Variational Autoencoder with Implicit Optimal Priors

Hiroshi Takahashi¹, Tomoharu Iwata², Yuki Yamanaka³, Masanori Yamada³, Satoshi Yagi¹

NTT Software Innovation Center
 NTT Communication Science Laboratories
 NTT Secure Platform Laboratories

{takahashi.hiroshi, iwata.tomoharu, yamanaka.yuki, yamada.m, yagi.satoshi}@lab.ntt.co.jp

Abstract

The variational autoencoder (VAE) is a powerful generative model that can estimate the probability of a data point by using latent variables. In the VAE, the posterior of the latent variable given the data point is regularized by the prior of the latent variable using Kullback Leibler (KL) divergence. Although the standard Gaussian distribution is usually used for the prior, this simple prior incurs over-regularization. As a sophisticated prior, the aggregated posterior has been introduced, which is the expectation of the posterior over the data distribution. This prior is optimal for the VAE in terms of maximizing the training objective function. However, KL divergence with the aggregated posterior cannot be calculated in a closed form, which prevents us from using this optimal prior. With the proposed method, we introduce the density ratio trick to estimate this KL divergence without modeling the aggregated posterior explicitly. Since the density ratio trick does not work well in high dimensions, we rewrite this KL divergence that contains the high-dimensional density ratio into the sum of the analytically calculable term and the lowdimensional density ratio term, to which the density ratio trick is applied. Experiments on various datasets show that the VAE with this implicit optimal prior achieves high density estimation performance.

1 Introduction

Estimating data distributions is one of the important challenges of machine learning. The variational autoencoder (VAE) (Kingma and Welling 2013; Rezende, Mohamed, and Wierstra 2014) was presented as a powerful generative model that can learn distributions by using latent variables and neural networks. Since the VAE can capture the high-dimensional complicated data distributions, it is widely applied to various data, such as images (Gulrajani et al. 2016), videos (Gregor et al. 2015), and audio and speech (Hsu, Zhang, and Glass 2017; van den Oord, Vinyals, and kavukcuoglu 2017).

The VAE is composed of three distributions: the encoder, the decoder, and the prior of the latent variable. The encoder and the decoder are conditional distributions, and neural networks are used to model these distributions. The encoder defines the posterior of the latent variable given the data point, whereas the decoder defines the distribution of the data point given the latent variable. The parameters of encoder and decoder neural networks are optimized by maximizing the sum

of the evidence lower bound of the log marginal likelihood. In the training of VAE, the prior regularizes the encoder by Kullback Leibler (KL) divergence. The standard Gaussian distribution is usually used for the prior since the KL divergence can be calculated in a closed form.

Recent research shows that the prior plays an important role in the density estimation (Hoffman and Johnson 2016). Although the standard Gaussian prior is usually used, this simple prior incurs over-regularization, which is one of the causes of the poor density estimation performance. This over-regularization is also known as the posterior-collapse phenomenon (van den Oord, Vinyals, and kavukcuoglu 2017). To improve the density estimation performance, the aggregated posterior prior has been introduced, which is the expectation of the encoder over the data distribution (Hoffman and Johnson 2016). The aggregated posterior is an optimal prior in terms of maximizing the training objective function of the VAE. However, KL divergence with the aggregated posterior cannot be calculated in a closed form, which prevents us from using this optimal prior. In previous work (Tomczak and Welling 2018), the aggregated posterior is modeled by using the finite mixture of encoders for calculating the KL divergence in a closed form. Nevertheless, it has sensitive hyperparameters such as the number of mixture components, which are difficult to tune.

In this paper, we propose the VAE with implicit optimal priors, where the aggregated posterior is used as the prior, but the KL divergence is directly estimated without modeling the aggregated posterior explicitly. This implicit modeling enables us to avoid the difficult hyperparameter tuning for the aggregated posterior model. We use the density ratio trick, which can estimate the density ratio between two distributions without modeling each distribution explicitly, since the KL divergence is the expectation of the density ratio between the encoder and aggregated posterior. Although the density ratio trick is powerful, it has been experimentally shown to work poorly in high dimensions (Sugiyama, Suzuki, and Kanamori 2012; Rosca, Lakshminarayanan, and Mohamed 2018). Unfortunately, with high-dimensional datasets, the density ratio between the encoder and the aggregated posterior also becomes high-dimensional. To avoid the density ratio estimation in high dimensions, we rewrite the KL divergence with the aggregated posterior to the sum of two terms. The first term is the KL divergence between the encoder and the standard Gaussian prior, which can be calculated in a closed form. The other term is the low-dimensional density ratio between the aggregated posterior and the standard Gaussian distribution, to which the density ratio trick is applied.

2 Preliminaries

2.1 Variational Autoencoder

First, we review the variational autoencoder (VAE) (Kingma and Welling 2013; Rezende, Mohamed, and Wierstra 2014). The VAE is a probabilistic latent variable model that relates an observed variable vector ${\bf z}$ to a low-dimensional latent variable vector ${\bf z}$ by a conditional distribution. The VAE models the probability of a data point ${\bf x}$ by

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z}) d\mathbf{z}, \tag{1}$$

where $p_{\lambda}(\mathbf{z})$ is a prior of the latent variable vector, and $p_{\theta}(\mathbf{x}|\mathbf{z})$ is the conditional distribution of \mathbf{x} given \mathbf{z} , which is modeled by neural networks with parameter θ . For example, if \mathbf{x} is binary, this distribution is modeled by a Bernoulli distribution $\mathcal{B}(\mathbf{x} \mid \mu_{\theta}(\mathbf{z}))$, where $\mu_{\theta}(\mathbf{z})$ is neural networks with parameter θ and input \mathbf{z} . These neural networks are called the decoder.

