

Auto-encoder and VAE

Thomas Gerald

April 16, 2026

Laboratoire Interdisciplinaire des Sciences du Numérique – LISN, CNRS

Review of some Neural Networks architecture

- Auto-encoder (semi supervised methods)
- Variational Auto-encoder (VAE) and the ELBO objective - **in depth**
- Generative Adversarial Networks (GAN) - **principle only**
- Diffusion Models - **principle only**

Auto-Encoder

Auto-encoder (AE)

Auto-encoder (AE) What is it ?

Objective: Learn a method for compressing data minimizing the loss.

Let be $x \in \mathbb{R}^n$ and $g \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ we would minimize

$$\sum_{x_i \in \mathcal{D}} \|g(f(x_i)) - x_i\|$$

With

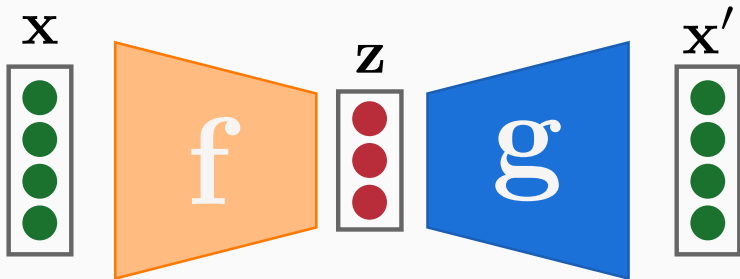
- $f : \mathbb{R}^n \rightarrow \mathbb{R}^d$
- $g : \mathbb{R}^d \rightarrow \mathbb{R}^n$

→ Most of the time $d \ll n$

What architecture for g and f ?

- Linear function
- MLP (Multi Layer Perceptron)
- Convolutional Model

Auto-encoder (AE)



Auto-encoder

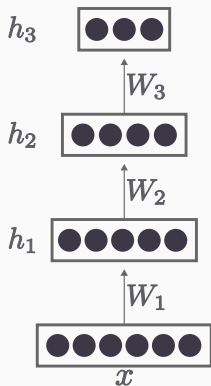
- f encode information
- g decode information

If the latent space (z space) is in low dimension we compressed the data efficiently (if the optimized cost is low enough)

Auto-encoder (AE)

In early Deep Learning

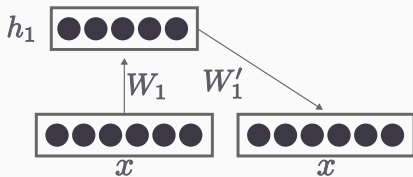
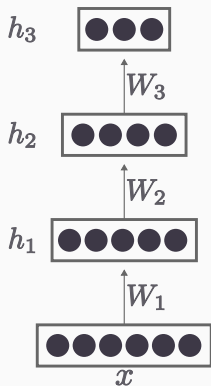
→ Used to initialize Deep Neural Networks



Auto-encoder (AE)

In early Deep Learning

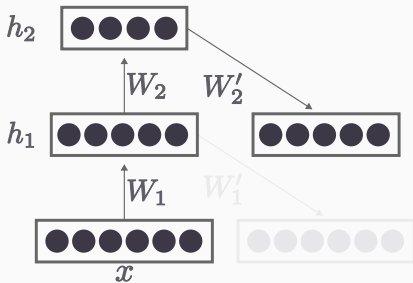
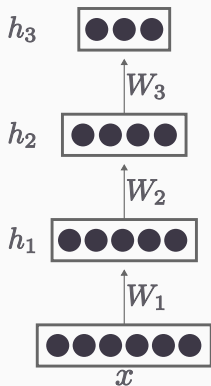
→ Used to initialize Deep Neural Networks



Auto-encoder (AE)

In early Deep Learning

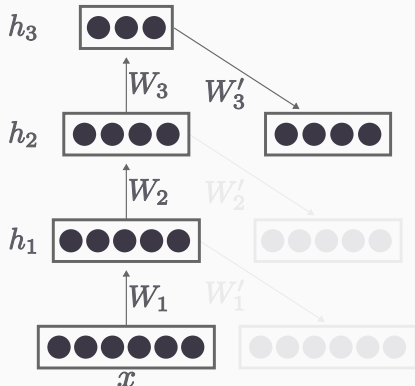
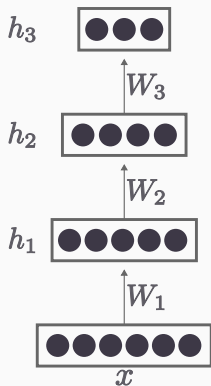
→ Used to initialize Deep Neural Networks



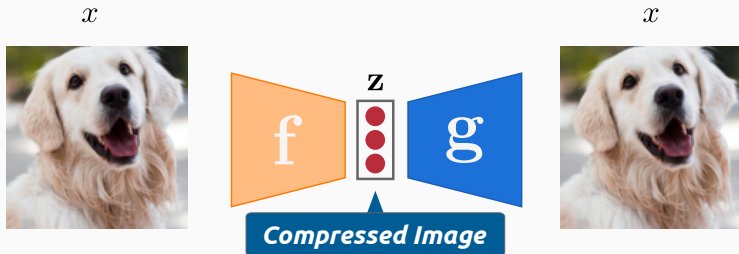
Auto-encoder (AE)

In early Deep Learning

→ Used to initialize Deep Neural Networks



Auto-encoder (AE)



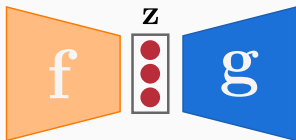
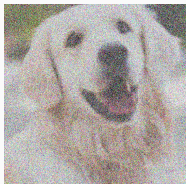
Auto-encoder

