

# Variational Autoencoders

by Paul Hand  
Northeastern University

## Outline :

Generative Models and Autoencoders

Variational Lower Bound (VLB)

Optimizing VLB + Variational Autoencoders

Resource: Kingma and Welling 2019, Chapter 2,  
"Introduction to Variational Autoencoders"

## Generative Models

A model that can sample from a  
learned distribution

6517814828  
9683960319  
3371368179  
8908691963  
8233331386  
6998616666  
9526651899  
9977872823  
0461232088  
9734934851

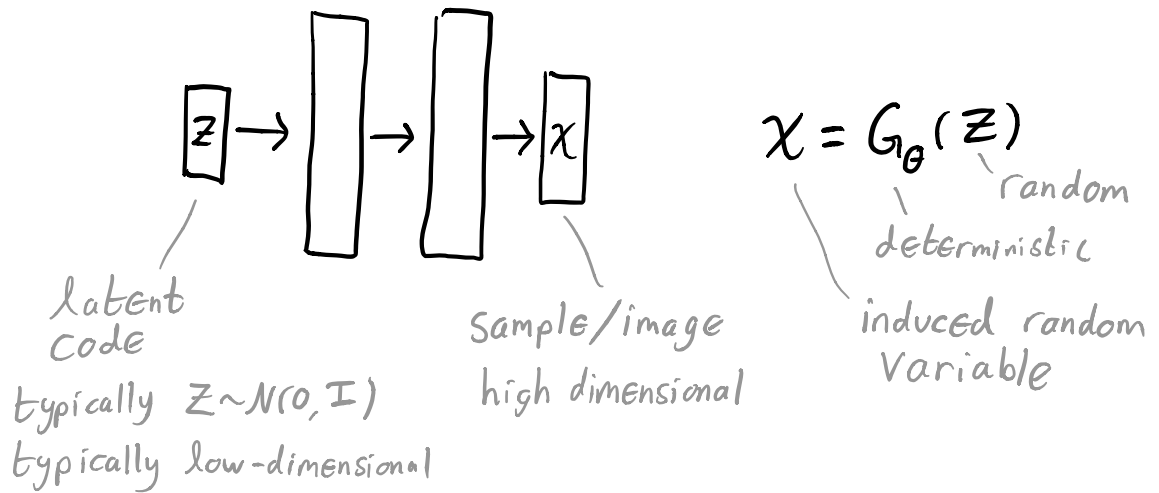
(a) 2-D latent space

(Kingma and Welling, 2014)



(Razavi et al. 2019)

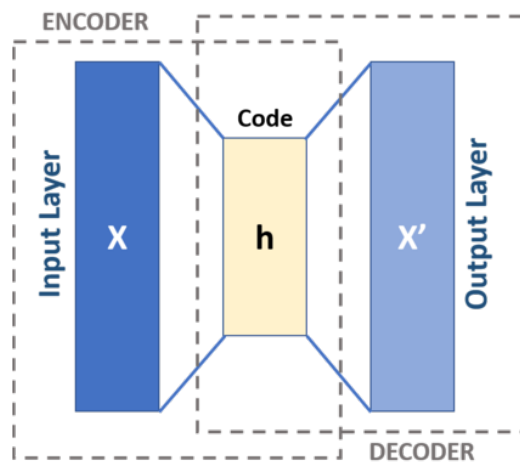
In some generative models, samples are generated by a net applied to a random latent code



## Autoencoders

Autoencoders attempt to reconstruct input signals/images by learning mappings to and from a code

Want:  $x' = D_{\theta}(E_{\phi}(x)) \approx x$



$$\min_{\theta, \phi} \sum_{i=1}^n \|D_{\theta}(E_{\phi}(x_i)) - x_i\|^2 \quad \text{w/ } \{x_i\}_{i=1 \dots n} \text{ is dataset}$$

A (plain) autoencoder is not a generative model as it does not define a distribution

