



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

Variational Autoencoder (VAE)

Cuong Dang, Nam Nguyen, Thanh Bui¹

¹CS Department
Ho Chi Minh City University Information of Technology (UIT)



Overview

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

I. Introduction



Introduction to Variational Autoencoders (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

• **Autoencoders (AE):**

- Autoencoders are deep learning models used for learning representations of data in a lower-dimensional latent space.
- They consist of two parts: the *Encoder* (which encodes data) and the *Decoder* (which reconstructs data).
- However, AE models are limited in their ability to generate new data in a smooth and natural way.



Introduction to Variational Autoencoders (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

● Autoencoders (AE):

- Autoencoders are deep learning models used for learning representations of data in a lower-dimensional latent space.
- They consist of two parts: the *Encoder* (which encodes data) and the *Decoder* (which reconstructs data).
- However, AE models are limited in their ability to generate new data in a smooth and natural way.

● Variational Autoencoders (VAE):

- VAE is an extension of Autoencoders where the latent space is modeled as a probabilistic distribution.
- The goal of VAE is to learn a probability distribution over the latent space, enabling the generation of new data samples.
- VAE uses *Variational Inference* methods to approximate the posterior distribution.



Differences Between Autoencoder (AE) and Variational Autoencoder (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- **Latent Space Representation:**

- **AE:** Encodes data into a fixed latent vector.
- **VAE:** Models the latent space as a probabilistic distribution (e.g., Gaussian distribution).



Differences Between Autoencoder (AE) and Variational Autoencoder (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

● Latent Space Representation:

- **AE:** Encodes data into a fixed latent vector.
- **VAE:** Models the latent space as a probabilistic distribution (e.g., Gaussian distribution).

● Goal:

- **AE:** Focuses on reconstructing input data as accurately as possible.
- **VAE:** Learns a distribution over the latent space for generating new data samples.



Differences Between Autoencoder (AE) and Variational Autoencoder (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- **Mathematical Framework:**

- **AE:** Uses a standard loss function (e.g., Mean Squared Error) for reconstruction.
- **VAE:** Combines a reconstruction loss and a regularization term (Kullback-Leibler divergence).



Differences Between Autoencoder (AE) and Variational Autoencoder (VAE)

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- **Mathematical Framework:**

- **AE:** Uses a standard loss function (e.g., Mean Squared Error) for reconstruction.
- **VAE:** Combines a reconstruction loss and a regularization term (Kullback-Leibler divergence).

- **Data Generation:**

- **AE:** Limited ability to generate new data.
- **VAE:** Generates smooth and realistic data by sampling from the learned distribution.



Why Choose Variational Autoencoders (VAE)?

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

● **Dimensionality Reduction:**

- One of the main advantages of autoencoders is their ability to reduce the dimensionality of data.
- They transform the data into a smaller latent representation while retaining the important information.



Why Choose Variational Autoencoders (VAE)?

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

● **Dimensionality Reduction:**

- One of the main advantages of autoencoders is their ability to reduce the dimensionality of data.
- They transform the data into a smaller latent representation while retaining the important information.

● **Feature Learning for Classification:**

- Autoencoders can serve as a feature learning tool.
- The latent space representation can be used as input for downstream classification tasks.



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

II. Structure of VAE



Basic Structure of VAE

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- **Encoder:**

- Takes the input data x and encodes it into a probabilistic representation in the latent space z .

- **Latent Space:**

- A space where data is represented probabilistically, capturing the underlying structure of the data.

- **Decoder:**

- Decodes the latent variable z back into the original data x , reconstructing the input.

- **Objective:**

- The goal is to learn to reconstruct the input data and optimize the distribution in the latent space.



