

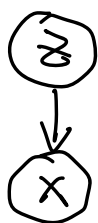
Why VAE:

traditional autoencoders encode the input into a single vector, which is a single point in the high-dim space. when generating new, unseen data, they fail to interpolate between seen data, so cannot generate those variations. In contrast,

VAE maps inputs into a distribution (Gaussian), which forces in the circle area around  $\mu$  to have data which makes data embedding in the high-dim space denser and easier to interpolate.

(VAE is like Gaussian Mixture Model, traditional autoencoder is like K-means)

## Derivation



What we want:

generate data  $x'$  which are similar to the training data  $X$

$\Rightarrow$  we need to model  $p(z|x)$  which we can sample  $z$  from

and generate  $x$ :  $p(x'|x) = \int p(z|x) p(x|z) dz$

$\Rightarrow$  how to model  $p(z|x)$

$$p(z|x) = \frac{p(x|z) p(z)}{p(x)}$$

$\Rightarrow$  problem:  $p(x) = \int p(x|z) p(z) dz$

is intractable due to continuous  $z$

$\therefore$  cannot iterate all  $z$

$\therefore$  cannot compute integral using computer.

⇒ solution: use another dist  
Easier  $q(z|x)$  to replace  $p(z|x)$  in  
distribution  
e.g. diagonal Gaussian.

⇒ Back to our goal: we  
need to generate  $x'$  from  $p(z|x)$   
which are similar to  $x$ . So  
we need  $q(z|x)$  is as close  
to  $p(z|x)$  as possible.

$$\begin{aligned} &\Rightarrow KL(q(z|x) || p(z|x)) \\ &= \int_z q(z|x) \log \frac{q(z|x)}{p(z|x)} dz \\ &= \int_z q(z|x) \log q(z|x) dz \\ &\quad - \int_z q(z|x) \log p(z|x) dz \\ &= \int_z q(z|x) \log \frac{p(x, z)}{p(x)} dz \\ &= \int_z q(z|x) \log p(x, z) dz \end{aligned}$$

$$- \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}) d\mathbf{z}$$

$$= \log p(\mathbf{x})$$

$$\therefore KL(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

$$= \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log q(\mathbf{z}|\mathbf{x}) d\mathbf{z}$$

$$- \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

$$+ \log p(\mathbf{x})$$

$\Rightarrow$  problem:  $p(\mathbf{x})$  cannot be computed !!  $\int p(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$ .

$$\Rightarrow \log p(\mathbf{x}) - KL(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$

what  $= \int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$

we want to maximize  $-\int_{\mathbf{z}} q(\mathbf{z}|\mathbf{x}) \log q(\mathbf{z}|\mathbf{x}) d\mathbf{z}$

$$\textcircled{1} \log p(\mathbf{x}) \uparrow \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z}|\mathbf{x})]$$

$\Rightarrow$  similar to neural data ELBO

$$\textcircled{2} \quad q(z|x) \Rightarrow p(z|x)$$

$\Rightarrow$  Equivalently, we need to maximize ELBO

Understand ELBO:

$$\text{ELBO} = \mathbb{E}_{q(z|x)} [\log p(x, z)] - \mathbb{E}_{q(z|x)} [\log q(z|x)]$$

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

$$- \mathbb{E}_{q(z|x)} [\log q(z|x)]$$

$$= \boxed{\mathbb{E}_{q(z|x)} [\log p(x|z)]}$$

encourage  $q(z|x)$  to put as much weights as possible on the mode of  $\log p(x|z)$  (or  $p(x|z)$ )

$$= \boxed{\text{KL}(q(z|x) \| \underbrace{p(z)}_{\mathcal{N}(0, I)})}$$

encourage  $q(z|x)$  to stretch a little bit to approximate  $p(z)$

If  $q(z|x)$  is so narrow, we  
can add more weights to  
 $KL(q(z|x) || p(z))$

In VAE:

Encoder Network to represent  
 $q_{\phi}(z|x)$

Generate (Decoder) Network to  
represent  $p_{\theta}(x|z)$

Make Everything differentiable:

Reparameterization Trick.

sample  $z \sim \mathcal{N}(\mu, \sigma^2)$

$\Leftarrow z = \mu + \sigma \odot \epsilon \quad \epsilon \in \mathcal{N}(0, I)$