

Variational Autoencoders

Elia Bronzo, Shoichi Yip

05/03/2024

- 1 VAEs from a probabilistic prospective
 - Learning probabilistic models
 - Latent variable models
 - Learning latent variable models with Expectation-Maximization
 - Learning latent variable models with Variational Auto-Encoders
- 2 VAEs from the perspective of neural networks
 - Autoencoder Neural Networks
 - Variational Autoencoder Neural Networks
 - Why use Neural Networks?
 - Finding the loss function
 - Computing the \mathcal{L}
 - The reparametrization trick
 - Loss function and NN structure

Learning probabilistic models

Let's start by fixing some notation:

- $\mathbf{x}^{(i)} \in \mathcal{D}$ that we assume to be N i.i.d. from $p^*(\mathbf{x})$ $i = 1, \dots, N$

Learning probabilistic models

Let's start by fixing some notation:

- $\mathbf{x}^{(i)} \in \mathcal{D}$ that we assume to be N i.i.d. from $p^*(\mathbf{x})$ $i = 1, \dots, N$
- $p^*(\mathbf{x})$ is the *True* generative process we want to learn

Learning probabilistic models

Let's start by fixing some notation:

- $\mathbf{x}^{(i)} \in \mathcal{D}$ that we assume to be N i.i.d. from $p^*(\mathbf{x})$ $i = 1, \dots, N$
- $p^*(\mathbf{x})$ is the *True* generative process we want to learn

Let's start considering a **fully observed model**

$$p_{\theta}(\mathbf{x}) \approx p^*(\mathbf{x})$$

Learning probabilistic models

Let's start by fixing some notation:

- $\mathbf{x}^{(i)} \in \mathcal{D}$ that we assume to be N i.i.d. from $p^*(\mathbf{x})$ $i = 1, \dots, N$
- $p^*(\mathbf{x})$ is the *True* generative process we want to learn

Let's start considering a **fully observed model**

$$p_{\theta}(\mathbf{x}) \approx p^*(\mathbf{x})$$

Maximum Likelihood (ML)

$$\begin{aligned}\theta_{ML} &= \arg \max_{\theta} p_{\theta}(\mathcal{D}) \\ &= \arg \max_{\theta} \sum_{i=1}^N p_{\theta}(\mathbf{x}^{(i)})\end{aligned}$$

Learning probabilistic models

Let's start by fixing some notation:

- $\mathbf{x}^{(i)} \in \mathcal{D}$ that we assume to be N i.i.d. from $p^*(\mathbf{x})$ $i = 1, \dots, N$
- $p^*(\mathbf{x})$ is the *True* generative process we want to learn

Let's start considering a **fully observed model**

$$p_{\theta}(\mathbf{x}) \approx p^*(\mathbf{x})$$

Maximum Likelihood (ML)

$$\begin{aligned}\theta_{ML} &= \arg \max_{\theta} p_{\theta}(\mathcal{D}) \\ &= \arg \max_{\theta} \sum_{i=1}^N p_{\theta}(\mathbf{x}^{(i)})\end{aligned}$$

Maximum a posteriori (MAP)

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\theta|\mathcal{D}) \\ &= \arg \max_{\theta} \frac{p_{\theta}(\mathcal{D})p(\theta)}{p(\mathcal{D})}\end{aligned}$$

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \underbrace{p_{\theta}(\mathbf{z})}_{\text{Prior}} p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad \text{Marginal likelihood}$$

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \underbrace{p_{\theta}(\mathbf{z})}_{\text{Prior}} p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad \text{Marginal likelihood}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a DNN \rightarrow **Deep Latent Variable model (DLVM)**

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \underbrace{p_{\theta}(\mathbf{z})}_{\text{Prior}} p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad \text{Marginal likelihood}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a DNN \rightarrow **Deep Latent Variable model (DLVM)**

Why?

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \underbrace{p_{\theta}(\mathbf{z})}_{\text{Prior}} p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad \text{Marginal likelihood}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a DNN \rightarrow **Deep Latent Variable model (DLVM)**

Why?

- high representation capacity: $p_{\theta}(\mathbf{x}) \rightarrow$ (almost) arbitrary dependencies even with simple (ex. Gaussian) priors
- easily trained with SGD

Latent variables models

$p_{\theta}(\mathbf{x}, \mathbf{z})$ where \mathbf{z} latent variables

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \underbrace{p_{\theta}(\mathbf{z})}_{\text{Prior}} p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad \text{Marginal likelihood}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a DNN \rightarrow **Deep Latent Variable model (DLVM)**

Why?

