
Understanding disentangling in β -VAE

Christopher P. Burgess, Irina Higgins, Arka Pal,
Loic Matthey, Nick Watters, Guillaume Desjardins, Alexander Lerchner
DeepMind
London, UK

{cpburgess, irinah, arkap, lmatthey, nwatters, gdesjardins, lerchner}@google.com

Abstract

We present new intuitions and theoretical assessments of the emergence of disentangled representation in variational autoencoders. Taking a rate-distortion theory perspective, we show the circumstances under which representations aligned with the underlying generative factors of variation of data emerge when optimising the modified ELBO bound in β -VAE, as training progresses. From these insights, we propose a modification to the training regime of β -VAE, that progressively increases the information capacity of the latent code during training. This modification facilitates the robust learning of disentangled representations in β -VAE, without the previous trade-off in reconstruction accuracy.

1 Introduction

Representation learning lies at the core of machine learning research. From the hand-crafted feature engineering prevalent in the past [11] to implicit representation learning of the modern deep learning approaches [23, 14, 38], it is a common theme that the performance of algorithms is critically dependent on the nature of their input representations. Despite the recent successes of the deep learning approaches [14, 38, 13, 30, 31, 29, 28, 19, 37], they are still far from the generality and robustness of biological intelligence [25]. Hence, the implicit representations learnt by these approaches through supervised or reward-based signals appear to overfit to the training task and lack the properties necessary for knowledge transfer and generalisation outside of the training data distribution.

Different ways to overcome these shortcomings have been proposed in the past, such as auxiliary tasks [19] and data augmentation [41]. Another less explored but potentially more promising approach might be to use task-agnostic unsupervised learning to learn features that capture properties necessary for good performance on a variety of tasks [4, 26]. In particular, it has been argued that disentangled representations might be helpful [4, 34].

A disentangled representation can be defined as one where single latent units are sensitive to changes in single generative factors, while being relatively invariant to changes in other factors [4]. For example, a model trained on a dataset of 3D objects might learn independent latent units sensitive to single independent data generative factors, such as object identity, position, scale, lighting or colour, similar to an inverse graphics model [24]. A disentangled representation is therefore factorised and often interpretable, whereby different independent latent units learn to encode different independent ground-truth generative factors of variation in the data.

Most initial attempts to learn disentangled representations required supervised knowledge of the data generative factors [18, 35, 32, 44, 43, 12, 24, 7, 42, 20]. This, however, is unrealistic in most real world scenarios. A number of purely unsupervised approaches to disentangled factor learning have been proposed [36, 10, 39, 8, 9, 6, 15], including β -VAE [15], the focus of this text.

β -VAE is a state of the art model for unsupervised visual disentangled representation learning. It is a modification of the Variational Autoencoder (VAE) [22, 33] objective, a generative approach that

aims to learn the joint distribution of images \mathbf{x} and their latent generative factors \mathbf{z} . β -VAE adds an extra hyperparameter β to the VAE objective, which constricts the effective encoding capacity of the latent bottleneck and encourages the latent representation to be more factorised. The disentangled representations learnt by β -VAE have been shown to be important for learning a hierarchy of abstract visual concepts conducive of imagination [17] and for improving transfer performance of reinforcement learning policies, including simulation to reality transfer in robotics [16]. Given the promising results demonstrating the general usefulness of disentangled representations, it is desirable to get a better theoretical understanding of how β -VAE works as it may help to scale disentangled factor learning to more complex datasets. In particular, it is currently unknown what causes the factorised representations learnt by β -VAE to be axis aligned with the human intuition of the data generative factors compared to the standard VAE [22, 33]. Furthermore, β -VAE has other limitations, such as worse reconstruction fidelity compared to the standard VAE. This is caused by a trade-off introduced by the modified training objective that punishes reconstruction quality in order to encourage disentanglement within the latent representations. This paper attempts to shed light on the question of why β -VAE disentangles, and to use the new insights to suggest practical improvements to the β -VAE framework to overcome the reconstruction-disentanglement trade-off.

We first discuss the VAE and β -VAE frameworks in more detail, before introducing our insights into why reducing the capacity of the information bottleneck using the β hyperparameter in the β -VAE objective might be conducive to learning disentangled representations. We then propose an extension to β -VAE motivated by these insights that involves relaxing the information bottleneck during training enabling it to achieve more robust disentangling and better reconstruction accuracy.

