

# $\beta$ -VAE

Benjamin Levy

December 10, 2020  
STAT185

## 1 Introduction

A  $\beta$ -variational autoencoder ( $\beta$ -VAE) is a family of unsupervised models whose goal is to infer a “disentangled,” low-dimensional representation of some high-dimensional data, often images [8]. In general, auto-encoders seek to reconstruct some input  $\vec{x} \in \mathbb{R}^d$  while simultaneously learning a latent representation,  $\vec{z} \in \mathbb{R}^k$ , which is typically chosen to be lower-dimensional than  $\vec{x}$ . This is accomplished by training two models: an encoder parameterized by  $\phi$ , denoted  $g_\phi(\vec{x})$ , and a decoder parameterized by  $\theta$ , denoted  $f_\theta(\vec{z})$ . The algorithm is, roughly:

1. Given example  $\vec{x}$ , compute  $\vec{z} = g_\phi(\vec{x})$
2. Compute reconstruction  $\hat{x} = f_\theta(\vec{z})$

A limitation of this deterministic autoencoder is that the latent space in which  $\vec{z}$  resides is typically quite sparse and/or irregular. This can make it difficult to use the latent representations for inference or other purposes. Variational auto-encoders (VAEs) improve upon this by replacing the decoder  $f_\theta$  with a stochastic sampling function. Consequently,  $g_\phi$  now outputs a mean  $\vec{\mu}$  and a covariance  $\Sigma$ , which parameterize a Gaussian distribution [3]. The algorithm then becomes:

1. Given example  $\vec{x}$ , compute  $(\vec{\mu}, \Sigma) = g_\phi(\vec{x})$
2. Sample latent variable  $\vec{z} \sim \mathcal{N}(\vec{\mu}, \Sigma)$
3. Compute reconstruction  $\hat{x} = f_\theta(\vec{z})$

To make this model tractable, we typically use isotropic Gaussians, meaning the covariance matrix  $\Sigma$  is diagonal. By replacing the deterministic algorithm with a stochastic one, the model is forced to learn latent representations that are more robust to Gaussian noise, since if small perturbations led to dramatically different outputs, the model would perform poorly at reconstructing the output. The loss function for the VAE can be expressed as a sum of two terms: (1) reconstruction error and (2) regularization (defined more rigorously later on in the methods section). The reconstruction term pushes the model to faithfully reconstruct the original data, while the regularization term encourages the model to learn a generalizable and efficient latent representation.

While the VAE model solves the problem of irregular latent spaces, there are no guarantees that the latent representations are *meaningful*. This gives rise to the idea of *disentanglement*, which is roughly the notion that the different components of the latent variable  $\vec{z}$  should correspond to distinct *generative factors* (i.e. the parameters of the generative process) [2]. For instance, if we are modelling an image dataset consisting of human faces, we might want one dimension to correspond to skin tone, another to correspond to eye shape, a further one to correspond to amount of hair, and so on. Introduced by Burgess et al. in 2018, the  $\beta$ -VAE is a model designed with this goal in mind [4]. An additional hyperparameter,  $\beta \in \mathbb{R}$  is added as a coefficient for the regularization term from the VAE loss. When  $\beta = 1$ , the loss function is identical to the original VAE loss function. When  $\beta > 1$ , the model is forced to learn a less expressive, and therefore more information-dense, latent encoding for the data. Furthermore, the representations learned by  $\beta$ -VAE

tend to be more disentangled, meaning that it is more possible to recover the original factors that may have generated the data.

In the landscape of dimensionality reduction techniques,  $\beta$ -VAE (and all types of VAE) are non-linear. The data domain (images) is such that the data are extremely high dimensional but the dimensions have highly complex relationships with one another, both geometrically and semantically. It is for this reason that simpler, more linear dimensionality reduction techniques, like PCA, or manifold-based techniques, like ISOMAP, are less likely to be able to perform well on image data, while VAE's can. However, this additional flexibility comes at the cost of reduced interpretability, since we no longer have any guarantees about the geometric relationship between the data space and latent space. This is where  $\beta$ -VAEs come in: by reducing the flexibility of the latent space, we can hopefully gain more interpretability in the low-dimensional representations.

## 2 Methods

In this section we first build up intuition and theory for vanilla (standard) VAEs before extending the math to  $\beta$ -VAE.

### 2.1 Variational auto-encoders

We have data  $\{\vec{x}_i\}_{i=1}^N$  and we wish to find some low-dimensional (latent) representation for each  $i$ :  $\{\vec{z}_i\}_{i=1}^N$ . We assume that the data are in fact generated stochastically conditioned on the latent representations, such that:

$$\vec{x}|\vec{z} \sim p(\vec{x}|\vec{z})$$

$p$  is commonly either a Gaussian distribution (in the case of normally distributed data) or a Bernoulli distribution (in the case of pixels). The goal of a VAE is to learn a good approximation to this distribution,  $p_\theta$ , with parameters  $\theta$ , to model our data  $\vec{x}$ , as well as the distribution  $p_\theta(\vec{z}|\vec{x})$  to model the hypothesized latent variables  $\vec{z}$  conditioned on the observed data. We can arrange related  $\vec{x}$  and  $\vec{z}$  using Bayes' rule:

$$p_\theta(\vec{x}|\vec{z}) = \frac{p_\theta(\vec{z}|\vec{x})p_\theta(\vec{x})}{p_\theta(\vec{z})}$$

Under this framework, we want to maximize the log-probability of data generated according to the process:  
(1) sample a point  $\vec{z} \sim p_\theta(\vec{z})$ , then sample  $\vec{x} \sim p_\theta(\vec{x}|\vec{z})$ :

$$\ell = \sum_{i=1}^N \log p_\theta(\vec{x}_i)$$

Marginalizing over all values of  $\vec{x}$ ,

$$\ell = \sum_{i=1}^N \log \int p_\theta(\vec{x}_i|\vec{z})p_\theta(\vec{z})d\vec{z}$$