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

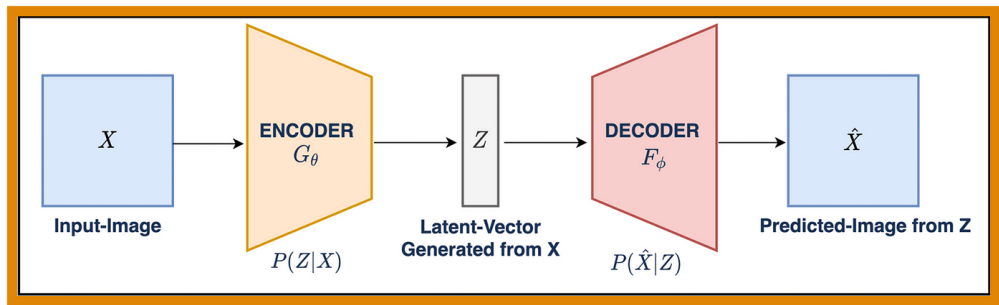


Figure 1: Structure of VAE



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

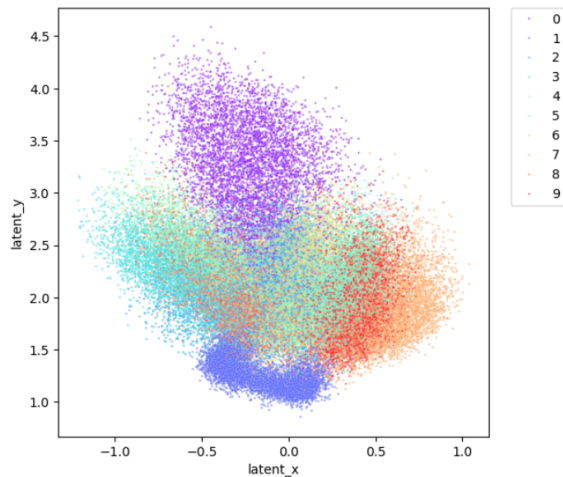


Figure 2: Latent space



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

III. Mathematical Symbols in VAE



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- $p(\mathbf{x})$: The true distribution of \mathbf{x} . The goal of diffusion models is to sample from this distribution, which is typically unknown.
- $p(z)$: The distribution of the latent variable, usually a zero-mean unit-variance Gaussian: $\mathcal{N}(0, \mathbf{I})$. This simplifies processing, as any distribution can be generated by mapping a Gaussian through a complex function.
- $p(z|\mathbf{x})$: The conditional distribution in the *encoder* — the likelihood of z given \mathbf{x} . This distribution is not directly accessible.
- $p(\mathbf{x}|z)$: The conditional distribution in the *decoder* — the probability of \mathbf{x} given z . This distribution is also not directly accessible.



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- $q_{\phi}(z|\mathbf{x})$: The approximate distribution for $p(z|\mathbf{x})$, parameterized by a neural network. Define:

$$(\mu, \sigma^2) = \text{EncoderNetwork}_{\phi}(\mathbf{x}),$$

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



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- $q_{\phi}(z|\mathbf{x})$: The approximate distribution for $p(z|\mathbf{x})$, parameterized by a neural network. Define:

$$(\mu, \sigma^2) = \text{EncoderNetwork}_{\phi}(\mathbf{x}),$$

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

- $p_{\theta}(\mathbf{x}|z)$: The approximate distribution for $p(\mathbf{x}|z)$, parameterized by a neural network. Define:

$$f_{\theta}(z) = \text{DecoderNetwork}_{\theta}(z),$$

$$p_{\theta}(\mathbf{x}|z) = \mathcal{N}(\mathbf{x} | f_{\theta}(z), \sigma_{\text{dec}}^2 \mathbf{I}).$$



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

The relationship between the input \mathbf{x} and the latent \mathbf{z} , as well as the conditional distributions, are summarized below

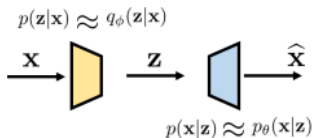


Figure 3: \mathbf{x} and \mathbf{z} are connected by $p(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{x}|\mathbf{z})$ in a VAE.



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

IV. Build construct in VAE



Encoding and Decoding

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Suppose that we have a random variable $x \in \mathbb{R}^d$ and a latent variable $z \in \mathbb{R}^d$ such that:

$$x \sim p(x) = \mathcal{N}(x | \mu, \sigma^2 \mathbf{I}),$$

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

We want to construct a Variational Autoencoder. By this, we mean that we want to build two mappings: the encoder $\text{Encoder}(\cdot)$ and the decoder $\text{Decoder}(\cdot)$. The encoder will take a sample x and map it to the latent variable z , whereas the decoder will take the latent variable z and map it to the generated variable \hat{x} .



If we know $p(x)$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- If we know the distribution $p(x)$, then there is a simple solution:

$$z = \frac{x - \mu}{\sigma}, \quad \hat{x} = \mu + \sigma z.$$

- In this case, the true distributions are expressed using delta functions:

$$p(x | z) = \delta(x - (\sigma z + \mu)),$$

$$p(z | x) = \delta\left(z - \frac{x - \mu}{\sigma}\right).$$



Dirac Delta Function

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Dirac Delta Function:

- The Dirac delta function is often represented as:

$$\delta_a(x) = \frac{1}{|a|\sqrt{\pi}} e^{-\left(\frac{x}{a}\right)^2}.$$

- For computational purposes, it can be approximated by a Gaussian distribution.

