

---

# VARIATIONAL AUTO ENCODER: A DERIVATION OF GAUSSIAN AND INVERSE-WISHART EVIDENCE LOWER BOUND'S (ELBO'S)

---

**Kenneth Ezukwoke\***

École des Mines de Saint-Etienne  
Univ. Clermont Auvergne, CNRS UMR 6158 LIMOS  
Mathematics and Industrial Engineering  
Henri FAYOL institute, F-42023, Saint-Etienne, France  
{ifeanyi.ezukwoke}@emse.fr

## ABSTRACT

Bayesian Inference has recently gained traction as an efficient method for text and image generation. Much of the success of this model is attributed to its rather unique loss function used for the generative process. In this paper we present the mathematical details behind its inference method and derive the Evidence lower bounds (ELBO) both for Gaussian distributions and Inverse-Wishart distributions.

**Keywords** Variational AutoEncoder; Evidence Lower Bound; Gaussian distribution; Inverse-Wishart distribution.

## 1 Introduction

## 2 Variational AutoEncoder (VAE)

VAE [1] is a generative model and a derivative of Bayesian inference. Unlike autoencoder, VAE are more suited for textual data and have been used for unsupervised abstractive sentence summarization [2] and long and coherent Text generation [3] amongst others. VAE specifies a prior distribution  $p(z)$  to which an encoder with posterior probability  $q_\phi(z|x)$  should map all  $x$  by maximizing a variational lower bound on a decoder posterior probability  $p_\theta(x|z)$ .  $\phi$  and  $\theta$  are the weights of the encoder and decoder respectively. Given an observation  $x_i$ , the marginal distribution over an unknown distribution is given as,

$$p_\theta(x_i) = \int_z p_\theta(x_i|z)p_\theta(z)dz, \quad (1)$$

The integral over  $z$  is computationally intractable for large number of  $z$  variables. We proceed by introducing the posterior distribution of the encoder in order to solve that challenge.

$$p_\theta(x_i) = \int_z p_\theta(x_i|z)p_\theta(z) \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} dz, \quad (2)$$

$$p_\theta(x_i) = \int_z q_\phi(z|x_i) \frac{p_\theta(x_i|z)p_\theta(z)}{q_\phi(z|x_i)} dz, \quad (3)$$

This can be transformed into an expectation form as,

$$p_\theta(x_i) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \frac{p_\theta(x_i|z)p_\theta(z)}{q_\phi(z|x_i)} \right], \quad (4)$$

---

\*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

We take the natural log of both sides as follows,

$$\ln p_\theta(x_i) = \ln \left( \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \frac{p_\theta(x_i|z)p_\theta(z)}{q_\phi(z|x_i)} \right] \right), \quad (5)$$

From Jensen inequality, we know that  $\mathbb{E}[\ln X] \leq \ln[\mathbb{E}(X)]$ , hence,

$$\ln p_\theta(x_i) \geq \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \frac{p_\theta(x_i|z)p_\theta(z)}{q_\phi(z|x_i)} \right], \quad (6)$$

$$\ln p_\theta(x_i) \geq - \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \frac{q_\phi(z|x_i)}{p_\theta(z)} \right] + \mathbb{E}_{z \sim q_\phi(z|x_i)} [\ln p_\theta(x_i|z)], \quad (7)$$

$$\ln p_\theta(x_i) \geq \mathbb{E}_{z \sim q_\phi(z|x_i)} [\ln p_\theta(x_i|z)] - D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)), \quad (8)$$

The R.H.S of equation 8 is referred to as the **Evidence Lower Bound (ELBO)** since it is upper-bounded by the logarithm of the evidence  $p_\theta(x_i)$ .  $\mathbb{E}[\ln p_\theta(x_i|z)]$  is the reconstruction loss while  $D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z))$  is a tractable regulation term known as Kullback-Leibler Divergence which ensures that the approximate posterior of the encoder  $q_\phi(z|x_i)$  does not deviate far from its prior  $p(z)$ .

the loss or objective function to be maximized is given by,

$$\mathcal{L}_{VAE}(\theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i|z)] - D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) \quad (9)$$

