

Variational Autoencoders with variance loss

Blaž Stojanovič

Cavendish Laboratory, Department of Physics, JJ Thomson Avenue, Cambridge. CB3 0HE

Abstract

Contents

1	Introduction	1
1.1	Generative modelling	1
1.2	Deep latent variable models	1
1.2.1	Learning in DLVMs	2
1.3	Variational Autoencoders	3
1.3.1	Gradient estimation	3
1.3.2	Variants and advances	4
2	Implementation of Variance loss	5
2.1	Continuous latent space	6
2.2	Discrete latent space	7
3	Experiments	7
3.1	Continuous latent space	7
3.1.1	Fashion MNIST	7
3.1.2	CELEBa	7
3.2	Discrete latent space	7
3.2.1	Binarized MNIST	7
3.2.2	Binarized Omniglot	7
4	Conclusion	7
	Bibliography	7

1. Introduction

1.1. Generative modelling

Depending on what we learn in machine learning we differentiate between *discriminative* and *generative* modelling. While there is some inconsistency in the terminology, at their core discriminative models only have the ability to discriminate between targets Y and inputs X while generative models "simulate" the underlying data generation process and have the ability to generate new samples from the underlying distribution of data. Generative modelling is very much similar to the scientific process of providing competing

models of how physical phenomena work and testing them through observations. The models are considered highly intuitive and interpretable as we can encode our prior knowledge into them while treating unknown details as noise [1]. While there has been a recent surge of popularity in generative modeling, with models like the *variational autoencoder* (VAE) [2, 3], *Generative Adversarial Network* (GAN) [4] and *Normalizing Flows* [5] generative modelling in ML has a rich history. Models such as Gaussian mixture models [6], Hidden Markov models [7], Autoregressive models, Latent Dirichlet allocation [8], Boltzmann and Helmholtz machines [9, 10] have long been known and used. With increase in computational power, modern generative models have achieved astounding results in speech synthesis [11], image super-resolution [12] and have been implemented to aid supervised learning [13].

Ian Goodfellow provided a helpful taxonomy of generative models in his tutorial on GANs [4]. He differentiates generative models that maximise likelihood between how they compute or approximate the likelihood and its gradients. The model can construct an explicit probability density $p_{\theta}(\mathbf{x})$ and consequently an explicit likelihood which is then optimised. This explicit density can either be tractable as is the case in *fully observed Bayesian networks*, or it could be intractable requiring either variational approaches as in the VAE, or MCMC for example in Boltzmann machines. On the other hand we have models that do not construct the probability distribution explicitly, but instead provide a less direct way to interact with it i.e. draw samples from it. GANs fall into this category.

1.2. Deep latent variable models

We assume that an observable value \mathbf{x} is generated by an unknown underlying random process with unknown probability distribution $p^*(\mathbf{x})$. This is the process we try to approximate using a generative model

$p_{\theta}(\mathbf{x})$, which is parameterised with parameters θ . The process of learning is finding parameters θ such that our model distribution $p_{\theta}(\mathbf{x})$ is as close to the underlying distribution $p^*(\mathbf{x})$. Normally though, this is the case for regression and classification problems, we are modelling the conditional distribution $p^*(y|\mathbf{x})$ with our model $p_{\theta}(y|\mathbf{x})$, where y is conditioned on the input \mathbf{x} . The approaches described and implemented in this written assignment are equally applicable to both conditional and unconditional modelling. In either case we want our model to be versatile enough to represent the underlying structure of the data. The last two decades have seen a renaissance of neural networks as models of complex processes [14], they are *universal function approximators*, generally made of compositions of smaller differentiable sub-units with nonlinear activation functions and trained using gradient-based approaches via backpropagation. They can be used to parameterize conditional distributions and form basis of probabilistic models thus combining the best of both worlds, a probabilistic description and *hopefully* better scalability, making use of gradient-based learning. The categorical distribution $p_{\theta}(y|\mathbf{x})$ in the case of classification, where y is the class, given the input \mathbf{x} , can be parameterised as

$$\begin{aligned} \mathbf{p} &= f(\mathbf{x}; \theta), \\ p_{\theta}(y|\mathbf{x}) &= \text{Categorical}(y; \mathbf{p}), \end{aligned} \quad (1)$$

where $f(\cdot, \theta) : \mathbb{R}^{n_{\text{in}}} \rightarrow \mathbb{R}^{n_{\text{out}}}$ is the neural network, with a final layer such that the probabilities $\sum_i^{n_{\text{out}}} p_i = 1$.

