

---

# $\rho$ -VAE: Autoregressive parametrization of the VAE encoder

---

Sohrab Ferdowsi, Shideh Rezaeifar, Maurits Diephuis, and Slava Voloshynovskiy

Department of Computer Science, University of Geneva, Switzerland  
{sohrab.ferdowsi, shideh.rezaeifar, maurits.diephuis, svolos}@unige.ch

## Abstract

We make a simple, but very effective plug-and-play alteration to the standard VAE. This is about parametrization of the approximate posterior of the latent space as a Gaussian distribution whose covariance matrix is constructed as an AR(1) process. We argue that this is a more natural prior for images to consider, as compared to the standard Gaussian distribution with a diagonal covariance matrix. While the standard diagonal parametrization consists of learning the mean vector, as well as the vector of diagonal values for each sample, our parametrization consists of the mean and two scalar values  $s$  and  $\rho$ , where the first one is scaling the distribution and the second one is a measure of correlation among the latent dimensions. Therefore, the performance boost which we show in our experiments to be consistently noticeable across different setups and variants of VAE models, comes at no cost. We even reduce the number of covariance parameters from  $d$ , the dimensionality of the latent space to only 2, yet providing a more flexible approximation of the intractable posterior.

## 1 Introduction

Arguably, one of the most successful approaches to representation and generative learning is that of “Auto-encoding variational Bayes” [1], a parametrization of the standard variational Bayes in the form of an autoencoder neural network with a recipe for end-to-end learning of its parameters, while providing effective approximation of the intractable posterior. This has then given rise to the very popular Variational AutoEncoder (VAE), a set of models

bla bla

bla bla bla

## 2 The VAE models

Here we briefly review the standard VAE model, highlighting aspects relevant to our work, as well as some of its further developments.

In a typical probabilistic model where a latent variable  $\mathbf{z} \in \mathbb{R}^d$  is the underlying factor to generate the observable samples  $\mathbf{x}'s \in \mathbb{R}^n$ , the standard variational Bayes [?] paradigm is concerned with finding an approximation  $q(\mathbf{z})$  for the intractable posterior  $p(\mathbf{z}|\mathbf{x})$ . This is achieved by minimizing the Kullback-Leibler divergence between these two distributions, i.e.,  $D_{\text{KL}}[q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})]$ . Standard

treatments of this quantity, along with its non-negativity property will then amount to the following inequality:

$$\log(p(\mathbf{x})) \leq \mathbb{E}_{q(\mathbf{z})} \left[ \log(p(\mathbf{x}|\mathbf{z})) \right] - D_{\text{KL}} \left[ q(\mathbf{z}) || p(\mathbf{z}) \right]. \quad (1)$$

## 2.1 The standard VAE

Autoencoding variational Bayes [1] is then constructing an explicit dependence of the latent variables to the  $i^{\text{th}}$  training sample by considering a parametrized distribution  $q_\phi(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})$  for the approximate posterior, whose construction resembles the encoder part of an autoencoder network with a set of learnable weights  $\phi$ . Furthermore the training samples can be decoded with  $p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i)})$ , another network with parameters symbolized as  $\theta$ .

Making this double-sided data dependence more explicit, and by summing over all  $N$  training samples results to the following inequality:

$$\frac{1}{N} \sum_{i=1}^N \log(p(\mathbf{x}^{(i)})) \leq \frac{1}{N} \sum_{i=1}^N \left[ \log(p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i)})) - D_{\text{KL}} [q_\phi(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}) || p(\mathbf{z})] \right]. \quad (2)$$

This, in fact, is highly relevant for generative modeling as the marginal log-likelihood of the training samples will be upper bounded by two terms, both of which amenable to mini-batch optimization with stochastic gradient descent.

During optimization, the first term of the LHS can be considered as a data fidelity term, minimized e.g., in the  $\ell_2$  sense, since a natural choice for the decoder is  $p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i)}) = \mathcal{N}(\mathbf{x}^{(i)}|\mathbf{z}^{(i)}, \sigma^2 \mathbf{I}_n)$ , where  $\mathbf{I}_n$  is the  $n$ -dimensional unity matrix.

The second term, from the other hand, can be interpreted as a regularization term, pushing the approximate posterior to a prior imposed on the latent space, most conveniently a simple  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Provided that the optimization is successful, and the inequality (??) is tight, one can generate random samples from this prior, pass it through the learned decoder and generate samples (non-trivially) similar to the underlying data.