For two multivariate normal distribution the closed-form for the objective loss function of VAE is given by,

$$\mathcal{L}_{VAE}(\theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i|z)] + \frac{1}{2} \left[ \ln \left| \frac{\Sigma_1}{\Sigma_0} \right| + m - \text{tr}(\Sigma_1 \Sigma_0^{-1}) - (\mu_1 - \mu_0)^T \Sigma_0^{-1} (\mu_1 - \mu_0) \right] \quad (10)$$

We train the VAE to find the optimal parameters

$$(\theta^*, \phi^*) = \text{argmax}_{(\theta, \phi)} \mathcal{L}_{VAE}(\theta, \phi), \quad (11)$$

Due to the lack of statistical independence of the observation in the latent space generated by VAE, we will adopt the Inverse-Wishart(IW) prior for latent space disentanglement, on the covariance matrix [4]. The inverse-Wishart distribution  $\mathcal{W}^{-1}(\Psi, \nu)$  is parameterized by a positive-definite scale matrix  $\phi \in \mathbb{R}^{p \times p}$  and a degrees of freedom (DoF)  $\nu \geq p - 1$ . An inverse-Wishart distributed random matrix  $X \in \mathbb{R}^{n \times m}$  has probability density function given by,

$$p(X) = \frac{\Psi^{\frac{\nu}{2}}}{2^{\frac{\nu p}{2}} \Gamma_p(\frac{\nu}{2})} X^{-\frac{(\nu+p+1)}{2}} \exp^{-\frac{1}{2} \text{tr}(\Psi X^{-1})}, \quad (12)$$

where  $\Gamma_p(\cdot)$  is the multivariate gamma function. For the inference update,

$$p(\sigma) \sim \mathcal{W}_p^{-1}(\Sigma|\Psi, \nu), \quad (13)$$

$$p(z|\sigma) \sim \mathcal{N}(z|0, \Sigma), \quad (14)$$

Where  $\Sigma$  is the covariance matrix of the matrix  $X$ . The

## 2.1 Transformer-Variational Autoencoder (Trans-VAE)

The objective of using the transformer-based variational autoencoder is to reduce the dimension of the original space (noise or non-keyword removal) of the data whilst preserving context.

We elaborate on the derivations of the VAE loss functions using Gaussian and Wishart prior and posteriors.

## 2.2 Loss function: Gaussian distribution

We recall the objective function for the VAE as,

$$\mathcal{L}_{VAE}(\theta, \phi) = - \underbrace{D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z))}_{\textcircled{1} \text{regularizer}} + \underbrace{E_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i|z)]}_{\textcircled{2} \text{reconstruction}} \quad (15)$$

We assume that the latent prior and the approximate posterior both follow a Gaussian distribution  $\mathcal{N}(\mu_0, \Sigma_0)$  and  $\mathcal{N}(\mu_1, \Sigma_1)$  respectively. We can compute the analytic closed form solutions for the two terms in the loss.

The Kullback-Liebler divergence  $D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z))$  is given as,

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \frac{q_\phi(z|x_i)}{p_\theta(z)} \right] \quad (16)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\ln q_\phi(z|x_i) - \ln p_\theta(z)] \quad (17)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\ln \mathcal{N}(\mu_1, \Sigma_1) - \ln \mathcal{N}(\mu_0, \Sigma_0)] \quad (18)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma_1|^{\frac{1}{2}}} e^{-\frac{1}{2} \|x - \mu_1\|_{\Sigma_1^{-1}}^2} \right] \quad (19)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma_0|^{\frac{1}{2}}} e^{-\frac{1}{2} \|x - \mu_0\|_{\Sigma_0^{-1}}^2} \right]$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_1| - \frac{1}{2} \|x - \mu_1\|_{\Sigma_1^{-1}}^2 \right]$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_0| - \frac{1}{2} \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (20)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \frac{1}{2} \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - \frac{1}{2} \|x - \mu_1\|_{\Sigma_1^{-1}}^2 + \frac{1}{2} \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (21)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (22)$$

