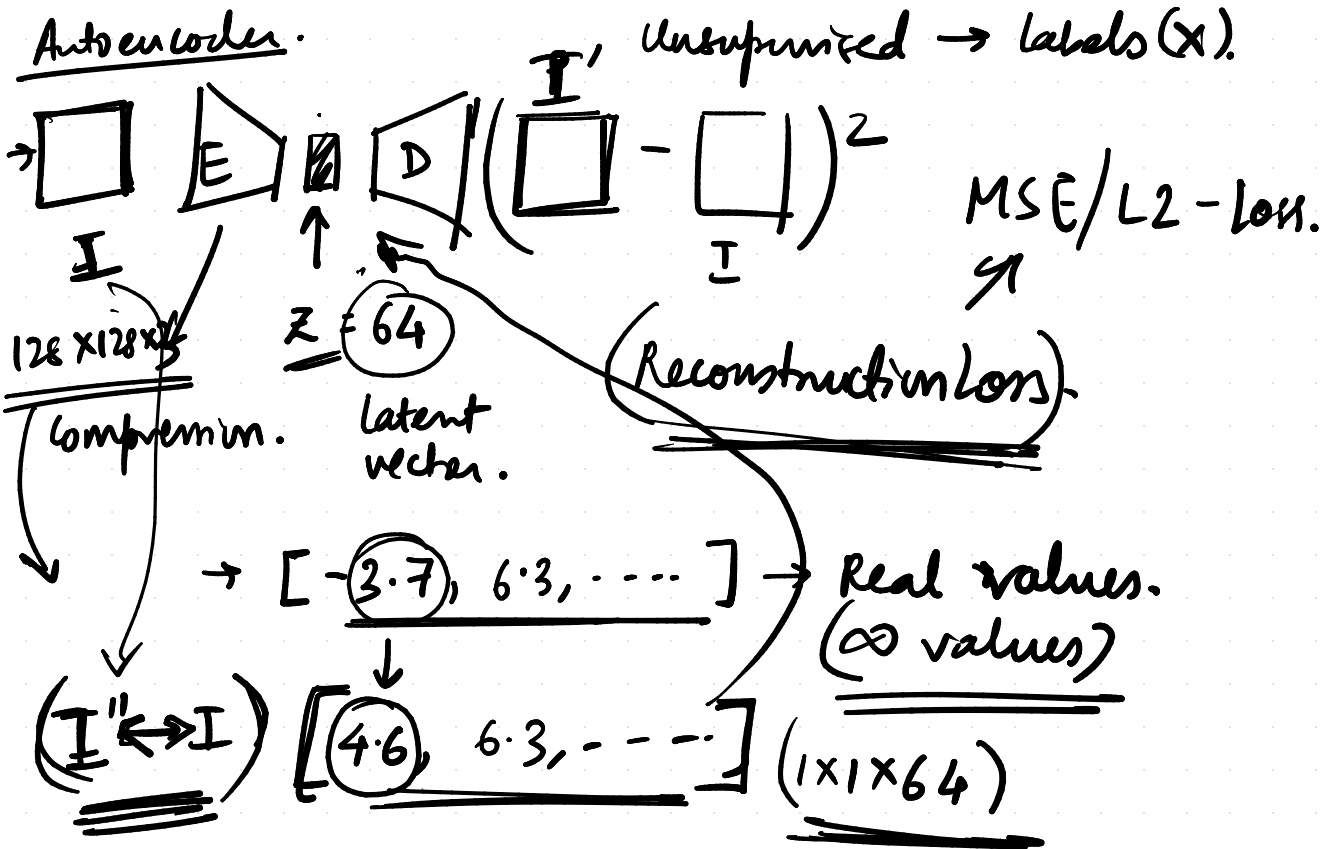
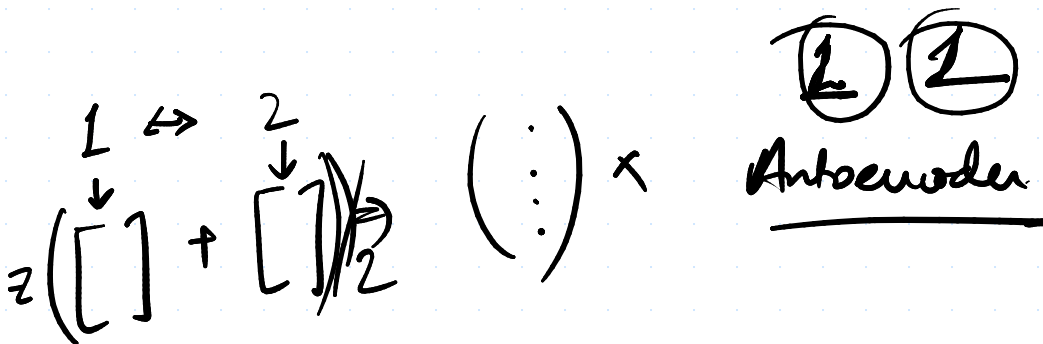
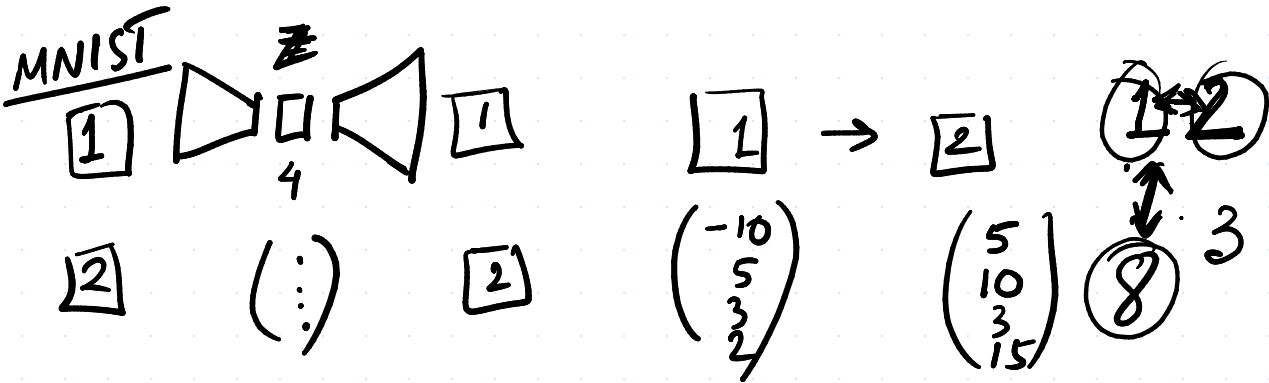


