

CS-E4891 Deep Generative Models

Lecture 3: Variational autoencoders

Harri Lähdesmäki

Department of Computer Science
Aalto University

March 31, 2025

Outline

- Variational autoencoders (recap)
- Applications of VAE (continued)
- VAE generalizations
- Identifiability of VAEs
- Amortization gap
- Reading: parts of Sec. 20-21 from (Murphy, 2023), Sec. 17 from (Prince, 2023), and other literature (see references at the end)

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

- 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

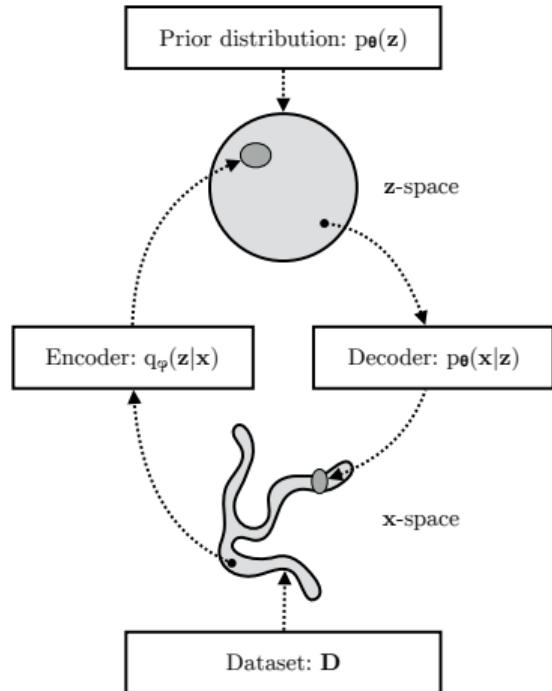


Figure from (Kingma & Welling, 2019)

Applications of variational autoencoders

Variational autoencoders can be used for various tasks

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

^aIn the previous lecture.

^bIn this lecture.

Variational autoencoder for reconstruction and data imputation

Consider a test data point $\mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$, where \mathbf{x}_o denotes observed features and \mathbf{x}_m denotes unobserved (missing) features

Observed and missing features can vary between different test data points \mathbf{x} and \mathbf{x}'

For a VAE model \mathbf{x}_o and \mathbf{x}_m are conditionally independent given \mathbf{z} (assuming diagonal likelihood covariance), therefore

$$\begin{aligned} p(\mathbf{x}_m \mid \mathbf{x}_o) &= \int p(\mathbf{x}_m \mid \mathbf{z}) p(\mathbf{z} \mid \mathbf{x}_o) d\mathbf{z} \\ &\approx \int p_{\theta}(\mathbf{x}_m \mid \mathbf{z}) q_{\phi}(\mathbf{z} \mid \mathbf{x}_o) d\mathbf{z} \end{aligned}$$

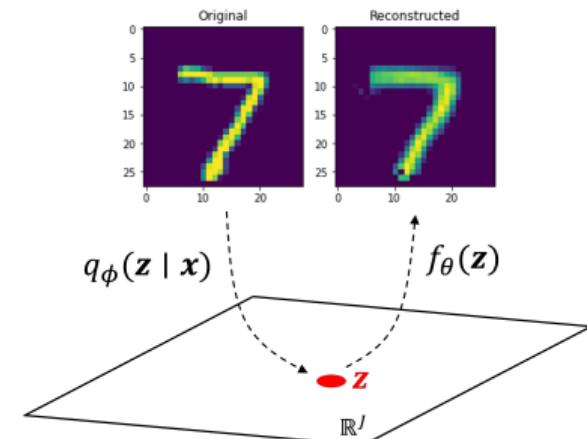
Variational autoencoder for reconstruction and data imputation

Consider a test data point $\mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$, where \mathbf{x}_o denotes observed features and \mathbf{x}_m denotes unobserved (missing) features

Observed and missing features can vary between different test data points \mathbf{x} and \mathbf{x}'

For a VAE model \mathbf{x}_o and \mathbf{x}_m are conditionally independent given \mathbf{z} (assuming diagonal likelihood covariance), therefore

$$\begin{aligned} p(\mathbf{x}_m | \mathbf{x}_o) &= \int p(\mathbf{x}_m | \mathbf{z}) p(\mathbf{z} | \mathbf{x}_o) d\mathbf{z} \\ &\approx \int p_{\theta}(\mathbf{x}_m | \mathbf{z}) q_{\phi}(\mathbf{z} | \mathbf{x}_o) d\mathbf{z} \end{aligned}$$



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

→ VAE architecture automatically implements a statistical data imputation method

Encoder model $q_{\phi}(\mathbf{z} | \mathbf{x})$ needs to be able to handle missing values in \mathbf{x}

Variational autoencoder for reconstruction and data imputation

Posterior matching technique (Strauss and Oliva, 2022)

- Artificially mask complete training data points as
 $\mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$
- Train an additional encoder model with parameters
 ϕ' such that $q_{\phi'}(\mathbf{z} \mid \mathbf{x}_o) \approx q_{\phi}(\mathbf{z} \mid \mathbf{x})$ by maximizing

$$\mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [\log q_{\phi'}(\mathbf{z} \mid \mathbf{x}_o)]$$

Variational autoencoder for reconstruction and data imputation

Posterior matching technique (Strauss and Oliva, 2022)

- Artificially mask complete training data points as
 $\mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$
- Train an additional encoder model with parameters
 ϕ' such that $q_{\phi'}(\mathbf{z} | \mathbf{x}_o) \approx q_{\phi}(\mathbf{z} | \mathbf{x})$ by maximizing

$$\mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log q_{\phi'}(\mathbf{z} | \mathbf{x}_o)]$$

Variationally optimal approach: use non-amortized variational inference to maximize (Agarwal et al., 2023)

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\psi} | \mathbf{x}) &= \mathbb{E}_{q_{\psi}(\mathbf{z})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_o | \mathbf{z})] \\ &\quad - D_{\text{KL}}(q_{\psi}(\mathbf{z}) || p_{\boldsymbol{\theta}}(\mathbf{z}))\end{aligned}$$

w.r.t. $\boldsymbol{\psi}$ and estimate with previously trained decoder

$$p(\mathbf{x}_m | \mathbf{x}_o) \approx \int p_{\boldsymbol{\theta}}(\mathbf{x}_m | \mathbf{z}) q_{\psi}(\mathbf{z}) d\mathbf{z}$$

Variational autoencoder for reconstruction and data imputation

Posterior matching technique (Strauss and Oliva, 2022)

- Artificially mask complete training data points as $\mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$
- Train an additional encoder model with parameters ϕ' such that $q_{\phi'}(\mathbf{z} | \mathbf{x}_o) \approx q_{\phi}(\mathbf{z} | \mathbf{x})$ by maximizing

$$\mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} [\log q_{\phi'}(\mathbf{z} | \mathbf{x}_o)]$$

Variationally optimal approach: use non-amortized variational inference to maximize (Agarwal et al., 2023)

$$\begin{aligned}\mathcal{L}(\theta, \psi | \mathbf{x}) &= \mathbb{E}_{q_{\psi}(\mathbf{z})} [\log p_{\theta}(\mathbf{x}_o | \mathbf{z})] \\ &\quad - D_{\text{KL}}(q_{\psi}(\mathbf{z}) || p_{\theta}(\mathbf{z}))\end{aligned}$$

w.r.t. ψ and estimate with previously trained decoder

$$p(\mathbf{x}_m | \mathbf{x}_o) \approx \int p_{\theta}(\mathbf{x}_m | \mathbf{z}) q_{\psi}(\mathbf{z}) d\mathbf{z}$$

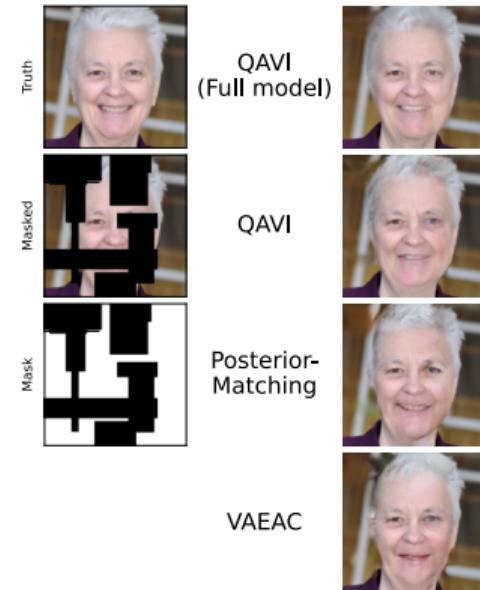


Figure from (Agarwal et al., 2023).