The log marginal likelihood $\ln p_{\theta}(\mathbf{x})$ is bounded below by the evidence lower bound (ELBO), which is derived from Jensen's inequality, as follows:

$$\ln p_{\theta}(\mathbf{x}) = \ln \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \right]$$

$$\geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln \frac{p_{\theta}(\mathbf{x} \mid \mathbf{z}) p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \right]$$

$$\equiv \mathcal{L}(\mathbf{x}; \theta, \phi), \tag{2}$$

where $\mathbb{E}[\cdot]$ represents the expectation, and $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ is the posterior of \mathbf{z} given \mathbf{x} , which are modeled by neural networks with parameter ϕ . $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ is usually modeled by a Gaussian distribution $\mathcal{N}(\mathbf{z} \mid \mu_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x}))$, where $\mu_{\phi}(\mathbf{x})$ and $\sigma_{\phi}^2(\mathbf{x})$ are neural networks with parameter ϕ and input \mathbf{x} . These neural networks are called the encoder.

The ELBO (Eq. (2)) can be also written as

$$\mathcal{L}(\mathbf{x}; \theta, \phi) = -D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || p_{\lambda}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\ln p_{\theta}(\mathbf{x} \mid \mathbf{z}) \right], \quad (3)$$

where $D_{KL}(P||Q)$ is the Kullback Leibler (KL) divergence between P and Q. The second expectation term in Eq. (3) is called the reconstruction term, which is also known as the negative reconstruction error.

The parameters of the encoder and decoder neural networks are optimized by maximizing the following expectation of the lower bound of the log marginal likelihood:

$$\max_{\theta,\phi} \int p_{\mathcal{D}}(\mathbf{x}) \mathcal{L}(\mathbf{x};\theta,\phi) d\mathbf{x}, \tag{4}$$

where $p_{\mathcal{D}}(\mathbf{x})$ is the data distribution.

2.2 Aggregated Posterior

The training of VAE is maximizing the reconstruction term with regularization by KL divergence between the encoder and the prior. The prior is usually modeled by a standard Gaussian distribution $\mathcal{N}(\mathbf{z}|\mathbf{0},\mathbf{I})$ (Kingma and Welling 2013). However, this is not an optimal prior for the VAE. This simple prior incurs over-regularization, which is one of the causes of the poor density estimation performance (Hoffman and Johnson 2016). This phenomenon is called the posterior-collapse (van den Oord, Vinyals, and kavukcuoglu 2017).

The optimal prior that maximizes the objective function of VAE (Eq. (4)) can be derived analytically. The maximization of Eq. (4) with respect to the prior $p_{\lambda}(\mathbf{z})$ is written as follows:

$$\arg \max_{p_{\lambda}(\mathbf{z})} \int p_{\mathcal{D}}(\mathbf{x}) \mathcal{L}(\mathbf{x}; \theta, \phi) d\mathbf{x}$$

$$= \arg \max_{p_{\lambda}(\mathbf{z})} \int p_{\mathcal{D}}(\mathbf{x}) \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\lambda}(\mathbf{z}) \right] d\mathbf{x}$$

$$= \arg \max_{p_{\lambda}(\mathbf{z})} \int \left\{ \int p_{\mathcal{D}}(\mathbf{x}) q_{\phi}(\mathbf{z} \mid \mathbf{x}) d\mathbf{x} \right\} \ln p_{\lambda}(\mathbf{z}) d\mathbf{z}$$

$$= \arg \max_{p_{\lambda}(\mathbf{z})} -H(\int p_{\mathcal{D}}(\mathbf{x}) q_{\phi}(\mathbf{z} \mid \mathbf{x}) d\mathbf{x}, p_{\lambda}(\mathbf{z})), \quad (5)$$

where -H(P,Q) is the negative cross entropy between P and Q. Since -H(P,Q) takes a maximum value when P is equal to Q, the optimal prior $p_{\lambda}^*(\mathbf{z})$ that maximizes Eq. (4) is

$$p_{\lambda}^{*}(\mathbf{z}) = \int p_{\mathcal{D}}(\mathbf{x}) q_{\phi}(\mathbf{z} \mid \mathbf{x}) d\mathbf{x} \equiv q_{\phi}(\mathbf{z}).$$
 (6)

This distribution $q_{\phi}(\mathbf{z})$ is called the aggregated posterior.

When we use the standard Gaussian prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0},\mathbf{I})$, the KL divergence $D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) \| p(\mathbf{z}))$ can be calculated in a closed form (Kingma and Welling 2013). However, when we use the aggregated posterior $q_{\phi}(\mathbf{z})$ as the prior, the KL divergence

$$D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || q_{\phi}(\mathbf{z})) = \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\ln \frac{q_{\phi}(\mathbf{z} \mid \mathbf{x})}{q_{\phi}(\mathbf{z})} \right]$$
(7)

cannot be calculated in a closed form, which prevents us from using the aggregated posterior as the prior.

2.3 Previous work: VampPrior

In previous work, the aggregated posterior is modeled by using the finite mixture of encoders to calculate the KL divergence. Given a dataset $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$, the aggregated posterior can be simply modeled by an empirical distribution:

$$q_{\phi}(\mathbf{z}) \simeq \frac{1}{N} \sum_{i=1}^{N} q_{\phi}(\mathbf{z} \mid \mathbf{x}^{(i)}).$$
 (8)

Nevertheless, this empirical distribution incurs over-fitting (Tomczak and Welling 2018). Thus, the VampPrior (Tomczak and Welling 2018) models the aggregated posterior by

$$q_{\phi}(\mathbf{z}) \simeq \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} \mid \mathbf{u}^{(k)}),$$
 (9)

where K is the number of mixtures, and $\mathbf{u}^{(k)}$ is the same dimensional vector as a data point. \mathbf{u} is regarded as the pseudo input for the encoder, and is optimized during the training of the VAE through the stochastic gradient descent (SGD). If $K \ll N$, the VampPrior can avoid over-fitting (Tomczak and Welling 2018). The KL divergence with the VampPrior can be calculated by the Monte Carlo approximation. The VAE with the VampPrior achieves better density estimation performance than the VAE with the standard Gaussian prior and the VAE with the Gaussian mixture prior (Dilokthanakul et al. 2016). However, this approach has a major drawback: it has sensitive hyperparameters such as the number of mixtures K, which are difficult to tune.

