

# VAE

A Detailed Introduction to Variational Autoencoder

Zhenya Liu

# Overview

---

1. Background

2. Variational Autoencoder

3. Experiments

# Autoencoder

---

**Autoencoder** is a neural network designed to learn an identity function in an unsupervised way to reconstruct the original input while compressing the data in the process so as to discover a more efficient and compressed representation.

This idea was originated from 1980s, and formally proposed by Hinton & Salakhutdinov, 2006 [1].

Use neural networks to define: encoder  $g_\phi$  and decoder  $f_\theta$ .

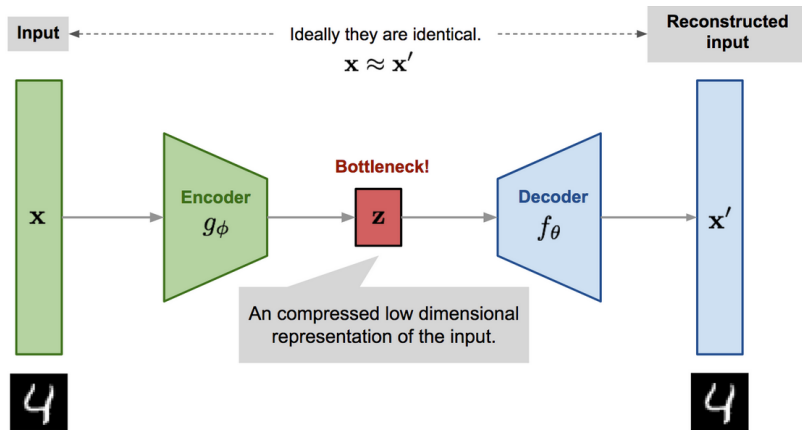


Figure: Illustration of autoencoder model architecture[4]

# Autoencoder

---

Training Objective:

$$L_{\text{AE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{x}^{(i)} - f_{\theta} \left( g_{\phi} \left( \mathbf{x}^{(i)} \right) \right) \right)^2$$

# Autoencoder

---

Unavoidable problem – Overfitting.

Several works on autoencoders were invented to train autoencoders with less overfitting and more robustness.

For instance, in **Denoising Autoencoder**(Vincent et al. 2008), the inputs are partially corrupted by masks or noises

$$\tilde{\mathbf{x}}^{(i)} \sim \mathcal{M}_{\mathcal{D}}(\tilde{\mathbf{x}}^{(i)} \mid \mathbf{x}^{(i)})$$

$$L_{\text{DAE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n \left( \mathbf{x}^{(i)} - f_{\theta} \left( g_{\phi} \left( \tilde{\mathbf{x}}^{(i)} \right) \right) \right)^2$$

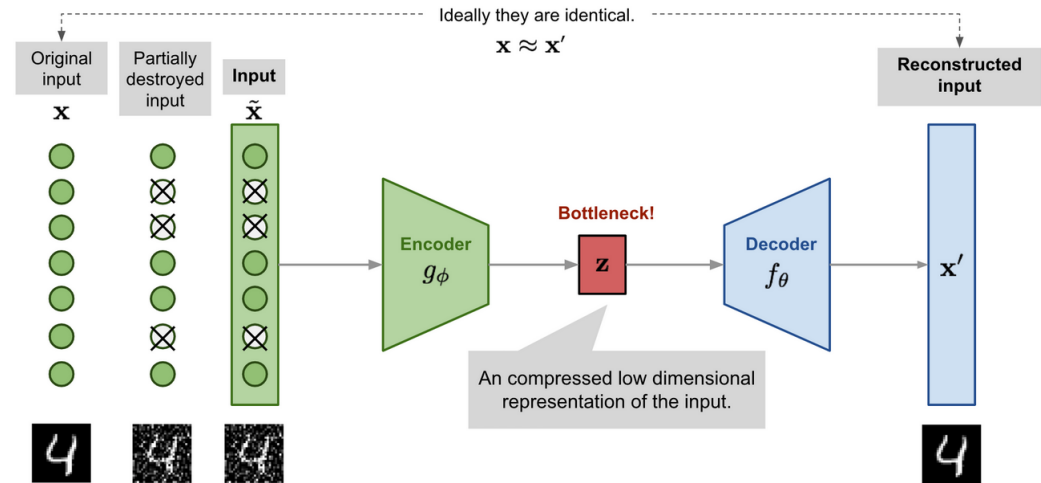


Figure: Illustration of denoising autoencoder model architecture[4]

# EM Algorithm

---

**Expectation Maximization** (EM) algorithm is an iterative method to find local maximum likelihood estimates (MLE) of parameters in statistical models, where the model depends on unobserved latent variables.



# EM Algorithm

---

## Problem Setup

Given one model parameterized by  $\theta$ , we aim to find  $\theta^*$  maximize the marginal likelihood of observed data  $\mathbf{X}$ , *i.e* MLE of  $P(\mathbf{X}|\theta)$ . We also have unobserved latent variable  $\mathbf{Z}$  such that  $P(\mathbf{Z})$  is unknown, but  $P(\mathbf{Z}|\theta)$ ,  $P(\mathbf{X}|\mathbf{Z}, \theta)$  is known.

$$L(\theta; \mathbf{X}) = p(\mathbf{X} | \theta) = \int p(\mathbf{X}, \mathbf{Z} | \theta) d\mathbf{Z} = \int p(\mathbf{X} | \mathbf{Z}, \theta) p(\mathbf{Z} | \theta) d\mathbf{Z}$$

In some cases, the above integral can be calculated, so we get one function over  $\theta$ . However, we always cannot find the closed form solution on  $\theta^*$  because solutions always have circular dependencies. Moreover, this integral is always non-concave, so we can only expect to find a local maximum.