- f encode information
- g decode information

If the latent space (z space) is in low dimension we compressed the data efficiently (if the optimized cost is low enough)

Auto-encoder (AE)

$$x + \mathcal{N}(\mu, \sigma^2)$$



$$x$$

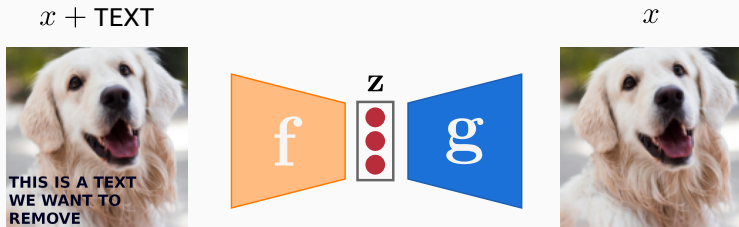


Auto-encoder

- f encode information
- g decode information

If the latent space (z space) is in low dimension we compressed the data efficiently (if the optimized cost is low enough)

Auto-encoder (AE)

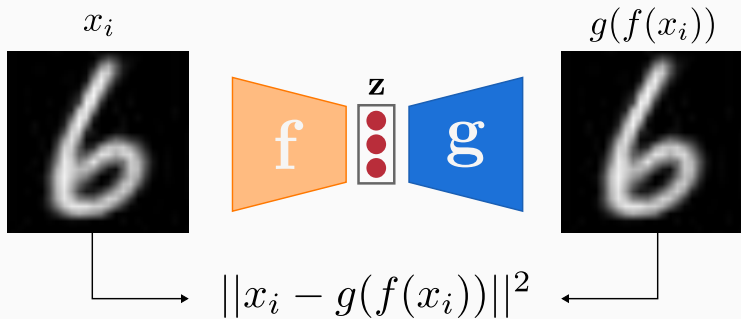


Auto-encoder

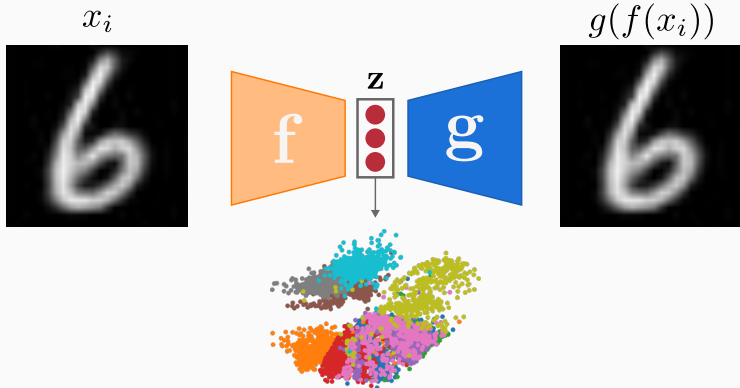
- f encode information
- g decode information

If the latent space (z space) is in low dimension we compressed the data efficiently (if the optimized cost is low enough)

Auto-encoder (AE): Latent representation



Auto-encoder (AE): Latent representation



- Using latent representation for classification
- Can we do something else on the latent space?

The latent space:

The latent space:

- Latent vectors encoding similar content close to each other?

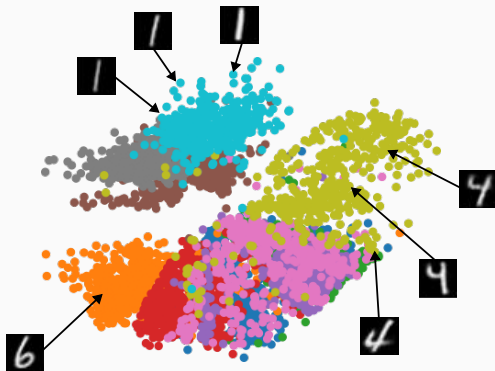
The latent space:

- Latent vectors encoding similar content close to each other?
- Can we sample in the latent space?

Auto-encoder (AE): Latent representation

The latent space:

- Latent vectors encoding similar content close to each other?
- Can we sample in the latent space?

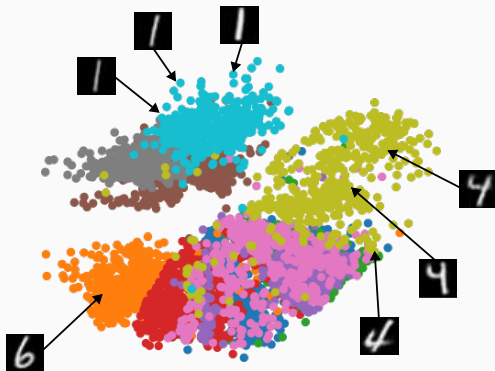


Generating content

Auto-encoder (AE): Latent representation

The latent space:

- Latent vectors encoding similar content close to each other?
- Can we sample in the latent space?



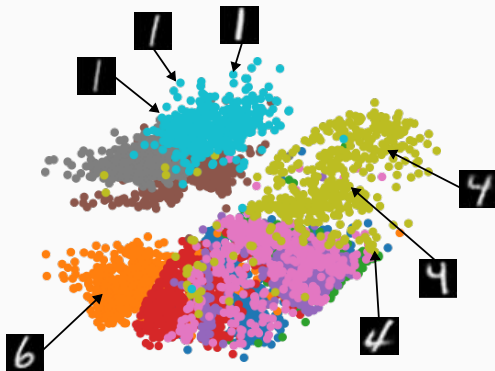
Generating content

- We do not know what is the underlying latent distribution $p_{\theta}(z)$

Auto-encoder (AE): Latent representation

The latent space:

- Latent vectors encoding similar content close to each other?
- Can we sample in the latent space?