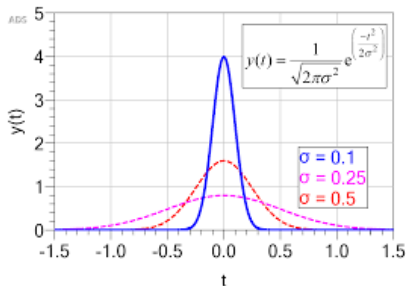


Figure 4: Dirac Delta Function



If we do not know $p(x)$

Let's first define the encoder. Our encoder in this example takes the input x and generates a pair of parameters $\hat{\mu}(x)$ and $\hat{\sigma}^2(x)$, denoting the parameters of a Gaussian. Then, we define $q_{\phi}(z | x)$ as a Gaussian:

$$(\hat{\mu}(x), \hat{\sigma}^2(x)) = \text{Encoder}_{\phi}(x),$$

$$q_{\phi}(z | x) = \mathcal{N}(z | \hat{\mu}(x), \hat{\sigma}^2(x)I).$$

For the purpose of discussion, we assume that $\hat{\mu}$ is an affine function of x such that: $\hat{\mu}(x) = ax + b$ for some parameters a and b . Similarly, we assume that:

$$\hat{\sigma}^2(x) = t^2,$$

for some scalar t . This will give us:

$$q_{\phi}(z | x) = \mathcal{N}(z | ax + b, t^2I).$$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE



If we do not know $p(x)$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

For the decoder, we deploy a similar structure by considering:

$$(\mu_e(z), \sigma_e^2(z)) = \text{Decoder}_\theta(z),$$

$$p_\theta(x | z) = \mathcal{N}(x | \mu_e(z), \sigma_e^2(z)I).$$

Again, for the purpose of discussion, we assume that μ_e is affine so that: $\mu_e(z) = cz + v$ for some parameters c and v , and:

$$\sigma_e^2(z) = s^2,$$

for some scalar s . Therefore, $p_\theta(x | z)$ takes the form of:

$$p_\theta(x | z) = \mathcal{N}(x | cz + v, s^2I).$$



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

V. Evidence Lower Bound



Why?

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Variational Autoencoders (VAEs) are designed to model observed data x through a latent variable z . The goal is to maximize the marginal likelihood of the data $p(x)$:

$$p(x) = \int p(x, z) dz,$$

where $p(x, z) = p(x | z)p(z)$.

Although using $p(x | z)$ and $p(z)$ simplifies the complex distribution $p(x, z)$, the integral is difficult to compute, especially when:

- The latent space z is very large or continuous.
- $p(x | z)$ and $p(z)$ are nonlinear or complex.

When the space of z is continuous, computing the integral over the entire latent space becomes impractical due to:

- The computational complexity.
- The high dimensionality of z , which makes the processing infeasible.



Define ELBO

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Theorem 1.1. Decomposition of Log-Likelihood. The log-likelihood $\log p(x)$ can be decomposed as

$$\log p(x) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] + D_{\text{KL}}(q_{\phi}(z|x) \| p(z|x)).$$



Proof

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

The trick is to use our magical proxy $q_{\phi}(z|x)$ to poke around $p(x)$ and derive the bound:

$$\begin{aligned}\log p(x) &= \log p(x) \times \underbrace{\int q_{\phi}(z|x) dz}_{=1} \\ &= \int \underbrace{\log p(x)}_{\text{constant w.r.t. } z} \times \underbrace{q_{\phi}(z|x)}_{\text{distribution in } z} dz \quad (\text{move } \log p(x) \text{ into the integral}) \\ &= \mathbb{E}_{q_{\phi}(z|x)}[\log p(x)].\end{aligned}$$



Kullback-Leibler Divergence and Derivation

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Define the Kullback-Leibler Divergence between two continuous distributions:

$$\begin{aligned} D_{\text{KL}}(P \parallel Q) &= \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx \\ &= \mathbb{E}_{p(x)} \left[\log \frac{p(x)}{q(x)} \right]. \end{aligned}$$



