

Variational auto-encoder

Variational autoencoder ideas

The original papers

- Rezende, D. J., Mohamed, S., & Wierstra, D. (2014, June). Stochastic backpropagation and approximate inference in deep generative models. In International conference on machine learning (pp. 1278-1286). PMLR.
- Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.

Main reference

- Diederik P. Kingma and Max Welling (2019), “An Introduction to Variational Autoencoders”, Foundations and Trends R in Machine Learning:

What is does

- generate realistic samples of data,
- allow for accurate imputations of missing data,
- high-dimensional data visualisation
- Clustering

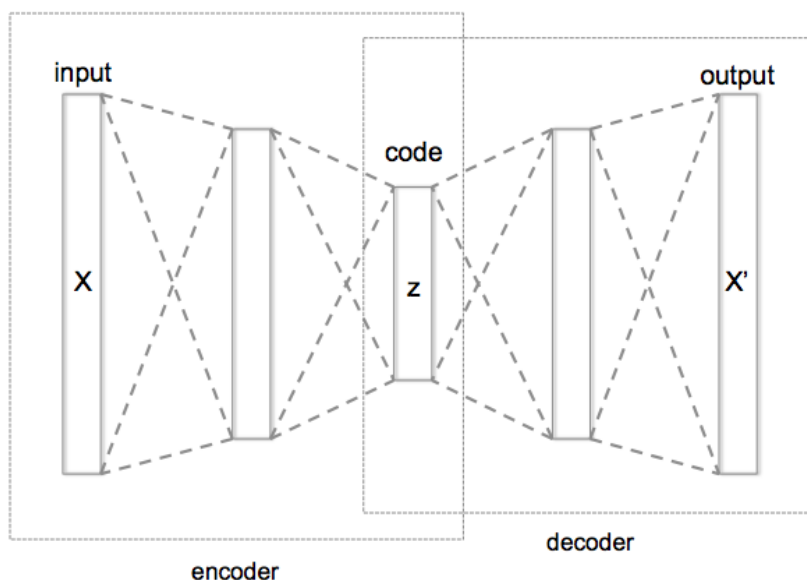
How it works

Latent variables models which marry ideas from

- approximate Bayesian inference
 - ELBO (Evidence Lower BOund)
 - Reparametrization
- deep neural networks
 - Stochastic Gradient Descent
 - Retropropagation of the Gradient

to represent an approximate posterior distribution through variational lower bound optimization

Auto-encoder Structure



Auto-encoder Structure

What is a VAE ?

Coupling of 2 parametric models

VAE is a latent variable vector \mathbf{z} and an observation $\mathbf{x} \in \mathbb{R}^D$

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

- encoder g (recognition model): $p(\mathbf{z}|\mathbf{x})$, which is approximated by $q_{\Phi}(\mathbf{z}|\mathbf{x})$
- decoder $h \approx g^{-1}$ (generative model): $p_{\Theta}(\mathbf{x}|\mathbf{z})$
- encoder and decoder could be neural networks

Optimization of ELBO via Stochastic Gradient Ascent

- The VAE ELBO approximates the likelihood of a **latent variable** model
- The Gradient computation uses the re-parametrization trick
- Each step of the gradient ascent augment the ELBO as an **EM iteration**