We can combine random variables $(\mathbf{x}_1, \dots, \mathbf{x}_M)$ and represent their conditional dependence structure as a network to construct a *probabilistic graphical model*. If the structure of the network is a directed acyclic graph, the joint distribution of the random variables is factorised as a product of conditionals and priors

$$p_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_M) = \prod_{j=1}^M p_{\theta}(\mathbf{x}_j | \text{pa}(\mathbf{x}_j)), \quad (2)$$

where $\text{pa}(\mathbf{x}_j)$ are the parent nodes of node j in the directed graph. A graphical model with this structure is termed a *Bayesian network* and we can parameterise the conditional distributions in a Bayesian network using neural networks

$$\begin{aligned} \boldsymbol{\eta} &= f(\text{pa}(\mathbf{x}); \theta) \\ p_{\theta}(\mathbf{x} | \text{pa}(\mathbf{x})) &= p_{\theta}(\mathbf{x} | \boldsymbol{\eta}). \end{aligned} \quad (3)$$

A model which only contains observable variables \mathbf{x}_i is a *fully observed* model and is straightforward to learn.

Given a dataset \mathcal{D} of *independently and identically distributed* observations \mathbf{x}_i we can perform *maximum likelihood estimation*, i.e. search for parameters θ such that the observed data is most probable, by simply estimating the gradients of the *log-likelihood*

$$\log p_{\theta}(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) \quad (4)$$

Moreover, *maximum a posteriori estimation* of which maximum likelihood is a special case, can be performed and full Bayesian treatment of the posterior is also possible [1]. A complication arises when we include *latent* variables into the model, these are variables which are not observable and do not form a part of the dataset \mathcal{D} , usually denoted with \mathbf{z} . For unconditional model learning the graphical model now represents a joint probability distribution $p_{\theta}(\mathbf{x}, \mathbf{z})$, which when marginalized gives the single point *model evidence* $p_{\theta}(\mathbf{x})$ as a function of the parameters θ . In the simplest deep latent variable models, the joint probability is factored into the prior distribution $p_{\theta}(\mathbf{z})$ and the conditional distribution, which is represented by a neural network $p_{\theta}(\mathbf{x}|\mathbf{z})$, as

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x} | \mathbf{z}). \quad (5)$$

Alternatively, we can have a model with many layers of latent variables \mathbf{z}_m ,

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z}_1) p(\mathbf{z}_1 | \mathbf{z}_2) \cdots p(\mathbf{z}_{m-1} | \mathbf{z}_m) p(\mathbf{z}_m), \quad (6)$$

which are then able to learn hierarchical latent representation. In either case learning the posterior $p_{\theta}(\mathbf{x}, \mathbf{z})$, i.e.

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}, \quad (7)$$

is intractable, meaning it is (without tricks) impossible to differentiate the marginal likelihood w.r.t parameters θ and perform gradient-based optimisation. Advantage of DLVMs is their flexibility, even using very simple priors and conditionals they can express complicated distributions because of the nonlinearity of the neural network representations.

1.2.1. Learning in DLVMs

There are two main issues with learning in deep latent variable models, succinctly illustrated with two optimisation approaches. The first is the Expectation-Maximization (EM) [15] algorithm, it consists of two steps which are repeated until convergence.

- **E-step:** For each datapoint \mathbf{x}_i compute the posterior $p_{\theta_i}(\mathbf{z}_i|\mathbf{x}_i)$.

- **M-step:** Estimate new parameters by maximising the expected log-likelihood $\theta_{t+1} = \arg \max_{\theta} \sum_{x \in D} \mathbb{E}_{z \sim p(z|x; \theta_t)} \log p(x, z; \theta)$.