- high representation capacity: $p_{\theta}(\mathbf{x}) \rightarrow$ (almost) arbitrary dependencies even with simple (ex. Gaussian) priors
- easily trained with SGD

VAE is used to learn this kind of models

Learning latent variable models

Learning a latent variable model means optimize the **marginal log-likelihood**:

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) d\mathbf{z}$$

Learning latent variable models

Learning a latent variable model means optimize the **marginal log-likelihood**:

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) d\mathbf{z}$$

The *classical* way of optimizing in this problems is **Expectation-Maximization**

Learning latent variable models

Learning a latent variable model means optimize the **marginal log-likelihood**:

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) d\mathbf{z}$$

The *classical* way of optimizing in this problems is **Expectation-Maximization**

We will first focus on how to optimize $l(\theta)$ in case of a dataset composed by only 1 point ($\mathbf{x}^{(1)}$) ($N=1$)

Evidence Lower Bound (ELBO)

$$\log p_{\theta}(\mathbf{x}^{(1)}) = \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right]$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right]\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right]\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \right] \right] + \mathbb{E}_q \left[\log \left[\frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right]\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \right] \right] + \mathbb{E}_q \left[\log \left[\frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \text{ELBO}(\mathbf{x}^{(1)}; q, \theta) + \mathcal{D}_{\text{KL}} \left[q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}) \right] \geq \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)\end{aligned}$$

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \right] \right] + \mathbb{E}_q \left[\log \left[\frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\ &= \text{ELBO}(\mathbf{x}^{(1)}; q, \theta) + \mathcal{D}_{\text{KL}} \left[q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}) \right] \geq \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)\end{aligned}$$

When does the ELBO is **tight**? (equality)

Evidence Lower Bound (ELBO)

$$\begin{aligned}
 \log p_{\theta}(\mathbf{x}^{(1)}) &= \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}^{(1)}) \right] \\
 &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\
 &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\
 &= \mathbb{E}_q \left[\log \left[\frac{p_{\theta}(\mathbf{x}^{(1)}, \mathbf{z})}{q(\mathbf{z})} \right] \right] + \mathbb{E}_q \left[\log \left[\frac{q(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})} \right] \right] \\
 &= \text{ELBO}(\mathbf{x}^{(1)}; q, \theta) + \mathcal{D}_{\text{KL}} \left[q(\mathbf{z}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x}) \right] \geq \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)
 \end{aligned}$$

When does the ELBO is **tight**? (equality)

$$\text{ELBO}(\mathbf{x}^{(1)}; q, \theta) = \log p_{\theta}(\mathbf{x}^{(1)}) \iff q(\mathbf{z}) = p_{\theta}(\mathbf{z}|\mathbf{x}^{(1)})$$

Expectation-Maximization Algorithm

- **Expectation** step: at fixed $\theta^{(t)}$, $q(\mathbf{z}) \leftarrow p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x}^{(1)})$

Expectation-Maximization Algorithm

- **Expectation** step: at fixed $\theta^{(t)}$, $q(\mathbf{z}) \leftarrow p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x}^{(1)})$
- **Maximization** step: at fixed q , $\theta^{(t+1)} \leftarrow \arg \max_{\theta} \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)$

Expectation-Maximization Algorithm

- **Expectation** step: at fixed $\theta^{(t)}$, $q(\mathbf{z}) \leftarrow p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x}^{(1)})$
- **Maximization** step: at fixed q , $\theta^{(t+1)} \leftarrow \arg \max_{\theta} \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)$

What about a dataset with more than 1 point? $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) \geq \sum_{i=1}^N \text{ELBO}(\mathbf{x}^{(i)}; q_i, \theta)$$

Expectation-Maximization Algorithm

- **Expectation** step: at fixed $\theta^{(t)}$, $q(\mathbf{z}) \leftarrow p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x}^{(1)})$
- **Maximization** step: at fixed q , $\theta^{(t+1)} \leftarrow \arg \max_{\theta} \text{ELBO}(\mathbf{x}^{(1)}; q, \theta)$

What about a dataset with more than 1 point? $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) \geq \sum_{i=1}^N \text{ELBO}(\mathbf{x}^{(i)}; q_i, \theta)$$

Repeat until convergence {

(E-step) **for each i:**

$$q_i(\mathbf{z}) \leftarrow p_{\theta^{(t)}}(\mathbf{z}|\mathbf{x}^{(i)})$$

(M-step)

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{i=1}^N \text{ELBO}(\mathbf{x}^{(i)}; q, \theta) \}$$

DLVM, intractabilities and problems with EM

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \underbrace{\int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) d\mathbf{z}}_{\text{intractable for a DNN}}$$

