# Disentangling Variational Autoencoders

Rafael Pastrana[a,*]

[a]*School of Architecture, Princeton University, United States of America*

**Abstract**

A variational autoencoder (VAE) is a probabilistic machine learning framework for posterior inference that projects an input set of high-dimensional data to a lower-dimensional, latent space. The latent space learned with a VAE offers exciting opportunities to develop new data-driven design processes in creative disciplines, in particular, to automate the generation of multiple novel designs that are aesthetically reminiscent of the input data but that were unseen during training. However, the learned latent space is typically disorganized and entangled: traversing the latent space along a single dimension does not result in changes to single visual attributes of the data. The lack of latent structure impedes designers from deliberately controlling the visual attributes of new designs generated from the latent space. This paper presents an experimental study that investigates latent space disentanglement. We implement three different VAE models from the literature and train them on a publicly available dataset of 60,000 images of hand-written digits. We perform a sensitivity analysis to find a small number of latent dimensions necessary to maximize a lower bound to the log marginal likelihood of the data. Furthermore, we investigate the trade-offs between the quality of the reconstruction of the decoded images and the level of disentanglement of the latent space. We are able to automatically align three latent dimensions with three interpretable visual properties of the digits: line weight, tilt and width. Our experiments suggest that i) increasing the contribution of the Kullback-Leibler divergence between the prior over the latents and the variational distribution to the evidence lower bound, and ii) conditioning input image class enhances the learning of a disentangled latent space with a VAE.

*Keywords:* variational autoencoder, generative model, disentanglement, latent space, machine learning, neural networks

## 1. Introduction

A *variational autoencoder* (VAE) is a posterior inference framework to fit latent variable models. A VAE-fitted model has two attributes relevant to creative disciplines like type and architectural design. First, the model can project the high-dimensional input data points **x** to a lower-dimensional latent space parametrized by a latent variable **z**. Second, the learned latent space is typically a smooth and continuous manifold that can be utilized for synthetic data generation. The latent space attains these properties because it is constrained to fit a variational probability distribution $q$ at training time. As a result, it is possible to take samples from the VAE-fitted latent space and then decode them to generate new data points $\hat{\mathbf{x}}$ that resemble the attributes of **x** but were never part of the input data.

Powered by neural networks, VAEs have gained state of the art prominence in creative design tasks that enable the generation of synthetic but photorealistic images of people who do not exist [1], the development of new music creation tools [2] and the design of novel dancing choreographies [3]. However, the latent space spanned by **z** and learned with a VAE is typically entangled and thus lacks two important factors to operationalize it in design tasks: *independence* and *interpretability* [4]. The goal of attaining a disentangled latent space, on the other hand, is to find a low-dimensional representations of the data **x** where single latent dimensions $z_i$ are sensitive to changes in single generative factors while being relatively invariant to changes in others [5]. To advance the use of machine learning-powered probabilistic tools in design, the desiderata is to develop and explore methodologies that produce a disentangled latent space that captures visual basic concepts imbued in the input data, such as scale, tilt, orientation or color. This way, one can have finer and deliberate control when generating novel data from that latent space.

In this paper, we investigate three different methodologies to disentangle the latent space learned by a VAE. First, we implement a standard VAE and replicate a portion of the numerical results in the original VAE paper on the autoencoding variational Bayes by Kingma and Welling [6]. Next, we explore the $\beta$-VAE [7] framework, an extension to the original VAE formulations that introduces a hyperparameter $\beta$ to hypothetically enforce a semantically richer yet efficient representation of the data in the latent space. Lastly, we investigate to what extent conditioning the dataset **x** on class labels **u** facilitates disentanglement with a $\beta$-VAE [8]. We perform multiple experiments to analyze the impact of $\beta$ on the ELBO and qualitatively assess the disentanglement of the latents.

This paper is accompanied by a blog post with interactive graphics and plots which the reader can access at https://bit.ly/3KyZOOf. For reproducibility purposes, we also make available a Pytorch [9] implementation of our work at https://github.com/arpastrana/neu_vae.

