

# CS-E4891 Deep Generative Models

## Lecture 2: Latent variable models, variational inference and variational autoencoders

Harri Lähdesmäki

Department of Computer Science  
Aalto University

March 24, 2025

## Outline

- Latent variable models (Section 6.2 from (Murphy, 2023))
- Linear latent variable models (Sec. 28.3)
- Deep latent variable models (Sec. 21.1)
- Variational inference: amortized, stochastic, gradient-based, reparameterized (Sec. 10.1-10.2)
- Variational autoencoders (Sec. 21)
- Reading: parts of Sections 6.2, 10.1-10.2, 21.1, 28.3 from (Murphy, 2023)

# Different types of latent variable models

## Latent variable model

A latent variable model (LVM) is a probabilistic model that contains variables that are always unobserved

Latent variables are random variable

Latent variables can also represent model parameters

Three types of latent variable models

- ① Global latent variables
- ② Local latent variables
- ③ Global and local latent variables

## Global latent variable models

- A model that contains global latent variables  $\theta$  that are shared across all  $N$  observations: e.g. model parameters
- E.g. the usual supervised learning setting

$$p(\mathbf{y}_{1:N}, \boldsymbol{\theta} | \mathbf{x}_{1:N}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta})$$

where  $\mathbf{y}_{1:N}$  is a shorthand for  $\mathbf{y}_1, \dots, \mathbf{y}_N$

- In Bayesian setting, the goal is to obtain the posterior

$$p(\boldsymbol{\theta} | \mathbf{x}_{1:N}, \mathbf{y}_{1:N})$$

- In maximum likelihood setting, the goal is to maximize the likelihood  $\prod_{n=1}^N p(\mathbf{y}_n | \mathbf{x}_n, \boldsymbol{\theta})$

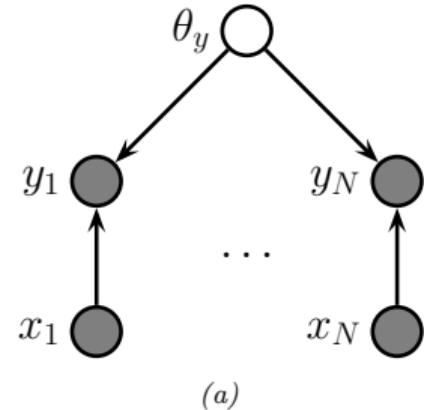


Figure 7.1a from (Murphy, 2023)

## Local latent variable models

- A model that contains local latent variables  $z_n$  ( $n \in \{1, \dots, N\}$ ) that are associated with individual observations  $x_n$
- Model parameters  $\theta = (\theta_z, \theta_x)$  are assumed to be deterministic or known
- The joint distribution typically factorizes

$$p(x_{1:N}, z_{1:N} | \theta) = \prod_{n=1}^N p(x_n | z_n, \theta_x) p(z_n | \theta_z)$$

- The goal is to obtain the posterior for each local latent

$$p(z_n | x_n, \theta)$$

- If parameters are not known, they can be treated as hyperparameters and optimized using e.g. the maximum likelihood

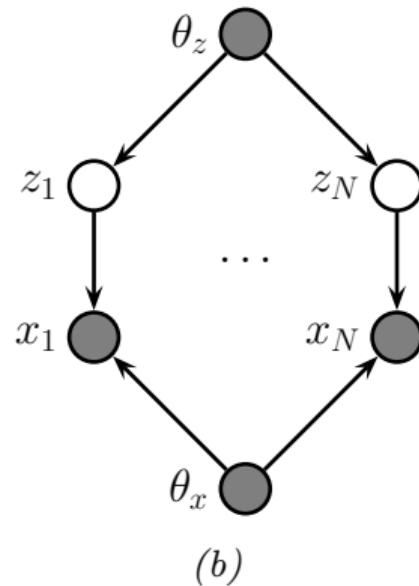


Figure 7.1b from (Murphy, 2023)

## Global and local latent variable models

- A model that contains both
  - global latent variables  $\theta = (\theta_z, \theta_x)$
  - local latent variables  $z_n$  ( $n \in \{1, \dots, N\}$ ), that are associated with individual observations  $x_n$
- The joint distribution typically factorizes

$$p(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}_z)p(\boldsymbol{\theta}_x) \prod_{n=1}^N p(\mathbf{x}_n | z_n, \boldsymbol{\theta}_x)p(z_n | \boldsymbol{\theta}_z)$$

- The goal is to obtain the posterior for unknowns

$$p(\mathbf{z}_{1:N}, \boldsymbol{\theta} | \mathbf{x}_{1:N})$$

- Benefits in considering  $\boldsymbol{\theta}$  as a random variable?

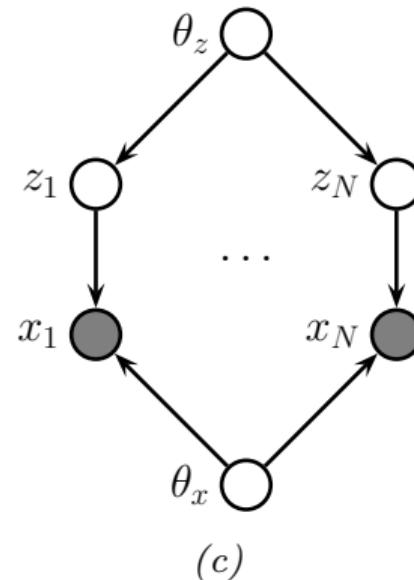


Figure 7.1c from (Murphy, 2023)

## Marginalization of latent variables

An indirect approach to define a probability distribution  $p(\mathbf{x} | \boldsymbol{\theta})$  via a joint distribution  $p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$

Target distribution  $p(\mathbf{x} | \boldsymbol{\theta})$  is obtained via marginalization

$$p(\mathbf{x} | \boldsymbol{\theta}) = \int p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}_x) p(\mathbf{z} | \boldsymbol{\theta}_z) d\mathbf{z}$$

## Marginalization of latent variables

An indirect approach to define a probability distribution  $p(\mathbf{x} | \boldsymbol{\theta})$  via a joint distribution  $p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$

Target distribution  $p(\mathbf{x} | \boldsymbol{\theta})$  is obtained via marginalization

$$p(\mathbf{x} | \boldsymbol{\theta}) = \int p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}_x) p(\mathbf{z} | \boldsymbol{\theta}_z) d\mathbf{z}$$

- Example: 1-D mixture of three Gaussians with

$$\begin{aligned} p(z = k) &= p_k, \quad k \in \{1, 2, 3\} \\ p(x | z = k) &= \mathcal{N}(x | \mu_k, \sigma_k^2) \end{aligned}$$

and

$$\begin{aligned} p(x) &= \sum_{k=1}^3 p(x | z_k) p(z_k) \\ &= \sum_{k=1}^3 p_k \cdot p(x | \mu_k, \sigma_k^2) \end{aligned}$$

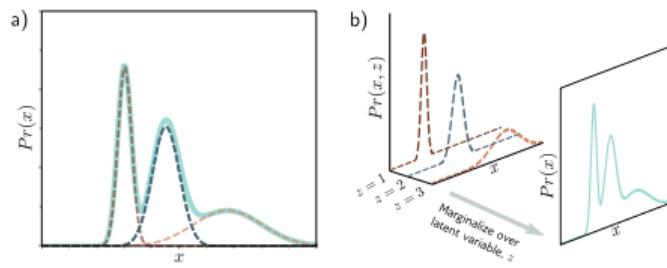


Figure 17.1 from (Prince, 2023)

## Generation

A new data point  $x^*$  can be generated from a latent variable model using the ancestral sampling

$$\begin{aligned}z^* &\sim p(z \mid \theta_z) \\x^* \mid z^* &\sim p(x \mid z^*, \theta_x)\end{aligned}$$

This is straightforward as long as distributions  $p(z \mid \theta_z)$  and  $p(x \mid z, \theta_x)$  are easy to sample

## Probabilistic factor analysis

Probabilistic factor analysis (FA) model is a linear (local) latent variable model defined as

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$$
$$p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{x} \mid \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

where

- $\mathbf{z} \in \mathbb{R}^L$  is a latent variable
- $\mathbf{x} \in \mathbb{R}^D$  denotes observed data
- $\mathbf{W} \in \mathbb{R}^{D \times L}$  is the factor loading matrix
- $\boldsymbol{\mu}$  is a mean offset
- $\boldsymbol{\Psi}$  is  $D$ -by- $D$  covariance matrix
- parameters are deterministic (not r.v.s)

## Probabilistic factor analysis

Probabilistic factor analysis (FA) model is a linear (local) latent variable model defined as

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$$
$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}),$$

where

- $\mathbf{z} \in \mathbb{R}^L$  is a latent variable
- $\mathbf{x} \in \mathbb{R}^D$  denotes observed data
- $\mathbf{W} \in \mathbb{R}^{D \times L}$  is the factor loading matrix
- $\boldsymbol{\mu}$  is a mean offset
- $\boldsymbol{\Psi}$  is  $D$ -by- $D$  covariance matrix
- parameters are deterministic (not r.v.s)

We can marginalize out the latent variable  $\mathbf{z}$  from  $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$  to get marginal likelihood

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \int \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})d\mathbf{z} \\ &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \boldsymbol{\Psi}) \end{aligned}$$

and

$$p(\mathbf{x}_{1:N}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \boldsymbol{\Psi})$$

