

# Variational Autoencoders (VAEs)

Allen C.

September 2023

## 1 Autoencoder

Illustration of an autoencoder:

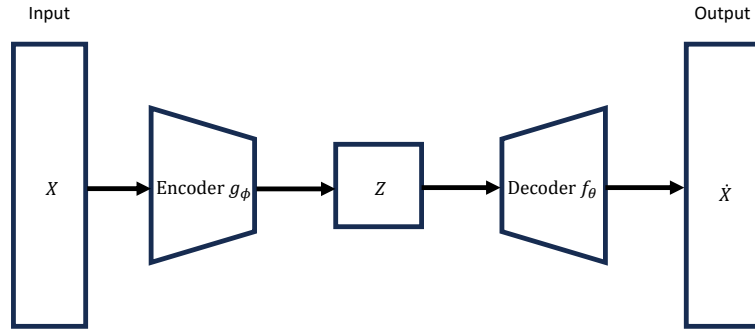


Figure 1: Autoencoder

As an applied mathematician, we understand the autoencoder model as follows:

- $X$  is a vector.
- $g_\phi$  can be a matrix and nonlinear functions (ReLU, tanh).
- $z$  is a reduced vector.
- $f_\theta$  can be a matrix and nonlinear functions (ReLU, tanh) which play the inverse role of  $g_\phi$ .
- $\hat{X}$  is a vector.

The loss function of the autoencoder can be defined as

$$L(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n [X^i - f_\theta(g_\phi(X^i))]^2 \quad (1)$$

Here, the loss function compares  $X$  and  $\hat{X}$  with the  $L^2$  norm.

## 2 Variational Autoencoder

Illustration of a Variational Autoencoder:

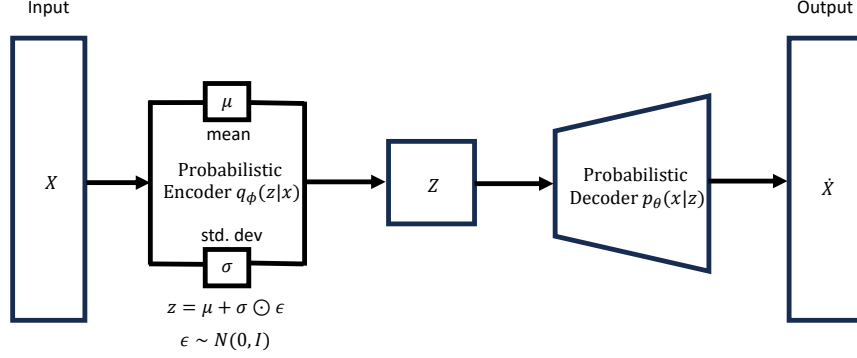


Figure 2: Variational Autoencoder

The loss function for the Variational Autoencoder model is given as follows, and it will be derived later:

$$L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + D_{KL}(q_\phi(z|x) || p_\theta(z)) \quad (2)$$

## 3 Probability

- $P(x)$ : defines the probability of random variable  $x$ .
- $P(x|y)$ : defines the probability of random variable  $x$  provided  $y$  has happened. Also called conditional probability.
- Bayes's Theorem:  $P(y|x) = \frac{P(x|y)P(y)}{P(x)}$

## 4 Kullback-Leibler Divergence

KL divergence is a measure of how one probability distribution is different from the other.

$$D_{KL}(P||Q) = \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) \quad (3)$$

Properties of KL divergence:

- $KL(P||Q) \geq 0$  or  $KL(Q||P) \geq 0$ .
- $KL(P||Q) \neq KL(Q||P)$

**Lemma 1.** Given a multivariate normal density function as follows:

$$N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) \quad (4)$$

where  $k$  is the dimension of  $x$  and  $\Sigma$ . With  $p(x) = N(x; \mu_1, \Sigma_1)$  and  $q(x) = N(x; \mu_2, \Sigma_2)$ , where  $\mu_1$  and  $\mu_2$  are the means and  $\Sigma_1$  and  $\Sigma_2$  are the covariance matrices, the KL divergence for  $p(x)$  and  $q(x)$  can be written as follows:

$$D_{KL}(p(x)||q(x)) = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - k + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_1 - \mu_2)^\top \Sigma_2^{-1}(\mu_1 - \mu_2) \right] \quad (5)$$

The above equation will be applied to rewrite the loss function of VAE.

