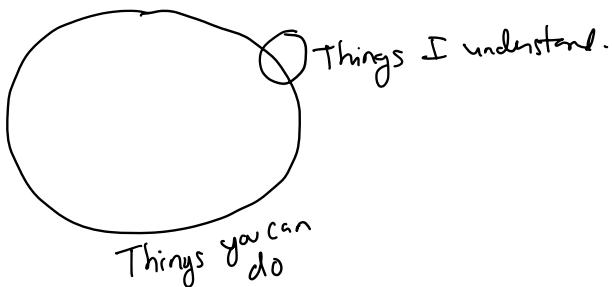


## TODAY

- deep generative models
- deep exponential families
- architectural elements: residuals, attention, convolution
- priors for neural networks



Ben Recht's Venn Diagram



## Deep Generative Model.

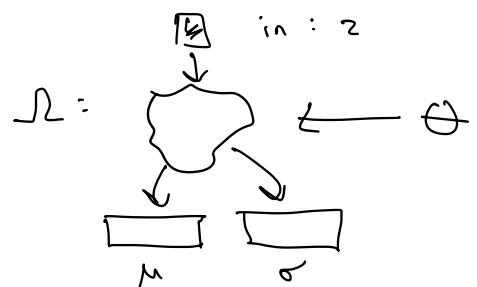
$$z_i \sim N_k(0, 1) \quad z_i \in \mathbb{R}^k$$

$$x_i \sim \text{expfam}(\mathcal{L}(z_i; \theta)) \quad x_i \in \mathbb{R}^d$$

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

local posterior:  $p(z_i | x_i, \theta)$   
 $\mathbb{E}[z | x_i, \theta] \stackrel{\Delta}{=} \hat{z}_i$

dimension reduction.



aside: Prob. PCA (Tipping + 1999)  
 $z_i \sim N_k(0, 1)$   
 $x_i | z_i \sim N(\beta \cdot z_i, \sigma^2)$   
 $d \times k$

Idea: DGM can capture many complex distributions of  $x$

$$\begin{aligned} p(x | \theta) &= \int p(z, x | \theta) dz && \text{marginal distribution of an obs.} \\ &= \int p(z) p(x | z, \theta) dz \\ &= \int_{\text{normal}_k(z; 0, 1)} \text{normal}(x; \mathcal{L}_\mu(z), \mathcal{L}_\sigma(z)^2) dz \end{aligned}$$

$$\mathcal{D} = \{x_i\}_{i=1}^n$$

$$\boxed{p(\theta | x_{1:n})}$$

$$p(\theta, z_{1:n}, x_{1:n}) = p(\theta) \prod_{i=1}^n p(z_i) p(x_i | z_i, \theta) \quad \leftarrow \text{"Factor model"}$$

$z_i$ : local latent variables

$\theta$ : global " "

Sc VI.

$$x_i \sim \text{NegBin}(\mathcal{L}(z_i; \theta))$$

RNA seq data: 10k-20k dim obs (counts), 20 dim  $z$ , goal: use  $\hat{z}_i$  to predict  $x_i$  if else.

Robots: 15 dim obs, 5 dim  $z$ , goal: sample from  $p(x)$

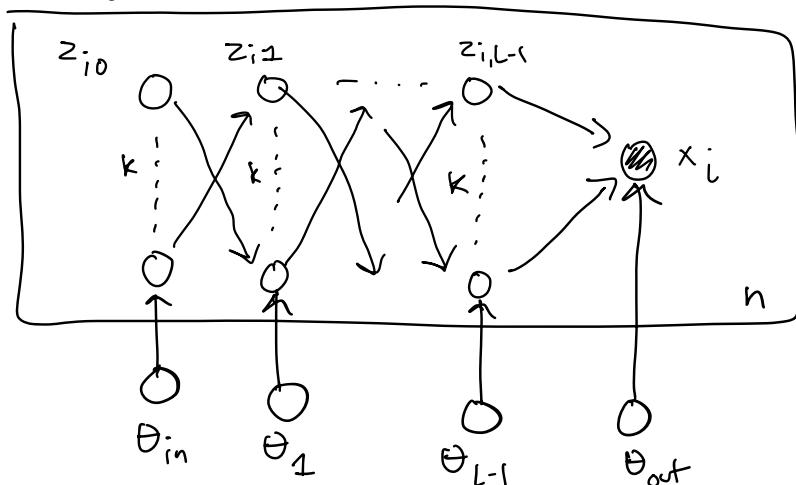
$$p(z_i | x_i, \hat{\theta}) = p(z_i) \quad \text{POSTERIOR COLLAPSE}$$

latent variable id. (Wang + 2020/1)

## Deep Exponential Families

Ranganath + 2013

everything becomes a random variable.