---

*Corresponding author

*Email address:* arpastrana@princeton.edu (Rafael Pastrana)
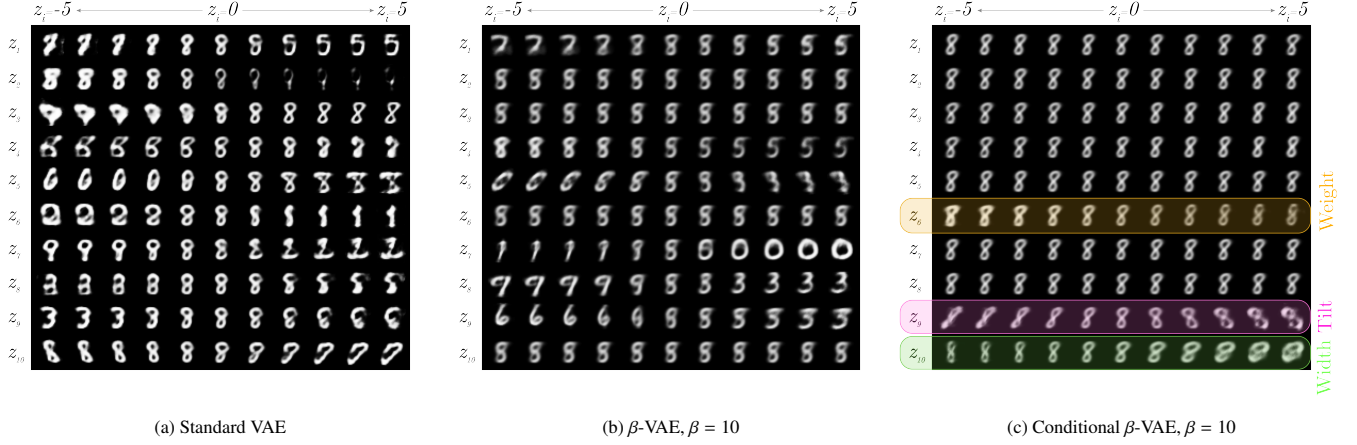
Figure 1: **(Dis)entangled latent spaces**. Latent space disentanglement comparison for three different VAEs that project the input data to a 10-dimensional latent space, $J = 10$. We feed in the same image to each of the three VAEs and traverse the latent space as we describe in Section 3.4. The standard VAE produces the sharpest digit reconstruction. Only the latent space learned by conditional $\beta$-VAE shows evidence of a disentangled latent space. Walking along latent dimensions $z_6$, $z_9$ and $z_{10}$ evidences that these dimensions align with three different visual properties: line weight, tilt and width.

## 2. Theoretical background

The observed data $\mathbf{x} \in \mathbb{R}^N$ is a random variable that follows a probability distribution $p(\mathbf{x})$. If we parametrize the density of the observed data with an unobserved variable $\mathbf{z} \in \mathbb{R}^J$, the log joint distribution of $\mathbf{x}$ and $\mathbf{z}$ is expressed as:

$$\log p(\mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \qquad (1)$$

Marginalizing the joint distribution over $\mathbf{z}$ to recover the log probability of the data $\log p(\mathbf{x})$ is generally intractable. Therefore, we maximize a lower bound to the log marginal likelihood of the data, the Evidence Lower Bound (ELBO), and fit instead the parameters of an variational distribution $q_\phi(\mathbf{z})$.

$$ELBO(\theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \qquad (2)$$

where $\phi$ and $\theta$ are the parameters of a VAE's encoder and decoder respectively. These parameters can be computed by non-linear function approximators such as neural networks.

We can use MCMC to compute the expectation of the conditional likelihood of the data under the variational distribution, $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$. In practice, a single sample per data point per batch suffices to approximate the expectation, which leads to the following updated expression for the ELBO:

$$ELBO(\theta, \phi) \approx \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \qquad (3)$$

The first term of this approximation can be understood as the reconstruction loss of e.g. Bernoulli- or Gaussian-decoded data points, while the second corresponds to a regularization term which penalizes the Kullback-Leibler (KL) divergence between the variational posterior $q_\phi(\mathbf{z}|\mathbf{x})$ and a prior $p(\mathbf{z})$ over the latent $\mathbf{z}$. To make the optimization of the ELBO tractable and differentiable, we assume that both the posterior $q_\phi(\mathbf{z}|\mathbf{x})$ and the prior $p(\mathbf{z})$ are Gaussian distributions $\mathcal{N}(\mu_j, \sigma_j^2)$ with mean $\mu_j$ and variance $\sigma_j^2$ per latent dimension $j$. The choice of a Gaussian allows the calculation of the KL divergence between

the variational posterior and the prior over the latents in closed form:

$$\frac{1}{2} \sum_{j=1}^{J} (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2) \qquad (4)$$

Furthermore, the Gaussian assumption enables gradient backpropagation using the *reparametrization trick*, where we can sample from the variational distribution $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mu, \sigma^2)$ via an affine transformation of an auxiliary noise variable $\epsilon$ sampled from a standard normal, $\epsilon \sim \mathcal{N}(0, 1)$:

$$\mathbf{z_j} = \mu_j + \epsilon \, \sigma_j \qquad (5)$$

### 2.1. $\beta$-VAE

To focus on learning statistically independent latent factors, the authors of [7] include an additional hyperparameter $\beta$ in the loss function proposed in the original formulation of the VAE. This hyperparameter scales the weight of the KL Divergence term in the ELBO:

$$ELBO(\theta, \phi) \approx \log p_\theta(\mathbf{x}|\mathbf{z}) - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \qquad (6)$$

