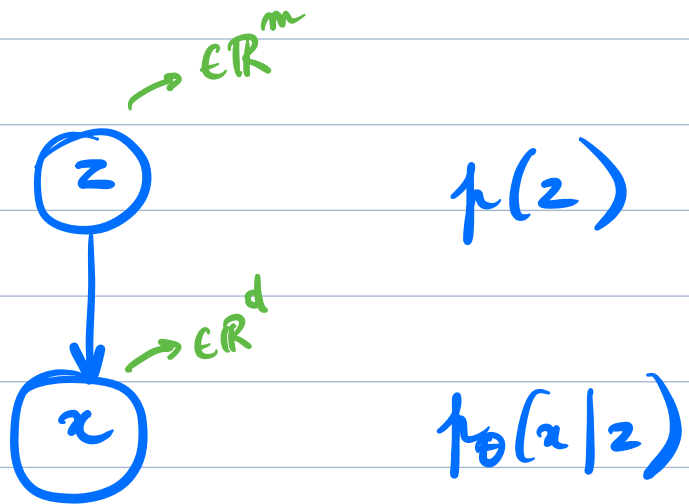
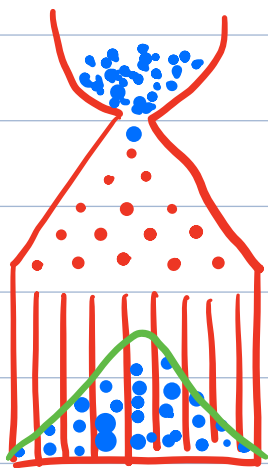


Deep Latent variable models



How to estimate θ ?

$$\begin{aligned} \max_{\theta} p_{\theta}(x) &= \int p_{\theta}(x, z) dz \\ &= \int p(z) p_{\theta}(x|z) dz \\ &= \mathbb{E}_{p(z)} [p_{\theta}(x|z)] \end{aligned}$$

$$\approx \frac{1}{K} \sum_k p_{\theta}(x|z_k)$$

MC approximation

Poor approximation when m is large due to the curse of dimensionality!

I. Variational inference

$$p_{\theta}(x) = \mathbb{E}_{p(z)} [p_{\theta}(x|z)]$$

! Inference

$$= \mathbb{E}_{q_\phi(z)} \left[\frac{p(z)}{q_\phi(z)} \log p_\theta(z|x) \right] \quad \downarrow \text{importance sampling}$$

$$\log p_\theta(x) = \log \mathbb{E}_{q_\phi(z)} \left[\frac{p(z)}{q_\phi(z)} \log p_\theta(z|x) \right]$$

Jensen's inequality $\mathbb{E}[\log] \geq \log \mathbb{E}$

$$\geq \mathbb{E}_{q_\phi(z)} \left[\log \frac{p(z)}{q_\phi(z)} \log p_\theta(z|x) \right]$$

ELBO_{VI}

$$\mathbb{E}_q \left[\log \frac{p(x,z)}{q(z)} \right] = \mathbb{E}_{q_\phi(z)} \left[\log p_\theta(z|x) \right] - \text{KL}(q_\phi(z) \| p(z))$$

$$= \mathbb{E}_{q_\phi(z)} \left[\log \frac{p(z)}{q_\phi(z)} \log p_\theta(z|x) \frac{p_\theta(z|x)}{p_\theta(z|x)} \right]$$

Reconstruction error

$$= \mathbb{E}_{q_\phi(z)} \left[\log \frac{p_\theta(z|x)}{q_\phi(z)} \log p_\theta(x) \right]$$

