## VAE structure:



$x$    input

$q_\phi(z|x)$   Encoder

$\theta$

$N(0,I) \rightarrow \varepsilon$

Sampling $\rightarrow z$

$\mu + \sigma \odot \varepsilon$

↑ Reparametrization trick

: Without $\sigma \odot \varepsilon$, model is not differentiable.

Decoder $g_\theta(z) = P_\theta(x|z)$

$\hat{x}$   output.

## Visualizing Reparametrization



: Without Reparametrization
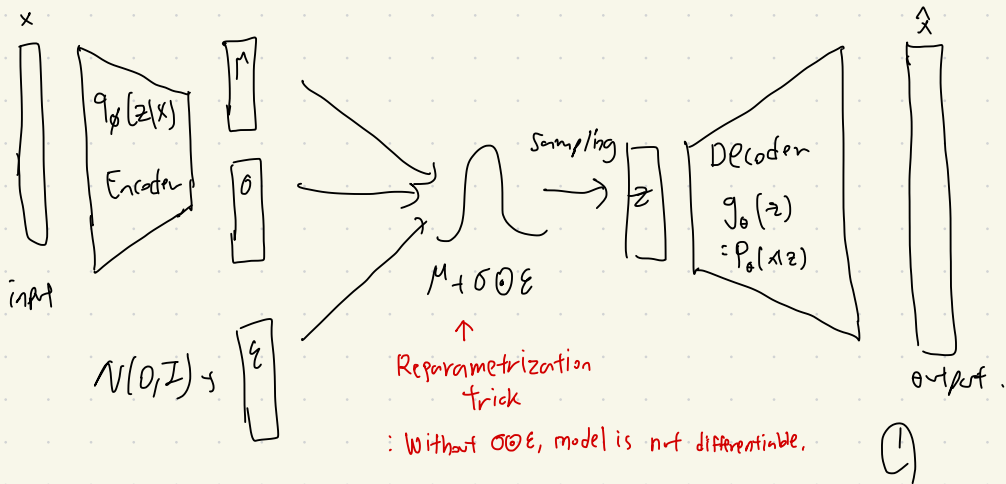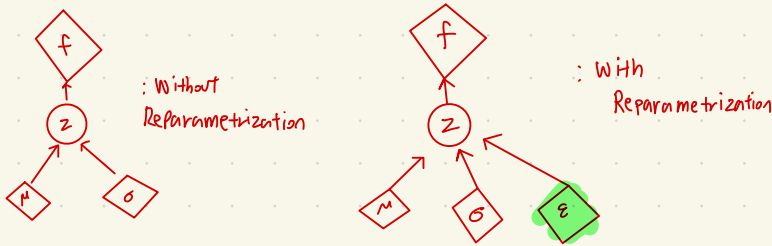
: With Reparametrization

Objective of VAE: Learn model parameters to maximize likelihood of training data: $P_\theta(x)$
And, $P_\theta(x) = \int P_\theta(z) P_\theta(x|z)\,dz$

Why? $\dfrac{P_\theta(x,z)}{P_\theta(z)} = P_\theta(x|z) \Rightarrow P_\theta(x,z) = P_\theta(z)P_\theta(x|z)$

And $\int P_\theta(x,z)\,dz = P_\theta(x)$.

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

However, $\int P_\theta(z)\,dz$ is intractable

It implies we would derive a lower bound on the data likelihood, which is tractable

Let's work on log data likelihood

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

$= \mathbb{E}_{z \sim q_\phi(z|x)}\left[\log \dfrac{P_\theta(x^{(i)}|z)\,P_\theta(z)}{P_\theta(z|x^{(i)})}\right]$

★ Use Bayes' Rule

$P(z|x) = \dfrac{P(x|z)P(z)}{P(x)} \Rightarrow P(x) = \dfrac{P(x|z)P(z)}{P(z|)}$

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \left( \frac{P_\theta(x^{(i)}|z) \, P_\theta(z)}{P_\theta(z|x^{(i)})} \cdot \boxed{\frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})}} \right) \right]$$

<span style="color:red">= Multiplying 1</span>

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \left( P_\theta(x^{(i)}|z) \cdot \frac{P_\theta(z)}{q_\phi(z|x^{(i)})} \cdot \frac{q_\phi(z|x^{(i)})}{P_\theta(z|x^{(i)})} \right) \right]$$

$$= \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \left( P_\theta(x^{(i)}|z) \right) - \log \left( \frac{q_\phi(z|x^{(i)})}{P_\theta(z)} \right) + \log \left( \frac{q_\phi(z|x^{(i)})}{P_\theta(z|x^{(i)})} \right) \right]$$

$$= \mathbb{E}_z \left[ \log P_\theta(x^{(i)}|z) \right] - \mathbb{E}_z \left[ \log \left( \frac{q_\phi(z|x^{(i)})}{P_\theta(z)} \right) \right] + \mathbb{E}_z \left[ \log \left( \frac{q_\phi(z|x^{(i)})}{P_\theta(z|x^{(i)})} \right) \right]$$

$$= \mathbb{E}_z \left[ \log P_\theta(x^{(i)}|z) - KL\left( q_\phi(z|x^{(i)}) \| P_\theta(z) \right) + KL\left( q_\phi(z|x^{(i)}) \| P_\theta(z|x^{(i)}) \right) \right.$$

<span style="color:red">$$\text{*} \quad \mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log \frac{q_\phi(z|x^{(i)})}{P_\theta(z)} \right] = \int q_\phi(z|x^{(i)}) \log \left( \frac{q_\phi(z|x^{(i)})}{P_\theta(z)} \right) dz = KL\left( P \| q \right)$$</span>