VAE - Theory Derivations & Insights.

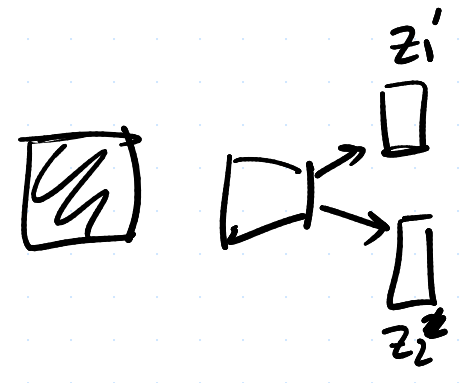
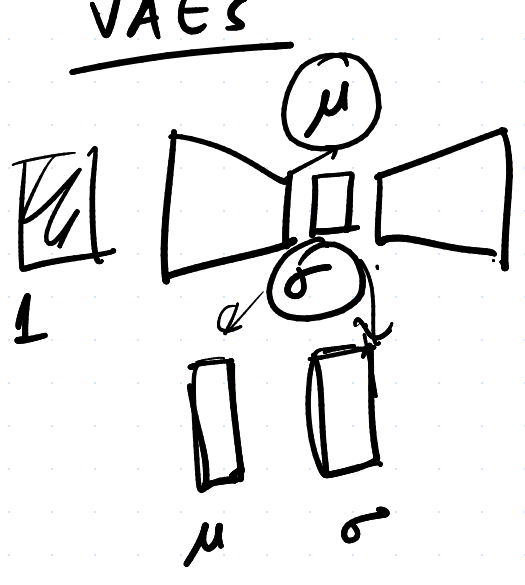
Autoencoder.



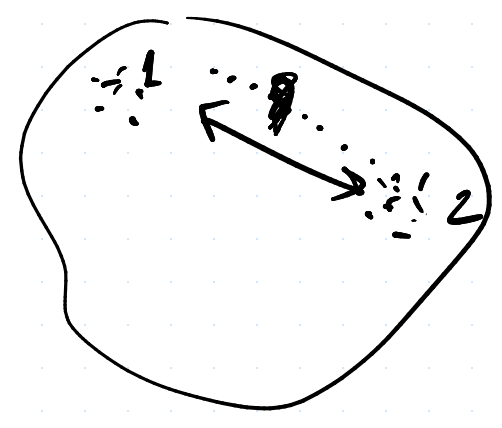
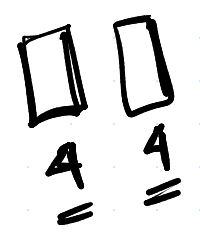
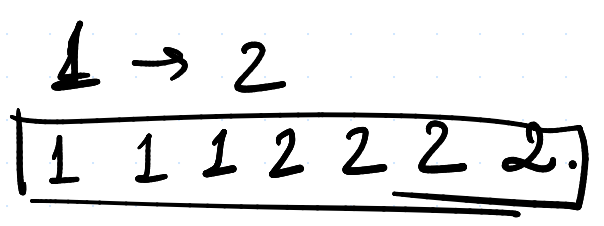
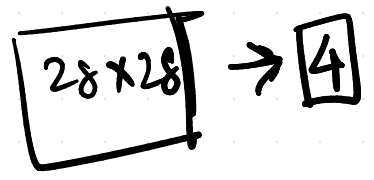
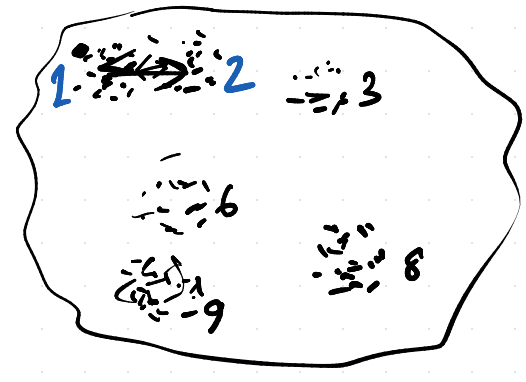
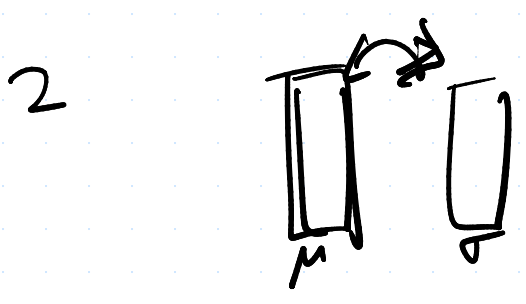
MNIST