An example of KL

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

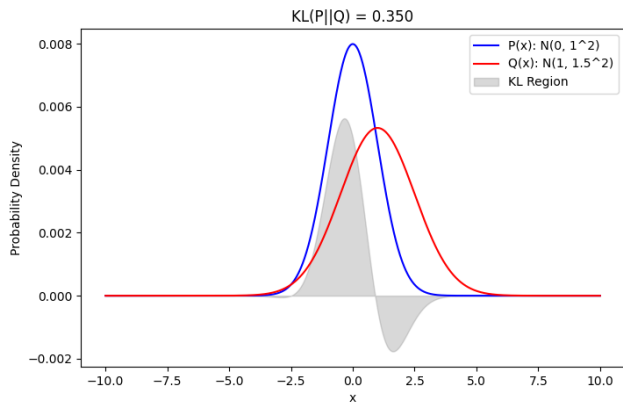


Figure 5: An example of KL



Deriving ELBO with Bayes' Theorem

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Using Bayes' theorem: $p(x, z) = p(z|x)p(x)$

$$\begin{aligned}\mathbb{E}_{q_{\phi}(z|x)}[\log p(x)] &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{p(z|x)} \right] \text{ (Bayes' theorem)} \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{p(z|x)} \times \frac{q_{\phi}(z|x)}{q_{\phi}(z|x)} \right] \text{ (Multiply and divide by } q_{\phi}(z|x)) \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] + \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{q_{\phi}(z|x)}{p(z|x)} \right] \\ &= \text{ELBO}(x) + D_{\text{KL}}(q_{\phi}(z|x) \| p(z|x)).\end{aligned}\tag{8}$$



Decompose ELBO

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Theorem 1.2

We can decompose ELBO as:

$$\text{ELBO}(x) = \mathbb{E}_{q_{\phi}(z|x)} \left[\underbrace{\log p_{\theta}(x|z)}_{\text{a Gaussian (how good of decoder)}} \right] - \underbrace{D_{\text{KL}}(q_{\phi}(z|x) \| p(z))}_{\text{a Gaussian (how good of encoder)}}.$$



Proof

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

$$\text{ELBO}(x) \stackrel{\text{def}}{=} \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right] \quad (\text{definition})$$

$$= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x|z)p(z)}{q_{\phi}(z|x)} \right] \quad (\text{since } p(x, z) = p(x|z)p(z))$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\log p(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(z)}{q_{\phi}(z|x)} \right] \quad (\text{split expectation})$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{(\text{KL})}(q_{\phi}(z|x) \| p(z)), \quad (\text{definition of KL})$$

where we replaced the inaccessible $p(x|z)$ by its proxy $p_{\theta}(x|z)$.



Key Objectives in Variational Autoencoders

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Reconstruction:

- The first term focuses on the decoder. We want the decoder to produce a good reconstruction of x when we feed a latent z into it.
- Our goal is to maximize $p_{\theta}(x|z)$, which is similar to maximum likelihood, where we optimize the model parameters to maximize the likelihood of observing the input x .



Key Objectives in Variational Autoencoders

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Prior Matching:

- The second term is the KL divergence for the encoder. We want the encoder to map x into a latent vector z such that z follows a chosen distribution:

$$z \sim \mathcal{N}(0, I).$$

- More generally, we define $p(z)$ as the target distribution. Since KL divergence measures the distance between two distributions, minimizing it ensures the latent space $q_\phi(z|x)$ becomes similar to $p(z)$.
- **Note:** A negative sign is used in the objective so that minimizing KL divergence aligns the distributions.



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

VI. Optimization in VAE



Definition

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

The optimization objective of the VAE is to maximize the ELBO:

$$(\phi, \theta) = \arg \max_{\phi, \theta} \sum_{x \in \mathcal{X}} \text{ELBO}(x),$$

where $\mathcal{X} = \{x^{(\ell)} \mid \ell = 1, \dots, L\}$ is the training dataset.

To be more precisely, we can show this objective function:

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



Gradient with parameter θ

Let's first look at θ . We can show that:

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

Using Monte Carlo approximation, we have:

$$\nabla_{\theta} \text{ELBO}(x) \approx \frac{1}{L} \sum_{\ell=1}^L \nabla_{\theta} \log p_{\theta}(x, z^{(\ell)}),$$

where $z^{(\ell)} \sim q_{\phi}(z|x)$.

