

Variational Autoencoders

Probabilistic Decision Making — Lecture 10

3rd December 2025

Robert Peharz

Institute of Machine Learning and Neural Computation

Graz University of Technology

Latent variable model:

$$\log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) = \log \int p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta}) d\mathbf{z}^{(i)}$$

Sample-wise **evidence lower bound (ELBO)**:

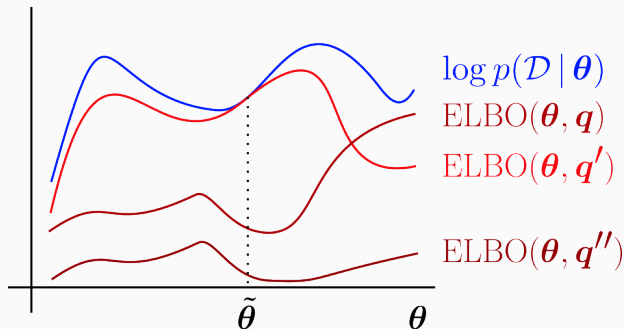
$$\begin{aligned} \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta}) &= \log \int \overbrace{\frac{q_i(\mathbf{z}^{(i)})}{q_i(\mathbf{z}^{(i)})}}^{=1} p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta}) d\mathbf{z}^{(i)} \\ &= \log \int q_i(\mathbf{z}^{(i)}) \frac{p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta})}{q_i(\mathbf{z}^{(i)})} d\mathbf{z}^{(i)} \\ &= \log \mathbb{E}_{q_i} \left[\frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \boldsymbol{\theta})}{q_i(\mathbf{Z}^{(i)})} \right] \\ &\stackrel{\text{Jensen}}{\geq} \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \boldsymbol{\theta})}{q_i(\mathbf{Z}^{(i)})} \right] =: \text{ELBO}(\boldsymbol{\theta}, q_i) \end{aligned}$$

Lower bound approximation of the log-likelihood.

(Total) ELBO:

$$\log p(\mathcal{D} | \theta) = \sum_i \log p(\mathbf{x}^{(i)} | \theta) \geq \sum_i \text{ELBO}(\theta, q_i) =: \text{ELBO}(\theta, \mathbf{q}),$$

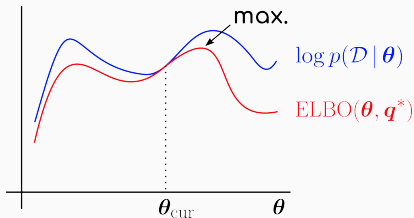
where $\mathbf{q} = \{q_i\}_{i=1}^N$ is the set of all **variational distributions**.



$$\log p(\mathcal{D} | \theta) = \text{ELBO}(\theta, \mathbf{q}) + \underbrace{\sum_i \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta))}_{\text{(total) variational gap}}$$

When $q_i^* \equiv p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta_{\text{cur}})$ for all $\mathbf{q}^* = \{q_i^*\}_{i=1}^N$, then ELBO is tight at θ_{cur} .

Maximizing $\text{ELBO}(\theta, \mathbf{q}^*)$ is a **minorization-maximization** method aka **expectation-maximization** (EM).



Variational Autoencoders

Variational Autoencoders (VAEs)

- **variational autoencoders (VAEs)** are “deep generative models,” combining mixture models and neural networks
- the probabilistic model, the **density network**, was already proposed by MacKay (1995), but learning and inference were slow and hard
- Kingma and Welling (2014) as well as Rezende et al. (2014) used several tricks to make them work:
 - **amortized variational inference**
 - **“re-parametrization” trick**
 - stochastic optimization
- this combination of model and learning/inference techniques can be understood as an autoencoder structure, hence the name “variational autoencoder”

Let \mathbf{Z} be a K -dimensional Gaussian with zero mean and identity covariance, a latent RV often called **code**:

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

Let a neural network with parameters θ be given, the so-called **decoder**, that takes a sample \mathbf{z} as input and returns two D -dimensional vectors

$$\mu_{\theta}(\mathbf{z}), \sigma_{\theta}^2(\mathbf{z})$$

They specify a Gaussian conditional distribution

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}^2(\mathbf{z}))),$$

where $\text{diag}(\mathbf{v})$ is the diagonal matrix with vector \mathbf{v} as diagonal.