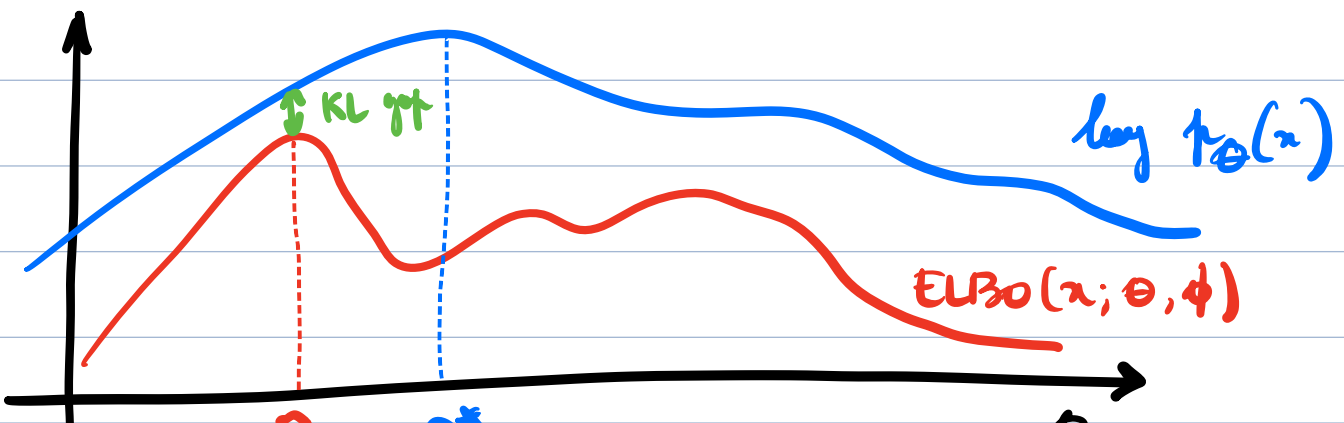
→ When $p(z|x)$ is Gaussian, this results in MSE

$$= \log p_\theta(x) - \text{KL}(q_\phi(z) \| p_\theta(z|x))$$

① + ②:

$$\text{ELBO}_{VI}(x; \theta, \phi) = \log p_\theta(x) - \text{KL}(q_\phi(z) \| p_\theta(z|x)) \geq 0$$

⇒ Maximizing $\log p_\theta(x)$ (is the same as max and min).



We want $KL(q_\phi(z) \parallel p_\theta(z|x)) \rightarrow 0$, otherwise the gap is large and $|\hat{\theta}^* - \hat{\theta}| \gg 0$.

\Rightarrow Put enough capacity in q_ϕ .

II. VAEs

Amortize inference for any x .

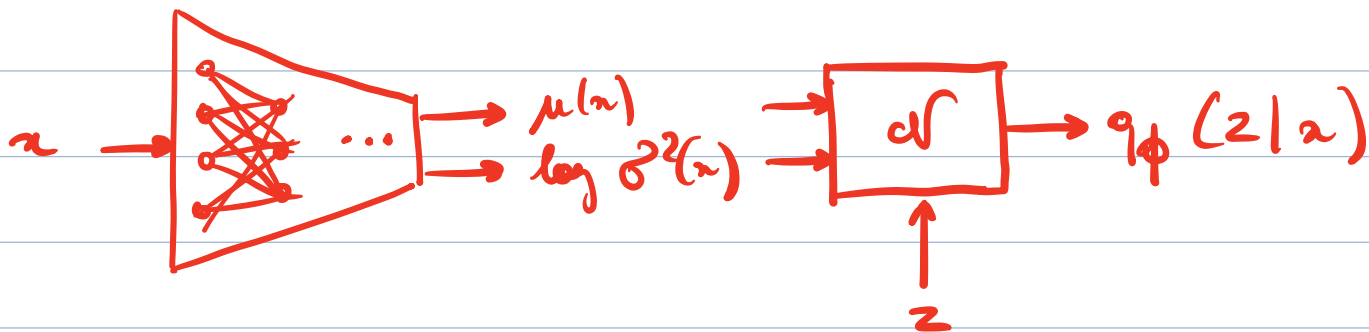
$q_\phi(z)$



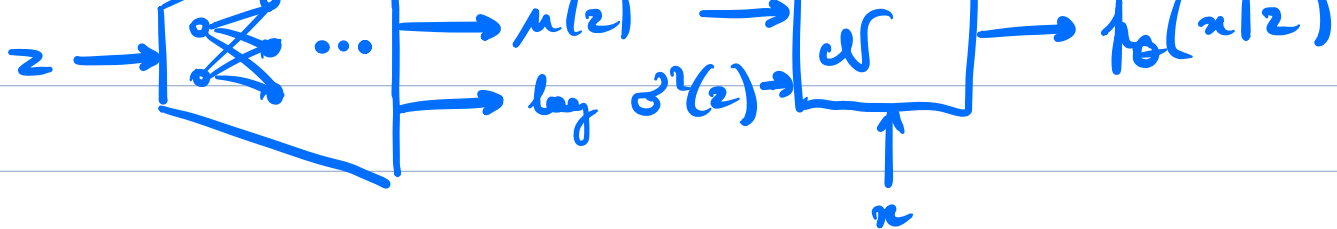
$q_\phi(z|x)$

Now a DNN that outputs the parameters of the variational dist.