Note: above results use additional tricks, such as deep/hierarchical VAEs, more expressive posteriors

Variational autoencoder for reconstruction and data imputation

Some VAE variants can explicitly model missingness together with x , e.g., by also modeling missingness mask m

Can model

- missing-completely-at-random (MCAR)
- missing-at-random (MAR)
- missing-not-at-random (MNAR)

Latent space interpolation

Explore the generative model via latent space interpolation

Consider two data points \mathbf{x}_1 and \mathbf{x}_2 and the mean of their embeddings $\mathbf{z}_1 = f_\phi(\mathbf{x}_1)$ and $\mathbf{z}_2 = f_\phi(\mathbf{x}_2)$

Generate new images via latent space interpolation

$$\mathbf{z} = \lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2, \quad 0 \leq \lambda \leq 1$$

and decoding $\mathbf{x} \sim p_\theta(\mathbf{x} | \mathbf{z})$

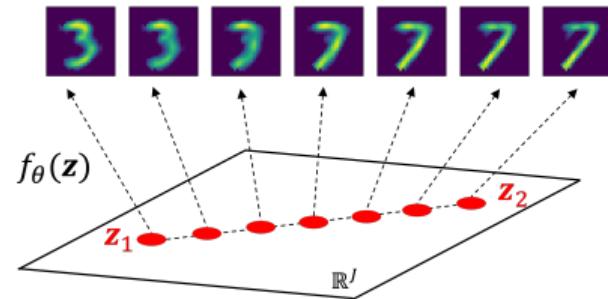


Figure from <https://mbernste.github.io/posts/vae/>

Note that the notation differs.

Latent space interpolation

Explore the generative model via latent space interpolation

Consider two data points \mathbf{x}_1 and \mathbf{x}_2 and the mean of their embeddings $\mathbf{z}_1 = f_\phi(\mathbf{x}_1)$ and $\mathbf{z}_2 = f_\phi(\mathbf{x}_2)$

Generate new images via latent space interpolation

$$\mathbf{z} = \lambda \mathbf{z}_1 + (1 - \lambda) \mathbf{z}_2, \quad 0 \leq \lambda \leq 1$$

and decoding $\mathbf{x} \sim p_\theta(\mathbf{x} | \mathbf{z})$

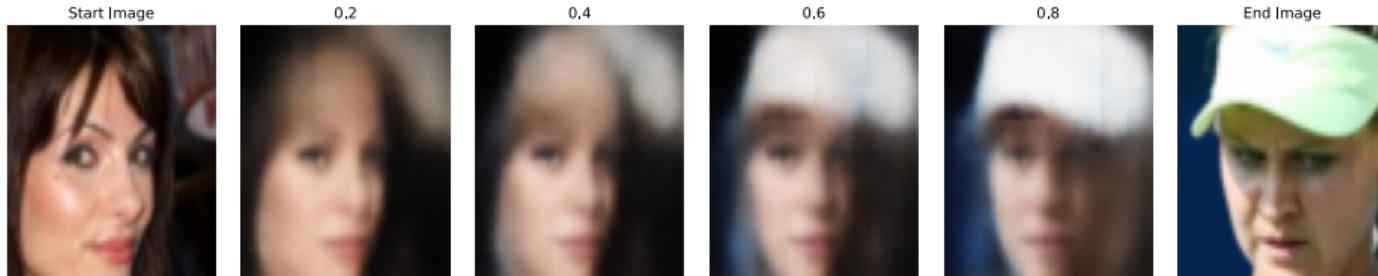


Figure 20.9 from (Murphy, 2023)

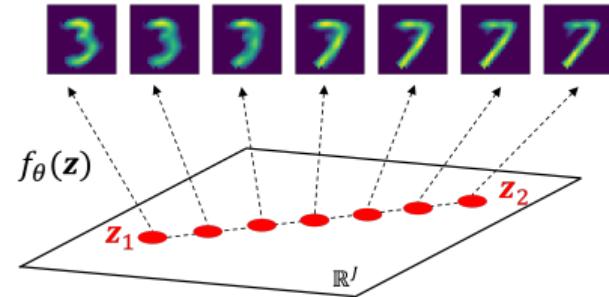


Figure from <https://mbernste.github.io/posts/vae/>

Note that the notation differs.

Latent space interpolation (and extrapolation)

Example: consider groups of images labeled as

- “neutral”: $\{x_i : i \in \mathcal{I}_{\text{neutral}}\}$
- “smiling”: $\{x_j : j \in \mathcal{I}_{\text{smiling}}\}$
- “mouth closed”: $\{x_k : k \in \mathcal{I}_{\text{closed}}\}$
- “mouth open”: $\{x_l : l \in \mathcal{I}_{\text{open}}\}$

and define

$$\Delta_1 = \sum_{j \in \mathcal{I}_{\text{smiling}}} f_\phi(x_j) - \sum_{i \in \mathcal{I}_{\text{neutral}}} f_\phi(x_i)$$

$$\Delta_2 = \sum_{l \in \mathcal{I}_{\text{open}}} f_\phi(x_l) - \sum_{k \in \mathcal{I}_{\text{closed}}} f_\phi(x_k)$$

A given image x can be manipulated via interpolation as

$$z = f_\phi(x)$$

$$z' = z + \lambda_1 \Delta_1 + \lambda_2 \Delta_2, \quad \lambda_1, \lambda_2 \in \mathbb{R}$$

$$x' \sim p_\theta(x | z')$$



Figure 17.13 from (Prince, 2023)

Note: interpolation above uses non-linear interpolation in spherical coordinates.