From the above discussion, the aggregated posterior seems to be difficult to model explicitly. In this paper, we estimate the KL divergence with the aggregated posterior without modeling the aggregated posterior explicitly.

3 Proposed Method

In this section, we propose the approximation method of the KL divergence with the aggregated posterior, and describe the optimization procedure of our approach.

3.1 Estimating the KL Divergence

As shown in Eq. (7), the KL divergence with the aggregated posterior is the expectation of the logarithm of the density ratio $q_{\phi}(\mathbf{z} \mid \mathbf{x})/q_{\phi}(\mathbf{z})$. In this paper, we introduce the density ratio trick (Sugiyama, Suzuki, and Kanamori 2012; Goodfellow et al. 2014), which can estimate the ratio of two distributions without modeling each distribution explicitly. Hence, there is no need to model the aggregated posterior explicitly. By using the density ratio trick, $q_{\phi}(\mathbf{z} \mid \mathbf{x})/q_{\phi}(\mathbf{z})$ can be estimated by using a probabilistic binary classifier $D(\mathbf{x}, \mathbf{z})$.

However, the density ratio trick has a serious drawback: it has been experimentally shown to work poorly in high dimensions (Sugiyama, Suzuki, and Kanamori 2012; Rosca, Lakshminarayanan, and Mohamed 2018). Unfortunately, if \mathbf{x} is high-dimensional, $q_{\phi}(\mathbf{z} \mid \mathbf{x})/q_{\phi}(\mathbf{z})$ also becomes a high-dimensional density ratio. The reason is as follows. Since the $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ is a conditional distribution of \mathbf{z} given \mathbf{x} , the density ratio trick has to use a probabilistic binary classifier $D(\mathbf{x}, \mathbf{z})$, which takes \mathbf{x} and \mathbf{z} jointly as an input. In fact, $D(\mathbf{x}, \mathbf{z})$ estimates the density ratio of joint distributions of \mathbf{x} and \mathbf{z} , which is a high-dimensional density ratio with high-dimensional \mathbf{x} (Mescheder, Nowozin, and Geiger 2017).

To avoid the density ratio estimation in high dimensions, we rewrite the KL divergence $D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || q_{\phi}(\mathbf{z}))$ as

follows:

$$D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || q_{\phi}(\mathbf{z}))$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\ln \frac{q_{\phi}(\mathbf{z} \mid \mathbf{x})}{q_{\phi}(\mathbf{z})} \right]$$

$$= \int q_{\phi}(\mathbf{z} \mid \mathbf{x}) \ln \frac{q_{\phi}(\mathbf{z} \mid \mathbf{x})}{p(\mathbf{z})} d\mathbf{z}$$

$$+ \int q_{\phi}(\mathbf{z} \mid \mathbf{x}) \ln \frac{p(\mathbf{z})}{q_{\phi}(\mathbf{z})} d\mathbf{z}$$

$$= D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || p(\mathbf{z})) - \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\ln \frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z})} \right]. \quad (10)$$

The first term in Eq. (10) is KL divergence between the encoder and standard Gaussian distribution, which can be calculated in a closed form. The second term is the expectation of the logarithm of the density ratio $q_{\phi}(\mathbf{z})/p(\mathbf{z})$. We estimate $q_{\phi}(\mathbf{z})/p(\mathbf{z})$ with the density ratio trick. Since the latent variable vector \mathbf{z} is low-dimensional, the density ratio trick works well.

We can estimate the density ratio $q_{\phi}(\mathbf{z})/p(\mathbf{z})$ as follows. First, we prepare the samples from $q_{\phi}(\mathbf{z})$ and samples from $p(\mathbf{z})$. We can sample from $p(\mathbf{z})$ and $q_{\phi}(\mathbf{z} \mid \mathbf{x})$ since these distributions are a Gaussian, and we can also sample from the aggregated posterior $q_{\phi}(\mathbf{z})$ by using ancestral sampling: we choose a data point \mathbf{x} from a dataset randomly and sample \mathbf{z} from the encoder given this data point \mathbf{x} . Second, we label y=1 to samples from $q_{\phi}(\mathbf{z})$ and y=0 to samples from $p(\mathbf{z})$. Then, we define $p^*(\mathbf{z} \mid y)$ as follows:

$$p^*(\mathbf{z} \mid y) \equiv \begin{cases} q_{\phi}(\mathbf{z}) & (y=1) \\ p(\mathbf{z}) & (y=0) \end{cases}$$
 (11)