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE



Gradient with parameter ϕ

The gradient with respect to ϕ is more difficult. We can show that:

$$\begin{aligned}\nabla_{\phi} \text{ELBO}(x) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] \\ &= \nabla_{\phi} \int [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] q_{\phi}(z|x) dz \\ &= \int \nabla_{\phi} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] q_{\phi}(z|x) dz \\ &\neq \int [\nabla_{\phi} (\log p_{\theta}(x, z) - \log q_{\phi}(z|x)) q_{\phi}(z|x)] dz \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\phi} (\log p_{\theta}(x, z) - \log q_{\phi}(z|x))] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\phi} (-\log q_{\phi}(z|x))] \\ &\approx \frac{1}{L} \sum_{\ell=1}^L \nabla_{\phi} \left(-\log q_{\phi}(z^{(\ell)}|x) \right), \text{ where } z^{(\ell)} \sim q_{\phi}(z|x).\end{aligned}$$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE



Reparameterization Trick

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

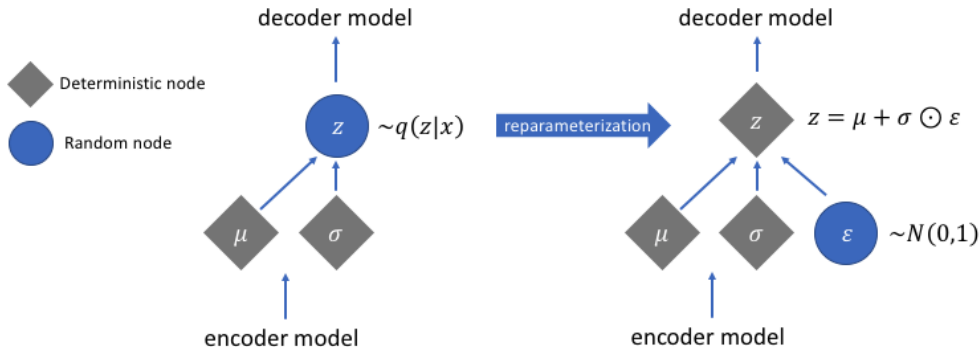


Figure 6: Comparison between Original form and Reparameterized form



Reparameterization Trick

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

As we can see, $z^{(\ell)} \sim q_{\phi}(z|x)$ represents a random node, symbolizing the stochastic sampling of the value z .

However, $q_{\phi}(z|x)$ is parameterized by a neural network (the Encoder network), which makes it difficult to compute a closed-form solution.

Thus, we need to adopt a new approach to address this problem: the **Reparameterization Trick**.



Define

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Suppose $z \sim q_\phi(z|x) \stackrel{\text{def}}{=} \mathcal{N}(z|\mu, \text{diag}(\sigma^2))$. We can define

$$z = g(\epsilon, \phi, x) \stackrel{\text{def}}{=} \epsilon \odot \sigma + \mu, \quad (1)$$

where $\epsilon \sim \mathcal{N}(0, I)$ and “ \odot ” means elementwise multiplication. The parameter ϕ is $\phi = (\mu, \sigma^2)$. For this choice of the distribution, we can show that by letting $\epsilon = \frac{z - \mu}{\sigma}$:

$$\begin{aligned} q_\phi(z|x) \cdot \left| \det \left(\frac{\partial z}{\partial \epsilon} \right) \right| &= \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left\{ -\frac{(z_i - \mu_i)^2}{2\sigma_i^2} \right\} \cdot \prod_{i=1}^d \sigma_i \\ &= \frac{1}{(2\pi)^{d/2}} \exp \left\{ -\frac{\|\epsilon\|^2}{2} \right\} = \mathcal{N}(0, I) = p(\epsilon). \end{aligned}$$



Jacobian of the Reparameterization Trick

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

$\frac{\partial \mathbf{z}}{\partial \epsilon}$ is the Jacobian, we can define that:

$$\frac{\partial \mathbf{z}}{\partial \epsilon} = \begin{bmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \frac{\partial z_1}{\partial \epsilon_2} & \dots & \frac{\partial z_1}{\partial \epsilon_d} \\ \frac{\partial z_2}{\partial \epsilon_1} & \frac{\partial z_2}{\partial \epsilon_2} & \dots & \frac{\partial z_2}{\partial \epsilon_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_d}{\partial \epsilon_1} & \frac{\partial z_d}{\partial \epsilon_2} & \dots & \frac{\partial z_d}{\partial \epsilon_d} \end{bmatrix}.$$

and expression in each component i is:

$$z_i = \sigma_i \epsilon_i + \mu_i$$



Gradient Computation with Reparameterization

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

We aim to compute the gradient $\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)}[f(z)]$ for some general function $f(z)$. (We will consider $f(z) = -\log q_{\phi}(z|x)$.)