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

$$p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{z}), \text{diag}(\boldsymbol{\sigma}_{\boldsymbol{\theta}}^2(\mathbf{z})))$$

prior
decoder

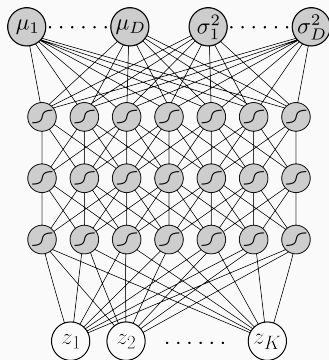
Prior and decoder give rise to a joint:

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

Treating \mathbf{Z} as latent variable we get the **density network** model:

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\theta}) &= \int p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) d\mathbf{z} \\ &= \int p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} \end{aligned}$$

Note that the decoder parameters $\boldsymbol{\theta}$ (i.e. neural network weights) are the only parameters of the model.

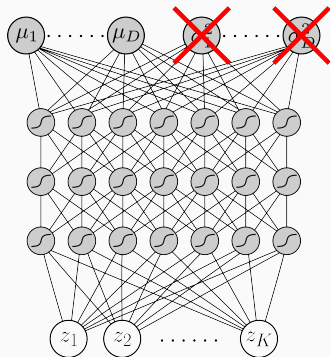


Density Network with Fixed Variances

Often, the density network is simplified to just output a mean, while the variances are fixed:

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \sigma^2 I),$$

where σ^2 is some hyperparameter (might be fixed or trained as well).



Density Networks vs. GMMs

Density networks and GMMs are close relatives. Both are Gaussian mixture models:

GMM

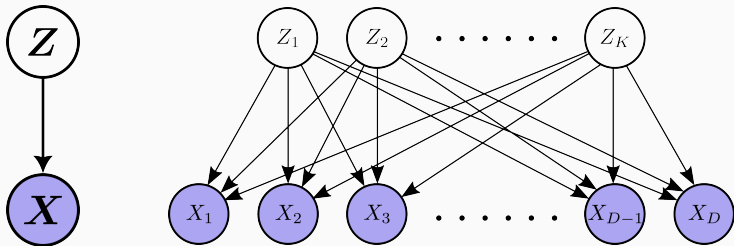


density network



The main difference is that GMMs are **finite** mixture models, while density networks are **infinite mixture models**.

Density Networks as Bayesian Network

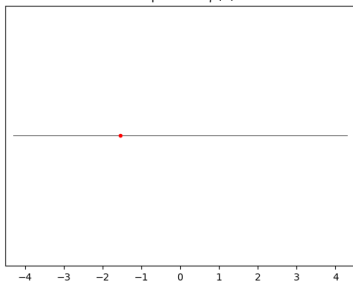


Left: density network as Bayesian network over vectors \mathbf{Z} and \mathbf{X}

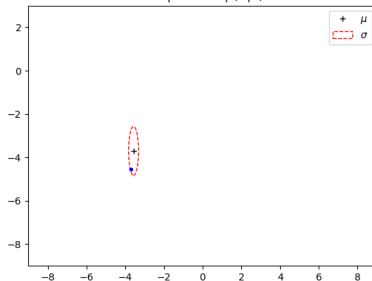
Right: density network as Bayesian network over single RVs

A toy density network with $K = 1$ (dimensionality of Z) and $D = 2$ (dimensionality of \mathbf{X}). Each sample of Z is transformed into a Gaussian over \mathbf{X} .

samples from $p(Z)$

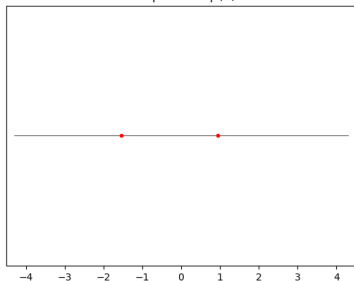


samples from $p(\mathbf{X}|Z)$

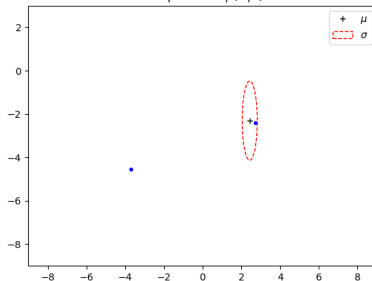


A toy density network with $K = 1$ (dimensionality of Z) and $D = 2$ (dimensionality of \mathbf{X}). Each sample of Z is transformed into a Gaussian over \mathbf{X} .