Introducing the trace operator, we have that,

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - \text{tr}[(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)] + \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (23)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - \text{tr}[(x - \mu_1)(x - \mu_1)^T \Sigma_1^{-1}] + \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (24)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - \text{tr}[\Sigma_1 \Sigma_1^{-1}] + \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (25)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - \text{tr}[I] + \|x - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (26)$$

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - m \right] + \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] \quad (27)$$

We expand the second term on the right hand side (R.H.S) as follows,

$$\frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[(x_i - \mu_0)^T \Sigma_0^{-1} (x_i - \mu_0)] \right] \quad (28)$$

$$\frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[(x_i - \mu_1 + \mu_1 - \mu_0)^T \Sigma_0^{-1} (x_i - \mu_1 + \mu_1 - \mu_0)] \right] \quad (29)$$

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) &= \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[(x_i - \mu_1)^T \Sigma_0^{-1} (x_i - \mu_1)] \right] \\ \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) &+ \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[(x_i - \mu_1)^T \Sigma_0^{-1} (\mu_0 - \mu_1)] \right] \\ \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) &+ \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[(\mu_1 - \mu_0)^T \Sigma_0^{-1} (\mu_1 - \mu_0)] \right] \end{aligned} \quad (30)$$

Since  $\mathbb{E}_{z \sim q_\phi(z|x_i)} [(x_i - \mu_1)]$  is zero, we have that

$$\begin{aligned} \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] \frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] &= \frac{1}{2} \text{tr} \left[ \mathbb{E}_{z \sim q_\phi(z|x_i)} [(x_i - \mu_1)(x_i - \mu_1)^T \Sigma_0^{-1}] \right] \\ &+ \frac{1}{2} \text{tr} \left[ \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \underbrace{(\mu_1 - \mu_0)^T \Sigma_0^{-1} (\mu_1 - \mu_0)}_{\text{constant}} \right] \right] \end{aligned} \quad (31)$$

$$\frac{1}{2} \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[ \text{tr}[\|x - \mu_0\|_{\Sigma_0^{-1}}^2] \right] = \frac{1}{2} [\text{tr}(\Sigma_1 \Sigma_0^{-1}) + \|\mu_1 - \mu_0\|_{\Sigma_0^{-1}}^2] \quad (32)$$

By substituting equation 32 into 27, the KL divergence can be rewritten as,

$$D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \left[ \ln \left| \frac{\Sigma_0}{\Sigma_1} \right| - m + \text{tr}(\Sigma_1 \Sigma_0^{-1}) + \|\mu_1 - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (33)$$

The negative KL is given as

$$-D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} \left[ \ln \left| \frac{\Sigma_1}{\Sigma_0} \right| + m - \text{tr}(\Sigma_1 \Sigma_0^{-1}) - \|\mu_1 - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (34)$$

However, if  $p_\theta(z) \sim \mathcal{N}(z; 0, I)$ , then the negative KL reduces to,

$$-D_{KL}(q_\phi(z|x_i) \parallel p_\theta(z)) = \frac{1}{2} [\ln |\Sigma_1| + m - \text{tr}(\Sigma_1) - (\mu_1)^T (\mu_1)] \quad (35)$$

Finally, the objective functions can be written as,

$$\mathcal{L}_{VAE}^G(\theta, \phi) = \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i|z)] + \frac{1}{2} \left[ \ln \left| \frac{\Sigma_1}{\Sigma_0} \right| + m - \text{tr}(\Sigma_1 \Sigma_0^{-1}) - \|\mu_1 - \mu_0\|_{\Sigma_0^{-1}}^2 \right] \quad (36)$$