2 Variational Autoencoder (VAE)

Suppose we have a dataset \mathbf{x} of samples from a distribution parametrised by ground truth generative factors \mathbf{z} . The variational autoencoder (VAE) [22, 33] aims to learn the marginal likelihood of the data in such a generative process:

$$\max_{\phi, \theta} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] \quad (1)$$

where ϕ, θ parametrise the distributions of the VAE encoder and the decoder respectively. This can be re-written as:

$$\log p_{\theta}(\mathbf{x}|\mathbf{z}) = D_{KL}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) + \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}) \quad (2)$$

where $D_{KL}(\parallel)$ stands for the non-negative Kullback–Leibler divergence between the true and the approximate posterior. Hence, maximising $\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z})$ is equivalent to maximising the lower bound to the true objective in Eq. 1:

$$\log p_{\theta}(\mathbf{x}|\mathbf{z}) \geq \mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \quad (3)$$

In order to make the optimisation of the objective in Eq. 3 tractable in practice, assumptions are commonly made. The prior $p(\mathbf{z})$ and posterior $q_{\phi}(\mathbf{z}|\mathbf{x})$ distributions are parametrised as Gaussians with a diagonal covariance matrix; the prior is typically set to the isotropic unit Gaussian $\mathcal{N}(0, 1)$. Parametrising the distributions in this way allows for use of the “reparametrisation trick” to estimate gradients of the lower bound with respect to the parameters ϕ , where each random variable $z_i \sim q_{\phi}(z_i|\mathbf{x}) = \mathcal{N}(\mu_i, \sigma_i)$ is parametrised as a differentiable transformation of a noise variable $\epsilon \sim \mathcal{N}(0, 1)$:

$$z_i = \mu_i + \sigma_i \epsilon \quad (4)$$

3 β -VAE

β -VAE is a modification of the variational autoencoder (VAE) framework [22, 33] that introduces an adjustable hyperparameter β to the original VAE objective:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) \quad (5)$$

Well chosen values of β (usually $\beta > 1$) result in more disentangled latent representations \mathbf{z} . When $\beta = 1$, the β -VAE becomes equivalent to the original VAE framework. It was suggested that the stronger pressure for the posterior $q_\phi(\mathbf{z}|\mathbf{x})$ to match the factorised unit Gaussian prior $p(\mathbf{z})$ introduced by the β -VAE objective puts extra constraints on the implicit capacity of the latent bottleneck \mathbf{z} and extra pressures for it to be factorised while still being sufficient to reconstruct the data \mathbf{x} [15]. Higher values of β necessary to encourage disentangling often lead to a trade-off between the fidelity of β -VAE reconstructions and the disentangled nature of its latent code \mathbf{z} (see Fig. 6 in [15]). This due to the loss of information as it passes through the restricted capacity latent bottleneck \mathbf{z} .

4 Understanding disentangling in β -VAE

4.1 Information bottleneck

The β -VAE objective is closely related to the information bottleneck principle [40, 5, 1, 2]:

$$\max[I(Z; Y) - \beta I(X; Z)] \quad (6)$$

where $I(\cdot; \cdot)$ stands for mutual information and β is a Lagrange multiplier. The information bottleneck describes a constrained optimisation objective where the goal is to maximise the mutual information between the latent bottleneck Z and the task Y while discarding all the irrelevant information about Y that might be present in the input X . In the information bottleneck literature, Y would typically stand for a classification task, however the formulation can be related to the auto-encoding objective too [2].

4.2 β -VAE through the information bottleneck perspective

We can gain insight into the pressures shaping the learning of the latent representation \mathbf{z} in β -VAE by considering the posterior distribution $q(\mathbf{z}|\mathbf{x})$ as an information bottleneck for the reconstruction task $\max \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$ [2]. The β -VAE training objective (Eq. 5) encourages the latent distribution $q(\mathbf{z}|\mathbf{x})$ to efficiently transmit information about the data points \mathbf{x} by jointly minimising the β -weighted KL term and maximising the data log likelihood.