$$\begin{aligned} \mathbb{E}[z_{i,l,k} | \mathbf{z}_{i,L-1}, \theta_{l,k}] \\ = g(\theta_{l,k}, \mathbf{z}_{i,L-1}) \end{aligned}$$

$$z_{i,l,k} | \mathbf{z}_{i,L-1} \sim \text{expfam}(\eta_z(g(\theta_{l,k}, z_{i,L-1})))$$

$$x_i | \mathbf{z}_{i,L-1} \sim \text{expfam}(\eta_x(h(\theta_{out}, \mathbf{z}_{i,L-1})))$$

$p(\mathbf{z}_{1:n}, \theta | \mathcal{D}) = \text{BBVI}$ .   
deep latent Gaussian model

- ① how do DEFs stack against deep NNs.
- ② explore the "natural DEF"
- ③ what about conjugate pairs?

$$z_{i,k} | z_{i,l-1} \sim N(g(\theta_{ek} \cdot z_{i,l-1}), \sigma^2), \mathbb{E}[z_{i,k} | z_{i,l-1}] = g(\theta \cdot z_{i,l-1})$$

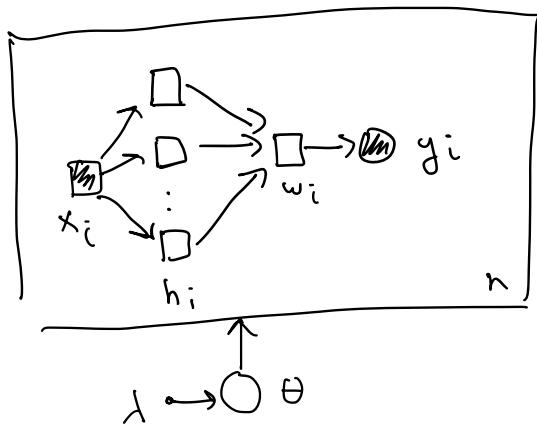
as  $\sigma^2 \rightarrow 0$ , DLGM approaches a NN.

### Sigmoid belief network

$$z_{i,k} | z_{i,l-1} \sim \text{Bernoulli}(\sigma(\theta_{ek} \cdot z_{i,l-1}))$$

$$\mathbb{E}[z_{i,k} | z_{i,l-1}] = \sigma(\cdot)$$

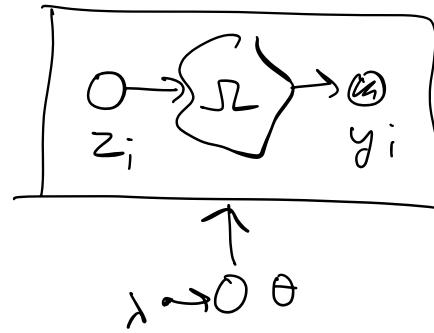
- deep exponential families
- discussion of DL
- architectural elements
- priors for neural networks
- BBVI ...



Supervised Bayesian NN  
MAP estimation. — Square loss  
x. entropy loss.

$$\mathcal{L}(x_i; \theta)$$

$$\nabla_{\theta} \mathcal{L}(x_i; \theta)$$



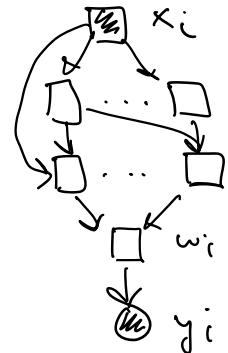
Deep generative model

- feature learners.
- Tensorflow, Pytorch, Jax.
- Automatic differentiation.

- activation function : ReLU
- intercept term :  $h_{i,ek} = g(\Theta_{ek} \cdot h_{i,l-1} + \alpha_{ek})$
- residual connections :

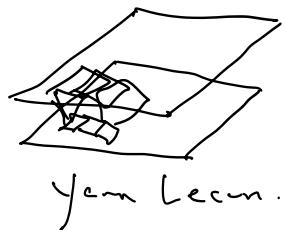
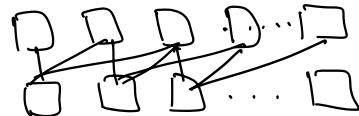
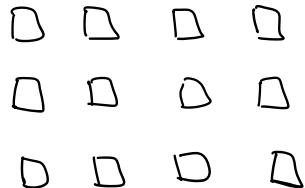
$$h_{i,ek} = g(\Theta_{ek} \cdot h_{i,l-1} + \gamma_{ek} \cdot x_i + \alpha_{ek})$$

you can implement hidden  $\rightarrow$  hidden residual connections too.  
"skip connections"

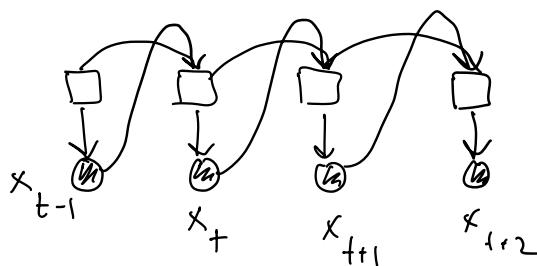


- parameter sharing.