A $\beta$-VAE with $\beta = 1$ corresponds to the original VAE formulation. Employing values of $\beta > 1$ hypothetically put pressure on the VAE bottleneck to match the prior $p(\mathbf{z})$ and thus promote the learning of a more efficient latent data representation. As reported in [7], disentanglement comes at the cost of a diminished reconstruction quality of $\mathbf{x}$. Moreover, too small or too large $\beta$ values may not necessarily lead to disentangled latents. Therefore, $\beta$ needs to be calibrated using qualitatively (e.g. using visual heuristics) or quantitatively approaches (e.g. developing a quantitative disentanglement metric after appending a simple linear classifier to a trained VAE).

### 2.2. Conditional β-VAE

Recent work in representation learning suggests that the latent space of fitted latent variable models cannot be disentangled without supervision [8, 10]. One way to fit disentangled latent spaces is thus to build a conditionally factorized prior over the latent variables $p(\mathbf{z}|\mathbf{u})$, which is possible by concurrently observing an auxiliary variable $\mathbf{u}$ that corresponds to the time index in a time series, previous data points, or class labels [10]. This means that the model takes as input a dataset with observation pairs $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{u}_i), ..., (\mathbf{x}_n, \mathbf{u}_n)\}$ instead of $\mathcal{D} = \{\mathbf{x}_i, ..., \mathbf{x}_n\}$. The total log density of the data $\mathbf{x}$ conditioned on labels $\mathbf{u}$ is found after marginalizing the latent variables $\mathbf{z}$:

$$\log p(\mathbf{x}|\mathbf{u}) = \int \log p(\mathbf{x}|\mathbf{z})p(\mathbf{z}|\mathbf{u})d\mathbf{z} \tag{7}$$

## 3. Method

### 3.1. Dataset

We train our model with the MNIST dataset [11]. This dataset consists of 60,000 images of ten different hand-written digits in the range $[0-9]$. The resolution per image is of $28 \times 28$ pixels. We use a randomized data loader to manage the train and validation sets that are utilized to fit the parameters of the VAE.

### 3.2. VAE Architecture

For consistency, we use the same neural architecture as in Kingma and Welling [6] in all our experiments unless otherwise noted. We use multilayer perceptrons (MLPs) in the encoder and decoder of the VAE. Every perceptron unit in a MLP is activated using the hyperbolic tangent function, except for the output layer in the decoder where we apply a sigmoid non-linearity.
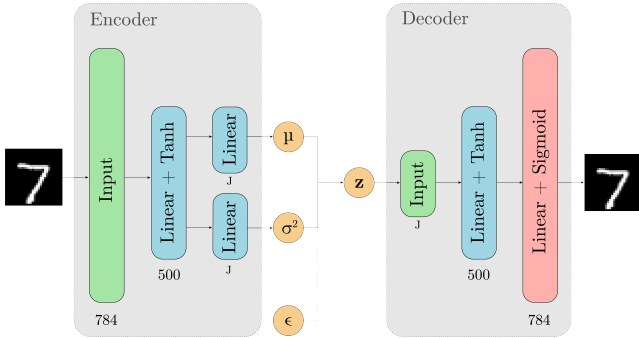


Figure 2: **VAE architecture**. The architecture of the VAE in all our experiments follows the structure reported in [6], except for when we work with a conditional β-VAE. The VAE consists of a three-layer encoder followed by a three-layer decoder. The encoder outputs the mean $\mu$ and the variance $\sigma^2$ of the variational posterior. We then sample $\epsilon$ from a standard normal $\mathcal{N}(0, I)$ and use the reparametrization trick to get the latent vector $\mathbf{z}$ that is fed to the decoder.

### 3.2.1. Gaussian encoder

The encoder consists of an input linear layer with 784 units (corresponding to a flattened $28 \times 28$ image) followed by a fully-connected, 500-units linear layer. To calculate the mean $\mu$ and the variance $\sigma$ of the variational distribution $q$, the output layer of the decoder is another linear layer whose number of hidden units is twice the dimensionality of $\mathbf{z}$, that is $2 \times J$.

### 3.2.2. Bernoulli decoder

After applying the reparametrization trick, we feed samples from $\mathbf{z}$ into the decoder's input linear layer which has the same number of units as in the size of the latent space, $J$. We complete the decoder with an intermediate 500-unit linear layer followed by a 784-unit output linear layer.