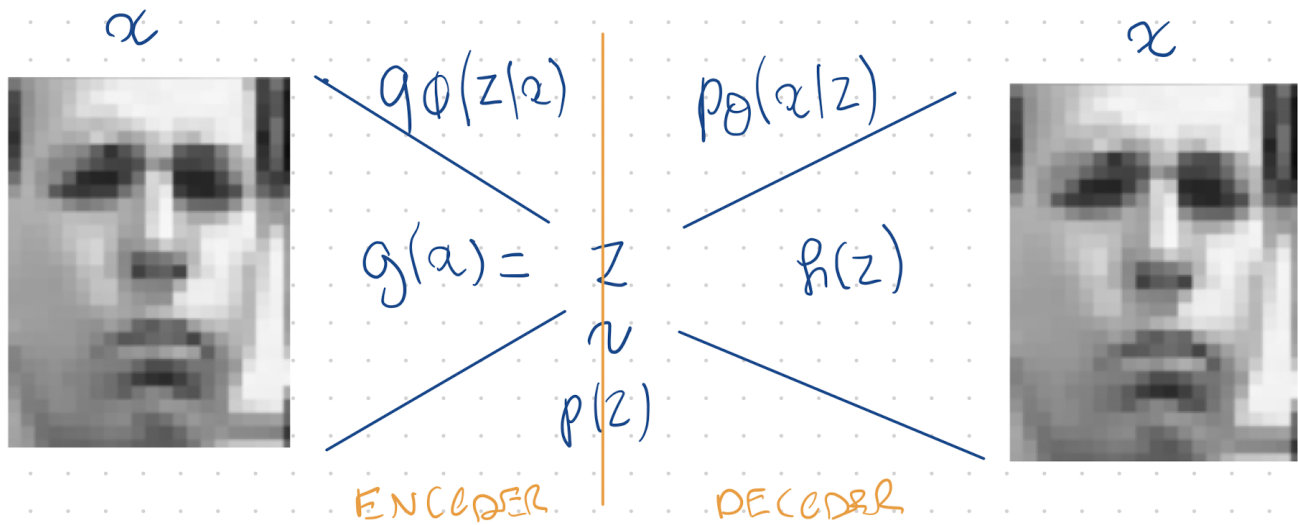
Example of use

Data

Left batch of original training set - right : random generation of images