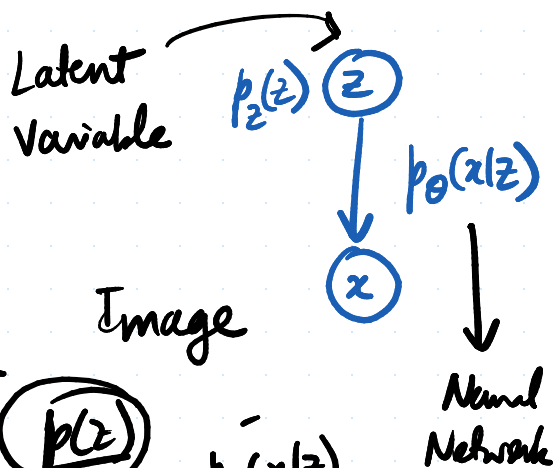
VAE's



$L1/L2 \quad \frac{\|z_1 - z_2\|}{\|z_1\| \|z_2\|} =$
 $z_1 \quad ([z_1^{(1)}] [z_2^{(1)}])$
 z_2



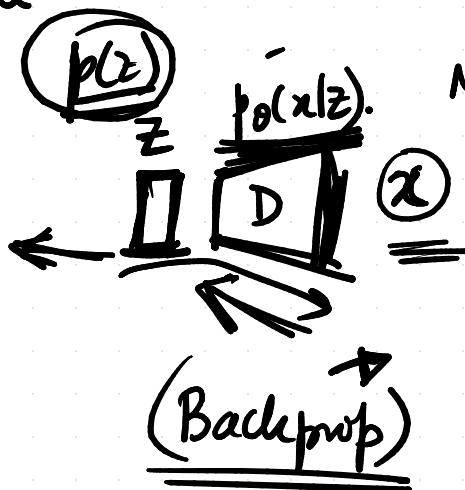
$$\left. \begin{aligned} x &\sim p_z(z) \\ x &\sim p_\theta(x|z) \end{aligned} \right\} \text{sample.}$$



Likelihood.

$$p_\theta(x) = \sum_z p_z(z) p_\theta(x|z)$$

easy to sample from
Marginalization



Training objective

$$z = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}$$

$$\max_{\theta} \sum_i \log p_\theta(x^{(i)})$$

$$= \sum_i \log \sum_z p_z(z) p_\theta(x^{(i)}|z)$$

$z \rightarrow$ can take many values.

Hard to sample z that leads to meaningful forms.

$$\sum_i \log \sum_z p_z(z) p_\theta(x^{(i)}|z)$$

$$\approx \sum_i \log \frac{1}{k} \sum_{k=1}^k p_\theta(x^{(i)}|z_k^{(i)})$$

$$z_k^{(i)} \sim p_z(z)$$

$z \rightarrow$ takes many values, hard to select samples that leads to meaningful tune.

Importance Weighted VAE (IWVAE)

$$\textcircled{z} \rightarrow \textcircled{x}$$

$z =$ discrete R.V.
(small domain).

But in general

$z \rightarrow$ CRV. & hence

the objective is not
tractable.

exact objective
is tractable.

$$\{ p_z(z=A) = p_z(z=B) = p_z(z=C) = \frac{1}{3}$$

May send
(Endsem)

$$p_\theta(x|z=k) = \frac{1}{(2\pi)^{n/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)\right)$$

Training objective :-

$$\max_{\theta} \sum_i \log p_\theta(x^{(i)})$$

$$= \sum_i \log \sum_z p_z(z=z) p_\theta(x|z=z)$$

$$= \frac{1}{3} p_z(z=A) p_\theta(x|z=A) + \frac{1}{3} p_z(z=B) p_\theta(x|z=B) + \frac{1}{3} p_z(z=C) p_\theta(x|z=C)$$

$$\begin{aligned}
&= \max_{\mu, \Sigma} \sum_i \log \left[\frac{1}{3} \frac{1}{(2\pi)^{n/2}} \frac{1}{|\Sigma_A|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_A)^T \Sigma_A^{-1} (x^{(i)} - \mu_A)\right) \right. \\
&\quad + \frac{1}{3} \frac{1}{(2\pi)^{n/2}} \frac{1}{|\Sigma_B|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_B)^T \Sigma_B^{-1} (x^{(i)} - \mu_B)\right) \\
&\quad \left. + \frac{1}{3} \frac{1}{(2\pi)^{n/2}} \frac{1}{|\Sigma_C|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_C)^T \Sigma_C^{-1} (x^{(i)} - \mu_C)\right) \right]
\end{aligned}$$