Approximating sample probability

Use a trained VAE to approximate sample probability for test data to

- Estimate the quality of the model
- Detect out-of-distribution samples

Approximating sample probability

Use a trained VAE to approximate sample probability for test data to

- Estimate the quality of the model
- Detect out-of-distribution samples

Instead of using the ELBO, we can directly estimate probability for a given data point \mathbf{x} using Monte Carlo

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \int p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{p_{\theta}(\mathbf{z})}[p_{\theta}(\mathbf{x} | \mathbf{z})] \\ &\approx \frac{1}{S} \sum_{s=1}^S p_{\theta}(\mathbf{x} | \mathbf{z}_s) \end{aligned}$$

where

$$\mathbf{z}_s \stackrel{\text{i.i.d.}}{\sim} p_{\theta}(\mathbf{z})$$

Approximating sample probability

Use a trained VAE to approximate sample probability for test data to

- Estimate the quality of the model
- Detect out-of-distribution samples

Instead of using the ELBO, we can directly estimate probability for a given data point \mathbf{x} using Monte Carlo

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \int p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_{p_{\theta}(\mathbf{z})}[p_{\theta}(\mathbf{x} | \mathbf{z})] \\ &\approx \frac{1}{S} \sum_{s=1}^S p_{\theta}(\mathbf{x} | \mathbf{z}_s) \end{aligned}$$

where

$$\mathbf{z}_s \stackrel{\text{i.i.d}}{\sim} p_{\theta}(\mathbf{z})$$

$$\begin{aligned} p_{\theta}(\mathbf{x}) &= \int p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z} \\ &= \int p_{\theta}(\mathbf{x} | \mathbf{z}) \frac{p_{\theta}(\mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} q_{\phi}(\mathbf{z} | \mathbf{x}) d\mathbf{z} \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[p_{\theta}(\mathbf{x} | \mathbf{z}) \frac{p_{\theta}(\mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \\ &\approx \frac{1}{S} \sum_{s=1}^S \tilde{w}_s p_{\theta}(\mathbf{x} | \mathbf{z}_s), \end{aligned}$$

where

$$\mathbf{z}_s \stackrel{\text{i.i.d}}{\sim} q_{\phi}(\mathbf{z} | \mathbf{x}), \quad \tilde{w}_s = \frac{p_{\theta}(\mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s | \mathbf{x})}, \quad \text{sup } p_{\theta} \subseteq \text{sup } q_{\phi}$$

Posterior collapse

- Consider a powerful decoder (a “super” deep neural net of some sort) such that for every value of latent variable $z \in \mathcal{Z}$ decoder outputs one of the data points
- Then approx. $p_{\theta}(x | z) = p_{\theta}(x)$ and $p_{\theta}(z | x) = p_{\theta}(z)$ because generative model does not really use the latent variable
- We can find inference network parameters ϕ^* such that $q_{\phi^*}(z | x) = p_{\theta}(z | x) = p_{\theta}(z)$

Posterior collapse

- Consider a powerful decoder (a “super” deep neural net of some sort) such that for every value of latent variable $z \in \mathcal{Z}$ decoder outputs one of the data points
- Then approx. $p_{\theta}(x | z) = p_{\theta}(x)$ and $p_{\theta}(z | x) = p_{\theta}(z)$ because generative model does not really use the latent variable
- We can find inference network parameters ϕ^* such that $q_{\phi^*}(z | x) = p_{\theta}(z | x) = p_{\theta}(z)$
- In this setting, ELBO is maximized because $D_{\text{KL}}(q_{\phi^*}(z | x) || p_{\theta}(z | x)) = 0$ but the model has not learned any meaningful latent representations
- Then also $D_{\text{KL}}(q_{\phi^*}(z | x) || p_{\theta}(z)) = 0$
- This is known as posterior collapse
- Posterior collapse is also due to other aspects, such local maxima in the complex objective function

β -VAE

VAE models can generate samples that may appear too smooth (e.g. blurry images)

β -VAE can improve sample quality (as well as help with posterior collapse) by reducing the penalty on the KL term

$$\mathcal{L}_\beta(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z} \mid \mathbf{x})} [\log p_\theta(\mathbf{x} \mid \mathbf{z})] - \beta D_{\text{KL}}(q_\phi(\mathbf{z} \mid \mathbf{x}) \parallel p_\theta(\mathbf{z})), \quad \beta \geq 0$$

If we set

- $\beta = 1$ we obtain the standard objective for VAEs
- $\beta < 1$ we can reconstruct data samples better
- $\beta > 1$ we get more compressed representations (in terms of rate distortion, see Murphy (2023))
- $\beta = 0$ we obtain an objective that is analogous for autoencoder models

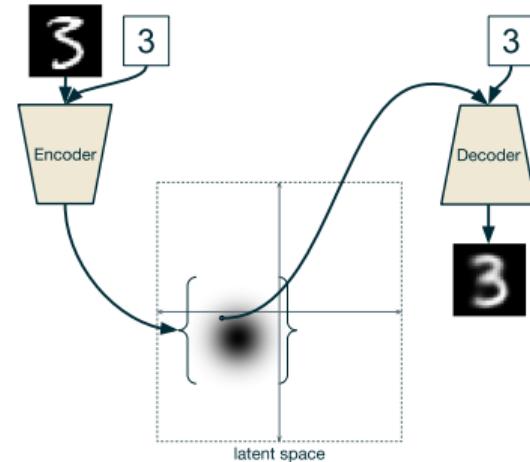
Conditional variational autoencoder

Conditional variational autoencoder (CVAE) is a neural architecture that consists of a conditional deep latent variable model of the form

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

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

where \mathbf{c} contains the explanatory covariates and the model is trained using (reparameterized) amortized variational inference with an inference model $q_{\phi}(z | \mathbf{x}, \mathbf{c})$.



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

Conditional variational autoencoder

Conditional variational autoencoder (CVAE) is a neural architecture that consists of a conditional deep latent variable model of the form

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

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

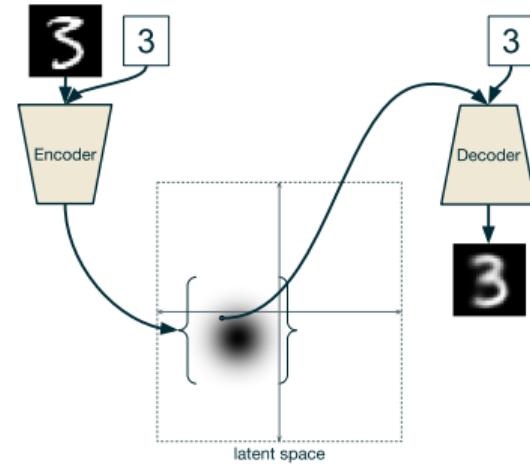
where \mathbf{c} contains the explanatory covariates and the model is trained using (reparameterized) amortized variational inference with an inference model $q_{\phi}(z | \mathbf{x}, \mathbf{c})$.

Alternative CVAE model definitions:

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

or

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



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

Conditional variational autoencoder

Conditional variational autoencoder (CVAE) is a neural architecture that consists of a conditional deep latent variable model of the form

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

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

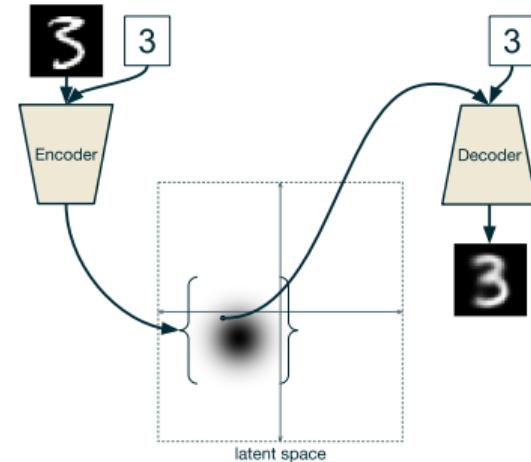
where c contains the explanatory covariates and the model is trained using (reparameterized) amortized variational inference with an inference model $q_{\phi}(z | \mathbf{x}, c)$.