Encoder



Decoder



Training

$$\max_{\theta, \phi} \mathbb{E}_{p(z)} [\text{ELBO}(z; \theta, \phi)]$$

$$= \mathbb{E}_{p(z)} \left[\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x) \| p(z)) \right]$$



Issue: posterior collapse
when $\text{KL} = 0$.

→ x does not bring
any info about z

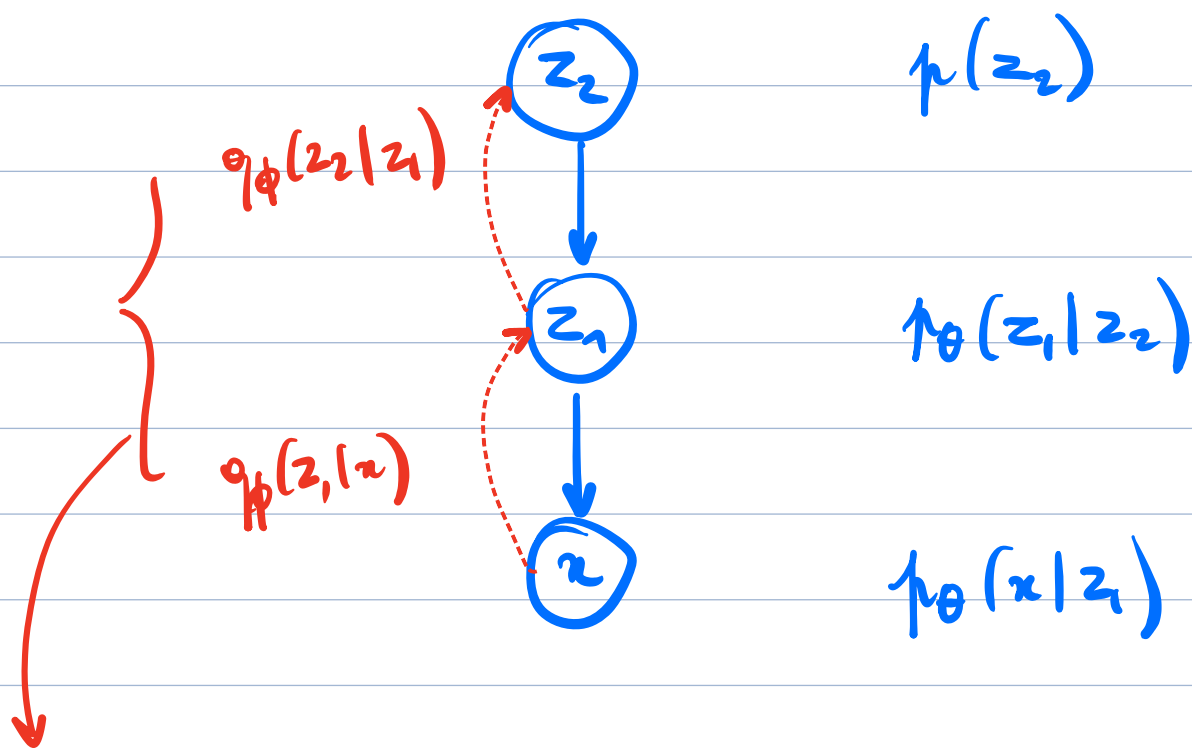
→ z is pure noise



Tension between $\text{KL}(q_{\phi}(z|x) \| p(z))$
and $\text{KL}(q_{\phi}(z|x) \| p(z|x))$

→ Slides + Code.