# EM Algorithm

---

Expectation step(E step):

$$Q\left(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}\right)=\mathbb{E}_{\mathbf{Z} \sim p(\cdot \mid \mathbf{X}, \boldsymbol{\theta}^{(t)})}[\log p(\mathbf{X}, \mathbf{Z}) \mid \boldsymbol{\theta}]$$

Maximization step(M step):

$$\boldsymbol{\theta}^{(t+1)}=\arg \max _{\boldsymbol{\theta}} Q\left(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)}\right)$$

Equivalently,

$$\boldsymbol{\theta}^{(t+1)}=\arg \max _{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{Z} \sim p(\cdot \mid \mathbf{X}, \boldsymbol{\theta}^{(t)})}[\log p(\mathbf{X}, \mathbf{Z}) \mid \boldsymbol{\theta}]$$

# EM Algorithm

---

## Example: Gaussian Mixture Model

Given observed set of datapoints  $\{X_i\}$ , and set of unobserved latent variables  $\{Z_i\}$

$$X_i \mid (Z_i = 1) \sim \mathcal{N}_d(\mu_1, \Sigma_1) \text{ and } X_i \mid (Z_i = 2) \sim \mathcal{N}_d(\mu_2, \Sigma_2),$$

where

$$P(Z_i = 1) = \tau_1 \text{ and } P(Z_i = 2) = \tau_2 = 1 - \tau_1.$$

We can use EM Algorithm to find

$$\theta = (\boldsymbol{\tau}, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \Sigma_1, \Sigma_2)$$

# Variational Inference

---

## Problem Setup

Given prior  $P(\mathbf{Z})$  and likelihood  $P(\mathbf{X} | \mathbf{Z})$ , aim to find the posterior distribution

$$P(\mathbf{Z} | \mathbf{X}) = \frac{P(\mathbf{X} | \mathbf{Z})P(\mathbf{Z})}{P(\mathbf{X})} = \frac{P(\mathbf{X} | \mathbf{Z})P(\mathbf{Z})}{\int_{\mathbf{Z}} P(\mathbf{X}, \mathbf{Z}') d\mathbf{Z}'}$$

Generally, the space of  $\mathbf{Z}$  is large so that the integral part is intractable, so we need to find  $Q(\mathbf{Z}) \approx P(\mathbf{Z} | \mathbf{X})$ .

Question:

how do we compare two distributions? how to find a tractable formula on  $Q(\mathbf{Z})$  to minimize the "distance" between them?

# KL Divergence

---

## Kullback–Leibler Divergence

For discrete probability distributions  $P$  and  $Q$  defined on the same sample space,  $\mathcal{X}$ ,

$$D_{\text{KL}}(P\|Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

KL divergence  $D_{\text{KL}}(P\|Q)$  measures how much information is lost if the distribution  $Q$  is used to represent  $P$ .

# KL Divergence

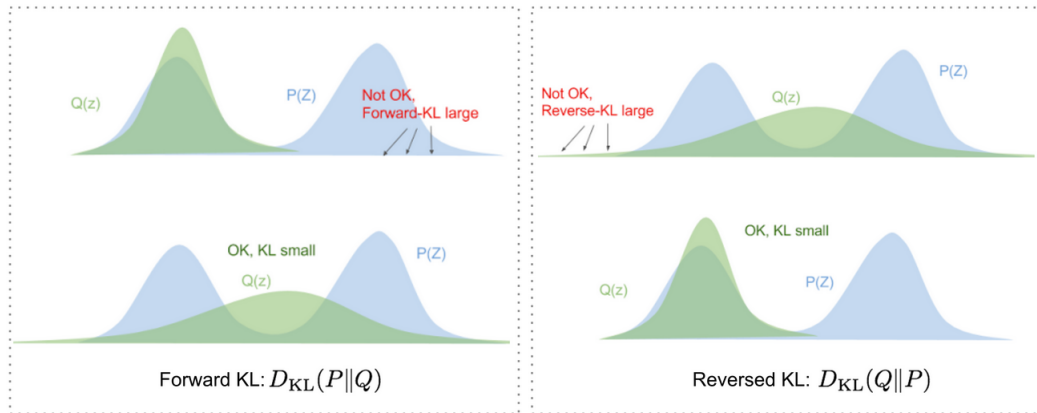


Figure: Forward and reversed KL divergence have different demands on how to match two distributions[4].

# KL Divergence

---

$$\begin{aligned} & D_{\text{KL}}(Q(\mathbf{Z}) \| P(\mathbf{Z}|\mathbf{X})) \\ &= \sum_{\mathbf{Z}} Q(\mathbf{Z}) \left[ \log \frac{Q(\mathbf{Z})}{P(\mathbf{Z}, \mathbf{X})} + \log P(\mathbf{X}) \right] \\ &= \sum_{\mathbf{Z}} Q(\mathbf{Z}) [\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] + \sum_{\mathbf{Z}} Q(\mathbf{Z}) [\log P(\mathbf{X})] \\ &= \mathbb{E}_{\mathbf{Q}}[\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] + \log P(\mathbf{X}) \end{aligned}$$

# KL Divergence

---

$$\begin{aligned}\log P(\mathbf{X}) &= D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X})) - \mathbb{E}_{\mathbf{Q}}[\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})] \\ &= D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X})) + \mathcal{L}(Q)\end{aligned}$$

Here  $\mathcal{L}(Q) = -\mathbb{E}_{\mathbf{Q}}[\log Q(\mathbf{Z}) - \log P(\mathbf{Z}, \mathbf{X})]$  is known as the **Evidence Lower Bound** (ELBO) of  $\log P(\mathbf{X})$ .