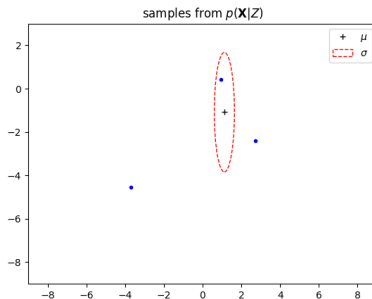
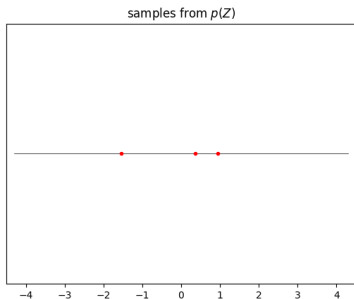
samples from $p(Z)$



samples from $p(\mathbf{X}|Z)$

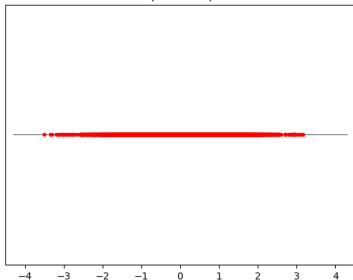


A toy density network with $K = 1$ (dimensionality of Z) and $D = 2$ (dimensionality of \mathbf{X}). Each sample of Z is transformed into a Gaussian over \mathbf{X} .

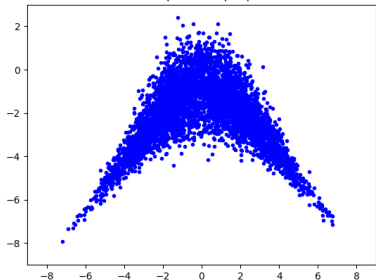


A toy density network with $K = 1$ (dimensionality of Z) and $D = 2$ (dimensionality of \mathbf{X}). Each sample of Z is transformed into a Gaussian over \mathbf{X} .

samples from $p(Z)$



samples from $p(\mathbf{X}|Z)$



How to Learn Density Networks?

Ideally, we would like to train by maximum likelihood:

$$\arg \max_{\theta} \log p(\mathcal{D} | \theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \theta)$$

However, the sample-wise likelihoods are defined as

$$p(\mathbf{x} | \theta) = \int \overbrace{p(\mathbf{x} | \mathbf{z}, \theta)}^{\mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}^2(\mathbf{z})))} p(\mathbf{z}) \mathrm{d}\mathbf{z},$$

a K -dimensional integral over a non-linear function $\mu_{\theta}, \sigma_{\theta}^2$.

This integral is analytically completely intractable.

Expectation-Maximization?

Can we employ the EM algorithm?

The key step is computing for each sample the posterior over latent variables given observed variable, using the current model:

$$p(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}) = \frac{p(\mathbf{x} \mid \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z})}{p(\mathbf{x} \mid \boldsymbol{\theta})} = \frac{p(\mathbf{x} \mid \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z})}{\int p(\mathbf{x} \mid \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z}}$$

Same integral as before, intractable.