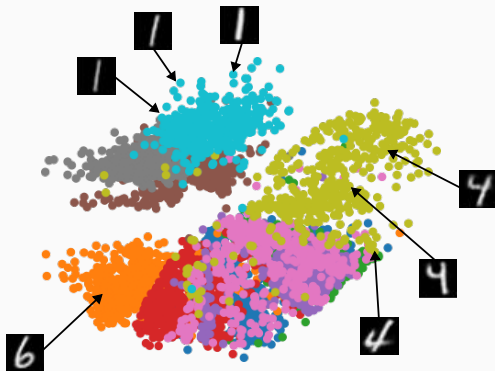
Generating content

- We do not know what is the underlying latent distribution $p_{\theta}(z)$
- Can we estimate this distribution ?

Auto-encoder (AE): Latent representation

The latent space:

- Latent vectors encoding similar content close to each other?
- Can we sample in the latent space?

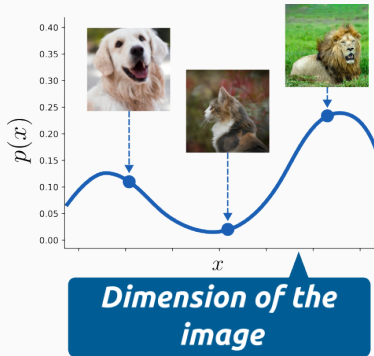


Generating content

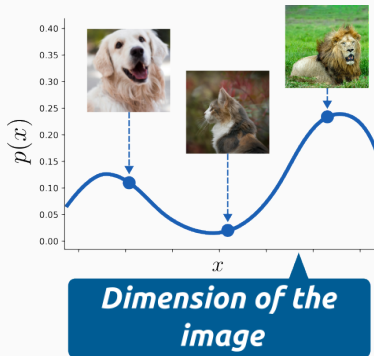
- We do not know what is the underlying latent distribution $p_{\theta}(z)$
- Can we estimate this distribution ?
- Can we force the latent variable to fit a chosen distribution ?

Generation

Generate data

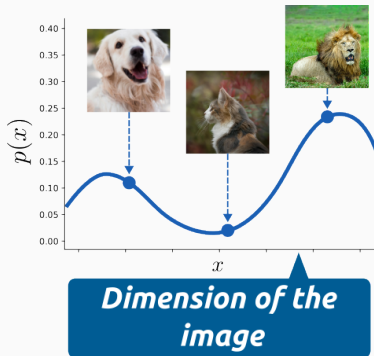


Generate data



There is an hidden distribution $p(x)$ over the data

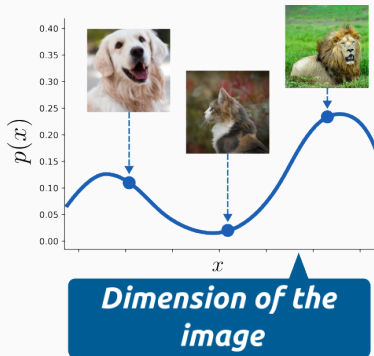
Generate data



There is an hidden distribution $p(x)$ over the data

- Can we use neural network ?

Generate data

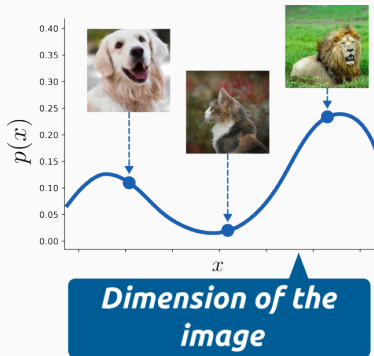


There is an hidden distribution $p(x)$ over the data

- Can we use neural network ?

$$\rightarrow p_{\theta}(x) \sim p(x)$$

Generate data



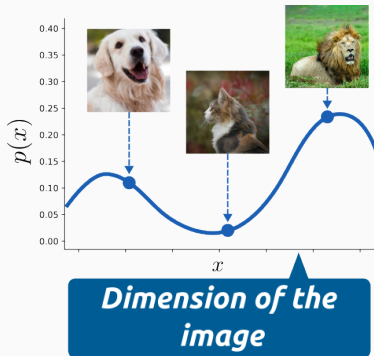
There is an hidden distribution $p(x)$ over the data

- Can we use neural network ?

$$\rightarrow p_{\theta}(x) \sim p(x)$$

- Can we sample from PDF $\rightarrow p_{\theta}(x)$?

Generate data



There is an hidden distribution $p(x)$ over the data

- Can we use neural network ?

$$\rightarrow p_{\theta}(x) \sim p(x)$$

- Can we sample from PDF $\rightarrow p_{\theta}(x)$?

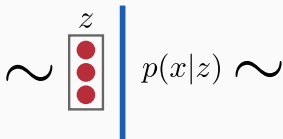
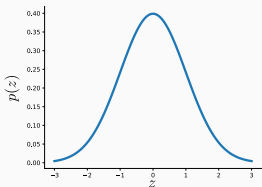
\rightarrow Difficult to optimize and sample...

Generate data using intermediate variable ?