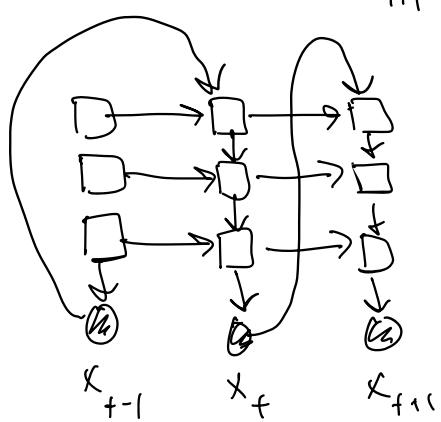
### ① convolution



- ② recurrence.



DNA data. methylation.  
text  
data

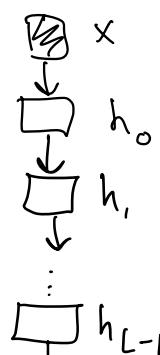


Sharing within each slice.

- Attention.

$$A(x_i; \alpha) : \mathbb{R}^P \rightarrow [0,1]^P$$

$$h_{i,ok} = g(\Theta_{in,k} \cdot (A(x_i; \alpha) \circ x_i))$$



$$x_i \rightarrow \begin{bmatrix} 0 \\ -1 \\ 0 \\ 0 \\ -1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\omega \downarrow y$$

$$x_i \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

The big dog ate my kid's gerbil  
and she was sad. "Why is she sad?"  
I said, she just ate a gerbil.

$$x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots \rightarrow x_T$$

=

$$y_i | x_i \sim N(\beta \cdot x_i, \sigma^2)$$

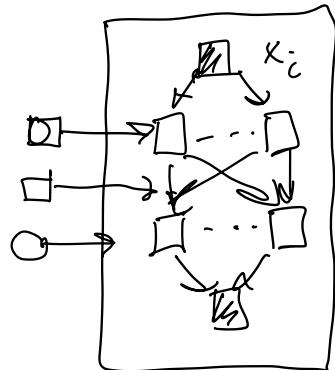
$$\sim N(\beta \cdot (\sigma(x \cdot x_i) \circ x_i), \sigma^2)$$

$\sigma(\cdot)$ : bank of logistic regressions.

$$\text{soft} \rightarrow \Delta^{P-1}$$

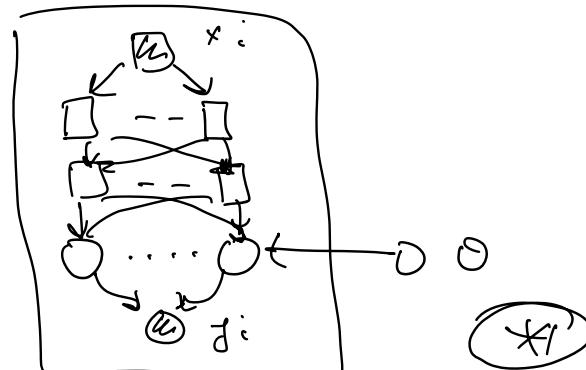
Bayesian N.N.

Automatic Relevance Determination.



feature learner

response model



## Black Box Variational Inference

complete conditionals  $\rightarrow$  Gibbs, CAVI

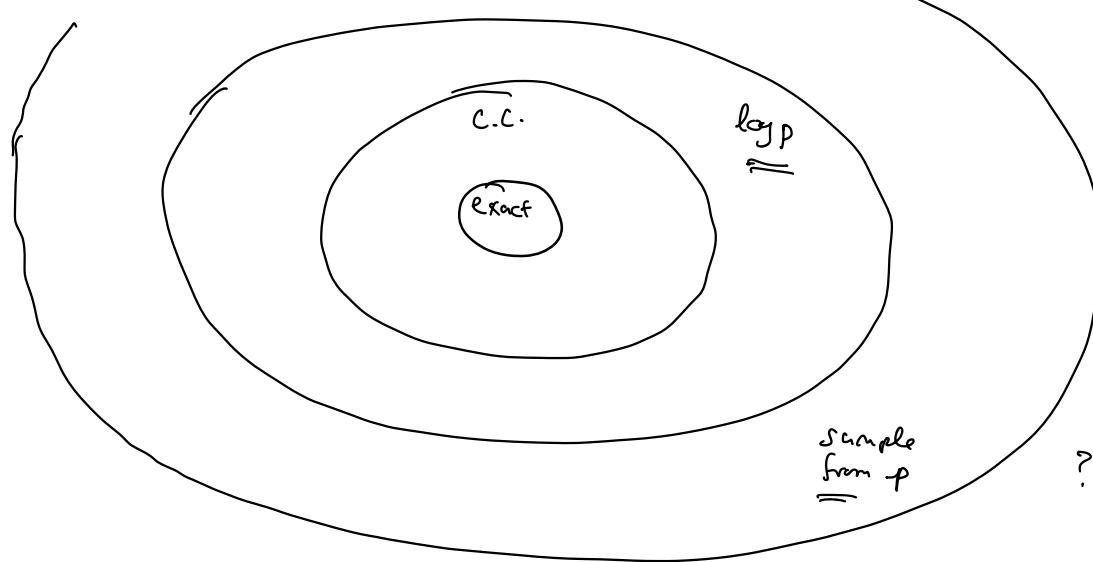
$$p(\theta, z_{1:n}, x_{1:n}) \quad \text{global } \theta \quad \text{local } z_{1:n}$$