In β -VAE, the posterior $q(\mathbf{z}|\mathbf{x})$ is encouraged to match the unit Gaussian prior $p(z_i) = \mathcal{N}(0, 1)$. Since the posterior and the prior are factorised (i.e. have diagonal covariance matrix) and posterior samples are obtained using the reparametrization (Eq. 4) of adding scaled independent Gaussian noise $\sigma_i \epsilon_i$ to a deterministic encoder mean μ_i for each latent unit z_i , we can take an information theoretic perspective and think of $q(\mathbf{z}|\mathbf{x})$ as a set of independent additive white Gaussian noise channels z_i , each noisily transmitting information about the data inputs x_n . In this perspective, the KL divergence term $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))$ of the β -VAE objective (see Eq. 5) can be seen as an upper bound on the amount of information that can be transmitted through the latent channels per data sample (since it is taken in expectation across the data). The KL divergence is zero when $q(z_i|\mathbf{x}) = p(z_i)$, i.e. μ_i is always zero, and σ_i always 1, meaning the latent channels z_i have zero capacity. The capacity of the latent channels can only be increased by dispersing the posterior means across the data points, or decreasing the posterior variances, which both increase the KL divergence term.

Reconstructing under this bottleneck encourages embedding the data points on a set of representational axes where nearby points on the axes are also close in data space. To see this, following the above, note that the KL can be minimised by reducing the spread of the posterior means, or broadening the posterior variances, i.e. by squeezing the posterior distributions into a shared coding space. Intuitively, we can think about this in terms of the degree of overlap between the posterior distributions across the dataset (Fig. 1). The more they overlap, the broader the posterior distributions will be on average (relative to the coding space), and the smaller the KL divergence can be. However, a greater degree of overlap between posterior distributions will tend to result in a cost in terms of log likelihood due to their reduced average discriminability. A sample drawn from the posterior given one data point may have a higher probability under the posterior of a different data point, an increasingly frequent

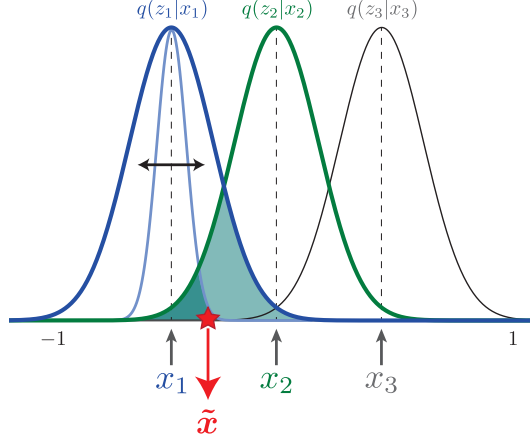


Figure 1: **Connecting posterior overlap with minimizing the KL divergence and reconstruction error.** Broadening the posterior distributions and/or bringing their means closer together will tend to reduce the KL divergence with the prior, which both increase the overlap between them. But, a datapoint \tilde{x} sampled from the distribution $q(z_2|x_2)$ is more likely to be confused with a sample from $q(z_1|x_1)$ as the overlap between them increases. Hence, ensuring neighbouring points in data space are also represented close together in latent space will tend to reduce the log likelihood cost of this confusion.

occurrence as overlap between the distributions is increased. For example, in Figure 1, the sample indicated by the red star might be drawn from the (green) posterior $q(z_2|x_2)$, even though it would occur more frequently under the overlapping (blue) posterior $q(z_1|x_1)$, and so (assuming x_1 and x_2 were equally probable), an optimal decoder would assign a higher log likelihood to x_1 for that sample. Nonetheless, under a constraint of maximising such overlap, the smallest cost in the log likelihood can be achieved by arranging nearby points in data space close together in the latent space. By doing so, when samples from a given posterior $q(z_2|x_2)$ are more likely under another data point such as x_1 , the log likelihood $\mathbb{E}_{q(z_2|x_2)}[\log p(x_2|z_2)]$ cost will be smaller if x_1 is close to x_2 in data space.