DLVM, intractabilities and problems with EM

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \underbrace{\int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}) d\mathbf{z}}_{\text{intractable for a DNN}}$$

Intractable means that we don't have neither an *analytical solution* nor an *efficient estimator* of the integral in case of a DNN

DLVM, intractabilities and problems with EM

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \underbrace{\int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) d\mathbf{z}}_{\text{intractable for a DNN}}$$

Intractable means that we don't have neither an *analytical solution* nor an *efficient estimator* of the integral in case of a DNN

Intractable $p_{\theta}(\mathbf{x}) \iff$ Intractable **posterior** $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})} = \frac{p(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x})}$

DLVM, intractabilities and problems with EM

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \underbrace{\int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) d\mathbf{z}}_{\text{intractable for a DNN}}$$

Intractable means that we don't have neither an *analytical solution* nor an *efficient estimator* of the integral in case of a DNN

Intractable $p_{\theta}(\mathbf{x}) \iff$ Intractable **posterior** $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})} = \frac{p(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x})}$

E-step: evaluation of an intractable posterior with MCMC **for each data-point** \rightarrow computationally expensive

DLVM, intractabilities and problems with EM

$$l(\theta) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)}) = \sum_{i=1}^N \underbrace{\int p(\mathbf{z}) p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) d\mathbf{z}}_{\text{intractable for a DNN}}$$

Intractable means that we don't have neither an *analytical solution* nor an *efficient estimator* of the integral in case of a DNN

Intractable $p_{\theta}(\mathbf{x}) \iff$ Intractable **posterior** $p_{\theta}(\mathbf{z}|\mathbf{x}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{p_{\theta}(\mathbf{x})} = \frac{p(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z})}{p_{\theta}(\mathbf{x})}$

E-step: evaluation of an intractable posterior with MCMC **for each data-point** \rightarrow computationally expensive

Solutions? \rightarrow **VAE**

Variational Auto-Encoders (VAE)

Model for the posterior $q_{\phi}(\mathbf{z}|\mathbf{x}) \rightarrow$ another DNN, called **Encoder**

Variational Auto-Encoders (VAE)

Model for the posterior $q_{\phi}(\mathbf{z}|\mathbf{x}) \rightarrow$ another DNN, called **Encoder**

While $p_{\theta}(\mathbf{x}|\mathbf{z})$ is called **Decoder**.

Variational Auto-Encoders (VAE)

Model for the posterior $q_{\phi}(\mathbf{z}|\mathbf{x}) \rightarrow$ another DNN, called **Encoder**

While $p_{\theta}(\mathbf{x}|\mathbf{z})$ is called **Decoder**.

$$\log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) = \underbrace{\mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right]}_{\text{loss}} = \text{ELBO}(\mathbf{x}; \phi, \theta)$$

We still optimize the ELBO, not with EM but directly with **SGD** exploiting back propagation.

Variational Auto-Encoders (VAE)

Model for the posterior $q_{\phi}(\mathbf{z}|\mathbf{x}) \rightarrow$ another DNN, called **Encoder**

While $p_{\theta}(\mathbf{x}|\mathbf{z})$ is called **Decoder**.

$$\log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) = \underbrace{\mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right]}_{\text{loss}} = \text{ELBO}(\mathbf{x}; \phi, \theta)$$

We still optimize the ELBO, not with EM but directly with **SGD** exploiting back propagation.

Difficulty: $\nabla_{\phi} \mathbf{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right] \rightarrow$ **Reparametrization Trick**

Why should we do this?

$$\text{ELBO}(\mathbf{x}; \phi, \theta) = \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Optimizing the ELBO so we do 2 things at the same time:

Why should we do this?

$$\text{ELBO}(\mathbf{x}; \phi, \theta) = \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Optimizing the ELBO so we do 2 things at the same time:

- 1 **Maximize** $p_{\theta}(\mathbf{x}) \rightarrow$ fitting the DLVM

Why should we do this?

$$\text{ELBO}(\mathbf{x}; \phi, \theta) = \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Optimizing the ELBO so we do 2 things at the same time:

- 1 **Maximize** $p_{\theta}(\mathbf{x}) \rightarrow$ fitting the DLVM
- 2 **Minimize** $D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \rightarrow$ learning a posterior over unobserved data \rightarrow typical task of **Variational Inference (VI)**

Why should we do this?

$$\text{ELBO}(\mathbf{x}; \phi, \theta) = \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Optimizing the ELBO so we do 2 things at the same time:

- 1 **Maximize** $p_{\theta}(\mathbf{x}) \rightarrow$ fitting the DLVM
- 2 **Minimize** $D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \rightarrow$ learning a posterior over unobserved data \rightarrow typical task of **Variational Inference (VI)**