### Training a low latent-dimensional generative model by likelihood

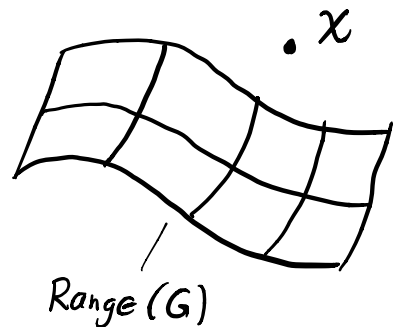
Given data  $\{x_i\}_{i=1 \dots n}$ , train a gen. model to maximize the likelihood of the observed data

If gen. model

$$G_{\theta} : \mathbb{R}^k \rightarrow \mathbb{R}^d \quad \text{w/ } k < d, \\ z \mapsto x$$

then  $p(x) = 0$  almost everywhere

So, can't directly optimize likelihood



To have nonzero likelihood everywhere,  
define noisy observation model

$$P_{\theta}(x|z) = \mathcal{N}(x; G_{\theta}(z), \eta I)$$

Under a simple prior  $p(z)$ , this induces

a joint distribution  $p_{\theta}(x, z)$

$$\text{Now } p(x) = \int p(z) p(x|z) dz$$

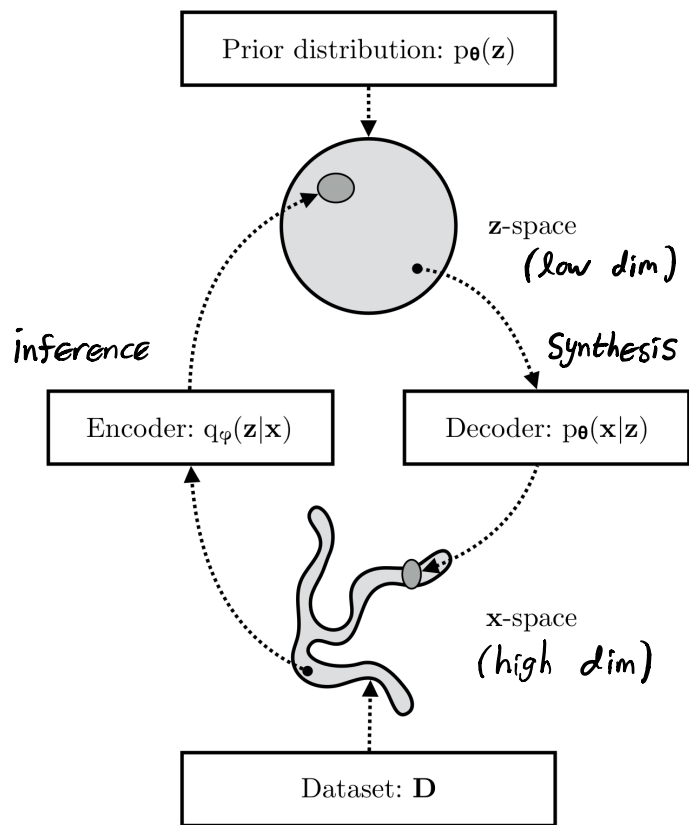
/  
intractable to evaluate at each iteration  
optimize a lower bound instead

## Variational Lower Bound

Setup:

Data generated by  
 $z \sim p(z)$  prior  
 $x \sim p_{\theta}(x|z)$

Use  $q_{\psi}(z|x)$  as  
proxy for  $p_{\theta}(z|x)$



(Kingma and Welling 2019)

Find a lower bound to  $P_\theta(x)$

$$\begin{aligned}\log P_\theta(x) &= \mathbb{E}_{z \sim q_\varphi(z|x)} \log P_\theta(x) = \mathbb{E}_{z \sim q_\varphi(z|x)} \log \frac{P_\theta(x, z)}{P_\theta(z|x)} \\ &= \mathbb{E}_{z \sim q_\varphi(z|x)} \log \left( \frac{P_\theta(x, z)}{q_\varphi(z|x)} \cdot \frac{q_\varphi(z|x)}{P_\theta(z|x)} \right) \\ &= \underbrace{\mathbb{E}_{z \sim q_\varphi(z|x)} \log \frac{P_\theta(x, z)}{q_\varphi(z|x)}}_{\mathcal{L}_{\theta, \varphi}(x)} + \underbrace{\mathbb{E}_{z \sim q_\varphi(z|x)} \log \frac{q_\varphi(z|x)}{P_\theta(z|x)}}_{D_{KL}(q_\varphi(z|x) \parallel P_\theta(z|x))}\end{aligned}$$