Alternative CVAE model definitions:

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

or

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



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

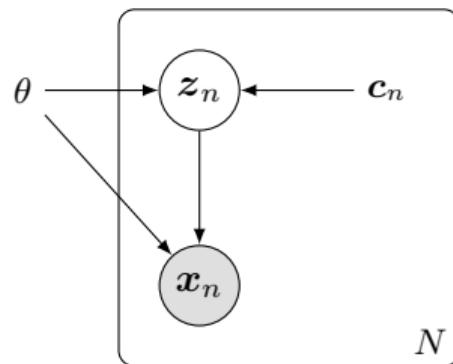
All methodological details remain essentially unchanged, including the ELBO objective

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

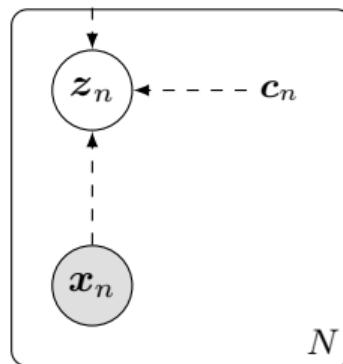
Conditional variational autoencoder: Probabilistic graphical model

Probabilistic graphical models for CVAE with $p_{\theta}(x, z | c) = p_{\theta}(x | z)p_{\theta}(z | c)$

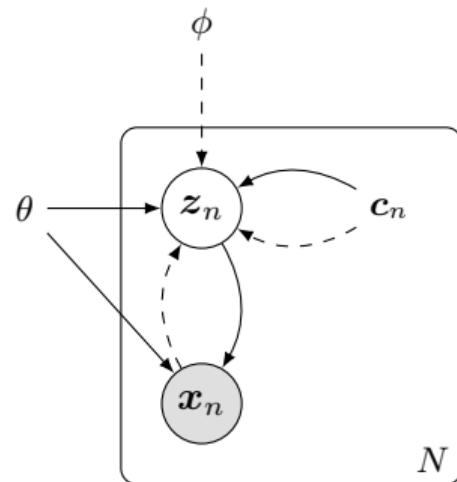
- Circled symbols denote random variables
- Non-circled symbols are hyperparameters
- White nodes (circles) are latent
- Shaded nodes are observed



(a) Generative model



(b) Inference model



(c) Combined

Multimodal variational autoencoder

- Multimodal VAE allows modeling joint distributions over multiple data modalities, e.g. image and text
- Assuming M modalities $\mathbf{x}_1 \in \mathcal{X}_1, \dots, \mathbf{x}_M \in \mathcal{X}_M$ and conditional independence given \mathbf{z}

$$p_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_M, \mathbf{z}) = p_{\theta}(\mathbf{z}) \prod_{m=1}^M p_{\theta}(\mathbf{x}_m | \mathbf{z})$$

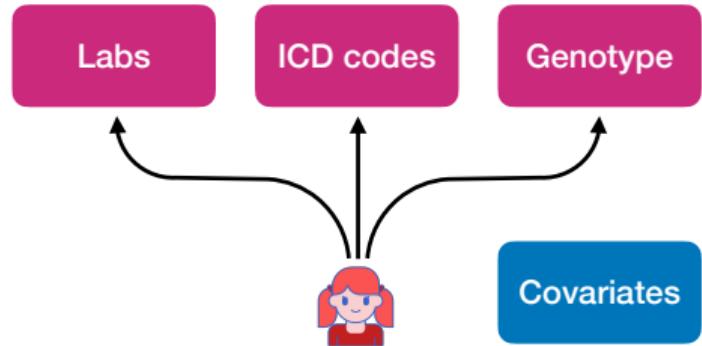


Illustration of multimodal data

Multimodal variational autoencoder

- Multimodal VAE allows modeling joint distributions over multiple data modalities, e.g. image and text
- Assuming M modalities $\mathbf{x}_1 \in \mathcal{X}_1, \dots, \mathbf{x}_M \in \mathcal{X}_M$ and conditional independence given \mathbf{z}

$$p_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_M, \mathbf{z}) = p_{\theta}(\mathbf{z}) \prod_{m=1}^M p_{\theta}(\mathbf{x}_m | \mathbf{z})$$

- The ELBO for $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is

$$\mathcal{L}(\theta, \phi | \mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{X})} \left[\sum_{m=1}^M \log p_{\theta}(\mathbf{x}_m | \mathbf{z}) \right] - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{X}) || p_{\theta}(\mathbf{z}))$$

- Likelihood terms may have different ranges
- Weighted ELBO with $\lambda_m \geq 0$

$$\mathcal{L}(\theta, \phi | \mathbf{X}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{X})} \left[\sum_{m=1}^M \lambda_m \log p_{\theta}(\mathbf{x}_m | \mathbf{z}) \right] - D_{\text{KL}}(q_{\phi}(\mathbf{z} | \mathbf{X}) || p_{\theta}(\mathbf{z}))$$

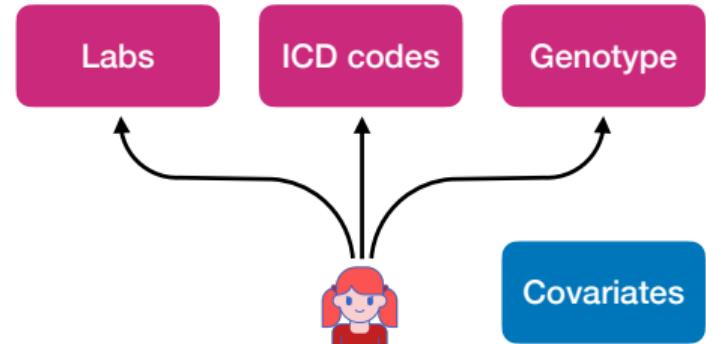


Illustration of multimodal data

Multimodal variational autoencoder for unpaired data

- Availability of paired multimodal data across all modalities is often limited
- For example, we may have lots of images, lots of text, but only few image-text pairs
- Let $O_m = 1$ denote that data modality m , \mathbf{x}_m , is observed, and $O_m = 0$ that modality m is not observed
- ELBO objective for partially paired data is

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{X}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{X})} \left[\sum_{m:O_m=1} \lambda_m \log p_{\boldsymbol{\theta}}(\mathbf{x}_m \mid \mathbf{z}) \right] - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{X}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}))$$

Multimodal variational autoencoder for unpaired data

- Availability of paired multimodal data across all modalities is often limited
- For example, we may have lots of images, lots of text, but only few image-text pairs
- Let $O_m = 1$ denote that data modality m , \mathbf{x}_m , is observed, and $O_m = 0$ that modality m is not observed
- ELBO objective for partially paired data is

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{X}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{X})} \left[\sum_{m:O_m=1} \lambda_m \log p_{\boldsymbol{\theta}}(\mathbf{x}_m \mid \mathbf{z}) \right] - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{X}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}))$$

- How to define the amortized VI such that it works for unpaired data?
- For example, if encoder $q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x}_1, \mathbf{x}_2)$ is trained on image-text pairs, how to compute posterior approximation using only images $q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x}_1)$ or text $q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x}_2)$, and vice versa?

Multimodal variational autoencoder with product of experts

From the conditional independence we get

$$\begin{aligned} p(\mathbf{z} \mid \mathbf{X}) &= \frac{p(\mathbf{z})p(\mathbf{x}_1, \dots, \mathbf{x}_m \mid \mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \\ &= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \prod_{m=1}^M p(\mathbf{x}_m \mid \mathbf{z}) \\ &= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \prod_{m=1}^M \frac{p(\mathbf{z} \mid \mathbf{x}_m)p(\mathbf{x}_m)}{p(\mathbf{z})} \\ &\propto p(\mathbf{z}) \prod_{m=1}^M \frac{p(\mathbf{z} \mid \mathbf{x}_m)}{p(\mathbf{z})} \\ &\approx p(\mathbf{z}) \prod_{m=1}^M q_m(\mathbf{z} \mid \mathbf{x}_m), \end{aligned}$$

where $q_m(\mathbf{z} \mid \mathbf{x}_m)$ is the “expert” for data modality m

Multimodal variational autoencoder with product of experts

From the conditional independence we get

$$\begin{aligned} p(\mathbf{z} \mid \mathbf{X}) &= \frac{p(\mathbf{z})p(\mathbf{x}_1, \dots, \mathbf{x}_m \mid \mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \\ &= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \prod_{m=1}^M p(\mathbf{x}_m \mid \mathbf{z}) \\ &= \frac{p(\mathbf{z})}{p(\mathbf{x}_1, \dots, \mathbf{x}_m)} \prod_{m=1}^M \frac{p(\mathbf{z} \mid \mathbf{x}_m)p(\mathbf{x}_m)}{p(\mathbf{z})} \\ &\propto p(\mathbf{z}) \prod_{m=1}^M \frac{p(\mathbf{z} \mid \mathbf{x}_m)}{p(\mathbf{z})} \\ &\approx p(\mathbf{z}) \prod_{m=1}^M q_m(\mathbf{z} \mid \mathbf{x}_m), \end{aligned}$$

where $q_m(\mathbf{z} \mid \mathbf{x}_m)$ is the “expert” for data modality m

This is known as the product of experts

Posterior approximation can be computed for any subset of modalities which are observed