EM

(E-step) **for each i:**

$$q_i(\mathbf{z}) \leftarrow p_{\theta(i)}(\mathbf{z}|\mathbf{x}^{(i)})$$

Why should we do this?

$$\text{ELBO}(\mathbf{x}; \phi, \theta) = \log p_{\theta}(\mathbf{x}) - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))$$

Optimizing the ELBO so we do 2 things at the same time:

- ① **Maximize** $p_{\theta}(\mathbf{x}) \rightarrow$ fitting the DLVM
- ② **Minimize** $D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})) \rightarrow$ learning a posterior over unobserved data \rightarrow typical task of **Variational Inference (VI)**

EM

(E-step) **for each i:**

$$q_i(\mathbf{z}) \leftarrow p_{\theta(i)}(\mathbf{z}|\mathbf{x}^{(i)})$$

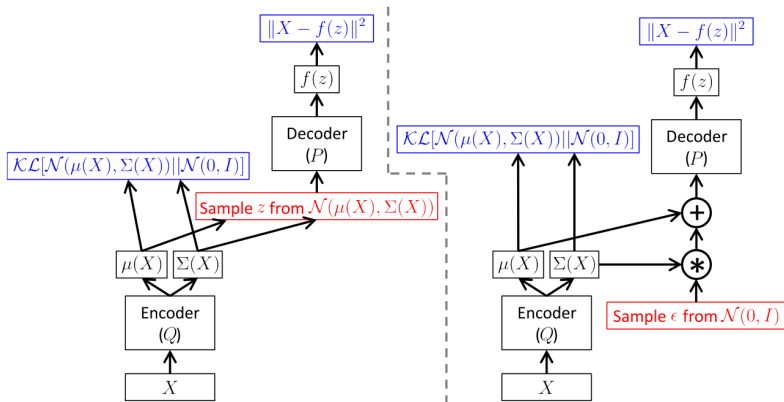
VAE:

Single set of parameters (ϕ) for all data-points

$$q_{\phi}(\mathbf{z}|\mathbf{x})$$

Amortized Inference

$$\underbrace{\mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \left[\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \right]}_{\text{loss}} = \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [p_\theta(\mathbf{x}|\mathbf{z})] - \mathbf{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \right]$$



Latent representations

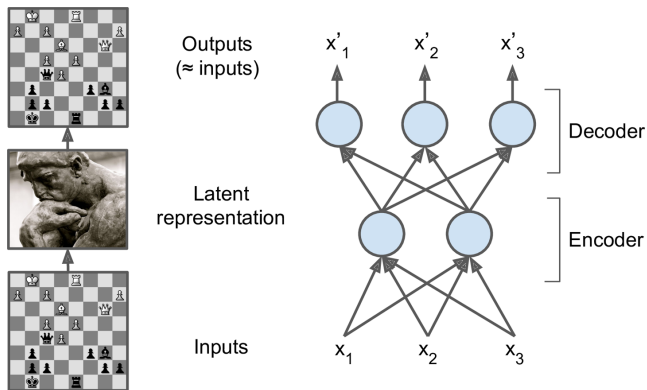


Figure: Chess memory experiment and structure of an autoencoder (from HOML).

Undercompleteness

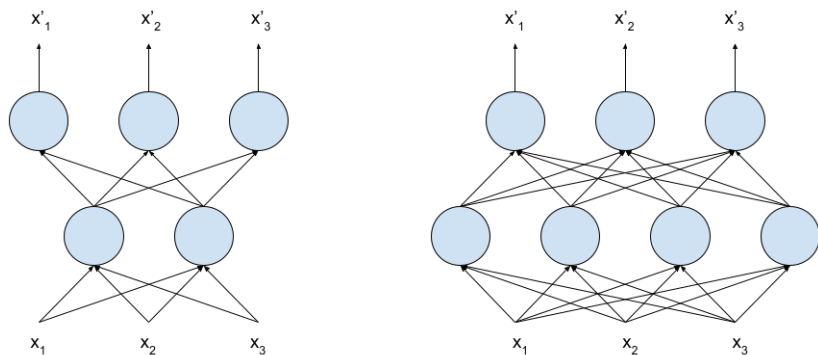


Figure: The typical structure of an autoencoder and of an MLP.

What is the autoencoder learning?

The autoencoder learns the **identity**. We expect the output to resemble to the input.

What is the autoencoder learning?

The autoencoder learns the **identity**. We expect the output to resemble to the input.

But it is characterised by the **information bottleneck** that yields **undercompleteness**.

What is the autoencoder learning?

The autoencoder learns the **identity**. We expect the output to resemble to the input.