Variational Lower Bound (VLB)

Evidence lower bound (ELBO)

Define:  $D_{KL}(q \parallel p) = \mathbb{E}_{z \sim q} \log \frac{q(z)}{p(z)}$

Note: •  $D_{KL}(q \parallel p) \neq D_{KL}(p \parallel q)$

• Measure of how far  $p$  is from  $q$

•  $D_{KL}(q \parallel p) \geq 0$  (and is 0 iff  $p=q$ )

So,  $\log P_\theta(x) \geq \mathbb{E}_{z \sim q_\varphi(z|x)} \log \frac{P_\theta(x, z)}{q_\varphi(z|x)} = \mathcal{L}_{\theta, \varphi}(x)$

## Interpretation of VLB

$$\begin{aligned} \mathcal{L}_{\theta, \psi}(\mathcal{X}) &= \mathbb{E}_{z \sim q_{\psi}(z|\mathcal{X})} \log \frac{p_{\theta}(\mathcal{X}, z)}{q_{\psi}(z|\mathcal{X})} = \mathbb{E}_{z \sim q_{\psi}(z|\mathcal{X})} \log p_{\theta}(\mathcal{X}|z) \frac{p(z)}{q_{\psi}(z|\mathcal{X})} \\ &= \underbrace{\mathbb{E}_{z \sim q_{\psi}(z|\mathcal{X})} \log p_{\theta}(\mathcal{X}|z)}_{\text{reconstruction error}} + \underbrace{\mathbb{E}_{z \sim q_{\psi}(z|\mathcal{X})} \log \frac{p(z)}{q_{\psi}(z|\mathcal{X})}}_{\text{regularization}} \end{aligned}$$

First term:  $p_{\theta}(\mathcal{X}|z) = \mathcal{N}(\mathcal{X}; G_{\theta}(z), \gamma I)$   
 $\Rightarrow \log p_{\theta}(\mathcal{X}|z) = -\frac{1}{2\gamma} \|\mathcal{X} - G_{\theta}(z)\|^2 + \text{constant}$

So  $\mathbb{E}_{z \sim q_{\psi}} \log p_{\theta}(\mathcal{X}|z)$  is expected  $\ell_2$  reconstruction error under the encoder model

Maximizing VLB encourages  $q_{\psi}$  to be point mass

Second term:  $\mathbb{E}_{z \sim q_{\psi}(z|\mathcal{X})} \log \frac{p(z)}{q_{\psi}(z|\mathcal{X})} = -D_{\text{KL}}(q_{\psi}(z|\mathcal{X}) \| p(z))$

So maximizing VLB  $\mathcal{L}_{\theta, \psi}$  pushes  $q_{\psi}(z|\mathcal{X})$  toward  $p(z)$ . Prevents  $q_{\psi}$  from being a point mass.

Makes  $q_{\psi}(z|\mathcal{X})$  more like standard normal

Maximizing VLB  $\mathcal{L}_{\theta, \psi}$  :

- Roughly maximizes  $P(x)$
- Minimizes KL divergence of  $q_{\psi}(z|x)$  and  $p_{\theta}(z|x)$ , making  $q_{\psi}$  better

Main

Idea: Instead of optimizing  $\sum_{i=1}^n \log p_{\theta}(x_i)$ ,  
optimize  $\sum_{i=1}^n \mathcal{L}_{\theta, \psi}(x_i)$

$$\text{w/ } \mathcal{L}_{\theta, \psi}(x) = \mathbb{E}_{z \sim q_{\psi}(z|x)} \log \frac{p_{\theta}(x, z)}{q_{\psi}(z|x)}$$

Optimizing Variational Lower Bound

$$\max_{\theta, \psi} \sum_{i=1}^n \mathcal{L}_{\theta, \psi}(x_i)$$

One possibility :

for each  $x_i$ , find best  $q_{\psi}(z|x_i)$  by multiple gradient steps in  $\psi$ . Then gradient ascend in  $\theta$ .