For simplicity, we write $g(\epsilon)$ instead of $g(\epsilon, \phi, x)$, even though g has three inputs.

Using the change of variables, we have:

$$\mathbb{E}_{q_{\phi}(z|x)}[f(z)] = \int f(z) \cdot q_{\phi}(z|x) dz = \int f(g(\epsilon)) \cdot q_{\phi}(g(\epsilon)|x) dg(\epsilon), \quad \text{where } z = g(\epsilon).$$

Incorporating the Jacobian determinant:

$$\mathbb{E}_{q_{\phi}(z|x)}[f(z)] = \int f(g(\epsilon)) \cdot q_{\phi}(g(\epsilon)|x) \cdot \det\left(\frac{\partial g(\epsilon)}{\partial \epsilon}\right) d\epsilon.$$

Simplification:

$$\mathbb{E}_{q_{\phi}(z|x)}[f(z)] = \int f(z) \cdot p(\epsilon) d\epsilon = \mathbb{E}_{p(\epsilon)}[f(z)].$$



Gradient Computation for ϕ

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

To compute the gradient with respect to ϕ , we write:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)}[f(z)] = \nabla_{\phi} \mathbb{E}_{p(\epsilon)}[f(z)] = \nabla_{\phi} \int f(z) \cdot p(\epsilon) d\epsilon.$$

Since $p(\epsilon)$ does not depend on ϕ , the gradient only acts on $f(z)$:

$$\nabla_{\phi} \mathbb{E}_{p(\epsilon)}[f(z)] = \int \nabla_{\phi} \{f(z) \cdot p(\epsilon)\} d\epsilon = \int \{\nabla_{\phi} f(z)\} \cdot p(\epsilon) d\epsilon.$$

Thus:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)}[f(z)] = \mathbb{E}_{p(\epsilon)}[\nabla_{\phi} f(z)].$$

Substitution for $f(z) = -\log q_{\phi}(z|x)$:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)}[-\log q_{\phi}(z|x)] = \mathbb{E}_{p(\epsilon)}[-\nabla_{\phi} \log q_{\phi}(z|x)].$$



Monte Carlo Approximation and Gradient Simplification

Monte Carlo Approximation:

$$\mathbb{E}_{p(\epsilon)}[-\nabla_{\phi} \log q_{\phi}(z|x)] \approx -\frac{1}{L} \sum_{\ell=1}^L \nabla_{\phi} \log q_{\phi}(z^{(\ell)}|x),$$

where $z^{(\ell)} = g(\epsilon^{(\ell)}, \phi, x)$.

Rewriting $q_{\phi}(z|x)$:

$$q_{\phi}(z|x) = p(\epsilon) : \det\left(\frac{\partial z}{\partial \epsilon}\right), \quad \log q_{\phi}(z|x) = \log p(\epsilon) - \log \det\left(\frac{\partial z}{\partial \epsilon}\right).$$

Gradient Simplification: Since $p(\epsilon)$ is independent of ϕ :

$$\nabla_{\phi} \log q_{\phi}(z|x) = \nabla_{\phi} \log \det\left(\frac{\partial z}{\partial \epsilon}\right).$$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE



Gradient Simplification for the Encoder

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Thus:

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)}[-\log q_{\phi}(z|x)] &\approx \frac{1}{L} \sum_{\ell=1}^L \nabla_{\phi} \log \det \left(\frac{\partial z^{(\ell)}}{\partial \epsilon^{(\ell)}} \right) \\ &= \frac{1}{L} \sum_{\ell=1}^L \nabla_{\phi} \left[\sum_{i=1}^d \log \sigma_i \right] \\ &= \nabla_{\phi} \left[\sum_{i=1}^d \log \sigma_i \right] \\ &= \frac{1}{\sigma} \odot \nabla_{\phi} [\sigma_{\phi}(x)].\end{aligned}$$

Note: $\sigma_{\phi}(x)$ is the output of the encoder, which is parameterized by a neural network.



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

VII. Specific Structure of VAE



Encoder

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

$$(\mu, \sigma^2) = \text{EncoderNetwork}_\phi(x)$$

$$q_\phi(z|x) = \mathcal{N}(z|\mu, \sigma^2 I)$$

The parameters μ and σ are technically neural networks because they are outputs of $\text{EncoderNetwork}_\phi(\cdot)$.