But it is characterised by the **information bottleneck** that yields **undercompleteness**.

Depending on the characteristics of the bottleneck, we get different autoencoders.

Stacked autoencoder

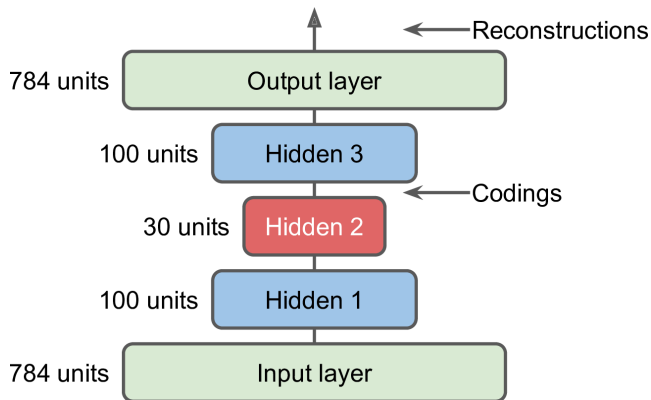


Figure: Structure of a stacked autoencoder (from HOML).

Reconstruction for the stacked autoencoder



Figure: Reconstruction of images on a stacked autoencoder (from HOML).

Denoising autoencoder

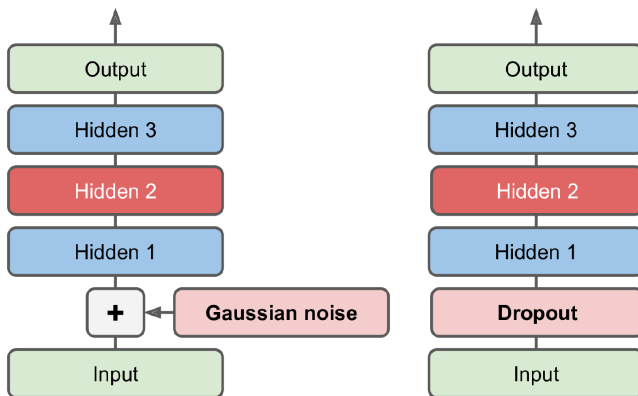


Figure: Structure of a denoising autoencoder (from HOML).

Reconstruction of noiseless images

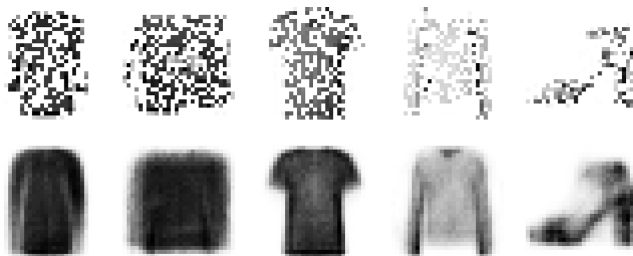


Figure: Reconstruction of noiseless images on a denoising autoencoder (from HOML).

Pros and Cons

What autoencoders **do**:

- learn **latent** representation

Pros and Cons

What autoencoders **do**:

- learn **latent** representation
- **semisupervised learning**

Pros and Cons

What autoencoders **do**:

- learn **latent** representation
- **semisupervised learning**
- **pretraining** using greedy layerwise approach

Pros and Cons

What autoencoders **do**:

- learn **latent** representation
- **semisupervised learning**
- **pretraining** using greedy layerwise approach
- **denoising** images

Pros and Cons

What autoencoders **do not do**:

- have a natural **interpolation** scheme

Pros and Cons

What autoencoders **do not do**:

- have a natural **interpolation** scheme
- have **generative** capacity

Variational autoencoder

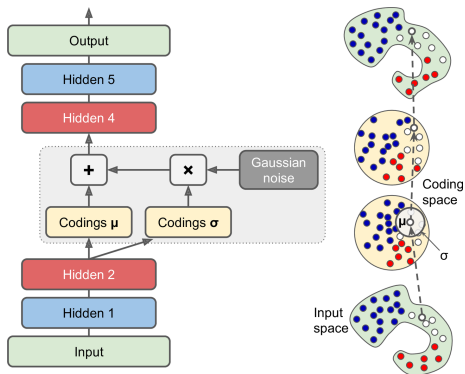


Figure: Structure of a variational autoencoder (from HOML).

The decoder

We have said that our model is **generative**.

¹This is justified by the **universal approximation theory** that states that deep ANNs are able to approximate any function.

The decoder

We have said that our model is **generative**.

The generative role of the variational autoencoder is played by the **decoder**.

¹This is justified by the **universal approximation theory** that states that deep ANNs are able to approximate any function.

The decoder

We have said that our model is **generative**.