Let introduce z a random variable where

- We know the distribution on z $p_{\theta}(z)$ (your choice)
- x depends on z (latent variable)

1. z_i is generated from a **prior distribution** $p_{\theta}(z)$
2. x_i is generated from the conditional probability $p_{\theta}(x|z)$

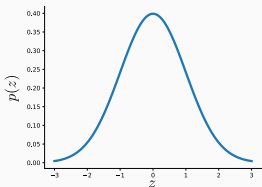


Generate data using intermediate variable ?

Let introduce z a random variable where

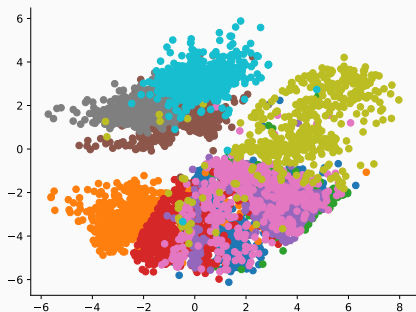
- We know the distribution on z $p_{\theta}(z)$ (your choice)
- x depends on z (latent variable)

1. z_i is generated from a **prior distribution** $p_{\theta}(z)$
2. x_i is generated from the conditional probability $p_{\theta}(x|z)$



\sim





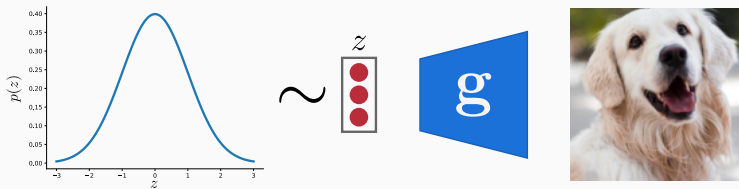
Estimate latent distribution

- What is $p(z)$?
- What prior?

Difficult to choose the probability function (Gaussian ?)

- Ensuring latent representation follow a prior distribution?
- Using a regularisation approach to force data follow a density function?

A first problem ?



→ Parameters (θ) is unknown

→ Latent variable is unknown

Objective:

Find θ^* such that we can generate x from z :

→ Maximize log-likelihood

$$\arg \max_{\theta} \mathbb{E}_{x \sim p(x)} [\log(p_{\theta}(x))]$$

What is $p_\theta(x)$?

$$\begin{aligned} p_\theta(x) &= \int_z p_\theta(x, z) dz \\ &= \int_z p_\theta(x|z) \cdot p_\theta(z) dz \\ &= \mathbb{E}_{z \sim p_\theta(z)} p_\theta(x|z) \end{aligned}$$

Additionally:

$$\begin{aligned} p(z|x) &= \frac{p(x|z)p(z)}{p(x)} \\ p(x|z) &= \frac{p(z|x)p(x)}{p(z)} \end{aligned}$$

→ Intractable (or very difficult to estimate)

Vocabulary

- $p(z|x)$ the posterior probability → we want to estimate

What is $p_{\theta}(x)$?

$$\begin{aligned} p_{\theta}(x) &= \int_z p_{\theta}(x, z) dz \\ &= \int_z p_{\theta}(x|z) \cdot p_{\theta}(z) dz \\ &= \mathbb{E}_{z \sim p_{\theta}(z)} p_{\theta}(x|z) \end{aligned}$$

Additionally:

$$\begin{aligned} p(z|x) &= \frac{p(x|z)p(z)}{p(x)} \\ p(x|z) &= \frac{p(z|x)p(x)}{p(z)} \end{aligned}$$

→ Intractable (or very difficult to estimate)

Vocabulary

- $p(z|x)$ the posterior probability → we want to estimate
- $p(z)$ the prior → we choose one (for instance gaussian)

What is $p_\theta(x)$?

$$\begin{aligned} p_\theta(x) &= \int_z p_\theta(x, z) dz \\ &= \int_z p_\theta(x|z) \cdot p_\theta(z) dz \\ &= \mathbb{E}_{z \sim p_\theta(z)} p_\theta(x|z) \end{aligned}$$

Additionally:

$$\begin{aligned} p(z|x) &= \frac{p(x|z)p(z)}{p(x)} \\ p(x|z) &= \frac{p(z|x)p(x)}{p(z)} \end{aligned}$$

→ Intractable (or very difficult to estimate)

Vocabulary

- $p(z|x)$ the posterior probability → we want to estimate
- $p(z)$ the prior → we choose one (for instance gaussian)
- $p(x|z)$ the likelihood → approximable (reconstruction)

What is $p_{\theta}(x)$?

$$\begin{aligned} p_{\theta}(x) &= \int_z p_{\theta}(x, z) dz \\ &= \int_z p_{\theta}(x|z) \cdot p_{\theta}(z) dz \\ &= \mathbb{E}_{z \sim p_{\theta}(z)} p_{\theta}(x|z) \end{aligned}$$

Additionally:

$$\begin{aligned} p(z|x) &= \frac{p(x|z)p(z)}{p(x)} \\ p(x|z) &= \frac{p(z|x)p(x)}{p(z)} \end{aligned}$$

→ Intractable (or very difficult to estimate)

Vocabulary

- $p(z|x)$ the posterior probability → we want to estimate
- $p(z)$ the prior → we choose one (for instance gaussian)
- $p(x|z)$ the likelihood → approximable (reconstruction)
- $p(x)$ the marginal probability/or evidence probability

What is $p_\theta(x)$?

$$\begin{aligned} p_\theta(x) &= \int_z p_\theta(x, z) dz \\ &= \int_z p_\theta(x|z) \cdot p_\theta(z) dz \\ &= \mathbb{E}_{z \sim p_\theta(z)} p_\theta(x|z) \end{aligned}$$