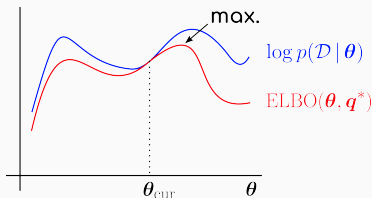
Variational EM

Recall that EM is just a minorization-maximization technique.

$$\log p(\mathcal{D} | \theta) = \text{ELBO}(\theta, \mathbf{q}) + \sum_i \text{KL}(q_i || p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta))$$

Tight ELBO: $q_i^* \equiv p(\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \theta_{\text{cur}})$

However, we might just give up on a tight ELBO and rather follow a **variational EM approach**:



$$\max_{\theta, \mathbf{q}} \text{ELBO}(\theta, \mathbf{q})$$

max w.r.t. \mathbf{q} : improve ELBO quality

max w.r.t. θ : learn model with ELBO as log-likelihood surrogate

How to represent the variational distributions?

$$\max_{\theta, \mathbf{q}} \text{ELBO}(\theta, \mathbf{q})$$

- θ are the neural network weights of the decoder
- $\mathbf{q} = \{q_i\}_{i=1}^N$ are the variational distributions

We might use parametric q_i 's, for example diagonal Gaussians.

However, we need **one q_i for each sample**—for large N this becomes quite a large number of variational parameters to be learned. Moreover, as soon as θ is changed, the \mathbf{q} need to be adapted as well.

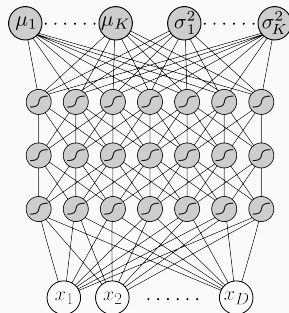
Is there a more compact way to represent \mathbf{q} ?

Before, we introduced the **decoder**, that represents a conditional $p(\mathbf{x} | \mathbf{z})$. Why not do this “trick” again and use a neural network for representing $q_i(\mathbf{z}^{(i)}) := q(\mathbf{z}^{(i)} | \mathbf{x}^{(i)})$?

Specifically, we might construct a neural network with **variational parameters** ϕ that takes \mathbf{x} as input and returns two K -dimensional vectors $\mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})$ specifying

$$p(\mathbf{z} | \mathbf{x}, \phi) = \mathcal{N}(\mu_\phi(\mathbf{x}), \text{diag}(\sigma_\phi^2(\mathbf{x})))$$

We call this network the **inference network** or **encoder**. The fact that we represent **all** variational distributions with a **single** neural net is called **amortized inference**.



Overview

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

prior

θ

decoder params

$$\mu_{\theta}(\mathbf{z}), \sigma_{\theta}^2(\mathbf{z})$$

decoder outputs

ϕ

encoder params

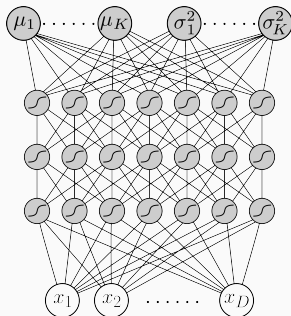
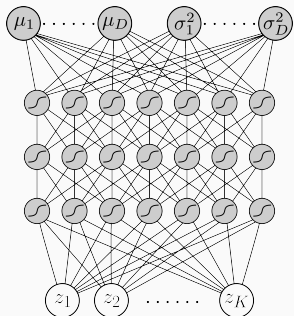
$$\mu_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})$$

encoder outputs

$$p(\mathbf{x} | \mathbf{z}, \theta) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \text{diag}(\sigma_{\theta}^2(\mathbf{z})))$$

$$q(\mathbf{z} | \mathbf{x}, \phi) = \mathcal{N}(\mu_{\phi}(\mathbf{x}), \text{diag}(\sigma_{\phi}^2(\mathbf{x})))$$

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



Learning Algorithm

Variational EM

The overall learning problem is now

$$\max_{\theta, \phi} \text{ELBO}(\theta, \phi)$$

- **maximization w.r.t. ϕ :** learn inference network, tighten lower bound
- **maximization w.r.t. θ :** learn generative model

Basic idea: apply **gradient ascent** simultaneously to ϕ and θ , meaning we need

$$\nabla_{\theta} \text{ELBO}(\theta, \phi) \quad \text{and} \quad \nabla_{\phi} \text{ELBO}(\theta, \phi)$$

Spoiler: we will get only an (unbiased) estimate of the gradient.