### 2.3 Loss function: Inverse-Wishart distribution

As previously shown for Gaussian distribution, we derive similarly the KL-divergence  $D_{KL}(\mathcal{W}^{-1}(\nu_0, \Phi_0) \parallel \mathcal{W}^{-1}(\nu_1, \Phi_1))$  for two Inverse-Wishart distributions. We note that for a random variable  $X$  which follows an Inverse-Wishart distribution,  $\mathbb{E}[\ln |X|] = \ln |\Phi| - m \ln 2 - \phi_m(\frac{\nu}{2})$  and  $\mathbb{E}[X] = \nu \Phi^{-1}$ . The KL is given as,

$$D_{KL}(\mathcal{W}_1^{-1}(\nu_1, \Phi_1) \parallel \mathcal{W}_0^{-1}(\nu_0, \Phi_0)) = \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} \left[ \ln \frac{\mathcal{W}^{-1}(\nu_1, \Phi_1)}{\mathcal{W}^{-1}(\nu_0, \Phi_0)} \right] \quad (37)$$

We denote  $D_{KL}(\mathcal{W}_1^{-1}(\nu_1, \Phi_1) \parallel \mathcal{W}_0^{-1}(\nu_0, \Phi_0))$  using the compact form  $D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1})$ , so that,

$$D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) = \underbrace{\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)]}_{\textcircled{1} \text{term}} - \underbrace{\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\mathcal{W}^{-1}(\nu_0, \Phi_0)]}_{\textcircled{2} \text{term}} \quad (38)$$

We expand the first term as follow,

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] &= \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} \left[ \frac{\nu_1}{2} \ln |\Phi_1| \right] - \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} \left[ \frac{\nu_1 + m + 1}{2} \ln |X| \right] \\ &\quad - \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} \left[ \frac{1}{2} \text{tr}(\Phi_1 X^{-1}) \right] - \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} \left[ \frac{\nu_1 m}{2} \ln 2 + \ln \Gamma_m\left(\frac{\nu_1}{2}\right) \right]\end{aligned}\quad (39)$$

Since the  $\mathbb{E}[k] = k$ , where  $k$  is a constant and  $\mathbb{E}(\text{tr}[M]) = \text{tr}(\mathbb{E}[M]) = \mu_m$  where  $\mu_m$  is the average of the entries of matrix  $M$ . Hence,

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] &= \frac{\nu_1}{2} \ln |\Phi_1| - \ln \Gamma_m\left(\frac{\nu_1}{2}\right) - \left( \frac{\nu_1 + m + 1}{2} \right) \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln |X|] \\ &\quad - \text{tr} \left[ \frac{1}{2} \left( \Phi_1 \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [X^{-1}] \right) \right]\end{aligned}\quad (40)$$

The expectation of the log of the determinant of a matrix is expanded as follow,

$$\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] = \frac{\nu_1}{2} \ln |\Phi_1| - \ln \Gamma_m\left(\frac{\nu_1}{2}\right) - \frac{\nu_1 + m + 1}{2} \left[ \ln \frac{|\Phi_1|}{2^m} - \psi_m\left(\frac{\nu_1}{2}\right) \right] \quad (41)$$

$$- \text{tr} \left[ \frac{1}{2} \left( \Phi_1 \mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [X^{-1}] \right) \right] \quad (42)$$

For a positive definite matrix  $X$ , its Inverse-Wishart expectation is given by  $\frac{\Phi}{\nu - m - 1}$  for  $\nu > m + 1$ , hence,

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] &= \frac{\nu_1}{2} \ln |\Phi_1| - \ln \Gamma_m\left(\frac{\nu_1}{2}\right) - \frac{\nu_1 + m + 1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right] \\ &\quad - \text{tr} \left[ \frac{1}{2} \left( \Phi_1 \frac{\Phi_1^{-1}}{(\nu_1 - m - 1)^{-1}} \right) \right]\end{aligned}\quad (43)$$

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] &= \frac{\nu_1}{2} \ln |\Phi_1| - \ln \Gamma_m\left(\frac{\nu_1}{2}\right) - \frac{\nu_1 + m + 1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right] \\ &\quad - \text{tr} \left[ \frac{1}{2} \left( \frac{I_m}{(\nu_1 - m - 1)^{-1}} \right) \right]\end{aligned}\quad (44)$$