If prior and experts are Gaussians

$$\begin{aligned} p(\mathbf{z}) &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0^{-1}) \\ q_m(\mathbf{z} \mid \mathbf{x}_m) &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1}), \end{aligned}$$

then the approximative posterior is also Gaussian

$$\begin{aligned} \prod_{m=0}^M \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m^{-1}) &\propto \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} = \left(\sum_{m=0}^M \boldsymbol{\Lambda}_m \right)^{-1} \quad \text{and} \quad \boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\sum_{m=0}^M \boldsymbol{\Lambda}_m \boldsymbol{\mu}_m \right) \end{aligned}$$

Multimodal variational autoencoder with product of experts: illustration

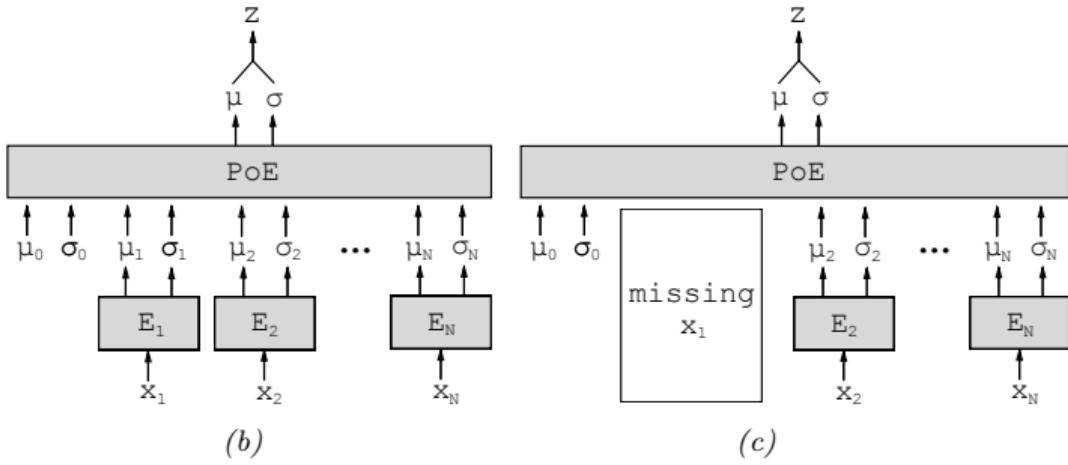
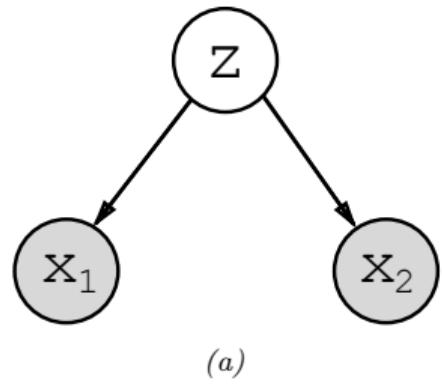


Figure 21.5 from (Murphy, 2023)

Semi-supervised VAEs

Semi-supervised learning methods can be useful when data consists of large amounts of unlabeled data $\mathcal{D}_U = \{(\mathbf{x}_n)\}$ and few labeled data points $\mathcal{D}_L = \{(\mathbf{x}_n, y_n)\}$

One of the earliest models, called M2, is

$$p_{\theta}(\mathbf{x}, y) = p_{\theta}(y)p_{\theta}(\mathbf{x} \mid y) = p_{\theta}(y) \int p_{\theta}(\mathbf{x} \mid y, \mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z},$$

where $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$, $p_{\theta}(y) = \text{Cat}(y \mid \boldsymbol{\pi})$ is trainable label prior, and $p_{\theta}(\mathbf{x} \mid y, \mathbf{z}) = p_{\theta}(\mathbf{x} \mid d_{\theta}(y, \mathbf{z}))$

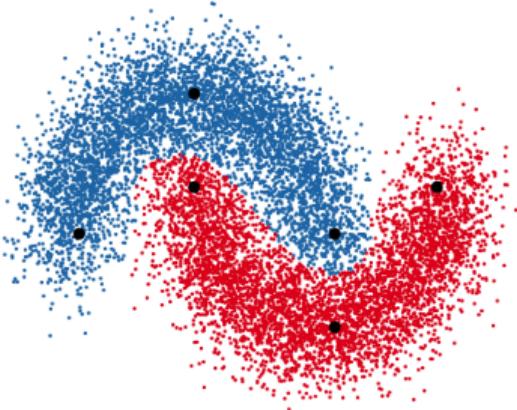


Figure from (Maaloe et al, 2016)

Comparison to CVAE:

- c is a fixed covariate/predictor
- here label y is a random variable

Semi-supervised VAEs

Variational approximation for labeled data

$$\begin{aligned} q_{\phi}(\mathbf{z} \mid \mathbf{x}, y) &= q(\mathbf{z} \mid f_{\phi}(\mathbf{x}, y)) \\ &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_{\phi}(\mathbf{x}, y), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x}, y))) \end{aligned}$$

Standard ELBO for labeled data (by marginalizing \mathbf{z})

$$\begin{aligned} \log p_{\theta}(\mathbf{x}, y) &\geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} \left[\log \frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} [\log p_{\theta}(\mathbf{x} \mid y, \mathbf{z}) + \log p_{\theta}(y) \\ &\quad + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)] \\ &= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}, y) \end{aligned}$$

Semi-supervised VAEs

Variational approximation for labeled data

$$\begin{aligned} q_{\phi}(\mathbf{z} \mid \mathbf{x}, y) &= q(\mathbf{z} \mid f_{\phi}(\mathbf{x}, y)) \\ &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_{\phi}(\mathbf{x}, y), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x}, y))) \end{aligned}$$

Standard ELBO for labeled data (by marginalizing \mathbf{z})

$$\begin{aligned} \log p_{\theta}(\mathbf{x}, y) &\geq \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} \left[\log \frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)} [\log p_{\theta}(\mathbf{x} \mid y, \mathbf{z}) + \log p_{\theta}(y) \\ &\quad + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z} \mid \mathbf{x}, y)] \\ &= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}, y) \end{aligned}$$

Variational approximation for unlabeled data

$$\begin{aligned} q_{\phi}(\mathbf{z}, y \mid \mathbf{x}) &= q_{\phi}(\mathbf{z} \mid \mathbf{x}) q_{\phi}(y \mid \mathbf{x}) \\ q_{\phi}(\mathbf{z} \mid \mathbf{x}) &= \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}_{\phi}(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_{\phi}(\mathbf{x}))) \\ q_{\phi}(y \mid \mathbf{x}) &= \text{Cat}(y \mid \boldsymbol{\pi}_{\phi}(\mathbf{x})) \end{aligned}$$