We denote them as:

$$\mu = \mu_\phi(x) \quad \text{and} \quad \sigma^2 = \sigma_\phi^2(x).$$

Given the ℓ -th training sample $x^{(\ell)}$, the latent variable $z^{(\ell)}$ is sampled as:

$$z^{(\ell)} \sim \mathcal{N}(z | \mu_\phi(x^{(\ell)}), \sigma_\phi^2(x^{(\ell)}) I).$$



Sampling with Reparameterization Trick

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

From the Gaussian, we draw a sample $z^{(\ell)}$:

$$z^{(\ell)} = \mu_\phi(x^{(\ell)}) + \sigma_\phi(x^{(\ell)})\epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$

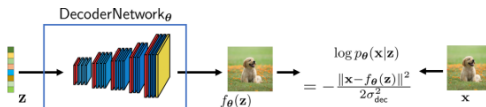


Figure 7: Implementation of a VAE encoder. The neural network takes x and estimates μ_ϕ and σ_ϕ^2 .



KL Divergence of Two Gaussians

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Theorem 1.3: KL-Divergence of Two Gaussians

The KL divergence for two d -dimensional Gaussian distributions $\mathcal{N}(\mu_0, \Sigma_0)$ and $\mathcal{N}(\mu_1, \Sigma_1)$ is:

$$D_{\text{KL}}(\mathcal{N}(\mu_0, \Sigma_0) \parallel \mathcal{N}(\mu_1, \Sigma_1)) = \frac{1}{2} \left[\text{Tr}(\Sigma_1^{-1} \Sigma_0) - d + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) + \log \frac{\det \Sigma_1}{\det \Sigma_0} \right].$$



KL Divergence in VAE

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Substituting Specific Distributions:

$$\mu_0 = \mu_\phi(x), \quad \Sigma_0 = \sigma_\phi^2(x)\mathbf{I}, \quad \mu_1 = 0, \quad \Sigma_1 = \mathbf{I}.$$

Simplified KL Divergence:

$$D_{\text{KL}}(q_\phi(z|x) \parallel p(z)) = \frac{1}{2} \left[\sigma_\phi^2(x)d - d + \|\mu_\phi(x)\|^2 - 2\log \sigma_\phi(x) \right].$$

Here, d is the dimension of the latent vector z .



Decoder

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

The decoder network, $\text{DecoderNetwork}_\theta(\cdot)$, maps a latent variable z to a generated image:

$$f_\theta(z) = \text{DecoderNetwork}_\theta(z).$$

The distribution $p_\theta(x | z)$ is:

$$p_\theta(x | z) = \mathcal{N}(x | f_\theta(z), \sigma_{\text{dec}}^2 I),$$

where σ_{dec} is a hyperparameter. Using the reparameterization trick:

$$\hat{x} = f_\theta(z) + \sigma_{\text{dec}}\epsilon, \quad \epsilon \sim \mathcal{N}(0, I).$$



Log-Likelihood of $p_{\theta}(x | z)$

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

The log-likelihood can be written as:

$$\begin{aligned}\log p_{\theta}(x|z) &= \log \mathcal{N}(x|f_{\theta}(z), \sigma_{\text{dec}}^2 I) \\ &= \log \left(\frac{1}{\sqrt{(2\pi\sigma_{\text{dec}}^2)^d}} \right) \exp \left(-\frac{\|x - f_{\theta}(z)\|^2}{2\sigma_{\text{dec}}^2} \right) \\ &= -\frac{\|x - f_{\theta}(z)\|^2}{2\sigma_{\text{dec}}^2} - \log \left((2\pi\sigma_{\text{dec}}^2)^d \right)\end{aligned}$$

The constant term $\log \sqrt{(2\pi\sigma_{\text{dec}}^2)^d}$ is independent of θ and can be dropped.