Expensive inference updates

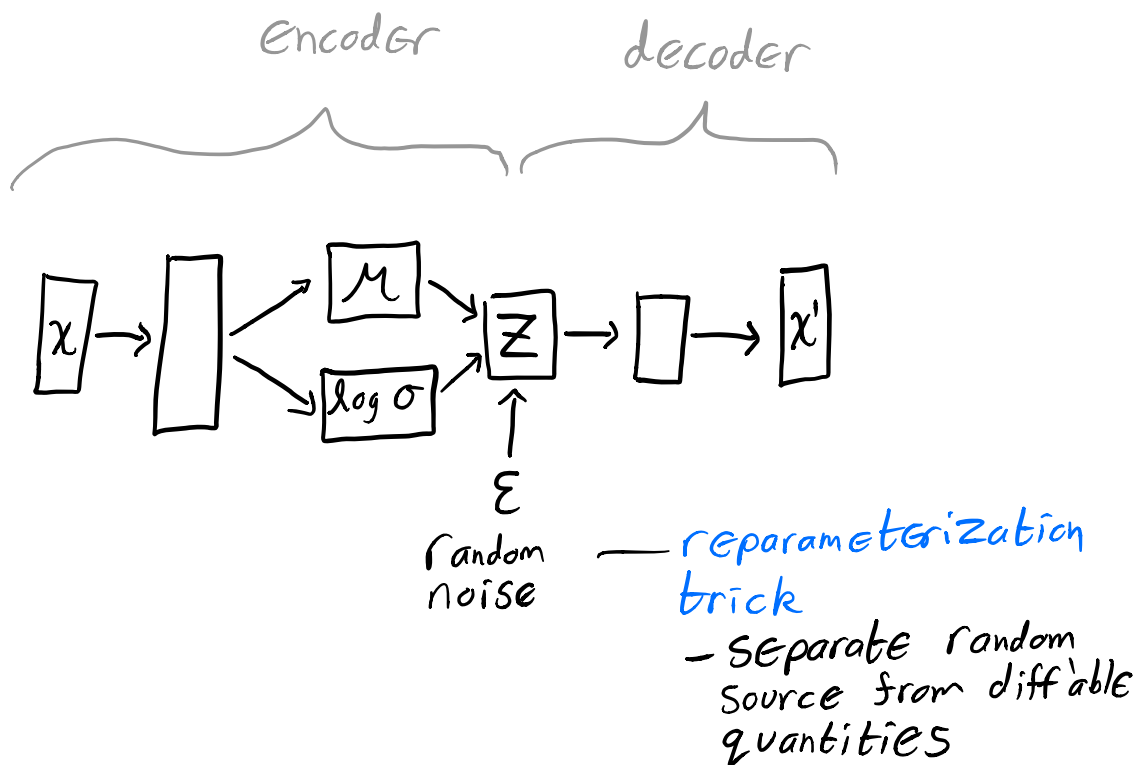
Instead<sup>o</sup>

amortize the inference costs  
by learning an inference net

$$x \mapsto (\mu, \Sigma) \text{ w/ } q_{\psi}(z|x) = \mathcal{N}(z; \underbrace{\mu(x)}_1, \underbrace{\Sigma(x)}_1) \\ \text{or } \sigma(x)\mathbf{I}$$

Parameters of inference model  
are shared between data points

## Variational Autoencoder architecture





# Stochastic Gradient Optimization of VLB

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

Solve  $\max_{\theta, \varphi} \sum_{x_i \in \mathcal{D}} \mathcal{L}_{\theta, \varphi}(x_i)$

$$\text{w/ } \mathcal{L}_{\theta, \varphi}(x) = \mathbb{E}_{z \sim q_{\varphi}(z|x)} \log \frac{p_{\theta}(x, z)}{q_{\varphi}(z|x)}$$