The algorithm falls apart in two ways, firstly the E-step requires computation of the posterior distribution which is intractable in DLVMs. This can be solved by employing a sampling technique (MCMC EM algorithm) but at a high computational cost, requiring a per datapoint MCMC step. Secondly, the M-step optimises over all datapoints requiring vast amounts of memory, although it can be done in mini-batches (online EM [16]). Hence, the first issue with learning in DLVMs is computational efficiency. The second algorithm is the *wake-sleep* algorithm [17], it employs a recognition model that approximates the posterior distribution and is applicable to both discrete and continuous latent variable models. While it is not computationally expensive relative to EM, its optimisation objectives do not coincide with optimising the marginal likelihood, thus giving a bad estimate of the posterior. Therefore the second issue is inaccurate estimates of the posterior. Both these issues were addressed in two seminal papers that constitute the framework of the variational autoencoder [2, 3].

1.3. Variational Autoencoders

The strength of variational methods lies in recasting inference as an optimisation problem, which scales better to large datasets than sampling methods. This is done by introducing an approximate posterior distribution $q_{\phi}(z|x)$, with parameters ϕ and optimising to maximise the similarity between $q_{\phi}(z|x)$ and $p_{\theta}(z|x)$. In the case of the VAE $q_{\phi}(z|x)$ is also parameterised using a neural network. The similarity between the two distributions is measured by the Kullback-Leibler divergence

$$D_{KL}(q||p) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}, \quad (8)$$

which has two important properties $D_{KL}(q||p) \geq 0$ for all q and p and it is equal to zero only when $q = p$. It can be rewritten as

$$\begin{aligned} D_{KL}(q_{\phi}(z|x)||p_{\theta}(z|x)) &= \\ &= \sum_z q_{\phi}(z|x) \left[\log \frac{q_{\phi}(z|x)}{p_{\theta}(z,x)} + \log p_{\theta}(x) \right] \\ &= \sum_z q_{\phi}(z|x) [\log q_{\phi}(z|x) - \log p_{\theta}(z,x)] + q_{\phi}(z|x) \log p_{\theta}(x), \end{aligned} \quad (9)$$

from which we can express the log likelihood as

$$\log p_{\theta}(x) = D_{KL}(q_{\phi}(z|x)||p_{\theta}(z|x)) + \mathcal{L}_{\theta, \phi}(x), \quad (10)$$

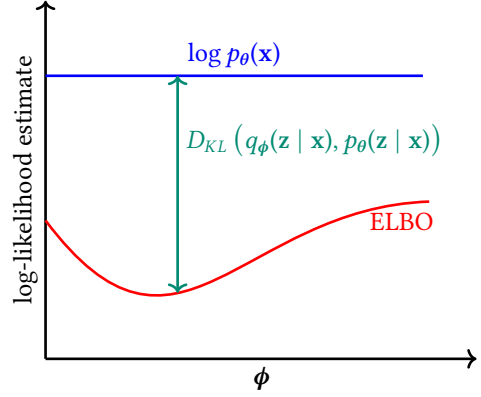


Figure 1: ELBO is a lower bound of the log likelihood, the difference between the two is the KL divergence. The better the approximation of the posterior $q_{\phi}(z|x)$, the smaller is KL and consequently the closer is ELBO to $\log p_{\theta}(x)$.

where $\mathcal{L}_{\theta, \phi}(x)$ is the *evidence lower bound* ELBO, defined as

$$\mathcal{L}_{\theta, \phi}(x) = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)]. \quad (11)$$

ELBO gives a lower bound on the log likelihood, because Kullback-Leibler divergence is non-negative. From equation (10) we see, that maximizing the ELBO with regards to both θ and ϕ does two things, it minimizes the KL divergence between the true and approximate posteriors, and consequently better approximates the log evidence, see Figure 1. This makes ELBO a good optimisation objective for learning in VAEs.

1.3.1. Gradient estimation

In order to perform gradient descent on the ELBO objective, we need to be able to evaluate its gradients with respect to parameters θ and ϕ . Taking the gradient w.r.t generative parameters θ is straightforward, because we can change the order of the expectation operator and the gradient, leaving us with

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\theta, \phi}(x) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] \\ &\approx \nabla_{\theta} (\log p_{\theta}(x, z) - \log q_{\phi}(z|x)) \\ &= \nabla_{\theta} (\log p_{\theta}(x, z)), \end{aligned} \quad (12)$$

where \approx denotes an unbiased estimator. This reversing of the order of operations is not possible when taking gradients w.r.t variational parameters ϕ , because the expectation $\mathbb{E}_{q_{\phi}(z|x)}$ is performed w.r.t the approximate posterior $q_{\phi}(z|x)$. The gradient could be estimated with a vanilla Monte Carlo estimator, but it has very high variance and is not practical [2].