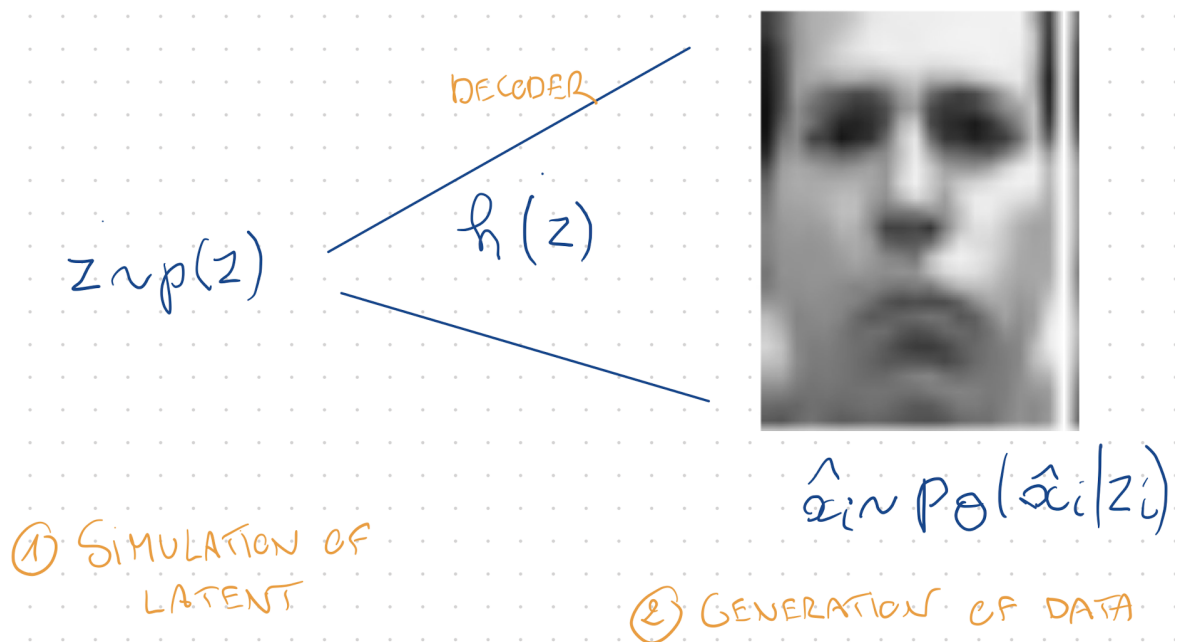
Learning from data



Frey image learning

51

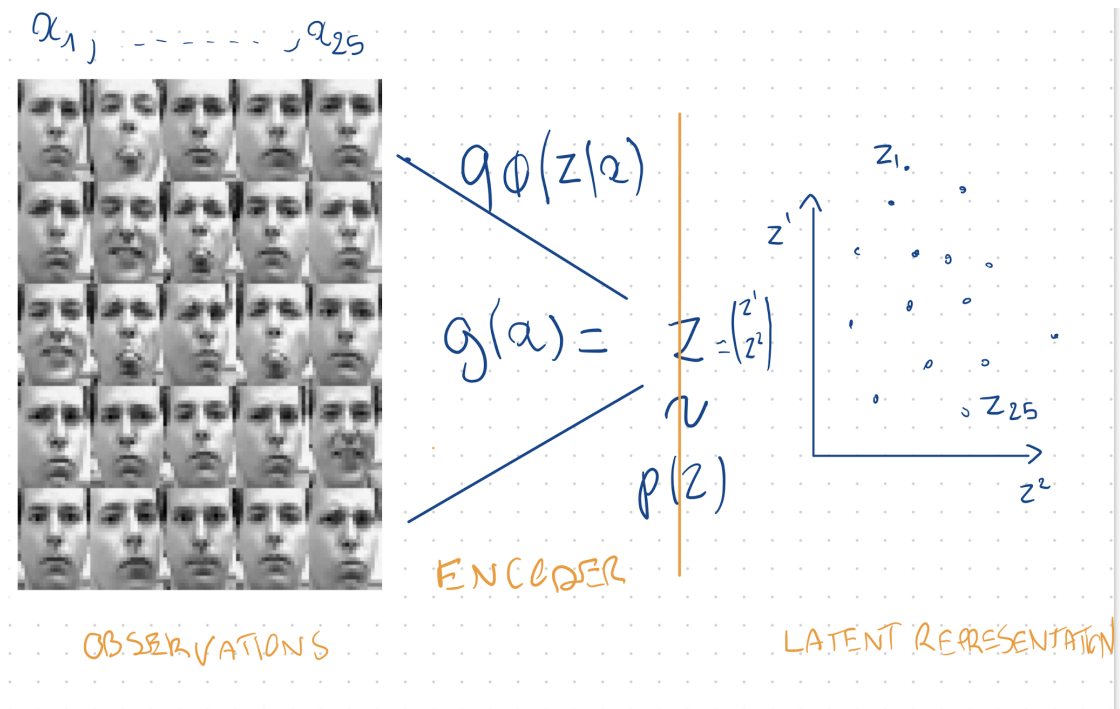
New Data generation from latent simulation