We now introduce an auxiliary distribution,  $q_\phi$ . This will allow us to massage this expression into one involving a KL-divergence and will ultimately be the distribution that we learn using the model.

$$\begin{aligned} \ell &= \sum_{i=1}^N \log \int p_\theta(\vec{x}_i|\vec{z}_i)p_\theta(\vec{z}_i) \frac{q_\phi(\vec{z}_i)}{q_\phi(\vec{z}_i)} d\vec{z}_i \\ &= \sum_{i=1}^N \log \int \frac{p_\theta(\vec{x}_i, \vec{z}_i)}{q_\phi(\vec{z}_i)} q_\phi(\vec{z}_i) d\vec{z}_i \\ &= \sum_{i=1}^N \log \mathbb{E} \left[ \frac{p_\theta(\vec{x}_i, \vec{z}_i)}{q_\phi(\vec{z}_i)} \right] \end{aligned}$$

Using Jensen's inequality, we get the evidence lower bound (ELBO)

$$\begin{aligned} &\geq \sum_{i=1}^N \mathbb{E}_{\vec{z} \sim q_\phi} \log \left[ \frac{p_\theta(\vec{x}_i, \vec{z}_i)}{q_\phi(\vec{z}_i)} \right] \\ \ell &= \sum_{i=1}^N (\mathbb{E}_{\vec{z} \sim q_\phi} [\log p_\theta(\vec{x}_i, \vec{z}_i)] - \mathbb{E}_{\vec{z} \sim q_\phi} [\log q_\phi(\vec{z})]) \end{aligned}$$

Recall that the definition of the KL divergence between  $q_\phi$  and  $p_\theta$  is:

$$D_{KL}(q_\phi(\vec{z}) || p_\theta(\vec{z}|\vec{x})) = \mathbb{E}_{\vec{z} \sim q_\phi} \left[ \log \frac{q_\phi(\vec{z})}{p_\theta(\vec{z}|\vec{x})} \right]$$

Since  $p_\theta(\vec{z}|\vec{x}) = \frac{p_\theta(\vec{z}, \vec{x})}{p_\theta(\vec{x})}$  (by definition of conditional probability),

$$D_{KL}(q_\phi(\vec{z}) || p_\theta(\vec{z}|\vec{x})) = -(\mathbb{E}_{\vec{z} \sim q_\phi} [\log p_\theta(\vec{z}, \vec{x})] - \mathbb{E}_{\vec{z} \sim q_\phi} [\log q_\phi(\vec{z})]) + \log p_\theta(\vec{x})$$

Plugging this back into our log-likelihood,

$$\ell = \sum_{i=1}^N (\log(p_\theta(\vec{x}_i)) - D_{KL}(q_\phi(\vec{z}_i) || p_\theta(\vec{z}_i|\vec{x}_i)))$$

The first term is maximized when the reconstruction  $\vec{x}$  is faithful to the original data, whereas the second term is minimized when  $p_\theta$  is close to  $q_\phi$ .

## 2.2 Choice of $q_\phi$

In theory, we could choose an arbitrarily expressive family of distributions  $q_\phi$  to model the latent variable  $\vec{z}$ . In practice, we make the simplifying **mean-field assumption**, which states that the dimensions of  $q_\phi$  are independent:

$$q_\phi(\vec{z}|\vec{x}) = \prod_{j=1}^d q_\phi(z_j|\vec{x}_j)$$

In our case, we operationalize this by saying that  $q_\phi$  is an isotropic Gaussian (diagonal covariance) [1]. Although this assumption reduces the expressiveness of the variational family  $q_\phi$ , we gain significantly in terms of computational efficiency when it comes time to learn the parameters of the distribution.

At this point, we are in effect trying to learn a unique distribution  $q_\phi(\vec{x}|\vec{z})$  for each  $\vec{x}$ . If we have  $N$  observations and a latent space of dimension  $d$ , then if the latent distribution is Gaussian, we would need to learn a mean and variance for each dimension of each latent variable, or  $2 \times N \times d$ . In other words,  $N \gg p$  (where  $p$  is the number of parameters in the model), which both dramatically increases the variance of the parameter estimates (they are highly dependent on the particular data sample at hand) and increases the cost of fitting the model. To get around this, we perform *amortised* variational inference: instead of directly estimating  $\mu_i$  and  $\Sigma_i$  for each  $i = 1, \dots, N$ , we assume that each  $\mu_i, \Sigma_i$  can be approximated by a function  $g_\phi(\vec{x})$ , typically a neural network [5]. We refer to this function interchangeably as the **encoder** or **inference network**. This is represented by the left rectangle in figure 1 and is denoted  $g_\phi$  here. Because the weights  $\phi$  are shared across all observations, we say that the cost of learning  $\phi$  is amortised across the entire dataset (hence, “amortised variational inference”). Symmetrically, the conditional posterior distribution  $p_\theta(\vec{x}|\vec{z})$  is implemented as a neural network that is typically the mirror image of the encoder, which is called the **decoder** or **generative network** and is denoted  $f_\theta$  here.