## Probabilistic factor analysis: illustration

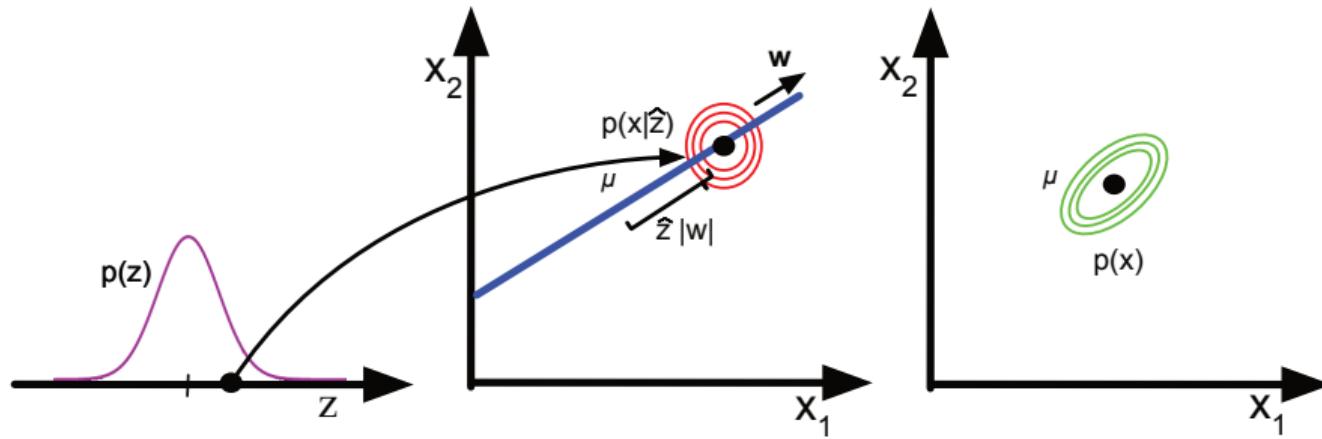


Figure 28.6: Illustration of the FA generative process, where we have  $L = 1$  latent dimension generating  $D = 2$  observed dimensions; we assume  $\Psi = \sigma^2 \mathbf{I}$ . The latent factor has value  $z \in \mathbb{R}$ , sampled from  $p(z)$ ; this gets mapped to a 2d offset  $\delta = zw$ , where  $w \in \mathbb{R}^2$ , which gets added to  $\mu$  to define a Gaussian  $p(\mathbf{x}|z) = \mathcal{N}(\mathbf{x}|\mu + \delta, \sigma^2 \mathbf{I})$ . By integrating over  $z$ , we “slide” this circular Gaussian “spray can” along the principal component axis  $w$ , which induces elliptical Gaussian contours in  $\mathbf{x}$  space centered on  $\mu$ . Adapted from Figure 12.9 of [Bis06].

Figure 28.3 from (Murphy, 2023)

## Probabilistic principle component analysis

Probabilistic principle component analysis (PPCA) is a special case of the probabilistic FA model where

- columns of  $\mathbf{W}$  are orthogonal
- $\Psi = \sigma^2 \mathbf{I}$ ,

and the generative model is

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x} | \mathbf{z}) = N(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

- The marginal likelihood

$$\begin{aligned} p(\mathbf{x} | \mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= \int \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) d\mathbf{z} \\ &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}) \end{aligned}$$

## Probabilistic principle component analysis: parameter learning

Given  $N$  observations  $\mathcal{D} = \mathbf{x}_{1:N}$ , we can estimate the model parameters by maximizing the marginal likelihood w.r.t.  $\mathbf{W}$ ,  $\boldsymbol{\mu}$  and  $\sigma^2$

$$L(\mathbf{W}, \boldsymbol{\mu}, \sigma^2 \mid \mathcal{D}) = p(\mathbf{x}_{1:N} \mid \mathbf{W}, \boldsymbol{\mu}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})$$

- Compute sample mean  $\boldsymbol{\mu} = \bar{\mathbf{x}}$  and covariance matrix  $\mathbf{S} = \frac{1}{N} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$
- Rewrite  $\mathbf{S}$  using the eigenvector-eigenvalue decomposition  $\mathbf{S} = \mathbf{U}\Lambda\mathbf{U}^\top$
- ML parameters are obtained using the  $L$  largest eigenvalues-vectors

$$\mathbf{W} = \mathbf{U}_L(\Lambda_L - \sigma^2\mathbf{I})^{\frac{1}{2}} \text{ (upto arbitrary rotation)}$$

$$\sigma^2 = \frac{1}{D-L} \sum_{i=L+1}^D \lambda_i$$

$$\boldsymbol{\mu} = \bar{\mathbf{x}}$$

## Probabilistic principle component analysis: posterior

- Given parameters  $\mathbf{W}$ ,  $\mu$ ,  $\sigma^2$  (fixed or optimized) and an observation  $\mathbf{x}$ , we want to know the distribution of the latent variable

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

- The posterior of the latent variable can be shown to have normal distribution

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x} - \mu), \sigma^2\mathbf{M}^{-1}),$$

where  $\mathbf{M} = \mathbf{W}^\top\mathbf{W} + \sigma^2\mathbf{I}$

## Probabilistic principle component analysis: posterior

- Given parameters  $\mathbf{W}$ ,  $\mu$ ,  $\sigma^2$  (fixed or optimized) and an observation  $\mathbf{x}$ , we want to know the distribution of the latent variable

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

- The posterior of the latent variable can be shown to have normal distribution

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x} - \mu), \sigma^2\mathbf{M}^{-1}),$$

where  $\mathbf{M} = \mathbf{W}^\top\mathbf{W} + \sigma^2\mathbf{I}$

- In the noise-free case of  $\sigma^2 = 0$  the PPCA and PCA are directly comparable

# Deep latent variable model (also called as non-linear latent variable models)

Deep latent variable model is a generative model

$$z \sim p_{\theta}(z)$$

$$x | z \sim p_{\theta}(x | z) \triangleq p(x | d_{\theta}(z)),$$

where

- $z \in \mathcal{Z} = \mathbb{R}^L$  is a latent variable
- $p_{\theta}(z)$  is a prior distribution, typically Gaussian
- $x \in \mathcal{X}$  denotes observed data, e.g.  $\mathcal{X} = \mathbb{R}^D$ , where typically  $D \gg L$
- $p(x | v)$  is a distribution from the exponential family with parameters  $v \in \Upsilon$  (e.g.  $\mu, \Sigma$  for Gaussian, or  $p = (p_1, \dots, p_D)$  for Bernoulli)
- $d_{\theta}(\cdot) : \mathcal{Z} \rightarrow \Upsilon$  is a non-linear function called decoder that maps the latent variable to parameters of the observation likelihood (e.g.  $(\mu, \Sigma) = d_{\theta}(z)$ )

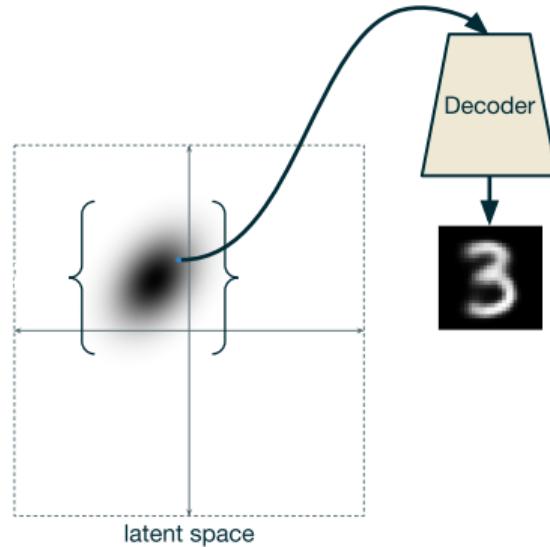


Figure adapted from <https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

# Deep latent variable model (also called as non-linear latent variable models)

Deep latent variable model is a generative model

$$z \sim p_{\theta}(z)$$

$$x | z \sim p_{\theta}(x | z) \triangleq p(x | d_{\theta}(z)),$$

where

- $z \in \mathcal{Z} = \mathbb{R}^L$  is a latent variable
- $p_{\theta}(z)$  is a prior distribution, typically Gaussian
- $x \in \mathcal{X}$  denotes observed data, e.g.  $\mathcal{X} = \mathbb{R}^D$ , where typically  $D \gg L$
- $p(x | v)$  is a distribution from the exponential family with parameters  $v \in \Upsilon$  (e.g.  $\mu, \Sigma$  for Gaussian, or  $p = (p_1, \dots, p_D)$  for Bernoulli)
- $d_{\theta}(\cdot) : \mathcal{Z} \rightarrow \Upsilon$  is a non-linear function called decoder that maps the latent variable to parameters of the observation likelihood (e.g.  $(\mu, \Sigma) = d_{\theta}(z)$ )

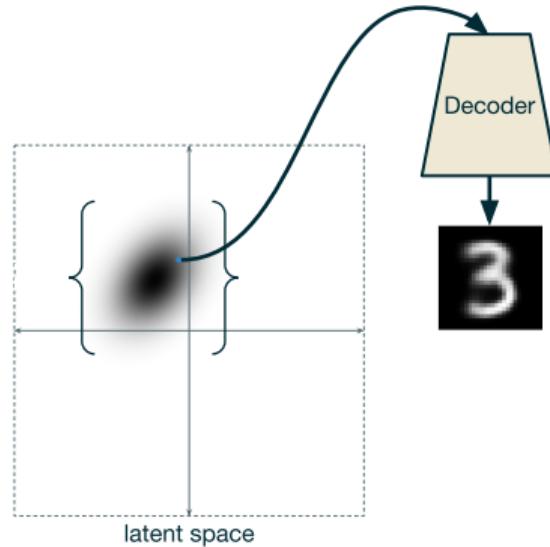


Figure adapted from <https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

Marg. likelihood  $p_{\theta}(x) = \int p(x | d_{\theta}(z))p_{\theta}(z)dz$  is generally intractable

## Variational inference

Assume a probabilistic model with

- latent variables  $z$
- observed variables  $x$
- fixed parameters  $\theta$  (if random variables, add  $\theta$  to  $z$ )

Interested in the posterior

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

where the evidence  $p_{\theta}(x) = \int p_{\theta}(x | z)p_{\theta}(z)dz$  is intractable

## Variational inference

Assume a probabilistic model with

- latent variables  $z$
- observed variables  $x$
- fixed parameters  $\theta$  (if random variables, add  $\theta$  to  $z$ )

Interested in the posterior

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

where the evidence  $p_{\theta}(x) = \int p_{\theta}(x | z)p_{\theta}(z)dz$  is intractable

### Variational inference

Variational inference (VI) approximates  $p_{\theta}(z | x)$  with a variational approximation  $q(z)$  by minimizing the reverse KL divergence over a variational family of distributions  $\mathcal{Q}$

$$q^* = \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q(z) || p_{\theta}(z | x)),$$

# Variational inference

Assume a probabilistic model with

- latent variables  $z$
- observed variables  $x$
- fixed parameters  $\theta$  (if random variables, add  $\theta$  to  $z$ )

Interested in the posterior

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

where the evidence  $p_{\theta}(x) = \int p_{\theta}(x | z)p_{\theta}(z)dz$  is intractable

## Variational inference

Variational inference (VI) approximates  $p_{\theta}(z | x)$  with a variational approximation  $q(z)$  by minimizing the reverse KL divergence over a variational family of distributions  $\mathcal{Q}$

$$\psi^* = \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q(z) || p_{\theta}(z | x)),$$