Additionally:

$$\begin{aligned} p(z|x) &= \frac{p(x|z)p(z)}{p(x)} \\ p(x|z) &= \frac{p(z|x)p(x)}{p(z)} \end{aligned}$$

→ Intractable (or very difficult to estimate)

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution, but sampling space too large)
- **Variational inference:** Choose an other distribution $q_\phi(z|x)$ to approximate $p_\theta(z|x)$

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution)
- **Variational inference:** choose an other distribution $q_\phi(z|x)$ to approximate $p_\theta(z|x)$

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution)
- **Variational inference:** choose an other distribution $q_\phi(z|x)$ to approximate $p_\theta(z|x)$

Variational Inference

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution)
- **Variational inference:** choose an other distribution $q_{\phi}(z|x)$ to approximate $p_{\theta}(z|x)$

Variational Inference

- Let consider $q_{\phi}(z|x)$ parametrized by ϕ (a neural network)

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution)
- **Variational inference:** choose an other distribution $q_\phi(z|x)$ to approximate $p_\theta(z|x)$

Variational Inference

- Let consider $q_\phi(z|x)$ parametrized by ϕ (a neural network)
- We would like $q_\phi(z|x)$ close to $p_\theta(z|x)$

A solution?

- Approximate using Monte-Carlo estimation (sampling to estimate distribution)
- **Variational inference:** choose an other distribution $q_\phi(z|x)$ to approximate $p_\theta(z|x)$

Variational Inference

- Let consider $q_\phi(z|x)$ parametrized by ϕ (a neural network)
- We would like $q_\phi(z|x)$ close to $p_\theta(z|x)$

→ Minimize a distance between $q_\phi(z|x)$ and $p_\theta(z|x)$

Kullback Leibler divergence

The KL-divergence is a metric between two distribution where the KL-divergence of two distribution is given by:

$$KL(p(x)||q(x)) = \int_x p(x) \log \frac{p(x)}{q(x)}$$

- When distributions are the same $KL(p(x)||q(x)) = 0$
- When distributions are not the same $0 < KL(p(x)||q(x))$

In our case we want to **minimize** $KL(q_\phi(z|x)||p(z|x))$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$KL(q_\phi(z|x)||p_\theta(z|x))$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \int_z q_\phi(z|x) \log p_\theta(x) \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \int_z q_\phi(z|x) \log p_\theta(x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \int_z q_\phi(z|x) \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \int_z q_\phi(z|x) \log p_\theta(x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \int_z q_\phi(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \end{aligned}$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \int_z q_\phi(z|x) \log p_\theta(x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \int_z q_\phi(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \end{aligned}$$

Rearranging terms we get :

$$\log p_\theta(x) = KL(q_\phi(z|x)||p_\theta(z|x)) - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z,x)]$$

Variational Auto-Encoder: $KL(q_\phi(z|x)||p_\theta(z|x))$

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ = & \int_z q_\phi(z|x) \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \\ = & \int_z q_\phi(z|x) [\log q_\phi(z|x) - \log(p_\theta(z|x))] \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log \frac{p_\theta(z,x)}{p_\theta(x)} \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \int_z q_\phi(z|x) \log p_\theta(x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \int_z q_\phi(z|x) \\ = & \int_z q_\phi(z|x) \log q_\phi(z|x) - \int_z q_\phi(z|x) \log p_\theta(z,x) + \log p_\theta(x) \end{aligned}$$

Rearranging terms we get :

$$\log p_\theta(x) = KL(q_\phi(z|x)||p_\theta(z|x)) - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z,x)]$$

**Maximizing the log-likelihood can be done maximizing the right expression
!!!**

Variational Auto-Encoder: The evidence lower bound

$$\log p_{\theta}(x) = KL(q_{\phi}(z|x)||p_{\theta}(z|x)) - \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log q_{\phi}(z|x)] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(z, x)]$$

Yes but... we don't know (or can estimate) $p_{\theta}(z|x)$

$KL(q_{\phi}(z|x)||p_{\theta}(z|x))$ is intractable however:

→ The Kullback Leibler divergence is positive

Variational Auto-Encoder: The evidence lower bound

$$\log p_{\theta}(x) = KL(q_{\phi}(z|x)||p_{\theta}(z|x)) - \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log q_{\phi}(z|x)] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(z, x)]$$

Yes but... we don't know (or can estimate) $p_{\theta}(z|x)$

$KL(q_{\phi}(z|x)||p_{\theta}(z|x))$ is intractable however:

→ The Kullback Leibler divergence is positive

$$\text{Thus } \log p_{\theta}(x) \geq - \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log q_{\phi}(z|x)] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(z, x)]$$

$$\text{Maximize } \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(z, x)] - \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log q_{\phi}(z|x)]$$

↓

$$\text{Maximize } \log p_{\theta}(x)$$

Variational Auto-Encoder objective

$$\text{Maximize } \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

This lower bound is called the Evidence Lower Bound (**ELBO**)

Variational Auto-Encoder objective

$$\text{Maximize } \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

This lower bound is called the Evidence Lower Bound (**ELBO**)

- $p_\theta(x)$ is often call the evidence (probability of input data)
- $\log p_\theta(x) \leq p_\theta(x)$
- The ELBO is a lower bound of the evidence

Variational Auto-Encoder objective

$$\text{Maximize } \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

This lower bound is called the Evidence Lower Bound (**ELBO**)

- $p_\theta(x)$ is often call the evidence (probability of input data)
- $\log p_\theta(x) \leq p_\theta(x)$
- The ELBO is a lower bound of the evidence

We have not finish yet:

Variational Auto-Encoder objective

$$\text{Maximize } \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

This lower bound is called the Evidence Lower Bound (**ELBO**)

- $p_\theta(x)$ is often call the evidence (probability of input data)
- $\log p_\theta(x) \leq p_\theta(x)$
- The ELBO is a lower bound of the evidence

We have not finish yet:

- How to estimate $\log p_\theta(z, x)$?