Decomposing ELBO

In the following, let $\mathbb{E}_{q_i}[\cdot]$ be short for $\mathbb{E}_{q(\mathbf{z}^{(i)}|\mathbf{x}^{(i)},\phi)}[\cdot]$

Recall that the overall ELBO is the sum of sample-wise ELBOs:

$$\text{ELBO}(\boldsymbol{\theta}, \phi) = \sum_{i=1}^N \text{ELBO}_i(\boldsymbol{\theta}, \phi) = \sum_{i=1}^N \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \boldsymbol{\theta})}{q(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \phi)} \right]$$

$$\begin{aligned} \text{ELBO}_i(\boldsymbol{\theta}, \phi) &= \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)}, \mathbf{Z}^{(i)} | \boldsymbol{\theta})}{q(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \phi)} \right] \\ &= \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{x}^{(i)} | \mathbf{Z}^{(i)}, \boldsymbol{\theta}) p(\mathbf{Z}^{(i)})}{q(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \phi)} \right] \\ &= \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] + \mathbb{E}_{q_i} \left[\log \frac{p(\mathbf{Z}^{(i)})}{q(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \phi)} \right] \\ &= \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] - \mathbb{E}_{q_i} \left[\log \frac{q(\mathbf{Z}^{(i)} | \mathbf{x}^{(i)}, \phi)}{p(\mathbf{Z}^{(i)})} \right] \\ &= \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] - \text{KL}(q_i || p(\mathbf{z}^{(i)})) \end{aligned}$$

Interlude: Different ELBO Decompositions

The last time we saw

$$\text{ELBO}_i(\theta, \phi) = \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \theta) \right] + H(q_i)$$

and

$$\text{ELBO}_i(\theta, \phi) = \log p(\mathbf{x}^{(i)} \mid \theta) - \text{KL}(q_i \parallel p(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}))$$

We just derived:

$$\text{ELBO}_i(\theta, \phi) = \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i)}, \theta) \right] - \text{KL}(q_i \parallel p(\mathbf{z}^{(i)}))$$

These decompositions just differ whether and which chain rule we apply to the joint model.

Deriving ELBO Gradient Estimate

In the following, we will use the third decomposition

$$\text{ELBO}_i(\theta, \phi) = \mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} \mid \mathbf{Z}^{(i)}, \theta) \right] - \text{KL}(q_i \parallel p(\mathbf{z}^{(i)}))$$

and derive gradient estimates w.r.t. θ and ϕ .

Since gradients of sums are the sum of the gradients, we can treat the two terms

- $\mathbb{E}_{q_i} [\log p(\mathbf{x}^{(i)} \mid \mathbf{Z}^{(i)}, \theta)]$
- $\text{KL}(q_i \parallel p(\mathbf{z}^{(i)}))$

independently. We will start with the second term.

Deriving ELBO Gradient Estimate cont'd

We start with $\text{KL}(q_i || p(\mathbf{z}^{(i)}))$

- q_i is Gaussian $\mathcal{N}(\mu_\phi(\mathbf{x}^{(i)}), \text{diag}(\sigma_\phi^2(\mathbf{x}^{(i)})))$
- $p(\mathbf{z}^{(i)})$ is Gaussian $\mathcal{N}(\mathbf{0}, I)$

There is a closed-form solution (Kingma and Welling, 2013)

$$\text{KL}(q_i || p(\mathbf{z}^{(i)})) = -\frac{1}{2} \sum_{j=1}^K \left[1 + \log \sigma_{\phi,j}^2(\mathbf{x}^{(i)}) - \mu_{\phi,j}^2(\mathbf{x}^{(i)}) - \sigma_{\phi,j}^2(\mathbf{x}^{(i)}) \right]$$

Recall that $\mu_{\phi,j}(\mathbf{x}^{(i)})$ and $\sigma_{\phi,j}^2(\mathbf{x}^{(i)})$ are just the outputs of the encoder.