The generative role of the variational autoencoder is played by the **decoder**. The role of the decoder is to generate new samples that belong to the input space \mathcal{X} starting from a given vector \mathbf{z} that belongs to the space of latent variables.

The role of the decoder network is to learn the **nonlinear mapping** $p_{\theta}(\mathbf{x}|\mathbf{z})$ that associates a value of \mathbf{x} given a point in the latent space \mathbf{z} .¹

¹This is justified by the **universal approximation theory** that states that deep ANNs are able to approximate any function.

The decoder

We can make the assumption that the $p(\mathbf{x}|\mathbf{z})$ is a multivariate Gaussian distribution, for example in the case of MNIST datasets.

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

The decoder

We can make the assumption that the $p(\mathbf{x}|\mathbf{z})$ is a multivariate Gaussian distribution, for example in the case of MNIST datasets.

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

Inputs in that case are, in fact, vectors of real valued variables (i.e. pixel brightness).

Back to the probabilistic picture

Remember that we are trying to maximise the **likelihood** $p(x_i)$ for each i -th data point.

Back to the probabilistic picture

Remember that we are trying to maximise the **likelihood** $p(x_i)$ for each i -th data point.

Given a new probability distribution $q_i(\mathbf{z})$ we have

$$\begin{aligned} p(x_i) &= \int p(x_i|\mathbf{z})p(\mathbf{z}) \, d\mathbf{z} \\ &= \int p(x_i|\mathbf{z})p(\mathbf{z}) \frac{q_i(\mathbf{z})}{q_i(\mathbf{z})} \, d\mathbf{z} \\ &= \left\langle \frac{p(x_i|\mathbf{z})p(\mathbf{z})}{q_i(\mathbf{z})} \right\rangle_{\mathbf{z} \sim q_i(\mathbf{z})} \end{aligned}$$

Back to the probabilistic picture

Per Jensen, we have that the **log likelihood** has the following bound

$$\begin{aligned}\log(p(x_i)) &\geq \left\langle \log \frac{p(x_i|\mathbf{z})p(\mathbf{z})}{q_i(\mathbf{z})} \right\rangle_{\mathbf{z} \sim q_i(\mathbf{z})} \\ &\geq \langle \log p(x_i|\mathbf{z})(p(\mathbf{z})) \rangle_{\mathbf{z} \sim q_i(\mathbf{z})} - H[q_i(\mathbf{z})]\end{aligned}$$

The RHS is what we have previously called ELBO (Evidence Lower Bound). We choose to maximise the ELBO instead of the log likelihood. But we are maximising against a **lot** of parameters: we have to optimise against θ and against each i -th $q_i(\mathbf{z})$!

Back to the probabilistic picture

Per Jensen, we have that the **log likelihood** has the following bound

$$\begin{aligned}\log(p(x_i)) &\geq \left\langle \log \frac{p(x_i|\mathbf{z})p(\mathbf{z})}{q_i(\mathbf{z})} \right\rangle_{\mathbf{z} \sim q_i(\mathbf{z})} \\ &\geq \langle \log p(x_i|\mathbf{z})(p(\mathbf{z})) \rangle_{\mathbf{z} \sim q_i(\mathbf{z})} - H[q_i(\mathbf{z})]\end{aligned}$$

The RHS is what we have previously called ELBO (Evidence Lower Bound). We choose to maximise the ELBO instead of the log likelihood. But we are maximising against a **lot** of parameters: we have to optimise against θ and against each i -th $q_i(\mathbf{z})$!

Can we reduce the number of parameters to optimise against?

Introducing: the encoder network

What if we add another network, but this time to perform the **inference**?
This means that we **choose** the $q_i(z)$ to be

$$q_i(z) = q_\phi(\mathbf{z}|\mathbf{x})$$

where we call $q_\phi(\mathbf{z}|\mathbf{x})$ the **encoder** and ϕ is the set of parameters of the variational inference.

Deep and amortised

Let us recap. In one go we

- 1 chose to perform the **deep inference**, which implies that we will train a neural network to **predict** the variational parameters ϕ instead of getting them by directly optimising the ELBO.
- 2 reduced the number of parameters that we have to optimise against since we do not have anymore N approximate posteriors $q_i(z)$'s, but one parameter set ϕ shared by all the data points. Since the task of inference is amortised across the entire dataset, this technique is called **amortised variational inference**.