In practice we use a parametric family  $\mathcal{Q}$  with variational parameters  $\psi$ ,  $q_{\psi}(z)$

- e.g.:  $q_{\psi}(z) = \mathcal{N}(z | \mu, \Sigma)$ ,  $\psi = (\mu, \Sigma)$

Variational inference converts posterior inference to optimization

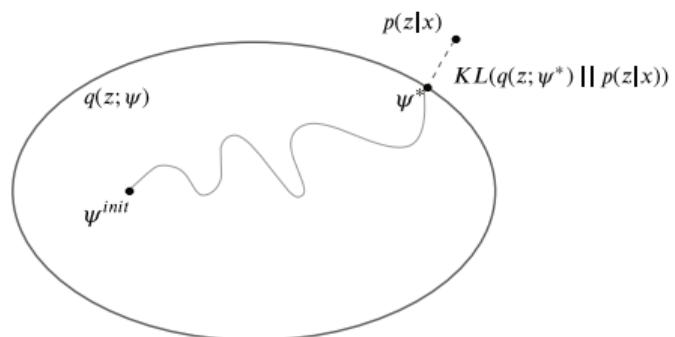


Figure 10.1 from (Murphy, 2023)

## Variational inference: objective function

The variational objective can be written as

$$\begin{aligned} D_{\text{KL}}(q_{\psi}(z) \parallel p_{\theta}(z \mid \mathbf{x})) &= \mathbb{E}_{q_{\psi}(z)} \left[ \log \left( \frac{q_{\psi}(z)}{p_{\theta}(z \mid \mathbf{x})} \right) \right] \\ &= \mathbb{E}_{q_{\psi}(z)} \left[ \log q_{\psi}(z) - \log \left( \frac{p_{\theta}(\mathbf{x} \mid z)p_{\theta}(z)}{p_{\theta}(\mathbf{x})} \right) \right] \\ &= \mathbb{E}_{q_{\psi}(z)} [\log q_{\psi}(z) - \log p_{\theta}(\mathbf{x} \mid z) - \log p_{\theta}(z)] + \log p_{\theta}(\mathbf{x}) \\ &= \mathbb{E}_{q_{\psi}(z)} [\log q_{\psi}(z) - \log p_{\theta}(\mathbf{x}, z)] + \log p_{\theta}(\mathbf{x}) \\ &= \underbrace{\mathbb{E}_{q_{\psi}(z)} \left[ \log \left( \frac{q_{\psi}(z)}{p_{\theta}(\mathbf{x}, z)} \right) \right]}_{\tilde{\mathcal{L}}(\theta, \psi \mid \mathbf{x})} + \log p_{\theta}(\mathbf{x}), \end{aligned}$$

where  $\log p_{\theta}(\mathbf{x})$  is intractable but independent of  $\psi$

We can minimize  $\tilde{\mathcal{L}}(\theta, \psi \mid \mathbf{x})$  w.r.t.  $\psi$

## Variational inference: evidence lower bound

Alternatively, we can derive a lower bound for the evidence (marg. likelihood)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \left( \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \right) \\ &= \log \left( \int q_{\psi}(\mathbf{z}) \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} d\mathbf{z} \right) \\ &= \log \left( \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right] \right) \\ &\stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right) \right] \\ &= \mathcal{L}(\theta, \psi \mid \mathbf{x})\end{aligned}$$

$\mathcal{L}(\theta, \psi \mid \mathbf{x})$  is called the evidence lower bound (ELBO): optimize w.r.t.  $\theta$  and  $\psi$

## Variational inference: evidence lower bound

Alternatively, we can derive a lower bound for the evidence (marg. likelihood)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \left( \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \right) \\ &= \log \left( \int q_{\psi}(\mathbf{z}) \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} d\mathbf{z} \right) \\ &= \log \left( \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right] \right) \\ &\stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right) \right] \\ &= \mathcal{L}(\theta, \psi | \mathbf{x})\end{aligned}$$

$\mathcal{L}(\theta, \psi | \mathbf{x})$  is called the evidence lower bound (ELBO): optimize w.r.t.  $\theta$  and  $\psi$

Notice that  $\mathcal{L}(\theta, \psi | \mathbf{x}) = -\tilde{\mathcal{L}}(\theta, \psi | \mathbf{x})$

→ Maximizing  $\mathcal{L}(\theta, \psi | \mathbf{x})$  equals minimizing  $\tilde{\mathcal{L}}(\theta, \psi | \mathbf{x})$

→ Minimizing  $D_{\text{KL}}(q_{\psi}(\mathbf{z}) || p_{\theta}(\mathbf{z} | \mathbf{x}))$  equals maximizing a lower bound for  $\log p_{\theta}(\mathbf{x})$

## ELBO as reconstruction loss minus KL divergence

We can also write the ELBO as

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) \\&= \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\psi}}(\mathbf{z})} \right) \right] \\&= \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} \left[ \log p_{\theta}(\mathbf{x} \mid \mathbf{z}) + \log \left( \frac{p_{\theta}(\mathbf{z})}{q_{\boldsymbol{\psi}}(\mathbf{z})} \right) \right] \\&= \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - \mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} \left[ \log \left( \frac{q_{\boldsymbol{\psi}}(\mathbf{z})}{p_{\theta}(\mathbf{z})} \right) \right] \\&= \underbrace{\mathbb{E}_{q_{\boldsymbol{\psi}}(\mathbf{z})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})]}_{\text{expected log likelihood}} - \underbrace{D_{\text{KL}}(q_{\boldsymbol{\psi}}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}))}_{\text{KL of the posterior from prior}}\end{aligned}$$

The first term measures the reconstruction loss

The second term can be interpreted as prior regularization

## ELBO as reconstruction loss minus KL divergence

We can also write the ELBO as

$$\begin{aligned}\log p_{\theta}(x) &\geq \mathcal{L}(\theta, \psi | x) \\&= \mathbb{E}_{q_{\psi}(z)} \left[ \log \left( \frac{p_{\theta}(x, z)}{q_{\psi}(z)} \right) \right] \\&= \mathbb{E}_{q_{\psi}(z)} \left[ \log p_{\theta}(x | z) + \log \left( \frac{p_{\theta}(z)}{q_{\psi}(z)} \right) \right] \\&= \mathbb{E}_{q_{\psi}(z)} [\log p_{\theta}(x | z)] - \mathbb{E}_{q_{\psi}(z)} \left[ \log \left( \frac{q_{\psi}(z)}{p_{\theta}(z)} \right) \right] \\&= \underbrace{\mathbb{E}_{q_{\psi}(z)} [\log p_{\theta}(x | z)]}_{\text{expected log likelihood}} - \underbrace{D_{\text{KL}}(q_{\psi}(z) || p_{\theta}(z))}_{\text{KL of the posterior from prior}}\end{aligned}$$

The first term measures the reconstruction loss

The second term can be interpreted as prior regularization

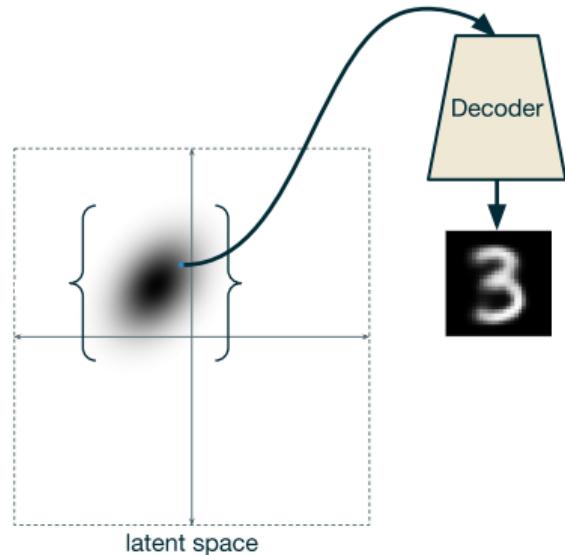


Figure adapted from

<https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

## ELBO: Monte Carlo estimate

ELBO involves expectations that generally cannot be computed analytically but can be estimated using Monte Carlo

For the standard ELBO formula

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right) \right] \\ &\approx \frac{1}{S} \left( \sum_{s=1}^S \log p_{\theta}(\mathbf{x}, \mathbf{z}_s) - \log q_{\psi}(\mathbf{z}_s) \right) \\ &= \frac{1}{S} \left( \sum_{s=1}^S \log p_{\theta}(\mathbf{x} \mid \mathbf{z}_s) + \log p_{\theta}(\mathbf{z}_s) \right. \\ &\quad \left. - \log q_{\psi}(\mathbf{z}_s) \right)\end{aligned}$$

where  $\mathbf{z}_s \stackrel{\text{i.i.d}}{\sim} q_{\psi}(\mathbf{z})$

## ELBO: Monte Carlo estimate

ELBO involves expectations that generally cannot be computed analytically but can be estimated using Monte Carlo

For the standard ELBO formula

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\psi}(\mathbf{z})} \left[ \log \left( \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\psi}(\mathbf{z})} \right) \right] \\ &\approx \frac{1}{S} \left( \sum_{s=1}^S \log p_{\theta}(\mathbf{x}, \mathbf{z}_s) - \log q_{\psi}(\mathbf{z}_s) \right) \\ &= \frac{1}{S} \left( \sum_{s=1}^S \log p_{\theta}(\mathbf{x} | \mathbf{z}_s) + \log p_{\theta}(\mathbf{z}_s) \right. \\ &\quad \left. - \log q_{\psi}(\mathbf{z}_s) \right)\end{aligned}$$

where  $\mathbf{z}_s \stackrel{\text{i.i.d}}{\sim} q_{\psi}(\mathbf{z})$

For the reconstruction minus KL formula

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\psi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - D_{\text{KL}}(q_{\psi}(\mathbf{z}) || p_{\theta}(\mathbf{z})) \\ &\approx \frac{1}{S} \sum_{s=1}^S \log p_{\theta}(\mathbf{x} | \mathbf{z}_s) - D_{\text{KL}}(q_{\psi}(\mathbf{z}) || p_{\theta}(\mathbf{z}))\end{aligned}$$

where  $\mathbf{z}_s \stackrel{\text{i.i.d}}{\sim} q_{\psi}(\mathbf{z})$

KL term has a closed-form solution for a pair of distributions from the same exponential family