Variational Auto-Encoder objective

$$\text{Maximize } \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

This lower bound is called the Evidence Lower Bound (**ELBO**)

- $p_\theta(x)$ is often call the evidence (probability of input data)
- $\log p_\theta(x) \leq p_\theta(x)$
- The ELBO is a lower bound of the evidence

We have not finish yet:

- How to estimate $\log p_\theta(z, x)$?
- Estimate update of $q_\phi(z|x)$ (especially because we sample on $q_\phi(z|x)$)

Variational Auto-Encoder: The evidence lower bound

$$ELBO = \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)]$$

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \end{aligned}$$

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \end{aligned}$$

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \end{aligned}$$

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p_\theta(z)) \end{aligned}$$

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p_\theta(z)) \end{aligned}$$

- We know $q_\phi(z|x) \rightarrow$ **Encoder**
- We know $p_\theta(z)$ (The distribution we want to have on the latent space)
- We know $p_\theta(x|z) \rightarrow$ **Decoder**
- We don't know the derivative when sampling $\mathbb{E}_{z \sim q_\phi(z|x)}$ (for backpropagation)

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p_\theta(z)) \end{aligned}$$

We can remark that we have two terms:

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p_\theta(z)) \end{aligned}$$

We can remark that we have two terms:

- $\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] \rightarrow$ fitting the data

Variational Auto-Encoder: The evidence lower bound

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(z, x)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [p_\theta(z) \log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} [\log q_\phi(z|x) - \log p_\theta(z)] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL(q_\phi(z|x) || p_\theta(z)) \end{aligned}$$

We can remark that we have two terms:

- $\mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] \rightarrow$ fitting the data
- $-KL(q_\phi(z|x) || p_\theta(z)) \rightarrow$ for the regularisation of the latent space

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using gradient descent:

- Have $\nabla_\phi \mathbb{E}_{z \sim q_\phi} [f(z)] = \mathbb{E}_{z \sim q_\phi} [f(z) \nabla_\phi \log(q_\phi(z))]$
- Estimate by Monte-Carlo: $\frac{1}{L} \sum_{l=1}^L [f(z^l) \nabla_\phi \log(q_\phi(z^l))]$ (L samples)

¹<https://towardsdatascience.com/policy-gradients-in-a-nutshell-8b72f9743c5d>

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using gradient descent:

- Have $\nabla_\phi \mathbb{E}_{z \sim q_\phi} [f(z)] = \mathbb{E}_{z \sim q_\phi} [f(z) \nabla_\phi \log(q_\phi(z))]$
- Estimate by Monte-Carlo: $\frac{1}{L} \sum_{l=1}^L [f(z^l) \nabla_\phi \log(q_\phi(z^l))]$ (L samples)

→ Bad estimator (high variance)

→ Connexion with Policy Gradient in reinforcement Learning¹

¹<https://towardsdatascience.com/policy-gradients-in-a-nutshell-8b72f9743c5d>

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

²**Original Paper:** Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

- Sample some noise $\epsilon \sim p(\epsilon)$

²Original Paper: Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

- Sample some noise $\epsilon \sim p(\epsilon)$
- Apply the transformation $\hat{z} = g_\phi(\epsilon, x)$

²Original Paper: Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

- Sample some noise $\epsilon \sim p(\epsilon)$
- Apply the transformation $\hat{z} = g_\phi(\epsilon, x)$
- f_ϕ should be differentiable

²Original Paper: Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

- Sample some noise $\epsilon \sim p(\epsilon)$
- Apply the transformation $\hat{z} = g_\phi(\epsilon, x)$
- f_ϕ should be differentiable
- Estimate with Monte-Carlo $\mathbb{E}_{z \sim q_\phi(z|x)} [g(z)] \approx \frac{1}{L} \sum_{l=1}^L [g(f_\theta(\epsilon^l, x))]$

²Original Paper: Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]?$

Using the reparametrization trick:

Estimate $z \sim q_\phi(z|x)$ by using a differentiable transformation:

- Sample some noise $\epsilon \sim p(\epsilon)$
- Apply the transformation $\hat{z} = g_\phi(\epsilon, x)$
- f_ϕ should be differentiable
- Estimate with Monte-Carlo $\mathbb{E}_{z \sim q_\phi(z|x)} [g(z)] \approx \frac{1}{L} \sum_{l=1}^L [g(f_\theta(\epsilon^l, x))]$

Notice that this is the Variational Auto-Encoder approach²

²Original Paper: Auto-Encoding Variational Bayes, Kingma And Welling, 2013

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization Gaussian case:

Estimate $z_i \sim q_\phi(z|x) = \mathcal{N}(\mu_i, \sigma_i^2)$ (in practice the encoder predict μ and σ for latent component):

- Sample some noise $\epsilon \sim \mathcal{N}(0, 1)$

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization Gaussian case:

Estimate $z_i \sim q_\phi(z|x) = \mathcal{N}(\mu_i, \sigma_i^2)$ (in practice the encoder predict μ and σ for latent component):

- Sample some noise $\epsilon \sim \mathcal{N}(0, 1)$
- Apply the transformation $\hat{z} = f_\phi(\epsilon, x) = \epsilon * \sigma_i + \mu_i$

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization Gaussian case:

Estimate $z_i \sim q_\phi(z|x) = \mathcal{N}(\mu_i, \sigma_i^2)$ (in practice the encoder predict μ and σ for latent component):

- Sample some noise $\epsilon \sim \mathcal{N}(0, 1)$
- Apply the transformation $\hat{z} = f_\phi(\epsilon, x) = \epsilon * \sigma_i + \mu_i$
- Estimate with Monte-Carlo $\mathbb{E}_{z \sim q_\phi(z|x)} [g(z)] \approx \frac{1}{L} \sum_{l=1}^L [g(\epsilon^{(l)} \sigma_i + \mu_i)]$

Variational Auto-Encoder: The reparametrization trick

Estimate ϕ for $\mathbb{E}_{z \sim q_\phi} [f(z)]$

Finding parameters ϕ of the distribution considering loss $\mathbb{E}_{z \sim q_\phi} [f(z)]$?

Using the reparametrization Gaussian case:

Estimate $z_i \sim q_\phi(z|x) = \mathcal{N}(\mu_i, \sigma_i^2)$ (in practice the encoder predict μ and σ for latent component):

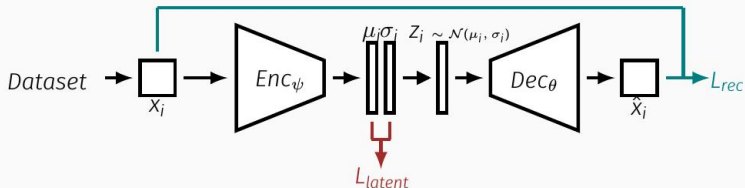
- Sample some noise $\epsilon \sim \mathcal{N}(0, 1)$
- Apply the transformation $\hat{z} = f_\phi(\epsilon, x) = \epsilon * \sigma_i + \mu_i$
- Estimate with Monte-Carlo $\mathbb{E}_{z \sim q_\phi(z|x)} [g(z)] \approx \frac{1}{L} \sum_{l=1}^L [g(\epsilon^{(l)} \sigma_i + \mu_i)]$

In most implementation, only one sample is used in the Monte-Carlo approximation

Variational Auto-Encoder: Summary

- Latent representation follow a distribution $p(z)$
- $q_{\phi}(z|x)$ encode data (It is a neural network)
- $p(x|z)$ decode data (It is also a neural network)
- We minimize the ELBO loss
- We sample in forward according to a distribution

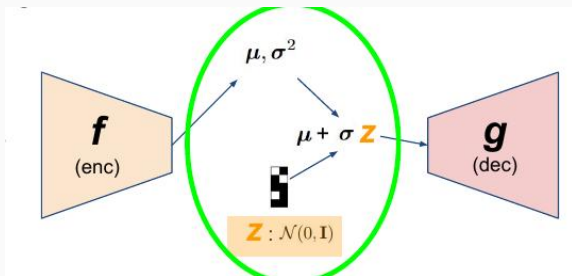
Variational Auto-Encoder: Framework



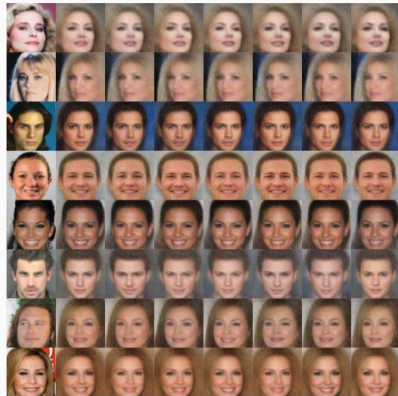
- encoding cost: $L_{latent} = \sum_i D_{KL} (Enc_\psi(x_i) || \mathcal{N}(0; 1))$
- reconstruction loss:

$$\begin{aligned} L_{rec} &= \sum_i \mathbb{E}_{z \sim Enc_\psi(x_i)} [-\log p_{Dec_\theta(z)}(x_i)] \\ &= \sum_i \mathbb{E}_{z \sim Enc_\psi(x_i)} \|Dec_\theta(z) - x_i\|^2 + cst. \end{aligned}$$

Variational Auto-Encoder: reparametrization trick



Variational Auto-Encoder: Examples



Variational Auto-Encoder: The latent space



Example of the latent space in the VAE

Variational Auto-Encoder: Code Examples

```
class MLPVAE(nn.Module):
    def __init__(self, input_size, inter_size, latent_size):
        super().__init__()
        self.latent_size = latent_size
        self.encoder = nn.Sequential(...)
        self.m_sigma = nn.Linear(inter_size, latent_size)
        self.m_mu = nn.Linear(inter_size, latent_size)
        self.decoder = nn.Sequential(...)

    def decode(self, mu, sigma):
        z = torch.randn(mu.shape) * sigma + mu
        return torch.sigmoid(self.decoder(z))

    def generate(self):
        z = torch.randn(1, latent_size)
        return torch.sigmoid(self.decoder(z))

    def forward(self, x):
        r = self.encoder(x)
        mu, log_sigma = self.m_mu(r), self.m_sigma(r)
        y = self.decode(mu, torch.exp(0.5 * log_sigma))
        return y, mu, log_sigma
```

- An Introduction to Variational Autoencoders
<https://arxiv.org/pdf/1906.02691>
- Auto-Encoding Variational Bayes <https://arxiv.org/abs/1312.6114>

