

A Proof of Theorem 1

First, we prepare some lemmas.

Lemma 2. Let $h : \mathbb{R}^{d_z \times d} \rightarrow \mathbb{R}^{d_z \times n}$ be a linear map defined by $h(\Phi) := \Phi G$. When the rank of G is n , there exists an orthogonal linear map $\tau : \mathbb{R}^{d_z \times d} \rightarrow \mathbb{R}^{d_z \times d}$ such that $\tau(\Phi) = [\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)}]$ satisfies $\ker \tilde{h} = \text{span} \left[\mathbf{0}_{d_z \times n}, \tilde{\Phi}^{(2)} \right]$, where $\tilde{h} := h \circ \tau^{-1}$, $\tilde{\Phi}^{(1)} := [\tilde{\phi}_1, \dots, \tilde{\phi}_n]$, $\tilde{\Phi}^{(2)} := [\tilde{\phi}_{n+1}, \dots, \tilde{\phi}_d]$ and $\tilde{\phi}_i \in \mathbb{R}^{d_z}$ for $i = 1, \dots, d$.

Proof. The singular value decomposition of G is represented by $TGT^\top = \begin{bmatrix} \Lambda \\ \mathbf{0}_{(d-n \times n)} \end{bmatrix}$, where Λ is a $n \times n$ diagonal matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, and T and T^\top are orthogonal matrices. Since the rank of G is n , $\lambda_1, \dots, \lambda_n$ are non-zero. When we set $\tau(\Phi) := \Phi T^\top$, we obtain

$$\begin{aligned} \tilde{h}(\tilde{\Phi}) &:= h(\tau^{-1}(\tilde{\Phi})) \\ &= \tilde{\Phi} T G = \tilde{\Phi} (T G T^\top) T \\ &= \tilde{\Phi} \begin{bmatrix} \Lambda \\ \mathbf{0}_{(d-n \times n)} \end{bmatrix} T. \end{aligned} \quad (21)$$

From the above equation, $\ker \tilde{h} = \text{span} \left[\mathbf{0}_{d' \times n}, \tilde{\Phi}^{(2)} \right]$ holds. \square

Lemma 3. For $\tilde{\Phi} \in \mathbb{R}^{d_z \times d}$, let Φ satisfy:

$$\tilde{\Phi} = [\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)}] = \tau(\Phi).$$

A map $\tilde{h}^{(1)} : \mathbb{R}^{d_z \times n} \rightarrow \mathbb{R}^{d_z \times n}$ defined by $\tilde{h}^{(1)}(\tilde{\Phi}^{(1)}) := \tilde{h}([\tilde{\Phi}^{(1)}, \mathbf{0}])$ satisfies $\Phi G = \tilde{h}^{(1)}(\tilde{\Phi}^{(1)})$ and is linear isomorphic.

Proof. From the definition, we have

$$\begin{aligned} \Phi G &= \tau^{-1}(\tilde{\Phi}) G = \tilde{\Phi} T G \\ &= \tilde{\Phi} (T G T^\top) T \\ &= [\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)}] \begin{bmatrix} \Lambda \\ \mathbf{0}_{(d-n \times n)} \end{bmatrix} T \\ &= [\tilde{\Phi}^{(1)}, \mathbf{0}] \begin{bmatrix} \Lambda \\ \mathbf{0}_{(d-n \times n)} \end{bmatrix} T \\ &= \tau^{-1}([\tilde{\Phi}^{(1)}, \mathbf{0}]) G \\ &= h \circ \tau^{-1}([\tilde{\Phi}^{(1)}, \mathbf{0}]) = \tilde{h}([\tilde{\Phi}^{(1)}, \mathbf{0}]) \\ &= \tilde{h}^{(1)}(\tilde{\Phi}^{(1)}). \end{aligned}$$