However, the above scenario comes with a major caution: the fact that sampling  $\mathbf{z}$  from  $q_\phi(\mathbf{z}|\mathbf{x})$  is a non-differentiable operation. The work-around for this issue is the wise ‘‘reparametrization trick’’, as proposed in [1].

The idea is to create the required randomness from a fixed distribution  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . The samples of the appropriate distribution can then be generated by injecting the learnable moments, e.g., using  $\mathbf{z}^{(i)} = \boldsymbol{\mu}^{(i)} + \tilde{\mathbf{C}}^{(i)} \epsilon$ , where  $\boldsymbol{\mu}^{(i)}$  is the mean vector of the posterior learned for the  $i^{\text{th}}$  sample and  $\tilde{\mathbf{C}}^{(i)}$  is the Choleskiy decomposition of the corresponding covariance matrix  $\mathbf{C}^{(i)}$ .

This then limits the practical choices for  $\mathbf{C}^{(i)}$  to have analytical Choleskiy decomposition forms, since both  $\mathbf{C}^{(i)}$  and  $\tilde{\mathbf{C}}^{(i)}$  participate in the optimization simultaneously.

Another issue to address is the calculation of  $D_{\text{KL}} [q_\phi(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}) || p(\mathbf{z})]$ , which we elaborate slightly more in section 2.2. In order to avoid many practical difficulties, the standard choice is to pick a closed-form expression for it, hence further limiting the choices of  $\mathbf{C}^{(i)}$ .

While the prior distribution is chosen as  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , considering the above two constraints, the standard choice widely adopted in many further variants for the sample-wise approximate posterior is to set  $\mathbf{C}_{(\mathbf{s})}^{(i)} = \text{diag}(\mathbf{s}^{(i)})$ . In other words,  $q_\phi(\mathbf{z}^{(i)}|\mathbf{x}^{(i)}) = \mathcal{N}(\boldsymbol{\mu}^{(i)}, \text{diag}(\mathbf{s}^{(i)}))$ , a diagonal Gaussian distribution parametrized by the pair  $(\boldsymbol{\mu}^{(i)}, \mathbf{s}^{(i)})$ .

Note that now, the reparametrization trick can run smoothly, since the Choleskiy decomposition has a closed expression as  $\tilde{\mathbf{C}}_{(\mathbf{s})}^{(i)} = \text{diag}(\sqrt{\mathbf{s}^{(i)}})$ . Furthermore, the regularization term  $D_{\text{KL}}[q_{\phi}(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})||p(\mathbf{z})]$  is also calculated analytically as:

$$D_{\text{KL}}[\mathcal{N}(\boldsymbol{\mu}^{(i)}, \text{diag}(\mathbf{s}^{(i)}))||\mathcal{N}(\mathbf{0}, \mathbf{I}_d)] = \frac{1}{2} [\mathbf{1}_d^T \mathbf{s}^{(i)} + \|\boldsymbol{\mu}^{(i)}\|_2^2 - d - \mathbf{1}_d^T \log(\mathbf{s}^{(i)})], \quad (3)$$

where  $\mathbf{1}_d$  is the unity vector of dimension  $d$ ,  $\|\cdot\|_2^2$  is the squared  $\ell_2$ -norm, and  $\log(\mathbf{s}^{(i)})$  is applied element-wise.

While this is a very practical choice, we argue in section 3 that it is too simplistic, as it disregards any correlation within dimensions.

## 2.2 Further variations

The literature around VAE is immense and still very active. Without aiming for any comprehensive literature review, here we still point out several of its variants.

As categorized in [2], VAE variants come in 3 categories.

## 3 The $\rho$ -VAE

We saw in section 2.1 that two considerations limit the choices of approximate posterior: the need for a parametric Choleskiy factorization of its covariance matrix, as well as closed-form expression for the regularization term of (2), which basically requires the expression of log-determinant of the covariance.

In spite of the general consensus to pick  $\mathbf{C}_{(\mathbf{s})}^{(i)} = \text{diag}(\mathbf{s}^{(i)})$ , which does not allow any correlation between the dimensions of the approximate posterior, this work proposes another parametrization that grants such freedom, satisfies the above-mentioned restrictions, and yet has less number of parameters.