Importance Sampling

$$\mathbb{E}_{z \sim p_z(z)} [f(z)] = \sum_z p_z(z) f(z)$$

$$= \sum_z \frac{q(z)}{q(z)} \cdot p_z(z) \cdot f(z).$$

$$= \mathbb{E}_{z \sim q(z)} \left[\frac{p_z(z) f(z)}{q(z)} \right]$$

Here we will sample from $q(z)$ rather than $p(z)$, sample based expectations

$$\approx \frac{1}{K} \sum_{k=1}^K \frac{p_z(z^{(k)})}{q(z^{(k)})} f(z^{(k)}) \quad \text{with } z^{(k)} \sim q(z).$$

We can sample from q to compute expectations w.r.t. p .

Train objective (new) :-

$$\max_{\theta} \sum_i \log p_{\theta}(x^{(i)}) = \sum_i \log \sum_z p_z(z) p_{\theta}(x^{(i)}|z)$$

$$\approx \sum_i \log \frac{1}{K} \sum_{k=1}^K \frac{p_z(z_k^{(i)})}{q(z_k^{(i)})} \underbrace{p_{\theta}(x^{(i)}|z_k^{(i)})}_{f(z)}$$

hence, we use importance sampling.

$$z_k^{(i)} \sim q(z_k^{(i)})$$

$$q(z) = p_{\theta}(z|x^{(i)}) = \frac{p_{\theta}(x^{(i)}|z) p_z(z)}{p_{\theta}(x^{(i)})}$$

Not easy to sample from.

$q(z) = \mathcal{N}(z; \mu, \sigma^2) \rightarrow$ easy to sample from.

Importance sampling \rightarrow Proposal Distribution

$$\Rightarrow \min_{q(z)} \text{KL}(q(z) \parallel p_{\theta}(z|x^{(i)})) \leftarrow \underline{\text{Distance}} \text{ b/w 2 dist.}$$

$$\Rightarrow \min_{q(z)} \sum_i q_i(z) \cdot \log\left(\frac{q(z)}{p_{\theta}(z|x^{(i)})}\right)$$

$$\Rightarrow \min_{q(z)} \mathbb{E}_{z \sim q(z)} \log\left(\frac{q(z)}{p_{\theta}(z|x^{(i)})}\right)$$

$$\Rightarrow \min_{q(z)} \mathbb{E}_{z \sim q(z)} \log\left(\frac{q(z)}{p_{\theta}(z|x^{(i)}) p_z(z)/p_{\theta}(x^{(i)})}\right) \quad q \rightarrow \text{gaussian}$$

$$\Rightarrow \min_{q(z)} \mathbb{E}_{z \sim q(z)} [\log q(z) - \log p_z(z) - \log p_\theta(x^{(i)}|z)]$$

constant independent of z .

$+ \log p_\theta(x^{(i)})$
doesn't depend on z



$$\Rightarrow \min_{q(z)} \mathbb{E}_{z \sim q(z)} [\log q(z) - \log p_z(z) - \log p_\theta(x^{(i)}|z)]$$

sample from q by design.

optimize to find q

easy to work with

Neural Network.

(SGLD)

Amortized Inference

$x^{(i)} \rightarrow$ want to solve. $N(\mu^{(i)}, \sigma^{(i)})$

$$\min_{q(z)} KL(q(z) \parallel p_\theta(z|x^{(i)}))$$

for all $x^{(i)}$, we want to solve.

$$\phi: x^{(i)} \rightarrow \mu^{(i)}, \sigma^{(i)}$$

Amortized formulation

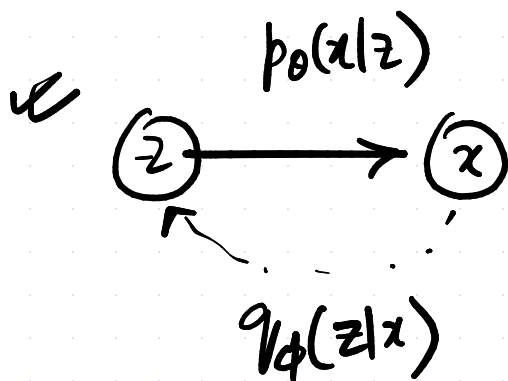
$$\min_{\phi} \sum_i \text{KL} (q_{\phi}(z|x^{(i)}) || p_{\theta}(z|x^{(i)}))$$

faster, regularization (not as precise).

Gaussian

$$\min_{\phi} \sum_i \text{KL} (q_{\phi}(z|x^{(i)}) || p_{\theta}(z|x^{(i)}))$$

↑
multiple posteriors.