Most importantly,  $\mathcal{L}(Q)$  is tractable and optimizable.



# KL Divergence

---

Here is another way to get ELBO by using *Jensen's Inequality*:

$$\begin{aligned}\log p(\mathbf{X}) &= \log \int_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}) \\ &= \log \int_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}) \frac{Q(\mathbf{Z})}{Q(\mathbf{Z})} \\ &= \log \left( \mathbb{E}_Q \left[ \frac{p(\mathbf{X}, \mathbf{Z})}{Q(\mathbf{Z})} \right] \right) \\ &\geq \mathbb{E}_Q[\log p(\mathbf{X}, \mathbf{Z})] - \mathbb{E}_Q[\log Q(\mathbf{Z})] \\ &= \mathcal{L}(Q)\end{aligned}$$

# KL Divergence

---

Back to what we get

$$\log P(\mathbf{X}) = D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X})) + \mathcal{L}(Q)$$

So

$$\mathcal{L}(Q) = \log P(\mathbf{X}) - D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X}))$$

Maximize  $\mathcal{L}(Q) \implies$  maximize  $\log P(\mathbf{X})$  and minimize  $D_{\text{KL}}(Q(\mathbf{Z})\|P(\mathbf{Z}|\mathbf{X}))$ .

# ELBO and EM Algorithm

---

Based on  $\mathcal{L}(Q)$ , we can prove the effectiveness of EM Algorithm.

Since

$$\log P(\mathbf{X}) = \mathcal{L}(Q) + D_{\text{KL}}(Q(\mathbf{Z}) \| P(\mathbf{Z} | \mathbf{X}))$$

When  $Q(\mathbf{Z}) = P(\mathbf{Z} | \mathbf{X})$ ,

$$\log P(\mathbf{X} | \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{P(\mathbf{Z}|\mathbf{X},\boldsymbol{\theta})}[\log P(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})]}_{\text{EM step maximizes it}} - \mathbb{E}_{P(\mathbf{Z}|\mathbf{X},\boldsymbol{\theta})}[\log P(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})]$$

For one iteration of EM Algorithm, let  $\boldsymbol{\theta}^{(t)}$  be the initial value, then we can also get

$$\log P(\mathbf{X} | \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{P(\mathbf{Z}|\mathbf{X},\boldsymbol{\theta}^{(t)})}[\log P(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})]}_{H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})} - \underbrace{\mathbb{E}_{P(\mathbf{Z}|\mathbf{X},\boldsymbol{\theta}^{(t)})}[\log P(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta})]}_{I(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})}$$

Question: why?

# ELBO and EM Algorithm

---

During this iteration, we get  $\boldsymbol{\theta}^{(t+1)}$ , and

$$\log p(\mathbf{X} \mid \boldsymbol{\theta}^{(t+1)}) = H(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) + I(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)})$$

By *Gibb's Inequality*,

$$\underbrace{I(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)})}_{\text{cross entropy}} - \underbrace{I(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)})}_{\text{entropy}} \geq 0$$

# ELBO and EM Algorithm

---

Thus, for each EM update,

$$\begin{aligned} & \log p(\mathbf{x} \mid \boldsymbol{\theta}^{(t+1)}) - \log p(\mathbf{x} \mid \boldsymbol{\theta}^{(t)}) \\ &= H(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) + I(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - I(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) \\ &\geq H(\boldsymbol{\theta}^{(t+1)} \mid \boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}^{(t)} \mid \boldsymbol{\theta}^{(t)}) \end{aligned}$$

This shows that updating  $\boldsymbol{\theta}$  to improve  $H(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{(t)})$  causes  $\log p(\mathbf{X} \mid \boldsymbol{\theta})$  to improve at least as much.

# Variational Inference

---

We have decided the optimization objective  $\mathcal{L}(Q)$ , then we need to decide how to represent  $Q(\mathbf{Z})$ . We expect to give  $Q(\mathbf{Z})$  high flexibility and tractable computability over expectation.

In variational inference, we pick a family of distributions over the latent variables with its own variational parameters, *i.e*

$$Q(\mathbf{Z}) = q(z_{1:m}|\nu)$$

# Mean Field Variational Inference

---

Assume that the variational family factorizes

$$q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j)$$

Here we assume all latent variables are independent. This also restricts our approximation to the true posterior distribution because latent variables are always dependent, *i.e* latent variables in Gaussian mixture models.

# Mean Field Variational Inference

---

Define  $\mathbf{X} = x_{1:n}$ , and choose one  $z_j$  from  $z_{1:m}$ . Refer  $-j$  to  $\{1, \dots, m\} \setminus \{j\}$ ,

$$E_{q(z_{1:m})} [\log P(z_{1:m}, x_{1:n})] = \int_{z_j} q_j(z_j) E_{-j} [\log P(x_{1:n}, z_{-j})] dz_j$$

Evaluate the expectation

$$E_{q(z_{1:m})} [\log q(z_{1:m})] = \sum_{j=1}^m E_j [\log q(z_j)] + C$$

where  $E_j$  represents the expectation with respect to  $q(z_j)$ , and  $C$  is some terms irrelevant with  $q(z_j)$ .