Third, we introduce a probabilistic binary classifier $D(\mathbf{z})$ that discriminates between the samples from $q_{\phi}(\mathbf{z})$ and samples from $p(\mathbf{z})$. If $D(\mathbf{z})$ can discriminate these samples perfectly, we can rewrite the density ratio $q_{\phi}(\mathbf{z})/p(\mathbf{z})$ by using Bayes theorem and $D(\mathbf{z})$ as follows:

$$\frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z})} = \frac{p^{*}(\mathbf{z} \mid y = 1)}{p^{*}(\mathbf{z} \mid y = 0)} = \frac{p^{*}(y = 0)p^{*}(y = 1 \mid \mathbf{z})}{p^{*}(y = 1)p^{*}(y = 0 \mid \mathbf{z})}$$

$$= \frac{p^{*}(y = 1 \mid \mathbf{z})}{p^{*}(y = 0 \mid \mathbf{z})} \equiv \frac{D(\mathbf{z})}{1 - D(\mathbf{z})}, \tag{12}$$

where $p^*(y=0)$ equals $p^*(y=1)$ since the number of samples is the same. We model $D(\mathbf{z})$ by $\sigma(T_{\psi}(\mathbf{z}))$, where $T_{\psi}(\mathbf{z})$ is a neural network with parameter ψ and input \mathbf{z} , and $\sigma(\cdot)$ is a sigmoid function. We train $T_{\psi}(\mathbf{z})$ to maximize the following objective function:

$$T^*(\mathbf{z}) = \max_{\psi} \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\ln(\sigma(T_{\psi}(\mathbf{z}))) \right] + \mathbb{E}_{p(\mathbf{z})} \left[\ln(1 - \sigma(T_{\psi}(\mathbf{z}))) \right]. \quad (13)$$

By using $T^*(\mathbf{z})$, we can estimate the density ratio $q_{\phi}(\mathbf{z})/p(\mathbf{z})$ as follows:

$$\frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z})} = \frac{\sigma(T^*(\mathbf{z}))}{1 - \sigma(T^*(\mathbf{z}))} \Leftrightarrow T^*(\mathbf{z}) = \ln \frac{q_{\phi}(\mathbf{z})}{p(\mathbf{z})}.$$
 (14)

Therefore, we can estimate the KL divergence with the aggregated posterior $D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || q_{\phi}(\mathbf{z}))$ by

$$D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || q_{\phi}(\mathbf{z}))$$

$$= D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || p(\mathbf{z})) - \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} [T^{*}(\mathbf{z})]. \quad (15)$$

3.2 Optimization Procedure

From the above discussion, we obtain the training objective function of the VAE with our implicit optimal prior:

$$\max_{\theta, \phi} \int p_{\mathcal{D}}(\mathbf{x}) \left\{ -D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || p(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x})} \left[\ln p_{\theta}(\mathbf{x} \mid \mathbf{z}) + T_{\psi}(\mathbf{z}) \right] \right\} d\mathbf{x}, \quad (16)$$

where $T_{\psi}(\mathbf{z})$ maximizes the Eq. (13). Given a dataset $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$, we optimize the Monte Carlo approximation of this objective:

$$\max_{\theta,\phi} \frac{1}{N} \sum_{i=1}^{N} \left\{ -D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}^{(i)}) || p(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z} \mid \mathbf{x}^{(i)})} \left[\ln p_{\theta}(\mathbf{x}^{(i)} \mid \mathbf{z}) + T_{\psi}(\mathbf{z}) \right] \right\}, \quad (17)$$

and we approximate the expectation term by the reparameterization trick (Kingma and Welling 2013):

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\ln p_{\theta}(\mathbf{x}^{(i)} \mid \mathbf{z}) + T_{\psi}(\mathbf{z}) \right]$$

$$\simeq \frac{1}{L} \sum_{\ell=1}^{L} \left\{ \ln p_{\theta}(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i,\ell)}) + T_{\psi}(\mathbf{z}^{(i,\ell)}) \right\}, \quad (18)$$

where $\mathbf{z}^{(i,\ell)} = \mu_{\phi}(\mathbf{x}^{(i)}) + \varepsilon^{(i,\ell)} \odot \sigma_{\phi}(\mathbf{x}^{(i)})$, $\varepsilon^{(i,\ell)}$ is a sample drawn from $\mathcal{N}(\mathbf{z}|\mathbf{0},\mathbf{I})$, \odot is the element-wise product, and L is the sample size of the reparameterization trick. Then, the resulting objective function is

$$\max_{\theta,\phi} \frac{1}{N} \sum_{i=1}^{N} \left[-D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}^{(i)}) || p(\mathbf{z})) + \frac{1}{L} \sum_{\ell=1}^{L} \left\{ \ln p_{\theta}(\mathbf{x}^{(i)} \mid \mathbf{z}^{(i,\ell)}) + T_{\psi}(\mathbf{z}^{(i,\ell)}) \right\} \right]. \quad (19)$$

We optimize this model with stochastic gradient descent (SGD) (Duchi, Hazan, and Singer 2011; Zeiler 2012; Tieleman and Hinton 2012; Kingma and Ba 2014) by iterating a two-step procedure: we first update θ and ϕ to maximize Eq. (19) with fixed ψ and next update ψ to maximize the Monte Carlo approximation of Eq. (13) with fixed θ and ϕ , as follows:

$$\max_{\psi} \frac{1}{M} \sum_{i=1}^{M} \ln(\sigma(T_{\psi}(\mathbf{z}_{1}^{(i)}))) + \frac{1}{M} \sum_{j=1}^{M} \ln(1 - \sigma(T_{\psi}(\mathbf{z}_{0}^{(j)}))), \quad (20)$$