### 3.3. Training

We use Adagrad [12] with a fixed learning rate of 0.01 in all our the experiments. We use an image batch size to 100 to compute gradients. We cap training at 200 epochs. In our work, the optimization objective well saturates within this budget of epochs. The optimization objective is to minimize the negative value of the ELBO. The reconstruction error is computed using the binary cross-entropy, and the KL divergence is calculated analytically since we assume Gaussians for the variational posterior and the prior.

### 3.4. Disentanglement evaluation

We traverse the $J$ distinct dimensions of the fitted $\mathbf{z}$, as in the work of Higgins, et al. [7], to visually examine the disentanglement quality of the learned latent representations.

First, we input a digit image into the VAE encoder to obtain the mean $\mu$ and the variance $\sigma$ of the variational posterior over the latent space. Next, we sample a latent variable $\mathbf{z}$ from the variational posterior using the reparametrization trick. We then alter each of the $J$ dimensions of the sampled $\mathbf{z}$: the value of every latent dimension $z_j$ is first zeroed and then traversed in the range $[-5, 5]$ in ten steps while keeping all the values of all the others latent dimensions $z_{\neq j}$ fixed. The updated $\mathbf{z}$ is then passed to the decoder to generate a new digit image that reflects the effect in image space of the perturbation we made to the data in latent space.

## 4. Results

### 4.1. Standard VAE

We reproduce a portion of the experiments on presented in [6]. We gradually increase the size of the latent variable $\mathbf{z}$, experimenting with dimension steps of $J \in \{3, 5, 10, 20, 200\}$, and observe what their impact was on the ELBO.

As Figure 3a shows, there are only minimal improvements to the decreasing the value of the objective function after 10 latent dimensions, $J > 10$. The close gap between the training and testing curves in within the allocated VAE's training epochs suggest the trained VAEs does not overfit, even when the dimensionality of the latent space scales up to $J = 200$. This finding is in accordance with Kingma and Welling who
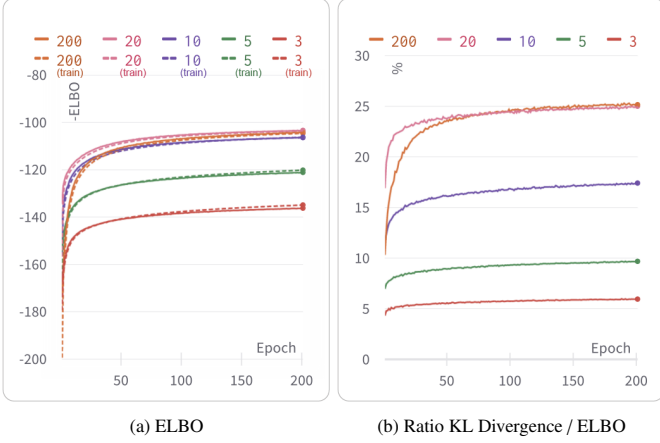
(a) ELBO

(b) Ratio KL Divergence / ELBO

Figure 3: **Standard VAE**. We reproduced of a portion of the experiments in [6] to validate our machine learning pipeline. We trained a standard VAE for 200 epochs and tested the effect of changing the number of dimensions of the latent space $\mathbf{z}$ from $J = 3$ until $J = 200$ in both train and test datasets. The plots show that the ELBO saturated well within the defined number of epochs in all four cases without overfitting.



(a) ELBO

(b) Ratio KL Divergence / ELBO

Figure 4: $\beta$-**VAE**. Sensitivity study on the effect that different values of $\beta$ have on the ELBO. In Figure 4b, the participation of the KL divergence in the ELBO is lower when $\beta = 10$ than when $\beta \in \{3, 5\}$.

attributed this phenomenon to the regularizing effect of the Gaussian prior and the KL divergence term in the ELBO [6]. However, latent variables with larger dimensionality raise the proportion that the KL divergence contributes to the ELBO. This contribution oscillates between about 6% when $J = 3$ and 25% when $J = 200$ (see Figure 3b). The increase in the contribution of the KL divergence term of the ELBO can be understood as a cumulative penalty accrued because the VAE is simultaneously fitting more dimensions of $\mathbf{z}$ to the Gaussian prior.