- no dependence on θ , hence $\nabla_{\theta} \text{KL}(q_i || p(\mathbf{z}^{(i)})) = \mathbf{0}$
- $\nabla_{\phi} \text{KL}(q_i || p(\mathbf{z}^{(i)}))$ is easily obtained by automatic differentiation (backprop)



Deriving ELBO Gradient Estimate cont'd

Let's turn to

$$\mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} \mid \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] = \int q(\mathbf{z}^{(i)} \mid \mathbf{x}^{(i)}, \phi) \log p(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i)}, \boldsymbol{\theta}) d\mathbf{z}^{(i)}$$

Note: depends on both $\boldsymbol{\theta}$ (through p) and on ϕ (through q)

Bad news: intractable to compute this integral

However: we might draw samples from q_i (Gaussian) and do a Monte Carlo estimate. One often just uses a **single** sample—which is still an unbiased (but noisy) estimator:

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

$$\mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} \mid \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] \approx \log p(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i)}, \boldsymbol{\theta})$$

Deriving ELBO Gradient Estimate cont'd

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

$$\mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \mathbf{Z}^{(i)}, \boldsymbol{\theta}) \right] \approx \log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}, \boldsymbol{\theta})$$

Taking the gradient w.r.t. $\boldsymbol{\theta}$ is no problem.

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}, \boldsymbol{\theta}) =$$

$$\nabla_{\boldsymbol{\theta}} \left(-\frac{1}{2} \sum_{d=1}^D \log(2\pi) + \log \sigma_{\boldsymbol{\theta},d}^2(\mathbf{z}^{(i)}) + \frac{(x_d^{(i)} - \mu_{\boldsymbol{\theta},d}(\mathbf{z}^{(i)}))^2}{\sigma_{\boldsymbol{\theta},d}^2(\mathbf{z}^{(i)})} \right)$$

Recall that $\mu_{\boldsymbol{\theta},d}(\mathbf{z}^{(i)})$ and $\sigma_{\boldsymbol{\theta},d}^2(\mathbf{z}^{(i)})$ are just the outputs of the decoder—computing the gradient w.r.t. $\boldsymbol{\theta}$ is easy with automatic differentiation (backprop).

Deriving ELBO Gradient Estimate cont'd

$$\mathbf{z}^{(i)} \sim \mathcal{N}(\mu_{\phi}(\mathbf{x}^{(i)}), \text{diag}(\sigma_{\phi}^2(\mathbf{x}^{(i)})))$$

$$\mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} \mid \mathbf{Z}^{(i)}, \theta) \right] \approx \log p(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i)}, \theta)$$

However, $\mathbf{z}^{(i)}$ is sampled from a Gaussian whose parameters depend on ϕ —how to “backprop” through the “sampling operation \sim ”?

The “Reparameterization Trick”


Note that sampling

$$\mathbf{z}^{(i)} \sim \mathcal{N}(\mu_{\phi}(\mathbf{x}^{(i)}), \text{diag}(\sigma_{\phi}^2(\mathbf{x}^{(i)})))$$

can be done like

$$\begin{aligned}\epsilon &\sim \mathcal{N}(\mathbf{0}, I) \\ \mathbf{z}^{(i)} &= \mu_{\phi}(\mathbf{x}^{(i)}) + \epsilon \cdot \sigma_{\phi}(\mathbf{x}^{(i)})\end{aligned}$$

where \cdot is element-wise multiplication.

ϵ is helper noise, sampled from a fixed distribution, independent of ϕ . Now differentiating $\mathbf{z}^{(i)}$ w.r.t. ϕ is straight-forward. 

This is called the **reparameterization trick**.

For a given sample $\mathbf{x}^{(i)}$ do

- pass $\mathbf{x}^{(i)}$ through encoder, yielding μ_ϕ and σ_ϕ^2
- compute $\text{KL} = -\frac{1}{2} \sum_{j=1}^K \left[1 + \log \sigma_{\phi,j}^2 - \mu_{\phi,j}^2 - \sigma_{\phi,j}^2 \right]$
- sample $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- let $\mathbf{z}^{(i)} = \mu_\phi(\mathbf{x}^{(i)}) + \epsilon \cdot \sigma_\phi(\mathbf{x}^{(i)})$
- pass $\mathbf{z}^{(i)}$ through decoder, yielding μ_θ and σ_θ^2
- compute $\log p = -\frac{1}{2} \sum_{d=1}^D \left[\log(2\pi) + \log \sigma_{\theta,d}^2 + \frac{(x_d^{(i)} - \mu_{\theta,d})^2}{\sigma_{\theta,d}^2} \right]$
- $\widehat{\text{ELBO}}_i = \log p - \text{KL}$
- compute $\nabla_\theta \widehat{\text{ELBO}}_i$ and $\nabla_\phi \widehat{\text{ELBO}}_i$ via autodiff