# Mean Field Variational Inference

---

Now calculate  $\mathcal{L}(Q)$

$$\begin{aligned}\mathcal{L} &= E_{q(z_{1:m})} [\log P(z_{1:m}, x_{1:n}) - \log q(z_{1:m})] \\ &= \int_{z_j} q_j(z_j) E_{-j} [\log P(x_{1:n}, z_{-j})] dz_j - \sum_{j=1}^m E_j [\log q(z_j)] - C\end{aligned}$$

From now on, we treat  $\mathcal{L}$  as one **functional** of  $q(z_j)$  while fixing all other  $q(z_{-j})$  and we want to find  $q^*(z_j)$  which maximizes  $\mathcal{L}$ .

How to find such  $q^*(z_j)$ ? Use methods in **Calculus of Variation**.

# Calculus of Variation

---

## Theorem (Euler Lagrange Equation)

*Given a functional  $J(f)$  defined on all functions  $f \in C^2[a, b]$  such that  $f(a) = A, f(b) = B$ ,*

$$J(f) = \int_a^b L(x, f, f') \, dx$$

*then if  $J(f)$  is a local extremum, then  $f$  satisfies*

$$\frac{d}{dx} \left( \frac{\partial L}{\partial f'} \right) - \frac{\partial L}{\partial f} = 0$$

# Calculus of Variation

---

## Example

$$s = \int_a^b \sqrt{dx^2 + dy^2} = \int_a^b \sqrt{1 + y'^2} \, dx$$

Let  $L(x, y, y') = \sqrt{1 + y'^2}$ ,

$$\frac{\partial L(x, y, y')}{\partial y'} = \frac{y'}{\sqrt{1 + y'^2}} \quad \text{and} \quad \frac{\partial L(x, y, y')}{\partial y} = 0$$

By **Euler Lagrange Equation**,

$$\begin{aligned} \frac{d}{dx} \frac{y'(x)}{\sqrt{1 + (y'(x))^2}} &= 0, \quad \frac{y'(x)}{\sqrt{1 + (y'(x))^2}} = C \\ \Rightarrow y'(x) &= \frac{C}{\sqrt{1 - C^2}} = A \\ \Rightarrow y(x) &= Ax + B \end{aligned}$$

# Mean Field Variational Inference

---

Now we drop out all irrelevant terms with  $q(z_j)$  in  $\mathcal{L}$ , and we get  $\mathcal{L}_j$ ,

$$\begin{aligned}\mathcal{L}_j &= \mathbb{E}_{-j} [\log p(z_j \mid z_{-j}, x_{1:n})] - \mathbb{E}_j [\log q(z_j)] \\ &= \int q(z_j) \mathbb{E}_{-j} [\log p(z_j \mid z_{-j}, x_{1:n})] - q(z_j) \log q(z_j) dz_j\end{aligned}$$

Define  $L$  be the inner part of above integral, by **Euler Lagrange Equation**,

$$\frac{dL}{dq(z_j)} = \mathbb{E}_{-j} [\log p(z_j \mid z_{-j}, x_{1:n})] - \log q(z_j) = 0$$

Then we get expected  $q^*(z_j)$  in **unconstrained** case.

# Mean Field Variational Inference

---

Note that since  $q(z_j)$  is one probability distribution, so there is one constrain functional  $\int q(z_j)dz_j = 1$ . We use **Lagrange Multiplier** to handle this situation. In this case, the solution  $q(z_j)$  must satisfy

$$\frac{dL}{dq(z_j)} - \lambda \cdot \frac{dq(z_j)}{dq(z_j)} = 0$$

so

$$\begin{aligned}\log q_j(z_j) &= E_{-j} [\log P(X, z_j, z_{-j})] + \lambda \\ q^*(z_j) &= \exp \{E_{-j} [\log p(z_j, z_{-j}, x_{1:n})]\} + \exp(\lambda)\end{aligned}$$

where  $\exp(\lambda)$  can be figured out when normalizing  $q^*(z_j)$ .

# Mean Field Variational Inference

---

So far we only consider maximize  $\mathcal{L}$  on one coordinate  $q(z_j)$  while fixing other coordinates. In practice, this is the **coordinate ascend** step of the algorithm.

We initialize all  $q(z_{1:m})$ , and iteratively update each coordinate  $q(z_j)$  by previous variational method. Finally we will converge to one local extreme  $\mathcal{L}$ .

This method is called **Mean Field Variational Inference**. The meaning of "Mean Field" is by the fact that  $q^*(z_j)$  is determined by the expectation over all  $q(z_{1:m})$ .

# Variational Autoencoder

---

## Problem Setting

Given  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ , we aim to model one distribution  $P$  parameterized by  $\theta$  with one continuous latent variable  $\mathbf{z}$  such that

- $\log p_{\theta}(\mathbf{X})$  is maximized
- $p_{\theta}(\mathbf{z}|\mathbf{X})$  can be efficiently approximated
- $p_{\theta}(\mathbf{X})$  can be efficiently approximated

# Variational Autoencoder

---

The only known variables are  $\mathbf{X}$  and the prior  $p(\mathbf{z})$ , and likelihood  $p_{\boldsymbol{\theta}}(\mathbf{X}|\mathbf{z})$  when  $\boldsymbol{\theta}$  is known.