ELBO for unlabeled data by (marginalizing \mathbf{z} and y)

$$\begin{aligned} \log p_{\theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\phi}(\mathbf{z}, y \mid \mathbf{x})} \left[\log \frac{p_{\theta}(\mathbf{x}, y, \mathbf{z})}{q_{\phi}(\mathbf{z}, y \mid \mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}, y \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x} \mid y, \mathbf{z}) + \log p_{\theta}(y) \\ &\quad + \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}, y \mid \mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}(y \mid \mathbf{x})} [\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}, y)] + \mathbb{H}(q_{\phi}(y \mid \mathbf{x})) \\ &= \mathcal{U}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) \end{aligned}$$

Semi-supervised VAEs

Variational approximation $q_\phi(y \mid \mathbf{x})$ can be considered as discriminative classifier

Because classifier is only used to compute log-likelihood of unlabeled data, a modified objective

$$\begin{aligned}\mathcal{L}'(\theta, \phi \mid \mathcal{D}_L, \mathcal{D}_U) = & \mathcal{L}(\theta, \phi \mid \mathcal{D}_L) + \mathcal{U}(\theta, \phi \mid \mathcal{D}_U) \\ & + \mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{D}_L} \log q_\phi(y \mid \mathbf{x}),\end{aligned}$$

where $(\mathbf{x}_n, y_n) \sim \mathcal{D}_L$ denotes labeled data distribution

N	NN	CNN	TSVM	CAE	MTC	AtlasRBF	M1+TSVM	M2	M1+M2
100	25.81	22.98	16.81	13.47	12.03	8.10 (± 0.95)	11.82 (± 0.25)	11.97 (± 1.71)	3.33 (± 0.14)
600	11.44	7.68	6.16	6.3	5.13	—	5.72 (± 0.049)	4.94 (± 0.13)	2.59 (± 0.05)
1000	10.7	6.45	5.38	4.77	3.64	3.68 (± 0.12)	4.24 (± 0.07)	3.60 (± 0.56)	2.40 (± 0.02)
3000	6.04	3.35	3.45	3.22	2.57	—	3.49 (± 0.04)	3.92 (± 0.63)	2.18 (± 0.04)

Figure from (Kingma et al, 2014)

Semi-supervised VAEs

Variational approximation $q_\phi(y | \mathbf{x})$ can be considered as discriminative classifier

Because classifier is only used to compute log-likelihood of unlabeled data, a modified objective

$$\begin{aligned}\mathcal{L}'(\theta, \phi | \mathcal{D}_L, \mathcal{D}_U) = & \mathcal{L}(\theta, \phi | \mathcal{D}_L) + \mathcal{U}(\theta, \phi | \mathcal{D}_U) \\ & + \mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{D}_L} \log q_\phi(y | \mathbf{x}),\end{aligned}$$

where $(\mathbf{x}_n, y_n) \sim \mathcal{D}_L$ denotes labeled data distribution

Methods have been further developed e.g. by adding additional latent variables (e.g. Maaloe et al, 2016) among others

N	NN	CNN	TSVM	CAE	MTC	AtlasRBF	M1+TSVM	M2	M1+M2
100	25.81	22.98	16.81	13.47	12.03	8.10 (± 0.95)	11.82 (± 0.25)	11.97 (± 1.71)	3.33 (± 0.14)
600	11.44	7.68	6.16	6.3	5.13	—	5.72 (± 0.049)	4.94 (± 0.13)	2.59 (± 0.05)
1000	10.7	6.45	5.38	4.77	3.64	3.68 (± 0.12)	4.24 (± 0.07)	3.60 (± 0.56)	2.40 (± 0.02)
3000	6.04	3.35	3.45	3.22	2.57	—	3.49 (± 0.04)	3.92 (± 0.63)	2.18 (± 0.04)

Figure from (Kingma et al, 2014)

Hierarchical VAEs

We can define hierarchical VAEs with K stochastic layers

$\mathbf{z}_{1:K} = (\mathbf{z}_1, \dots, \mathbf{z}_K)$ as follows

$$p_{\theta}(\mathbf{x}, \mathbf{z}_{1:K}) = p_{\theta}(\mathbf{z}_K) \prod_{k=1}^K p_{\theta}(\mathbf{z}_k | \mathbf{z}_{k+1}) p_{\theta}(\mathbf{x} | \mathbf{z}_1)$$

Generative model can be extended into non-Markovian model

$$p_{\theta}(\mathbf{x}, \mathbf{z}_{1:K}) = p_{\theta}(\mathbf{z}_K) \prod_{k=1}^K p_{\theta}(\mathbf{z}_k | \mathbf{z}_{k+1:K}) p_{\theta}(\mathbf{x} | \mathbf{z}_{1:K})$$

Similar to using skip connection architecture

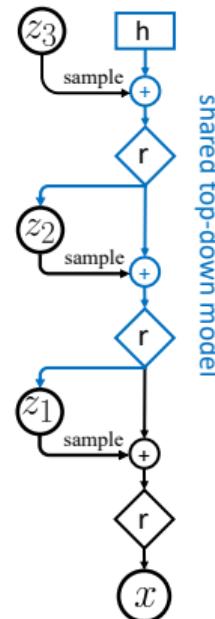


Figure 21.11 from (Murphy et al, 2023)

Hierarchical VAEs: inference model

Variational inference model can be define as bottom-up

$$q_{\phi}(\mathbf{z}_{1:K}) = q_{\phi}(\mathbf{z}_1 \mid \mathbf{x}) \prod_{k=2}^K q_{\phi}(\mathbf{z}_k \mid \mathbf{x}, \mathbf{z}_{1:k-1})$$

Hierarchical VAEs: inference model

Variational inference model can be defined as bottom-up

$$q_{\phi}(\mathbf{z}_{1:K}) = q_{\phi}(\mathbf{z}_1 \mid \mathbf{x}) \prod_{k=2}^K q_{\phi}(\mathbf{z}_k \mid \mathbf{x}, \mathbf{z}_{1:k-1})$$

or as top-down

$$q_{\phi}(\mathbf{z}_{1:K}) = q_{\phi}(\mathbf{z}_K \mid \mathbf{x}) \prod_{k=K-1}^1 q_{\phi}(\mathbf{z}_k \mid \mathbf{x}, \mathbf{z}_{k+1:K})$$

Hierarchical VAEs: inference model

Variational inference model can be defined as bottom-up

$$q_{\phi}(\mathbf{z}_{1:K}) = q_{\phi}(\mathbf{z}_1 | \mathbf{x}) \prod_{k=2}^K q_{\phi}(\mathbf{z}_k | \mathbf{x}, \mathbf{z}_{1:k-1})$$

or as top-down

$$q_{\phi}(\mathbf{z}_{1:K}) = q_{\phi}(\mathbf{z}_K | \mathbf{x}) \prod_{k=K-1}^1 q_{\phi}(\mathbf{z}_k | \mathbf{x}, \mathbf{z}_{k+1:K})$$

or both

Top-down inference model is generally considered better because it can be shown to more closely approximate the true posterior of a given layer

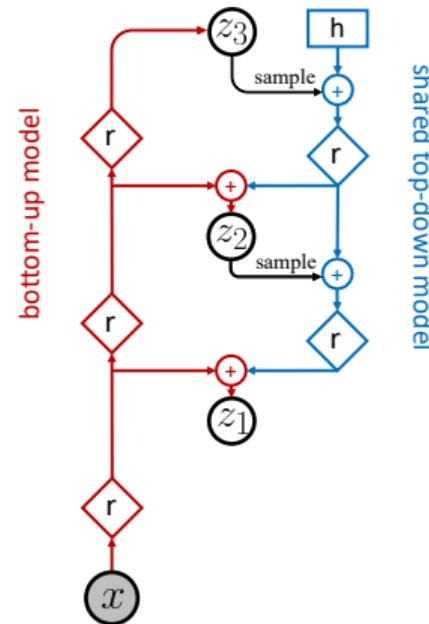


Figure 21.11 from (Murphy et al, 2023)