Final Learning Algorithm: Stochastic Gradient Ascent

In each iteration:

- draw mini-batch $\mathcal{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(L)}\}$ of samples
- for each $\mathbf{x}^{(i)} \in \mathcal{B}$, compute

$$\nabla_{\theta} \widehat{\text{ELBO}}_i \quad \text{and} \quad \nabla_{\phi} \widehat{\text{ELBO}}_i$$

using the procedure on the previous slide. Note that this can be parallelized over the mini-batch

- compute

$$\nabla_{\theta} \widehat{\text{ELBO}} = \sum_{l=1}^L \nabla_{\theta} \widehat{\text{ELBO}}_i \quad \text{and} \quad \nabla_{\phi} \widehat{\text{ELBO}} = \sum_{l=1}^L \nabla_{\phi} \widehat{\text{ELBO}}_i$$

and perform a gradient ascent step.

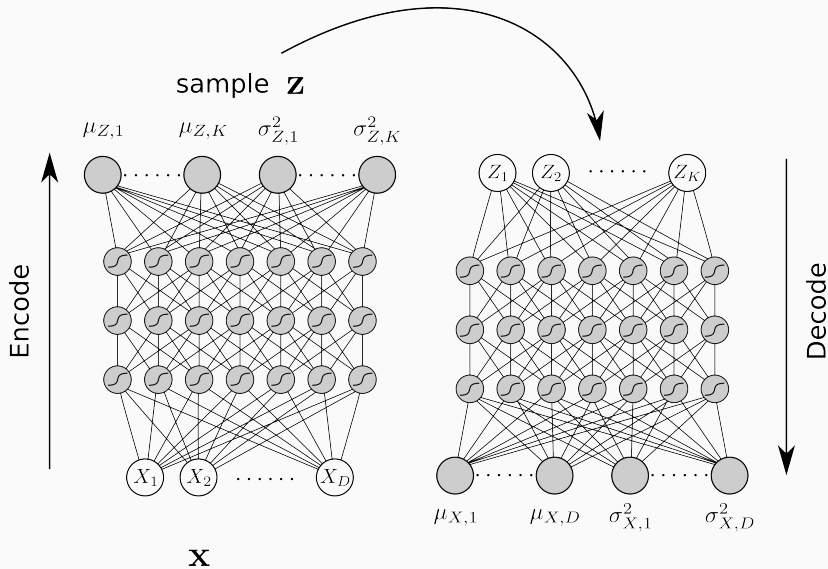
Note that this learning scheme resembles **classical autoencoders**

- encode sample $\mathbf{x}^{(i)}$ into a latent representation $\mathbf{z}^{(i)}$
- decode $\mathbf{z}^{(i)}$ into a (probabilistic) reconstruction $\log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)})$
- in this interpretation the ELBO can be understood as

$$\underbrace{\mathbb{E}_{q_i} \left[\log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}, \theta) \right]}_{\text{reconstruction loss}} - \underbrace{\text{KL}(q_i || p(\mathbf{z}^{(i)}))}_{\text{regularizer}}$$

- nice side effect: encoder delivers a **representation** (code, embedding) $\mathbf{z}^{(i)}$ for data $\mathbf{x}^{(i)}$
- **but: only the decoder specifies the probabilistic model—the encoder is merely an inference tool!**

Variational Autoencoder



VAE on MNIST Data

MNIST Data

28×28 gray scale images (784-dimensional vectors)