Frey image generation

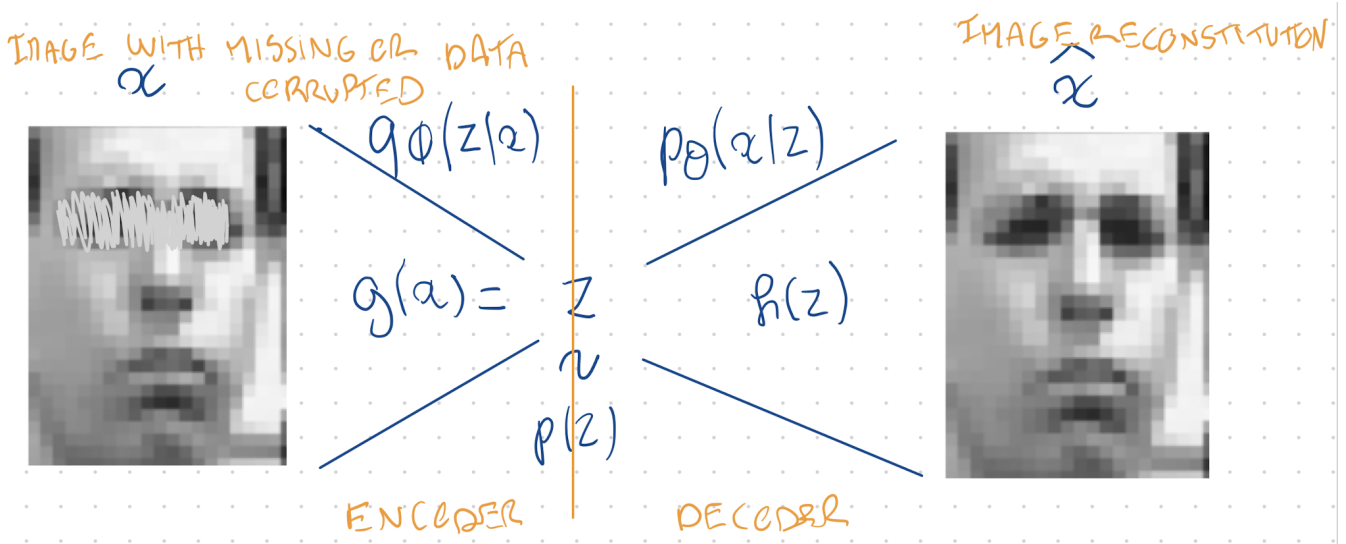
52

Data representation in latent space



Frey image representation

Missing data imputation after learning



Frey image pixel imputation

Ingredient: Parameterization of conditional distributions with Neural Networks

56

Modeling joint distribution

- \mathbf{x} : Observed random variables
- $p^*(\mathbf{x})$: underlying unknown distribution
- $p_\theta(\mathbf{x})$: model distribution
- Goal: $p_\theta(\mathbf{x}) \approx p^*(\mathbf{x})$

We wish flexible $p_\theta(\mathbf{x})$

57

Modeling Conditional distribution

Classification and regression

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

58

Parameterization of conditional distributions with Neural Networks

Classification

$$\theta = \text{NeuralNet}(\mathbf{x})$$

$$p_{\theta}(y|\mathbf{x}) = \text{Categorical}(y, \theta)$$

59

Ingredient: Stochastic Gradient

61

What differences between Oja's rule and VAE

What is a VAE ?

Coupling of 2 **parametric models**

- decoder (generative model): $\mathbf{x} = h_{\Theta}(\mathbf{z})$ where $p_{\Theta}(\mathbf{x}|\mathbf{z})$
- encoder (recognition model): $\mathbf{z} = g_{\Phi}(\mathbf{x})$ where $p(\mathbf{z}|\mathbf{x})$ is approximated by $q_{\Phi}(\mathbf{z}|\mathbf{x})$

In Oja's rule

62

Factor analysis generalizes Oja and is closer to a VAE

Factor analysis considers the observation $\mathbf{x} \in \mathbb{R}^D$

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

where

- the noise $\boldsymbol{\epsilon} \sim \mathcal{N}_D(\mathbf{0}, \boldsymbol{\Psi})$
- the hidden (latent) vector $\mathbf{z} \sim \mathcal{N}_L(\mathbf{0}, \mathbf{I}_L)$

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

the mean is a linear function of the (hidden) inputs

- \mathbf{W} is a $D \times L$ matrix, known as the factor loading matrix,
- $\boldsymbol{\Psi}$ is a $D \times D$ covariance matrix that we take to be diagonal

The special case in which $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ is called probabilistic principal components analysis or PPCA.

Ingredient : Evidence Lower BOund Minimization

65

Missing data

In a missing data framework the log-likelihood of the parameters is advantageously expressed as

$$\log P_{\Theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{z}|\mathbf{x}} \left[\log \frac{P(\mathbf{x}, \mathbf{z})}{P(\mathbf{z}|\mathbf{x})} \right]$$

Approximation

When the distribution of $\mathbf{z}|\mathbf{x}$ is intractable, an **approximation** $q_{\Phi}(\mathbf{z}|\mathbf{x})$ is used

$q_{\Phi}(\mathbf{z}|\mathbf{x})$ is the inverse function for $p_{\Theta}(\mathbf{x}|\mathbf{z})$ in a **Bayes sense**