By the definition, $\tilde{h}^{(1)}$ is linear. Here, $\tilde{h}^{(1)}$ is injective, since $\ker \tilde{h} = \text{span} \left[\mathbf{0}_{d' \times n}, \tilde{\Phi}^{(2)} \right]$, and hence, $\dim(\text{Im } \tilde{h}^{(1)}) \geq d_z \times n$. Since $\text{Im } \tilde{h}^{(1)} \subset \mathbb{R}^{d_z \times n}$, $\tilde{h}^{(1)}$ is surjective. \square

Lemma 4. For $V : \mathbb{R}^D \ni \phi \mapsto V(\Phi) = U(\mathbf{X}, f_{\mathbf{z}|\mathbf{x}}(\mathbf{X}; \Phi)) := \sum_{i=1}^n U(\mathbf{x}^{(i)}, f_{\mathbf{z}|\mathbf{x}}(\mathbf{x}^{(i)}; \Phi)) \in \mathbb{R}$, $\tilde{\Phi} = [\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)}] := \tau(\Phi)$, $\tilde{V} := V \circ \tau^{-1}$ and $\tilde{V}^{(1)}(\tilde{\Phi}^{(1)}) := \tilde{V}([\tilde{\Phi}^{(1)}, \mathbf{0}_{d_z \times (d-n)}])$, Eq. (12) is equivalent to

$$d\tilde{\Phi}^{(1)} = -\nabla_{\tilde{\Phi}^{(1)}} \tilde{V}^{(1)}(\tilde{\Phi}^{(1)}) dt + \sqrt{2} dB, \quad (22)$$

$$d\tilde{\Phi}^{(2)} = \sqrt{2} dB. \quad (23)$$

Proof. By direct calculation, we obtain

$$\begin{aligned}
& \tilde{V} \left(\left[\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)} \right] \right) & (24) \\
& = V \circ \tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)} \right] \right) \\
& = U \left(\mathbf{X}, f_{\mathbf{z}|\mathbf{x}} \left(\mathbf{X}; \tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)} \right] \right) \right) \right) \\
& = U \left(\mathbf{X}, h \left(\tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)} \right] \right) \right) \right) \\
& = U \left(\mathbf{X}, h \circ \tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \mathbf{0} \right] \right) + h \circ \tau^{-1} \left(\left[\mathbf{0}, \tilde{\Phi}^{(2)} \right] \right) \right) \\
& = U \left(\mathbf{X}, h \left(\tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \mathbf{0} \right] \right) \right) \right) \\
& = U \left(\mathbf{X}, f_{\mathbf{z}|\mathbf{x}} \left(\mathbf{X}; \tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \mathbf{0} \right] \right) \right) \right) \\
& = V \circ \tau^{-1} \left(\left[\tilde{\Phi}^{(1)}, \mathbf{0} \right] \right) \\
& = \tilde{V} \left(\left[\tilde{\Phi}^{(1)}, \mathbf{0} \right] \right). & (25)
\end{aligned}$$

Then, the following equivalence holds:

$$\begin{aligned}
& d\Phi = -\nabla_{\Phi} V(\Phi) dt + \sqrt{2} dB, \\
& \Leftrightarrow d\tau^{-1}(\Phi) = -\tau^{\top} \left(\nabla_{\tilde{\Phi}} \tilde{V}(\tilde{\Phi}) \right) dt + \sqrt{2} dB \\
& \Leftrightarrow d\tilde{\Phi} = -\tau \circ \tau^{\top} \left(\nabla_{\tilde{\Phi}} \tilde{V}(\tilde{\Phi}) \right) dt + \sqrt{2} d\tau(B) \\
& = -\nabla_{\tilde{\Phi}} \tilde{V}(\tilde{\Phi}) dt + \sqrt{2} dB, & (26)
\end{aligned}$$

where we used $\tau \circ \tau^{\top} = \text{id}$ because τ is orthogonal. From Eq. (25), the dynamics in Eq. (26) is equivalent to Eq. (22) and Eq. (23). \square

In the following, we prove Theorem 1 using the above lemmas.