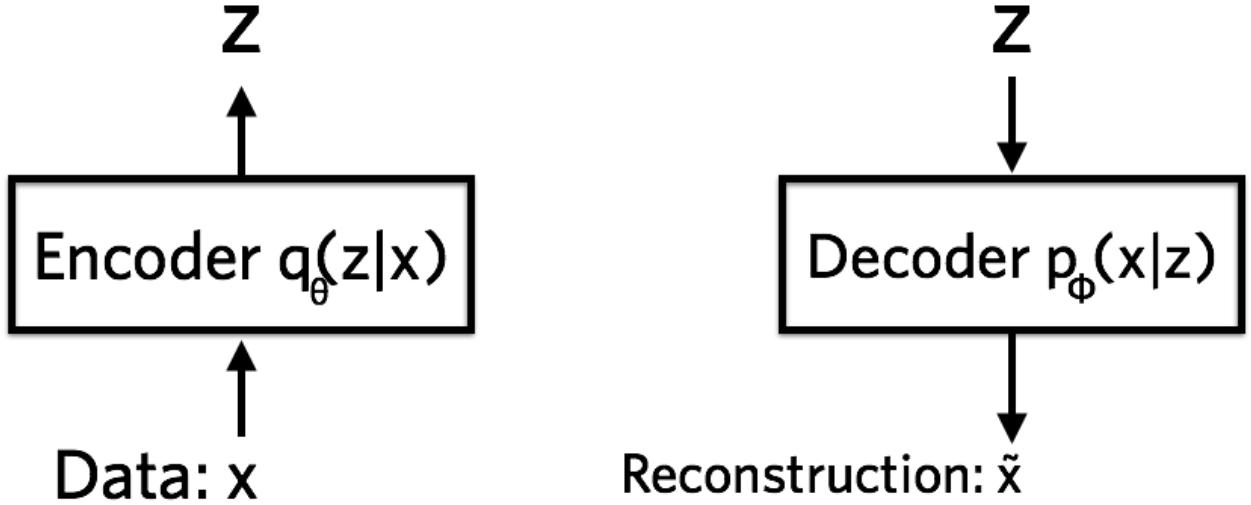


Figure 1: Simplified schematic of the encoder and decoder, where the left and right rectangles labelled  $q(z|x)$  and  $p(x|z)$  refer to the inferential and generative networks, respectively. First,  $\vec{z}$  is generated from  $\vec{x}$  using the encoder and then  $\vec{x}$  is reconstructed using  $\vec{z}$  [5]. Note: the notation in this figure is opposite to what we use in this paper:  $\theta$  and  $\phi$  have been flipped such that  $\theta$  parameterizes  $q$  and  $\phi$  parameterizes  $p$ . Otherwise, the setup is the same to what we have been discussing.

### 2.3 The reparameterization trick

A further advantage of variational inference is that it allows us to convert this problem from one of sampling from a posterior distribution – which can require the use of computationally intensive or mathematically complex methods such as expectation-maximization – to a simple non-convex optimization problem. That is, all the terms in the ELBO loss function can be expressed in terms of  $\phi$  and  $\theta$ , meaning that if we can differentiate through all the operations involved, then this problem can be solved with gradient descent. To revisit, the encoder and decoder form a single algorithm:

1. Given  $\vec{x}$ , compute  $\vec{\mu}, \Sigma = g_\phi(\vec{x})$
2. Sample  $\vec{z} \sim \mathcal{N}(\vec{\mu}, \Sigma)$
3. Compute reconstruction  $\hat{x} = f_\theta(\vec{z})$

The one operation that does not allow for easy differentiation is the sampling step (2). In order to get around this, we use the “reparameterization trick”. Recall that a  $d$ -dimensional isotropic Gaussian random variable can be converted into a sample from the  $d$ -dimensional multivariate normal distribution:

$$\vec{\epsilon} = (\vec{z} - \vec{\mu}) \odot \frac{1}{\vec{\sigma}} + \vec{\mu}$$

Where  $\vec{\epsilon}$  is a standard normal random variable (not to be confused with  $\vec{z}$ , which is our latent variable with mean  $\mu$  and covariance  $\Sigma$ ) and  $\vec{\sigma}$  is the  $d$ -dimensional vector of standard deviations for the components of  $\vec{z}$ .  $\odot$  represents element-wise multiplication. Using simple algebra, we can rearrange this:

$$\vec{z} = (\vec{\epsilon} \odot \vec{\sigma}) + \vec{\mu}$$

This means that sampling  $\vec{z} \sim \mathcal{N}(\vec{\mu}, \Sigma)$  is equivalent to sampling  $\vec{\epsilon} \sim \mathcal{N}(0, I)$  and then transforming it into  $\vec{z}$  as above. The gradient operation is well-defined for sampling from the standard normal, so we can make the entire VAE algorithm differentiable by adding this reparameterization step:

1. Given  $\vec{x}$ , compute  $\vec{\mu}, \Sigma = g_\phi(\vec{x})$ . Let  $\vec{\sigma} = \sqrt{\text{diag}(\Sigma)}$

2. Sample  $\vec{\epsilon} \sim \mathcal{N}(0, I)$
3. Let  $\vec{z} = \vec{\epsilon} \odot \vec{\sigma} + \vec{\mu}$
4. Compute reconstruction  $\hat{x} = f_{\theta}(\vec{z})$

## 2.4 $\beta$ -VAE

In the original paper, the authors provide an independent derivation of the  $\beta$ -VAE loss using constrained optimization [4]. Given data  $\vec{x}$  distributed according to the data distribution  $\mathcal{D}$ , the goal is to jointly learn generative parameters  $\theta$  and latent parameters  $\phi$  for the distributions  $p_{\theta}$  and  $q_{\phi}$ , respectively. We want to find  $\theta$  and  $\phi$  such that the posterior log-probability  $p_{\theta}(\vec{x}|\vec{z})$  where  $\vec{z} \sim q_{\phi}$  is maximized:

$$\theta^*, \phi^* = \underset{\theta, \phi}{\operatorname{argmax}} \mathbb{E}_{\vec{x} \sim \mathcal{D}} [\mathbb{E}_{\vec{z} \sim q_{\phi}} [\log p_{\theta}(\vec{x}|\vec{z})]]$$

Subject to constraint

$$D_{KL}(q_{\phi}(\vec{z}|\vec{x})||p(\vec{z})) < \delta$$

Where  $\delta$  is a small positive constant (since KL divergence is nonnegative). The constraint is imposed to narrow the "information bottleneck" by encouraging the latent posterior  $q_{\phi}(\vec{z}|\vec{x})$  to be close to a simple prior  $p(\vec{z})$ , which in this case is a standard gaussian,  $\mathcal{N}(0, I)$ . The idea is that we want to make  $q_{\phi}$  less expressive, which in turn forces it to be as efficient as possible in which latent factors are represented. Hopefully, the efficiency is maximized when the latent factors line up with the original factors that generated the data.