## Tightness of the ELBO

From the variational bound

$$D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) = \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + \log p_{\theta}(\mathbf{x})$$

we get

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= -\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) \\ &= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) \\ &\geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}),\end{aligned}$$

because KL divergence is non-negative

## Tightness of the ELBO

From the variational bound

$$D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) = \tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + \log p_{\theta}(\mathbf{x})$$

we get

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= -\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) \\ &= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}) + D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) \\ &\geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} \mid \mathbf{x}),\end{aligned}$$

because KL divergence is non-negative

KL divergence  $D_{\text{KL}}(q_{\psi}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x}))$  defines us two “distances”

- ➊ The KL divergence of the variational approximation from the true posterior (i.e., the VI objective)
- ➋ The tightness of the bound, i.e., the gap between the ELBO and  $\log p_{\theta}(\mathbf{x})$

## Variational inference: mean field model

For a standard local (deep) latent variable model with  $N$  observations  $\mathcal{D} = \mathbf{x}_{1:N}$  and deterministic parameters  $\theta$

$$\begin{aligned} p_{\theta}(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}) &= p_{\theta}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p_{\theta}(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\theta}(\mathbf{z}_n)) p_{\theta}(\mathbf{z}_n) \end{aligned}$$

The posterior factorize accordingly

$$\begin{aligned} p_{\theta}(\mathbf{z}_{1:N} \mid \mathbf{x}_{1:N}) &\propto p_{\theta}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\theta}(\mathbf{z}_n)) p_{\theta}(\mathbf{z}_n) \end{aligned}$$

## Variational inference: mean field model

For a standard local (deep) latent variable model with  $N$  observations  $\mathcal{D} = \mathbf{x}_{1:N}$  and deterministic parameters  $\boldsymbol{\theta}$

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}) &= p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p_{\boldsymbol{\theta}}(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\boldsymbol{\theta}}(\mathbf{z}_n)) p_{\boldsymbol{\theta}}(\mathbf{z}_n) \end{aligned}$$

The posterior factorize accordingly

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{z}_{1:N} \mid \mathbf{x}_{1:N}) &\propto p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\boldsymbol{\theta}}(\mathbf{z}_n)) p_{\boldsymbol{\theta}}(\mathbf{z}_n) \end{aligned}$$

It is common to use the so-called mean field variational approximation that factorizes accordingly

$$q(\mathbf{z}_{1:N} \mid \boldsymbol{\psi}_{1:N}) = \prod_{n=1}^N q_{\boldsymbol{\psi}_n}(\mathbf{z}_n)$$

Note: each latent variable  $\mathbf{z}_n$  has its own variational parameters  $\boldsymbol{\psi}_n$

## Variational inference: mean field model

For a standard local (deep) latent variable model with  $N$  observations  $\mathcal{D} = \mathbf{x}_{1:N}$  and deterministic parameters  $\boldsymbol{\theta}$

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N}, \mathbf{z}_{1:N}) &= p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p_{\boldsymbol{\theta}}(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\boldsymbol{\theta}}(\mathbf{z}_n)) p_{\boldsymbol{\theta}}(\mathbf{z}_n) \end{aligned}$$

The posterior factorize accordingly

$$\begin{aligned} p_{\boldsymbol{\theta}}(\mathbf{z}_{1:N} \mid \mathbf{x}_{1:N}) &\propto p_{\boldsymbol{\theta}}(\mathbf{x}_{1:N} \mid \mathbf{z}_{1:N}) p(\mathbf{z}_{1:N}) \\ &= \prod_{n=1}^N p(\mathbf{x}_n \mid d_{\boldsymbol{\theta}}(\mathbf{z}_n)) p_{\boldsymbol{\theta}}(\mathbf{z}_n) \end{aligned}$$

It is common to use the so-called mean field variational approximation that factorizes accordingly

$$q(\mathbf{z}_{1:N} \mid \boldsymbol{\psi}_{1:N}) = \prod_{n=1}^N q_{\boldsymbol{\psi}_n}(\mathbf{z}_n)$$

Note: each latent variable  $\mathbf{z}_n$  has its own variational parameters  $\boldsymbol{\psi}_n$

Because

- the above mean field approximation factorizes
- log likelihood  $p_{\boldsymbol{\theta}}(\mathbf{x}_{1:n} \mid \mathbf{z}_{1:N})$  factorizes
- KL for independent r.v.s. factorizes

the ELBO for  $\log p_{\boldsymbol{\theta}}(\mathcal{D})$  factorizes as well

$$\log p_{\boldsymbol{\theta}}(\mathcal{D}) \geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_{1:N} \mid \mathcal{D}) = \sum_{n=1}^N \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_n \mid \mathbf{x}_n)$$

## Variational inference and parameter estimation

- If generative model parameter  $\theta$  are unknown, we generally optimize them by maximizing the (intractable) log marginal likelihood  $\log p_\theta(\mathcal{D})$
- Again, the ELBO provides a lower bound:  $\log p_\theta(\mathcal{D}) \geq \mathcal{L}(\theta, \psi_{1:N} \mid \mathcal{D})$
- **Variational EM** algorithm iteratively optimizes  $\psi_{1:N}$  and  $\theta$ 
  - E step: maximize  $\mathcal{L}(\theta, \psi_{1:N} \mid \mathcal{D})$  w.r.t.  $\psi_{1:N}$
  - M step: maximize  $\mathcal{L}(\theta, \psi_{1:N} \mid \mathcal{D})$  w.r.t.  $\theta$
- Similar approach, called **variational Bayes**, if generative model parameters  $\theta$  are random variables  $p(\theta \mid \xi)$

## Variational inference: stochastic VI

- Computing and minimizing the ELBO for large datasets (e.g. millions of data points) can be slow
- An unbiased stochastic approximation with a minibatch  $\mathcal{B} \subseteq \mathcal{D}$  of size  $B = |\mathcal{B}|$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_{1:N} \mid \mathcal{D}) &= \sum_{n=1}^N \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi}_n \mid \mathbf{x}_n) \\ &\approx \frac{N}{B} \sum_{\mathbf{x}_n \in \mathcal{B}} \left[ \mathbb{E}_{q_{\boldsymbol{\psi}_n}(\mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_n \mid \mathbf{z}_n)] - D_{\text{KL}}(q_{\boldsymbol{\psi}_n}(\mathbf{z}_n) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}_n)) \right]\end{aligned}$$

- This corresponds to stochastic gradient descent (SGD) optimization algorithm and scales VI to large datasets

## Variational inference: amortized VI

- For large datasets, optimizing  $N$  local variational parameters  $\psi_{1:N}$  can be slow
- Amortized VI trains a separate deterministic mapping, so-called encoder or inference network, that predicts  $\psi_n$  from  $x_n$

$$\psi_n = f_\phi(x_n)$$

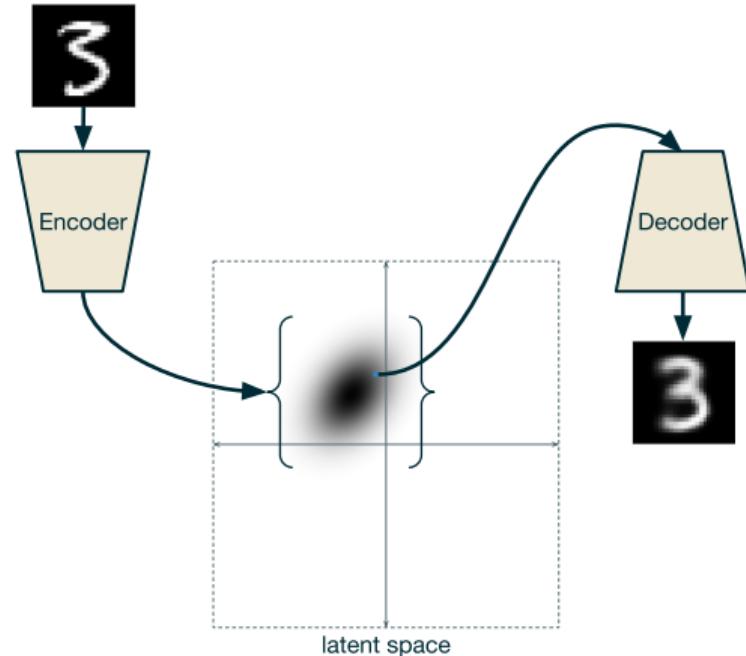
- Parameters  $\phi$  are shared across samples  $N$
- Example: assume

$$q_{\psi_n}(z_n) = \mathcal{N}(z_n | \mu_n, \text{diag}(\sigma_n^2))$$

$$\psi_n = (\mu_n, \log \sigma_n^2)$$

then

$$\psi_n = (\mu_n, \log \sigma^2) = f_\phi(x_n)$$



<https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

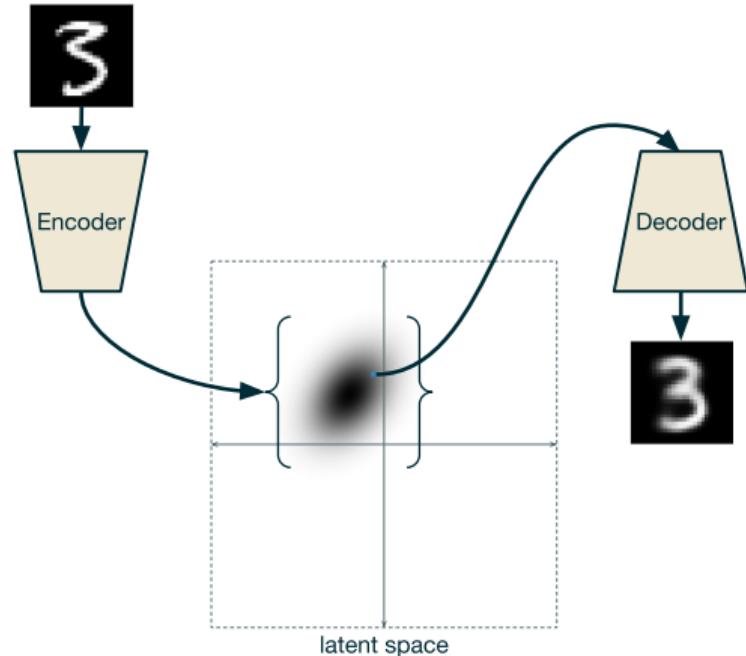
## Variational inference: amortized VI

- The amortized version of the standard variational approximation  $q_{\psi_n}(\mathbf{z}_n)$  is denoted as

$$q(\mathbf{z}_n \mid f_{\phi}(\mathbf{x}_n)) = q_{\phi}(\mathbf{z}_n \mid \mathbf{x}_n)$$