Training Data: 60,000 samples

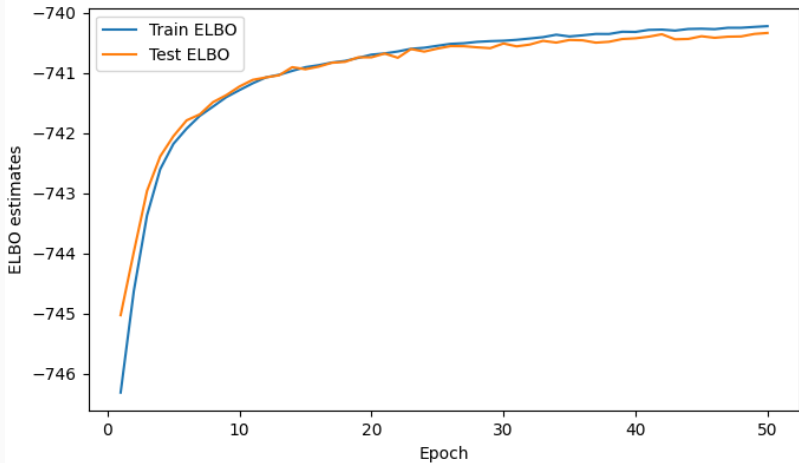
Test Data: 10,000 samples



VAE Architecture and Training Details

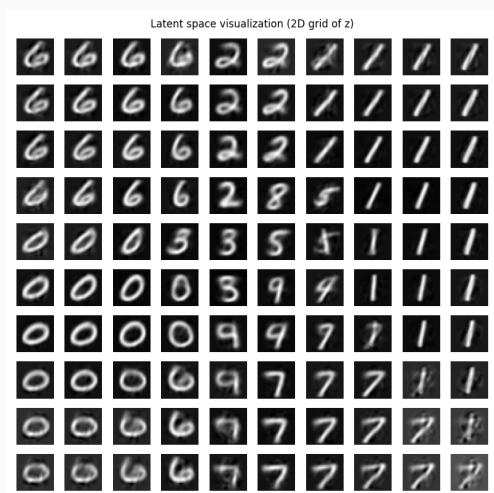
- $K = 2$ latent dimensions (for simplicity and visualization)
- encoder and decoder are both 2-layer MLPs with 200 units, using relu-activation functions
- train for 50 epochs with Adam as optimizer
- stepsize 0.001 and batch-size 128
- the decoder used a fixed variance of 1, i.e. only means are predicted

Training Curves



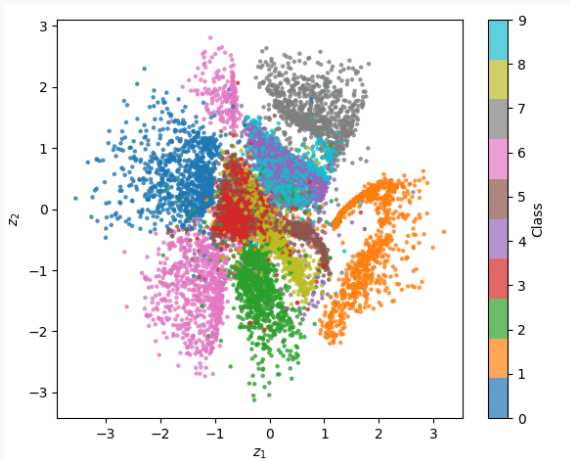
Latent Space Visualization

Sampling z on a grid from -3 to 3 , followed by decoding, we can visualize the learned latent space (only decoded means are shown)



Embedding

The encoder delivers for each sample an embedding vector $\mathbf{z}^{(i)}$, that can be used as **dimensionality reduction** technique. Here, we color-code the classes (class information was not used in training):



VAE: Remarks and Extensions

- VAEs are one of **the** most prominent deep generative models
- Kingma & Welling: cited over **50,000 times today**
- **historical note:** actually the model is a density network (MacKay, 1995), but the VAE tricks made them fly (stochastic gradient, amortized inference, etc.)
- many extensions exist:
 - combinations with convnets for image data
 - combinations with graph networks, e.g. for molecule data
 - using multiple samples for \mathbf{Z} (importance weighted autoencoder)
 - applications to time series and dynamical systems
 - trading off reconstruction and KL term (β -VAE)
 - learning the prior (VampPrior)