In particular, we chose a first-order autoregressive covariance which is characterized by a scaling factor  $s$ , and another scalar  $\rho$  to control the level of correlation, hence the term  $\rho$ -VAE. This has the form of a simple symmetric Toeplitz matrix as the following:

$$\mathbf{C}_{(\rho, s)} = s \times \text{Toeplitz}([1, \rho, \rho^2, \dots, \rho^{d-1}]) = s \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 & \dots & \rho^{d-1} \\ \rho & 1 & \rho & \rho^2 & \dots & \rho^{d-2} \\ \rho^2 & \rho & 1 & \rho & \dots & \rho^{d-3} \\ \rho^3 & \rho^2 & \rho & 1 & \dots & \rho^{d-4} \\ \vdots & & & & \ddots & \vdots \\ \rho^{d-1} & \dots & \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}, \quad (4)$$

where  $s$  is a positive scalar, and the correlation parameter is bounded as  $-1 < \rho < +1$ .

The determinant for this matrix can be calculated as [3]:

$$\det(\mathbf{C}_{(\rho, s)}) = s^d (1 - \rho^2)^{d-1}, \quad (5)$$

based on which we can derive the regularization term of the loss function as:

$$D_{\text{KL}}[\mathcal{N}(\boldsymbol{\mu}^{(i)}, \mathbf{C}_{(\rho, s)})||\mathcal{N}(\mathbf{0}, \mathbf{I}_d)] = \frac{1}{2} [\|\boldsymbol{\mu}^{(i)}\|_2^2 + d(s - 1 - \log(s)) - (d-1) \log(1 - \rho^2)]. \quad (6)$$

As far as the reparametrization trick is concerned, the Choleskiy decomposition of our choice of covariance matrix has the following lower triangular form:

$$\tilde{C}_{(\rho,s)} = \frac{1}{\sqrt{s}} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ \rho & \sqrt{1-\rho^2} & 0 & 0 & \dots & 0 \\ \rho^2 & \rho\sqrt{1-\rho^2} & \sqrt{1-\rho^2} & 0 & \dots & 0 \\ \rho^3 & \rho^2\sqrt{1-\rho^2} & \rho\sqrt{1-\rho^2} & \sqrt{1-\rho^2} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^d & \dots & \rho^3\sqrt{1-\rho^2} & \rho^2\sqrt{1-\rho^2} & \rho\sqrt{1-\rho^2} & \sqrt{1-\rho^2} \end{bmatrix}, \quad (7)$$

which can be used to generate the latent codes as  $\mathbf{z}^{(i)} = \boldsymbol{\mu}^{(i)} + \tilde{C}_{(\rho,s)}^{(i)} \boldsymbol{\epsilon}$ , which can be constructed also as the element-wise product of  $C_{(\rho,s)}$  with another highly structured matrix.

Otherwise, if depending on the choice of the deep learning framework used, the realization of Toeplitz matrices is not straightforward, one can generate AR(1) samples directly from their definition, i.e.,  $\mathbf{z}^{(i)}[j] = \boldsymbol{\mu}^{(i)}[j] + \sqrt{s}\boldsymbol{\epsilon}[j] + \rho\mathbf{z}^{(i)}[j-1]$ , for  $1 < j \leq d$ .

Although it has less number of parameters than the standard choice and is hence more resilient towards over-fitting, this structure for the approximate posterior is more natural to consider, since correlation will somehow be represented.

Note that the fact that the prior is chosen as a white Gaussian by design, i.e.,  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , does not obviate the need for the per-sample approximate posterior to account for correlation. In fact, the per-sample posterior can be correlated, yet the aggregation of all samples can be a white Gaussian matching the prior.

Furthermore, the need for correlation does not solely stem from the natural signals like images being correlated. As a matter of fact, another requirement for the success of the VAE-based generative modeling is the tightness of the bound in (2), which is controlled by  $D_{\text{KL}}[q_\phi(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z}^{(i)}|\mathbf{x}^{(i)})]$ .

In other words, to guarantee a successful training, the approximate posterior should have enough capacity to match the unknown and intractable posterior. In VAE models, however, it is usually only ‘‘hoped’’ that this will be the case. We believe (albeit without providing quantitative evidence), that accounting for correlation may help reduce this gap.

Next we will show the effectiveness of our proposition. We show that the simple alterations to the standard approach, without the need for any sort of hyper-parameter tuning, will noticeably and consistently improve the performance under all variations considered and for all setups.

## 4 Experiments

### References

- [1] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [2] Michael Tschannen, Olivier Bachem, and Mario Lucic. Recent advances in autoencoder-based representation learning. *arXiv preprint arXiv:1812.05069*, 2018.
- [3] Harold Widom. On the eigenvalues of certain hermitian operators. *Transactions of the American Mathematical Society*, 88(2):491–522, 1958.