The problem of stochastic gradient estimation of an expectation of a function is a well studied problem that transcends machine learning and has a variety of applications [18, 19]. Different estimators differ in from and their properties, variance being one of the most important. In their review [20] Mohamed et al. categorise MC gradient estimators into three categories

1. *Score-function estimator*: The score function is a logarithm of a probability distribution w.r.t to distributional parameters. It can be used as a gradient estimator

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(\mathbf{x})}[f(\mathbf{x})] &= \nabla_{\theta} \int p_{\theta}(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}_{p_{\theta}(\mathbf{x})} [f(\mathbf{x}) \nabla_{\theta} \log p_{\theta}(\mathbf{x})].\end{aligned}\quad (13)$$

The score-function estimator is compatible with any cost function, it requires that the measure $p_{\theta}(\mathbf{x})$ is differentiable and easy to sample. Very importantly it is applicable to both discrete and continuous distribution, but has a drawback of having high variance.

2. *Pathwise estimator*: Continuous distributions can be sampled either directly by generating samples from the distribution $p_{\theta}(\mathbf{x})$ or indirectly, by sampling from a simpler base distribution $p(\epsilon)$ and transforming the variate through a deterministic path $g_{\theta}(\epsilon)$. Using this, it is possible to move the source of randomness in such a way that the objective is differentiable. In essence this approach pushes the parameters of the measure into the cost function which is then differentiated. The estimator is

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(\mathbf{x})}[f(\mathbf{x})] &= \nabla_{\theta} \int p_{\theta}(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} \\ &= \nabla_{\theta} \int p(\epsilon) f(g_{\theta}(\epsilon)) d\epsilon \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_{\theta} f(g_{\theta}(\epsilon))].\end{aligned}\quad (14)$$

This was the gradient estimator originally used in the VAE implementation [2] there named as the *reparametrization trick*, see also Figure 2. In many cases the transformation paths are so simple they can be implemented in one line of code, referred to as *one-liners*. The pathwise-estimator can only be used on differentiable cost functions, but is easy to implement and crucially has lower variance than the score-function estimator.

3. *Measure-valued gradient estimator*, which exploits the properties of signed-measures is beyond the scope of this report.

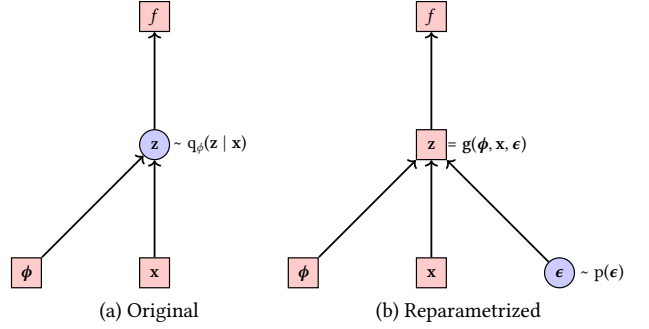


Figure 2: The reparametrization trick, adapted from [1]. The stochasticity of the z node is pushed out into a separate input to the same node, resulting in deterministic gradients w.r.t ϕ through the node.

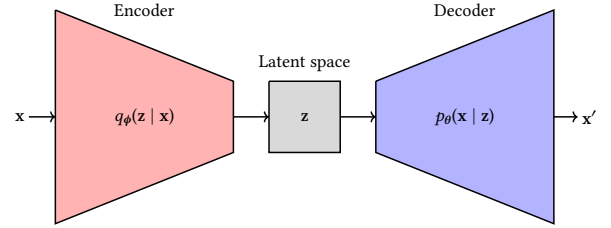


Figure 3: The variational autoencoder.

Unbiased gradient estimation makes training on the ELBO objective possible. During training the VAE learns a mapping between a complicated observed space of \mathbf{x} and the latent space of \mathbf{z} , which is usually defined to be relatively simple. The *decoder* network or the *generative* part of the VAE $p_{\theta}(\mathbf{x}|\mathbf{z})$ learns to map from the latent space to the data space, while the *inference* or *encoder* network $q_{\phi}(\mathbf{z}|\mathbf{x})$ approximates the true posterior $p_{\theta}(\mathbf{z}|\mathbf{x})$ of this process, see Figure 3.