II. Hierarchical VAEs



$$q_\phi(z_1, z_2 | x) = q_\phi(z_1 | x) q_\phi(z_2 | z_1)$$

Training

$$\begin{aligned} \max_{\phi, \theta} \mathbb{E}_{p(x)} [\text{ELBO}(x)] & \quad \rightarrow p(x|z_1) p(z_1|z_2) p(z_2) \\ &= \mathbb{E}_{p(x)} \mathbb{E}_{q_\phi(z_1, z_2 | x)} \left[\log \frac{p(x, z_1, z_2)}{q(z_1, z_2 | x)} \right] \\ &\stackrel{\text{do it!}}{=} \mathbb{E}_{p(x)} \mathbb{E}_{q_\phi(z_1, z_2 | x)} \left[\log p_0(x|z_1) \right. \\ &\quad \left. - \text{KL}(q_\phi(z_1|x) \| p(z_1|z_2)) \right. \\ &\quad \left. - \text{KL}(q_\phi(z_2|z_1) \| p(z_2)) \right] \end{aligned}$$

Some as done but with this new term

→ 0 when $q_\phi(z_2|z_1)$ has too much capacity
 $\Rightarrow q_\phi(z_2|z_1) = p(z_2) = \mathcal{U}$

no info about x in z_2

\Rightarrow the second layer is not used!
 \Rightarrow Some \Rightarrow regular VAEs!