Algorithm 1 VAE with Implicit Optimal Priors

```
1: while not converged do
          for J_1 steps do
               Sample minibatch \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}\} from X
 3:
               Compute the gradients of Eq. (19) w.r.t. \theta and \phi
 4:
               Update \theta and \phi with their gradients
 5:
 6:
          end for
          for J_2 steps do
 7:
               Sample minibatch \left\{\mathbf{z}_0^{(1)},\dots,\mathbf{z}_0^{(K)}\right\} from p(\mathbf{z})
 8:
               Sample minibatch \left\{\mathbf{z}_1^{(1)}, \dots, \mathbf{z}_1^{(K)}\right\} from q_{\phi}(\mathbf{z})
 9:
10:
               Compute the gradient of Eq. (20) w.r.t. \psi
11:
               Update \psi with its gradient
12:
          end for
13: end while
```

where $\mathbf{z}_1^{(i)}$ is a sample drawn from $q_{\phi}(\mathbf{z})$, $\mathbf{z}_0^{(j)}$ is a sample drawn from $p(\mathbf{z})$, and M is the sampling size of Monte Carlo approximation. Note that we need to compute the gradient of $T_{\psi}(\mathbf{z})$ with respect to ϕ in the optimization of Eq. (19) since $T_{\psi}(\mathbf{z})$ models $\ln q_{\phi}(\mathbf{z})/p(\mathbf{z})$. However, when $T_{\psi}(\mathbf{z})$ equals $T^*(\mathbf{z})$, the expectation of this gradient becomes zero, as follows:

$$\mathbb{E}_{p_{\mathcal{D}}(\mathbf{x})q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\nabla_{\phi} T^{*}(\mathbf{z}) \right] = \mathbb{E}_{q_{\phi}(\mathbf{z})} \left[\nabla_{\phi} \ln q_{\phi}(\mathbf{z}) \right]$$

$$= \int q_{\phi}(\mathbf{z}) \frac{\nabla_{\phi} q_{\phi}(\mathbf{z})}{q_{\phi}(\mathbf{z})} d\mathbf{z} = \nabla_{\phi} \int q_{\phi}(\mathbf{z}) d\mathbf{z} = \nabla_{\phi} 1 = 0.$$
(21)

Therefore, we ignore this gradient in the optimization 1 . We also note that $T_{\psi}(\mathbf{z})$ is likely to overfit to the log density ratio between the empirical aggregated posterior (Eq. (8)) and the standard Gaussian distribution. As mentioned in Section 2.3, this over-fitting also incurs over-fitting of the VAE (Tomczak and Welling 2018). Therefore, we use the regularization techniques such as dropout (Srivastava et al. 2014) for $T_{\psi}(\mathbf{z})$, which prevents it from over-fitting. We train ψ more than θ and ϕ : if we update θ and ϕ for J_1 steps, we update ψ for J_2 steps, where J_2 is larger than J_1 . Algorithm 1 shows the pseudo code of the optimization procedure of this model, where K is the minibatch size of SGD.

4 Related Work

For improving the density estimation performance of the VAE, numerous works have focused on the regularization effect of the KL divergence between the encoder and the prior. These works improve either the encoder or the prior.

First, we focus on the works about the prior. Although the optimal prior for the VAE is the aggregated posterior, the KL divergence with the aggregated posterior cannot be calculated in a closed form. As described in Section 2.3, the VampPrior (Tomczak and Welling 2018) has been presented to solve this problem. However, it has sensitive hyperparameters such as the number of mixtures K. Since the

¹There is almost the same discussion in (Mescheder, Nowozin, and Geiger 2017).

VampPrior requires a heavy computational cost, these hyperparameters are difficult to tune. In contrast to this, our approach can estimate the KL divergence more easily and robustly than the VampPrior since it does not need to model the aggregated posterior explicitly. In addition, since the computational cost of our approach is much more lightweight than that of VampPrior, the hyperparameters of our approach are easier to tune than those of VampPrior.

There are approaches on improving the prior other than the aggregated posterior. For example, non-parametric Bayesian distribution (Nalisnick and Smyth 2017) and hyperspherical distribution (Davidson et al. 2018) are used for the prior. These approaches aim to obtain the useful and interpretable latent representation rather than improving the density estimation performance, which is opposite to our purpose. We should mention the disadvantage of our approach compared with these approaches. Since our prior is implicit, we cannot sample from our prior directly. Instead, we can sample from the aggregated posterior, which our implicit prior models, by using ancestral sampling. That is, when we sample from the prior, we need to prepare a data point.

Next, we focus on the works about the encoder. To improve the density estimation performance, these works increase the flexibility of the encoder. The normalizing flow (Rezende and Mohamed 2015; Kingma et al. 2016; Huang et al. 2018) is one of the main approaches, which applies a sequence of invertible transformations to the latent variable vector until a desired level of flexibility is attained. Our approach is orthogonal to the normalizing flow and can be used together with it.

The similar approaches to ours are the adversarial variational Bayes (AVB) (Mescheder, Nowozin, and Geiger 2017) and the adversarial autoencoders (AAE) (Makhzani et al. 2015; Tolstikhin et al. 2017). These approaches use the implicit encoder network, which takes as input a data point x and Gaussian random noise and produces a latent variable vector z. Since the implicit encoder does not assume the distribution type, it can become a very flexible distribution. In these approaches, the standard Gaussian distribution is used for the prior. Although the KL divergence between the implicit encoder and the standard Gaussian prior $D_{KL}(q_{\phi}(\mathbf{z} \mid \mathbf{x}) || p(\mathbf{z}))$ cannot be calculated in a closed form, the AVB estimates this KL divergence by using the density ratio trick. However, this estimation does not work well with high-dimensional datasets since this KL divergence also becomes a high-dimensional density ratio (Rosca, Lakshminarayanan, and Mohamed 2018). Our approach can avoid this problem since we use the density ratio trick in a low dimension. The AAE is an expansion of the Autoencoder rather than the VAE. The AAE regularizes the aggregated posterior to be close to the standard Gaussian prior by minimizing the KL divergence $D_{KL}(q_{\phi}(\mathbf{z})||p(\mathbf{z}))$. The AAE also uses the density ratio trick to estimate this KL divergence, and this works well since this KL divergence is a low-dimensional density ratio. However, the AAE cannot estimate the probability of a data point. Our approach is based on the VAE, and can estimate the probability of a data point.