Hierarchical VAEs: ELBO

Evidence lower-bound for hierarchical VAE with top-down inference model is

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x} \mid z_{1:K})] - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}_K \mid \mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}_K)) \\ &\quad - \sum_{k=K-1}^1 \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{k+1:K} \mid \mathbf{x})} D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x}, \mathbf{z}_{k+1:K}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}_k \mid \mathbf{z}_{k+1:K})),\end{aligned}$$

where $q_{\boldsymbol{\phi}}(\mathbf{z}_{k+1:K} \mid \mathbf{x}) = \prod_{i=k+1}^K q_{\boldsymbol{\phi}}(\mathbf{z}_i \mid \mathbf{x}, \mathbf{z}_{i+1:K})$

Hierarchical VAE example

So-called very deep VAE (VDVAE) model uses hierarchical VAE with bidirectional CNN architecture and nearest neighbour upsampling

Low-resolution (top) latents capture global structure, while the bottom of the hierarchy fills in details

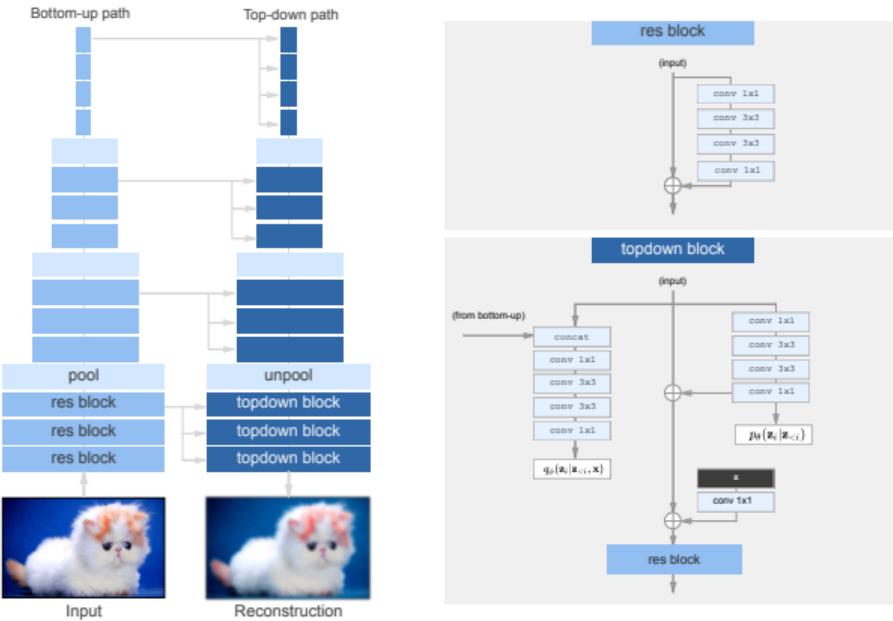


Figure 21.12 from (Murphy et al, 2023)

Hierarchical VAE example

So-called very deep VAE (VDVAE) model uses hierarchical VAE with bidirectional CNN architecture and nearest neighbour upsampling

Low-resolution (top) latents capture global structure, while the bottom of the hierarchy fills in details



Figure 21.13 from (Murphy et al, 2023)

Recap: Marginal likelihood, ELBO, importance sampling

Marginal likelihood and ELBO

$$\log p_{\theta}(\mathbf{x}) = \log \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right) \right] = \mathcal{L}(\theta, \phi | \mathbf{x})$$

and Monte Carlo estimate

$$\mathcal{L}(\theta, \phi | \mathbf{x}) \approx \frac{1}{S} \left(\sum_{s=1}^S \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s | \mathbf{x})} \right) \quad \text{where} \quad \mathbf{z}_s \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z} | \mathbf{x})$$

Recap: Marginal likelihood, ELBO, importance sampling

Marginal likelihood and ELBO

$$\log p_{\theta}(\mathbf{x}) = \log \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right) \right] = \mathcal{L}(\theta, \phi | \mathbf{x})$$

and Monte Carlo estimate

$$\mathcal{L}(\theta, \phi | \mathbf{x}) \approx \frac{1}{S} \left(\sum_{s=1}^S \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s | \mathbf{x})} \right) \quad \text{where} \quad \mathbf{z}_s \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z} | \mathbf{x})$$

Importance sampling estimate of $p_{\theta}(\mathbf{x})$ and $\log p_{\theta}(\mathbf{x})$

$$p_{\theta}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z} | \mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \approx \frac{1}{S} \sum_{s=1}^S \tilde{w}_s \quad \text{and} \quad \log p_{\theta}(\mathbf{x}) \approx \log \left(\frac{1}{S} \sum_{s=1}^S \tilde{w}_s \right)$$

where

$$\tilde{w}_s = \frac{p_{\theta}(\mathbf{x}, \mathbf{z}_s)}{q_{\phi}(\mathbf{z}_s | \mathbf{x})}, \quad \mathbf{z}_s \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z} | \mathbf{x}), \quad \sup p_{\theta} \subseteq \sup q_{\phi}$$

Importance weighted variational autoencoder

Empirical estimate of the importance weighted variational autoencoder (IWAE) objective

$$\tilde{\mathcal{L}}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \log \left(\frac{1}{K} \sum_{k=1}^K \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right) = \log \left(\frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right) \approx \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

where

$$\mathbf{z}_k \stackrel{\text{i.i.d.}}{\sim} q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})$$

Importance weighted variational autoencoder

Empirical estimate of the importance weighted variational autoencoder (IWAE) objective

$$\tilde{\mathcal{L}}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \log \left(\frac{1}{K} \sum_{k=1}^K \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right) = \log \left(\frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right) \approx \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

where

$$\mathbf{z}_k \stackrel{\text{i.i.d.}}{\sim} q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})$$

By defining $q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x}) = \prod_{k=1}^K q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})$, IWAE objective (also called as the multiple-sample ELBO) is

$$\mathcal{L}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\log \frac{1}{K} \sum_{k=1}^K \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\log \frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right],$$

where

$$\tilde{w}_k = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})}$$

Importance weighted variational autoencoder: properties

- ① IWAE objective reduces to the standard single sample ELBO when $K = 1$

$$\mathcal{L}_1^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:1} \mid \mathbf{x})} \left[\log \frac{1}{1} \sum_{k=1}^1 \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \right] = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x})$$

Importance weighted variational autoencoder: properties

- ① IWAE objective reduces to the standard single sample ELBO when $K = 1$

$$\mathcal{L}_1^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:1} \mid \mathbf{x})} \left[\log \frac{1}{1} \sum_{k=1}^1 \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \right] = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x})$$

- ② IWAE objective is a lower bound for $\log p_{\boldsymbol{\theta}}(\mathbf{x})$

$$\mathcal{L}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\log \frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] \stackrel{\text{Jensen}}{\leq} \log \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] = \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Importance weighted variational autoencoder: properties

- ① IWAE objective reduces to the standard single sample ELBO when $K = 1$

$$\mathcal{L}_1^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:1} \mid \mathbf{x})} \left[\log \frac{1}{1} \sum_{k=1}^1 \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \right] = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x})$$

- ② IWAE objective is a lower bound for $\log p_{\boldsymbol{\theta}}(\mathbf{x})$

$$\mathcal{L}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\log \frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] \stackrel{\text{Jensen}}{\leq} \log \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] = \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