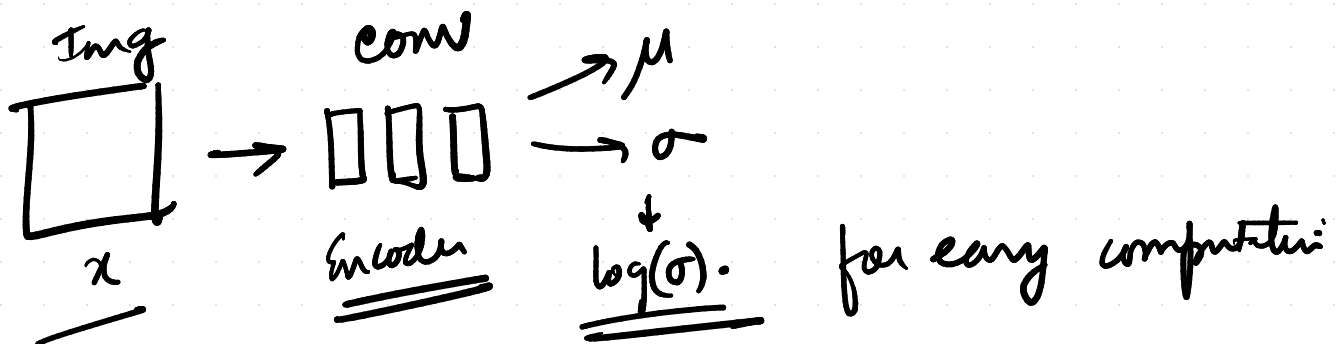


→ part of our inference model, not a part of generative model.

$$q_{\phi}(z|x) = \mathcal{N}(\mu_{\phi}(x), \sigma_{\phi}^2(x))$$

$$\mathbb{Z} = \underline{\underline{\mu_{\phi} + \varepsilon \sigma_{\phi}(x)}} \quad \text{with } (\varepsilon \sim \underline{\underline{\mathcal{N}(0, I)}})$$

$$z \sim q(z|x) = \mathcal{N}(z; \mu_{\phi}^{(x)}, \sigma_{\phi}^{(x)})$$



Importance - Weighted A.E. (IWAE)

$$\text{objective} \quad \sum_i \log \frac{1}{k} \sum_{k=1}^k \frac{p_z(z_k^{(i)})}{q(z_k^{(i)})} p_{\theta}(x^{(i)}|z_k^{(i)})$$

$$\text{with } z_k^{(i)} \sim q(z_k^{(i)})$$

$$\min_{\phi} \sum_i \text{KL}(q_{\phi}(z|x^{(i)}) \parallel p_{\theta}(z|x^{(i)}))$$

maximize $(\text{term 1} - \text{term 2}) = \mathcal{L}_k$ - Likelihood.
 θ, ϕ

For all k , the lower bounds satisfy :-

$$\text{real objective} \leftarrow \log p(x) \geq \mathcal{L}_{k+1} \geq \mathcal{L}_k.$$

Moreover, if $p(z|x)/q(z|x)$ is bounded

then $\mathcal{L}_k \rightarrow \log p(x)$ (approaches)

as $k \rightarrow \infty$.

As we have more terms inside the importance sampling, we will do better.

$$\mathcal{L}_k = \mathbb{E}_{\underline{h_1, h_2, \dots, h_k}} \left[\log \frac{1}{k} \sum_{i=1}^k \frac{p(x, h_i)}{q(h_i|x)} \right]$$

$$= \mathbb{E}_{h_1, h_2, \dots, h_k} \left[\log \mathbb{E}_{\mathbf{I} = \{i_1, i_2, \dots, i_m\}} \left[\frac{1}{m} \sum_{j=1}^m \frac{p(x, h_{ij})}{q(h_{ij}|x)} \right] \right]$$

$$\geq \mathbb{E}_{h_1, h_2, \dots, h_k} \left[\mathbb{E}_{\mathbf{I} = \{i_1, i_2, \dots, i_m\}} \left[\log \frac{1}{m} \sum_{j=1}^m \frac{p(x, h_{ij})}{q(h_{ij}|x)} \right] \right]$$

$$= \mathbb{E}_{h_1, h_2, \dots, h_m} \left[\log \frac{1}{m} \sum_{i=1}^m \frac{p(x, h_i)}{q(h_i|x)} \right] = \mathcal{L}_m$$

Variational Lower Bound Derivation $\rightarrow 1$

(ELBO)

(Using Jensen)

$$\max_{\theta} \sum_i \log p_{\theta}(x^{(i)})$$

$$= \max_{\theta} \sum_i \log \left(\sum_z p_z(z) p_{\theta}(x^{(i)}|z) \right)$$

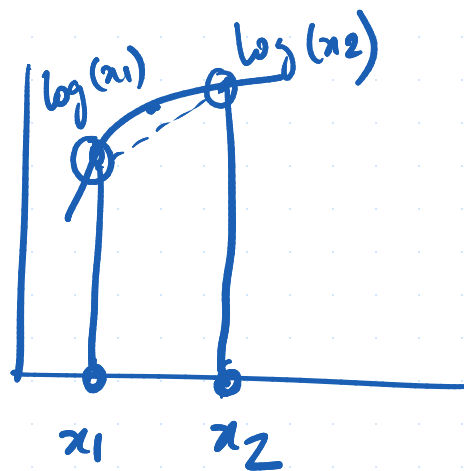
$$= \max_{\theta} \sum_i \log \left(\sum_z \frac{q_i(z)}{q(z)} \underbrace{p_z(z) \cdot p_{\theta}(x^{(i)}|z)} \right)$$