Table 1: Number and dimensions of datasets

	Dimension	Train size	Valid size	Test size
OneHot	4	1,000	100	1,000
MNIST	784	50,000	10,000	10,000
OMNIGLOT	784	23,000	1,345	8,070
FreyFaces	560	1,565	200	200
Histopathology	784	6,800	2,000	2,000

5 Experiments

In this section, we experimentally evaluate the density estimation performance of our approach.

5.1 Data

We used five datasets: OneHot (Mescheder, Nowozin, and Geiger 2017), MNIST (Salakhutdinov and Murray 2008), OMNIGLOT (Burda, Grosse, and Salakhutdinov 2015), FreyFaces², and Histopathology (Tomczak and Welling 2016). OneHot consists of only four-dimensional one hot vectors: $(1,0,0,0)^{\mathrm{T}}$, $(0,1,0,0)^{\mathrm{T}}$, $(0,0,1,0)^{\mathrm{T}}$, and $(0,0,0,1)^{\mathrm{T}}$. This simple dataset is useful for observing the posterior of the latent variable, which is used in (Mescheder, Nowozin, and Geiger 2017). MNIST and OMNIGLOT are binary image datasets, and FreyFaces and Histopathology are grayscale image datasets. These image datasets are useful for measuring the density estimation performance, which are used in (Tomczak and Welling 2018). The number and the dimensions of data points of the five datasets are listed in Table 1.

5.2 Setup

We compared our implicit optimal prior with standard Gaussian prior and VampPrior. We set the dimensions of the latent variable vector to 2 for OneHot, and 40 for other datasets. We used two-layer neural networks (500 hidden units per layer) for the encoder, the decoder, and the density ratio estimator. We used the gating mechanism (Dauphin et al. 2016) for the encoder and the decoder and used a hyperbolic tangent as the activation function for the density ratio estimator. We initialized the weights of these neural networks in accordance with the method in (Glorot and Bengio 2010). We used a Gaussian distribution as the encoder. As the decoder, we used a Bernoulli distribution for OneHot, MNIST, and OMNIGLOT and used a Gaussian distribution for Frey-Faces and Histopathology, means of which were constrained to the interval [0, 1] by using a sigmoid function. We trained all methods by using Adam (Kingma and Ba 2014) with a mini-batch size of 100 and learning rate in $[10^{-4}, 10^{-3}]$. We set the maximum number of epochs to 1,000 and used earlystopping (Goodfellow, Bengio, and Courville 2016) on the basis of validation data. We set the sample size of the reparameterization trick to L=1. In addition, we used warmup (Bowman et al. 2015) for the first 100 epochs of Adam.

²This dataset is available at https://cs.nyu.edu/~roweis/data/frey_rawface.mat

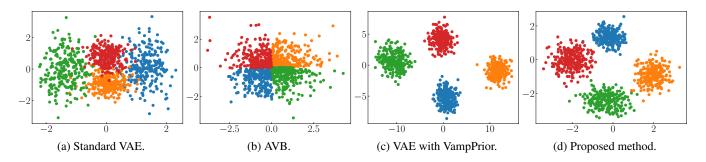


Figure 1: Comparison of posteriors of latent variable on OneHot. We plotted samples drawn from $q_{\phi}(\mathbf{z} \mid \mathbf{x})$, where \mathbf{x} is a one hot vector: $(1,0,0,0)^{\mathrm{T}}$, $(0,1,0,0)^{\mathrm{T}}$, $(0,0,1,0)^{\mathrm{T}}$, or $(0,0,0,1)^{\mathrm{T}}$. We used test data for this sampling. Samples in each color correspond to each latent representation of one hot vectors. (a) Standard VAE (VAE with standard Gaussian prior). (b) AVB. (c) VAE with VampPrior. (d) Proposed method.

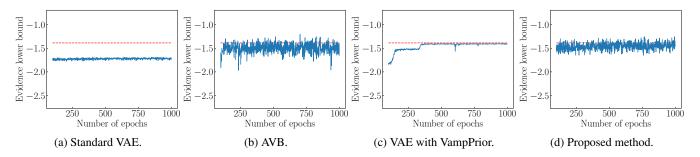


Figure 2: Comparison of the evidence lower bound (ELBO) with validation data on OneHot. We plotted the ELBO from 100 to 1,000 epochs since we used warm-up for the first 100 epochs. The optimal log-likelihood on this dataset is $-\ln(4) \approx -1.386$. We plotted this value by a dashed line for comparison. (a) Standard VAE (VAE with standard Gaussian prior). (b) AVB. (c) VAE with VampPrior. (d) Proposed method.

For MNIST and OMNIGLOT, we used dynamic binarization (Salakhutdinov and Murray 2008) during the training of VAE to avoid over-fitting. For image datasets, we calculated the log marginal likelihood of the test data by using the importance sampling (Burda, Grosse, and Salakhutdinov 2015). We set the sample size of the importance sampling to 10. We ran all experiments eight times each.

With VampPrior, we set the number of mixtures K to 50 for OneHot, 500 for MNIST, FreyFaces, and Histopathology, and 1,000 for OMNIGLOT. In addition, for image datasets, we used a clipped relu function that equals $\min(\max(x,0),1)$ to scale the pseudo inputs in [0,1] since the range of data points of these datasets is [0,1]³.

With our approach, we used dropout (Srivastava et al. 2014) in the training of the density ratio estimator since it is likely to over-fit. We set the keep probability of dropout to 50%. We updated the parameter of the density ratio estimator: ψ for 10 epochs during the updating of the parameters of VAE: θ and ϕ for one epoch. We set the sampling size of Monte Carlo approximation in Eq. (20) to M=N.

In addition, we compared our approach with adversarial variational Bayes (AVB) on OneHot. We set the dimension of the Gaussian random noise input of AVB to 10, and other

settings are almost the same as those for our approach.

5.3 Results

Figures 1a–1d show the posteriors of latent variable of each approach on OneHot, and Figures 2a–2d show the evidence lower bound of each approach on OneHot.

