

# CSC420: Tutorial 4

## VAE and Diffusion Model

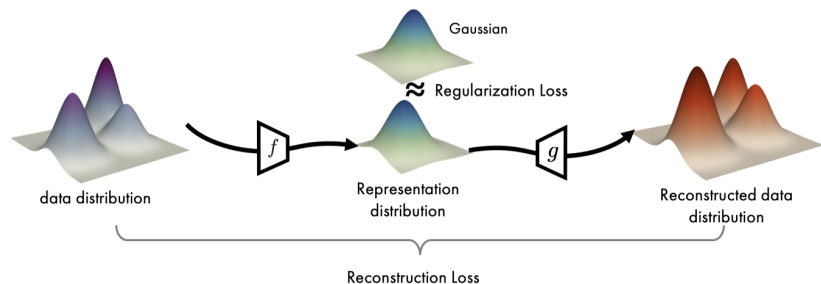
Michael Neumayr

September 28, 2025

- ▶ P1: VAE and Variational Inference Background
- ▶ P2: DDPM From Scratch

# vae overview

VAE conceptually:



train with ELBO objective:

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \underbrace{-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction term}} + \underbrace{D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))}_{\text{regularization term}}$$

# how did we get here

## goals:

- ▶ **encoder:** we want to map high-dimensional data  $\mathbf{x}$  to a lower-dimensional latent space  $\mathbf{z}$
- ▶ **decoder:** we want to be able to sample the latent space  $\mathbf{z}$  to generate new data  $\mathbf{x}$

use **latent variable model** (in our case, a VAE):

- ▶ assume that a few lower-dimensional latent factors  $\mathbf{z}$  can explain our complex data  $\mathbf{x}$
- ▶ exploit low-dim. latent structure to facilitate modeling and sampling of distribution  $p_{\theta}(\mathbf{x})$  over high-dim. data  $\mathbf{x}$

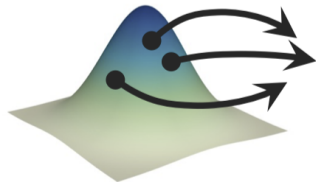
# how did we get here

**main idea:** true distribution  $p(\mathbf{x})$  is “complex” (e.g. images), but the conditional distribution  $p(\mathbf{x}|\mathbf{z})$  is “simple” (e.g. Gaussian)

**generate data** in two steps:

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}) \quad (\text{sample latent space } \mathbf{z})$$

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



since  $p_{\theta}(\mathbf{z})$  standard normal  
(no parameters  $\theta$ ), use  $p(\mathbf{z})$

# how did we get here

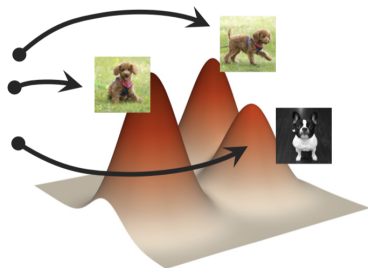
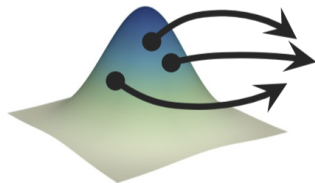
**main idea:** true distribution  $p(\mathbf{x})$  is “complex” (e.g. images), but the conditional distribution  $p(\mathbf{x}|\mathbf{z})$  is “simple” (e.g. Gaussian)

**generate data** in two steps:

$\mathbf{z} \sim p(\mathbf{z})$  (sample latent space  $\mathbf{z}$ )

$\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z})$  (generate data conditional on  $\mathbf{z}$ )

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



# how did we get here

**generate data** in two steps:

$$\mathbf{z} \sim p(\mathbf{z}) \quad (\text{sample latent variable } \mathbf{z})$$

$$\mathbf{x} \sim p_{\theta}(\mathbf{x}|\mathbf{z}) \quad (\text{generate data conditional on } \mathbf{z})$$

together, we get the joint distribution

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})$$

and model the full distribution  $p_{\theta}(\mathbf{x})$  by **marginalizing over  $\mathbf{z}$** :

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} = \int p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z}) \, d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[p_{\theta}(\mathbf{x} | \mathbf{z})] \quad (1)$$