A short introduction to Generative Adversarial Network

Objective

Find a generative model

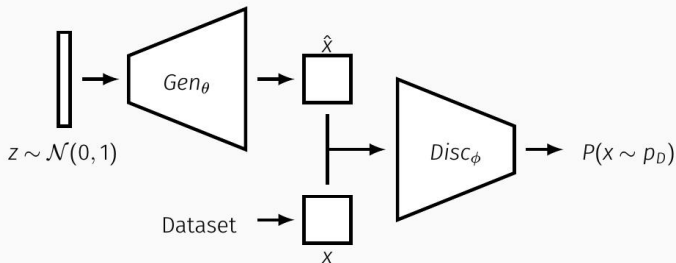
- Classical approach: learn a distribution
- Idea: Sampling directly on latent space (z) and evaluate if the generated images is likely to be a correct image

Principle

- Find a good generative model where generated samples cannot be discriminated from real samples

Generative Adversarial Network: Principle

- A dataset of true samples x (real)
- A generator G (a variational decoder), that sample from $p_g(z)$ and generate data (fake data)
- A discriminator D that discriminates real data from generated ones
 - Generator $G_\theta : \mathcal{L} \rightarrow \mathcal{D}$
 - Discriminator $D_\phi : \mathcal{D} \rightarrow [0, 1]$



Training GAN

- Alternate Optimisation
- Optimize D to discriminate fake from generated
- Optimize G to desceive D

Objective:

$$\text{Min}_G \text{Max}_D \mathbb{E}_{x \sim \text{data}} [\log D(x)] + \mathbb{E}_{z \sim p_g(z)} [\log(1 - D(G(z)))]$$

Algorithm 1 Training GAN

for n iterations **do**

for k steps **do**

 ▷ Discriminator update loop

- Sample $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ from prior noise $p_g(z)$
- Sample $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ from data distribution $p_{data}(x)$
- Update D by gradient descent:

$$\nabla_{\theta_D} - \frac{1}{m} \sum_{i=1}^m \left[\log D(x^i) + \log(1 - D(G(z^{(i)}))) \right]$$

end for

 ▷ Generator update

- Sample $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$ from prior noise $p_g(z)$
- Update D by gradient descent:

$$\nabla_{\theta_D} \frac{1}{m} \sum_{i=1}^m \left[\log(1 - D(G(z^{(i)}))) \right]$$

end for

Generative Adversarial Network: GAN vs VAE

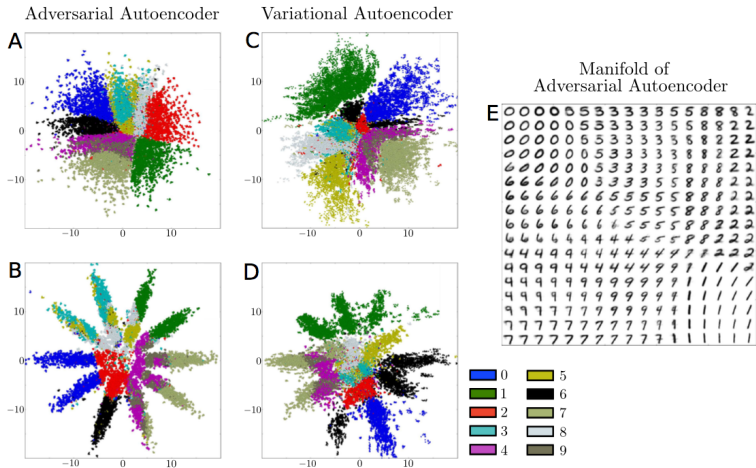


Figure 2: Comparison of adversarial and variational autoencoder on MNIST. The hidden code z of the *hold-out* images for an adversarial autoencoder fit to (a) a 2-D Gaussian and (b) a mixture of 10 2-D Gaussians. Each color represents the associated label. Same for variational autoencoder with (c) a 2-D Gaussian and (d) a mixture of 10 2-D Gaussians. (e) Images generated by uniformly sampling the Gaussian percentiles along each hidden code dimension z in the 2-D Gaussian adversarial autoencoder.

Generative Adversarial Network: examples



Stability issues

If P_r (probability of real images) and P_g (on generated images) are not in the same manifolds :

- It exists a perfect discriminator
- End of optimization (but bad generator)

Stability issues

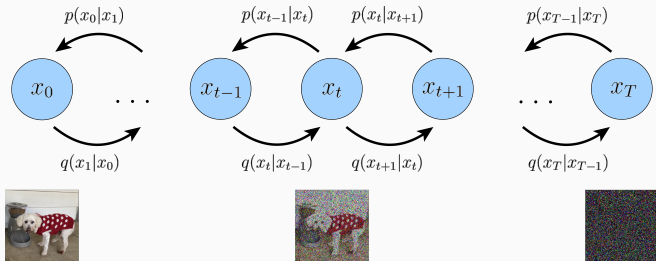
If P_r (probability of real images) and P_g (on generated images) are not in the same manifolds :

- It exists a perfect discriminator
- End of optimization (but bad generator)

Stabilizing training through regularization

- Improved training of Wasserstein GANs (Gulrajani, Ahmed, Arjovsky, Dumoulin, Courville 17)
→ using regularisation based on the wasserstein distance (with weights clipping)
- Stabilizing Training of Generative Adversarial Networks through Regularization (Roth, Lucchi, Nowozin, Hofmann, 17)

Diffusion models



Intuition beyond diffusion model

- Sampling gaussian noise
- Learn a denoising function $p(x_{t-1}|x_t)$
- Repeat the denoising step

What models are based on

- Stable diffusion
- Dall-E (since second version, first version was a GAN)
- Midjourney

²The image was taken from <https://calvinyluo.com/2022/08/26/diffusion-tutorial.html>