$$= p(\theta) \prod_{i=1}^n p(z_i | \theta) p(x_i | z_i, \theta)$$

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

Black box criteria:

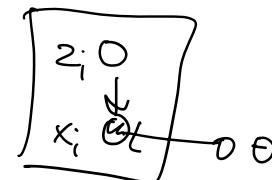
- compute the log joint (compute  $\nabla_{\theta, z} \log p(z, \theta, x)$ )
- compute things about  $q$
- sample from  $q$



Recall: Deep generative model.

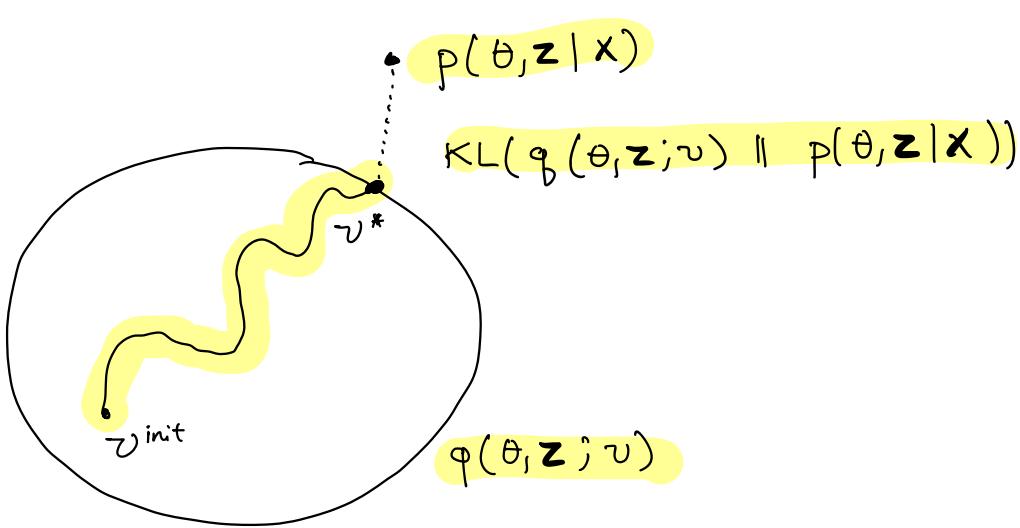
$$z_i \sim N_p(0, 1)$$

$$x_i | z_i, \theta \sim N(\Sigma_\mu(z_i), \Sigma_\sigma^2(z_i))$$



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

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



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  $\mathbb{E}_q [ \dots ]$
- ② Take a Monte Carlo approximation of  $\nabla_v \mathcal{L} \approx \frac{1}{B} \sum_b \dots_b$
- ③ Use stochastic optimization

Score gradient

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

$$\begin{aligned} & \nabla_v \log q(\theta, Z; v) \\ & \nabla_{\theta, Z} \log q(\theta, Z; v) \\ & \nabla \log q(\theta, Z; v) \end{aligned} \quad \left. \right\} \text{"Score function"}$$

## 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

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

$$q(\theta, z; v) = q(\theta; v_0) \prod_{i=1}^n q(z_i; v_i)$$

↑                      ↑  
 N                      N

In practice:

Variance is high. ∫ Andy Miller, Justin Domke, ...  
control variates.

Let  $g(\theta, z; v) \stackrel{\text{stochastic}}{\triangleq} \text{Score gradient}, \mathbb{E}[g(\theta, z; v)] = \nabla_v \mathcal{L}(v)$

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

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

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

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

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

FACT:

$$\nabla_v q(z; v) = q(z; v) \nabla_v \log q(z; v)$$

## Reparameterization gradient

transform the variables in  $q$ .

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

e.g. for a normal:

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

## Reparameterization gradient

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

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

where:  $\theta = t(\varepsilon_\theta, v_\theta)$ ,  $z = t(\varepsilon_z, v_z)$

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

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