EM algorithm does not work because we do not know the posterior  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{X})$ .

Mean field variational inference does not work here to approximate the posterior because  $\boldsymbol{\theta}$  is unknown, and even it is known, the necessary integral is intractable when  $\mathbf{z}$  is continuous.

**Question:** What should we rely on?



# Neural Networks

**Variational Autoencoder(VAE)**(Kingma & Welling, 2014) is a generative model using neural networks to train on ELBO.

# Variational Autoencoder

---

Construct neural networks  $p_{\theta}(\mathbf{x}|\mathbf{z})$  (*probabilistic decoder*) for likelihood distribution and  $q_{\phi}(\mathbf{z}|\mathbf{x})$  (*probabilistic encoder*) for marginal distribution.

**Objection:** Maximize  $\log p_{\theta}(\mathbf{x})$ , and approximate  $p_{\theta}(\mathbf{z}|\mathbf{x})$  by  $q_{\phi}(\mathbf{z}|\mathbf{x})$ .

Recall:

Maximize  $\mathcal{L}(Q) \implies$  maximize  $\log P(\mathbf{X})$  and minimize  $D_{\text{KL}}(Q(\mathbf{Z})|P(\mathbf{Z}|\mathbf{X}))$ .

so the **ELBO** is a sufficient objective for jointly training  $\theta, \phi$ .

# Variational Autoencoder

---

Marginal log likelihood for individual datapoint

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}),$$

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [-\log q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}) + \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = -D_{KL}(q_{\boldsymbol{\phi}}(\mathbf{z} | \mathbf{x}^{(i)}) \| p_{\boldsymbol{\theta}}(\mathbf{z})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} | \mathbf{z})]$$

The *KL* divergence term on **RHS** can always be analytically calculated, but the expectation term (reconstruction error) must be approximated by samples.

# Variational Autoencoder

---

We cannot directly calculate the gradient of  $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})]$  with respect to  $\theta, \phi$ . The usual Monte Carlo gradient estimator for this type of problem is

$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} \left[ f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z}) \right] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$$

where  $\mathbf{z}^{(l)} \sim q_\phi(\mathbf{z} | \mathbf{x}^{(i)})$ .

This sampling estimator always exhibits a high variance.

# Reparameterization

---

Under certain mild conditions, for  $\tilde{\mathbf{z}} \sim q_\phi(\mathbf{z} \mid \mathbf{x})$ , we can reparameterize it with one differentiable transformation  $g_\phi(\epsilon, \mathbf{x})$  and noise variable  $\epsilon$ ,

$$\tilde{\mathbf{z}} = g_\phi(\epsilon, \mathbf{x}) \quad \text{with} \quad \epsilon \sim p(\epsilon)$$

Then we apply Monte Carlo gradient estimator on sampling  $\epsilon$ ,

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})}[f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} \left[ f \left( g_\phi \left( \epsilon, \mathbf{x}^{(i)} \right) \right) \right] \simeq \frac{1}{L} \sum_{l=1}^L f \left( g_\phi \left( \epsilon^{(l)}, \mathbf{x}^{(i)} \right) \right)$$

where  $\epsilon^{(l)} \sim p(\epsilon)$

# Reparameterization

---

This reparameterization trick yields Stochastic Gradient Variational Bayes (SGVB) estimator for the **ELBO**.

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \simeq \tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = -D_{KL} \left( q_{\phi}(\mathbf{z} \mid \mathbf{x}^{(i)}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}) \right) + \frac{1}{L} \sum_{l=1}^L \left( \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i,l)}) \right)$$

where  $\mathbf{z}^{(i,l)} = g_{\phi}(\epsilon^{(i,l)}, \mathbf{x}^{(i)})$  and  $\epsilon^{(l)} \sim p(\epsilon)$

# Variational Autoencoder

---

Consider the dataset  $\mathbf{X}$  with  $N$  samples, we approximate the total **ELBO** loss based on **minibatch**  $\mathbf{X}^M == \{\mathbf{x}^{(i)}\}_{i=1}^M$ ,

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{X}) \simeq \tilde{\mathcal{L}}^M(\boldsymbol{\theta}, \phi; \mathbf{X}^M) = \frac{N}{M} \sum_{i=1}^M \tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)})$$

# Variational Autoencoder

---

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$

---

Figure: VAE Algorithm [2]



# Variational Autoencoder

---

## Example

Define prior  $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ .

Both  $p_{\theta}(\mathbf{x} \mid \mathbf{z})$  and  $q_{\phi}(\mathbf{z} \mid \mathbf{x})$  are multivariate Gaussian distribution with diagonal covariance, *i.e*

$$\log p_{\theta}(\mathbf{x} \mid \mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\sigma}_{\theta}^2 \mathbf{I})$$

$$\log q_{\phi}(\mathbf{z} \mid \mathbf{x}) = \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}, \boldsymbol{\sigma}_{\phi}^2 \mathbf{I})$$

# Variational Autoencoder

---

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}^{(i)}) = -D_{KL} \left( q_{\boldsymbol{\phi}} \left( \mathbf{z} \mid \mathbf{x}^{(i)} \right) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}) \right) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z} \mid \mathbf{x}^{(i)})} \left[ \log p_{\boldsymbol{\theta}} \left( \mathbf{x}^{(i)} \mid \mathbf{z} \right) \right]$$