- ③ For any $K \geq 1$, IWAE lower bound satisfies

$$\mathcal{L}_K^{\text{IWAE}} \leq \mathcal{L}_{K+1}^{\text{IWAE}} \leq \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Importance weighted variational autoencoder: properties

- ① IWAE objective reduces to the standard single sample ELBO when $K = 1$

$$\mathcal{L}_1^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:1} \mid \mathbf{x})} \left[\log \frac{1}{1} \sum_{k=1}^1 \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_k)}{q_{\boldsymbol{\phi}}(\mathbf{z}_k \mid \mathbf{x})} \right] = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \left[\log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x})} \right] = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x})$$

- ② IWAE objective is a lower bound for $\log p_{\boldsymbol{\theta}}(\mathbf{x})$

$$\mathcal{L}_K^{\text{IWAE}}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\log \frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] \stackrel{\text{Jensen}}{\leq} \log \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_{1:K} \mid \mathbf{x})} \left[\frac{1}{K} \sum_{k=1}^K \tilde{w}_k \right] = \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

- ③ For any $K \geq 1$, IWAE lower bound satisfies

$$\mathcal{L}_K^{\text{IWAE}} \leq \mathcal{L}_{K+1}^{\text{IWAE}} \leq \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

- ④ IWAE lower bound approaches the true log marginal likelihood in the limit

$$\lim_{K \rightarrow \infty} \mathcal{L}_K^{\text{IWAE}} = \log p_{\boldsymbol{\theta}}(\mathbf{x})$$

Identifiability of VAEs

Consider VAE model with linear parametrization

$$\begin{aligned} p(\mathbf{x} \mid \mathbf{z}) &= \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \\ q(\mathbf{z} \mid \mathbf{x}) &= \mathcal{N}(\mathbf{V}(\mathbf{x} - \boldsymbol{\mu}), \mathbf{D}), \end{aligned}$$

where \mathbf{D} is diagonal and all parameters are shared with all samples

VAE models are trained via ELBO $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) \leq \log p_{\boldsymbol{\theta}}(\mathbf{x})$

Probabilistic principle component analysis (PPCA) is equivalent with the above **decoder model**

- Marginal likelihood $\log p(\mathbf{x})$ has closed-form
- Posterior $p(\mathbf{z} \mid \mathbf{x})$ has closed-form
- Maximum likelihood parameter estimates have also closed form (up to rotation), i.e., identifiable

Identifiability of VAEs

Consider VAE model with linear parametrization

$$p(\mathbf{x} \mid \mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$
$$q(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{V}(\mathbf{x} - \boldsymbol{\mu}), \mathbf{D}),$$

where \mathbf{D} is diagonal and all parameters are shared with all samples

VAE models are trained via ELBO $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi} \mid \mathbf{x}) \leq \log p_{\boldsymbol{\theta}}(\mathbf{x})$

Probabilistic principle component analysis (PPCA) is equivalent with the above **decoder model**

- Marginal likelihood $\log p(\mathbf{x})$ has closed-form
- Posterior $p(\mathbf{z} \mid \mathbf{x})$ has closed-form
- Maximum likelihood parameter estimates have also closed form (up to rotation), i.e., identifiable

Identifiability of linear VAE (Lucas et al., 2019)

The global maximum of the ELBO for the linear VAE

- Is identical to the global maximum for the log marginal likelihood of PPCA
- Has the scaled principal components as the columns of the decoder network

Identifiability of VAEs

An unconstrained deep latent variable model $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x} | \mathbf{z})p_{\theta}(\mathbf{z})$ is not in general identifiable
Model identifiability can be proved by assuming specific sparsity in $p_{\theta}(\mathbf{x} | \mathbf{z})$ and additional technical assumptions (Moran et al., 2022)

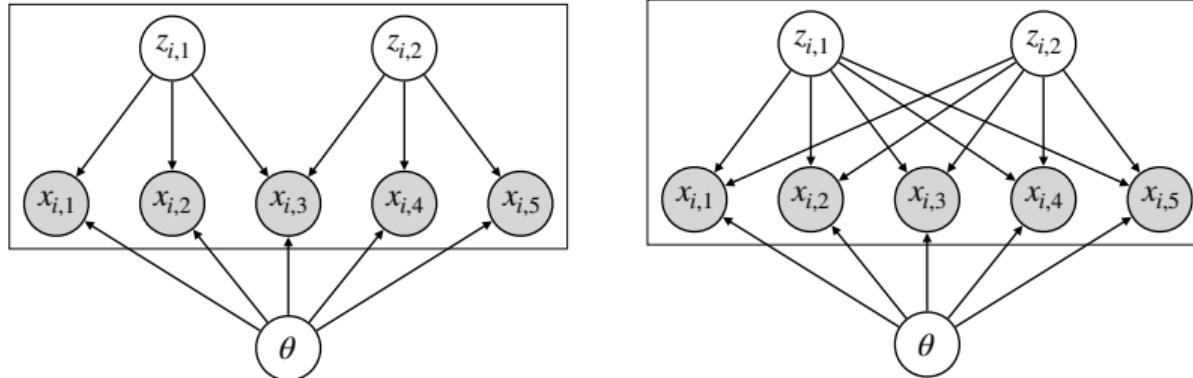


Figure from (Moran et al, 2022)

Identifiability of CVAEs

Consider a CVAE with conditional generative model $p_{\theta}(\mathbf{x}, \mathbf{z} | \mathbf{c}) = p_{\theta}(\mathbf{x} | \mathbf{z})p_{\theta}(\mathbf{z} | \mathbf{c})$

- Assume prior $p_{\theta}(\mathbf{z} | \mathbf{c})$ is mean field with exponential family of distributions
- Covariate \mathbf{c} has at least $\dim(\mathbf{z}) \cdot k + 1$ distinct values, where k denotes the number of sufficient statistics of the prior
- Amortized variational approximation $q_{\phi}(\mathbf{z} | \mathbf{x}, \mathbf{c})$ includes the true posterior
- + Additional technical assumptions (Khemakem et al., 2020)

Then the parameters of the CVAE are **asymptotically** identifiable (Khemakem et al., 2020)

Amortized VI: amortization gap

References

- Agarwal S, Hope G, Younis A, Sudderth EB, A decoder suffices for query-adaptive variational inference, *Uncertainty in Artificial Intelligence*, 2023.
- Khemakhem I, Kingma D, Monti R, Hyvärinen A, Variational autoencoders and nonlinear ICA: A unifying framework, *International Conference on Artificial Intelligence and Statistics*, 2020.
- Kingma DP and Welling M, An Introduction to variational autoencoders, *Foundations and Trends in Machine Learning*, Vol. 12, No. 4, pp. 307-392.
- Kingma DP, Rezende DJ, Mohamed S, Welling M, Semi-supervised learning with deep generative models, *Neural information processing systems*, 2014.
- Lucas J, Tucker G, Grosse R, Norouzi M, Don't blame the ELBO! A linear VAE perspective on posterior collapse, *Neural Information Processing Systems*, 2019.
- Maaloe L, Sonderby CK, Sønderby SK, Winther O, Auxiliary deep generative models, *International Conference on Machine Learning*, PMLR 48:1445-1453, 2016.
- Moran GE, Sridhar D, Wang Y, Blei DM, Identifiable deep generative models via sparse decoding, *Transactions on Machine Learning Research*, 2022.
- Murphy K, Probabilistic Machine Learning: Advanced Topics, 2023.
- Prince SJD, Understanding Deep Learning, The MIT Press, 2023.
- Strauss R, Oliva J, Posterior matching for arbitrary conditioning, *Advances in Neural Information Processing Systems*, 2022.