We can rewrite this constrained optimization problem as a Lagrangian subject to the KKT condition, with Lagrange multiplier  $\beta$  and Lagrangian objective  $\mathcal{F}$ :

$$\begin{aligned} \mathcal{F}(\theta, \phi, \beta; \vec{x}, \vec{z}) &= \mathbb{E}_{\vec{z} \sim q_{\phi}} [\log p_{\theta}(\vec{x}|\vec{z})] - \beta (D_{KL}(q_{\phi}(\vec{z}|\vec{x})||p(\vec{z})) - \delta) \\ &= \mathbb{E}_{\vec{z} \sim q_{\phi}} [\log p_{\theta}(\vec{x}|\vec{z})] - \beta (D_{KL}(q_{\phi}(\vec{z}|\vec{x})||p(\vec{z}))) + \beta \delta \\ &\geq \mathbb{E}_{\vec{z} \sim q_{\phi}} [\log p_{\theta}(\vec{x}|\vec{z})] - \beta (D_{KL}(q_{\phi}(\vec{z}|\vec{x})||p(\vec{z}))) \end{aligned}$$

Immediately, we can recognize that this is nearly identical to the vanilla VAE loss, except for the hyper-parameter  $\beta$ . When  $\beta = 1$ , this is equivalent to the vanilla VAE loss. When  $\beta = 0$ , the model becomes similar to a non-variational auto-encoder, since the only objective is to faithfully reconstruct the data. As  $\beta$  increases, the bottleneck in the latent space becomes narrower, promoting a sparser latent representation (i.e. a latent space that more closely resembles a standard Gaussian).

### 2.4.1 Measuring disentanglement

The authors of the original  $\beta$ -VAE paper go a step further and propose a "disentanglement metric" to quantify how well the model's latent representations align with the hypothesized generative factors for the data [4].

Formally, we assume the data  $\vec{x}$  are drawn from a dataset  $\mathcal{D} = \{X, V\}$ , where  $X$  is the data matrix, and  $V$  is a matrix of hypothesized *generative factors* for each data point. For instance, a factor for a dataset of faces might be *skin tone* or *glasses*. An image is created using a simulator:  $\vec{x} \sim \text{Sim}(\vec{v})$ . Then,

1. Uniformly sample factor  $y$  from the set of  $K$  possible generative factors
2. Repeat  $L$  times:
  - (a) Sample two sets of generative latent representations,  $\vec{v}_{1,l}$  and  $\vec{v}_{2,l}$  such that both  $\vec{v}_{1,l}$  and  $\vec{v}_{2,l}$  have the same value for factor  $y$  and are free to differ in all other generative factors
  - (b) Simulate an image  $x_{1,l} \sim \text{Sim}(v_{1,l})$  and compute the latent representation  $\vec{z}_{1,l} \sim q(\vec{z}|\vec{x}_{1,l})$  using the encoder  $g_{\phi}(\vec{x}_{1,l})$ . Repeat for  $v_{1,l}$

- (c) Compute a difference vector  $z_{\text{diff}}^l = |z_{1,l} - z_{2,l}|_1$
- 3. Given all  $z_{\text{diff}}^l$  computed as above, calculate the average difference vector  $\bar{z}_{\text{diff}} = \frac{1}{L} \sum_{l=1}^L z_{\text{diff}}^l$
- 4. Use  $\bar{z}_{\text{diff}}$  as the input for a linear classifier (i.e. logistic regression), where the goal is to predict the factor  $y$  that was kept constant between all pairs of latent factors

This algorithm can be seen in figure 2. The intuition behind this metric is that when a pair of data points  $(x_{1,l}, x_{2,l})$  has the same value for some generative factor, then a well-disentangled latent representation will be such that a small subset of the latent dimensions will correspond to that factor and will therefore have small difference, averaged over many such pairs. Thus, a linear classifier trained on supervised pairs of  $(\bar{z}_{\text{diff}}, y)$  will be able to map the averaged difference vectors to the corresponding generative factor that was held constant across the pairs.

## 2.5 Implementation

Depending on the type and complexity of the data, various architectures for the inference and generative networks are used. In the simplest case, each is a feed-forward neural network, which is common for either numeric data or small-sized image data (e.g. MNIST, which is  $28 \times 28 \times 1$ ). For larger image data, such as the CelebA dataset ( $64 \times 64 \times 3$ ), convolutional neural networks (CNNs) are often used, with the structure of the encoder neural network mirroring the structure of the decoder neural network. The weights of the encoder and decoder are found using well-known gradient descent algorithms (e.g. the authors use Adam, RMSProp, and Adagrad for different experiments).

## 3 Discussion

### 3.1 Strengths

The authors show that  $\beta$ -VAE is able to outperform other comparable deep methods of nonlinear dimensionality reduction in terms of learning a latent space that aligns with interpretable factors. They compare  $\beta$ -VAE to three other methods qualitatively: DC-IGN, InfoGAN, and a vanilla VAE (see fig. 3).

On its own,  $\beta$ -VAE can lead to some interesting results. For instance, the authors find that the model can learn effective latent representations for (a) skin colour, (b) age/gender, and (c) image saturation on the CelebA dataset (fig. 4).

### 3.2 Disentanglement and information bottleneck

Burgess et al. [7] take an information-theoretic approach to understanding how  $\beta$ -VAE improves disentanglement in the latent representations  $\vec{z}$ . Both the prior  $p(\vec{z})$  and the posterior  $q_\phi(\vec{z}|\vec{x})$  are factorized, as per the mean-field assumption, so we can conceptualize  $q(\vec{z}|\vec{x})$  as a set of noisy channels, each transmitting some information about the input  $\vec{x}_i$ . The KL-divergence term in the loss,  $D_{KL}(q_\phi(\vec{z}|\vec{x})||p(\vec{z}))$  can be seen as a bound on the information contained in  $q_\phi$ .  $D_{KL}$  will be 0 only if  $\vec{\mu}_i = 0, \Sigma_i = I$ , meaning that  $q_\phi$  is simply white noise and contains no useful information about  $\vec{x}_i$ . As  $D_{KL}$  increases,  $q_\phi$  can deviate further from  $\mathcal{N}(0, I)$  and contain more information that is particular to  $\vec{x}_i$ . As can be seen in figure 5, when  $\beta$  is increased, the sparsity of the latent encoding is greater, since only two out of the six possible latent dimensions differ noticeably from the prior, while the reconstructions are indistinguishable from the vanilla VAE reconstructions.