These results show the difference between these approaches. We can see that the evidence lower bound (ELBO) of the standard VAE (VAE with standard Gaussian prior) on OneHot was worse than the optimal log-likelihood on this dataset: $-\ln(4)\approx -1.386.$ The over-regularization incurred by the standard Gaussian prior can be given as a reason. The posteriors were overlapped, and it became difficult to discriminate between samples from these posteriors. Hence, the decoder became confused when reconstructing. This caused the poor density estimation performance.

On the other hand, the ELBOs of AVB, VAE with Vamp-Prior, and our approach are much closer to the optimal log-likelihood than the standard VAE. We note that the ELBOs of the AVB and our approach are the estimated values, and that these approaches may overestimate the ELBO on One-Hot since the training data and validation data of OneHot are the same. First, we focus on the AVB. Although there is still the strong regularization by the standard Gaussian prior, the posteriors barely overlapped, and the data point was easy to reconstruct from the latent representation. The reason is

³We referred to https://github.com/jmtomczak/ vae_vampprior

Table 2: Comparison of test log-likelihoods on four image datasets.

	MNIST	OMNIGLOT	FreyFaces	Histopathology
Standard VAE	-85.84 ± 0.07	-111.39 ± 0.11	1382.53 ± 3.57	1081.53 ± 0.70
VAE with VampPrior	-83.90 ± 0.08	-110.53 ± 0.09	1392.62 ± 6.25	1083.11 ± 2.10
Proposed method	$pprox$ -83.21 \pm 0.13	$pprox$ -108.48 \pm 0.16	$pprox$ 1396.27 \pm 2.75	$\approx \textbf{1087.42} \pm \textbf{0.60}$

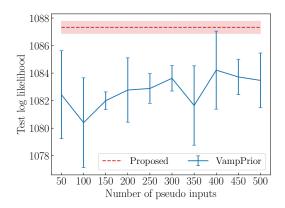


Figure 3: Relationship between the test log-likelihoods and number of pseudo inputs of VampPrior on Histopathology. We plotted the test log-likelihoods of our approach by a dashed line for comparison. The semi-transparent area and error bar represent standard deviations.

that the implicit encoder network of AVB can learn complex posterior distributions. Next, we focus on the VAE with VampPrior and our approach. The VampPrior and our implicit optimal prior model the aggregated posterior that is the optimal prior for the VAE. These priors made the posteriors of these approaches different from each other, and the data point was easy to reconstruct from the latent representation.

Table 2 compares the test log-likelihoods on four image datasets. We used bold to highlight the best result and the results that are not statistically different from the best result according to a pair-wise t-test. We used 5% as the p-value. We did not compare with AVB since the estimated log marginal likelihood of AVB with high-dimensional datasets such as images is not accurate (Rosca, Lakshminarayanan, and Mohamed 2018).

First, we focus on the VampPrior. We can see that test log-likelihoods of VampPrior are better than those of standard VAE. However, we found two drawbacks with the VampPrior. One is that the pseudo inputs of VampPrior are difficult to optimize. For example, the pseudo inputs have an initial value dependence. Although the warm-up helps in solving this problem, it seems difficult to solve completely. The other is that the number of mixtures K is a sensitive hyperparameter. Figure 3 shows the test log-likelihoods with various K on Histopathology. The high standard deviation of the VampPrior indicates its high dependence of the pseudo input initial values. In addition, even though we choose the optimal K, the test log-likelihood of the VampPrior is worse than that of our approach.

Next, we focus on our approach. Our approach obtained the equal to or better density estimation performance than the VampPrior. Since our approach models the aggregated posterior implicitly, it can estimate the KL divergence more easily and robustly than the VampPrior. In addition, it has a much more lightweight computational cost than the VampPrior. In the training phase on MNIST, our approach was almost 2.83 times faster than the VampPrior. Therefore, although our approach has as many hyperparameters, like the neural architecture of the density ratio estimator, as the VampPrior, these hyperparameters are easier to tune than those of the VampPrior.

These results indicate that our implicit optimal prior is a good alternative to the VampPrior: our implicit optimal prior can be optimized easily and robustly, and its density estimation performance is equal to or better than that of the VAE with the VampPrior.

6 Conclusion

In this paper, we proposed the variational autoencoder (VAE) with implicit optimal priors. Although the standard Gaussian distribution is usually used for the prior, this simple prior incurs over-regularization, which is one of the causes of poor density estimation performance. To improve the density estimation performance, the aggregated posterior has been introduced as a sophisticated prior, which is optimal in terms of maximizing the training objective function of VAE. However, Kullback Leibler (KL) divergence between the encoder and the aggregated posterior cannot be calculated in a closed form, which prevents us from using this optimal prior. Even though explicit modeling of the aggregated posterior has been tried, this optimal prior is difficult to model explicitly.

With the proposed method, we introduced the density ratio trick for estimating this KL divergence directly. Since the density ratio trick can estimate the density ratio between two distributions without modeling each distribution explicitly, there is no need to model the aggregated posterior explicitly. Although the density ratio trick is useful, it does not work well in a high dimension. Unfortunately, the KL divergence between the encoder and the aggregated posterior is highdimensional. Hence, we rewrite the KL divergence into the sum of two terms: the KL divergence between the encoder and the standard Gaussian distribution that can be calculated in a closed form, and the low-dimensional density ratio between the aggregated posterior and the standard Gaussian distribution, to which the density ratio trick is applied. We experimentally showed the high density estimation performance of the VAE with this implicit optimal prior.

