

# Enhancing Variational Autoencoders with Mutual Information Neural Estimation for Text Generation

## 1 Appendix

We define the *variational joint distribution* over  $\mathbf{x}$  and  $\mathbf{z}$  as:

$$q_\phi(\mathbf{x}, \mathbf{z}) = \frac{e^{f_\psi(\mathbf{x}, \mathbf{z})}}{Z_\psi} q(\mathbf{x}) q_\phi(\mathbf{z}) \quad (1)$$

where  $f_\psi(\mathbf{x}, \mathbf{z})$  is an energy function with parameters  $\psi$  and the partition function  $Z_\psi$  is defined as  $\mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[e^{f_\psi(\mathbf{x}, \mathbf{z})}]$ . In particular,  $f_\psi(\mathbf{x}, \mathbf{z})$  is approximated by neural networks.

The lower bound on the mutual information  $\mathbb{I}_q[\mathbf{x}, \mathbf{z}]$  is given by:

$$\begin{aligned} \mathbb{I}_q[\mathbf{x}, \mathbf{z}] &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})} \left[ \log \frac{q_\phi(\mathbf{x}, \mathbf{z})}{q(\mathbf{x})q_\phi(\mathbf{z})} \right] \\ &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\log Z_\psi] \\ &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi(\mathbf{x}, \mathbf{z})] - \log Z_\psi \\ &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi(\mathbf{x}, \mathbf{z})] - \log[\mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[e^{f_\psi(\mathbf{x}, \mathbf{z})}]] \\ &\geq \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi(\mathbf{x}, \mathbf{z})] - \xi \cdot \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[e^{f_\psi(\mathbf{x}, \mathbf{z})}] + \log(\xi) + 1 \end{aligned} \quad (2)$$

**Proposition 1.** With the fixed  $\phi$  and  $\xi$ , the optimal energy function  $f_\psi^*(\mathbf{x}, \mathbf{z})$  according to the objective in Eq.(2) is given by

$$\begin{aligned} f_\psi^*(\mathbf{x}, \mathbf{z}) &= \operatorname{argmax}_{f_\psi(\mathbf{x}, \mathbf{z})} \left\{ \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi(\mathbf{x}, \mathbf{z})] - \xi \cdot \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[e^{f_\psi(\mathbf{x}, \mathbf{z})}] \right\} \\ \nabla_\psi f_\psi(\mathbf{x}, \mathbf{z}) &= q_\phi(\mathbf{x}, \mathbf{z}) - \xi \cdot [q(\mathbf{x})q_\phi(\mathbf{z})]e^{f_\psi(\mathbf{x}, \mathbf{z})} = 0 \\ f_\psi^*(\mathbf{x}, \mathbf{z}) &= \log q_\phi(\mathbf{x}, \mathbf{z}) - \log[q(\mathbf{x})q_\phi(\mathbf{z})] - \log(\xi) \end{aligned} \quad (3)$$

When  $\xi = 1$ ,  $f_\psi^*(\mathbf{x}, \mathbf{z})$  becomes essentially pointwise mutual information. This means that the energy function assigns zero probability to the samples independently from  $q(\mathbf{x})q_\phi(\mathbf{z})$ .

With the optimal function  $f_\psi^*(\mathbf{x}, \mathbf{z})$  defined, the max-max objective in Equa-

tion (2) can be reformulated as:

$$\begin{aligned}
C(\phi, \psi^*) &= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[f_\psi^*(\mathbf{x}, \mathbf{z})] - \xi \cdot \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}\left[e^{f_\psi^*(\mathbf{x}, \mathbf{z})}\right] + \log(\xi) + 1 \\
&= \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}\left[\log \frac{q_\phi(\mathbf{x}, \mathbf{z})}{\xi \cdot q(\mathbf{x})q_\phi(\mathbf{z})}\right] - \xi \cdot \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}\left[\frac{q_\phi(\mathbf{x}, \mathbf{z})}{\xi \cdot q(\mathbf{x})q_\phi(\mathbf{z})}\right] + \log(\xi) + 1 \\
&= -\log(\xi) + \mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}\left[\log \frac{q_\phi(\mathbf{x}, \mathbf{z})}{q(\mathbf{x})q_\phi(\mathbf{z})}\right] - \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}\left[\frac{q_\phi(\mathbf{x}, \mathbf{z})}{q(\mathbf{x})q_\phi(\mathbf{z})}\right] + \log(\xi) + 1 \\
&= \text{KL}[q_\phi(\mathbf{x}, \mathbf{z})||q(\mathbf{x})q_\phi(\mathbf{z})] \tag{4}
\end{aligned}$$

Thus, maximizing the lower bound on the MI with respect to  $\phi$  and  $\psi$  is equivalent to maximizing the KL divergence between the joint distribution and two marginal distributions.

We can then compute the gradients of Eq.(4) with respect to  $\phi$  and  $\theta$  for the optimization. While it is easy to compute the gradient with respect to  $\theta$ , the gradient with respect to  $\phi$  is hard to compute since  $C(\phi, \psi^*)$  itself depends on  $\phi$ . Actually, when the function  $f_\psi(\mathbf{x}, \mathbf{z})$  is optimal, the expectation of the gradients becomes zero, that is

$$\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\nabla_\phi f_\psi^*(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[\nabla_\phi e^{f_\psi^*(\mathbf{x}, \mathbf{z})}] = 0 \tag{5}$$

Thus, we can ignore gradients in the optimization when  $f_\psi(\mathbf{x}, \mathbf{z})$  is optimal.

**Proof.**

$$\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\nabla_\phi f_\psi^*(\mathbf{x}, \mathbf{z})] = \underbrace{\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})]}_{\circledast} - \underbrace{\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\nabla_\phi \log q_\phi(\mathbf{z})]}_{\circledcirc} \tag{6}$$

For the  $\circledast$  part,

$$\begin{aligned}
\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\nabla_\phi \log q_\phi(\mathbf{z}|\mathbf{x})] &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\frac{1}{q_\phi(\mathbf{z}|\mathbf{x})}\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})\right] \\
&= \int_{\mathbf{z}} \nabla_\phi q_\phi(\mathbf{z}|\mathbf{x}) \\
&= 0 \tag{7}
\end{aligned}$$

For the  $\circledcirc$  part,

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\nabla_\phi \log q_\phi(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\frac{1}{q_\phi(\mathbf{z})}\nabla_\phi q_\phi(\mathbf{z})\right] \tag{8}$$

$$\begin{aligned}
\mathbb{E}_{q_\phi(\mathbf{z})}[\nabla_\phi e^{f_\psi^*(\mathbf{x}, \mathbf{z})}] &= \mathbb{E}_{q_\phi(\mathbf{z})}\left[\nabla_\phi \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z})}\right] \\
&= \mathbb{E}_{q_\phi(\mathbf{z})}\left[\frac{\nabla_\phi q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z})} - \frac{q_\phi(\mathbf{z}|\mathbf{x})\nabla_\phi q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})^2}\right] \\
&= -\mathbb{E}_{q_\phi(\mathbf{z})}\left[\frac{q_\phi(\mathbf{z}|\mathbf{x})\nabla_\phi q_\phi(\mathbf{z})}{q_\phi(\mathbf{z})^2}\right] \\
&= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\frac{1}{q_\phi(\mathbf{z})}\nabla_\phi q_\phi(\mathbf{z})\right] \tag{9}
\end{aligned}$$

Therefore,

$$\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\nabla_\phi f_\psi^*(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q(\mathbf{x})q_\phi(\mathbf{z})}[\nabla_\phi e^{f_\psi^*(\mathbf{x}, \mathbf{z})}] = 0 \quad (10)$$