### 3.3 Dimensionality reduction perspective

We can also analyze  $\beta$ -VAE more explicitly from the perspective of dimensionality reduction. Just like a vanilla VAE,  $\beta$ -VAE is a nonlinear dimensionality reduction technique. More specifically, because the relationship between the latent probability space  $q_\phi(\vec{z}|\vec{x})$  and the data space  $p_\theta(\vec{x}|\vec{z})$  is parameterised by neural networks, there is no easily interpretable geometric similarity between the two spaces. This separates

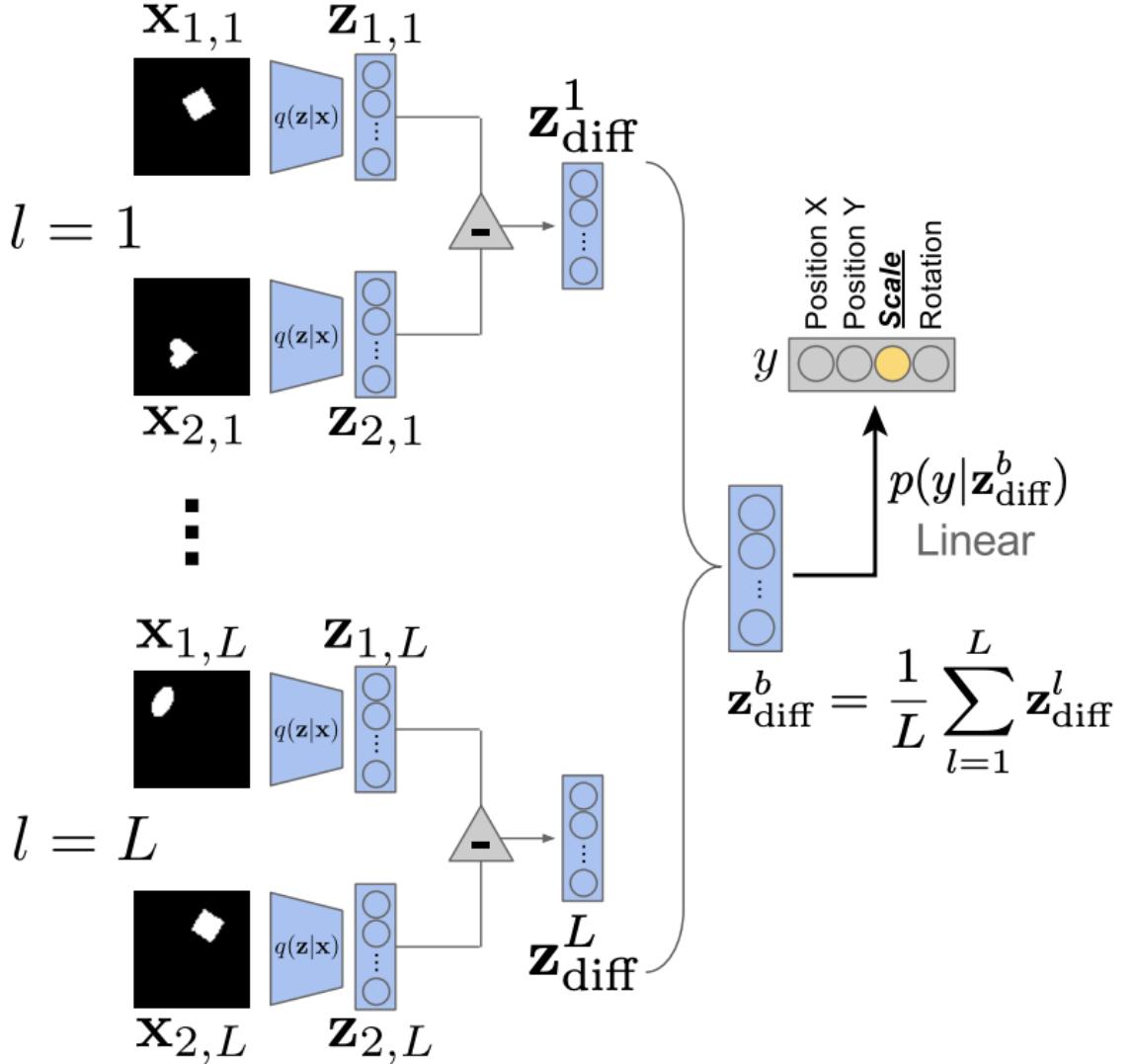


Figure 2: A schematic of the proposed disentanglement metric from [4]. Each vector  $\mathbf{z}_{\text{diff}}^l$  in the process is generated from pairs of latent representations  $z_{1,l}, z_{2,l}$  with the exact same value for a chose generative factor  $y$ . All of the  $L$  pairs in batch  $b$  may differ in terms of the value of the generative factor, but the choice of  $y$  is constant between pairs. The linear classifier,  $p(y|\mathbf{z}_{\text{diff}}^b)$  is a multivariate classification algorithm whose goal is to predict which generative factor was used to generate the average difference vector  $\bar{\mathbf{z}}_{\text{diff}}^b$  (in the text, we refer to this average difference vector as  $\bar{\mathbf{z}}_{\text{diff}}$ ).