Because the latent variables $\mathbf{Z} := \tilde{h}^{(1)} \left(\tilde{\Phi}^{(1)} \right)$ are independent of $\tilde{\Phi}^{(2)}$, the stationary distribution $q(\mathbf{Z} | \mathbf{X})$ is given by $\left(\tilde{h}^{(1)} \right)_{\#} \left(p_{*}^{(1)} \right)(\mathbf{Z})$, which is the pushforward measure of the probability distribution $p^{(1)} \left(\tilde{\Phi} \right)$ by $\tilde{h}^{(1)}$. Then, we have

$$\begin{aligned}
& q(\mathbf{Z} | \mathbf{X}) \\
& = \left(\tilde{h}^{(1)} \right)_{\#} \left(p_{*}^{(1)} \right)(\mathbf{Z}) \\
& = p^{(1)} \left(\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}) \right) \left| \det \frac{d\left(\tilde{h}^{(1)} \right)^{-1}}{d\mathbf{Z}} \right| \\
& = p^{(1)} \left(\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}) \right) \left| \det \frac{d\tilde{h}^{(1)}}{d\tilde{\Phi}^{(1)}} \right|^{-1} \\
& = p^{(1)} \left(\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}) \right) \times \left| \det \tilde{h}^{(1)} \right|^{-1} \\
& \propto \exp \left(-\tilde{V} \left(\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}) \right) \right) \\
& = \exp \left(-V \left(\tau^{-1} \left(\left[\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}), \mathbf{0} \right] \right) \right) \right) \\
& = \exp \left(-U \left(\mathbf{X}, \tau^{-1} \left(\left[\left(\tilde{h}^{(1)} \right)^{-1}(\mathbf{Z}), \mathbf{0} \right] \mathbf{G} \right) \right) \right) \\
& = \exp(-U(\mathbf{X}, \mathbf{Z})),
\end{aligned}$$

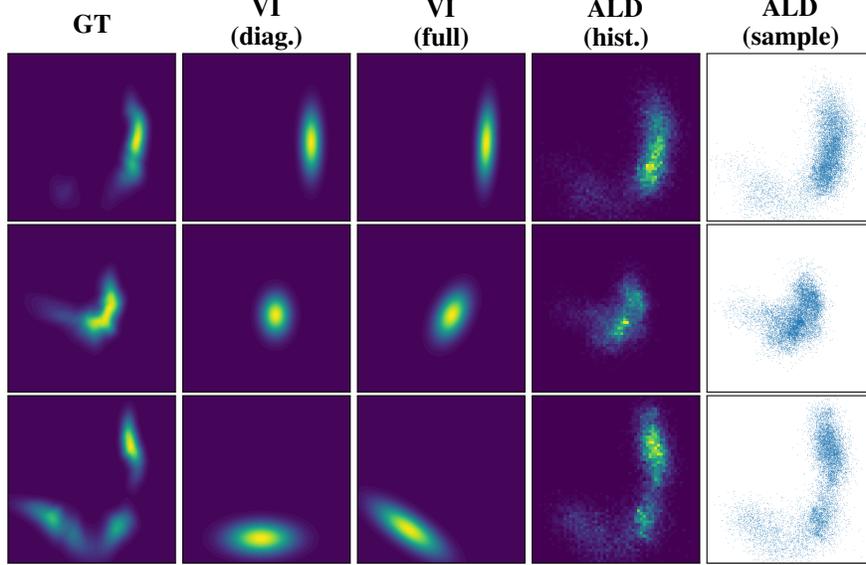


Figure 5: Neural likelihood experiments.

where we used that $\frac{d\tilde{h}^{(1)}}{d\tilde{\Phi}^{(1)}} = \tilde{h}^{(1)}$ because of the linearity of $\tilde{h}^{(1)}$ and is constant with respect to \mathbf{Z} . The last equation is derived as follows. From Lemma 3, $\Phi \mathbf{G} = \tilde{h}^{(1)}(\tilde{\Phi}^{(1)})$ holds when $\tilde{\Phi} = [\tilde{\Phi}^{(1)}, \tilde{\Phi}^{(2)}] = \tau(\Phi)$. Thus, when $\Phi = \tau^{-1}([\tilde{\Phi}^{(1)}, \mathbf{0}])$, we obtain $\tilde{h}^{(1)}(\tilde{\Phi}^{(1)}) = \Phi \mathbf{G} = \tau^{-1}([\tilde{\Phi}^{(1)}, \mathbf{0}]) \mathbf{G}$. In particular, for $\tilde{\Phi}^{(1)} = (\tilde{h}^{(1)})^{-1}(\mathbf{Z})$, we have