Importance Sampling

$$= \max_{\theta} \sum_i \log E_{z \sim q_i(z)} \frac{p_z(z)}{q_i(z)} p_{\theta}(x^{(i)}|z).$$

$$\geq \max_{\theta} \sum_i E_{z \sim q_i(z)} \left\{ \log \frac{p_z(z)}{q_i(z)} \cdot p_{\theta}(x^{(i)}|z) \right\}$$

constant, then Equality



$$\log E_z(a) \geq \underline{E_z(\log(x))}$$

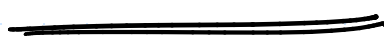
$$= \max_{\theta} \sum_i \mathbb{E}_{z \sim q(z)} \log p_z(z) +$$

$$\mathbb{E}_{z \sim q(z)} \log p_{\theta}(x^{(i)}|z)$$

$$- \mathbb{E}_{z \sim q(z)} \log q(z).$$

$q \rightarrow$ something we can choose,

$$q(z) \propto p_z(z) p_{\theta}(x^{(i)}|z)$$



\hookrightarrow un-normalized dist.

$p_{\theta}(z|x^{(i)}) \rightarrow$ normalized.



$$\Rightarrow \max_{\theta, q} \sum_i \mathbb{E}_{z \sim q(z)} \log p_i(z) +$$

$$\mathbb{E}_{z \sim q(z)} \log p_\theta(z^{(i)} | z)$$

$$\mathbb{E}_{z \sim q(z)} \log q(z).$$

Evidence Lower Bound (ELBO)
 Variational Lower Bound (VLB)

Derivation using KL

$$D_{KL} [q_\alpha(z) \parallel p(z|\alpha)] = \mathbb{E}_{z \sim q_\alpha(z)} \left[\log q_\alpha(z) - \right.$$

conditional $\left. \log p(z|\alpha) \right]$

$$= \mathbb{E}_{z \sim q_\alpha(z)} \left[\log q_\alpha(z) - \log \frac{p(z, \alpha)}{p(\alpha)} \right]$$

↓ joint

$$= \mathbb{E}_{z \sim q_{\alpha}(z)} \left[\log q_{\alpha}(z) - \log p(z) - \log p(x|z) + \log p(x) \right]$$

$$= \mathbb{E}_{z \sim q_{\alpha}(z)} \left[\log q_{\alpha}(z) - \log p(z) - \log p(x|z) \right]$$

only this part depends
on z , (VLB)

+ $\log p(x)$

$$D_{KL} [q_{\alpha}(z) \parallel p(z|x)] = \mathbb{E}_{z \sim q_{\alpha}(z)} \left[\log q_{\alpha}(z) - \log p(z) - \log p(x|z) \right] + \log p(x)$$

$$\log p(x) = \mathbb{E}_{z \sim q_{\alpha}(z)} \left[-\log q_{\alpha}(z) + \log p(z) \right]$$

$$\underline{\underline{VLB}} \left\{ + \log p(x|z) \right\} + \underbrace{D_{KL} [q_{\alpha}(z) \parallel p(z|x)]}_{\geq 0}$$

Same with Jensen's, but now we know
the gap = KL.

Stochastic Optimization (Objective amenable
to this)

$$\log p(x) = -\mathbb{E}_{z \sim q_x(z)} [\log q_x(z) - \log p(z) - \log p(x|z)]$$

$$+ D_{KL} [q_x(z) \parallel p(z|x)]$$

$$= \mathbb{E}_{z \sim q_x(z)} [\log p(z) + \log p(x|z) - \log q_x(z)] +$$

VLB.

$$D_{KL} [q_x(z) \parallel p(z|x)]$$

only work with this

$$= 0, \text{ we set this to 0.}$$

$q_z(z) \rightarrow$ optimal of VLB is $p(z|x)$ at which point, VLB is tight, $\log p(x)$.

$x \sim p_{\text{data}}$, we can now train the generative model by maximizing the VLB under data distribution.

$$\text{VLB} = \mathbb{E}_{x \sim p_{\text{data}}} \left[\mathbb{E}_{z \sim q_z(z)} \left[\log p(z) + \log p(x|z) - \log q_z(z) \right] \right]$$

$$\leq \mathbb{E}_{x \sim p_{\text{data}}} [\log p(x)] \uparrow$$

Likelihood Ratio Gradient

$$\max_{\phi} \mathbb{E}_{z \sim q_{\phi}(z)} [\underline{f(z)}]$$

$$z^{(i)} \sim q_{\phi}(z)$$

~~$$\nabla_{\phi} \frac{1}{K} \sum_{i=1}^K f(z^{(i)})$$~~

Empirical expectation

$$\max_{\phi} \sum_z q_{\phi}(z) f(z).$$

$$\Rightarrow \nabla_{\phi} \left(\sum_z q_{\phi}(z) f(z) \right)$$

↓
change dist q
the expectations
will change
for most $f(z)$.

$$\Rightarrow \sum_z \nabla_{\phi} q_{\phi}(z) f(z).$$

$$\Rightarrow \sum_z \frac{q_{\phi}(z)}{q_{\phi}(z)} \nabla_{\phi} q_{\phi}(z) f(z)$$

$$\mathbb{E}_{z \sim q_{\phi}(z)} \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)} f(z)$$

$$\Rightarrow \mathbb{E}_{z \sim q_{\phi}(z)} \left[\nabla_{\phi} \log q_{\phi}(z) f(z) \right]$$

$$\approx \mathbb{E}_{z \sim q_{\phi}(z)} \left[\nabla_{\phi} \log q_{\phi}(z) f(z) \right]$$

$$\approx \frac{1}{k} \sum_{i=1}^k \nabla_{\phi} \log q_{\phi}(z^{(i)}) f(z^{(i)})$$

↑ ↪ gradient

(Increasing the log probability of something we sampled higher.)