*Proof.* With the definition of KL divergence, we have

$$KL(p(x)||q(x)) = \sum_x p(x) \log \left( \frac{p(x)}{q(x)} \right)$$

The probability density functions  $p(x)$  and  $q(x)$  are given by

$$p(x) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_1|}} \exp \left( -\frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) \right)$$

$$q(x) = \frac{1}{\sqrt{(2\pi)^k |\Sigma_2|}} \exp \left( -\frac{1}{2} (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2) \right)$$

Then we have

$$\log p(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1)$$

$$\log q(x) = -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2)$$

Then KL divergence can be rewritten as

$$\begin{aligned} KL(p(x)||q(x)) &= \sum_x p(x) (\log p(x) - \log q(x)) \\ &= \sum_x p(x) \left[ -\frac{k}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) \right. \\ &\quad \left. + \frac{k}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma_2| + \frac{1}{2} (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2) \right] \\ &= \sum_x p(x) \left[ \frac{1}{2} \log \left( \frac{|\Sigma_2|}{|\Sigma_1|} \right) + \frac{1}{2} (x - \mu_2)^\top \Sigma_2^{-1} (x - \mu_2) - \frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) \right] \end{aligned} \quad (6)$$

Consider the above right hand side part by part:

$$\sum_x p(x) \frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) = E_{p(x)} \left[ \frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1) \right]$$

To prove the lemma, we need to use the following trace and expectation identity:

$$\begin{aligned} E(X^\top A X) &= E(\text{tr}(X^\top A X)) \\ &= E(\text{tr}(A X X^\top)) \\ &= \text{tr}(E(A X X^\top)) \end{aligned}$$

Then we have

$$\begin{aligned}
& \frac{1}{2}E_{p(x)}\left[(x - \mu_1)^\top \Sigma_1^{-1}(x - \mu_1)\right] \\
&= E_{p(x)}\left[\text{tr}\left(\frac{1}{2}(x - \mu_1)^\top \Sigma_1^{-1}(x - \mu_1)\right)\right] \\
&= E_{p(x)}\left[\text{tr}\left(\frac{1}{2}(x - \mu_1)(x - \mu_1)^\top \Sigma_1^{-1}\right)\right] \\
&= \text{tr}\left[E_{p(x)}\left(\frac{1}{2}(x - \mu_1)(x - \mu_1)^\top \Sigma_1^{-1}\right)\right] \\
&= \text{tr}\left[E_{p(x)}\left((x - \mu_1)(x - \mu_1)^\top\right)\frac{1}{2}\Sigma_1^{-1}\right] \\
&= \text{tr}\left[\Sigma_1\frac{1}{2}\Sigma_1^{-1}\right] \\
&= \frac{k}{2}
\end{aligned}$$

Then for the second term of the right hand side of Equation (6) can be rewritten as:

$$\begin{aligned}
& \sum_x p(x) \left[\frac{1}{2}(x - \mu_2)^\top \Sigma_2^{-1}(x - \mu_2)\right] \\
&= \sum_x p(x) \left[\frac{1}{2}[(x - \mu_1) + (\mu_1 - \mu_2)]^\top \Sigma_2^{-1}[(x - \mu_1) + (\mu_1 - \mu_2)]\right] \\
&= \sum_x p(x) \left[\frac{1}{2}(x - \mu_1)^\top \Sigma_2^{-1}(x - \mu_1) + (x - \mu_1)^\top \Sigma_2^{-1}(\mu_1 - \mu_2) + \frac{1}{2}(\mu_1 - \mu_2)^\top \Sigma_2^{-1}(\mu_1 - \mu_2)\right] \\
&= E_{p(x)}\left[\frac{1}{2}(x - \mu_1)^\top \Sigma_2^{-1}(x - \mu_1) + (x - \mu_1)^\top \Sigma_2^{-1}(\mu_1 - \mu_2) + \frac{1}{2}(\mu_1 - \mu_2)^\top \Sigma_2^{-1}(\mu_1 - \mu_2)\right]
\end{aligned}$$

Since we have the following identity:

$$\begin{aligned}
& E_{p(x)}\left[(x - \mu_1)^\top \Sigma_2^{-1}(\mu_1 - \mu_2)\right] \\
&= \left(E_{p(x)}x - \mu_1\right)^\top \Sigma_2^{-1}(\mu_1 - \mu_2) \\
&= (\mu_1 - \mu_1)^\top \Sigma_2^{-1}(\mu_1 - \mu_2) \\
&= 0
\end{aligned}$$

Then we have

$$\begin{aligned}
& \sum_x p(x) \left[\frac{1}{2}(x - \mu_2)^\top \Sigma_2^{-1}(x - \mu_2)\right] \\
&= \text{tr}\left(\frac{\Sigma_2^{-1}\Sigma_1}{2}\right) + (\mu_1 - \mu_2)^\top \Sigma_2^{-1}(\mu_1 - \mu_2)
\end{aligned}$$