- The ELBO objective becomes

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathcal{D}) = & \sum_{n=1}^N \left[ \mathbb{E}_{q_{\phi}(\mathbf{z}_n \mid \mathbf{x}_n)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_n \mid \mathbf{z}_n)] \right. \\ & \left. - D_{\text{KL}}(q_{\phi}(\mathbf{z}_n \mid \mathbf{x}_n) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}_n)) \right]\end{aligned}$$



<https://ijdykeman.github.io/ml/2016/12/21/cvae.html>

## Variational inference: gradient-based VI

- The gradient of the ELBO w.r.t. generative model parameters

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &\stackrel{\text{MC}}{\approx} \sum_{s=1}^S \nabla_{\theta} \log p_{\theta}(\mathbf{x} \mid \mathbf{z}_s) - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z}))\end{aligned}$$

where  $\mathbf{z}_s \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z} \mid \mathbf{x})$  and assuming the KL divergence has a closed-form solution (computation does not involve any random variables)

## Variational inference: gradient-based VI

- The gradient of the ELBO w.r.t. generative model parameters

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(z \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid z)] - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(z \mid \mathbf{x}) \parallel p_{\theta}(z)) \\ &= \mathbb{E}_{q_{\phi}(z \mid \mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x} \mid z)] - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(z \mid \mathbf{x}) \parallel p_{\theta}(z)) \\ &\stackrel{\text{MC}}{\approx} \sum_{s=1}^S \nabla_{\theta} \log p_{\theta}(\mathbf{x} \mid z_s) - \nabla_{\theta} D_{\text{KL}}(q_{\phi}(z \mid \mathbf{x}) \parallel p_{\theta}(z))\end{aligned}$$

where  $z_s \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(z \mid \mathbf{x})$  and assuming the KL divergence has a closed-form solution (computation does not involve any random variables)

- The gradient of the ELBO w.r.t. the inference network parameters  $\phi$  is trickier

$$\begin{aligned}\nabla_{\phi} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(z \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid z)] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(z \mid \mathbf{x}) \parallel p_{\theta}(z)) \\ &\neq \mathbb{E}_{q_{\phi}(z \mid \mathbf{x})} [\nabla_{\phi} \log p_{\theta}(\mathbf{x} \mid z)] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(z \mid \mathbf{x}) \parallel p_{\theta}(z))\end{aligned}$$

## Reparametrization trick

The reparametrization trick is to rewrite a random variable  $\mathbf{z} \sim q_\phi(\mathbf{z} \mid \mathbf{x})$  using a differentiable and invertible transformation  $g$  of another random variable  $\epsilon \sim p(\epsilon)$  that does not depend on  $\phi$ , i.e.,

$$\begin{aligned}\epsilon &\sim p(\epsilon) \\ \mathbf{z} &= g(\phi, \mathbf{x}, \epsilon)\end{aligned}$$

An example: if  $\mathbf{z} \sim q_\phi(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$  where  $(\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) = f_\phi(\mathbf{x})$ , we can use

$$\begin{aligned}\epsilon &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{z} &= \boldsymbol{\mu} + \text{diag}(\boldsymbol{\sigma})\epsilon\end{aligned}$$

## Reparametrization trick

The reparametrization trick is to rewrite a random variable  $\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x})$  using a differentiable and invertible transformation  $g$  of another random variable  $\epsilon \sim p(\epsilon)$  that does not depend on  $\phi$ , i.e.,

$$\begin{aligned}\epsilon &\sim p(\epsilon) \\ \mathbf{z} &= g(\phi, \mathbf{x}, \epsilon)\end{aligned}$$

An example: if  $\mathbf{z} \sim q_\phi(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$  where  $(\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) = f_\phi(\mathbf{x})$ , we can use

$$\begin{aligned}\epsilon &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{z} &= \boldsymbol{\mu} + \text{diag}(\boldsymbol{\sigma})\epsilon\end{aligned}$$

Expectation of a function  $h(\mathbf{z})$  w.r.t.  $q_\phi(\mathbf{z} | \mathbf{x})$  can be written as

$$\mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [h(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} [h(\mathbf{z}')] = \mathbb{E}_{p(\epsilon)} [h(g(\phi, \mathbf{x}, \epsilon))]$$

where  $\mathbf{z}' = g(\phi, \mathbf{x}, \epsilon)$ , and gradient accordingly

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z} | \mathbf{x})} [h(\mathbf{z})] = \nabla_\phi \mathbb{E}_{p(\epsilon)} [h(\mathbf{z}')] = \mathbb{E}_{p(\epsilon)} [\nabla_\phi h(\mathbf{z}')$$

## Reparametrization: illustration

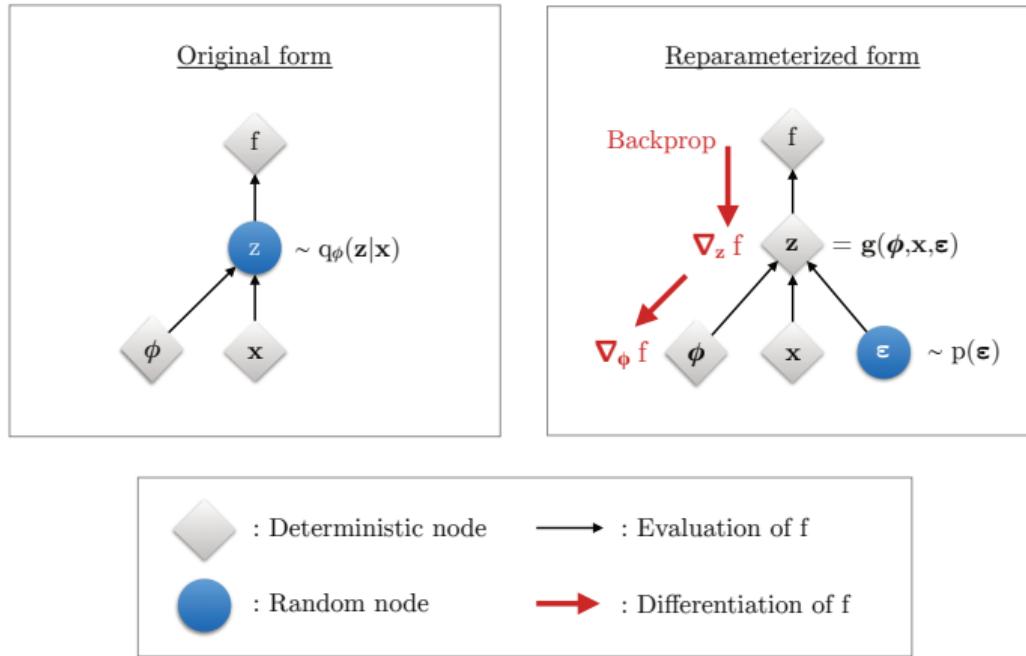


Figure 10.4 from (Murphy, 2023), originally from (Kingma & Welling, 2019)

## Variational inference: reparameterized VI

Using reparametrization, the gradient of the ELBO w.r.t. the inference network parameters  $\phi$  is

$$\begin{aligned}\nabla_{\phi} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &\stackrel{\text{reparam.}}{=} \nabla_{\phi} \mathbb{E}_{p(\boldsymbol{\epsilon})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z}')] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &= \mathbb{E}_{p(\boldsymbol{\epsilon})} [\nabla_{\phi} \log p_{\theta}(\mathbf{x} \mid \mathbf{z}')] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &\stackrel{\text{MC}}{\approx} \sum_{s=1}^S \nabla_{\phi} \log p_{\theta}(\mathbf{x} \mid \mathbf{z}'_s) - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z}))\end{aligned}$$

where  $\mathbf{z}' = g(\phi, \mathbf{x}, \boldsymbol{\epsilon})$ ,  $\mathbf{z}'_s = g(\phi, \mathbf{x}, \boldsymbol{\epsilon}_s)$ , and  $\boldsymbol{\epsilon}_s \stackrel{\text{i.i.d.}}{\sim} p(\boldsymbol{\epsilon})$

## Variational inference: reparameterized VI

Using reparametrization, the gradient of the ELBO w.r.t. the inference network parameters  $\phi$  is

$$\begin{aligned}\nabla_{\phi} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z})] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &\stackrel{\text{reparam.}}{=} \nabla_{\phi} \mathbb{E}_{p(\epsilon)} [\log p_{\theta}(\mathbf{x} \mid \mathbf{z}')] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_{\phi} \log p_{\theta}(\mathbf{x} \mid \mathbf{z}')] - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z})) \\ &\stackrel{\text{MC}}{\approx} \sum_{s=1}^S \nabla_{\phi} \log p_{\theta}(\mathbf{x} \mid \mathbf{z}'_s) - \nabla_{\phi} D_{\text{KL}}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \parallel p_{\theta}(\mathbf{z}))\end{aligned}$$

where  $\mathbf{z}' = g(\phi, \mathbf{x}, \epsilon)$ ,  $\mathbf{z}'_s = g(\phi, \mathbf{x}, \epsilon_s)$ , and  $\epsilon_s \stackrel{\text{i.i.d.}}{\sim} p(\epsilon)$

Reparameterized gradient is an unbiased estimate of the exact gradient

If the KL does not have a closed form solution, reparameterized gradients may be computed from

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} \mid \mathbf{x})] \\ \nabla_{\phi} \mathcal{L}(\theta, \phi \mid \mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} \mid \mathbf{x})]\end{aligned}$$

## Variational autoencoder

Variational autoencoder (VAE) is a neural architecture that consists of a deep latent variable model (deep generative model) of the form

$$\begin{aligned} z &\sim p_{\theta}(z) \\ x \mid z &\sim p(x \mid d_{\theta}(z)) \end{aligned}$$

that is trained using (reparameterized) amortized variational inference with an inference model  $q_{\phi}(z \mid x)$

- Data  $x$  can be continuous or discrete, or both, typically modeled using distribution from the exponential family
- Decoder  $d_{\theta}(\cdot)$  and encoder (inference network)  $f_{\phi}(\cdot)$  are commonly neural networks
- For general statistical latent variable models, this can be called as autoencoding variational Bayes method
- Variational approximation  $q_{\phi}(z \mid x)$  needs to admit reparametrization