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_1, \Phi_1)] &= \frac{\nu_1}{2} \ln |\Phi_1| - \frac{\nu_1 p}{2} - \ln \Gamma_m\left(\frac{\nu_1}{2}\right) - \frac{\nu_1 + m + 1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right] \\ &\quad - \frac{(\nu_1 - m - 1)}{2} m\end{aligned}\quad (45)$$

Similarly, we expand the second term to obtain,

$$\begin{aligned}\mathbb{E}_{z \sim \mathcal{W}_1^{-1}} [\ln \mathcal{W}^{-1}(\nu_0, \Phi_0)] &= \frac{\nu_0}{2} \ln |\Phi_0| - \ln \Gamma_m\left(\frac{\nu_0}{2}\right) - \frac{\nu_1 + m + 1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right] \\ &\quad - \frac{1}{2} \text{tr} \left( \frac{\Phi_0 \Phi_1^{-1}}{(\nu_0 - m - 1)^{-1}} \right)\end{aligned}\quad (46)$$

If  $\nu_1 = \nu_0$  and the degree of freedom is very large indicating an informative prior i.e  $\nu \gg m$ , then we can combine equations (45) and (46) to obtain the following,

$$\begin{aligned}D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) &= \frac{\nu_0}{2} \ln |\Phi_1 \Phi_0^{-1}| + \frac{(\nu_0 - m - 1)}{2} \text{tr}(\Phi_0 \Phi_1^{-1}) - \frac{(\nu_1 - m - 1)}{2} m + \ln \frac{\Gamma_m(\frac{\nu_0}{2})}{\Gamma_m(\frac{\nu_1}{2})} \\ &\quad - \frac{\nu_1 + m + 1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right]\end{aligned}\quad (47)$$

$$\begin{aligned}D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) &= \frac{\nu_0}{2} \ln |\Phi_1 \Phi_0^{-1}| + \frac{\nu_1}{2} (\text{tr}(\Phi_0 \Phi_1^{-1}) - m) + \underbrace{\frac{\Gamma_m(\frac{\nu_0}{2})}{\Gamma_m(\frac{\nu_1}{2})} - \frac{\nu_1}{2} \left[ \ln |\Phi_1| - \psi_m\left(\frac{\nu_1}{2}\right) \right]}_{\text{constant}}\end{aligned}\quad (48)$$

and since the last two terms are constant, we can easily reduce the KL to

$$D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) = \frac{\nu_0}{2} \ln |\Phi_1 \Phi_0^{-1}| + \frac{\nu_1}{2} (tr(\Phi_0 \Phi_1^{-1}) - m) \quad (49)$$

We can further expand the R.H.S of equation (49) as follow,

$$D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) = \frac{\nu_0}{2} \ln |\Phi_1| + \frac{\nu_1}{2} tr(\Phi_0 \Phi_1^{-1}) - \underbrace{\frac{\nu_1}{2} \ln |\Phi_0^{-1}| - \frac{\nu_1}{2} m}_{constant} \quad (50)$$

$$D_{KL}(\mathcal{W}_1^{-1} \parallel \mathcal{W}_0^{-1}) = \frac{\nu_0}{2} \ln |\Phi_1| + \frac{\nu_1}{2} tr(\Phi_0 \Phi_1^{-1}) \quad (51)$$

Finally, the loss function for the Inverse-Wishart VAE can be given as,

$$\mathcal{L}_{VAE}^{\mathcal{W}^{-1}} = E_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i|z)] - \frac{\nu_0}{2} \ln |\Phi_1| - \frac{\nu_1}{2} tr(\Phi_0 \Phi_1^{-1}) \quad (52)$$

## References

- [1] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- [2] Raphael Schumann. Unsupervised abstractive sentence summarization using length controlled variational autoencoder, 2018.
- [3] Dinghan Shen, Asli Celikyilmaz, Yizhe Zhang, Liquan Chen, Xin Wang, Jianfeng Gao, and Lawrence Carin. Towards generating long and coherent text with multi-level latent variable models, 2019.
- [4] Abdul Fatir Ansari and Harold Soh. Hyperprior induced unsupervised disentanglement of latent representations, 2019.