### 4.2. $\beta$-VAE

We carry out a sensitivity analysis on the hyperparameter $\beta$ by training a $\beta$-VAE with different values of $\beta \in \{1, 3, 5, 10, 20\}$. We fix the number of latent dimensions to $J = 10$ in all the experiments hereafter for the $\beta$-VAE as further increases to $J$ did not lead to decreasing the negative value of the ELBO as we report in Section 4.1). Figure 4a shows higher values of $\beta$ lead to larger negative ELBO. The value of the ELBO after 200 epochs decreases by almost 70% when $\beta = 10$ compared to when $\beta = 1$. This result suggests that further increasing $\beta$ would lead to a poorer log-likelihood estimate of the data. In contrast, we observe that the contribution of the KL divergence between the variational posterior $q_\phi(\mathbf{z}|\mathbf{x})$ and the Gaussian prior $p(\mathbf{z})$ to the ELBO diminishes as $\beta$ is larger, except for when $\beta = 10$ (see Figure 4b ).

We also study the effect that increasing $\beta$ has on the visual quality of the digits reconstructed by the decoder of the $\beta$-VAE (see Figure 5). In contrast to the standard VAE (i.e. equivalent to setting $\beta = 1$), the decoded images are consistently blurrier but still preserved the key features. However, digits visually unidentifiable when $\beta = 20$. Figure 1b shows that visual patterns begin to emerge when we traverse the latents with $\beta = 10$. For example, the latent dimensions $z_2, z_3, z_6, z_{10}$ do not experience any changes after traversing them, which suggests
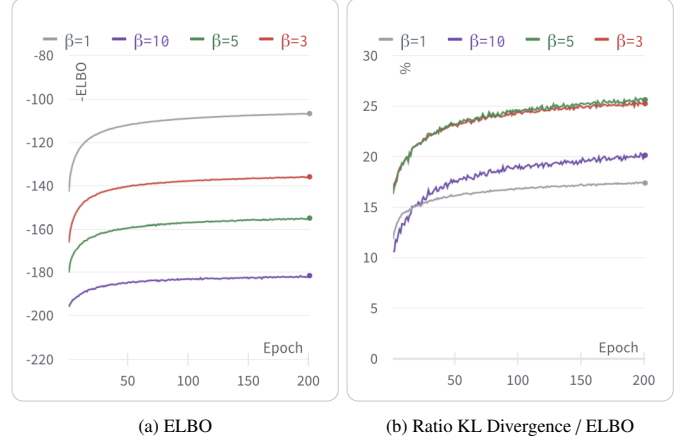
that the $\beta$-VAE model started to concentrate data variance over particular latent directions and neutralize information along others. However, the visual recognition of single generative factors corresponding to single feature dimensions is still unclear. For example, in Figure 1b, traversing single latent dimensions results in inter-digit transformations instead of modifying any intrinsic visual property: digit 8 transitioned between digit 6 ($z_9 = -5$) when and digit 3 ($z_9 = 5$) when traversing the latent space along dimension $z_9$.

### 4.3. Conditional $\beta$-VAE

We study the effect that using the class labels had on disentangling the latent space learned by a $\beta$-VAE. We set the number of dimensions of the latent vector $\mathbf{z}$ in all the experiments we describe in this section to $J = 10$, but we monotonically increase the value of $\beta$ in the range $\beta \in \{1, 3, 5, 10\}$. The VAE ingests data pairs $(\mathbf{x}, \mathbf{u})$, where $\mathbf{u}$ is the digit class label. The class label $\mathbf{u}$ is encoded as a one-hot vector of size ten that is concatenated with the 784-dimensional original input vector. We also concatenate the class one-hot vector to the sampled
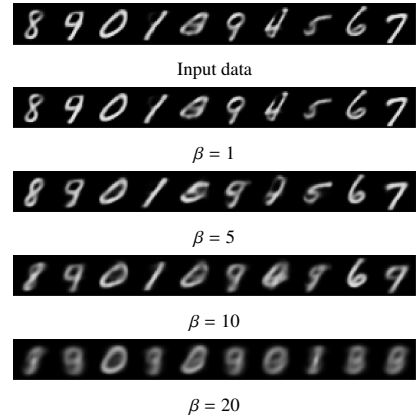


Input data

$\beta = 1$

$\beta = 5$

$\beta = 10$

$\beta = 20$

Figure 5: **Digit reconstruction**. How does the choice of $\beta$ affect the visual aspect of the hand-written digits decoded by a $\beta$-VAE? Higher values of $\beta$ are conducive to blurrier reconstructions.
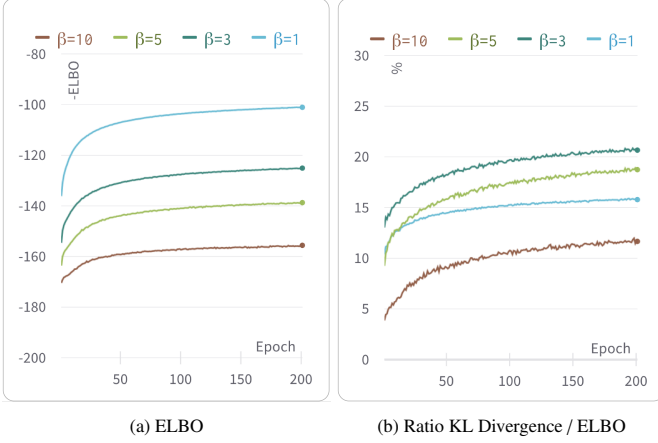
(a) ELBO