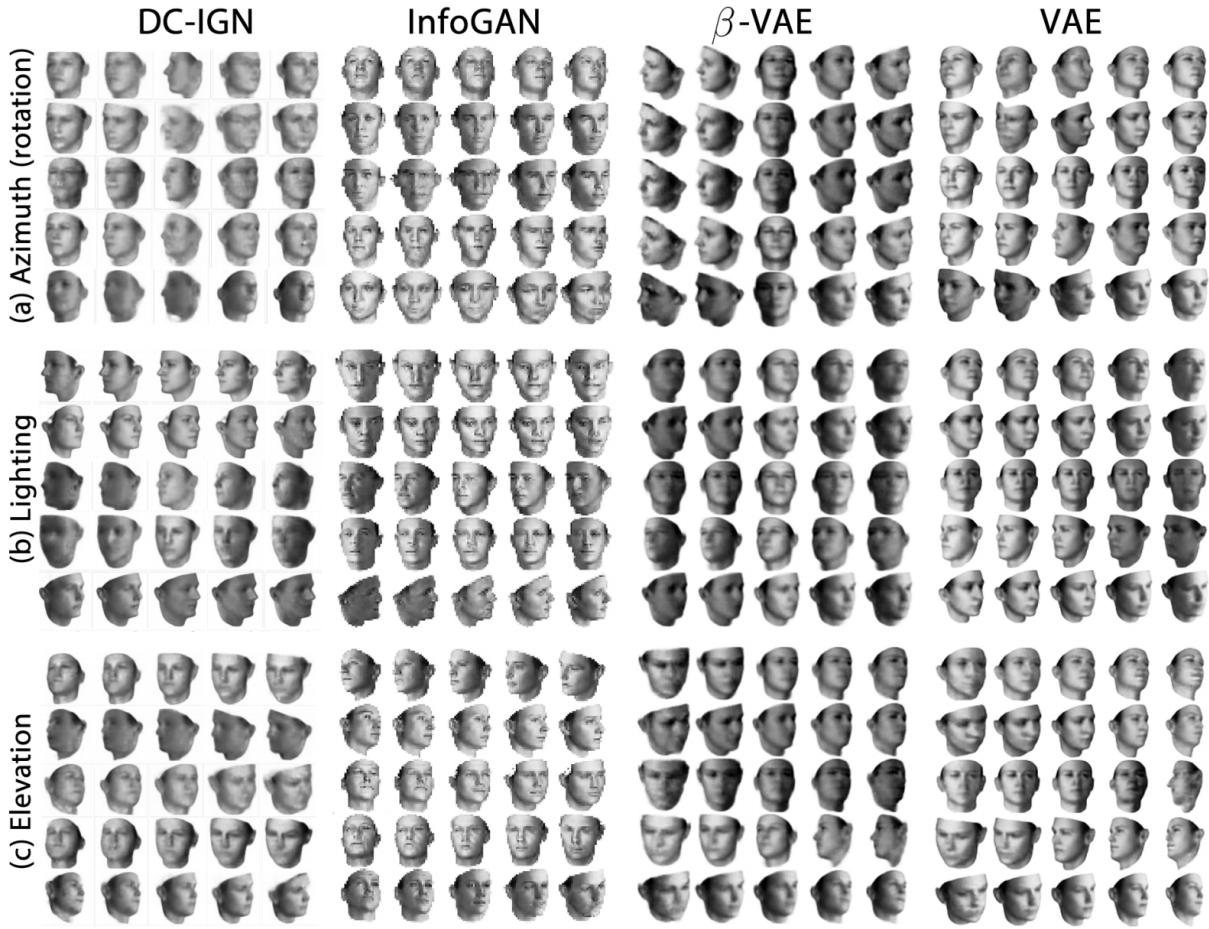


Figure 3: Qualitative comparison of three dimensionality reduction techniques with  $\beta$ -VAE ( $\beta$ -VAE) on the 3D faces dataset. Although all models are relatively good at disentangling (b) lighting and (c) elevation,  $\beta$ -VAE far outperforms the other three in terms of disentangling azimuth.

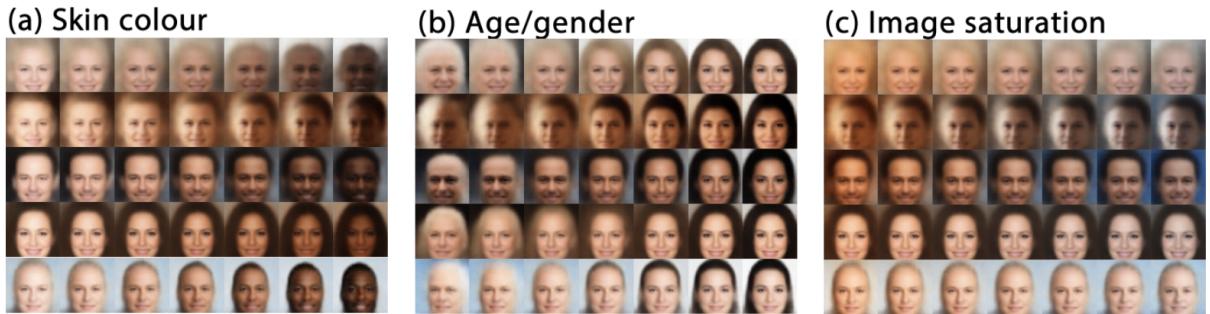


Figure 4:  $\beta$ -VAE latent factors for the CelebA dataset. Each panel shows a traversal of one of the latent dimensions with all others fixed. The model is able to learn clear disentangled representations for (a) skin colour, (b) age/gender, and (c) image saturation.