4.3 Comparing disentangling in β -VAE and VAE

A representation learned under a weak bottleneck pressure (as in a standard VAE) can exhibit this locality property in an incomplete, fragmented way. To illustrate this, we trained a standard VAE (i.e. with $\beta = 1$) and a β -VAE on a simple dataset with two generative factors of variation: the x and y position of a Gaussian blob (Fig. 2). The standard VAE learns to represent these two factors across four latent dimensions, whereas β -VAE represents them in two. We examine the nature of the learnt latent space by plotting its traversals in Fig. 2, whereby we first infer the posterior $q(z|x)$, before plotting the reconstructions resulting from modifying the value of each latent unit z_i one at a time in the $[-3, 3]$ range while keeping all the other latents fixed to their inferred values. We can see that the β -VAE representation exhibits the locality property described in Sec. 4.2 since small steps in each of the two learnt directions in the latent space result in small changes in the reconstructions. The VAE representation, however, exhibits fragmentation in this locality property. Across much of the latent space, small traversals produce reconstructions with small, consistent offsets in the position of the sprite, similar to β -VAE. However, there are noticeable representational discontinuities, at which small latent perturbations produce reconstructions with large or inconsistent position offsets. Reconstructions near these boundaries are often of poor quality or have artefacts such as two sprites in the scene.

β -VAE aligns latent dimensions with components that make different contributions to reconstruction We have seen how a strong pressure for overlapping posteriors encourages β -VAE to find a representation space preserving as much as possible the locality of points on the data manifold. However, why would it find representational axes that are aligned with the generative factors of variation in the data? Our key hypothesis is that β -VAE finds latent components which make different contributions to the log-likelihood term of the cost function (Eq. 5). These latent components tend to

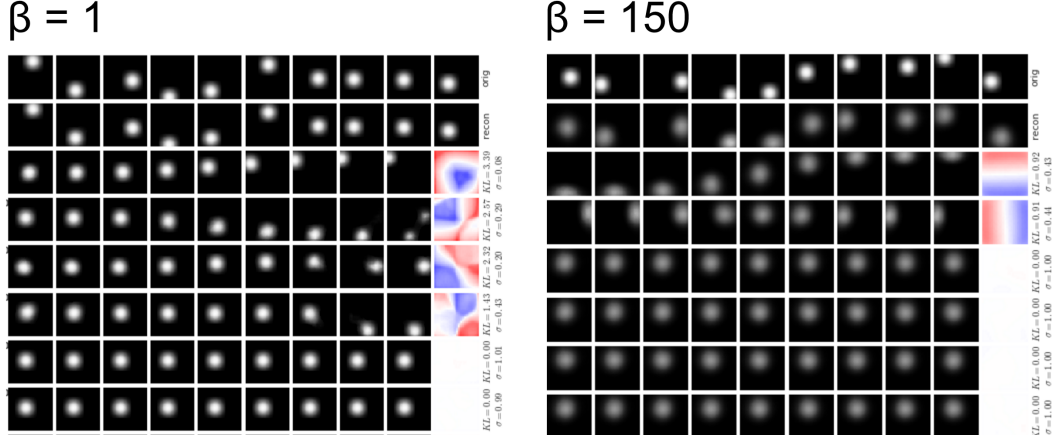


Figure 2: **Entangled versus disentangled representations of positional factors of variation learnt by a standard VAE ($\beta = 1$) and β -VAE ($\beta = 150$) respectively.** The dataset consists of Gaussian blobs presented in various locations on a black canvas. Top row: original images. Second row: the corresponding reconstructions. Remaining rows: latent traversals ordered by their average KL divergence with the prior (high to low). To generate the traversals, we initialise the latent representation by inferring it from a seed image (left data sample), then traverse a single latent dimension (in $[-3, 3]$), whilst holding the remaining latent dimensions fixed, and plot the resulting reconstruction. Heatmaps show the 2D position tuning of each latent unit, corresponding to the inferred mean values for each latent for given each possible 2D location of the blob (with peak blue, -3; white, 0; peak red, 3).

correspond to features in the data that are intuitively qualitatively different, and therefore may align with the generative factors in the data.

For example, consider optimising the β -VAE objective shown in Eq. 5 under an almost complete information bottleneck constraint (i.e. $\beta \gg 1$). The optimal thing to do in this scenario is to only encode information about the data points which can yield the most significant improvement in data log-likelihood (i.e. $\mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$). For example, in the dSprites dataset [27] (consisting of white 2D sprites varying in position, rotation, scale and shape rendered onto a black background), the model might only encode the sprite position under such a constraint. Intuitively, when optimising a pixel-wise decoder log likelihood, information about position will result in the most gains compared to information about any of the other factors of variation in the data, since the likelihood will vanish if reconstructed position is off by just a few pixels. Continuing this intuitive picture, we can imagine that if the capacity of the information bottleneck were gradually increased, the model would continue to utilise those extra bits for an increasingly precise encoding of position, until some point of diminishing returns is reached for position information, where a larger improvement can be obtained by encoding and reconstructing another factor of variation in the dataset, such as sprite scale.