The Kullback-Leibler divergence

We are going to compute the **Kullback-Leibler divergence** for a particular case, in which the real and approximate posteriors have take the **multinomial Gaussian** form.

$$\mathcal{D}_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z})] = \int d\mathbf{z} q_{\phi}(\mathbf{z}|\mathbf{x}) [\log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \quad (1)$$

First term of D_{KL}

$$\begin{aligned} & \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{z}) \, d\mathbf{z} = \\ & = \int d\mathbf{z} \frac{1}{\sqrt{\det(\Sigma)} \sqrt{2\pi}^L} \prod_{l=1}^L \exp\left\{-\frac{(z_l - \mu_l)^2}{2\sigma_l^2}\right\} \log \frac{1}{\sqrt{2\pi}^L} e^{-\frac{|\mathbf{z}|^2}{2}} \end{aligned}$$

First term of D_{KL}

$$\begin{aligned} \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{z}) \, d\mathbf{z} &= \\ &= \int d\mathbf{z} \frac{1}{\sqrt{\det(\Sigma)} \sqrt{2\pi}^L} \prod_{l=1}^L \exp\left\{-\frac{(z_l - \mu_l)^2}{2\sigma_l^2}\right\} \left[-\frac{L}{2} \log 2\pi - \frac{|\mathbf{z}|^2}{2}\right] \\ &= -\frac{L}{2} \log 2\pi - \frac{1}{2} \sum_{l=1}^L (\sigma_l^2 + \mu_l^2) \end{aligned}$$

Second term of D_{KL}

$$\begin{aligned} & \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log q_{\phi}(\mathbf{z}|\mathbf{x}) \, d\mathbf{z} \\ &= \int d\mathbf{z} \frac{1}{\sqrt{\det(\Sigma)} \sqrt{2\pi}^L} \prod_{l=1}^L \exp\left\{-\frac{(z_l - \mu_l)^2}{2\sigma_l^2}\right\} \\ & \quad \log \frac{1}{\sqrt{\det(\Sigma)} \sqrt{2\pi}^L} \exp\left\{-\frac{(z_l - \mu_l)^2}{2\sigma_l^2}\right\} \end{aligned}$$

Second term of D_{KL}

$$\begin{aligned} & \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log q_{\phi}(\mathbf{z}|\mathbf{x}) \, d\mathbf{z} \\ &= \int d\mathbf{z} \frac{1}{\sqrt{\det(\Sigma)} \sqrt{2\pi}^L} \prod_{l=1}^L \exp \left\{ -\frac{(z_l - \mu_l)^2}{2\sigma_l^2} \right\} \\ & \quad \left[-\frac{L}{2} \log 2\pi - \log \sigma_l^2 - \frac{(z_l - \mu_l)^2}{2\sigma_l^2} \right] \\ &= -\frac{L}{2} \log 2\pi - \frac{1}{2} \sum_{l=1}^L (1 + \log \sigma_l^2) \end{aligned}$$

Final expression for the Kullback-Leibler divergence

KL divergence for the Gaussian case

The KL divergence in the Gaussian case is given by

$$\mathcal{D}_{\text{KL}}[q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z})] = \frac{1}{2} \sum_{l=1}^L [1 + \log \sigma_l^2(\mathbf{x}) - \sigma_l^2(\mathbf{x}) - \mu_l^2(\mathbf{x})] \quad (2)$$

We have hence obtained the **tightness** term of our loss function given by the **ELBO**.

The loss function

Loss function for the Gaussian case

The loss function for the Gaussian case is given by

$$\mathcal{L}_{\theta,\phi}(\mathbf{x}^{(i)}) = \frac{1}{2}(\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)})^2 - \frac{1}{2} \sum_{l=1}^L \left[1 + \log \sigma_l^2(\mathbf{x}^{(i)}) - \sigma_l^2(\mathbf{x}^{(i)}) - \mu_l^2(\mathbf{x}^{(i)}) \right] \quad (3)$$

In this simplified case, hence, we can get an closed explicit form for the loss function. It can be directly implemented in the neural network.

A nondeterministic component

Remember the structure of the variational autoencoder.

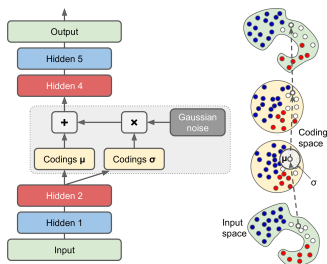


Figure: Structure of a variational autoencoder (from HOML).

The part in grey is a **nondeterministic** component, since it has a noise term. This makes backpropagation **unfeasible**! How shall we deal with it?