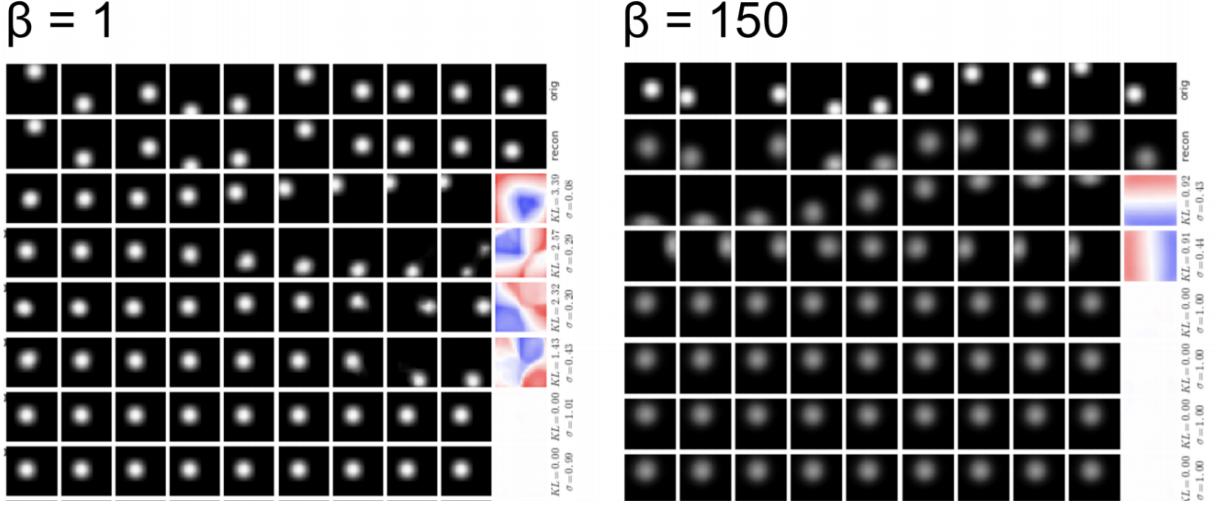


Figure 5: Comparison of vanilla VAE ( $\beta = 1$ ) with strongly disentangled  $\beta$ -VAE ( $\beta = 150$ ). The dataset is a set of Gaussian blobs located in various locations in the cartesian plane (top row of both panels). The second row of each panel is the reconstructions of the original images. The remaining rows are traversals of one of the latent dimensions. The latent dimensions are arranged from bottom to top in order of highest to lowest  $D_{KL}$  from the prior standard Gaussian.

Model	Disentanglement metric score
<i>Ground truth</i>	100%
Raw pixels	$45.75 \pm 0.8\%$
PCA	$84.9 \pm 0.4\%$
ICA	$42.03 \pm 10.6\%$
DC-IGN	<b><math>99.3 \pm 0.1\%</math></b>
InfoGAN	$73.5 \pm 0.9\%$
VAE untrained	$44.14 \pm 2.5\%$
VAE	$61.58 \pm 0.5\%$
<b><math>\beta</math>-VAE</b>	<b><math>99.23 \pm 0.1\%</math></b>

Figure 6: Results from the [4] comparing the performance of different dimensionality reduction techniques vis-a-vis the proposed disentanglement metric.  $\beta$ -VAE and DC-IGN both reach nearly perfect accuracy, whereas a vanilla VAE is only able to successfully classify just over 60% of generative factors based on the latent representations.

$(\beta)$ -VAE from simple linear dimensionality reduction techniques, such as PCA and CCA, but also from nonlinear manifold or kernel techniques such as LLE, ISOMAP, MDS, and kernel PCA.

However, VAE-type models are still similar to other dimensionality reduction techniques in the sense that the goal is to embed the data in a latent space where the geometry is more comprehensible to a human or simple classifier. Further,  $\beta$ -VAE’s augmentation over vanilla VAE is to make the axes of this lower-dimensional space align with human-interpretable generative factors. This contrasts with a method like PCA, where the goal is to find a lower-dimensional space that maximizes variance along each axis, with no regard to the relationship between these low-dimensional axes and the factors that generated the data.

### 3.4 Weaknesses

Despite the promise of learning a disentangled latent representation of complex image or other high dimensional data, the  $\beta$ -VAE is not without significant weaknesses. Most glaring is the arbitrariness of the choice of  $\beta$ . The authors state that if even a small amount of labelled generative factors are available for the data,  $\beta$  can be chosen through a hyperparameter sweep (e.g. a gridsearch or some other hyperparameter optimization algorithm). If such labels are not available, then  $\beta$  can be chosen through “visual inspection of what effect the traversal of each single latent unit  $z_m$  has on the generated images  $(\vec{x}|\vec{z})$  in pixel space” [4]. In other words,  $\beta$  is chosen through trial-and-error. In most cases, few if any labels are available, meaning that the visual inspection technique, along with some “rules of thumb” will be all that is available to the analyst.

Furthermore, there are few guarantees that this technique will work outside of the toy datasets employed in the paper. The major assumption underlying the technique is that the *generative factors exist*. That is, we have no reason to believe that there genuinely is a data-generating process that creates each image from a random choice of a small set of discrete or continuous generative factors. This is crucial, since if there is no set of generative factors to be recovered, then no choice of  $\beta$  or careful optimization procedure will allow these phantom factors to be found. As well, even if the factors exist, there is no guarantee that they are in any way human-interpretable. I.e. we may be able to recover some kind of generative factors, but we might not know them when we see them.

## 4 Experiment

In this section we explore the results of a Python implementation of  $\beta$ -VAE on a toy EMNIST (extension to MNIST) dataset of handwritten digits and characters. The code can be found here: <https://github.com/benlevyx/beta-vae>.

### 4.1 Model architecture and dataset

We use the EMNIST dataset, which is an image dataset similar to MNIST, except that it also includes handwritten letters [6]. Each digit is  $28 \times 28$  pixels with a single grayscale channel. The dataset includes handwritten digits from 0 to 9 and letters from a to z, including lowercase and uppercase. For simplicity, we take a subset of the data comprising the letters A through D (uppercase and lowercase).

The model encoder and decoder both consist of feed-forward neural networks. The 2D images are flattened to a vector of 784, which are then fed into the network in batches of 256 for computational efficiency. Since the goal of this experiment is to be able to plot the low-dimensional representations of the data for different levels of  $\beta$ , we use a latent dimension  $d = 2$ . We train models with  $\beta = 0.0, 1.0, 5.0$  for 15 epochs each using the Adam optimizer with learning rate  $10^{-3}$ .