66

ELBO

$$\begin{aligned} \log P_{\Theta}(\mathbf{x}) &= \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{P_{\Theta}(\mathbf{z}|\mathbf{x})} \right] \\ &= \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z}) q_{\Phi}(\mathbf{z}|\mathbf{x})}{P_{\Theta}(\mathbf{z}|\mathbf{x}) q_{\Phi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{q_{\Phi}(\mathbf{z}|\mathbf{x})} \right]}_{ELBO} + \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\Phi}(\mathbf{z}|\mathbf{x})}{p_{\Theta}(\mathbf{z}|\mathbf{x})} \right]}_{D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}|\mathbf{x}))} \end{aligned}$$

where

$$D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}|\mathbf{x})) = -\mathbb{E}_q[\log \frac{p}{q}] \geq -\log \mathbb{E}_q[\frac{p}{q}] \geq 0$$

from Jensen

67

Two for one

- VAE finds parameters which approximately maximize the marginal likelihood $P_{\Theta}(\mathbf{x})$ (good generative function)
- VAE finds the approximation of the recognition model which minimizes the KL divergence

68

Alternative formulation of ELBO

ELBO can be rewritten as

$$ELBO = E_{q_{\Phi}(\mathbf{z}|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}))$$

For a suited choice of $p(\mathbf{z})$ and $q(\mathbf{z}|\mathbf{x})$, $D_{KL}(q_{\Phi}(\mathbf{z}|\mathbf{x})|p_{\Theta}(\mathbf{z}))$ can be calculated in closed form.

Exercices

1. Show that the ELBO can be rewritten as above
2. Compute the KL divergence between two multivariate Gaussians

69

The ELBO is maximized by Stochastic Gradient

Let \mathbf{X} be a i.i.d sample of random vector \mathbf{x}

$$ELBO(\mathbf{X}) = \sum_{\mathbf{x} \in \mathbf{X}} \mathcal{L}(\Theta, \Phi; \mathbf{x})$$

Gradient

The Gradient can be separated into 2 parts

1. $\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x})$
2. $\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

70

Decoder Gradient: $\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

Given some usually verified conditions and a Monte Carlo Approximation

$$\begin{aligned}\nabla_{\Theta} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Theta} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} \left[\log \frac{P_{\Theta}(\mathbf{x}, z)}{q_{\Phi}(z|\mathbf{x})} \right] \\ &= \nabla_{\Theta} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}, z)] \\ &= \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\nabla_{\Theta} \log P_{\Theta}(\mathbf{x}, z)] \\ &\approx \nabla_{\Theta} \log P_{\Theta}(\mathbf{x}, z)\end{aligned}$$

71

Encoder Gradient: $\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x})$

Encoder Gradient is more difficult to compute since in general

$$\begin{aligned}\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Phi} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} \left[\log \frac{P_{\Theta}(\mathbf{x}, z)}{q_{\Phi}(z|\mathbf{x})} \right] \\ &= \nabla_{\Phi} \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\log P_{\Theta}(\mathbf{x}, z) - \log q_{\Phi}(z|\mathbf{x})] \\ &\neq \mathbb{E}_{q_{\Phi}(z|\mathbf{x})} [\nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, z) - \nabla_{\Phi} \log q_{\Phi}(z|\mathbf{x})]\end{aligned}$$

A reparametrization (variable change) trick allows a workaround

72

Third ingredient: Encoder approximation using Reparametrization and Monte Carlo

74

Encoder function

Let us rewrite the decoder function with a random vector ϵ whose distribution is not parametrized by Φ :