Since

$$\begin{aligned} \int q_{\boldsymbol{\phi}}(\mathbf{z}) \log q_{\boldsymbol{\phi}}(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\boldsymbol{\phi}}, \boldsymbol{\sigma}_{\boldsymbol{\phi}}^2 \mathbf{I}) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\boldsymbol{\phi}}, \boldsymbol{\sigma}_{\boldsymbol{\phi}}^2 \mathbf{I}) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2) \end{aligned}$$

$$\begin{aligned} -D_{KL} (q_{\boldsymbol{\phi}}(\mathbf{z}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z})) &= \int q_{\boldsymbol{\phi}}(\mathbf{z}) (\log p_{\boldsymbol{\theta}}(\mathbf{z}) - \log q_{\boldsymbol{\phi}}(\mathbf{z})) d\mathbf{z} \\ &= \frac{1}{2} \sum_{j=1}^J \left( 1 + \log \left( (\sigma_j)^2 \right) - (\mu_j)^2 - (\sigma_j)^2 \right) \end{aligned}$$

# Variational Autoencoder

---

We reparameterize  $\mathbf{z}^{(i,l)} \sim q_\phi(\mathbf{z} \mid \mathbf{x}^{(i)})$  as  
 $\mathbf{z}^{(i,l)} = g_\phi(\mathbf{x}^{(i)}, \boldsymbol{\epsilon}^{(l)}) = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\epsilon}^{(l)}$  where  $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \simeq \frac{1}{2} \sum_{j=1}^J \left( 1 + \log \left( \left( \sigma_j^{(i)} \right)^2 \right) - \left( \mu_j^{(i)} \right)^2 - \left( \sigma_j^{(i)} \right)^2 \right) + \frac{1}{L} \sum_{l=1}^L \log p_{\boldsymbol{\theta}} \left( \mathbf{x}^{(i)} \mid \mathbf{z}^{(i,l)} \right)$$

# Experiments

---

## Encoder

```
1 # Network parameters
2 original_dim = 784
3 intermediate_dim = 256
4 latent_dim = 2
5 # Encoder
6 inputs = Input(shape=(original_dim,))
7 h = Dense(intermediate_dim, activation='relu')(inputs)
8 z_mean = Dense(latent_dim)(h)
9 z_log_var = Dense(latent_dim)(h)
```

# Experiments

---

## Reparameterization trick

```
1  # Reparameterization trick
2  def sampling(args):
3      z_mean, z_log_var = args
4      batch = K.shape(z_mean)[0]
5      dim = K.int_shape(z_mean)[1]
6      epsilon = K.random_normal(shape=(batch, dim))
7      return z_mean + K.exp(0.5 * z_log_var) * epsilon
8
9  z = Lambda(sampling, output_shape=(latent_dim,))([z_mean, z_log_var])
```

# Experiments

---

## Decoder

```
1 # Decoder
2 decoder_h = Dense(intermediate_dim, activation='relu')
3 decoder_mean = Dense(original_dim, activation='sigmoid')
4 h_decoded = decoder_h(z)
5 x_decoded_mean = decoder_mean(h_decoded)
```

# Experiments

---

## End-to-end Autoencoder

```
1 # End-to-end autoencoder
2 vae = Model(inputs, x_decoded_mean)
3
4 # Encoder, from inputs to latent space
5 encoder = Model(inputs, z_mean)
6
7 # Generator, from latent space to reconstructed inputs
8 decoder_input = Input(shape=(latent_dim,))
9 _h_decoded = decoder_h(decoder_input)
10 _x_decoded_mean = decoder_mean(_h_decoded)
11 generator = Model(decoder_input, _x_decoded_mean)
```

# Experiments

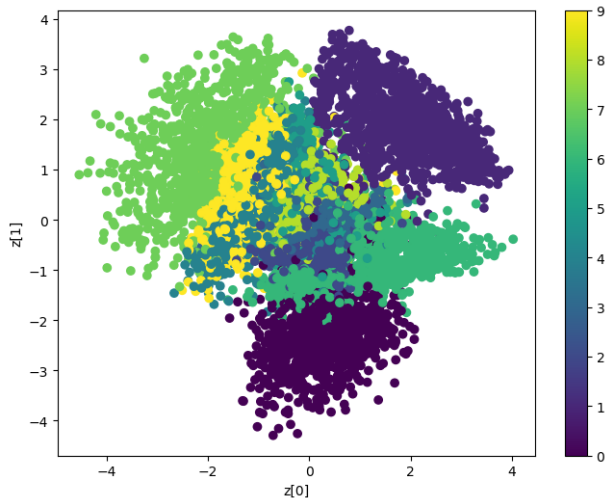
---

## Loss function & Train

```
1 # Loss
2 xent_loss = original_dim * binary_crossentropy(inputs, x_decoded_mean)
3 kl_loss = - 0.5 * K.sum(1 + z_log_var - K.square(z_mean) - K.exp(
    z_log_var), axis=-1)
4 vae_loss = K.mean(xent_loss + kl_loss)
5
6 vae.add_loss(vae_loss)
7 vae.compile(optimizer='rmsprop')
8
9 # Train
10 vae.fit(x_train, epochs=50, batch_size=256, validation_data=(x_test,
    None))
```



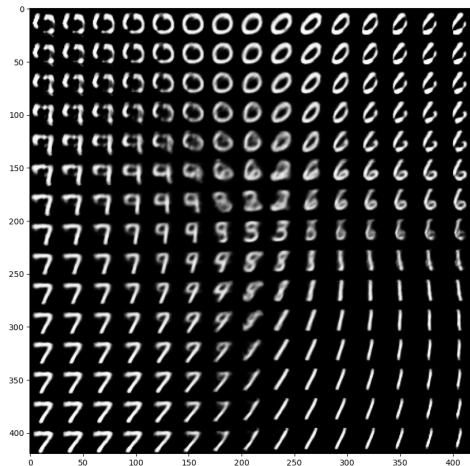
VAE implicitly maps one class of number to one neighbor in latent space ( $d = 2$ )



**Figure:** The test data is plotted in the latent space. Each point corresponds to an MNIST digit, and the color of the point indicates the digit's label.