Top-down VAEs

! Swap the directional dependencies between the latents

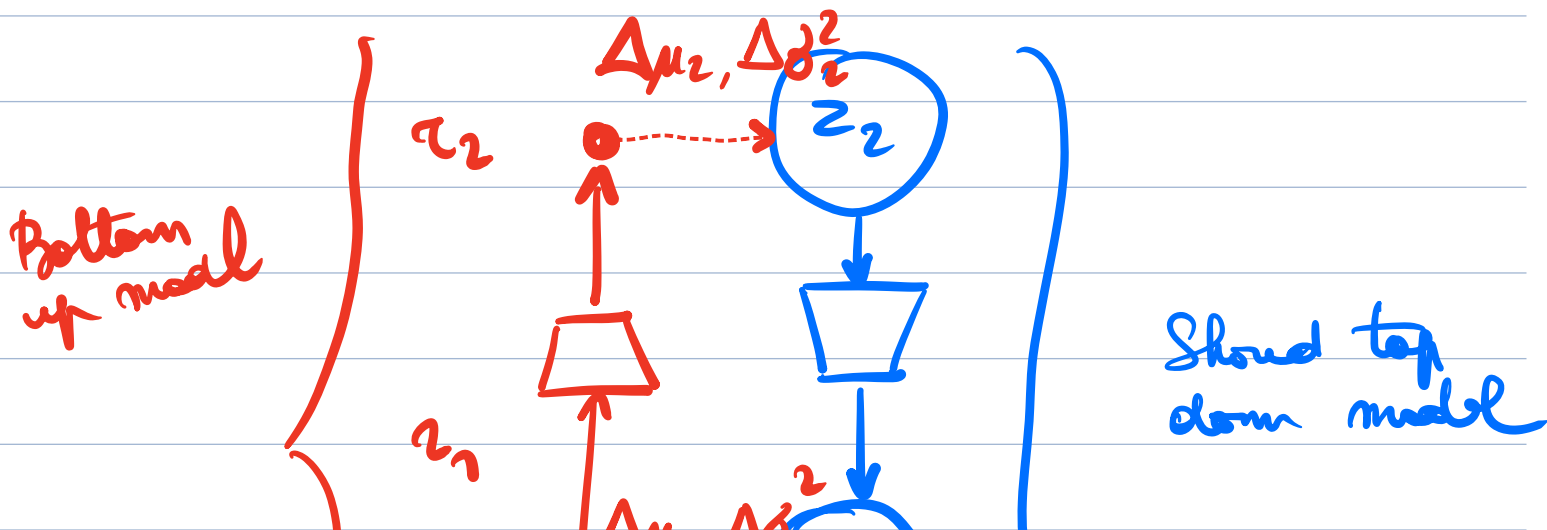
force injecting info about x in z_2

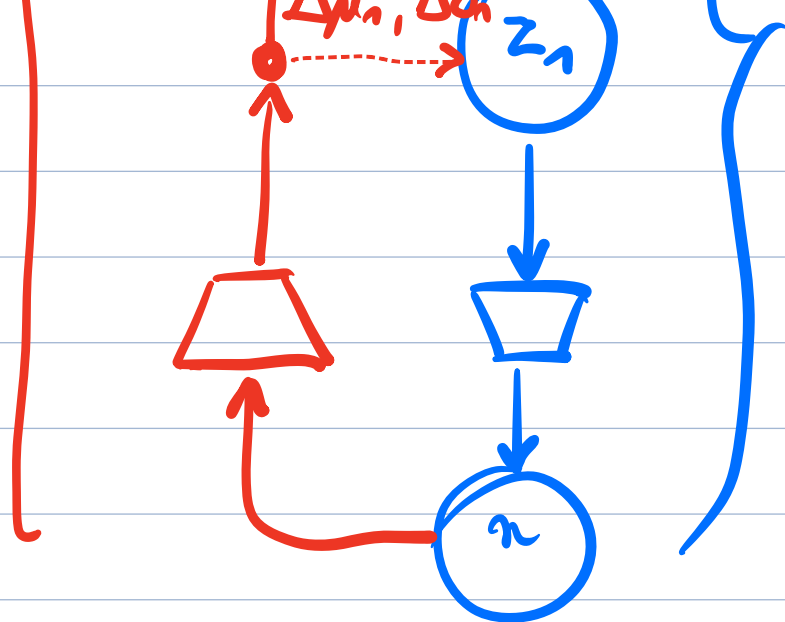
$$q_{\phi}(z_1, z_2 | x) = q_{\phi}(z_2 | x) q_{\phi}(z_1 | z_2, x)$$

! Some directed dependencies as in $p_{\theta}(z_1 | z_2)$

INDUCTIVE BIAS

\Rightarrow Show a common top-down path.





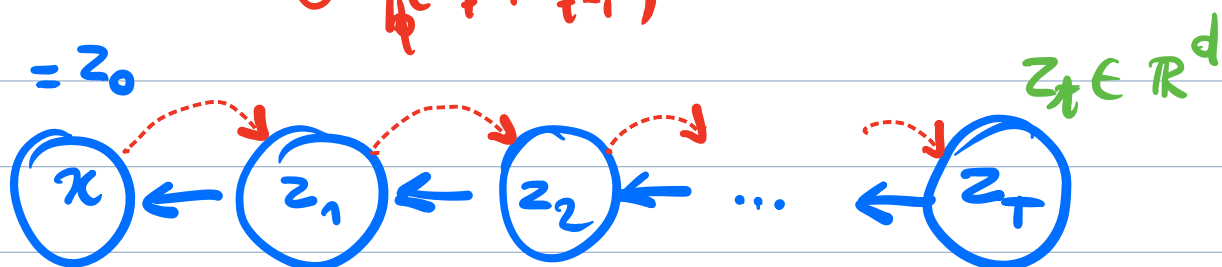
Focus on close connection between q_ϕ and p_θ and helps encode information about x in z_2, z_1 .

→ Slides.

IV. Deep diffusion probabilistic models

Diffusion

① W encoder
 \uparrow
 $q_\phi(z_t | z_{t-1})$



$z \in \mathbb{R}^k$

$$p_\theta(z_t | z_{t+1}) \rightarrow$$

DNN decoder

$$p(z_T) := \mathcal{U}$$
$$T = 0(1000)$$



Big first-order Markov chain

← Generation = Reverse diffusion

→ Slow diffusion.

$$p_\theta(z_{0:T}) = \left[\prod_{t=0}^{T-1} p_\theta(z_t | z_{t+1}) \right] p(z_T)$$

$$q_\phi(z_{1:T} | z_0) = \prod_{t=1}^T q_\phi(z_t | z_{t-1})$$

$$q_\phi(z_t | z_{t-1}) = \mathcal{U}(z_t | \sqrt{1 - \beta_t} z_{t-1}; \beta_t \mathbb{I})$$

$$\Leftrightarrow z_t := \sqrt{1 - \beta_t} z_{t-1} + \beta_t \epsilon, \quad \epsilon \sim \mathcal{U}(0, 1)$$

Training

$$\max_{\phi, \theta} \mathbb{E}_{p(z_0)} [\text{ELBO}(z)]$$

$$= \mathbb{E}_{p(z_0)} \mathbb{E}_{q_\phi(z_{1:T} | z_0)} \left[\log \frac{p_\theta(z_{0:T})}{q_\phi(z_{1:T} | z_0)} \right]$$

$$= \mathbb{E}_p(z_0) \mathbb{E}_{q_\beta(z_{1:T}|z_0)} \left[\log p_0(z_0|z_1) - \sum_{t=1}^T \text{KL}(q_\beta(z_t|z_{t-1}) \| p(z_t|z_{t+1})) - \text{KL}(q_\beta(z_T|z_{T-1}) \| p(z_T)) \right]$$

He et al, 2020:

[SKIP]
if time is short.