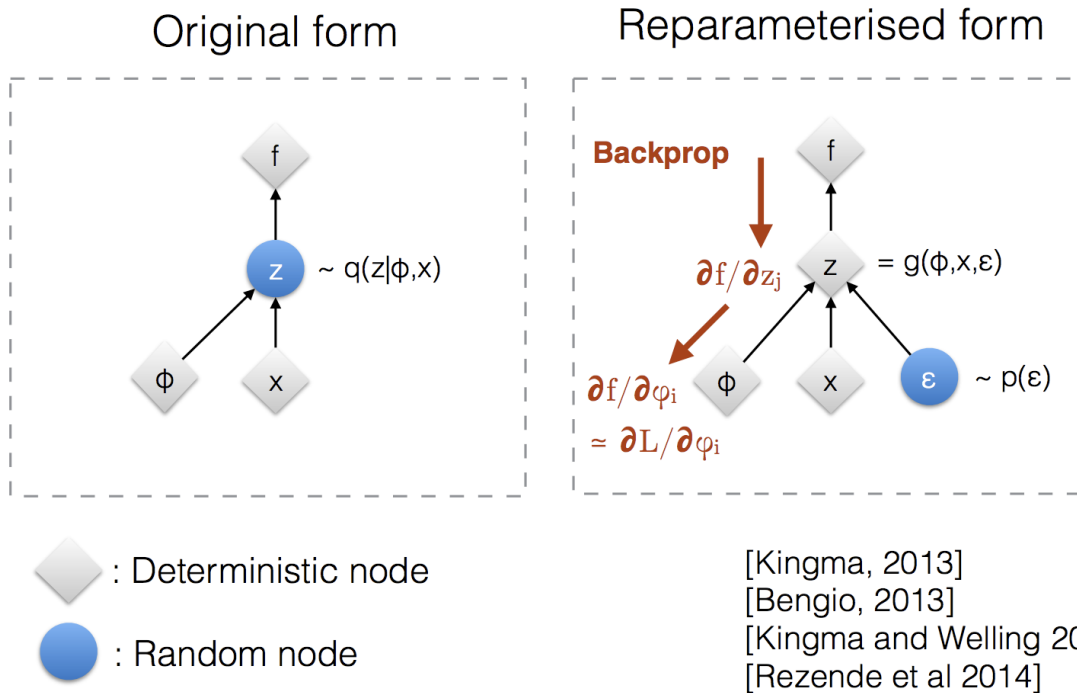
$$\mathbf{z} = g(\mathbf{x}, \epsilon; \Phi)$$

$$\begin{aligned}\nabla_{\Phi} \mathcal{L}(\Theta, \Phi; \mathbf{x}) &= \nabla_{\Phi} \mathbb{E}_{p(\epsilon)} \left[\log \frac{P_{\Theta}(\mathbf{x}, \mathbf{z})}{q_{\Phi}(\mathbf{z}|\mathbf{x})} \right] \\ &= \nabla_{\Phi} \mathbb{E}_{p(\epsilon)} [\log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \log q_{\Phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \nabla_{\Phi} \log q_{\Phi}(\mathbf{z}|\mathbf{x})] \\ &\approx \nabla_{\Phi} \log P_{\Theta}(\mathbf{x}, \mathbf{z}) - \nabla_{\Phi} \log q_{\Phi}(\mathbf{z}|\mathbf{x})\end{aligned}$$

We just have to compute $\log q_{\Phi}(\mathbf{z}|\mathbf{x})$ after the change of variable

75

Reparametrization trick



76

Computing $\log q_{\Phi}(\mathbf{z}|\mathbf{x})$ with a change of variable

$$\log q_{\Phi}(\mathbf{z}|\mathbf{x}) = \log p(\boldsymbol{\epsilon}) - \log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon})$$

where the second term is the log of the absolute value of the determinant of the Jacobian matrix:

$$\log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon}) = \log \left| \det \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \cdots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \cdots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \right|$$

77

Factorized Gaussian Posterior

Model

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = \text{EncoderNN}_{\Phi}(\mathbf{x})$$

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \prod_i q_{\Phi}(z_i|\mathbf{x}) = \prod_i \mathcal{N}(z_i | \text{EncoderNN}_{\Phi}(\mathbf{x})) = \mathbf{I}$$