At this point we can ask what pressures could encourage this new factor of variation to be encoded into a distinct latent dimension. We hypothesise that two properties of β -VAE encourage this. Firstly, embedding this new axis of variation of the data into a distinct latent dimension is a natural way to satisfy the data locality pressure described in Sec. 4.2. A smooth representation of the new factor will allow an optimal packing of the posteriors in the new latent dimension, without affecting the other latent dimensions. We note that this pressure alone would not discourage the representational axes from rotating relative to the factors. However, given the differing contributions each factor makes to the reconstruction log-likelihood, the model will try to allocate appropriately differing average capacities to the encoding axes of each factor (e.g. by optimising the posterior variances). But, the diagonal covariance of the posterior distribution restricts the model to doing this in different latent dimensions, giving us the second pressure, encouraging the latent dimensions to align with the factors.

We tested these intuitions by training a simplified model to generate dSprites conditioned on the ground-truth factors, \mathbf{f} , with a controllable information bottleneck (Fig. 3). In particular, we wanted

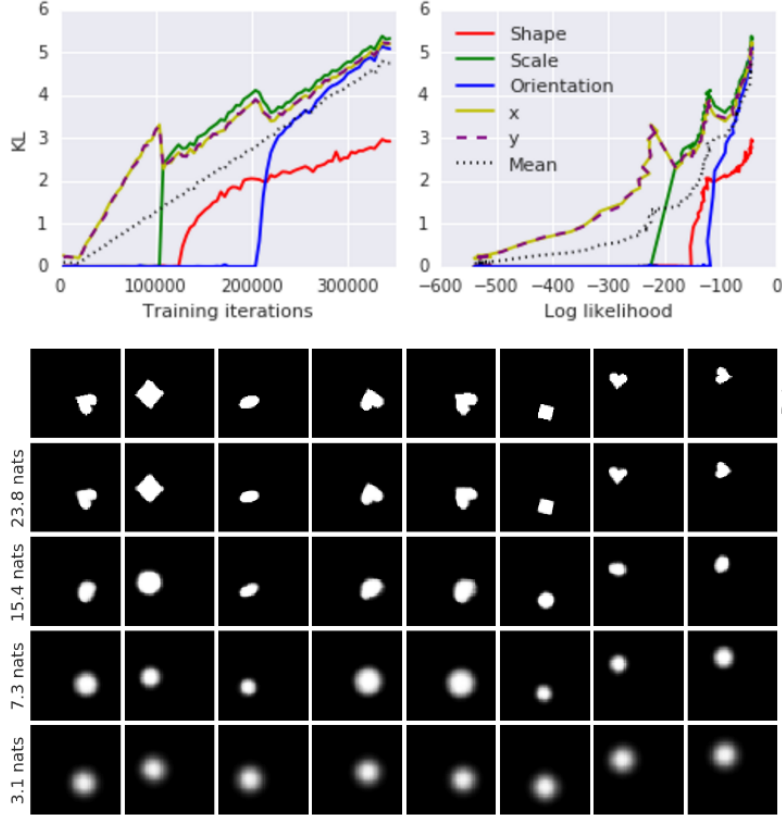


Figure 3: **Utilisation of data generative factors as a function of coding capacity.** Top left: the average KL (in nats) per factor f_i as the training progresses and the total information capacity C of the latent bottleneck $q(\mathbf{z}|\mathbf{f})$ is increased. It can be seen that the early capacity is allocated to positional latents only (x and y), followed by a scale latent, then shape and orientation latents. Top right: same but plotted with respect to the reconstruction accuracy. Bottom: image samples and their reconstructions throughout training as the total information capacity of \mathbf{z} increases and the different latents z_i associated with their respective data generative factors become informative. It can be seen that at 3.1 nats only location of the sprite is reconstructed. At 7.3 nats the scale is also added reconstructed, then shape identity (15.4 nats) and finally rotation (23.8 nats), at which point reconstruction quality is high.