### 4.2 Effect of $\beta$ on latent space

In figure 7, we can see the effect of  $\beta$  on the latent space. The starker difference is between  $\beta = 0$  and the cases where  $\beta > 0$ . In the former, we see that the latent space is extended along a line such that each point is essentially assigned its own location in that latent space. This resembles the latent space for a non-Bayesian

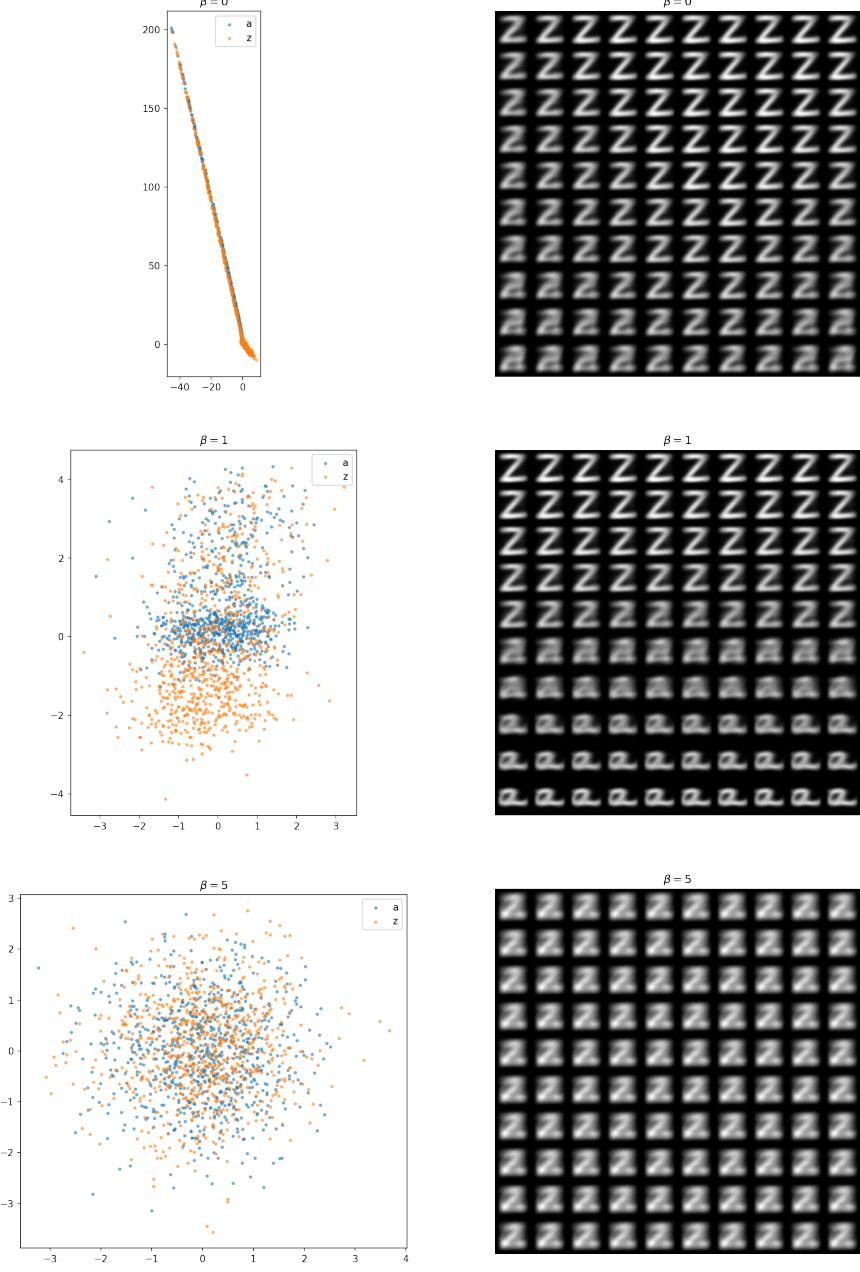


Figure 7: Results from experiment with different levels of  $\beta$ . Each row represents a model trained with a different level of  $\beta$ . The left column plots the 2-dimensional representation of the evaluation data, colour-coded according to whether the point is an ‘A’ or ‘Z’. The right column shows a traversal of the latent space in the fixed range  $(-4, 4)$  for both latent dimensions.

autoencoder, since the model is only incentivized to reconstruct the data faithfully and is not encouraged to learn a smooth latent space.

As  $\beta$  increases, we see that the latent space becomes closer to a standard Gaussian distribution. This is a double-edged sword: although the latent space is smoother and perhaps generalizes better, we completely lose the ability to discriminate between the two classes.

Of course, this is an extreme case meant to illustrate the effect of strong and weak regularization on the

geometry of the latent space. In the real world, we would implement a  $\beta$ -VAE with a larger bottleneck dimension (2 is quite extreme), which would also allow for more fine-tuning of  $\beta$ .

## References

- [1] David Blei. *Variational Inference*. Lecture notes for COS597c. 2011. URL: <https://www.cs.princeton.edu/courses/archive/fall11/cos597C/lectures/variational-inference-i.pdf>.
- [2] Yoshua Bengio, Aaron Courville, and Pascal Vincent. “Representation Learning: A Review and New Perspectives. ArXiv e-prints”. In: *arXiv preprint arXiv:1206.5538* (2012).
- [3] Diederik P Kingma and Max Welling. “Auto-encoding variational bayes”. In: *arXiv preprint arXiv:1312.6114* (2013).
- [4] Irina Higgins et al. “beta-vae: Learning basic visual concepts with a constrained variational framework”. In: (2016).
- [5] <https://jaan.io/what-is-variational-autoencoder-vae-tutorial/>. *Tutorial - What is a variational autoencoder?* 2016.
- [6] Gregory Cohen et al. “EMNIST: Extending MNIST to handwritten letters”. In: *2017 International Joint Conference on Neural Networks (IJCNN)*. IEEE. 2017, pp. 2921–2926.
- [7] Christopher P Burgess et al. “Understanding disentangling in  $\beta$ -VAE”. In: *arXiv preprint arXiv:1804.03599* (2018).
- [8] Lilian Weng. “From Autoencoder to Beta-VAE”. In: *lilianweng.github.io/lil-log* (2018). URL: <http://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>.