Reparametrization

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

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$

78

Factorized Gaussian Posterior

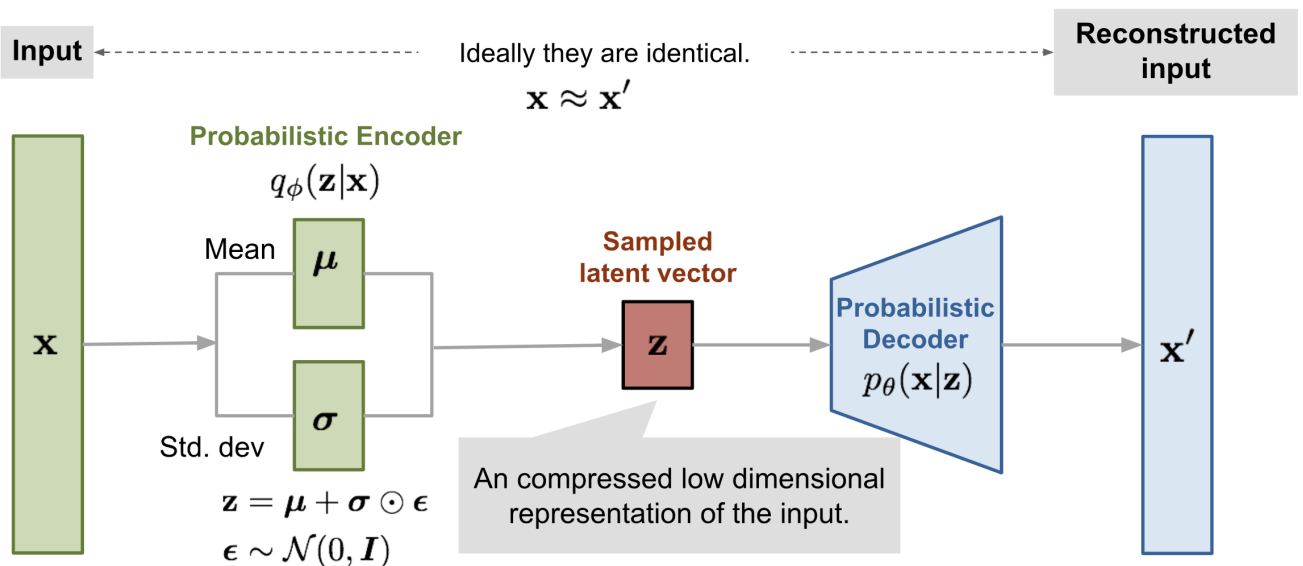
The Jacobian of the transformation is

$$\log d_{\Phi}(\mathbf{x}, \boldsymbol{\epsilon}) = \log \left| \det \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \cdots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \cdots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix} \right| = \log \prod_i \sigma_i$$

The log posterior density is

$$\begin{aligned} \log q_{\Phi}(\mathbf{z}|\mathbf{x}) &= \log p(\boldsymbol{\epsilon}) - \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \\ &= \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log \sigma_i \end{aligned}$$

Reparametrization trick



Full Gaussian posterior

Model

$$q_{\Phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Reparametrization

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

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}$$

Where \mathbf{L} is a lower triangular matrix obtained from a Cholesky decomposition of $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$

81

Full Gaussian posterior

The Jacobian has a simple form

$$\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \mathbf{L}$$

As the determinant of a triangular matrix is the product of its diagonal terms,

$$\begin{aligned} \log q_{\Phi}(\mathbf{z}|\mathbf{x}) &= \log p(\boldsymbol{\epsilon}) - \log \left| \det \left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| \\ &= \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log L_{ii} \end{aligned}$$

82

Variational Auto Encoder in short

84

Algorithm input and output

Input:

- \mathbf{X} : Dataset
- $h(\mathbf{z})$ decoding function, with dist. $p_{\Theta}(\mathbf{x}, \mathbf{z})$
- $g(\mathbf{x})$ encoding function with dist. $q_{\Psi}(\mathbf{z}|\mathbf{x})$

Output:

- Θ
- Φ

85

Algorithm

Initialisation of Θ and Φ

While SGD not converged do

- Draw a random minibatch $\mathbf{X}^M \in \mathbf{X}$
- $\epsilon \sim p(\epsilon)$ (Random noise for every datapoint in \mathbf{X}^M)
- Compute

→

$$\tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \epsilon) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \left(\log p_{\Theta}(\mathbf{x}, \mathbf{z}) - \underbrace{\log q_{\Phi}}_{\log p(\epsilon) - \log} \right)$$

and

→ its gradients

$$\Rightarrow \nabla_{\Theta} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \epsilon) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \frac{\partial \log p_{\Theta}(\mathbf{x}, \mathbf{z})}{\partial \Theta}$$

$$\Rightarrow \nabla_{\Phi} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \epsilon) = \frac{1}{m} \sum_{\mathbf{x} \in \mathbf{X}^M} \frac{\log \partial d_{\Phi}(\mathbf{x}, \epsilon)}{\partial \Phi}$$

$$\left(\Theta \right) \quad \left(\Theta \right) \quad \left(\nabla_{\Theta} \tilde{\mathcal{L}}_{\Theta, \Phi}(\mathbf{X}^M, \epsilon) \right)$$