$$\begin{aligned} \mathbf{Z} &= \tilde{h}^{(1)}\left(\left(\tilde{h}^{(1)}\right)^{-1}(\mathbf{Z})\right) \\ &= \tau^{-1}\left(\left[\left(\tilde{h}^{(1)}\right)^{-1}(\mathbf{Z}), \mathbf{0}\right]\right) \mathbf{G}. \end{aligned}$$

□

B Experimental Settings

B.1 Neural likelihood example

We perform an experiment with a complex posterior, wherein the likelihood is defined with a randomly initialized neural network f_θ . Particularly, we parameterize f_θ by four fully-connected layers of 128 units with ReLU activation and two dimensional outputs like $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(f_\theta(\mathbf{z}), \sigma_x^2 \mathbf{I})$. We initialize the weight and bias parameters with $\mathcal{N}(0, 0.2\mathbf{I})$ and $\mathcal{N}(0, 0.1\mathbf{I})$, respectively. In addition, we set the observation variance σ_x to 0.25. We used the same neural network architecture for the encoder f_ϕ . Other settings are same as the conjugate Gaussian experiment described in Section 5.1.

The results are shown in Figure 5. The left three columns show the density visualizations of the ground truth or approximation posteriors of VI methods; the right two columns show the visualizations of 2D histograms and samples obtained using ALD. For VI method, we use two different models. One uses diagonal Gaussians, i.e., $\mathcal{N}(\mu(\mathbf{x}; \phi), \text{diag}(\sigma^2(\mathbf{x}; \phi)))$, for the variational distribution, and the other uses Gaussians with full covariance $\mathcal{N}(\mu(\mathbf{x}; \phi), \Sigma(\mathbf{x}; \phi))$. From the density visualization of GT, the true posterior is multimodal and skewed; this leads to the failure of the Gaussian VI methods notwithstanding considering covariance. In contrast, the samples of ALD accurately capture such a complex distribution, because ALD does not need to assume any tractable distributions for approximating the true posteriors. The samples of ALD capture well the multimodal and skewed posterior, while Gaussian VI methods fail it even when considering covariance.

B.2 Image Generation

For the image generation experiments, we use a standard Gaussian distribution $\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ for the latent prior. The latent dimensionality is set to 8 for MNIST, 16 for SVHN, and 32 for CIFAR-10 and CelebA. The raw images, which take the values in $\{0, 1, \dots, 255\}$, are scaled into the range from -1 to 1 via preprocessing. Because the values of the preprocessed images are not continuous in a precise sense due to the quantization, it is not desirable to use continuous distributions, e.g., Gaussians, for the likelihood function $p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})$. Therefore, we define the likelihood using a discretized logistic distribution [Salimans et al., 2017] as follows:

$$\begin{aligned}
 p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) &= \prod_i^{d_{\mathbf{x}}} \int_{a_i}^{b_i} \text{Logistic}(\tilde{x}_i; \mu_i, s) d\tilde{x}_i, \\
 &= \prod_i^{d_{\mathbf{x}}} \left(\sigma\left(\frac{b_i - \mu_i}{s}\right) - \sigma\left(\frac{a_i - \mu_i}{s}\right) \right), \\
 a_i &= \begin{cases} -\infty & x = -1 \\ x_i - \frac{1}{255} & \text{otherwise} \end{cases}, \\
 b_i &= \begin{cases} \infty & x = 1 \\ x_i + \frac{1}{255} & \text{otherwise} \end{cases},
 \end{aligned} \tag{27}$$