This leads to perfect gradient when we pick only many samples, hence, we do a reparameterization trick.

Reparameterization trick

Pathwise Derivative

$$\mathbb{E}_{z \sim q_{\phi}(z)} [f(z)]$$

$$q_{\phi}(z) = \mathcal{N}(\mu, \sigma^2)$$

↓

$$= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} [f(\mu + \epsilon \sigma)]$$

$$z = \mu + \epsilon \cdot \sigma$$

$$\epsilon \sim \mathcal{N}(0, I)$$

↑
no, ϕ here

$$= \frac{1}{K} \sum_{i=1}^K f(\mu + \epsilon^{(i)} \sigma)$$

$\nabla_{\mu, \sigma} (\longrightarrow)$

\hookrightarrow much lower
variance

$$\max_{\theta, \phi} \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_z(z) + \log p_{\theta}(x^{(i)}|z) - \log q_{\phi}(z|x^{(i)})]$$

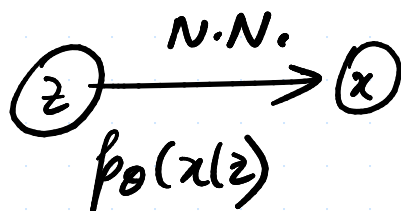
Differentiate
w.r.t. θ

$$\nabla_{\theta} = \nabla_{\theta} \mathbb{E}_{z \sim q_{\phi}} [\log p_{\theta}(x^{(i)}|z)]$$

$$= \nabla_{\theta} \sum_{k=1}^K \log p_{\theta}(x^{(i)}|z^{(k)})$$

(optimizing)

$$z^{(k)} \sim q_{\phi}(z|x^{(i)})$$



$$\nabla_{\phi} = \nabla_{\phi} q_{\phi}(z^{(k)}) \left[\log p_z(z^{(k)}) + \log p_{\theta}(x^{(i)}|z^{(k)}) - \log q_{\phi}(z^{(k)}|x^{(i)}) \right]$$

$$\Rightarrow \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[-\nabla_{\phi} \log q_{\phi}(z|x^{(i)}) \right]$$

$$\Rightarrow \mathbb{E}_{z \sim q_{\phi}} - \frac{\nabla_{\phi} q_{\phi}(z|x^{(i)})}{q_{\phi}(z|x^{(i)})}$$

$$\Rightarrow \sum_z \frac{\cancel{q_{\phi}(z|x^{(i)})} \nabla_{\phi} q_{\phi}(z|x^{(i)})}{\cancel{q_{\phi}(z|x^{(i)})}}$$

$$\Rightarrow \nabla_{\phi} \left(\sum_z q_{\phi}(z|x^{(i)}) \right) = 0 \rightarrow \text{distribution, hence the gradient contribution is } \underline{\underline{0}}.$$

Likelihood Ratio Estimator.

We are interested in :-

$$\arg \max_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(z)]$$

How do we compute,

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(z)]$$

$$\nabla_{\phi} \sum_z q_{\phi}(z|x) f(z) = \sum_z \nabla_{\phi} q_{\phi}(z|x) f(z)$$

$$= \sum_z \nabla_{\phi} \frac{q_{\phi}(z|x)}{q_{\phi}(z|x)} f(z) q_{\phi}(z|x)$$

$$= \nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [f(z)]$$

KL.

$$\Rightarrow \sum_z (\nabla_{\phi} \log q_{\phi}(z|x) f(z)) q_{\phi}(z|x)$$

if you derivative, then
this will be the (reverse).

$$\Rightarrow \sum_{\mathbf{z}} (\nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x}) f(\mathbf{z})) q_{\phi}(\mathbf{z}|\mathbf{x})$$

$$\Rightarrow \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\cdot|\mathbf{x})} [\nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x}) f(\mathbf{z})]$$

$$\phi \leftarrow \phi + \alpha \nabla_{\phi} \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\cdot|\mathbf{x})} [\nabla_{\phi} \log q_{\phi}(\mathbf{z}|\mathbf{x}) f(\mathbf{z})]}_{\text{}} \quad \uparrow$$

Pathwise derivative \Leftarrow

$\mathbf{z} \sim$ continuous \rightarrow cast \mathbf{z} as a function of a sample fixed noise, such as a standard Gaussian.

$$\mathbf{z} = g(\epsilon, \phi) \Rightarrow \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

$$\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\cdot|\mathbf{x})} [f(\mathbf{z})] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} [f(g(\epsilon, \phi))]$$

When f is differentiable :-

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(\cdot|x)} [f(z)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbb{I})} [\nabla_{\phi} f(g(\epsilon, \phi))]$$

Pathwise derivative applied to

Variational Inference \rightarrow Variational A-E.