1.3.2. Variants and advances

In the previous section we have described the basic variational autoencoder as proposed in the original paper [2]. But as with any new paradigm/approach in machine learning, the variational autoencoder was quickly extended in a variety of ways. A large line of work tries to address the issues that VAEs have with learning, this was prompted by the discovery that ELBO optimisation can get stuck in suboptimal equilibria [21]. Suggestions include annealing the latent cost $D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}))$ from 0 to 1 [21], or modifying the ELBO objective to ensure that a certain minimum number of bits is encoded in each latent variable [22]. The β -VAE [23] also modifies the ELBO objective by introducing the β hyperparameter to balance the regularization and reconstruction terms. One thing that the

β -VAE does particularly well is learn disentangled latent representations, which are latent representations where each component is sensitive to changes in one generative factor and is not much affected by others. Further work has been done in understanding how introducing the β parameter leads to these representations [24].

Alternative objectives to the evidence lower bound have also been proposed, instead of using the KL divergence one can use the chi-square divergence to arrive at an upper bound of the log likelihood (CUBO) and optimise that [25]. In addition to the objective of the VAE, any other component can be, and has been modified and studied. There has been work done on improving the variational posteriors by using conditional flow-based models [22, 5], implementing all sorts of architectures of the decoder, including convolutional neural networks, residual networks and recurrent neural networks [26]. Much in the same vein, variational autoencoders have been adapted for use with other data types than the original application to images. Among other uses, they have been used for text generation [27], molecular graph generation [28] and video compression [29]. Variational autoencoders can also be used in tandem with a classifier in a procedure called semi-supervised learning [13], allowing to learn from large datasets which have only partially been labeled.

We have alluded in section 1.2 that VAEs can have hierarchical latent space structure, such models have very recently been implemented and show promising results. The Bidirectional-Inference Variational Autoencoder (BIVA) [30] uses a skip-connected generative model and an inference network formed by a bidirectional stochastic inference path to learn a hierarchy of latent variables. Another hierarchical VAE is the Nouveau-VAE (NVAE) [31] which uses depth-wise separable convolutions and batch normalization. Moreover, very deep VAE's [32] can represent autoregressive models and beat them on likelihood benchmarks.

Given the reliance of most variational autoencoders on the reparametrization trick, one of their main flaws is their inability to treat discrete latent variables. Current workarounds to this issue include using a continuous relaxation of the discrete latent variables [33], or pairing each discrete latent variable with an auxiliary continuous one and performing the reparametrization trick on the marginal distributions [34]. Current state of the art is the vector quantised VAE (VQ-VAE) which uses vector quantisation for the latent space and learns the prior [35].

Nevertheless, there are some problems with VAEs that have not yet been fully addressed. Perhaps the

most striking issue, one that pertains to all deep generative models not just VAEs, is the *out-of-distribution problem*. Deep generative models have trouble distinguishing between a sample from the distribution they were trained on and another sample from a different distribution. This was discovered by Nalisnick et al. [36], when they noticed that the models could not distinguish common objects (CIFAR-10) from house numbers (SVHN). The second problem is the *posterior collapse*. Because of the variance of the stochastic gradient approximation the posterior approximation $q_\phi(\mathbf{z}|\mathbf{x})$ may provide too weak of a signal

$$q_\phi(\mathbf{z}|\mathbf{x}) \approx q_\phi(\mathbf{z}), \quad (15)$$

resulting in the decoder ignoring \mathbf{z} samples from the posterior. Proposed solutions include skip connections [25] and aggressively optimising the inference model before performing each model update. [37]. There is also the *hole problem* [38], which occurs when there are discrepancies between the prior $p(\mathbf{z})$ and the aggregate posterior $q_\phi(\mathbf{z}) = \frac{1}{N} \sum_n q_\phi(\mathbf{z}|\mathbf{x}_n)$, where \mathbf{x}_n is a training example. Sampling from these discrepancies or *holes* provides poor latent variable samples which are then badly reconstructed by the generative network. Finally there is the *balancing problem*, the issue of finding a good balance between the sampling quality and the reconstruction quality.