where $\boldsymbol{\mu} := f_{\mathbf{x}|\mathbf{z}}(\mathbf{z}; \boldsymbol{\theta})$, $f_{\mathbf{x}|\mathbf{z}} : \mathbb{R}^{d_{\mathbf{z}}} \rightarrow \mathbb{R}^{d_{\mathbf{x}}}$. $\text{Logistic}(\cdot; \mu, s)$ is the density function of a logistic distribution with the location parameter μ and the scale parameter s , and σ is the logistic sigmoid function. We use a neural network with four fully-connected layers for the decoder function $f_{\mathbf{x}|\mathbf{z}}$. The number of hidden units are set to 1,024, and ReLU is used for the activation function. Before each activation, we apply the layer normalization [Ba et al., 2016] to stabilize training. The scale parameter s is also optimized in the training. Because it has a constraint of $s > 0$, we parameterize $s = \zeta(b)^{-1/2}$, where ζ is the softplus function, and treat b as a learnable parameter instead. When the model has sufficiently high expressive power, b may diverge to infinity [Rezende and Viola, 2018], so we add a regularization term of $(b + 2\zeta(-b))/m$ to the loss function, where m is the number of training examples. This regularization corresponds to assuming a standard logistic distribution $\text{Logistic}(b; 0, 1)$ for the prior distribution of b . We optimize the models using stochastic gradient ascent with the learning rate of 1×10^{-4} and the batch size of 100. We run two steps of ALD iterations, i.e., $T = 2$ in Algorithm 2, and the step size η is set to 1×10^{-4} . We use the same experimental settings for the baseline models. We run the training iterations for 50 epochs for MNIST, SVHN, CIFAR-10 and 20 epochs for CelebA. The implementation is available at <https://github.com/iShohei220/LAE>.

B.3 Datasets

All the dataset we use in the experiment is public for non-commercial research purposes. MNIST [LeCun et al., 1998], SVHN [Netzer et al., 2011], CIFAR-10 [Krizhevsky et al., 2009], CelebA [Liu et al., 2015] are available at <http://yann.lecun.com/exdb/mnist/>, <http://ufldl.stanford.edu/housenumber>, <https://www.cs.toronto.edu/~kriz/cifar.html>, <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>, and https://github.com/tkarras/progressive_growing_of_gans, respectively. The images of CelebA are resized to 32×32 in the experiment. We use the default settings of data splits for all datasets.

B.4 Computational Resources

We run all the experiments on AI Bridging Cloud Infrastructure (ABCI), which is a large-scale computing infrastructure provided by National Institute of Advanced Industrial Science and Technology (AIST). The experiments are performed on Computing Nodes (V) of ABCI, each of which has four NVIDIA V100 GPU accelerators, two Intel Xeon Gold 6148, one NVMe SSD, 384GiB memory, two InfiniBand EDR ports (100Gbps each). Please see <https://abci.ai> for more details.

Table 2: Effects of the number of MCMC iterations of Hoffman [2017]. We report the mean and standard deviation of the negative evidence lower bound per data dimension in three different seeds. Lower is better.

		MNIST	SVHN	CIFAR-10	CelebA
Hoffman [2017]	$(T = 0)$	1.189 ± 0.002	4.442 ± 0.003	4.820 ± 0.005	4.671 ± 0.001
Hoffman [2017]	$(T = 2)$	1.189 ± 0.002	4.440 ± 0.007	4.831 ± 0.005	4.662 ± 0.011
Hoffman [2017]	$(T = 10)$	1.188 ± 0.001	4.437 ± 0.009	4.832 ± 0.006	4.664 ± 0.004
LAE	$(T = 2)$	1.177 ± 0.001	4.412 ± 0.002	4.773 ± 0.003	4.636 ± 0.003

C Additional Experiments

In the main result in Section 5, we fix the number of MCMC iterations (i.e., T) for the model of Hoffman [2017]. In this additional experiment, we further investigate the effect of T by changing it from 0 to 10. Note that when $T = 0$, the model is identical to the normal VAE. The result is shown in Table 2. It can be seen that the effect is relatively small, and our LAE (with $T = 2$) shows better performance than all cases.