Framing the problem

Formally, what bothers us is that when we are backpropagating weight updates through the neural network, we want to be able to compute

$$\theta^{\text{new}} = \theta^{\text{old}} + \eta \nabla_{\theta} \text{ELBO}_{\theta, \phi}$$

$$\phi^{\text{new}} = \phi^{\text{old}} + \eta \nabla_{\phi} \text{ELBO}_{\theta, \phi}$$

where

$$\text{ELBO}_{\theta, \phi} = \langle \log p(\mathbf{x}|\mathbf{z})(p(\mathbf{z})) \rangle_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} - H[q_{\phi}(\mathbf{z}|\mathbf{x})]$$

Framing the problem

Formally, what bothers us is that when we are backpropagating weight updates through the neural network, we want to be able to compute

$$\theta^{\text{new}} = \theta^{\text{old}} + \eta \nabla_{\theta} \text{ELBO}_{\theta, \phi}$$

$$\phi^{\text{new}} = \phi^{\text{old}} + \eta \nabla_{\phi} \text{ELBO}_{\theta, \phi}$$

where

$$\text{ELBO}_{\theta, \phi} = \langle \log p(\mathbf{x}|\mathbf{z})(p(\mathbf{z})) \rangle_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} - H[q_{\phi}(\mathbf{z}|\mathbf{x})]$$

First derivation is **OK!**

Framing the problem

Formally, what bothers us is that when we are backpropagating weight updates through the neural network, we want to be able to compute

$$\theta^{\text{new}} = \theta^{\text{old}} + \eta \nabla_{\theta} \text{ELBO}_{\theta, \phi}$$

$$\phi^{\text{new}} = \phi^{\text{old}} + \eta \nabla_{\phi} \text{ELBO}_{\theta, \phi}$$

where

$$\text{ELBO}_{\theta, \phi} = \langle \log p(\mathbf{x}|\mathbf{z})(p(\mathbf{z})) \rangle_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} - H[q_{\phi}(\mathbf{z}|\mathbf{x})]$$

First derivation is **OK!**

Second derivation is **NOT OK!** Expectation value and gradient do not commute in this case. What can we do?

Let it commute!

Let us **make it commute!** In fact, in particularly simple settings such as when the approximate posteriors $q_\phi(\mathbf{z}|\mathbf{x})$ have Gaussian form

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\sigma}(\mathbf{x}))$$

the $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ can be rewritten as

$$\mathbf{z} = \boldsymbol{\mu}(\mathbf{x}) + \boldsymbol{\epsilon} \cdot \boldsymbol{\sigma}(\mathbf{x})$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$.

The reparametrization trick

Instead of sampling from the approximate variational posterior we sample ϵ from a base noise distribution and apply a **deterministic transformation**.

This amounts to replacing $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ with $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ and defining the new deterministic $\mathbf{z} = \boldsymbol{\mu} + \epsilon\sigma$. It is called **reparametrization trick**

$$\begin{aligned}\nabla_\phi \text{ELBO}_{\theta, \phi} &= \nabla_\phi \langle \log p(\mathbf{x}, \boldsymbol{\mu} + \epsilon\sigma) \rangle_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})} - \nabla_\phi H[q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \langle \nabla_\phi \log p(\mathbf{x}, \boldsymbol{\mu} + \epsilon\sigma) \rangle_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})} - \nabla_\phi H[q_\phi(\mathbf{z}|\mathbf{x})]\end{aligned}$$

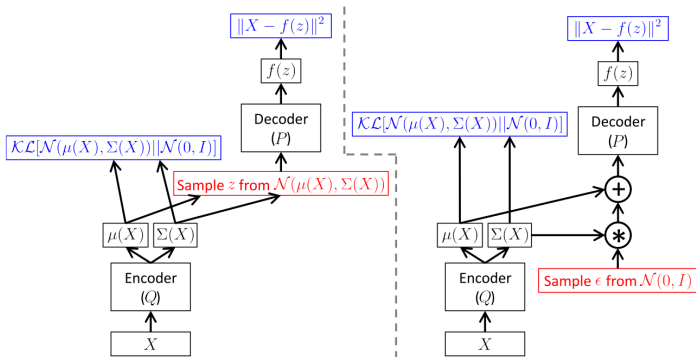


Figure: Structure of NN with and without reparametrization trick.

Structure of our network

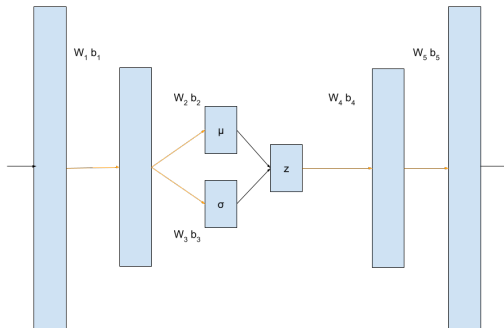


Figure: Structure of our variational autoencoder

Let's get to work

The Colab notebook is available at the following URL
https://bit.ly/pcsm1_vae