2. Implementation of Variance loss

The evidence lower bound (ELBO) pathwise gradient estimator (SGVB), described in section 1.3.1, has found great success but is limited to models with continuous latent variables. It would be more convenient to have a loss function that would deal with both discrete and continuous latent variables within the same framework. Recent work by Richter et al. [39] analysed a score-function gradient estimator for ELBO and showed that it can be obtained from a different divergence measure, which they termed the *log-variance loss*

$$\mathcal{L}_{LV}(q_\phi(\mathbf{z})\|p_\theta(\mathbf{z}|\mathbf{x})) = \frac{1}{2} \mathbb{V}_{\mathbf{z} \sim r(\mathbf{z})} \left[\log \left(\frac{q_\phi(\mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \right) \right], \quad (16)$$

where $r(\mathbf{z})$ is some arbitrary distribution. Moreover, they showed that the gradient of \mathcal{L}_{LV} w.r.t ϕ corresponds to that of the KL divergence if $r(\mathbf{z}) = q_\phi(\mathbf{z})$ and provided a gradient estimator of the log-variance, thus suggested an optimisation procedure that works for both discrete and continuous latent distributions. In this work we will proceed in a similar fashion but

with a different loss function. Motivated by the fact that when the posterior $q(z|x) = p(z|x)$ is exact we have

$$\frac{p(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{q(\mathbf{z} | \mathbf{x})} = p(\mathbf{x}). \quad (17)$$

The variance of $p(\mathbf{x})$ with regards to samples drawn from an arbitrary probability distribution $r(\mathbf{z})$ which suggests the variance of the left hand side

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) = \mathbb{V}_{\mathbf{z} \sim r(\mathbf{z})} \left[\log \left(\frac{p_{\theta}(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right) \right] \quad (18)$$

as an objective function. We can estimate the loss with sample variance using S Monte Carlo samples

$$\mathcal{L}_{\phi, \theta}(\mathbf{x}) \approx \frac{1}{(S-1)} \sum_{s=1}^S \left(f_{\phi, \theta}(\mathbf{z}^{(s)}, \mathbf{x}) - \bar{f}_{\phi, \theta}(\mathbf{x}) \right)^2, \quad (19)$$

$\mathbf{z}^{(s)} \stackrel{\text{i.i.d.}}{\sim} r(\mathbf{z}),$

where

$$f_{\phi, \theta} = \log \left(\frac{p_{\theta}(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right) \quad (20)$$

and $\bar{f}_{\phi, \theta}$ is the average

$$\bar{f}_{\phi, \theta}(\mathbf{x}) = \frac{1}{S} \sum_{s=1}^S f_{\phi, \theta}(\mathbf{z}^{(s)}, \mathbf{x}). \quad (21)$$

The gradient of expression (19) can be computed using automatic differentiation without any issue. The form of the gradient estimator of $\mathcal{L}_{\phi, \theta}(\mathbf{x})$ outlined above, is the same as the leave-one-out estimator described in [39], which means that it is an unbiased estimator of the gradient. The procedure of calculating the gradients for a given example \mathbf{x} and parameters ϕ, θ is given in Algorithm 1. Importantly the samples $\mathbf{z}^{(s)}$, which will be for simplicity sampled from the model probability $\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})$, must be detached from the computational graph so that backpropagation is possible.

To test the applicability of the loss function we will test two variational autoencoders, one with continuous latent space and one with discrete latent space. The architecture of the neural networks that were used to parameterise the encoder and decoder will be discussed in section 3.

2.1. Continuous latent space

For the continuous case we will recreate the model from the original VAE paper [2]. The inference and generative models will both be parameterised by fully connected multilayer perceptrons. The approximate

Algorithm 1: Gradients of $\mathcal{L}_{\phi, \theta}(\mathbf{x})$

Input: Encoder parameters ϕ , decoder parameters θ , example \mathbf{x}

Result: Gradients w.r.t ϕ, θ of the objective function 1

for $s = 1, \dots, S$ **do**

$\mathbf{z}^{(s)} \leftarrow \text{sample}(q_{\phi}(\cdot | \mathbf{x}))$

$\mathbf{z}^{(s)} \leftarrow \text{stop_gradient}(\mathbf{z}^{(s)})$

$f_{\phi, \theta}^{(s)} \leftarrow \log p_{\theta}(\mathbf{x} | \mathbf{z}^{(s)}) p(\mathbf{z}^{(s)}) - \log q_{\phi}(\mathbf{z}^{(s)})$

end

$\hat{\mathcal{L}} \leftarrow \text{Variance} \left(\left\{ f_{\phi, \theta}^{(s)} \right\}_{s=1}^S \right)$

return $\text{grad}(\hat{\mathcal{L}})$

posterior $q_{\phi}(\mathbf{z} | \mathbf{x})$ will be a multivariate Gaussian with a diagonal covariance matrix