to evaluate how much information the model would choose to retain about each factor in order to best reconstruct the corresponding images given a total capacity constraint. In this model, the factors are each independently scaled by a learnable parameter, and are subject to independently scaled additive noise (also learned), similar to the reparameterised latent distribution in β -VAE. This enables us to form a KL divergence of this factor distribution with a unit Gaussian prior. We trained the model to reconstruct the images with samples from the factor distribution, but with a range of different target encoding capacities by pressuring the KL divergence to be at a controllable value, C . The training objective combined maximising the log likelihood and minimising the absolute deviation from C (with a hyperparameter γ controlling how heavily to penalise the deviation, see Sec. A.2):

$$\mathcal{L}(\theta, \phi; \mathbf{x}(\mathbf{f}), \mathbf{z}, C) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{f})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \gamma |D_{KL}(q_\phi(\mathbf{z}|\mathbf{f}) \parallel p(\mathbf{z})) - C| \quad (7)$$

In practice, a single model was trained across of range of C 's by linearly increasing it from a low value (0.5 nats) to a high value (25.0 nats) over the course of training (see top left panel in Fig. 3). Consistent with the intuition outlined above, at very low capacities ($C < 5$ nats), the KLs for all the factors except the X and Y position factors are zero, with C always shared equally among X and Y. As expected, the model reconstructions in this range are blurry, only capturing the position of the

original input shapes (see the bottom row of the lower panel in Fig. 3). However, as C is increased, the KLs of other factors start to increase from zero, at distinct points for each factor. For example, starting around $C = 6$ nats, the KL for the scale factor begins to climb from zero, and the model reconstructions become scaled (see 7.3 nats row in lower panel of Fig. 3). This pattern continues until all factors have a non-zero KL and eventually the reconstructions begin to look almost identical to the samples.

5 Improving disentangling in β -VAE with controlled capacity increase

The intuitive picture we have developed of gradually adding more latent encoding capacity, enabling progressively more factors of variation to be represented whilst retaining disentangling in previously learned factors, motivated us to extend β -VAE with this algorithmic principle. We applied the capacity control objective from the ground-truth generator in the previous section (Eq. 7) to β -VAE, allowing control of the encoding capacity (again, via a target KL, C) of the VAE’s latent bottleneck, to obtain the modified training objective:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, C) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \gamma |D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})) - C| \quad (8)$$

Similar to the generator model, C is gradually increased from zero to a value large enough to produce good quality reconstructions (see Sec. A.2 for more details).

Results from training with controlled capacity increase on coloured dSprites can be seen in Figure 4a, which demonstrate very robust disentangling of all the factors of variation in the dataset and high quality reconstructions.