# our tasks in this framework

**encoder:** we want to map high-dimensional data  $\mathbf{x}$  to a lower-dimensional latent space  $\mathbf{z}$

- ▶ **inference:** with sample  $\mathbf{x}$ , find posterior distribution over  $\mathbf{z}$

$$p_{\theta}(\mathbf{z} \mid \mathbf{x}) = \frac{p_{\theta}(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})}{p_{\theta}(\mathbf{x})}$$

- ▶ but we need to model the distribution  $p_{\theta}(\mathbf{x})$  for our data first



# our tasks in this framework

**learning:** given empirical dataset  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  (assume i.i.d.), find parameters  $\theta$  that best explain data

- ▶ typically, we **maximize the log-likelihood** over the dataset

$$\max_{\theta} \log p_{\theta}(\mathbf{X}) = \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

- ▶ for simplicity, look at a single data point  $\mathbf{x}$ , using 1 we get

$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \max_{\theta} \log \int p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (2)$$

- ▶ what is the problem?

# our tasks in this framework

**learning:** given empirical dataset  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  (assume i.i.d.), find parameters  $\theta$  that best explain data

- ▶ typically, we **maximize the log-likelihood** over the dataset

$$\max_{\theta} \log p_{\theta}(\mathbf{X}) = \max_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

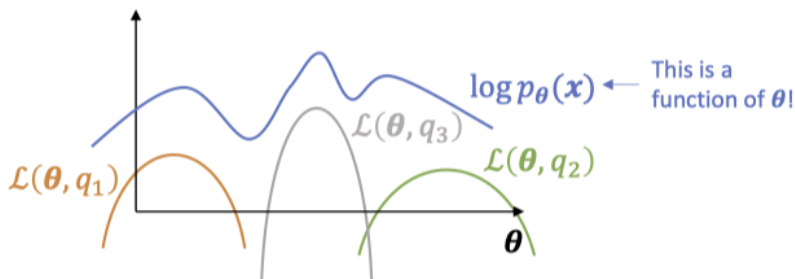
- ▶ for simplicity, look at a single data point  $\mathbf{x}$ , using 1 we get

$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \max_{\theta} \log \underbrace{\int p_{\theta}(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}}_{f(\theta)} \quad (2)$$

- ▶  $f(\theta)$  is **intractable** to evaluate: no closed-form solution and numerical integration is infeasible ( $\mathbf{z}$  still high-dimensional)

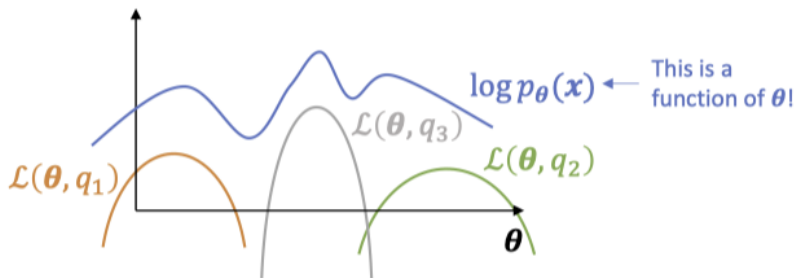
# log-likelihood maximization

- ▶  $f(\theta)$  is **intractable** to evaluate: no closed-form solution and numerical integration is infeasible ( $\mathbf{z}$  still high-dimensional)
- ▶ instead, find a **lower bound** on the log-likelihood using easy-to-evaluate functions (for us multivariate Gaussians)



# log-likelihood maximization

- ▶  $f(\theta)$  is **intractable** to evaluate: no closed-form solution and numerical integration is infeasible ( $\mathbf{z}$  still high-dimensional)
- ▶ instead, find a lower bound on the log-likelihood using easy-to-evaluate functions (for us multivariate Gaussians)



→

variational

optimize over functions

inference

infer latent features  $\mathbf{z}$  from data  $\mathbf{x}$

# ELBO derivation

- ▶ still looking for objective to tractably optimize  $\theta$  by substituting  $p_{\theta}(\mathbf{x})$  maximization with lower bound in optim
- ▶ let  $q(\mathbf{z})$  be an arbitrary distribution over  $\mathbf{z}$  (choose later)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log p_{\theta}(\mathbf{x})] \\&= \int q(\mathbf{z}) \log p_{\theta}(\mathbf{x}) \, d\mathbf{z} && | \text{ def. of expectation} \\&= \int q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} \, d\mathbf{z} && | \text{ def. of cond. probability} \\&= \int q(\mathbf{z}) \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} \cdot \frac{q(\mathbf{z})}{q(\mathbf{z})} \right) \, d\mathbf{z} && | \text{ Id. trick} \\&= \int q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \, d\mathbf{z} + \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} \, d\mathbf{z} \\&= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] + D_{\text{KL}}(q(\mathbf{z}) \| p_{\theta}(\mathbf{z} | \mathbf{x})) && | \text{ defs}\end{aligned}$$