(b) Ratio KL Divergence / ELBO

Figure 6: **Conditional $\beta$-VAE**. How does the addition of labels to the input dataset affect the learning process of a $\beta$-VAE? Overall, the minimization of the negative ELBO follows a similar downward trend than that of the unconditioned $\beta$-VAE for the four values of $\beta$ we test, $\beta \in \{1, 3, 5, 10\}$.



(a) ELBO

(b) Ratio KL Divergence / ELBO

Figure 7: **Standard VAE vs. $\beta$-VAE vs. conditional $\beta$-VAE**. Metrics comparison between a standard VAE ($\beta = 1$) [6] and a $\beta$-VAE ($\beta = 10$) [7] with unlabeled and labeled data. $J = 10$ in the all four tests.

latent vector $\mathbf{z}$ that is input to the VAE decoder. Therefore, we adjust the architecture of the VAE accordingly and add ten more units to the input layers of both the encoder and decoder: the number of of units in the first layer of the encoder increases from 784 to 794 whereas the number of units in the first layer of the decoder is $J + 10 = 20$.

Figure 6a shows that the ELBO and the KL divergence contribution fluctuates in agreement with the unconditional $\beta$-VAE, regardless of the choice of $\beta$: higher $\beta$ values lead to a lower ELBO estimate. The KL Divergence over ELBO ratio was lowest when $\beta = 10$, even lower than the case where $\beta = 1$.

The figure also helps identifying that adding the class labels increase the fit of the model to the data. The negative of the approximation to the log marginal likelihood is consistently 15% lower for the conditioned $\beta$-VAE than it is in the unconditioned case. For example, when $\beta = 5$, the negative of the ELBO is -160 in the former case (Figure 4a), whereas this figure decreases to -140 in the latter (Figure 6a). Moreover, when $\beta \in \{1, 3, 5\}$, Figure 6b tells us that the participation of the KL divergence in the ELBO increases steadily from 15% until over to 20%, but it interestingly sees a sharp decline when $\beta = 10$.

After training the $\beta$-VAE on the labeled dataset, we observe that the disentanglement of the latent space is qualitatively clearer compared to that produced by the $\beta$-VAE framework trained on the unlabeled dataset with $\beta = 10$. As we show in Figure 1c, one of the most significant findings is that after running inference on images of the digits, seven of the ten latent dimensions $z_i$ are not changed by the latent traversals. The second most significant result is that dimensions $z_6$, $z_9$, and $z_{10}$ concentrate all the information related to the digit's generative factors, which we associate to the latent directions that control the line-weight, lateral tilt, and width of the hand-written digits. These disentangled latent directions remain consistent even as we traverse the latent space of the conditional $\beta$-VAE for all nine digits (see Figure 8).
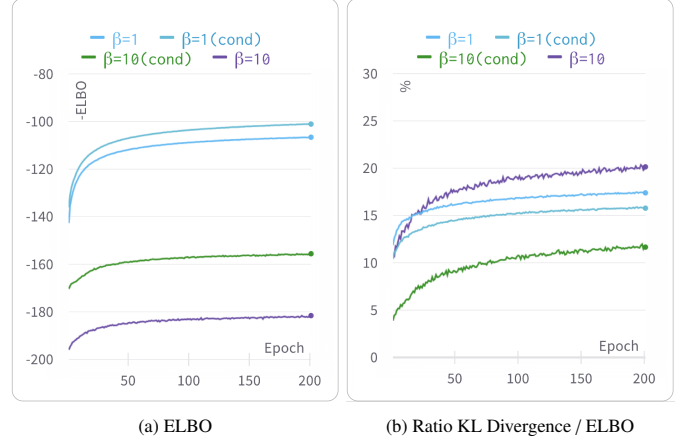
### 4.4. Comparison: $\beta$-VAE and conditional $\beta$-VAE

We examine and compare the behavior of $\beta$-VAE, in both the conditioned and unconditioned case with two different reference values $\beta = 1$ and $\beta = 10$. In terms of the magnitude of the ELBO, the conditioned dataset with $\beta = 1$ exhibits the highest value and the unconditioned case with $\beta = 10$ the lowest, as depicted in Figure 7a. This concur with our previous experimental results that show that higher values of $\beta$ lead to lower values of the ELBO.