$q_{\phi}(z|x) \rightarrow$ modelled as a Gaussian

with params μ & σ a DNN encoder (param ϕ) of x . DNN decoder $p_{\theta}(x|z)$ is differentiable.

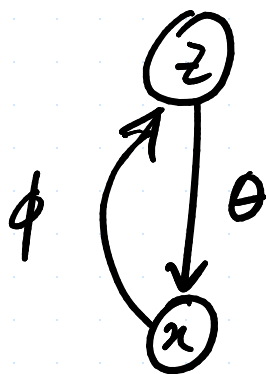
$$z = \sum^{1/2} (x; \phi) \epsilon + \mu(x; \phi).$$

$$VLB = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbb{I})} [\log p_{\theta}(x|z) - \log q_{\phi}(z|x) + \log p(z)]$$

$$= \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \Sigma)} [\log p_{\theta}(x|z)] \\ - \text{KL}(q_{\phi}(z|x) \parallel p(z))$$

$\nabla_{\theta} [\text{VLB}] \triangleq \nabla_{\phi} [\text{VLB}] \rightarrow$ effectively computed using SGD.

$$q_{\phi}(z|x) = \mathcal{N}(\mu_{\phi}(x), \sigma_{\phi}(x))$$



$$p_{\theta}(x|z) = \mathcal{N}(\mu_{\theta}(z), \sigma_{\theta}(z))$$

$$x^{(1)}, x^{(2)}, \dots, x^{(m)}$$

$$\underline{\underline{q_{\phi}(z|x^{(1)}) \rightarrow z^{(1)}}}$$

$\text{ELBO}(z^{(1)}) \rightarrow$ do a gradient descent using that, update & repeat.

Why is it called Auto-Encoder?

Gaussian Prior $\rightarrow p(z)$.

Approximate posterior $\rightarrow q_\phi(z|x)$.

reconstruction loss.

$$\log p_\theta(x) \geq \mathbb{E}_{z \sim q_\phi(z|x)} \log p_\theta(x|z)$$

$$- \text{KL}(q_\phi(z|x) \parallel p(z))$$

Auto Encoder part.

Regularization

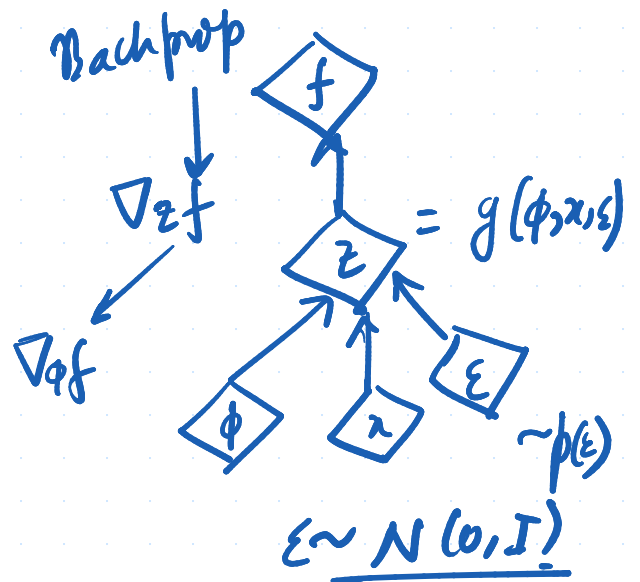
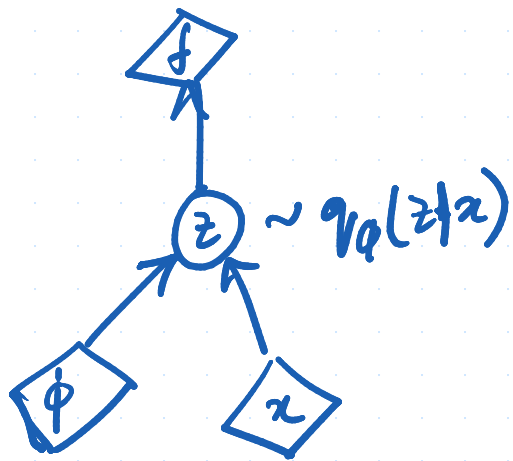
How likely $p_\theta(x|z)$ to sample the x & reconstructed, we will get the original one back.

$$x \rightarrow \square \rightarrow \square \xrightarrow{\mu} \mu \xrightarrow{\sigma} \sigma$$

Decoder $p_\theta(x|z)$

$\mathcal{L}(\theta, \phi)$ - VAE objective

My posterior of $z|x$ has to be some-what close to $p(z)$



The variational parameter ϕ affects the objective f through the R.V. $z \sim q_\phi(z|x)$. We wish to compute gradients $\nabla_\phi f$ to optimize the objective with SGD.

Original from (left) we cannot differentiate f w.r.t. ϕ , because we cannot directly backprop gradients through R.V. (z).

Externalize the randomness in z by re-parameterizing the variable as a

deterministic and differentiable function ϕ , and a newly introduced R-V. ϵ .

This allows us to backprop through z & compute $\nabla \phi$.