To conclude, we prove the lemma:

$$D_{K|L}(p(x)||q(x)) = \frac{1}{2}\left[\log\frac{|\Sigma_2|}{|\Sigma_1|} - k + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_1 - \mu_2)^\top \Sigma_2^{-1}(\mu_1 - \mu_2)\right]$$

□

In the following, we will derive the loss function of VAE:

$$L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + D_{K|L}(q_\phi(z|x)||p_\theta(z))$$

We have Baye's theorem:

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)} \quad (7)$$

Here  $p(x)$  can be written as:

$$p(x) = \int_z p(x, z) dz \quad (8)$$

which is difficult to integrate if  $z$  is in high dimension. To obtain  $p(z|x)$ , the alternative way is to approximate  $p(z|x)$  by another distribution  $q(z|x)$  which is defined in a way that has tractable solution, such as Gaussian distribution. As discussed above, we use KL divergence to measure the 'distance' between  $p(z|x)$  and  $q(z|x)$ :

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ &= \sum_z q_\phi(z|x) \log \left( \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \\ &= E_{z \sim q_\phi(z|x)} \left[ \log \left( \frac{q_\phi(z|x)}{p_\theta(z|x)} \right) \right] \\ &= E_{z \sim q_\phi(z|x)} \left[ \log q_\phi(z|x) - \log p_\theta(z|x) \right] \end{aligned} \quad (9)$$

Substituting Equation (7) into Equation (9) gives:

$$\begin{aligned} & KL(q_\phi(z|x)||p_\theta(z|x)) \\ &= E_{z \sim q_\phi(z|x)} \left[ \log q_\phi(z|x) - \log \frac{p_\theta(x|z)p_\theta(z)}{p_\theta(x)} \right] \\ &= E_{z \sim q_\phi(z|x)} \left[ \log q_\phi(z|x) - \log p_\theta(x|z) - \log p_\theta(z) + \log p_\theta(x) \right] \end{aligned} \quad (10)$$

Then we have

$$KL(q_\phi(z|x)||p_\theta(z|x)) - \log p_\theta(x) = E_{z \sim q_\phi(z|x)} \left[ \log q_\phi(z|x) - \log p_\theta(x|z) - \log p_\theta(z) \right] \quad (11)$$

Multiplying -1 at both sides of Equation (11) gives:

$$\begin{aligned} & \log p_\theta(x) - KL(q_\phi(z|x)||p_\theta(z|x)) \\ &= E_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - E_{z \sim q_\phi(z|x)} \left[ \log q_\phi(z|x) - \log p_\theta(z) \right] \\ &= E_{z \sim q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - D_{KL}[q_\phi(z|x)||p_\theta(z)] \end{aligned} \quad (12)$$

If we optimize the right hand side of Equation (12), it implies to optimize  $\log p_\theta(x)$  and minimize  $KL(q_\phi(z|x)||p_\theta(z|x))$ . Then the loss function can be defined as:

$$L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + D_{KL}(q_\phi(z|x)||p_\theta(z)) \quad (13)$$

where  $q_\phi(z|x)$  is a neural network - encoder;  $p_\theta(x|z)$  is a neural network - decoder (it could be a normal distribution).  $p_\theta(z)$  is the latent variable distribution. The simple choice is  $N(0, 1)$  (it can also be any distribution: flower shaped distribution). Then  $D_{KL}$  can be simplified as:

$$\begin{aligned} & D_{KL}(N(\mu_\phi(x), \Sigma_\phi(x))||N(0, 1)) \\ &= \frac{1}{2} \left[ \text{tr}(\Sigma_\phi(x)) + \mu_\phi(x)^\top \mu_\phi(x) - k - \log |\Sigma_\phi(x)| \right] \end{aligned}$$

where  $k$  is the dimension of the Gaussian.  $tr(\Sigma_\phi(x))$  is the trace function, which is the sum of the diagonal matrix of  $\Sigma_\phi(x)$ . Furthermore the determinant of  $\Sigma_\phi(x)$  is the product of a diagonal matrix. So the above equation can be written as

$$\begin{aligned} & D_{KL}(N(\mu_\phi(x), \Sigma_\phi(x)) || N(0, 1)) \\ &= \frac{1}{2} \left( \sum_k \Sigma_\phi(x) + \sum_k \mu_\phi(x)^2 + \sum_k 1 - \log \left( \prod_k \Sigma_\phi(x) \right) \right) \\ &= \frac{1}{2} \sum_k \left( \Sigma_\phi(x) + \mu_\phi(x)^2 - 1 - \log \Sigma_\phi(x) \right) \end{aligned}$$