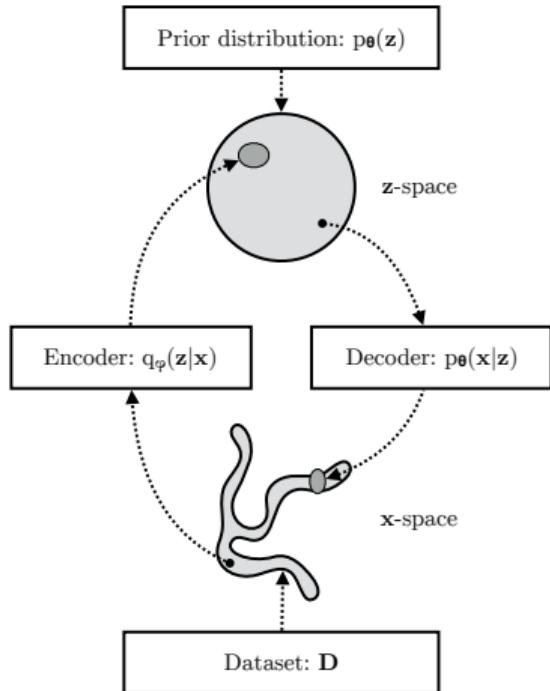
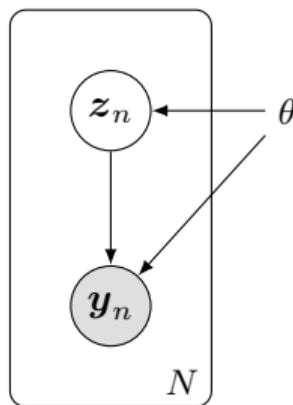


Figure from (Kingma & Welling, 2019)

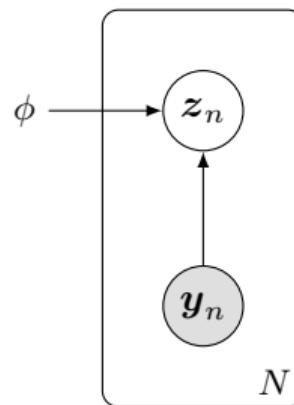
# Variational autoencoder: Probabilistic graphical model

## Probabilistic graphical models for VAE

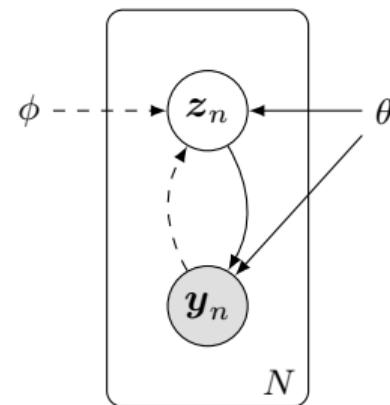
- Circled symbols denote random variables
- Non-circled symbols are hyperparameters
- White nodes (circles) are latent
- Shaded nodes are observed



(a) Generative model



(b) Inference model



(c) Combined

## Learning algorithm for VAE

Algorithm for obtaining maximum likelihood estimate for deep latent variable model using reparameterized, stochastic and amortized variational inference, here with

- a minibatch sample size of one
- Monte Carlo estimate with sample size one

---

**Algorithm 10.3:** Reparameterized amortized SVI for MLE of an LVM

---

```
1 Initialize  $\theta, \phi$ 
2 repeat
3   Sample  $x_n \sim p_{\mathcal{D}}$ 
4   Sample  $\epsilon_n \sim q_0$ 
5   Compute  $z_n = g(\phi, x_n, \epsilon_n)$ 
6   Compute  $\mathcal{L}(\theta, \phi | x_n, z_n) = -\log p_{\theta}(x_n, z_n) + \log q_{\phi}(z_n | x_n)$ 
7   Update  $\theta := \theta - \eta \nabla_{\theta} \mathcal{L}(\phi, \theta | x_n, z_n)$ 
8   Update  $\phi := \phi - \eta \nabla_{\phi} \mathcal{L}(\phi, \theta | x_n, z_n)$ 
9 until converged
```

---

Algorithm 10.3 from (Murphy, 2023)  
Note notation  $\tilde{\mathcal{L}}$  vs.  $\mathcal{L}$

## Variational autoencoder: neural architecture training pipeline

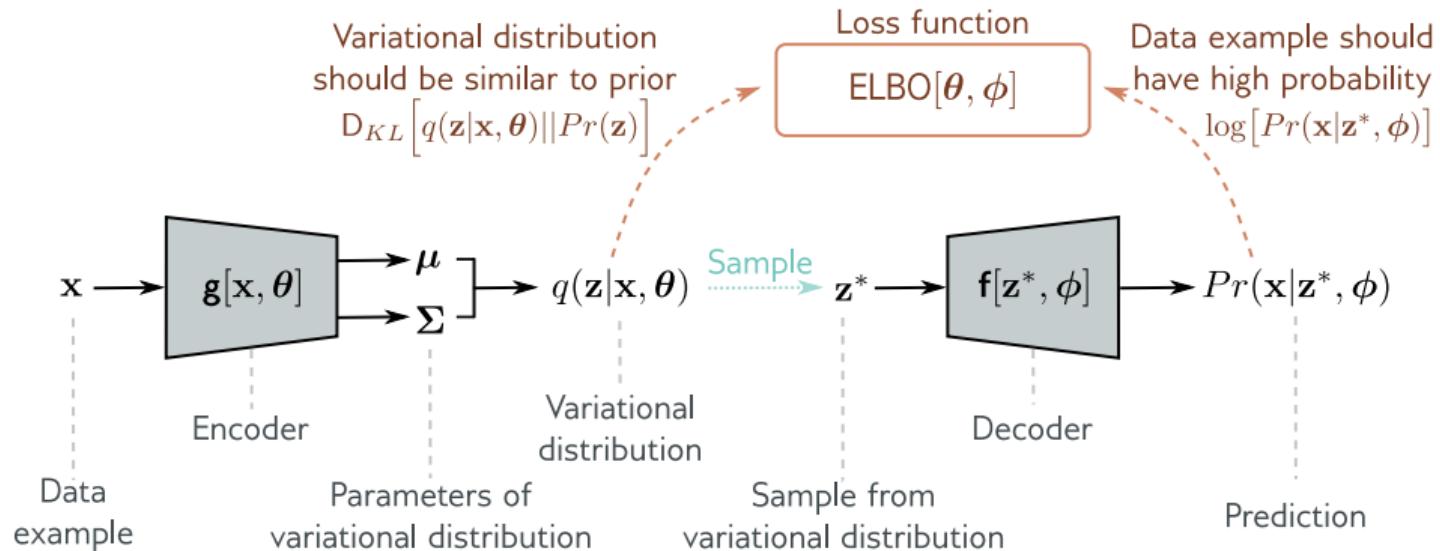


Figure 17.9 from (Prince, 2023). Note that notation differs.

## Variational autoencoder: parametrization, NN architectures

In standard VAE models:

- Prior  $p_\theta(z)$  is typically Gaussian
- Likelihood  $p(x | d_\theta(z))$  depends on application: typically Gaussian, Bernoulli, etc.
- Variational approximation  $q_\phi(z | x)$  is typically Gaussian (but can be more expressive, such as normalizing flows etc.)

Decoder  $d_\theta(\cdot)$  and encoder  $f_\phi(\cdot)$  neural network architectures are application dependent:

- MLPs for tabular data
- CNNs for images
- RNNs and attention for sequential data

## Autoencoder vs. variational autoencoder

Autoencoder (AE) is a deterministic neural network architecture that

- Maps input  $x$  into a lower dimensional representation  $z = f_e(x)$
- Attempts to reconstruct the original input by  $f_d(z)$ , i.e.,  $x \approx \hat{x} = f_d(f_e(x))$

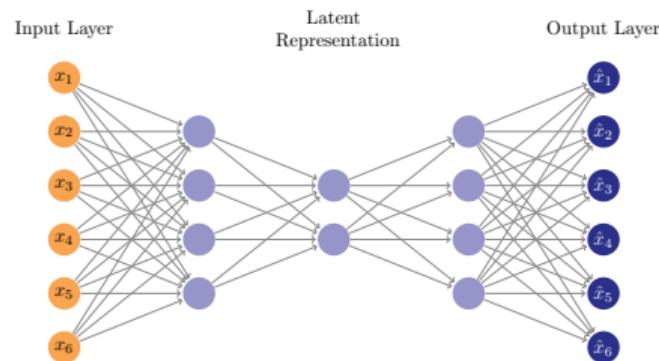


Figure 16.13 from (Murphy, 2023)

## Autoencoder vs. variational autoencoder

Autoencoder (AE) is a deterministic neural network architecture that

- Maps input  $x$  into a lower dimensional representation  $z = f_e(x)$
- Attempts to reconstruct the original input by  $f_d(z)$ , i.e.,  $x \approx \hat{x} = f_d(f_e(x))$

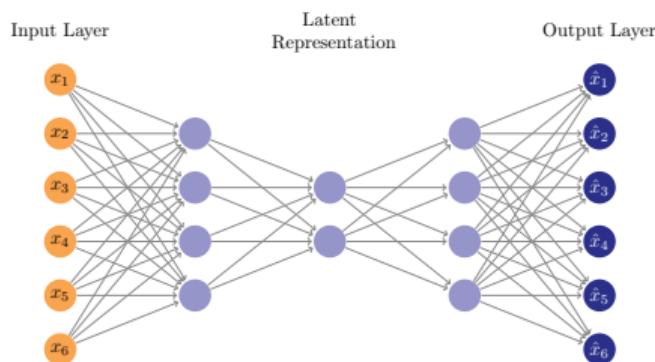
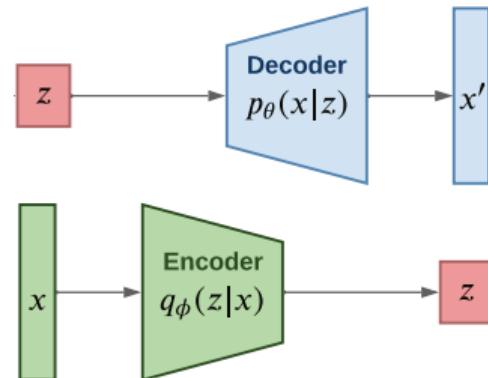


Figure 16.13 from (Murphy, 2023)

Variational autoencoder (VAE) is a neural architecture that consists of a deep generative model that is trained using an inference model