# ELBO derivation

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log p_{\theta}(\mathbf{x})] \\&= \int q(\mathbf{z}) \log p_{\theta}(\mathbf{x}) d\mathbf{z} && | \text{ def. of expectation} \\&= \int q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} d\mathbf{z} && | \text{ def. of cond. probability} \\&= \int q(\mathbf{z}) \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} \cdot \frac{q(\mathbf{z})}{q(\mathbf{z})} \right) d\mathbf{z} && | \text{ Id. trick} \\&= \int q(\mathbf{z}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} d\mathbf{z} + \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x})} d\mathbf{z} \\&= \underbrace{\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]}_{=: \text{ELBO}(\mathbf{x})} + \underbrace{D_{\text{KL}}(q(\mathbf{z}) \| p_{\theta}(\mathbf{z} | \mathbf{x}))}_{\geq 0 \text{ (gap)}} && | \text{ defs}\end{aligned}$$

► for a fixed  $\theta$ , we want  $q_{\phi}(\mathbf{z})$  to be close to  $p_{\theta}(\mathbf{z} | \mathbf{x})$

## exercises around ELBO to deepen your understanding:

1. warm up: write out ELBO from

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[ \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$$

to

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z}))$$

2. prove  $\log p_{\theta}(\mathbf{x}) \geq \text{ELBO}(\mathbf{x})$  (it lower bounds log-likelihood) using Jensen's inequality  $\rightarrow$  Q1 (a) in assignment 2
3. prove  $\log p_{\theta}(\mathbf{x}) - \text{ELBO}(\mathbf{x}) = D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x}))$  (it is gap between log-likelihood and KL divergence between approximate and true posterior)  $\rightarrow$  Q1 (b) in assignment 2

# optimizing with elbo

optimizing  $\text{ELBO}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z}))$

- ▶ given  $\mathbf{x}$ , we want to optimize the  $\text{ELBO}(\mathbf{x})$  wrt  $\theta$  and  $q(\mathbf{z})$
- ▶ but what are we optimizing over in the case of  $q(\mathbf{z})$ ?
- ▶ choose set of tractable, parametric distributions  $\mathcal{Q}$
- ▶ every distribution  $q_{\phi}(\mathbf{z})$  is specified by its parameters  $\phi$
- ▶ best distribution  $q_{\phi}(\mathbf{z}) \iff$  best parameters  $\phi$
- ▶ instead of optimizing  $\phi_{\text{optimal}}^{(i)}$  for every data point  $\mathbf{x}^{(i)}$ , we learn a neural network that maps **every**  $\mathbf{x}_i$  to  $\phi_{\text{optimal}}^{(i)}$   
→ **amortized inference**
- ▶ as  $\mathbf{x}$  is network input, write  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  instead of  $q_{\phi}(\mathbf{z})$ :

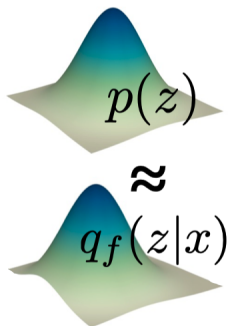
$$\text{ELBO}(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p(\mathbf{z}))$$



## default choices (keep simple and standard)

**prior:**  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  (no learnable params)

- ▶  $D_{\text{KL}}(q_\phi(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))$  term regularizes  $q_\phi(\mathbf{z} | \mathbf{x})$  to learn a distribution close to the prior  $p(\mathbf{z})$

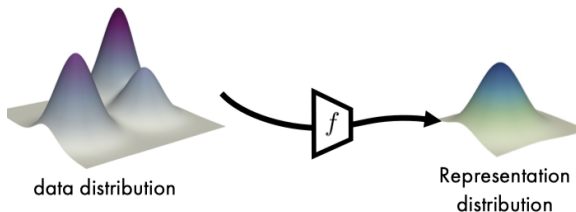


# default choices (keep simple and standard)

**encoder** (approximate posterior):

$$q_{\phi}(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x})))$$

- ▶ neural network outputs *parameters*  $\boldsymbol{\mu}_{\phi}(\mathbf{x}), \boldsymbol{\sigma}_{\phi}^2(\mathbf{x})$
- ▶ given a k-dimensional latent vector  $\mathbf{z}$ , what dimension is the output tensor of the encoder?