① Since q is linear Gaussian, we have

$$q(z_t|z_0) = \mathcal{N}(z_t | \sqrt{\bar{\alpha}_t} z_0, 1 - \bar{\alpha}_t \mathbf{I})$$

where $\alpha_t = 1 - \beta_t$
 $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$

$\Rightarrow z_t$ can be sampled without all intermediate steps!

$$\textcircled{2} \quad q(z_{t-1}|z_t, z_0) = \frac{q(z_t|z_{t-1}, z_0) q(z_{t-1}|z_0)}{q(z_t|z_0)}$$

\Downarrow \mathcal{N} (conjugate prior) \Downarrow \mathcal{N}

$$= \mathcal{N}(z_{t-1} | \tilde{\mu}_t(z_t, z_0), \tilde{\beta}_t \mathbf{I}) \quad \frac{1 - \bar{\alpha}_{t+1} \beta_t}{1 - \bar{\alpha}_t} \beta_t$$

$$\frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} z_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} z_t$$

Then, the ELBO can be expressed as

$$\mathbb{E}_{p(z_0)} \mathbb{E}_{q(z_{1:T}|z_0)} \left[\sum_T \log p_\theta(z_0|z_1) \rightarrow z_0 \rightarrow z_t \right. \\ \left. - \sum_{t=1} \text{KL}(q_\phi(z_{t+1}|z_t, z_0) \parallel p_\theta(z_{t+1}|z_t)) \rightarrow z_t \right. \\ \left. - \text{KL}(q_\phi(z_T|z_0) \parallel p(z_T)) \right]$$

Training

$\text{KL}(d \parallel d)$
 cloud form!

$$\mathbb{E}_{t \sim [1..T]} \mathbb{E}_{p(z_0)} \mathbb{E}_{q(z_t|z_0)} [\mathcal{L}_t]$$

- Update the layers one at a time!
- much more memory efficient
- Scale to $T = 0$ (1000)

found potential mem

→ slides

$$\mathcal{L}_{t-1} = \mathbb{E}_q \left[\frac{1}{2\sigma_t^2} \left\| \tilde{\mu}_t(z_t, z_0) - \mu_\theta(z_t, t) \right\|^2 \right]$$

DMT

V. VAE prior

$$ELBO = \mathbb{E}_{p(z)} \mathbb{E}_{q_{\phi}(z|x)} \left[\underbrace{\log p_{\theta}(x|z)}_{\text{Reconstruction}} + \underbrace{\log \frac{p(z)}{q_{\phi}(z|x)}}_{\Omega} \right]$$

$$\begin{aligned} \Omega &= \mathbb{E}_{p(x)} \mathbb{E}_{q_{\phi}(z|x)} \left[\log p(z) - \log q_{\phi}(z|x) \right] \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_{\phi}(z|x_n)} \left[\log p(z) - \log q_{\phi}(z|x_n) \right] \\ &= \mathbb{E}_{p(z)} \left[\log p(z) \right] + \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{q_{\phi}(z|x_n)} \left[-\log q_{\phi}(z|x_n) \right] \\ &= \underbrace{-CE(q_{\phi}(z) || p(z))}_{N \text{ non-entropies}} + \underbrace{H_{p(z)}[q_{\phi}(z|x)]}_{\text{entropies}} \end{aligned}$$

$$q_{\phi}(z) = \frac{1}{N} \sum_{n=1}^N q_{\phi}(z|x_n)$$

aggregated posterior

Make the aggregated posterior match with the prior.

Make the posterior narrower $\rightarrow \infty$, but counter-balanced with RE.

hole!

Difficult when