## Autoencoder vs. variational autoencoder

Autoencoder (AE) is a deterministic neural network architecture that

- Maps input  $x$  into a lower dimensional representation  $z = f_e(x)$
- Attempts to reconstruct the original input by  $f_d(z)$ , i.e.,  $x \approx \hat{x} = f_d(f_e(x))$

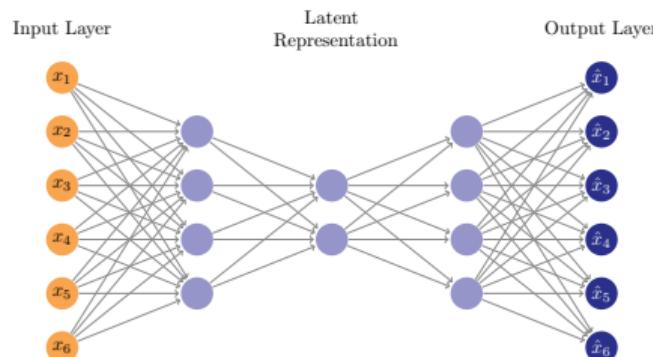
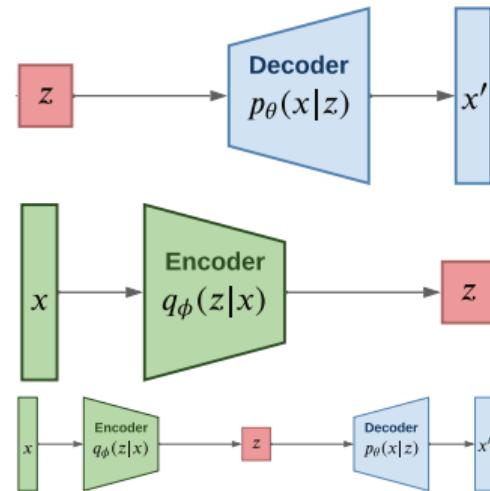


Figure 16.13 from (Murphy, 2023)

Variational autoencoder (VAE) is a neural architecture that consists of a deep generative model that is trained using an inference model



Adapted from Figure 20.1 from (Murphy, 2023)

# Applications of variational autoencoders

Variational autoencoders can be used for various tasks

- Statistical modeling<sup>a</sup>
- Generation<sup>a</sup>, conditional generation<sup>b</sup>
- Amortized variational inference<sup>a</sup>
- Representation learning<sup>a</sup>
- Reconstruction, missing value estimation<sup>b</sup>
- Latent space interpolation<sup>b</sup>
- Out-of-distribution detection<sup>b</sup>

---

<sup>a</sup>In this lecture.

<sup>b</sup>In the next lecture.

## Variational autoencoder: MNIST image example

MNIST is a large image dataset  $\mathcal{D} = \mathbf{x}_{1:N}$  of handwritten digits



[https://riccardo-cantini.netlify.app/post/cnn\\_vae\\_mnist/](https://riccardo-cantini.netlify.app/post/cnn_vae_mnist/)

## Variational autoencoder: MNIST image example

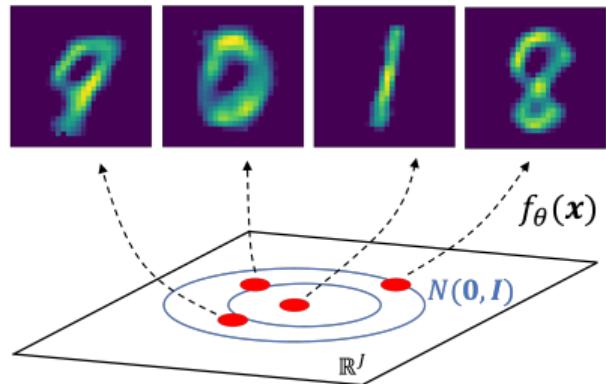
MNIST is a large image dataset  $\mathcal{D} = \mathbf{x}_{1:N}$  of handwritten digits



[https://riccardo-cantini.netlify.app/post/cnn\\_vae\\_mnist/](https://riccardo-cantini.netlify.app/post/cnn_vae_mnist/)

Train a VAE model and sample from the trained generative model to generate new images

$$\begin{aligned} z^* &\sim p_{\theta}(z) \\ x^* \mid z^* &\sim p(x \mid d_{\theta}(z^*)) \end{aligned}$$



<https://mbernste.github.io/posts/vae/>

## Sample generation from VAE

As before, a new data point  $\mathbf{x}^*$  can be generated from a latent variable model using the ancestral sampling

$$\mathbf{z}^* \sim p_{\theta}(\mathbf{z})$$

$$\mathbf{x}^* \sim p_{\theta}(\mathbf{x} \mid \mathbf{z}^*)$$

For example in image tasks, so-called decoded mean samples can have higher quality

$$\mathbf{z}^* \sim p_{\theta}(\mathbf{z})$$

$$\hat{\mathbf{x}}^* = d_{\theta}(\mathbf{z}^*)$$

## Sample generation from VAE

As before, a new data point  $\mathbf{x}^*$  can be generated from a latent variable model using the ancestral sampling

$$\mathbf{z}^* \sim p_{\theta}(\mathbf{z})$$

$$\mathbf{x}^* \sim p_{\theta}(\mathbf{x} \mid \mathbf{z}^*)$$

For example in image tasks, so-called decoded mean samples can have higher quality

$$\mathbf{z}^* \sim p_{\theta}(\mathbf{z})$$

$$\hat{\mathbf{x}}^* = d_{\theta}(\mathbf{z}^*)$$

Another trick is to define a joint inference distribution

$$q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) = p_{\mathcal{D}}(\mathbf{x})q_{\phi}(\mathbf{z} \mid \mathbf{x}),$$

where

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n),$$

and then sample latent variables from the aggregated posterior (marginal inference distribution)  $q_{\mathcal{D}, \phi}(\mathbf{z})$

$$\begin{aligned} q_{\mathcal{D}, \phi}(\mathbf{z}) &= \int q_{\mathcal{D}, \phi}(\mathbf{x}, \mathbf{z}) d\mathbf{x} \\ &= \int \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n) q_{\phi}(\mathbf{z} \mid \mathbf{x}) d\mathbf{x} \\ &= \frac{1}{N} \sum_{n=1}^N q_{\phi}(\mathbf{z} \mid \mathbf{x}_n) \end{aligned}$$

## Variational autoencoder: CELEBA image example

Samples from VAEs that are trained on CELEBA dataset

a) Random samples:

$$z^* \sim p_{\theta}(z)$$

$$x^* | z^* \sim p_{\theta}(x | d_{\theta}(z^*))$$

b) Decoded mean samples:

$$z^* \sim p_{\theta}(z)$$

$$\hat{x}^* | z^* = d_{\theta}(z^*)$$

c) "Noise":  $\hat{\epsilon} = x^* - \hat{x}^*$

d) High-quality generation from a VAE with hierarchical prior and specialized architectures (details later)

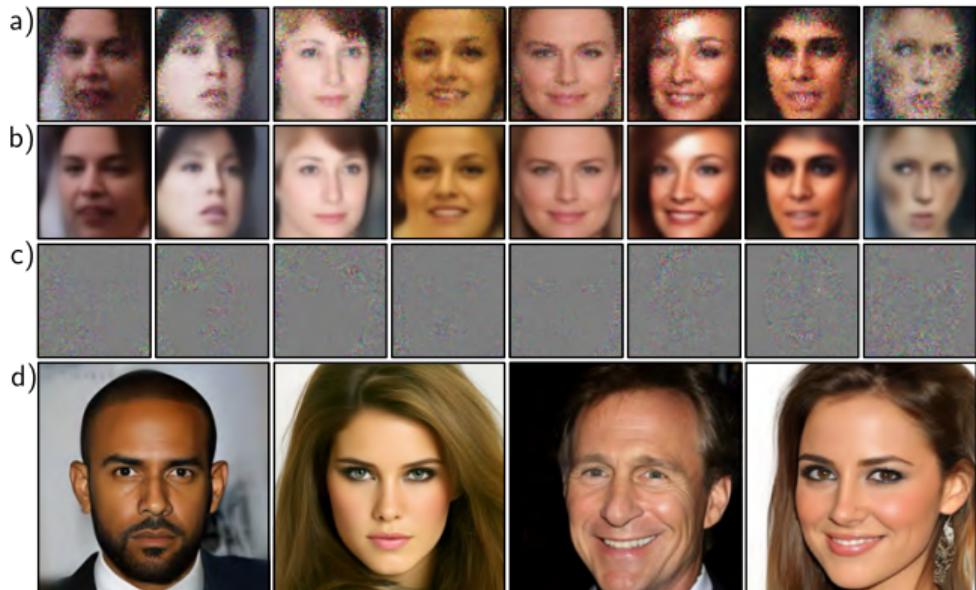


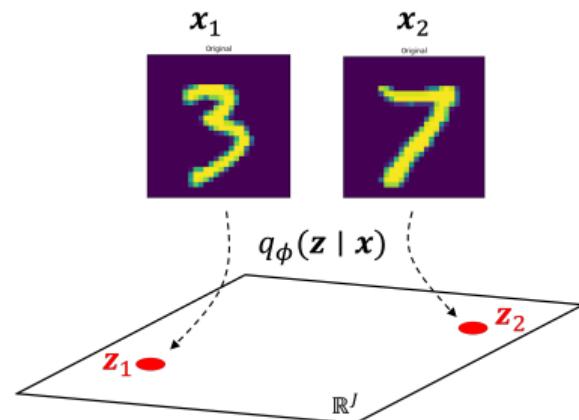
Figure 17.12 from (Prince, 2023).