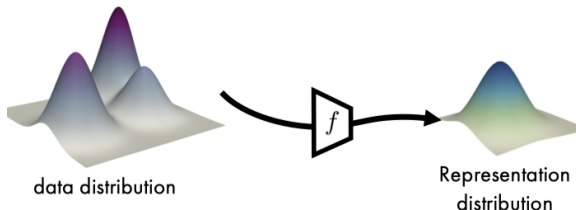


# default choices (keep simple and standard)

**encoder** (approximate posterior):

$$q_{\phi}(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x})))$$

- ▶ neural network outputs *parameters*  $\boldsymbol{\mu}_{\phi}(\mathbf{x}), \boldsymbol{\sigma}_{\phi}^2(\mathbf{x})$
- ▶ given a k-dimensional latent vector  $\mathbf{z}$ , what dimension is the output tensor of the encoder?
- ▶ for every image, the encoder outputs a latent vector with dimension (k, 2): per-latent dimension mean and variance



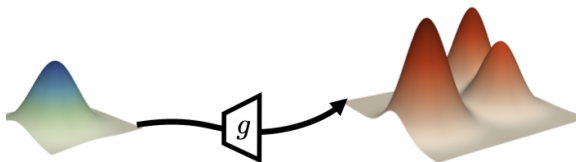
# default choices (keep simple and standard)

**decoder** (likelihood):

- ▶ pick by data type, for example:

$$p_{\theta}(\mathbf{x} \mid \mathbf{z}) = \begin{cases} \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{z}), \mathbf{I}), & \mathbf{x} \in \mathbb{R}^D \\ \text{Bernoulli}(\boldsymbol{\pi}_{\theta}(\mathbf{z})), & \mathbf{x} \in \{0, 1\}^D \end{cases}$$

- ▶ neural network outputs *parameters*  $\boldsymbol{\mu}_{\theta}(\mathbf{z})$  or  $\boldsymbol{\pi}_{\theta}(\mathbf{z})$



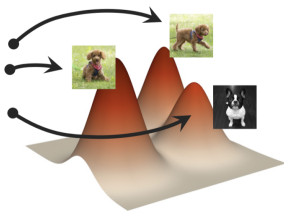
# default choices (keep simple and standard)

**decoder** (likelihood):

- ▶ pick by data type, for example:

$$p_{\theta}(\mathbf{x} \mid \mathbf{z}) = \begin{cases} \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{z}), \mathbf{I}), & \mathbf{x} \in \mathbb{R}^D \\ \text{Bernoulli}(\boldsymbol{\pi}_{\theta}(\mathbf{z})), & \mathbf{x} \in \{0, 1\}^D \end{cases}$$

- ▶ neural network outputs *parameters*  $\boldsymbol{\mu}_{\theta}(\mathbf{z})$  or  $\boldsymbol{\pi}_{\theta}(\mathbf{z})$
- ▶ retrieve samples from the distribution using the parameters



## additional slides

# bayes terminology recap

$$\underbrace{p_{\theta}(\mathbf{z} \mid \mathbf{x})}_{\text{posterior}} = \frac{\overbrace{p_{\theta}(\mathbf{x} \mid \mathbf{z})}^{\text{likelihood}} \overbrace{p(\mathbf{z})}^{\text{prior}}}{\underbrace{p_{\theta}(\mathbf{x})}_{\text{(evidence)}}}$$

# Kullback-Leibler divergence recap

- ▶ KL divergence from  $q(\mathbf{z})$  to  $p(\mathbf{z})$  is defined as

$$D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z})) := \int q(\mathbf{z}) \log \frac{q(\mathbf{z})}{p(\mathbf{z})} d\mathbf{z}$$

- ▶ properties:
  - ▶ asymmetric,  $D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z})) \neq D_{\text{KL}}(p(\mathbf{z}) \parallel q(\mathbf{z}))$  in general
  - ▶ nonnegative,  $D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z})) \geq 0$
  - ▶  $D_{\text{KL}}(q(\mathbf{z}) \parallel p(\mathbf{z})) = 0 \iff p = q$  almost everywhere



# high level understanding