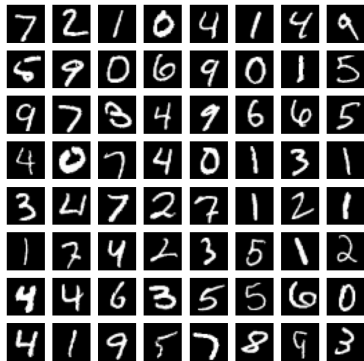
Grid sampling from the normal distribution.

Recall that  $q_\phi(\mathbf{z}|x) \approx p_\theta(\mathbf{z}) \sim \mathcal{N}(0, \mathbf{I})$

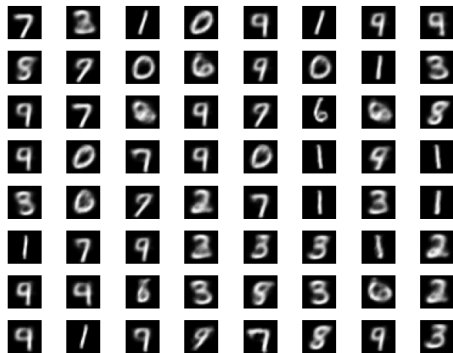


**Figure:** Each image is generated by sampling a point in the latent space and then using the generator (decoder) to produce an image.

## Reconstruction results



(a) sample



(b) reconstruction = 128

Generations from random noises

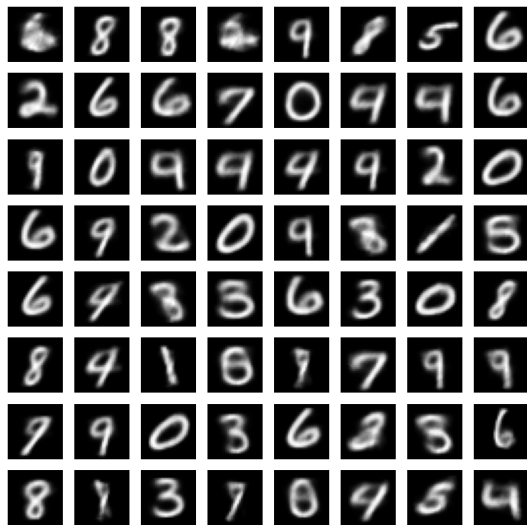


Figure: 64 samples from  $\mathcal{N}(0, I)$

Try with CIFAR-10,

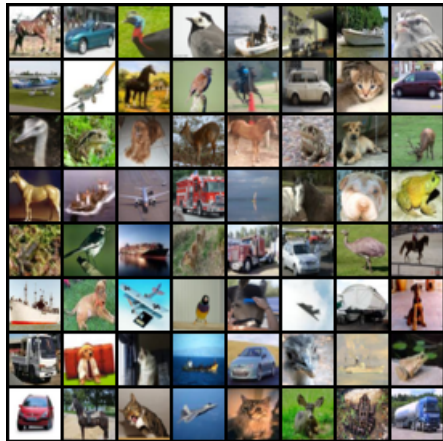


Figure: Randomly chosen 64 images from CIFAR-10

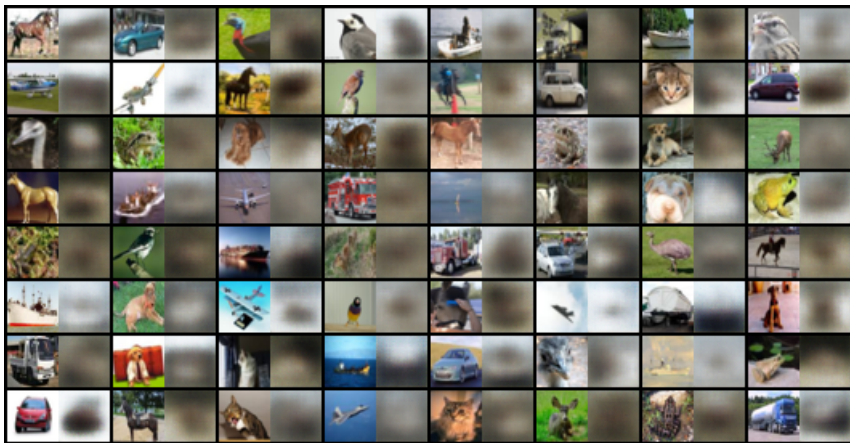


Figure: Images V.S. Reconstructions

## FID Score

---

1. (FID) Given two distributions  $P$  and  $Q$  on  $R^d$ , we assume they have finite first and second moments:  $\mu_X = \int x dP < \infty$ ,  $\mu_Y = \int y dQ < \infty$ , and  $\Sigma_X = \int x x^t dP < \infty$ ,  $\Sigma_Y = \int y y^t dQ < \infty$ .

$$FID(P, Q) = \inf \{E|X - Y|^2; (X, Y) : X \sim P, Y \sim Q\}$$

This is the minimum expected square distance between two random variables in  $R^d$  with some joint distribution such that the marginal distribution of  $X$  is  $P$  and the marginal distribution of  $Y$  is  $Q$ .