One of the most interesting findings pertains to the KL divergence to ELBO ratio, where conditioning the dataset leads to changes in the trend we observed in previous experiments. In Figure 7b, for example, the contribution of the KL divergence to the ELBO is almost cut from 20% down to almost 10% when $\beta = 10$, Since the class-labeled dataset with this value of $\beta$ produces the best disentanglement results so far, we hypothesize whether minimizing the participation of the KL divergence in the ELBO calculation while preserving a reasonable reconstruction error is conducive to good latent space disentanglement. A more detailed investigation on the matter is left to future work.

### 5. Conclusion

In this paper, we implemented and trained three different VAEs to study latent space disentanglement on a dataset of 60,000 images of hand-written digits: a standard VAE [6], a $\beta$-VAE [7] and a conditional $\beta$-VAE [8]. We analyzed the effect of varying the magnitude of the hyperparameter $\beta$, and that of conditioning images on discrete class labels $\mathbf{u}$ on the disentanglement of the latents.

We found that using $\beta = 10$ together with label-conditioned data led to the clearest level of disentanglement through our experiments, revealing single latent directions that allowed for individual control of three generative factors in the images supplied to the VAE. These generative factors corresponded to the line weight, the tilt and the width of the digit images. This finding supports the idea that good disentanglement is
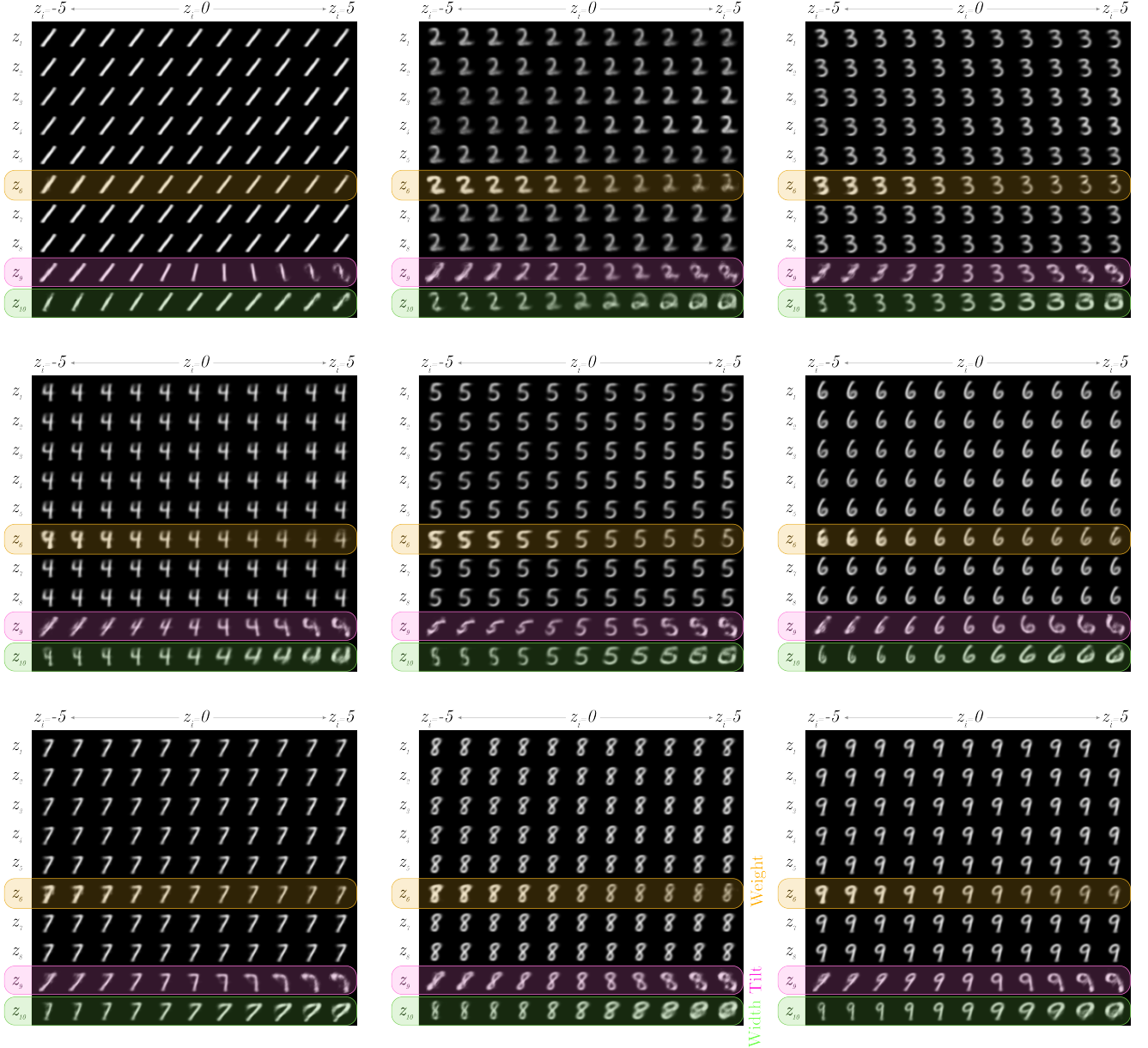
Figure 8: **Disentangled latent space**. The alignment between line weight, tilt and digit width and the latent dimensions $z_6$, $z_9$ and $z_{10}$ we found after traversing the latent space learned by the conditional $\beta$-VAE remained consistent across all nine digits.