ELBO with Monte Carlo Approximation

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

Using the reparameterization trick, the expectation is approximated as:

$$\begin{aligned}\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] &= - \int \frac{\|x - f_{\theta}(z)\|^2}{2\sigma_{\text{dec}}^2} \cdot \mathcal{N}(z | \mu_{\phi}(x), \sigma_{\phi}(x)^2) dz \\ &\approx - \frac{1}{M} \sum_{m=1}^M \frac{\|x - f_{\theta}(z^{(m)})\|^2}{2\sigma_{\text{dec}}^2},\end{aligned}$$

where:

$$z^{(m)} = \mu_{\phi}(x) + \sigma_{\phi}(x)\epsilon^{(m)}, \quad \epsilon^{(m)} \sim \mathcal{N}(0, \mathbf{I}).$$



VAE Training: Optimization Objective

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

To train a VAE, we need to solve the optimization problem:

$$\arg \max_{\theta, \phi} \sum_{x \in \mathcal{X}} \text{ELBO}_{\phi, \theta}(x),$$

where:

$$\text{ELBO}_{\phi, \theta}(x) = \text{Reconstruction Term} + \text{KL Divergence Term}.$$

Reconstruction Term:

$$-\frac{1}{M} \sum_{m=1}^M \frac{\|x - f_{\theta}(\mu_{\phi}(x) + \sigma_{\phi}(x)\epsilon^{(m)})\|^2}{2\sigma_{\text{dec}}^2}.$$



VAE Training: ELBO Expression

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

KL Divergence Term:

$$+\frac{1}{2}\left(\sigma_{\phi}^2(x)d - d + \|\mu_{\phi}(x)\|^2 - 2\log\sigma_{\phi}(x)\right).$$

Explanation:

- $\mu_{\phi}(x)$: Mean output of the encoder.
- $\sigma_{\phi}(x)$: Standard deviation output of the encoder.
- d : Dimension of the latent variable z .



I. Introduction

II. Structure of
VAE

III.
Mathematical
Symbols in VAE

IV. Build
construct in
VAE

V. Evidence
Lower Bound

VI. Optimization
in VAE

VII. Specific
Structure of
VAE

VIII.
Applications of
VAE

VIII. Applications of VAE



Applications of VAE: Image Generation

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- VAE can generate new images similar to a given dataset.
- Example: Generating new handwritten digits using the MNIST dataset.
- VAE encodes the input into a latent space and decodes it back to generate new images.



I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

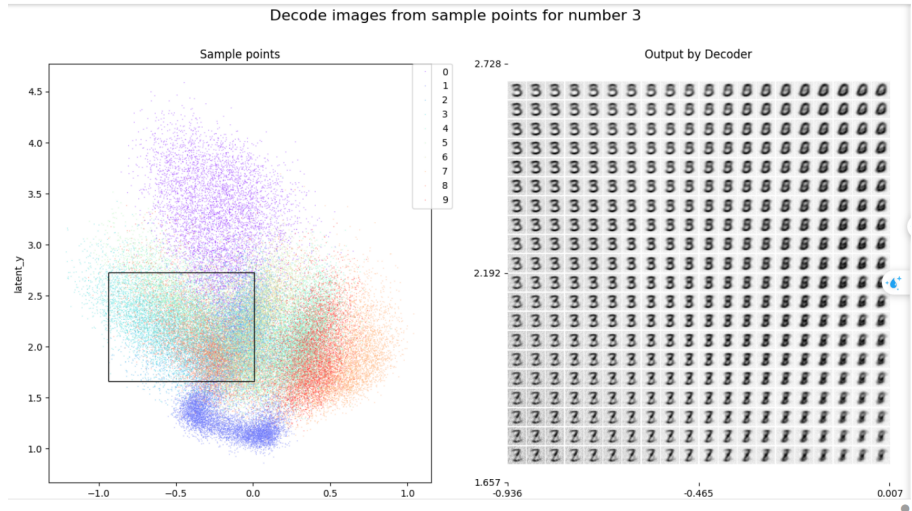
IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE





Applications of VAE: Anomaly Detection

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- VAE learns the normal distribution of data.
- Any data point with low reconstruction probability is flagged as an anomaly.
- Applications:
 - Fraud detection.
 - Industrial defect detection.



Applications of VAE: Data Compression

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- The latent space z provides a compressed representation of the data.
- Applications:
 - Reducing storage requirements while retaining meaningful features.
 - Efficient transmission of high-dimensional data.



Applications of VAE: Medicine and Biology

I. Introduction

II. Structure of VAE

III. Mathematical Symbols in VAE

IV. Build construct in VAE

V. Evidence Lower Bound

VI. Optimization in VAE

VII. Specific Structure of VAE

VIII. Applications of VAE

- Generating molecular structures or drug compounds.
- Applications:
 - Drug discovery.
 - Protein structure prediction.