References

- [Bowman et al. 2015] Bowman, S. R.; Vilnis, L.; Vinyals, O.; Dai, A. M.; Jozefowicz, R.; and Bengio, S. 2015. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*.
- [Burda, Grosse, and Salakhutdinov 2015] Burda, Y.; Grosse, R.; and Salakhutdinov, R. 2015. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*.
- [Dauphin et al. 2016] Dauphin, Y. N.; Fan, A.; Auli, M.; and Grangier, D. 2016. Language modeling with gated convolutional networks. *arXiv* preprint arXiv:1612.08083.
- [Davidson et al. 2018] Davidson, T. R.; Falorsi, L.; De Cao, N.; Kipf, T.; and Tomczak, J. M. 2018. Hyperspherical variational auto-encoders. *arXiv preprint arXiv:1804.00891*.
- [Dilokthanakul et al. 2016] Dilokthanakul, N.; Mediano, P. A.; Garnelo, M.; Lee, M. C.; Salimbeni, H.; Arulkumaran, K.; and Shanahan, M. 2016. Deep unsupervised clustering with Gaussian mixture variational autoencoders. *arXiv* preprint arXiv:1611.02648.
- [Duchi, Hazan, and Singer 2011] Duchi, J.; Hazan, E.; and Singer, Y. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12(Jul):2121–2159.
- [Glorot and Bengio 2010] Glorot, X., and Bengio, Y. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 249–256.
- [Goodfellow, Bengio, and Courville 2016] Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep Learning*. MIT Press. http://www.deeplearningbook.org.
- [Goodfellow et al. 2014] Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. In *Advances in neural information processing systems*, 2672–2680.
- [Gregor et al. 2015] Gregor, K.; Danihelka, I.; Graves, A.; Rezende, D.; and Wierstra, D. 2015. DRAW: A recurrent neural network for image generation. In *Proceedings of the 32nd International Conference on Machine Learning*, 1462–1471.
- [Gulrajani et al. 2016] Gulrajani, I.; Kumar, K.; Ahmed, F.; Taiga, A. A.; Visin, F.; Vazquez, D.; and Courville, A. 2016. PixelVAE: A latent variable model for natural images. *arXiv* preprint arXiv:1611.05013.
- [Hoffman and Johnson 2016] Hoffman, M. D., and Johnson, M. J. 2016. ELBO surgery: yet another way to carve up the variational evidence lower bound. In *Workshop in Advances in Approximate Bayesian Inference*, NIPS.
- [Hsu, Zhang, and Glass 2017] Hsu, W.-N.; Zhang, Y.; and Glass, J. 2017. Learning latent representations for speech generation and transformation. *Proc. Interspeech* 2017 1273–1277.
- [Huang et al. 2018] Huang, C.-W.; Krueger, D.; Lacoste, A.; and Courville, A. 2018. Neural autoregressive flows. In *Pro-*

- ceedings of the 35th International Conference on Machine Learning, 2078–2087.
- [Kingma and Ba 2014] Kingma, D., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Kingma and Welling 2013] Kingma, D. P., and Welling, M. 2013. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*.
- [Kingma et al. 2016] Kingma, D. P.; Salimans, T.; Jozefowicz, R.; Chen, X.; Sutskever, I.; and Welling, M. 2016. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, 4743–4751.
- [Makhzani et al. 2015] Makhzani, A.; Shlens, J.; Jaitly, N.; Goodfellow, I.; and Frey, B. 2015. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*.
- [Mescheder, Nowozin, and Geiger 2017] Mescheder, L.; Nowozin, S.; and Geiger, A. 2017. Adversarial Variational Bayes: Unifying variational autoencoders and generative adversarial networks. In *International Conference on Machine Learning*, 2391–2400.
- [Nalisnick and Smyth 2017] Nalisnick, E., and Smyth, P. 2017. Stick-breaking variational autoencoders. In *International Conference on Learning Representations (ICLR)*.
- [Rezende and Mohamed 2015] Rezende, D., and Mohamed, S. 2015. Variational inference with normalizing flows. In *International Conference on Machine Learning*, 1530–1538.
- [Rezende, Mohamed, and Wierstra 2014] Rezende, D. J.; Mohamed, S.; and Wierstra, D. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, 1278–1286.
- [Rosca, Lakshminarayanan, and Mohamed 2018] Rosca, M.; Lakshminarayanan, B.; and Mohamed, S. 2018. Distribution matching in variational inference. *arXiv preprint arXiv:1802.06847*.
- [Salakhutdinov and Murray 2008] Salakhutdinov, R., and Murray, I. 2008. On the quantitative analysis of deep belief networks. In *Proceedings of the 25th international conference on Machine learning*, 872–879. ACM.
- [Srivastava et al. 2014] Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research* 15(1):1929–1958.
- [Sugiyama, Suzuki, and Kanamori 2012] Sugiyama, M.; Suzuki, T.; and Kanamori, T. 2012. *Density ratio estimation in machine learning*. Cambridge University Press.
- [Tieleman and Hinton 2012] Tieleman, T., and Hinton, G. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning* 4(2):26–31.
- [Tolstikhin et al. 2017] Tolstikhin, I.; Bousquet, O.; Gelly, S.; and Schoelkopf, B. 2017. Wasserstein auto-encoders. *arXiv preprint arXiv:1711.01558*.

- [Tomczak and Welling 2016] Tomczak, J. M., and Welling, M. 2016. Improving variational auto-encoders using householder flow. *arXiv preprint arXiv:1611.09630*.
- [Tomczak and Welling 2018] Tomczak, J. M., and Welling, M. 2018. VAE with a VampPrior. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, 1214–1223.
- [van den Oord, Vinyals, and kavukcuoglu 2017] van den Oord, A.; Vinyals, O.; and kavukcuoglu, k. 2017. Neural discrete representation learning. In *Advances in Neural Information Processing Systems*, 6309–6318.
- [Zeiler 2012] Zeiler, M. D. 2012. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*.