Original example from Kingma: Gaussian model with MLP parametrization

88

Multivariate Gaussian decoder with a diagonal covariance structure

Decoder $p_{\Theta}(\mathbf{x}|\mathbf{z})$ or decoder (just swap \mathbf{x} and \mathbf{z}) are assumed to have multivariate Gaussian dist. with a diagonal covariance structure:

Decoder

- $\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 I)$ where $\boldsymbol{\mu} = W_4 \mathbf{h} + b_4$
- $\log \sigma^2 = W_5 \mathbf{h} + b_5$
- $\mathbf{h} = \tanh(W_3 \mathbf{z} + b_3)$

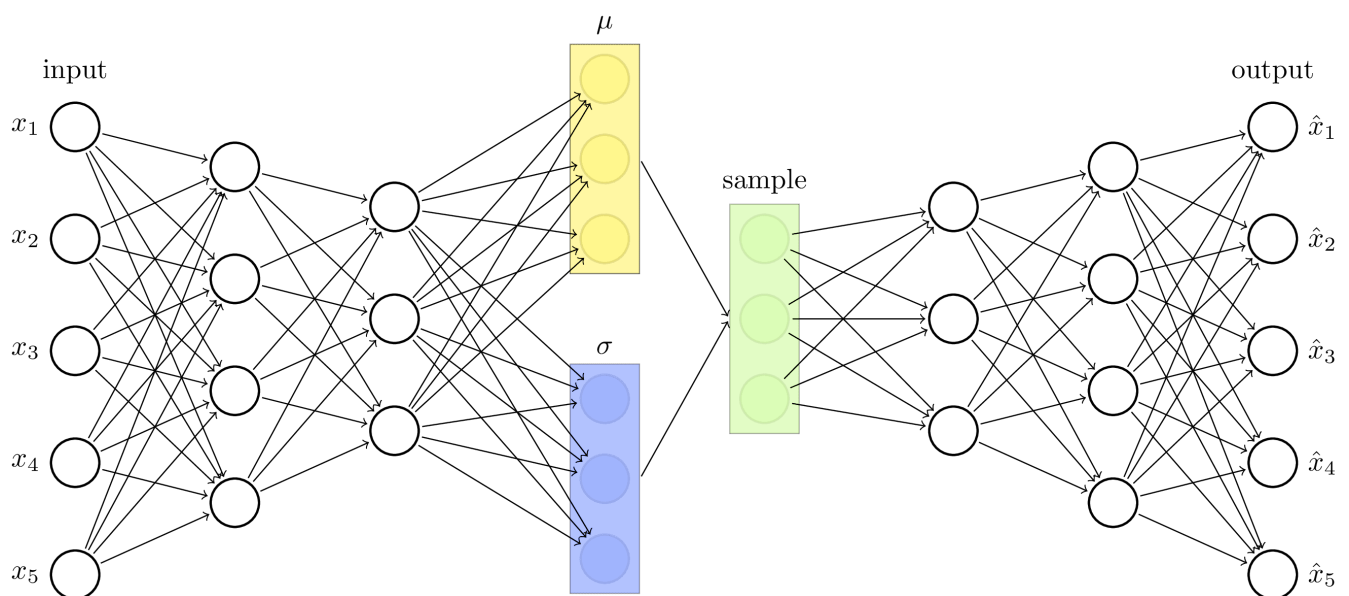
where $\{W_3, W_4, W_5, b_3, b_4, b_5\}$ are the weights and biases of the MLP and part of Θ when used as decoder.

Encoder

- Let us consider the prior $p_{\Theta}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$
- swap \mathbf{x} and \mathbf{z} in the decoder above to get $q_{\Phi}(\mathbf{z}|\mathbf{x})$

89

Gaussian VAE illustrated



90