In practice, we use  $\exp(\Sigma_\phi(x))$  instead of  $\Sigma_\phi(x)$ . Then the KL divergence can be written as:

$$\begin{aligned} & D_{KL}(N(\mu_\phi(x), \Sigma_\phi(x)) || N(0, 1)) \\ &= \frac{1}{2} \sum_k \left( \exp(\Sigma_\phi(x)) + \mu_\phi(x)^2 - 1 - \Sigma_\phi(x) \right) \end{aligned}$$

Correspondingly the loss function is written as

$$L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + \frac{1}{2} \sum_k \left( \exp(\Sigma_\phi(x)) + \mu_\phi(x)^2 - 1 - \Sigma_\phi(x) \right) \quad (14)$$

## 5 Optimization of the loss function

Find  $\theta^*, \phi^* = \operatorname{argmin}_{\theta, \phi} L(\theta, \phi)$ . In variational Bayesian method, this loss function is known as the variaional lower bound or evidence lower bound (ELBD). We know that

$$L(\theta, \phi) \leq \log p_\theta(x) \quad (15)$$

Recall the loss function

$$L(\theta, \phi) = -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) + \frac{1}{2} \sum_k \left( \exp(\Sigma_\phi(x)) + \mu_\phi(x)^2 - 1 - \Sigma_\phi(x) \right) \quad (16)$$

$$\theta^* = \theta - \eta \nabla_\theta L(\theta, \phi)$$

$$\phi^* = \phi - \eta \nabla_\phi L(\theta, \phi)$$

where

$$\begin{aligned} & \nabla_\theta L(\theta, \phi) \\ &= \nabla_\theta \left[ -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) \right] + \nabla_\theta \left[ \frac{1}{2} \sum_k \left( \exp(\Sigma_\phi(x)) + \mu_\phi(x)^2 - 1 - \Sigma_\phi(x) \right) \right] \\ &= \frac{1}{L} \sum_{l=1}^L \nabla_\theta \log p_\theta(x|z^l) (\text{Monte Carlo estimation}) + 0 \end{aligned}$$

where  $z^l \sim q_\phi(z|x)$ . For the derivative with respect to  $\phi$ ,  $\nabla_\phi L(\theta, \phi)$  can be written as

$$\nabla_\phi L(\theta, \phi) = \nabla_\phi \left[ -E_{z \sim q_\phi(z|x)} \log(p_\theta(x|z)) \right] + \nabla_\phi \left[ \frac{1}{2} \sum_k \left( \exp(\Sigma_\phi(x)) + \mu_\phi(x)^2 - 1 - \Sigma_\phi(x) \right) \right]$$

For the second term of the right hand side, it is easy to calculate. However for the first term of the right hand side, we face some challenge, since the derivative for the first term is difficult to take to the sample points.

$$\nabla_{\phi} E_{z \sim q_{\phi}(z|x)} f(z) \neq E_{z \sim q_{\phi}(z|x)} \nabla_{\phi} f(z)$$

To overcome this, the expectation is rewritten by linear transformation:

$$E_{z \sim q_{\phi}(z|x)} f(z) = E_{p(\epsilon)} f(g_{\phi}(\epsilon, x))$$

where  $z = g_{\phi}(\epsilon, x)$ , with  $\epsilon \sim N(0, 1)$  and

$$g_{\phi}(\epsilon, x) = \mu_{\phi}(x) + \epsilon \odot \Sigma_{\phi}^{1/2}(x) = z \sim N(\mu(x), \Sigma(x))$$

This transformation makes sure that the distribution is independent of  $\phi$ :

$$\begin{aligned} & \nabla_{\phi} L(\theta, \phi) \\ &= \nabla_{\phi} \left[ - E_{z^l \sim p(\epsilon)} \log(p_{\theta}(x|z^l)) \right] + \nabla_{\phi} \left[ \frac{1}{2} \sum_k \left( \exp(\Sigma_{\phi}(x)) + \mu_{\phi}(x)^2 - 1 - \Sigma_{\phi}(x) \right) \right] \\ &= - \frac{1}{S} \sum_{l=1}^S \nabla_{\phi} \log p_{\theta}(x|z^l) + \nabla_{\phi} \left[ \frac{1}{2} \sum_k \left( \exp(\Sigma_{\phi}(x)) + \mu_{\phi}(x)^2 - 1 - \Sigma_{\phi}(x) \right) \right] \end{aligned}$$

If the loss function exclusively incorporates the log entropy term, VAEs will generate only one value (vector) for each input. Conversely, if the loss function solely incorporates the KL divergence term, VAEs will attempt to approximate  $q_{\phi}(z|x)$  to be as close to the normal distribution  $N(0, 1)$  as possible. This can potentially result in image blurriness and hinder effective learning from the training data.