Theorem: Assume that  $\Sigma_X, \Sigma_Y$  are non-singular, i.e. positive definite. Then

$$FID(P, Q) = |\mu_X - \mu_Y|^2 + \text{tr} \left[ \Sigma_X + \Sigma_Y - 2 \left( \Sigma_X^{1/2} \Sigma_Y \Sigma_X^{1/2} \right)^{1/2} \right]$$

# Experiments

---

**FID-Calculation** The FID score we get is around 8.04

```
1 def fid_score(mu1, sigma1, mu2, sigma2):
2     """Compute FID score given mean and covariance of two datasets."""
3     diff = mu1 - mu2
4     # covmean = sqrtm(sqrtm(sigma1) @ sigma2 @ sqrtm(sigma1))
5     covmean = stable_covmean(sigma1, sigma2)
6     # Real part is taken to avoid any potential imaginary number due to
7     numerical issues
8     fid = diff.dot(diff) + np.trace(sigma1 + sigma2 - 2 * covmean)
9     return fid
```



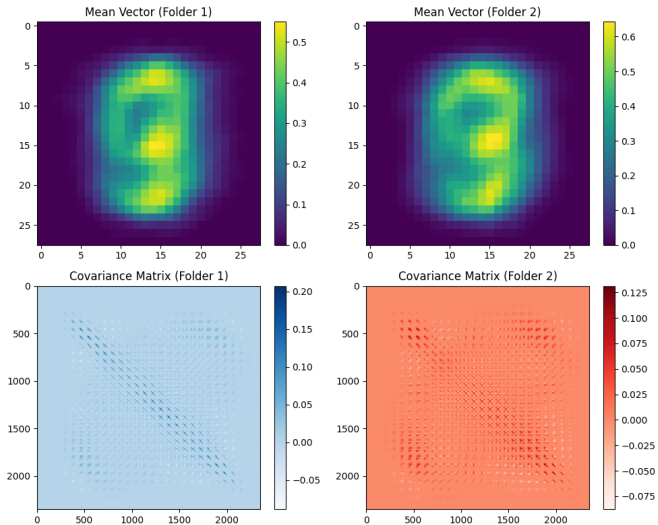


Figure: Computed means and Variances

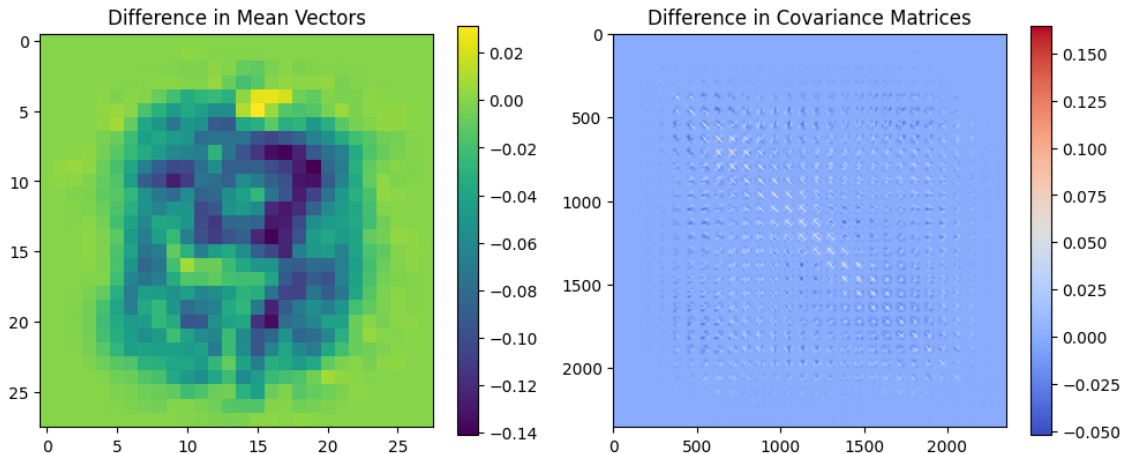


Figure: Differences in means and Variances

# References

---



Hinton, G. & Salakhutdinov, R. Reducing the dimensionality of data with neural networks. *Science*. **313**, 504-507 (2006)



Kingma, D. & Welling, M. Auto-encoding variational bayes. *ArXiv Preprint ArXiv:1312.6114*. (2013)



David M. Blei, A. & McAuliffe, J. Variational Inference: A Review for Statisticians. *Journal Of The American Statistical Association*. **112**, 859-877 (2017),



Weng, L. From Autoencoder to Beta-VAE. *Lilianweng.github.io*. (2018),  
<https://lilianweng.github.io/posts/2018-08-12-vae/>



Variational Bayesian methods. Wikipedia.  
[https : //en.wikipedia.org/wiki/Variational\\_Bayesian\\_methods](https://en.wikipedia.org/wiki/Variational_Bayesian_methods)



Expectation–maximization algorithm. Wikipedia. [https : //en.wikipedia.org/wiki/Expectation](https://en.wikipedia.org/wiki/Expectation)



ethanluoyc, pytorch-vae, 2022, GitHub repository, <https://github.com/ethanluoyc/pytorch-vae>

**End**