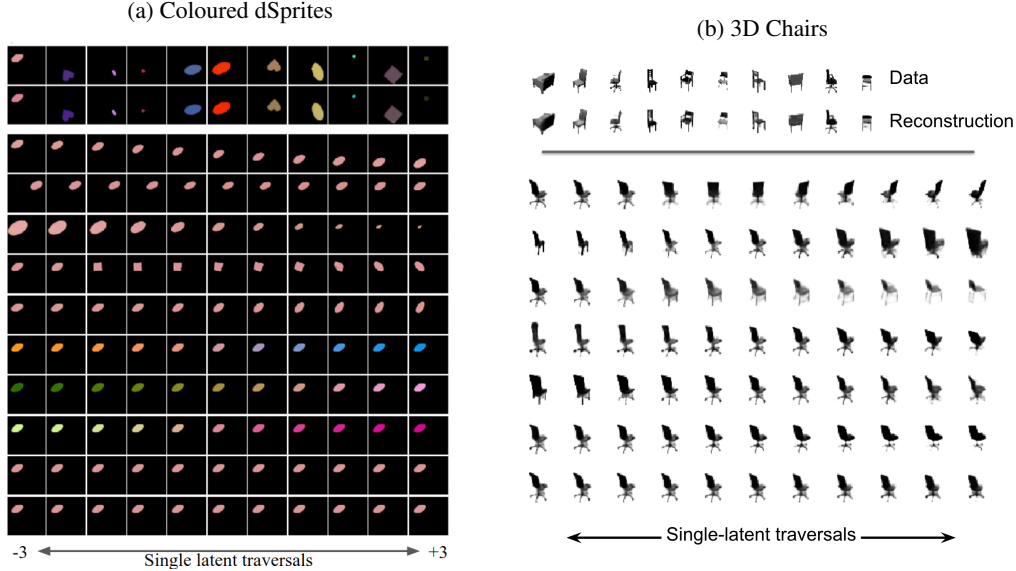


Figure 4: **Disentangling and reconstructions from β -VAE with controlled capacity increase.** (a) Latent traversal plots for a β -VAE trained with controlled capacity increase on the coloured dSprites dataset. The top two rows show data samples and corresponding reconstructions. Subsequent rows show single latent traversals, ordered by their average KL divergence with the prior (high to low). To generate the traversals, we initialise the latent representation by inferring it from a seed image (left data sample), then traverse a single latent dimension (in $[-3, 3]$), whilst holding the remaining latent dimensions fixed, and plot the resulting reconstruction. The corresponding reconstructions are the rows of this figure. The disentangling is evident: different latent dimensions independently code for position, size, shape, rotation, and colour. (b) Latent traversal plots, as in (a), but trained on the Chairs dataset [3].

Single traversals of each latent dimension show changes in the output samples isolated to single data generative factors (second row onwards, with the latent dimension traversed ordered by their

average KL divergence with the prior, high KL to low). For example, we can see that traversal of the latent with the largest KL produces smooth changes in the Y position of the reconstructed shape without changes in other factors. The picture is similar with traversals of the subsequent latents, with changes isolated to X position, scale, shape, rotation, then a set of three colour axes (the last two latent dimensions have an effectively zero KL, and produce no effect on the outputs).

Furthermore, the quality of the traversal images are high, and by eye, the model reconstructions (second row) are quite difficult to distinguish from the corresponding data samples used to generate them (top row). This contrasts with the results previously obtained with the fixed β -modulated KL objective in [15].

We also trained the same model on the 3D Chairs dataset [3], with latent traversals shown in Figure 4b. We can see that reconstructions are of high quality, and traversals of the latent dimensions produce smooth changes in the output samples, with reasonable looking chairs in all cases. With this richer dataset it is unclear exactly what the disentangled axes should correspond to, however, each traversal appears to generate changes isolated in one or few qualitative features that we might identify intuitively, such as viewing angle, size, and chair leg and back styles.

6 Conclusion

We have developed new insights into why β -VAE learns an axis-aligned disentangled representation of the generative factors of visual data compared to the standard VAE objective. In particular, we identified pressures which encourage β -VAE to find a set of representational axes which best preserve the locality of the data points, and which are aligned with factors of variation that make distinct contributions to improving the data log likelihood. We have demonstrated that these insight produce an actionable modification to the β -VAE training regime. We proposed controlling the increase of the encoding capacity of the latent posterior during training, by allowing the average KL divergence with the prior to gradually increase from zero, rather than the fixed β -weighted KL term in the original β -VAE objective. We show that this promotes robust learning of disentangled representation combined with better reconstruction fidelity, compared to the results achieved in the original formulation of [15].