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\sigma^2))$$

$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \prod_i q_{\phi}(z_i | \mathbf{x}) = \prod_i \mathcal{N}(z_i; \mu_i, \sigma_i^2) \quad (22)$$

where

$$(\boldsymbol{\mu}, \log \sigma) = \text{EncoderNeuralNet}_{\phi}(\mathbf{x}). \quad (23)$$

The choice of decoder depends on the type of data. For real-valued outputs the generative $p_{\theta}(\mathbf{z} | \mathbf{x})$ model is again a Gaussian, with mean and variance parameterised by a neural network

$$(\boldsymbol{\mu}, \log \sigma) = \text{DecoderNeuralNet}_{\theta}(\mathbf{z}). \quad (24)$$

Alternatively if the outputs are binary we use the Bernoulli distribution for the logarithm of $p_{\theta}(\mathbf{x} | \mathbf{z})$

$$\begin{aligned} \log p_{\theta}(\mathbf{x} | \mathbf{z}) &= \sum_{j=1}^D \log p(x_j | z_j) \\ &= \sum_{j=1}^D \log \text{Bernoulli}(x_j; y_j) \\ &= \sum_{j=1}^D x_j \log y_j + (1 - x_j) \log (1 - y_j). \end{aligned} \quad (25)$$

Where D is the dimension of \mathbf{x} and \mathbf{y} is the output of a sigmoid activation in the last layer of the NN

$$\mathbf{y} = f_{\sigma}(\text{DecoderNeuralNet}_{\theta}(\mathbf{z})). \quad (26)$$

The prior $p(\mathbf{z})$ will be a unit isotropic Gaussian in either case $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.

2.2. Discrete latent space

3. Experiments

3.1. Continuous latent space

We compare the performance of the latent variable model described in section 2.1 and compare the efficiency of different gradient estimators/objectives.