# Variational autoencoder: dimension reduction, visualization, clustering

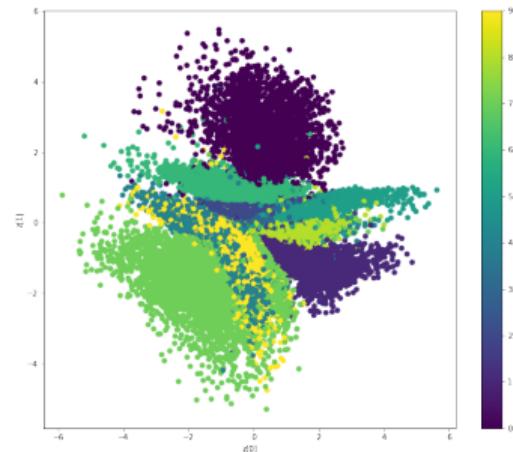
Example: Train a VAE model on the MNIST images

Use the encoder (amortized VI) for dimension reduction and visualization by plotting the approximative posteriors  $q_\phi(z | x)$  (or their means)

It is common to use the learned representations  $q_\phi(z | x)$  for clustering or other downstream tasks



<https://mbernste.github.io/posts/vae/>



[https://riccardo-cantini.netlify.app/post/cnn\\_vae\\_mnist/](https://riccardo-cantini.netlify.app/post/cnn_vae_mnist/)

# Variational autoencoder: modern molecular biology example

Single-cell sequencing is a modern technology that allows quantifying transcriptomic profiles of individual cells

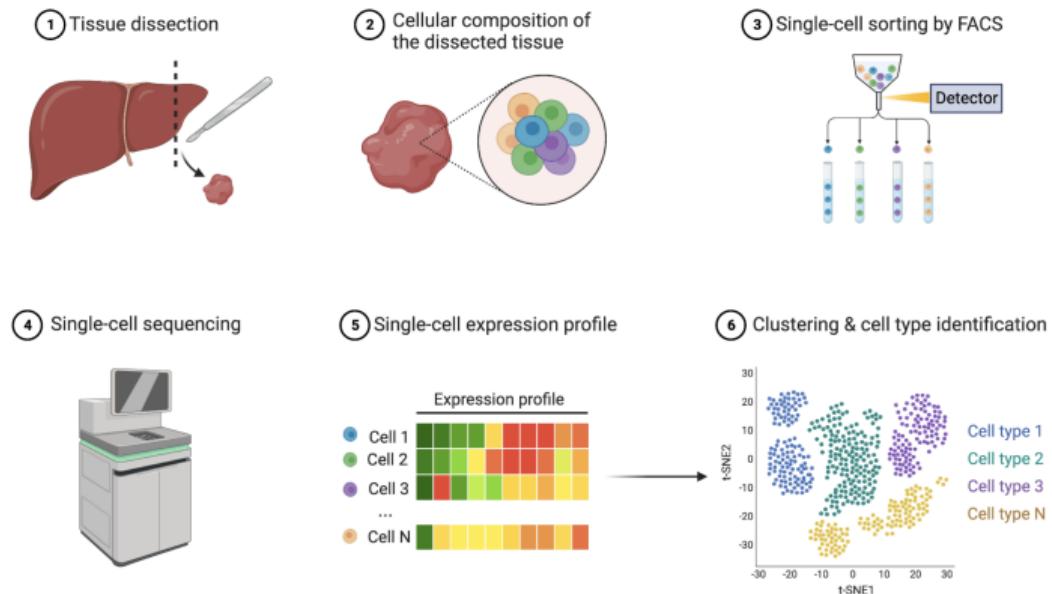


Figure from bioRender (<https://www.biorender.com/template/single-cell-sequencing>)

## Variational autoencoder: modern molecular biology example

Challenging to specify a well-motivated probabilistic model for the data

→ Incorporate biological inductive bias into neural network parameterized generative model

$$z_n \sim \text{Normal}(0, I)$$

$$\ell_n \sim \text{LogNormal}(\ell_\mu, \ell_\sigma^2)$$

$$\rho_n = f_w(z_n, s_n)$$

$$w_{ng} \sim \text{Gamma}(\rho_n^g, \theta)$$

$$y_{ng} \sim \text{Poisson}(\ell_n w_{ng})$$

$$h_{ng} \sim \text{Bernoulli}(f_h^g(z_n, s_n))$$

$$x_{ng} = \begin{cases} y_{ng} & \text{if } h_{ng} = 0, \\ 0 & \text{otherwise.} \end{cases}$$

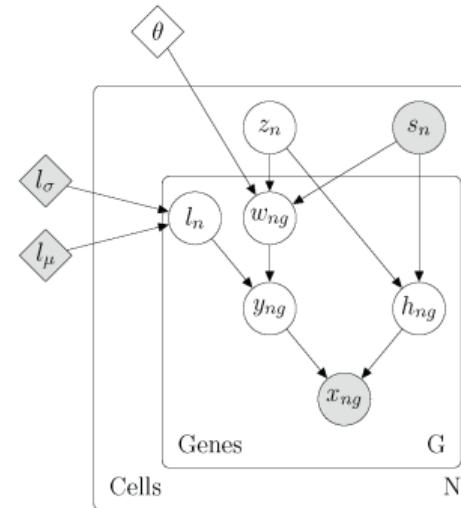


Figure from (Lopez et al, 2019)

# Variational autoencoder: modern molecular biology example

Variational autoencoder architecture (generative model and amortized VI) with deep neural networks

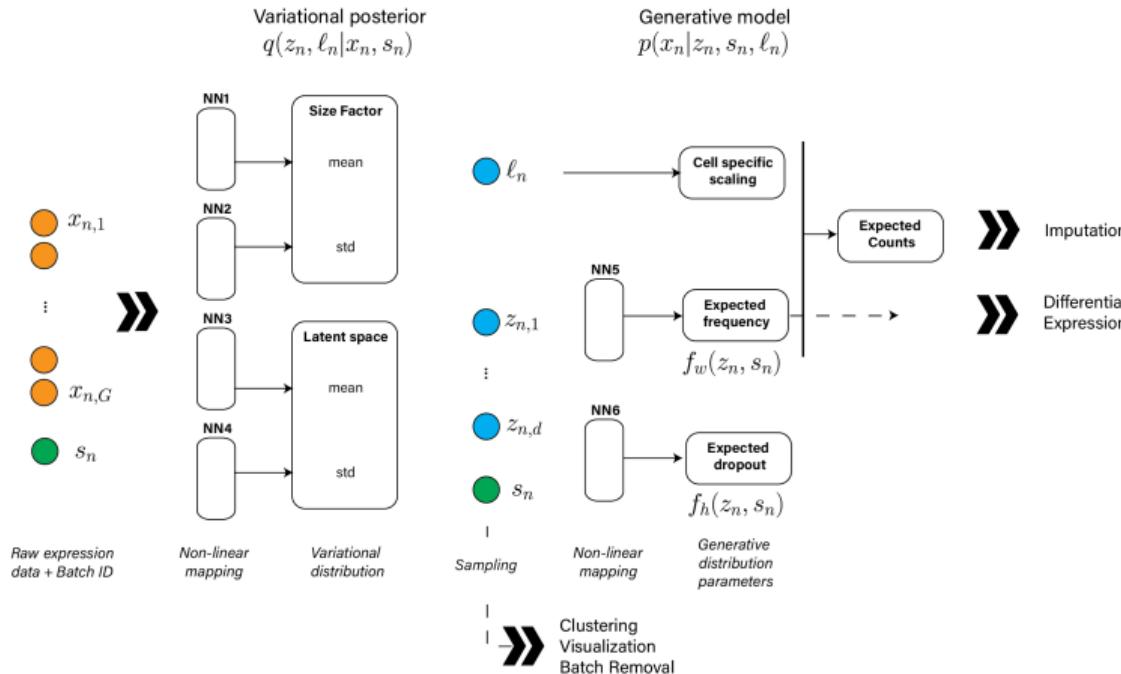


Figure from (Lopez et al, 2019)

# Variational autoencoder: modern molecular biology example

Encoder (amortized VI) provides representations for visualization, clustering, and other downstream tasks

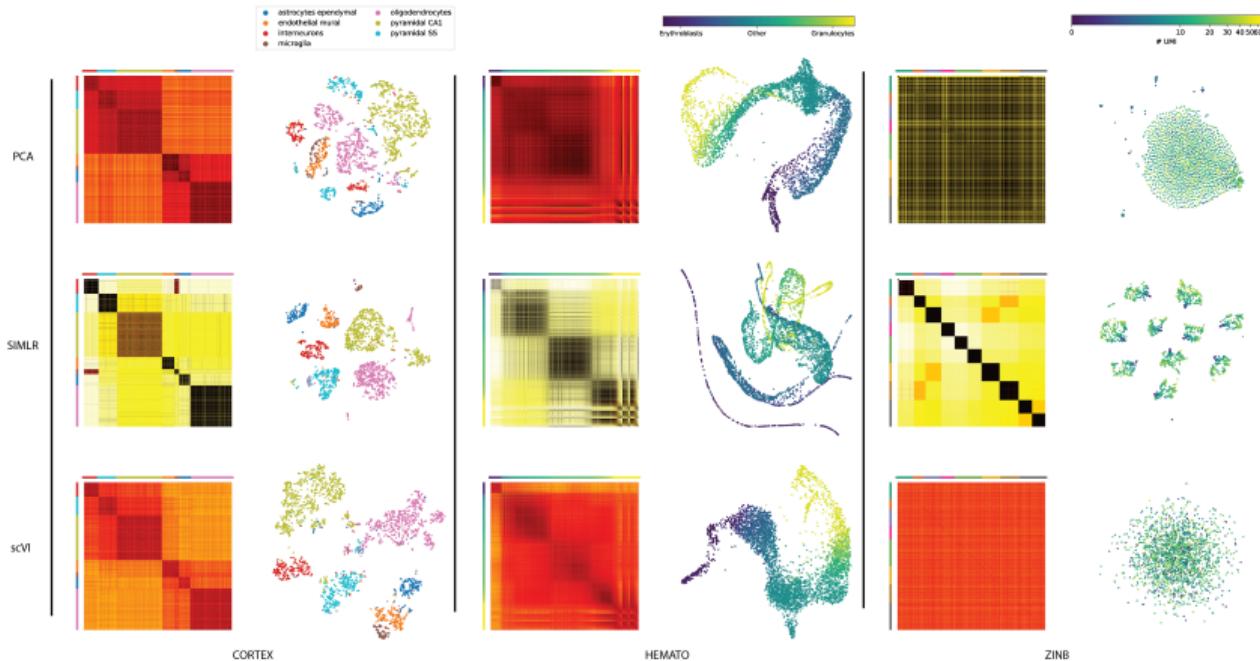


Figure from (Lopez et al, 2019)

## References

- Kingma DP and Welling M, An Introduction to Variational Autoencoders, Foundations and Trends in Machine Learning, Vol. 12, No. 4, pp. 307-392.
- R. Lopez, J. Regier, MB. Cole, M. Jordan, N. Yosef, Deep generative modeling for single-cell transcriptomics, *Nature Methods*, 2019.
- Murphy K, Probabilistic Machine Learning: Advanced Topics, The MIT Press, 2023.
- Prince SJD, Understanding Deep Learning, The MIT Press, 2023.