<span style="color:red">$$\text{Since } KL(P \| Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$</span>

However, we know $P_\theta(z|x^{(i)})$ is intractable. However KL always $\geq 0$.
Hence,

$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[ \log P_\theta(x^{(i)}|z) \right] - KL\left( q_\phi(z|x^{(i)}) \| P_\theta(z) \right) \quad \text{is our ELBO.}$$

$$\mathcal{L}(x^{(i)}, \theta, \phi) : \text{Tractable lower bound.}$$

<span style="color:blue">Hence, VAE Objective : $\theta^*, \phi^* = \underset{\theta, \phi}{\arg\max} \sum_{i=1}^{N} \mathcal{L}(x^{(i)}, \theta, \phi)$</span>

<span style="color:red">Reconstruction Error</span>          <span style="color:red">Regularization</span>

$$\theta^*, \phi^* = \underset{\theta, \phi}{\arg\min} \sum_{i=1}^{N} - \mathcal{L}(x^{(i)}, \theta, \phi) = \sum_{i=1}^{N} - \log P_\theta\left( x^{(i)} | g_\phi(z) \right) + KL\left( q_\phi(z|x^{(i)}) \| P(z) \right)$$

<span style="color:red">$\Downarrow$</span>        where $P(z) = \mathcal{N}(0, I)$

<span style="color:red">$$P(x | g_\phi(z)) = P_\theta(x|z)$$</span>

Regularization:
Assumption 1: $q_\phi(z|x^{(i)}) \sim \mathcal{N}(\mu^{(i)}, \sigma^{(i)} I)$
Assumption 2: $P_\theta(z) \sim \mathcal{N}(0, I)$
            Cause it makes VAE simple and reasonable.

$$D_{KL}(N_0 \| N_1) = \frac{1}{2}\left(tr\left[\Sigma_1^{-1}\Sigma_0\right] + (M_1 - M_0)^T \Sigma_1^{-1}(M_1 - M_0) - k + \ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right)\right)$$

Since our $N_1 = N(0, I)$, $\Sigma_1 = I$, $M_1 = 0$

Then, $KL(q_\phi(z|x^{(i)}) \| P_\theta(z)) = \frac{1}{2}\sum_{j=1}^{J} M_{i,j}^2 + \sigma_{i,j}^2 - \ln(\sigma_{i,j}^2) - 1)$

Reconstruction Error : MSE or Cross Entropy

MSE : When Decoder uses Gaussian Distribution

Cross-Entropy: When Decoder uses Bernoulli Distribution

Therefore, Goal of VAE is

$$\theta^*, \phi^* = \underset{\theta, \phi}{argmin} \sum_i -\mathbb{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log(P_\theta(x^{(i)}|z))\right] + KL(q_\phi(z|x^{(i)}) \| P(z))$$

<span style="color:red">Reconstruction Error<br>(MSE or Cross-Entropy)</span>

<span style="color:blue">Regularization<br>: Encouraging a structured and well-behaved latent space</span>

β-VAE (Disentangled VAE)

Recall VAE Loss: $L_{VAE} = -\log P_\theta(x|z) + D_{KL}(q_\phi(z|x) \| P_\theta(z))$

<span style="color:red">improving decoder<br>$z \to x$</span>   <span style="color:blue">improving encoder<br>(better representation of z)</span>

Another way of writing VAE objective:

$$\underset{\phi, \theta}{max} \; \mathbb{E}_{x \sim D}\left[\mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x|z)\right]$$
$$s.t \; D_{KL}(q_\phi(z|x) \| P_\theta(z)) < \delta$$

: Maximizing probability of generating real data, while keeping distance between real and approximate posterior distribution $(q_\phi(z|x))$ small. (under small constant $\delta$)

• VAE maximization objective can then be rewritten as a Lagrangian with a Lagrangian multiplier β under KKT condition

$$: \mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x|z) - \beta(D_{KL}(q_\phi(z|x) \| P_\theta(z)) - \delta)$$
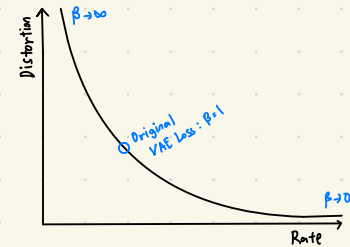$$= \mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x|z) - \beta(D_{KL}(q_\phi(z|x) \| P_\theta(z)) + \beta\delta$$

β-VAE Loss hence given by:

$$\mathcal{L}_{Beta}(\phi, \beta) = -\mathbb{E}_{z \sim q_\phi(z|x)} \log P_\theta(x|z) + \beta D_{KL}(q_\phi(z|x) \| P_\theta(z))$$

- When β=1 ⇒ Standard VAE
- When β>1 ⇒ Stronger constraint on latent bottleneck, follow generative process and thus encourage disentanglement
  ↳ Increasing representation capacity of z. Generally, "stronger disentanglement" means more Gaussian like for most linear definition.
- Disadvantage: Trade-off between disentanglement and reconstruction capability.

★ Distortion: $d(x, \hat{x})$ where $\hat{x}$ reconstruction/estimate of x given z. It's a reconstruction loss. e.g. $-\log(P(x|z)) = MSE[x|z]$
★ Rate: Regularization: $KL(q_\phi(z|x) \| P(z))$



As β → 0, Optimizer Prioritizes minimizing distortion: Rate ↑, Distortion ↓
β → ∞, Optimizer prioritizes minimizing rate over maximizing distortion
Rate ↓, Distortion ↑