3.1.1. Fashion MNIST

3.1.2. CELEBa

3.2. Discrete latent space

3.2.1. Binarized MNIST

3.2.2. Binarized Omniglot

4. Conclusion

Bibliography

- [1] D. P. Kingma, Variational inference & deep learning: A new synthesis (2017).
- [2] D. P. Kingma, M. Welling, Auto-encoding variational bayes, arXiv preprint arXiv:1312.6114 (2013).
- [3] D. J. Rezende, S. Mohamed, D. Wierstra, Stochastic backpropagation and approximate inference in deep generative models, in: International conference on machine learning, PMLR, 2014, pp. 1278–1286.
- [4] I. Goodfellow, Nips 2016 tutorial: Generative adversarial networks, arXiv preprint arXiv:1701.00160 (2016).
- [5] D. Rezende, S. Mohamed, Variational inference with normalizing flows, in: International Conference on Machine Learning, PMLR, 2015, pp. 1530–1538.
- [6] B. S. Everitt, Finite mixture distributions, Wiley StatsRef: Statistics Reference Online (2014).
- [7] L. Rabiner, B. Juang, An introduction to hidden markov models, *IEEE ASSP Magazine* 3 (1) (1986) 4–16.
- [8] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent dirichlet allocation, the *Journal of machine Learning research* 3 (2003) 993–1022.
- [9] D. H. Ackley, G. E. Hinton, T. J. Sejnowski, A learning algorithm for boltzmann machines, *Cognitive science* 9 (1) (1985) 147–169.
- [10] P. Dayan, G. E. Hinton, R. M. Neal, R. S. Zemel, The helmholtz machine, *Neural computation* 7 (5) (1995) 889–904.
- [11] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, Wavenet: A generative model for raw audio, arXiv preprint arXiv:1609.03499 (2016).
- [12] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al., Photo-realistic single image super-resolution using a generative adversarial network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4681–4690.
- [13] D. P. Kingma, D. J. Rezende, S. Mohamed, M. Welling, Semi-supervised learning with deep generative models, arXiv preprint arXiv:1406.5298 (2014).
- [14] I. Goodfellow, Y. Bengio, A. Courville, Y. Bengio, Deep learning, Vol. 1, MIT press Cambridge, 2016.
- [15] A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from incomplete data via the em algorithm, *Journal of the Royal Statistical Society: Series B (Methodological)* 39 (1) (1977) 1–22.
- [16] M. Hoffman, F. R. Bach, D. M. Blei, Online learning for latent dirichlet allocation, in: advances in neural information processing systems, Citeseer, 2010, pp. 856–864.
- [17] G. E. Hinton, P. Dayan, B. J. Frey, R. M. Neal, The “wake-sleep” algorithm for unsupervised neural networks, *Science* 268 (5214) (1995) 1158–1161.
- [18] N. A. Chriss, N. Chriss, Black Scholes and beyond: option pricing models, McGraw-Hill, 1997.
- [19] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, et al., Mastering atari, go, chess and shogi by planning with a learned model, *Nature* 588 (7839) (2020) 604–609.
- [20] S. Mohamed, M. Rosca, M. Figurnov, A. Mnih, Monte carlo gradient estimation in machine learning, *Journal of Machine Learning Research* 21 (132) (2020) 1–62.
- [21] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, S. Bengio, Generating sentences from a continuous space, arXiv preprint arXiv:1511.06349 (2015).
- [22] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, M. Welling, Improving variational inference with inverse autoregressive flow, arXiv preprint arXiv:1606.04934 (2016).
- [23] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, A. Lerchner, beta-vae: Learning basic visual concepts with a constrained variational framework (2016).
- [24] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, A. Lerchner, Understanding disentangling in β -vae, arXiv preprint arXiv:1804.03599 (2018).
- [25] A. B. Dieng, D. Tran, R. Ranganath, J. Paisley, D. M. Blei, Variational inference via χ -upper bound minimization, arXiv preprint arXiv:1611.00328 (2016).
- [26] K. Gregor, I. Danihelka, A. Graves, D. Rezende, D. Wierstra, Draw: A recurrent neural network for image generation, in: International Conference on Machine Learning, PMLR, 2015, pp. 1462–1471.
- [27] S. Semeniuta, A. Severyn, E. Barth, A hybrid convolutional variational autoencoder for text generation, arXiv preprint arXiv:1702.02390 (2017).
- [28] W. Jin, R. Barzilay, T. Jaakkola, Junction tree variational autoencoder for molecular graph generation, in: International Conference on Machine Learning, PMLR, 2018, pp. 2323–2332.
- [29] A. Habibi, T. v. Rozendaal, J. M. Tomczak, T. S. Cohen, Video compression with rate-distortion autoencoders, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 7033–7042.
- [30] L. Maaløe, M. Fraccaro, V. Liévin, O. Winther, Biva: A very deep hierarchy of latent variables for generative modeling, arXiv preprint arXiv:1902.02102 (2019).
- [31] A. Vahdat, J. Kautz, Nvae: A deep hierarchical variational autoencoder, arXiv preprint arXiv:2007.03898 (2020).
- [32] R. Child, Very deep vae’s generalize autoregressive models and can outperform them on images, arXiv preprint arXiv:2011.10650 (2020).
- [33] C. J. Maddison, A. Mnih, Y. W. Teh, The concrete distribution: A continuous relaxation of discrete random variables, *CoRR abs/1611.00712* (2016). arXiv:1611.00712.
- [34] J. T. Rolfe, Discrete variational autoencoders, arXiv preprint arXiv:1609.02200 (2016).
- [35] A. van den Oord, O. Vinyals, K. Kavukcuoglu, Neural discrete representation learning, *CoRR abs/1711.00937* (2017). arXiv:1711.00937.
- [36] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, B. Lakshminarayanan, Do deep generative models know what they don’t know?, arXiv preprint arXiv:1810.09136 (2018).
- [37] J. He, D. Spokoyny, G. Neubig, T. Berg-Kirkpatrick, Lagging inference networks and posterior collapse in variational autoencoders, arXiv preprint arXiv:1901.05534 (2019).
- [38] D. J. Rezende, F. Viola, Taming vae’s, arXiv preprint

arXiv:1810.00597 (2018).

- [39] L. Richter, A. Boustati, N. Nüsken, F. J. Ruiz, Ö. D. Akyildiz, Vargrad: A low-variance gradient estimator for variational inference, arXiv preprint arXiv:2010.10436 (2020).