Computing  $\nabla_{\theta, \varphi} \mathcal{L}_{\theta, \varphi}(x_i)$  is intractable,  
but there are unbiased estimators

Easy to get unbiased  $\nabla_{\theta} \mathcal{L}_{\theta, \varphi}$ :

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\theta, \varphi}(x) &= \nabla_{\theta} \mathbb{E}_{z \sim q_{\varphi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\varphi}(z|x)] \\ &= \mathbb{E}_{z \sim q_{\varphi}} \nabla_{\theta} \log p_{\theta}(x, z) \\ &\approx \nabla_{\theta} \log p_{\theta}(x, z) \text{ w/ } z \sim q_{\varphi}(z|x) \\ &\quad \text{unbiased estimate} \end{aligned}$$

Not as easy to get unbiased  $\nabla_{\psi} \mathcal{L}_{\theta, \psi}^{\circ}$

$$\begin{aligned}\nabla_{\psi} \mathcal{L}_{\theta, \psi}(\chi) &= \nabla_{\psi} \mathbb{E}_{z \sim q_{\psi}(z|\chi)} [\log P_{\theta}(\chi, z) - \log q_{\psi}(z|\chi)] \\ &\neq \mathbb{E}_{z \sim q_{\psi}} \nabla_{\psi} (\log P_{\theta}(\chi, z) - \log q_{\psi}(z|\chi))\end{aligned}$$

Recall,  $q_{\psi}(z|\chi) = \mathcal{N}(z; \mu(\chi), \sigma(\chi)\mathbf{I})$   
 $= \mu(\chi) + \sigma(\chi) \cdot \varepsilon$ , w/  $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$

$$\begin{aligned}\mathcal{L}_{\theta, \psi}(\chi) &= \mathbb{E}_{z \sim q_{\psi}(z|\chi)} (\log P_{\theta}(\chi, z) - \log q_{\psi}(z|\chi)) \\ &= \mathbb{E}_{\varepsilon \sim P(\varepsilon)} (\log P_{\theta}(\chi, z) - \log q_{\psi}(z|\chi))\end{aligned}$$

Form estimator of  $\mathcal{L}_{\theta, \psi}(\chi)$  as  $\tilde{\mathcal{L}}_{\theta, \psi}(\chi)$  by:

$$\begin{aligned}\varepsilon &\sim P(\varepsilon) \\ z &= \mu_{\psi}(\chi) + \sigma_{\psi}(\chi) \varepsilon = g(\psi, \chi, \varepsilon) \\ \tilde{\mathcal{L}}_{\theta, \psi}(\chi) &= \log P_{\theta}(\chi, z) - \log q_{\psi}(z|\chi)\end{aligned}$$

Unbiased estimate of  $\nabla_{\psi} \mathcal{L}_{\theta, \psi}(\mathcal{X})$

$$\nabla_{\psi} \hat{\mathcal{L}}_{\theta, \psi}(\mathcal{X})$$

Note:  $\mathbb{E}_{\mathcal{E} \sim p(\mathcal{E})} \hat{\mathcal{L}}_{\theta, \psi}(\mathcal{X}) = \mathcal{L}_{\theta, \psi}(\mathcal{X})$

So  $\mathbb{E}_{\mathcal{E} \sim p(\mathcal{E})} \nabla_{\psi} \hat{\mathcal{L}}_{\theta, \psi}(\mathcal{X}) = \nabla_{\psi} \mathbb{E}_{\mathcal{E} \sim p(\mathcal{E})} \hat{\mathcal{L}}_{\theta, \psi} = \nabla_{\psi} \mathcal{L}_{\theta, \psi}$

Optimize VAE parameters w/ stochastic gradients.

Can extend models to be more sophisticated,  
eg  $\Sigma(\mathcal{X})$  vs  $\sigma \mathbb{I}(\mathcal{X})$  as inference model.

Key points:

- optimize a lower bound to likelihood
- lower bound has terms for reconstruction and for regularization
- maintain an inference model for  $Z|\mathcal{X}$  in place of intractable true distribution

- reparameterization trick allows backpropagating on mean and variance of inference model
- VAEs have been trained with photorealistic outputs