contingent upon some a priori level of supervision on the input data [8]. However, our experiments also suggest that using a moderate value of $\beta$ to scale the KL divergence in the ELBO is also instrumental to arrive at a disentangled latent space.

Disentangled representation learning is an active area of research with plenty of interesting challenges ahead. Several routes to extend our work are thus outlined. First, we plan to develop a consistent and robust quantitative disentanglement metric –resorting to qualitative visual heuristics to evaluate latent space disentanglement may be cumbersome. We are interested in following a methodology similar to the prediction-based measurements exposed in [7] to this end. An alternative approach to disentanglement we would like to explore later is to directly perform *independent component analysis* (ICA) on the learned latents **z** and to evaluate whether ICA facilitates finding disentangled directions without having to use the hyperparameter $\beta$ or any class labels, as suggested in recent work [10].

We also aim to leverage other neural networks in the encoder and the decoder components of the VAE to learn better-quality latents by exploiting the symmetries in the data, in particular, the translational invariance of image-structured data that we used herein. One concrete example is to use convolutional neural networks instead of multilayer perceptrons. Finally, we anticipate working with richer and structured prior distributions for the variational distribution, other than Gaussians, which may ultimately better capture the hidden and disentangled generative structure of the input data.

# References

[1] A. Vahdat, J. Kautz, NVAE: A Deep Hierarchical Variational Autoencoder, arXiv:2007.03898 [cs, stat] (Oct. 2020). arXiv:2007.03898.

[2] A. Roberts, J. Engel, D. Eck, Hierarchical Variational Autoencoders for Music, in: 31st Conference on Neural Information Processing Systems, Long Beach, CA, USA, 2017, p. 6.

[3] M. Pettee, C. Shimmin, D. Duhaime, I. Vidrin, Beyond Imitation: Generative and Variational Choreography via Machine Learning, arXiv:1907.05297 [cs, stat] (Jul. 2019). arXiv:1907.05297.

[4] I. Higgins, D. Amos, D. Pfau, S. Racaniere, L. Matthey, D. Rezende, A. Lerchner, Towards a Definition of Disentangled Representations, arXiv:1812.02230 [cs, stat] (Dec. 2018). arXiv:1812.02230.

[5] Y. Bengio, A. Courville, P. Vincent, Representation Learning: A Review and New Perspectives, IEEE Transactions on Pattern Analysis and Machine Intelligence 35 (8) (2013) 1798–1828. doi:10.1109/TPAMI.2013.50.

[6] D. P. Kingma, M. Welling, Auto-Encoding Variational Bayes, arXiv:1312.6114 [cs, stat] (May 2014). arXiv:1312.6114.

[7] 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 (Nov. 2016).

[8] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, O. Bachem, Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations, in: K. Chaudhuri, R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning, Vol. 97 of Proceedings of Machine Learning Research, PMLR, Long Beach, California, USA, 2019, pp. 4114–4124.

[9] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, PyTorch: An imperative style, high-performance deep learning library, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems 32, Curran Associates, Inc., 2019, pp. 8024–8035.

[10] I. Khemakhem, D. Kingma, R. Monti, A. Hyvarinen, Variational Autoencoders and Nonlinear ICA: A Unifying Framework, in: S. Chiappa, R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics, Vol. 108 of Proceedings of Machine Learning Research, PMLR, Online, 2020, pp. 2207–2217.

[11] Y. LeCun, C. Cortes, C. Burges, MNIST handwritten digit database, ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist 2 (2010).

[12] J. Duchi, E. Hazan, Y. Singer, Adaptive Subgradient Methods for Online Learning and Stochastic Optimization 39.