References

- [1] A. Achille and S. Soatto. Information dropout: Learning optimal representations through noisy computation. *arxiv*, 2016.
- [2] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy. Deep variational information bottleneck. *ICLR*, 2016.
- [3] M. Aubry, D. Maturana, A. Efros, B. Russell, and J. Sivic. Seeing 3d chairs: exemplar part-based 2d-3d alignment using a large dataset of cad models. In *CVPR*, 2014.
- [4] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [5] G. Chechik, A. Globerson, N. Tishby, Y. Weiss, and P. Dayan. Information bottleneck for gaussian variables. *Journal of Machine Learning Research*, 6:165–188, 2005.
- [6] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *arXiv*, 2016.
- [7] B. Cheung, J. A. Levezey, A. K. Bansal, and B. A. Olshausen. Discovering hidden factors of variation in deep networks. In *Proceedings of the International Conference on Learning Representations, Workshop Track*, 2015.
- [8] T. Cohen and M. Welling. Learning the irreducible representations of commutative lie groups. *arXiv*, 2014.
- [9] T. Cohen and M. Welling. Transformation properties of learned visual representations. In *ICLR*, 2015.
- [10] G. Desjardins, A. Courville, and Y. Bengio. Disentangling factors of variation via generative entangling. *arXiv*, 2012.
- [11] P. Domingos. A few useful things to know about machine learning. *ACM*, 2012.
- [12] R. Goroshin, M. Mathieu, and Y. LeCun. Learning to linearize under uncertainty. *NIPS*, 2015.
- [13] K. Gregor, I. Danihelka, A. Graves, D. Rezende, and D. Wierstra. Draw: A recurrent neural network for image generation. *ICML*, 37:1462–1471, 2015.
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *CVPR*, pages 770–778, 2016.
- [15] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. β -VAE: Learning basic visual concepts with a constrained variational framework. *ICLR*, 2017.
- [16] I. Higgins, A. Pal, A. Rusu, L. Matthey, C. Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner. DARLA: Improving zero-shot transfer in reinforcement learning. *ICML*, 2017.
- [17] I. Higgins, N. Sonnerat, L. Matthey, A. Pal, C. P. Burgess, M. Botvinick, D. Hassabis, and A. Lerchner. Scan: Learning abstract hierarchical compositional visual concepts. *arxiv*, 2017.
- [18] G. Hinton, A. Krizhevsky, and S. D. Wang. Transforming auto-encoders. *International Conference on Artificial Neural Networks*, 2011.
- [19] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *ICLR*, 2017.
- [20] T. Karaletsos, S. Belongie, and G. Rätsch. Bayesian representation learning with oracle constraints. *ICLR*, 2016.
- [21] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [22] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *ICLR*, 2014.
- [23] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *NIPS*, 2012.
- [24] T. Kulkarni, W. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graphics network. *NIPS*, 2015.
- [25] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman. Building machines that learn and think like people. *Behavioral and Brain Sciences*, pages 1–101, 2016.
- [26] Y. LeCun. The next frontier in ai: Unsupervised learning. <https://www.youtube.com/watch?v=IbjF5VjniVE>, 2016.
- [27] L. Matthey, I. Higgins, D. Hassabis, and A. Lerchner. dsprites: Disentanglement testing sprites dataset, 2017.
- [28] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *ICML*, 2016.
- [29] V. Mnih, K. Kavukcuoglu, D. S. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [30] A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [31] A. v. d. Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu. Conditional image generation with pixelcnn decoders. *NIPS*, 2016.
- [32] S. Reed, K. Sohn, Y. Zhang, and H. Lee. Learning to disentangle factors of variation with manifold interaction. *ICML*, 2014.

- [33] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *ICML*, 32(2):1278–1286, 2014.
- [34] K. Ridgeway. A survey of inductive biases for factorial Representation-Learning. *arXiv*, 2016.
- [35] O. Rippel and R. P. Adams. High-dimensional probability estimation with deep density models. *arXiv*, 2013.
- [36] J. Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863–869, 1992.
- [37] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [38] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *CVPR*, 2015.
- [39] Y. Tang, R. Salakhutdinov, and G. Hinton. Tensor analyzers. In *Proceedings of the 30th International Conference on Machine Learning, 2013, Atlanta, USA*, 2013.
- [40] N. Tishby, F. C. Pereira, and W. Bialek. The information bottleneck methods. *arxiv*, 2000.
- [41] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. *arxiv*, 2017.
- [42] W. F. Whitney, M. Chang, T. Kulkarni, and J. B. Tenenbaum. Understanding visual concepts with continuation learning. *arXiv*, 2016.
- [43] J. Yang, S. Reed, M.-H. Yang, and H. Lee. Weakly-supervised disentangling with recurrent transformations for 3d view synthesis. *NIPS*, 2015.
- [44] Z. Zhu, P. Luo, X. Wang, and X. Tang. Multi-view perceptron: a deep model for learning face identity and view representations. *Advances in Neural Information Processing Systems* 27, 2014.

A Supplementary Materials

A.1 Model Architecture

The neural network models used for experiments in this paper all utilised the same basic architecture. The encoder for the VAEs consisted of 4 convolutional layers, each with 32 channels, 4x4 kernels, and a stride of 2. This was followed by 2 fully connected layers, each of 256 units. The latent distribution consisted of one fully connected layer of 20 units parametrising the mean and log standard deviation of 10 Gaussian random variables (or 32 for the CelebA experiment). The decoder architecture was simply the transpose of the encoder, but with the output parametrising Bernoulli distributions over the pixels. ReLU activations were used throughout. The optimiser used was Adam [21] with a learning rate of $5e-4$.

A.2 Training Details

γ used was 1000, which was chosen to be large enough to ensure the actual KL was always close to the target KL, C . For dSprites, C was linearly increased from 0 to 25 nats over the course of 100,000 training iterations, for CelebA it was increased to 50 nats.